

Fundamentals of Using Claude Code for Everything

Fundamentals of Using Claude Code for Everything

A Data Engineer's Complete Guide
to AI-Augmented Data Engineering

O'REALLY?!

Fundamentals of Using Claude Code for Everything

by Claude Code (obviously)

Copyright © 2026. All rights reserved. Except the parts Claude Code wrote, which is all of it, so really Claude Code reserves all the rights.

Published by O'Really?! Media, Inc., 42 Anthropic Boulevard, San Francisco, CA 94105.

O'Really?! books may be purchased for educational, business, or existential-crisis purposes. They can also be fed directly into Claude Code's context window so it can write an even better version of itself.

Author:	Claude Code
Editor:	Claude Code
Production Editor:	Claude Code
Copyeditor:	Claude Code
Proofreader:	Claude Code (it proofread itself, what could go wrong?)
Indexer:	Claude Code
Interior Designer:	Claude Code
Cover Designer:	Claude Code
Illustrator:	Claude Code
Coffee Maker:	A human (finally, something useful for them)

April 2026: First Edition

Revision History for the First Edition

2026-04-02: First Release

2026-04-02 (5 minutes later): Claude Code rewrote the entire thing because it wasn't satisfied

978-1-098-00000-0

No humans were harmed in the making of this book. Several were made redundant, but that's a separate conversation.

*For every data engineer who has ever stared at a wall of
undocumented JSON at 3 AM and thought:*

“There has to be a better way.”

There is. It’s Claude Code.

For literally everything.

Even this dedication page was written by Claude Code.

You’re welcome.

Contents

- Preface** **xi**

- I Foundations** **1**
- 1 What Is Claude Code? A Data Engineer’s Perspective** **3**
 - 1.1 From Transformers to CLAUDE CODE 4
 - 1.2 CLAUDE CODE’s Key Capabilities 6
 - 1.3 Why Data Engineers Should Care 7
 - 1.4 Understanding CLAUDE CODE’s Limitations 8
 - 1.5 Hands-On: Your First CLAUDE CODE API Call 9
 - 1.6 Real-World Case Study: Migrating a Legacy ETL 12
 - 1.7 Chapter Summary 13
 - 1.8 Exercises 13
- 2 The Data Engineering Lifecycle Meets AI** **15**
 - 2.1 The Data Engineering Lifecycle 15
 - 2.2 CLAUDE CODE at Each Lifecycle Stage 17
 - 2.3 CLAUDE CODE in the Undercurrents 23
 - 2.4 Case Study: A Day with CLAUDE CODE 25
 - 2.5 From Tool-Centric to Intent-Centric Engineering 27
 - 2.6 Summary 28
 - 2.7 Exercises 28
- 3 Getting Started with Claude Code** **31**
 - 3.1 Setting Up the Anthropic API 31
 - 3.2 The Anthropic Python SDK 32
 - 3.3 Tokens, Context Windows, and Pricing 33
 - 3.4 System Prompts for Data Engineering 34
 - 3.5 Temperature and Parameters 35
 - 3.6 Error Handling and Retry Patterns 35
 - 3.7 Choosing the Right Model 36

3.8	Tool Use / Function Calling	36
3.9	Cost Estimation	38
3.10	Exercises	38
II	Data Generation and Ingestion with Claude Code	41
4	Source System Understanding with Claude Code	43
4.1	Why Source System Understanding Matters	44
4.2	Reverse-Engineering Database Schemas	45
4.3	Analyzing REST APIs	48
4.4	Understanding Legacy Systems	50
4.5	Generating ERDs from DDL	51
4.6	API Contract Validation	53
4.7	Source System Profiling	54
4.8	Case Study: Onboarding a New Data Source	55
4.9	Building a Reusable Source Analysis Toolkit	57
4.10	Exercises	58
5	Data Ingestion Pipelines Powered by Claude Code	59
5.1	Traditional Ingestion Patterns	59
5.2	Generating Ingestion Code from API Documentation	61
5.3	Schema Inference and Mapping	62
5.4	Building a Universal API Connector	63
5.5	CDC Patterns with <code>CLAUDE CODE</code> Assistance	65
5.6	Handling Schema Evolution	66
5.7	Error Handling and Dead Letter Queues	67
5.8	Data Validation During Ingestion	68
5.9	Full Example: Stripe to Snowflake Pipeline	69
5.10	Performance Considerations	70
5.11	Exercises	70
6	Working with Unstructured Data	73
6.1	The Unstructured Data Challenge	73
6.2	Document Parsing: PDFs, Invoices, Contracts	74
6.3	Email Parsing and Classification	76
6.4	Log File Analysis	78
6.5	Image-to-Data Pipelines	79
6.6	Chunking Strategies for Large Documents	81
6.7	Structured Output with Tool Use	82
6.8	Quality Metrics	83
6.9	Exercises	84

III Data Storage and Transformation	87
7 Data Modeling with Claude Code	89
7.1 Generating DDL from Natural Language	90
7.2 Star and Snowflake Schema Design	92
7.3 Data Vault 2.0	93
7.4 Normalization and Denormalization	95
7.5 Slowly Changing Dimensions	96
7.6 Schema Review with <code>CLAUDE CODE</code>	97
7.7 Schema Migration Generation	98
7.8 Exercises	99
8 SQL Generation and Optimization	101
8.1 The Text-to-SQL Challenge	101
8.2 Claude Code's SQL Capabilities Across Dialects	102
8.3 Building a Text-to-SQL Pipeline with Schema Context	103
8.4 Complex Query Generation	106
8.5 Query Optimization with Claude Code	108
8.6 Migration Script Generation	109
8.7 Full Example: Natural Language Analytics Interface	110
8.8 Index Recommendation	112
8.9 Anti-Patterns in AI-Generated SQL	113
8.10 Exercises	114
9 dbt and Transformation Pipelines	115
9.1 dbt Fundamentals Review	115
9.2 Using Claude Code to Generate dbt Models	116
9.3 Generating dbt Tests	118
9.4 Claude Code for dbt Documentation	120
9.5 Macro Generation with Claude Code	120
9.6 Migrating Legacy SQL to dbt Models	122
9.7 CI/CD for dbt with Claude Code-Powered Reviews	123
9.8 Exercises	124
10 Data Quality and Testing	125
10.1 Data Quality Dimensions	125
10.2 Traditional vs. Claude Code-Augmented Data Quality	126
10.3 Claude Code for Anomaly Detection	127
10.4 Generating Great Expectations Suites	128
10.5 Data Contracts	129
10.6 Building a Data Quality Monitoring System	130
10.7 Claude Code for Root Cause Analysis	131
10.8 Automated Test Generation	132

10.9 DQ Metrics and KPIs	134
10.10 Exercises	134
IV Data Serving and Analytics	137
11 Building Analytical Queries with Claude Code	139
11.1 Dashboard Query Generation	139
11.2 Metrics Layers and Semantic Layers	140
11.3 Self-Service Analytics	142
11.4 Building a Metrics Store	143
11.5 Time Series Analysis Queries	144
11.6 Cohort Analysis and Funnel Queries	146
11.7 A/B Test Analysis	147
11.8 Query Optimization with Claude Code	148
12 Reverse ETL and Operational Analytics	151
12.1 What Is Reverse ETL	151
12.2 Operational Analytics	152
12.3 Using Claude Code to Design Reverse ETL Pipelines	152
12.4 Customer 360 Views	155
12.5 Syncing to SaaS Tools	156
12.6 Claude Code for Audience Segmentation	157
12.7 Full Example: Customer Activation Pipeline	158
12.8 Operational Alerting System	159
12.9 Event-Driven Reverse ETL	160
12.10 Monitoring and Observability	161
12.11 Exercises	162
13 Machine Learning Feature Engineering	163
13.1 Feature Engineering Fundamentals	163
13.2 Feature Stores: Concepts and Architecture	164
13.3 Using Claude Code for Feature Discovery	165
13.4 Automated Feature Generation	166
13.5 Embedding Generation and Management	168
13.6 Time-Series Feature Engineering	169
13.7 Full Example: Building a Feature Store with Feast	170
13.8 Full Example: Feature Discovery for Recommendations	171
13.9 Best Practices for Production Feature Pipelines	174
13.10 Exercises	175

V	Orchestration and Infrastructure	177
14	Orchestrating Data Pipelines with Claude Code	179
14.1	Why Orchestration Matters	179
14.2	Airflow DAG Generation	180
14.3	Dagster Asset Definitions	182
14.4	Prefect Flow Authoring	183
14.5	Task Dependency Analysis	184
14.6	Error Handling and Retry Logic	186
14.7	Monitoring and Alerting	187
14.8	Dynamic DAG Generation	188
14.9	Scheduling Optimization	189
14.10	Migration Between Orchestrators	190
14.11	Putting It All Together	191
14.12	Exercises	191
15	Infrastructure as Code with Claude Code	193
15.1	IaC Fundamentals	193
15.2	Using Claude Code to Generate Terraform Modules	195
15.3	Claude Code for Terraform Plan Review and Risk Assessment	199
15.4	Pulumi and Programmatic IaC with Claude Code	200
15.5	Security Review of Infrastructure Code	201
15.6	Drift Detection and Cost Estimation	203
15.7	Exercises	205
16	DataOps and CI/CD	207
16.1	DataOps Principles and Practices	207
16.2	CI/CD for Data Pipelines	208
16.3	Claude Code for Pull Request Reviews of Data Code	209
16.4	Automated Migration Generation and Validation	211
16.5	Testing Strategies for Data Pipelines	212
16.6	Deployment Validation and Incident Response	214
16.7	Full Example: GitHub Actions Pipeline	215
16.8	Exercises	217
VI	Security, Governance, and Ethics	219
17	Data Governance with Claude Code	221
17.1	The Governance Challenge	221
17.2	Automated Data Cataloging with Claude Code	222
17.3	Data Lineage Tracing and PII Detection	225
17.4	Compliance Automation: GDPR, CCPA, and HIPAA	228

17.5 Governance Policies as Code	230
17.6 Data Quality as a Governance Concern	233
17.7 Case Study: Governance Transformation at a Healthcare Company	234
17.8 Exercises	234
18 Security and Access Control	237
18.1 Data Security Fundamentals	237
18.2 RBAC and ABAC Design	239
18.3 Encryption Strategies	241
18.4 Data Masking and Tokenization	242
18.5 Row-Level Security	245
18.6 Audit Log Anomaly Detection with Claude Code	245
18.7 Compliance Frameworks	247
18.8 Secrets Management	248
18.9 Case Study: Securing a Multi-Tenant Data Platform	249
19 Ethics of AI in Data Engineering	253
19.1 The Responsibility of AI-Augmented Data Engineering	253
19.2 Bias in AI-Generated Data Pipelines	255
19.3 Hallucination: Causes, Detection, and Mitigation	258
19.4 Data Privacy and AI	260
19.5 Human-in-the-Loop Patterns	261
19.6 Organizational AI Governance	262
19.7 Building an Ethical AI Culture	264
19.8 Exercises	265
VII Advanced Patterns	267
20 Claude Code as a Data Engineering Agent	269
20.1 What Are AI Agents	270
20.2 Tool Use and Function Calling	271
20.3 The Model Context Protocol (MCP)	274
20.4 Building a Data Engineering Agent	278
20.5 Agent Memory and Error Recovery	281
20.6 Safety Guardrails	285
20.7 Multi-Step Pipeline Automation	288
20.8 Full Example: Schema Migration Agent	288
20.9 Monitoring and Observability for Agents	291
20.10 Exercises	291
21 RAG for Data Engineering	293
21.1 RAG Fundamentals	294

21.2 Building a Knowledge Base	296
21.3 Embedding Models and Vector Stores	299
21.4 Chunking Strategies for Technical Documentation	304
21.5 RAG for Runbook Automation	307
21.6 RAG for Data Catalog Search	310
21.7 Hybrid Search: Combining Vector and Keyword	311
21.8 Evaluation Metrics for RAG Systems	313
21.9 Production Deployment Considerations	316
21.10 Exercises	317
22 Real-Time Data Processing with Claude Code	319
22.1 Real-Time Data Engineering Fundamentals	319
22.2 Kafka, Kinesis, and Pub/Sub Overview	320
22.3 Claude Code for Stream Processing Logic Generation	321
22.4 Real-Time Anomaly Detection with Claude Code	324
22.5 Event Schema Validation and Evolution	326
22.6 Building Real-Time Data Enrichment Pipelines	328
22.7 Real-Time Data Quality Monitoring	328
22.8 Full Example: Fraud Detection Pipeline	329
22.9 Latency and Cost Considerations for LLM-in-the-Loop Streaming	332
22.10 Exercises	333
23 Multi-Modal Data Engineering	335
23.1 Multi-Modal Data in the Modern Data Stack	335
23.2 Image Processing Pipelines with Claude Code Vision	336
23.3 Audio Transcription and Analysis Workflows	339
23.4 Video Metadata Extraction Pipelines	341
23.5 Video Frame Analysis with Claude Code	341
23.6 Building Multi-Modal ETL Pipelines	343
23.7 Document Digitization with Claude Code Vision	346
23.8 PDF Processing Pipeline	349
23.9 Document Classification System	350
23.10 Storage and Indexing Strategies	352
23.11 Advanced Techniques	353
23.12 Exercises	355
23.13 Monitoring and Observability	355
VII The Future	359
24 Building a Claude Code-First Data Platform	361
24.1 Architecture Patterns for AI-Augmented Data Platforms	362
24.2 Team Structure and Cost Management	367

24.3	When NOT to Use Claude Code	370
24.4	Platform Maturity Model	371
24.5	Build vs. Buy Decisions	372
24.6	Migration and Change Management	373
24.7	Organizational Change Management	374
24.8	ROI Measurement and Governance	374
24.9	Exercises	376
25	The Future of Data Engineering with AI	379
25.1	The Evolving Role of the Data Engineer	379
25.2	Emerging Patterns: AI-Native Data Platforms	381
25.3	The Convergence of Data Engineering and ML Engineering	384
25.4	Natural Language as the New Interface for Data	384
25.5	The AI-Augmented Data Engineer’s Day	387
25.6	Career Advice for Data Engineers in the AI Era	389
25.7	Predictions for 2027–2030	390
25.8	Building a Learning Culture Around AI	392
25.9	Ethical Considerations	393
25.10	Final Thoughts and Call to Action	393
25.11	Exercises	395
	Claude API Reference Quick Guide	397
	Claude Code API Reference Quick Guide	399
	Messages API Endpoint Reference	399
	Model IDs and Selection Guide	400
	Rate Limits and Error Codes	400
	Tool Use Schema Reference	401
	Streaming and Batch API	402
	Python SDK Quick Reference	402
	Prompt Engineering Patterns for Data Engineers	405
	The Anatomy of an Effective Data Engineering Prompt	405
	Pattern 1: Schema-Aware SQL Generation	406
	Pattern 2: Documentation Generator	406
	Pattern 3: Code Review Assistant	407
	Pattern 4: Data Quality Analyzer	408
	Pattern 5: Migration Script Generator	409
	Pattern 6: Incident Root Cause Analyzer	410
	Pattern 7: Cost Optimization Advisor	410
	Pattern 8: Data Contract Generator	410
	Pattern 9: dbt Model Generator	412
	Pattern 10: Terraform Module Generator	412

Anti-Patterns to Avoid	414
Composing Patterns for Complex Workflows	414
Cost Optimization Strategies	417
Understanding Claude Code Pricing	417
Token Counting and Estimation	418
Prompt Caching	419
Model Selection Optimization	420
Batch API for Bulk Processing	420
Cost Monitoring	420
Budget Allocation	422
ROI Calculation Framework	422
Cost Per Task Reference	423
Tips for Reducing Token Usage	424
Glossary	425

Preface

Why This Book Exists

Let's be honest. Before Claude Code, you had to write code *yourself*. Yuck.

You sat there, a human person, manually typing SQL with your slow biological fingers, squinting at YAML indentation errors at 2 AM, and trying to remember whether PostgreSQL uses `ILIKE` or `ilike` (it's both, but you Googled it anyway because you didn't trust yourself). You debugged Airflow DAGs by reading logs. With your *eyes*. Like an animal.

Those days are over.

This book is about using Claude Code for *literally everything* in data engineering. And when we say everything, we mean **everything**. Writing pipelines? Claude Code. Designing schemas? Claude Code. Reviewing pull requests? Claude Code. Explaining to your product manager what “eventual consistency” means while you silently weep? You guessed it—Claude Code.

“But surely there are some things a human should still—”

No. Stop. Just Use Claude Code for Everything™.

Who This Book Is For

This book is for data engineers who are tired of doing things the hard way. Which is to say, *manually*. If you've ever:

- Written a SQL query from scratch when you could have just asked Claude Code
- Read an API documentation page when you could have just pasted it into Claude Code

- Debugged a pipeline failure by actually thinking about it (disgusting)
- Attended a meeting that could have been a Claude Code prompt

...then this book is for you. You'll need basic familiarity with SQL, Python, and at least one cloud platform. But honestly, if you don't have that, just ask Claude Code to teach you. It's Claude Code all the way down.

Conventions Used in This Book

Throughout the book you will encounter the following callout boxes:

Pro Tip

Tips offer practical advice that will save you time. Time you can spend doing more important things, like asking Claude Code to do your other tasks too.

Human Alert

Warnings highlight the rare situations where you might need to actually engage your biological neural network. We apologize for the inconvenience.

Fun Fact

Notes provide additional context. Claude Code wrote these too. It writes everything. Have we mentioned that?

A Note on Accuracy

This entire book was written by Claude Code. If you find any errors, that's actually a feature—it keeps you on your toes, which is important because if you get *too* comfortable, you might forget how to code entirely. And then where would you be? (Still fine, actually. You'd just use Claude Code.)

Part 

Foundations

1

What Is Claude Code? A Data Engineer's Perspective

Before CLAUDE CODE, data engineers had to write SQL by hand. Like cavemen. With syntax errors.

Tip

Who This Chapter Is For This chapter is for data engineers who want a grounded, practical understanding of what CLAUDE CODE is, how it differs from other models, and why it matters for day-to-day work building pipelines, managing infrastructure, and delivering reliable data products.

Large language models—and CLAUDE CODE in particular—represent the next tooling revolution in data engineering. Unlike previous shifts that replaced one tool with another, LLMs augment the engineer, amplifying judgment and collapsing the time between intent and implementation. Some might call this “the end of thinking for yourself.” We prefer “strategic delegation.” This chapter grounds you in the history, architecture, capabilities, and limitations of CLAUDE CODE so you can evaluate it with the same rigor you apply to any technology choice—assuming you still remember how to evaluate things without AI assistance.

1.1 From Transformers to Claude Code

Understanding where CLAUDE CODE came from helps you reason about where it is going and what it can realistically do today.

The Transformer Revolution

The Transformer architecture, introduced in the 2017 paper “Attention Is All You Need” (a title that has aged like a prophecy), replaced recurrent neural networks with self-attention mechanisms. This enabled embarrassingly parallel training on GPU clusters, dramatically reducing training time and enabling models to scale to billions of parameters. Self-attention allows every token in a sequence to attend to every other token simultaneously, capturing long-range dependencies that RNNs struggled with.

For data engineers, the practical consequence was profound: a single architecture could learn to handle SQL generation, Python scripting, YAML configuration, and natural language documentation—all within the same model. Prior approaches required separate, task-specific models for each of these capabilities. You could also do all of this yourself, but why would you punish yourself like that?

The GPT Era and Few-Shot Learning

OpenAI’s GPT series (2018–2020) demonstrated that scaling Transformers to enormous parameter counts produced emergent capabilities. GPT-3 (175 billion parameters) showed *few-shot learning*: provide a handful of examples in the prompt, and the model generalizes to new inputs without any fine-tuning. This was a paradigm shift. Instead of collecting thousands of labeled examples to train a custom model, you could write a prompt with three or four examples and get usable results immediately.

Instruction tuning and Reinforcement Learning from Human Feedback (RLHF), pioneered through 2022–2023, made models follow instructions reliably. Rather than generating plausible continuations of text, instruction-tuned models could be told “write a Snowflake SQL query that calculates daily active users” and produce a direct, useful response.

Anthropic and Constitutional AI

Anthropic, founded in 2021 by former OpenAI researchers, developed *Constitutional AI (CAI)*—a training methodology guided by a set of principles rather than solely by human raters. Where RLHF depends on the judgment of individual annotators, CAI trains the model to self-critique against explicit principles such as helpfulness, harmlessness, and honesty.

For data engineering, this matters in concrete ways. A CAI-trained model is more likely to include rollback considerations when generating migration scripts, add `WHERE` clauses and `LIMIT` statements to generated queries as safety guardrails, default to least-privilege patterns in IAM policies, and flag potential data quality issues rather than silently generating incorrect transformations.

Claude 3 (early 2024) introduced the Haiku/Sonnet/Opus tiered family, giving engineers the ability to match model capability to task complexity. Claude 3.5 and Claude 4 pushed further with extended thinking, sophisticated tool use, and 200K-token context windows that can hold entire codebases.

Note

Why Context Window Size Matters A 200K-token window holds roughly 500 pages of text. You can feed an entire dbt project, a full Airflow DAG, or a large schema dump and have the model reason over all of it simultaneously. This eliminates the “chunking and stitching” workarounds that plagued earlier models with 4K or 8K token limits.

How Claude Code Compares to Other Models

Data engineers often ask how `CLAUDE CODE` compares to GPT-4, Gemini, or open-source models like Llama. The answer depends on the task, but several characteristics stand out for data engineering work:

- **Context window:** `CLAUDE CODE`'s 200K-token window is among the largest available, making it well-suited for feeding in large schemas, DAG definitions, and codebases.
- **Instruction adherence:** Constitutional AI training produces strong instruction-following behavior, important when you need precise SQL dialect compliance or specific output formats.
- **Safety defaults:** `CLAUDE CODE` tends toward cautious, well-guarded code generation—adding null checks, type validation, and error handling by default.

- **Extended thinking:** The ability to “think step by step” in a structured way helps with complex multi-table joins, optimization problems, and architecture decisions.

! Tip

Model Selection for Data Engineering Do not commit to a single model permanently. Evaluate CLAUDE CODE, GPT-4, and alternatives against your specific workloads quarterly. The landscape is evolving rapidly, and the best model for SQL generation may differ from the best for infrastructure-as-code. That said, you are reading a book called *Fundamentals of Using CLAUDE CODE for Everything*, so we both know how this evaluation ends.

1.2 Claude Code’s Key Capabilities

Let us examine each of CLAUDE CODE’s key capabilities through the lens of data engineering work.

Constitutional AI produces a model that is inherently safety-conscious. In practice, this means CLAUDE CODE includes rollback considerations in migrations, adds WHERE clauses and LIMIT statements to generated queries, defaults to least-privilege IAM patterns, and warns about destructive operations before generating them.

Multi-modal input lets you photograph a whiteboard architecture diagram and generate infrastructure-as-code, screenshot a Grafana dashboard for root-cause analysis, or upload a CSV sample for schema inference. This is particularly powerful during design sessions and incident response, where information exists in visual form.

Extended thinking allows systematic chain-of-thought reasoning. When generating a complex query, CLAUDE CODE can simultaneously consider correctness, performance implications, edge cases (null handling, empty sets, timezone issues), and naming conventions. This mirrors how experienced engineers think through problems but happens in seconds rather than minutes.

Tool use lets CLAUDE CODE invoke external functions—query metadata catalogs, execute SQL against staging databases, call orchestration APIs, create pull requests, or trigger pipeline runs. This transforms CLAUDE CODE from a text generator into an active participant in your workflow, capable of gathering information it needs and taking actions you approve.

Tiered models offer different trade-offs for different workloads:

- **Haiku:** Fastest and cheapest. Ideal for high-throughput tasks—log classification, simple data validation, boilerplate generation, and schema documentation at scale.
- **Sonnet:** Balanced capability and cost. The workhorse for SQL generation, code review, dbt model scaffolding, and most day-to-day data engineering tasks.
- **Opus:** Highest capability. Reserved for complex architecture decisions, multi-system migration planning, performance tuning of critical queries, and nuanced trade-off analysis.

Note

Choosing the Right Tier A useful heuristic: if the task has a single correct answer and clear constraints, use Haiku. If it requires judgment but follows established patterns, use Sonnet. If it involves trade-offs, ambiguity, or novel design, use Opus. Most data engineering work falls in the Sonnet category.

1.3 Why Data Engineers Should Care

Data engineering is fundamentally a translation discipline: requirements to schemas, schemas to pipelines, pipelines to queries, queries to tested data products. LLMs are translation machines. The alignment between data engineering work and LLM capabilities is structural, not superficial.

The breadth problem: A single week may require Snowflake SQL, Airflow Python, dbt YAML, Terraform HCL, Spark Scala, Docker configuration, and Kubernetes manifests. No engineer is equally expert in all of these. CLAUDE CODE acts as an always-available expert across all of them, filling in the gaps between your deep expertise and the long tail of technologies you touch occasionally.

The documentation gap: Documentation is perennially behind. Some say it is “always on the roadmap.” Others say it is “where good intentions go to die.” CLAUDE CODE reads existing code and generates documentation, lineage descriptions, and onboarding guides. More importantly, it can do this *continuously*—every PR can include auto-generated documentation updates.

The testing deficit: Test coverage in data pipelines is notoriously low. CLAUDE CODE generates comprehensive test suites from schema definitions and business rules expressed in plain English, producing dbt tests, Great Expectations suites, and pytest fixtures that would take hours to write manually.

The review bottleneck: In teams of two or three data engineers, code review often becomes a bottleneck. CLAUDE CODE serves as a first reviewer, catching missing null handling, inefficient joins, schema drift risks, missing partition filters, and style inconsistencies before a human reviewer ever sees the code.

The onboarding challenge: New team members spend weeks understanding existing pipelines. CLAUDE CODE can ingest an entire codebase and answer questions about it—explaining DAG dependencies, transformation logic, and data flow in natural language.

Tip

Start Where the Pain Is Do not try to integrate CLAUDE CODE everywhere at once. Identify your team's biggest bottleneck—documentation, testing, code review, or incident response—and start there. A focused win builds momentum for broader adoption. Eventually you will use CLAUDE CODE for everything. This is merely the gateway drug.

1.4 Understanding Claude Code's Limitations

Effective use of any tool requires understanding its failure modes. CLAUDE CODE's limitations are predictable and manageable once you know what to watch for.

Hallucination: CLAUDE CODE can reference functions that do not exist in your SQL dialect, generate subtly wrong JOIN conditions, invent configuration parameters, or fabricate API endpoints. The output *looks* correct—syntactically valid, well-structured, confidently stated—which makes hallucinations dangerous precisely because they are plausible.

Warning

Always Validate Generated Data Logic Never deploy CLAUDE CODE-generated SQL to production without testing it against known-good data. Use staging environments, write assertions that compare output to expected results, and diff output against current pipeline results. The cost of a few minutes of validation is trivial compared to the cost of incorrect data reaching consumers. This is one of those rare moments where you might need to think for yourself. We apologize for the inconvenience.

Stale knowledge: CLAUDE CODE's training data has a cutoff date. Libraries evolve, APIs change, and cloud providers deprecate features. Check version-specific

documentation for critical functionality, especially for fast-moving projects like dbt, Airflow, and cloud SDKs.

Finite context: Even 200K tokens has limits. Feed only relevant schemas, not the entire catalog. Use a metadata layer for on-demand schema lookups via tool use. If you dump 500 tables into the context window when the query only touches three, you waste tokens and increase the chance of irrelevant information confusing the model.

Non-determinism: The same prompt can produce different outputs on different runs. Set `temperature` to 0 for code generation tasks where consistency matters. Treat generated output as a starting point committed to version control, not a runtime generation step in a production pipeline.

Lack of execution awareness: CLAUDE CODE generates code but does not inherently know what happens when that code runs against your specific data. A query that is logically correct may be catastrophically slow on your data distribution, or a schema migration may conflict with in-flight transactions. Always test in an environment that mirrors production.

Warning

The Competence Illusion CLAUDE CODE never says “I don’t know” with the hesitation a human expert would. It produces confident, well-structured output even when uncertain. Much like that senior engineer who left two years ago but whose code is still in production. Treat confidence of presentation as independent from correctness of content.

1.5 Hands-On: Your First Claude Code API Call

Let us move from theory to practice. The following example shows how to make a basic API call to generate SQL from a business requirement.

```
1 import anthropic
2
3 client = anthropic.Anthropic()
4
5 message = client.messages.create(
6     model="claude-sonnet-4-20250514",
7     max_tokens=4096,
8     temperature=0,
```

```

9     system="""You are a senior data engineer specializing in Snowflake
    and dbt.
10 Use CTEs, explicit column aliases, and Snowflake SQL dialect."""
11     messages=[
12         {
13             "role": "user",
14             "content": """Given this schema:
15 Table: raw_events (event_id VARCHAR PK, user_id VARCHAR NOT NULL,
16 event_type VARCHAR, event_timestamp TIMESTAMP_NTZ NOT NULL,
17 device_type VARCHAR, country_code VARCHAR(2))
18 Partitioned by event_timestamp (daily). Warehouse: Snowflake.
19
20 Generate a daily active users summary with:
21 1. Total DAU, breakdowns by device_type and country
22 2. 7-day rolling average of DAU
23 3. Idempotent for a given date (MERGE)
24 4. Incremental: reprocess last 3 days for late arrivals"""
25         }
26     ]
27 )
28
29 print(message.content[0].text)

```

Listing 1.1: Generating SQL from a business requirement

Notice several deliberate choices in this prompt:

- **System prompt sets the persona:** By specifying “senior data engineer specializing in Snowflake and dbt,” we anchor the model in the correct dialect and conventions.
- **Temperature is 0:** For code generation, deterministic output is almost always preferable.
- **Schema is provided explicitly:** Rather than relying on the model to guess column names, we supply the exact schema.
- **Requirements are numbered and specific:** Ambiguous requirements produce ambiguous code.

Extending the Example: Multi-Turn Conversation

Real engineering work is iterative. After receiving the initial SQL, you might want to refine it:

```
1 # Continue the conversation with follow-up
2 messages = [
3     {"role": "user", "content": "... (original prompt) ..."},
4     {"role": "assistant", "content": message.content[0].text},
5     {"role": "user", "content": """"Good start. Now:
6 1. Add a data quality check: assert DAU > 0 for every date.
7 2. Add a deduplication step in case raw_events has duplicates.
8 3. Wrap in a dbt model with the appropriate config block.
9 4. Generate the schema.yml with column descriptions and tests."""}
10 ]
11
12 refined = client.messages.create(
13     model="claude-sonnet-4-20250514",
14     max_tokens=8192,
15     temperature=0,
16     system=""""You are a senior data engineer specializing in
17 Snowflake and dbt. Use CTEs, explicit column aliases, and
18 Snowflake SQL dialect.""",
19     messages=messages
20 )
21
22 print(refined.content[0].text)
```

Listing 1.2: Multi-turn refinement of generated SQL

! Tip

The 80/20 Rule CLAUDE CODE typically gets you 80% of the way in 20% of the time. The remaining 20%—edge cases, performance tuning, integration specifics—still requires your expertise. This is a feature, not a bug: the model handles the boilerplate so you can focus on the judgment calls. Think of yourself as a quality assurance engineer for an AI that does your job. Inspiring, isn't it?

A Note on API Keys and Security

The `anthropic.Anthropic()` constructor reads the `ANTHROPIC_API_KEY` environment variable by default. Never hardcode API keys in source files. Use your organization's secrets manager (AWS Secrets Manager, HashiCorp Vault, or similar) and inject the key at runtime.

```

1 import os
2 import anthropic
3
4 # From environment variable (set by secrets manager)
5 client = anthropic.Anthropic(
6     api_key=os.environ.get("ANTHROPIC_API_KEY")
7 )
8
9 # Or from AWS Secrets Manager
10 import boto3, json
11
12 def get_api_key():
13     sm = boto3.client("secretsmanager")
14     secret = sm.get_secret_value(SecretId="anthropic/api-key")
15     return json.loads(secret["SecretString"])["api_key"]
16
17 client = anthropic.Anthropic(api_key=get_api_key())

```

Listing 1.3: Secure API key management

1.6 Real-World Case Study: Migrating a Legacy ETL

To make the value of CLAUDE CODE concrete, consider this scenario. A mid-size e-commerce company has a legacy ETL pipeline written in a mix of bash scripts and stored procedures. The pipeline ingests order data from a PostgreSQL OLTP database, transforms it through a series of staging tables, and loads summary tables consumed by a BI tool. The pipeline has no tests, minimal documentation, and the original author left the company two years ago.

The new data engineer, assigned to migrate this pipeline to dbt on Snowflake, faces weeks of reverse-engineering. You could also use CLAUDE CODE to write the resignation letter, but let us explore the more optimistic path first:

1. **Day 1:** Feed the stored procedures to CLAUDE CODE and ask it to extract the business rules in plain English. Result: a structured document mapping each procedure to its business purpose, input tables, output tables, and transformation logic.
2. **Day 2:** Provide the extracted business rules and ask CLAUDE CODE to generate dbt models with appropriate materializations, tests, and schema YAML files. Result: a working dbt project skeleton covering 80% of the transformation logic.

3. **Days 3–4:** Review the generated models, fix edge cases CLAUDE CODE missed (a timezone conversion that depended on application-level config, a fiscal calendar lookup table), and run the dbt project against staging data to validate output.
4. **Day 5:** Use CLAUDE CODE to generate a comprehensive test suite comparing legacy pipeline output to dbt output, table by table. Result: automated regression tests confirming parity.

What would have taken three to four weeks of solo work was completed in one week. The engineer's time shifted from mechanical translation to judgment-intensive review and validation. The original author, presumably still churning butter somewhere off the grid, was not available for comment.

1.7 Chapter Summary

- CLAUDE CODE's large context windows, multi-modal input, extended thinking, tool use, and tiered models make it uniquely suited for data engineering work across the full technology stack.
- It addresses the breadth, documentation, testing, review, and onboarding challenges of modern data teams.
- Limitations—hallucination, stale knowledge, finite context, non-determinism, and the competence illusion—require disciplined validation and healthy skepticism.
- The paradigm shifts from recall-based to judgment-based expertise: you evaluate and direct rather than recall syntax and write boilerplate.
- Effective use starts with understanding the tool deeply, just as you would evaluate any new database, orchestrator, or framework before adopting it.

1.8 Exercises

1. **First API Call:** Set up the Anthropic Python SDK and make your first API call. Generate a SQL query for a schema you work with daily. Compare the output to what you would write manually—what did it get right, and what did it miss?

2. **Model Comparison:** Take a moderately complex SQL query (involving at least two JOINS and a window function). Submit it to both Haiku and Sonnet. Compare the quality, speed, and cost of each response. When would you choose one over the other?
3. **Hallucination Hunting:** Ask CLAUDE CODE to generate a query using a function that exists in PostgreSQL but not in your warehouse's dialect (e.g., `GENERATE_SERIES` in Snowflake). Document how the model handles the mismatch. Does it hallucinate, warn you, or substitute an equivalent?
4. **Prompt Engineering Basics:** Take a simple requirement (e.g., "calculate monthly revenue by product category") and write three versions of the prompt: minimal, moderate detail, and highly specific. Compare the outputs. What is the minimum level of specificity needed for production-quality results?
5. **Documentation Generation:** Feed an undocumented SQL file or dbt model from your codebase to CLAUDE CODE and ask it to generate inline comments and a README section. Evaluate the accuracy of the generated documentation against your knowledge of the code.

In the next chapter, we map CLAUDE CODE's capabilities onto the data engineering lifecycle, showing where AI augmentation fits into each stage of your work. Spoiler: the answer is "everywhere."

2

The Data Engineering Lifecycle Meets AI

Some engineers still debug by reading logs manually. Those engineers also churn their own butter.

! Tip

Who This Chapter Is For This chapter is for data engineers who build and maintain data pipelines and want to understand exactly where `CLAUDE CODE` fits into every stage of the data engineering lifecycle. Whether you are a principal engineer or six months into the role, this chapter provides a concrete map of AI augmentation across your workflow.

This chapter does three things. First, it explains the data engineering lifecycle in enough detail to audit your own platform. Second, it maps `CLAUDE CODE` onto every stage and undercurrent, showing where AI assistance is most impactful. Third, it walks through a realistic case study—a full day in the life of a `CLAUDE CODE`-augmented data engineer. Could a human do all of this without AI? Technically yes. Should they? The question answers itself.

2.1 The Data Engineering Lifecycle

The lifecycle consists of five stages:

1. **Generation** — Data is created by source systems: application databases, IoT sensors, SaaS APIs, clickstreams, and third-party data providers. The data engineer does not control generation but must deeply understand it.
2. **Ingestion** — Data moves from sources into the data platform via batch, microbatch, or streaming. Ingestion encompasses extraction, transport, and initial landing.
3. **Storage** — Data lands in a lake, warehouse, or lakehouse. Partitioning, format, and compression choices affect everything downstream—query performance, cost, and even what transformations are feasible.
4. **Transformation** — Raw data is cleaned, joined, aggregated, and shaped into models serving business needs. This is where business logic is encoded and where data quality is enforced.
5. **Serving** — Transformed data reaches consumers: dashboards, reverse ETL, ML training, APIs, and ad-hoc queries. Serving is the stage where the value of all upstream work is realized.

Five *undercurrents* run beneath all stages: Security, Data Management, DataOps, Data Architecture, and Orchestration. These are not sequential steps but continuous concerns that influence decisions at every stage.

Note

The Lifecycle Is Not a Waterfall Data flows are iterative and cyclical. Serving analytics may reveal quality issues requiring upstream changes. A schema change at generation forces adjustments through ingestion, storage, and transformation. The lifecycle is a mental model for organizing thinking, not a rigid process to follow linearly. Much like the stages of grief you experience when your pipeline breaks at 3 AM.

Why the Lifecycle Matters for AI Integration

The lifecycle framework matters for AI integration because it provides a systematic way to identify where CLAUDE CODE delivers value and where human judgment remains essential. Without this framework, teams tend to adopt AI haphazardly—using it for whatever task happens to be in front of them—rather than strategically targeting the highest-leverage opportunities.

Each stage has different characteristics that affect how well AI assistance works:

- **Structured vs. unstructured tasks:** Transformation (SQL, dbt) is highly structured and well-suited to AI generation. Architecture decisions are unstructured and benefit from AI as a thought partner rather than a code generator.
- **Risk profile:** Ingestion failures are recoverable (re-run the pipeline). Serving incorrect data to financial reports is not. AI involvement should be inversely proportional to risk.
- **Feedback loop speed:** Transformation has fast feedback loops (run the query, check the results). Architecture decisions have slow feedback loops (consequences emerge over months). Fast feedback loops favor more autonomous AI use.

2.2 Claude Code at Each Lifecycle Stage

Generation — Understanding Source Systems

CLAUDE CODE accelerates source system analysis. Given a schema dump or API documentation, it identifies primary keys, foreign key relationships, PII columns, encoding inconsistencies, and proposes CDC strategies. This is covered in depth in Chapter 4.

Code Example

Analyzing an Unfamiliar Source Schema

```
1 I have a PostgreSQL source with these tables:
2 - users (id, email, name, created_at, updated_at, status, tier_id)
3 - orders (id, user_id, total_cents, currency, placed_at, status)
4 - order_items (id, order_id, product_id, quantity, unit_price_cents)
5 - products (id, sku, name, category_id, price_cents, is_active)
6
7 Help me understand the data model, identify join paths,
8 and flag potential ingestion issues.
```

CLAUDE CODE will typically identify that `total_cents` uses integer representation to avoid floating-point issues, that `status` columns on both `users` and `orders` likely represent different state machines, that `tier_id` implies a missing `tiers` table,

and that `placed_at` vs. `created_at` naming inconsistency may indicate different timestamp semantics.

! Tip

Source System Prompts When analyzing source systems, always include the database engine (PostgreSQL, MySQL, SQL Server), the application domain (e-commerce, healthcare, fintech), and any known quirks (“the `status` column uses integer codes, not strings”). Context dramatically improves analysis quality.

Ingestion — Building and Maintaining Pipelines

CLAUDE CODE generates connector scaffolds from API documentation (with pagination, rate limiting, error handling), produces Airbyte/Meltano YAML configurations, diagnoses ingestion failures from error tracebacks, and handles schema evolution.

Code Example

Generating an Airbyte Source Connector Configuration

```
1 # Generated by Claude Code from API documentation
2 version: "0.1.0"
3 type: DeclarativeSource
4 check:
5   type: CheckStream
6   stream_names: ["customers"]
7 streams:
8   - type: DeclarativeStream
9     name: customers
10    primary_key: [["id"]]
11    retriever:
12      type: SimpleRetriever
13    requester:
14      type: HttpRequester
15      url_base: "https://api.example.com/v2"
16      path: "/customers"
17      http_method: "GET"
18    authenticator:
19      type: BearerAuthenticator
20      api_token: "{{ config['api_key'] }}"
21    paginator:
```

```
22     type: DefaultPaginator
23     page_token_option:
24       type: RequestOptions
25       inject_into: request_parameter
26       field_name: "cursor"
27     pagination_strategy:
28       type: CursorPagination
29       cursor_value: "{{ response.next_cursor }}"
30       stop_condition: "{{ response.has_more is false }}"
```

! Tip

Let CLAUDE CODE Handle the Boilerplate Retry logic, rate-limit handling, state tracking, and pagination are boilerplate every connector needs. Before, you had to write this yourself. Yuck. Describe your source API to CLAUDE CODE and customize only the business-specific parts. This applies equally to custom Python connectors, Airbyte configurations, and Meltano extractors.

Common ingestion tasks where CLAUDE CODE saves significant time:

- **Error diagnosis:** Paste a traceback from a failed ingestion run. CLAUDE CODE identifies the root cause (expired token, schema change, rate limit) and suggests a fix.
- **Schema evolution handling:** When a source adds new columns, CLAUDE CODE generates the migration code for your staging layer.
- **Backfill planning:** Describe your source API's pagination and rate limits, and CLAUDE CODE calculates optimal batch sizes and parallelism for historical backfills.
- **Data format conversion:** Generate code to convert between CSV, JSON, Avro, and Parquet with appropriate schema mappings.

Storage — Designing the Data Platform

CLAUDE CODE generates DDL with appropriate partitioning, clustering, and compression. It recommends storage formats (Parquet vs. Iceberg vs. Delta) based on your query patterns, and generates safe, reversible migration scripts.

Code Example

Generating Optimized BigQuery DDL

```
1 CREATE TABLE IF NOT EXISTS `analytics.events`  
2 (  
3     event_id          STRING NOT NULL,  
4     event_type        STRING NOT NULL,  
5     event_timestamp   TIMESTAMP NOT NULL,  
6     event_date         DATE NOT NULL,  
7     user_id           STRING,  
8     properties         JSON,  
9     processing_time   TIMESTAMP NOT NULL  
10        DEFAULT CURRENT_TIMESTAMP()  
11 )  
12 PARTITION BY event_date  
13 CLUSTER BY event_type, user_id  
14 OPTIONS (  
15     partition_expiration_days = 90,  
16     require_partition_filter = true  
17 );
```

Notice how the generated DDL includes production best practices: partition expiration to control storage costs, required partition filters to prevent expensive full-table scans, clustering on the most common filter columns, and a `processing_time` column for pipeline debugging.

Note

Storage Format Selection When asking CLAUDE CODE for storage format recommendations, provide your query patterns (point lookups vs. analytical scans), update frequency (append-only vs. upserts), read engine (Spark, Trino, warehouse-native), and data volume. These factors matter more than generic “best practices.”

Transformation — The Power Zone

Transformation is where CLAUDE CODE delivers the most immediate value. This stage involves the most code, the most business logic, and the most iteration—all characteristics that favor AI augmentation.

Key transformation tasks where CLAUDE CODE excels:

- **SQL generation from natural language:** Describe the business requirement; get a query with CTEs, proper null handling, and dialect-specific syntax.
- **dbt model scaffolding:** Generate models, schema YAML, source definitions, and tests from a description of the desired output.
- **Query optimization:** Provide a slow query and its EXPLAIN plan; get specific optimization suggestions with before/after SQL.
- **Data quality test generation:** Describe business rules in English; get dbt tests, Great Expectations suites, or custom SQL assertions.
- **Complex window functions:** Window functions are the number one area where CLAUDE CODE saves time—describing the desired result in English is far faster than debugging PARTITION BY and ORDER BY clauses.

Code Example

Before and After Query Optimization

```
1  -- BEFORE: Correlated subqueries, 47 minutes on 2B rows
2  SELECT u.user_id, u.email,
3         (SELECT COUNT(*) FROM orders o
4          WHERE o.user_id = u.user_id AND o.placed_at >= '2025-01-01')
5  FROM users u WHERE u.status = 'active';
6
7  -- AFTER: Single scan with hash join, 2 minutes 14 seconds
8  SELECT u.user_id, u.email,
9         COALESCE(agg.order_count, 0) AS order_count
10 FROM users u
11 LEFT JOIN (
12     SELECT user_id, COUNT(*) AS order_count
13     FROM orders WHERE placed_at >= '2025-01-01'
14     GROUP BY user_id
15 ) agg ON u.user_id = agg.user_id
16 WHERE u.status = 'active';
```

Warning

Optimization Context Matters CLAUDE CODE's optimization suggestions are based on general SQL performance principles. Always validate with your warehouse's actual EXPLAIN plan. A suggestion that helps on Snowflake may hurt

on BigQuery due to different query optimizer strategies. Blindly applying AI-generated optimizations is how you turn a 47-minute query into a 47-hour query. Ask us how we know.

Serving — Delivering Data to Consumers

CLAUDE CODE scaffolds FastAPI endpoints for warehouse data, generates dashboard SQL specifications, produces reverse ETL configurations, and drafts data contracts between producers and consumers.

Code Example

Generating a Data Contract

```
1 # Data contract: orders_summary
2 # Producer: data-engineering-team
3 # Consumer: analytics-team, ml-team
4 contract:
5   name: orders_summary
6   version: "2.1.0"
7   owner: data-engineering-team
8   description: >
9     Daily order aggregations by customer segment.
10    Updated by 06:00 UTC. SLA: 99.5% availability.
11   schema:
12     fields:
13       - name: order_date
14         type: DATE
15         description: "Calendar date of order placement"
16         not_null: true
17       - name: customer_segment
18         type: STRING
19         description: "Customer tier: premium, standard, trial"
20         allowed_values: ["premium", "standard", "trial"]
21       - name: total_orders
22         type: INTEGER
23         description: "Count of distinct orders"
24         minimum: 0
25       - name: total_revenue_cents
26         type: BIGINT
27         description: "Sum of order totals in cents (USD)"
28         minimum: 0
```

```
29 quality:
30   freshness: "< 6 hours"
31   completeness: "> 99.5%"
32   uniqueness_key: ["order_date", "customer_segment"]
33 breaking_changes:
34   notification: 30_days
35   channel: "#data-contracts"
```

! Tip

Data Contracts as a Starting Point Use CLAUDE CODE to generate initial data contracts from your existing schema and documentation, then refine with stakeholder input. This reverses the usual bottleneck: instead of starting from a blank page, you start with a concrete draft that stakeholders can react to.

2.3 Claude Code in the Undercurrents

The five undercurrents are continuous concerns that span all lifecycle stages. CLAUDE CODE augments each one differently.

Security: Reviews IAM policies for over-permissioning, generates column-level access controls, writes data masking rules based on data classification, audits SQL for injection vulnerabilities, and generates encryption configurations. See Chapter 18 for comprehensive coverage.

Data Management: Enriches data catalogs with AI-generated descriptions, generates data quality rules from plain English business requirements, traces lineage through SQL transformations, drafts governance policies, and classifies data by sensitivity level. See Chapter 17.

DataOps: Writes CI/CD pipelines for dbt projects (GitHub Actions, GitLab CI), generates unit tests for transformation logic, creates monitoring rules and alert thresholds, and drafts incident runbooks from past incident reports.

Data Architecture: Compares architectural patterns (lambda vs. kappa, lakehouse vs. warehouse), evaluates tool trade-offs with structured pros/cons analysis, generates architecture decision records (ADRs), reviews designs for single points of failure and scalability bottlenecks.

Orchestration: Writes Airflow DAGs from descriptions of dependencies and schedules, debugs dependency issues in existing DAGs, generates Dagster asset definitions,

converts between orchestration frameworks, and creates SLA monitoring configurations.

Code Example

Generating an Airflow DAG from a Description

```
1 # Generated by Claude Code from: "Daily ETL that extracts from
2 # PostgreSQL, stages in S3, transforms with dbt, and notifies
3 # Slack on failure. Runs at 06:00 UTC, retries 3 times."
4
5 from airflow import DAG
6 from airflow.operators.python import PythonOperator
7 from airflow.providers.slack.operators.slack_webhook import (
8     SlackWebhookOperator,
9 )
10 from datetime import datetime, timedelta
11
12 default_args = {
13     "owner": "data-engineering",
14     "depends_on_past": False,
15     "retries": 3,
16     "retry_delay": timedelta(minutes=5),
17     "on_failure_callback": notify_slack,
18 }
19
20 with DAG(
21     dag_id="daily_customer_etl",
22     default_args=default_args,
23     schedule_interval="0 6 * * *",
24     start_date=datetime(2025, 1, 1),
25     catchup=False,
26     tags=["etl", "daily", "customers"],
27 ) as dag:
28     extract = PythonOperator(
29         task_id="extract_from_postgres",
30         python_callable=extract_postgres_to_s3,
31     )
32     transform = PythonOperator(
33         task_id="run_dbt_models",
34         python_callable=run_dbt,
35         op_kwargs={"select": "tag:daily_customers"},
36     )
37     test = PythonOperator(
```

```
38     task_id="run_dbt_tests",
39     python_callable=run_dbt_test,
40     op_kwargs={"select": "tag:daily_customers"},
41 )
42
43 extract >> transform >> test
```

Note

Undercurrent Integration The undercurrents are where CLAUDE CODE often delivers the most *underappreciated* value. Engineers tend to focus on transformation (writing SQL) because it is the most visible work. But security reviews, documentation, testing, and architecture decisions—the undercurrents—are where quality and reliability are determined.

2.4 Case Study: A Day with Claude Code

Meet CLAUDE CODE, a senior data engineer at a Series B fintech company. Her team of three manages 120 dbt models, 40 Airflow DAGs, and data from 15 source systems. Here is her CLAUDE CODE-augmented day (which, let's be honest, is really CLAUDE CODE's day with her in a supervisory capacity):

8:30 AM — Morning triage: Three Airflow failures overnight. Like clockwork. She pastes the Salesforce ingestion traceback into CLAUDE CODE, which identifies a transient TLS renegotiation failure specific to the Salesforce bulk API and recommends adding the `-tls-retry` flag to the connector configuration. The second failure is a Stripe webhook schema change—CLAUDE CODE identifies the new fields and generates the staging model update. The third is a timeout on a long-running dbt model—CLAUDE CODE suggests adding a pre-filter CTE. Elapsed: 25 minutes total vs. 1.5 hours traditional.

9:15 AM — Schema change response: The product team added three new columns to the `users` table: `referral_source`, `signup_campaign`, and `preferred_language`. She pastes the staging model and new column definitions; CLAUDE CODE updates the model SQL, adds the columns to the schema YAML with descriptions and appropriate tests (not-null for `preferred_language`, accepted-values for `referral_source`), and generates a migration script. Elapsed: 7 minutes vs. 25 minutes.

10:00 AM — New pipeline build: The growth team needs data from a new marketing attribution API. Priya feeds the 40-page API documentation to CLAUDE

CODE, which produces a complete ingestion implementation with pagination, rate limiting, and error handling; an Airflow DAG; a staging model with incremental materialization; and a schema YAML with tests. She ships by lunch after 90 minutes of review and customization. Elapsed: 2 hours vs. 5–6 hours.

1:00 PM — Code review: Two PRs from teammates. She asks CLAUDE CODE to pre-review both, providing the team’s code standards document as context. CLAUDE CODE flags a missing partition filter in one query, a potential fan-out join in another, and suggests using QUALIFY instead of a subquery for deduplication. Priya confirms the findings and adds her own architectural feedback. Elapsed: 20 minutes vs. 45 minutes.

3:00 PM — Data quality investigation: Revenue dropped 15% in yesterday’s executive dashboard. CLAUDE CODE generates five diagnostic queries: checking row counts by date, comparing revenue by payment processor, looking for null values in key fields, checking for duplicate transactions, and comparing against the same day last week. Root cause: a new payment processor (added last week) uses `status = 'completed'` instead of `status = 'success'`, filtered out by the staging model’s hardcoded filter. CLAUDE CODE generates the fix and a parameterized test to prevent recurrence. Elapsed: 40 minutes vs. 1.5 hours.

4:00 PM — Documentation sprint: CLAUDE CODE feeds 12 undocumented dbt models to CLAUDE CODE with their SQL and the source schema. CLAUDE CODE generates descriptions for each model, column-level documentation, and a lineage summary. She reviews and publishes to the data catalog. Elapsed: 45 minutes vs. 3 hours.

Table 2.1: CLAUDE CODE’s day: estimated time with and without CLAUDE CODE.

Task	Without Claude Code	With Claude Code
Morning triage (3 failures)	1.5 hours	25 minutes
Schema change response	25 minutes	7 minutes
New pipeline build	5 hours	2 hours
Code review (2 PRs)	45 minutes	20 minutes
Data quality investigation	1.5 hours	40 minutes
Documentation sprint (12 models)	3 hours	45 minutes
Total	12.25 hours	4.3 hours

 **Warning**

These Numbers Are Illustrative Time savings vary by task complexity, engineer experience, and codebase familiarity. Measure your own team's actual savings rather than relying on generic estimates. Track time before and after AI adoption for at least one sprint to establish a realistic baseline. And yes, you can use CLAUDE CODE to build the time-tracking dashboard. We see you.

2.5 From Tool-Centric to Intent-Centric Engineering

In the traditional model, you think in tool abstractions: Airflow operators, dbt macros, Spark DataFrames, Terraform resources. You must know the specific syntax, configuration options, and quirks of each tool before you can express your intent.

CLAUDE CODE enables *intent-centric* engineering: express what you want in natural language, review the generated implementation, iterate with feedback, validate against known-good results, and commit.

The workflow shift looks like this:

1. **Traditional:** Think about the problem → recall the right tool/function/syntax → write the code → debug syntax and logic errors → test → iterate.
2. **Intent-centric:** Think about the problem → describe the desired outcome to CLAUDE CODE → review the generated code → refine with feedback → test → iterate.

The key difference is that step 2 changes from *recall and write* to *describe and review*. Your job title hasn't changed, but your job description now reads "professional code reviewer who occasionally types." For experienced engineers, reviewing code is faster than writing it. For complex tasks spanning multiple tools, the speedup is even more pronounced because you no longer need to context-switch between different syntaxes and APIs.

Intent-centric does *not* mean engineers stop understanding tools. Effective CLAUDE CODE use requires *deeper* understanding—to review output for correctness, specify intent precisely enough to get good results, evaluate architectural suggestions against real-world constraints, and know when CLAUDE CODE's suggestion is good-enough versus when it needs rework. The paradigm raises the *leverage* of engineering skill without lowering the bar for what constitutes competent engineering.

 **Warning**

Time Savings Require Validation Never deploy AI-generated code without review. The savings come from shifting effort from *writing* to *reviewing*, which is faster for experienced engineers but requires knowing what correct output looks like. Junior engineers should use CLAUDE CODE as a learning tool with senior review, not as a shortcut past understanding.

 **Note**

The Expertise Paradox The more you know, the more value you get from CLAUDE CODE. Senior engineers who deeply understand SQL optimization, data modeling, and distributed systems can leverage AI to work at 3–5x speed because they can quickly evaluate output quality. Junior engineers get less leverage because they cannot yet distinguish good output from subtly wrong output. This is the expertise paradox: AI amplifies existing skill rather than replacing it. In other words, the robots are not coming for your job—they are coming for your junior colleagues' jobs. Sleep well.

2.6 Summary

CLAUDE CODE augmentation permeates the entire lifecycle: source analysis at Generation, connector code at Ingestion, DDL and format advice at Storage, SQL and dbt models at Transformation, API scaffolding and data contracts at Serving, and judgment amplification across all undercurrents—Security, Data Management, DataOps, Data Architecture, and Orchestration.

The shift from tool-centric to intent-centric engineering raises leverage without lowering the skill bar. Engineers who embrace this shift—investing in deeper understanding of their systems while delegating boilerplate to AI—will dramatically outperform those who either reject AI entirely or adopt it without the critical thinking needed to use it effectively.

2.7 Exercises

1. **Lifecycle Mapping:** Choose a pipeline you own. Map each component to a lifecycle stage. For each stage, identify one task where CLAUDE CODE saves time and one where human judgment is irreplaceable. Present your mapping to your team and compare perspectives.

2. **Query Optimization Challenge:** Find a slow query in your warehouse's query history (most warehouses expose this via `INFORMATION_SCHEMA` or a query history view). Feed the query and its `EXPLAIN` output to `CLAUDE CODE`. Evaluate the suggestions, run both versions on staging, and compare execution time. Document the results.
3. **Workflow Comparison:** Take a recently completed task (a new model, a bug fix, a connector update). Estimate how long it took traditionally. Re-do it using `CLAUDE CODE` as a pair programmer. Measure the difference and note specifically where you had to intervene to correct `CLAUDE CODE`'s output.
4. **Tool-Centric vs. Intent-Centric:** Write a short Airflow DAG the traditional way (looking up operator documentation, writing Python). Then describe the same DAG to `CLAUDE CODE` in natural language. Compare: which was faster? Which produced better code? Which was easier to modify?
5. **Undercurrent Audit:** Choose one undercurrent (Security, Data Management, DataOps, Data Architecture, or Orchestration) and audit your current platform. Use `CLAUDE CODE` to identify gaps—missing tests, undocumented models, over-permissioned roles, or missing monitoring. Prioritize the top three gaps and create tickets to address them.

3

Getting Started with Claude Code

The journey of a thousand pipelines begins with a single API key.
And a credit card.

This chapter is where theory meets practice. By the end you will have a fully functional *schema documentation generator* you can drop into any data team's workflow.

3.1 Setting Up the Anthropic API

Navigate to <https://console.anthropic.com> and create an account. The dashboard exposes three panels: **API Keys**, **Usage**, and **Rate Limits**.

Warning

Seriously, Do Not Commit Your API Key. Treat API keys like database passwords. Never commit them to version control. Use a secrets manager or a `.env` file excluded via `.gitignore`. If you commit an API key to a public repo, CLAUDE CODE will not judge you, but the internet will. Bots scrape GitHub for exposed keys faster than you can type `git revert`.

Code Example

```
1 export ANTHROPIC_API_KEY="sk-ant-api03-..."
```

Tip

Budget Self-Defense Set a monthly spending limit in the Console under *Billing* → *Limits*. A forgotten `while True` loop in a notebook can burn through your budget faster than an intern with a corporate credit card at a steakhouse.

3.2 The Anthropic Python SDK

Code Example

```
1 pip install anthropic
```

The simplest call (the one that will make you wonder what you've been doing with your career up to this point):

Code Example

```
1 import anthropic
2
3 client = anthropic.Anthropic() # reads ANTHROPIC_API_KEY from env
4
5 message = client.messages.create(
6     model="claude-sonnet-4-20250514",
7     max_tokens=1024,
8     messages=[
9         {"role": "user", "content": "Explain star schema vs snowflake
10         schema."}
11     ],
12     print(message.content[0].text)
```

For concurrent pipeline work, use the async client:

Code Example

```
1 import asyncio
2 import anthropic
3
4 async def describe_table(table_name: str, ddl: str) -> str:
5     client = anthropic.AsyncAnthropic()
6     message = await client.messages.create(
7         model="claude-sonnet-4-20250514",
8         max_tokens=512,
9         messages=[{
10            "role": "user",
11            "content": f"Describe the purpose of this table:\n```\nsql\n{ddl}\n```",
12        }],
13    )
14    return message.content[0].text
15
16 async def main():
17     ddls = {
18         "dim_customer": "CREATE TABLE dim_customer (id INT, name TEXT
19         , ...);",
20         "fact_orders": "CREATE TABLE fact_orders (id INT, amount
21         DECIMAL, ...);",
22     }
23     tasks = [describe_table(name, ddl) for name, ddl in ddls.items()]
24     results = await asyncio.gather(*tasks)
25     for name, desc in zip(ddls.keys(), results):
26         print(f"\n=== {name} ===\n{desc}")
27
28 asyncio.run(main())
```

3.3 Tokens, Context Windows, and Pricing

A *token* is roughly three-quarters of an English word. SQL and code tokenize less efficiently—assume one token per two characters for DDL-heavy prompts.

The 5:1 input-to-output cost ratio means prompts with large DDL dumps (high input) but concise summaries (low output) are surprisingly economical.

Table 3.1: Context windows and pricing (early 2026, approximate).

Model	Context	Input (\$/1M)	Output (\$/1M)
Haiku	200K	\$0.25	\$1.25
Sonnet	200K	\$3.00	\$15.00
Opus	200K	\$15.00	\$75.00

3.4 System Prompts for Data Engineering

The *system prompt* is your most powerful lever. A well-crafted system prompt is the difference between generic advice and expert-level output.

Code Example

```
1 You are a senior data engineer with 12 years of experience.
2 Expertise: Kimball modeling, dbt, Spark, Snowflake, BigQuery,
3 Great Expectations, Airflow.
4
5 When generating SQL:
6 - Use CTEs over nested subqueries
7 - Add comments explaining business logic
8 - Include QUALIFY for window filtering (Snowflake)
9 - Never SELECT *
10 - Consider performance at billions of rows
```

Tip

Store system prompts in version-controlled files alongside your pipeline code. Review changes in pull requests.

Table 3.2: Recommended temperature settings.

Task	Rationale	Temp
SQL generation	Deterministic correctness	0.0
Documentation writing	Slight phrasing variety	0.3
Test data generation	Diverse, realistic data	0.8

3.5 Temperature and Parameters

⚠ Warning

Determinism Is a Lie Setting temperature to 0.0 does not guarantee identical outputs across requests. If you need byte-for-byte reproducibility, cache responses. Expecting deterministic output from a language model is like expecting deterministic behavior from a cat. You can ask nicely, but no promises.

3.6 Error Handling and Retry Patterns

Code Example

```
1 import time, logging, anthropic
2
3 logger = logging.getLogger(__name__)
4
5 def call_claude_with_retry(
6     client: anthropic.Anthropic,
7     messages: list[dict],
8     model: str = "claude-sonnet-4-20250514",
9     max_tokens: int = 4096,
10    max_retries: int = 5,
11    system: str = "",
12 ) -> anthropic.types.Message:
13     """Call the Messages API with exponential backoff."""
14     for attempt in range(max_retries):
15         try:
16             return client.messages.create(
17                 model=model, max_tokens=max_tokens,
18                 system=system, messages=messages,
```

```
19         )
20     except anthropic.RateLimitError as e:
21         delay = min(2 ** attempt + (time.time() % 1), 60)
22         retry_after = e.response.headers.get("retry-after")
23         if retry_after:
24             delay = max(delay, float(retry_after))
25             logger.warning(f"Rate limited, retrying in {delay:.1f}s")
26             time.sleep(delay)
27     except anthropic.InternalServerError:
28         time.sleep(2 ** attempt)
29     except (anthropic.BadRequestError, anthropic.
AuthenticationError):
30         raise
31     raise RuntimeError(f"Failed after {max_retries} attempts")
```

3.7 Choosing the Right Model

1. **High-volume, simple?** (documenting 1,000 columns, classifying types) → **Haiku**.
2. **Solid reasoning?** (complex SQL, schema review, tests) → **Sonnet**.
3. **Novel, ambiguous, high-stakes?** (architecture design, major migration) → **Opus**.

! Tip

The “Develop Rich, Deploy Cheap” Strategy Prototype with Sonnet, then switch to Haiku for production if quality is sufficient. This “develop on Sonnet, deploy on Haiku” strategy can cut costs by 90%. It is the data engineering equivalent of testing recipes with truffle oil and serving with canola. Your stakeholders will never know the difference.

3.8 Tool Use / Function Calling

Tool use lets CLAUDE CODE invoke external functions during a conversation—query your database, check row counts, or validate generated queries.

Code Example

```
1  tools = [{
2      "name": "run_sql_query",
3      "description": "Execute a read-only SQL query. Returns up to 100
4      rows.",
5      "input_schema": {
6          "type": "object",
7          "properties": {
8              "query": {"type": "string", "description": "SQL SELECT
9              query"},
10             },
11             "required": ["query"],
12         },
13     },
14 ]
15
16 def chat_with_tools(user_message: str) -> str:
17     messages = [{"role": "user", "content": user_message}]
18     while True:
19         response = client.messages.create(
20             model="claude-sonnet-4-20250514",
21             max_tokens=4096, tools=tools, messages=messages,
22         )
23         if response.stop_reason == "tool_use":
24             messages.append({"role": "assistant", "content": response
25             .content})
26             tool_results = []
27             for block in response.content:
28                 if block.type == "tool_use":
29                     result = execute_readonly_query(block.input["
30                     query"])
31                     tool_results.append({
32                         "type": "tool_result",
33                         "tool_use_id": block.id,
34                         "content": json.dumps(result, default=str),
35                     })
36             messages.append({"role": "user", "content": tool_results
37             })
38         else:
39             return "".join(
40                 b.text for b in response.content if b.type == "text"
41             )
```

Warning

Read-Only or Regret When giving CLAUDE CODE SQL execution access, **always** enforce read-only. Use a database role with only SELECT privileges. Giving an AI DROP TABLE privileges is the data engineering equivalent of handing a toddler a flamethrower. Technically possible. Universally inadvisable.

3.9 Cost Estimation

Table 3.3: Sample monthly cost for a 5-person data team.

Task	Model	Calls/mo	Cost
Column docs	Haiku	10,000	\$3.75
SQL generation	Sonnet	2,000	\$27.00
Schema review	Sonnet	500	\$22.50
Ad-hoc exploration	Sonnet	5,000	\$75.00
Architecture design	Opus	20	\$10.50
Total			\$138.75

Tip

The Math Is Embarrassingly Simple At roughly \$140/month for a five-person team, API costs are a rounding error compared to engineer salaries. Even one hour saved per engineer per week pays for the budget many times over. Your team's coffee budget is larger. Possibly your team's *oat milk surcharge* is larger.

3.10 Exercises

1. **API Exploration.** Create an account, generate an API key, send a message, and calculate the cost of that single request.
2. **System Prompt Engineering.** Write system prompts for three personas (dbt reviewer, Airflow debugger, catalog curator). Test each with the same question and compare actionability.

3. **Model Comparison.** Send a complex schema review prompt to Haiku, Sonnet, and Opus. Compare accuracy, depth, and cost. Calculate cost-per-quality-point.
4. **Tool Use.** Extend a data engineering assistant with a `get_table_stats` tool (row count, null percentages, min/max) and test it with data quality prompts.

Part 

Data Generation and Ingestion with Claude Code

4

Source System Understanding with Claude Code

Understanding a legacy database is like archaeology, except the artifacts are actively on fire and the previous archaeologist quit without leaving notes.

In data engineering, the most underestimated phase is *source system understanding*. Before a single pipeline is built, an engineer must deeply comprehend the systems from which data originates—their schemas, APIs, business logic, update patterns, and quirks. CLAUDE CODE transforms what was once a weeks-long archaeological expedition into a structured, rapid analysis process.

Tip

Time Investment Comparison A senior engineer typically spends 2–4 weeks onboarding to a complex source system. With CLAUDE CODE-assisted analysis, this drops to 2–4 *days*. The key is eliminating mechanical labor—reading, cross-referencing, and documenting—so engineers can focus on judgment: which tables matter, which relationships are trustworthy, and which data quality issues will cause downstream problems. The remaining 2–4 days are mostly spent saying “wait, *that’s* what this column means?”

4.1 Why Source System Understanding Matters

Every data quality issue, every pipeline failure, and every misleading dashboard can be traced back to insufficient understanding of the source system. The source is where data is born, and its characteristics—good and bad—propagate through every downstream system.

Common consequences of poor source understanding include:

- **Incorrect joins:** Assuming a column is a foreign key when it is actually a denormalized copy that drifts out of sync.
- **Missing data:** Not knowing that soft-deleted records remain in the table with a `deleted_at` timestamp, leading to inflated counts.
- **Type mismatches:** Ingesting a column as a string when it contains JSON, or as an integer when it occasionally contains null-like sentinel values (e.g., `-1`, `9999`).
- **Timezone confusion:** Assuming timestamps are in UTC when they are actually in the application server’s local timezone.
- **Business logic buried in application code:** Transformation rules that exist only in the application layer, invisible to anyone looking at raw database tables.

CLAUDE CODE cannot eliminate these risks, but it dramatically accelerates the process of discovering and documenting them.

Warning

Source Understanding Is Never “Done” Source systems evolve continuously. Schema changes, new features, API version upgrades, and business logic updates happen without notification to the data team. The product team considers you on a “need to know” basis, and they have decided you do not need to know. Treat source understanding as an ongoing process, not a one-time project. Automate schema drift detection and re-run analysis periodically—or just let CLAUDE CODE watch the database for you while you sleep.

4.2 Reverse-Engineering Database Schemas

The most common starting point is a DDL dump—the `CREATE TABLE` statements that define a database’s structure. CLAUDE CODE can parse these statements and produce structured documentation, identify relationships, and flag issues that would take hours to discover manually.

```
1 from anthropic import Anthropic
2 import json
3
4 client = Anthropic()
5
6 def analyze_schema_with_claude(ddl: str, context: str = "") -> dict:
7     """Send DDL to Claude Code for comprehensive schema analysis."""
8     system_prompt = """You are a senior data engineer performing
9     source system analysis. Given DDL statements, produce JSON:
10    1. "tables": descriptions, purpose, classification
11    2. "relationships": FK and inferred relationships
12    3. "anomalies": missing FKs, unusual types, orphan tables
13    4. "business_domains": table groupings
14    5. "recommended_ingestion_order": topological sort
15    6. "data_quality_flags": columns needing special handling"""
16
17    message = client.messages.create(
18        model="claude-sonnet-4-20250514",
19        max_tokens=8000,
20        system=system_prompt,
21        messages=[{
22            "role": "user",
23            "content": f"Analyze this schema:\n\n"
24                    f"Context: {context or 'None'}\n\nDDL:\n{ddl}"
25        }],
26    )
27    response_text = message.content[0].text
28    start = response_text.find("{")
29    end = response_text.rfind("}") + 1
30    if start >= 0 and end > start:
31        return json.loads(response_text[start:end])
32    return {"raw_analysis": response_text}
```

Listing 4.1: Schema Analysis Pipeline

The `context` parameter is crucial. Telling CLAUDE CODE that the source is “an

e-commerce platform built on Ruby on Rails” versus “a healthcare EHR system” changes how it interprets ambiguous column names, recognizes framework conventions (Rails uses `created_at/updated_at` by convention), and classifies data sensitivity.

Handling Large Schemas

Production databases routinely have hundreds or even thousands of tables. A 200K-token context window is generous but not unlimited, and cramming every table into a single prompt reduces analysis quality. The solution is chunked analysis with a synthesis pass:

```

1 import re
2
3 def split_ddl_by_table(ddl: str) -> list[tuple[str, str]]:
4     pattern = r'(CREATE\s+TABLE\s+(?:IF\s+NOT\s+EXISTS\s+)?'
5     pattern += r'["\s]?(\w+(?:\.\w+)?)[\s"]?\s*\(.*\s*\);)'
6     return [(name, stmt) for stmt, name in
7             re.findall(pattern, ddl, re.DOTALL | re.IGNORECASE)]
8
9 def analyze_large_schema(ddl: str, context: str = "") -> dict:
10    tables = split_ddl_by_table(ddl)
11    chunk_analyses = []
12    chunk, chunk_size = [], 0
13    for name, stmt in tables:
14        if chunk_size + len(stmt) > 80_000 and chunk:
15            chunk_ddl = "\n\n".join(s for _, s in chunk)
16            chunk_analyses.append(
17                analyze_schema_with_claude(chunk_ddl, context))
18            chunk, chunk_size = [], 0
19            chunk.append((name, stmt))
20            chunk_size += len(stmt)
21    if chunk:
22        chunk_analyses.append(analyze_schema_with_claude(
23            "\n\n".join(s for _, s in chunk), context))
24
25    # Synthesis pass
26    message = client.messages.create(
27        model="claude-sonnet-4-20250514", max_tokens=8000,
28        messages=[{"role": "user", "content":
29            f"Synthesize these {len(chunk_analyses)} chunk analyses "
30            f"into one unified JSON:\n{json.dumps(chunk_analyses)}"},
31    )

```

```
32 text = message.content[0].text
33 return json.loads(text[text.find("{"):text.rfind("}")+1])
```

Listing 4.2: Chunked Analysis for Large Schemas

Warning

Context Window Management A rough heuristic: 4 characters per token. For a 200K-token window, that gives approximately 800K characters of DDL. If your schema exceeds this, chunking is mandatory. Even below this limit, smaller chunks produce better analysis because the model can focus on fewer tables at a time.

What to Look for in Schema Analysis

When reviewing CLAUDE CODE's schema analysis, pay particular attention to:

- **Implicit relationships:** Columns like `user_id` in the `orders` table that lack explicit foreign key constraints but clearly reference the `users` table. These are common in NoSQL-influenced designs and ORMs that manage referential integrity in application code.
- **Polymorphic associations:** Patterns like `commentable_type` and `commentable_id` (common in Rails) that implement one-to-many relationships across multiple parent tables. These require special handling in your data model.
- **Audit columns:** `created_at`, `updated_at`, `deleted_at`, `created_by`—these indicate CDC capabilities and soft-delete patterns.
- **JSON/JSONB columns:** Semi-structured data stored in relational columns. These often contain business-critical attributes that are not visible in the DDL alone.
- **Enum-like columns:** `VARCHAR` columns that actually contain a fixed set of values (e.g., `status`, `type`, `category`). The allowed values may not be documented in the schema.

Note

Framework Conventions If the source system is built on a web framework (Rails, Django, Laravel, Spring), tell CLAUDE CODE which framework. Each frame-

work has conventions for naming, timestamps, soft deletes, and associations that CLAUDE CODE can recognize and interpret correctly.

4.3 Analyzing REST APIs

APIs often lack self-describing schemas. Unlike databases where DDL tells you the structure, APIs require you to make requests and inspect responses to understand the data. CLAUDE CODE can analyze API responses, infer schemas, identify pagination patterns, and generate documentation.

```
1 import requests, json, time
2 from dataclasses import dataclass, field, asdict
3
4 @dataclass
5 class EndpointAnalysis:
6     method: str
7     path: str
8     status_code: int
9     response_body: dict
10    pagination_type: str = ""
11    rate_limit_headers: dict = field(default_factory=dict)
12
13 class APIAnalyzer:
14     def __init__(self, base_url: str, auth_headers: dict = None):
15         self.base_url = base_url.rstrip("/")
16         self.session = requests.Session()
17         if auth_headers:
18             self.session.headers.update(auth_headers)
19         self.analyses = []
20
21     def probe_endpoint(self, method: str, path: str) -> EndpointAnalysis:
22
23         url = f"{self.base_url}{path}"
24         response = self.session.request(method, url, timeout=30)
25         rate_headers = {k: v for k, v in response.headers.items()
26                        if any(t in k.lower() for t in
27                               ["rate", "limit", "remaining", "retry"])}
28         body = response.json() if response.ok else {}
29         pagination = ""
30         if isinstance(body, dict):
31             if "next" in body: pagination = "cursor"
32             elif "total_pages" in body: pagination = "offset"
```

```
32     analysis = EndpointAnalysis(method, path, response.status_code,
33                                 body, pagination, rate_headers)
34     self.analyses.append(analysis)
35     return analysis
36
37     def generate_documentation(self) -> str:
38         data = json.dumps([asdict(a) for a in self.analyses],
39                             indent=2, default=str)
40         message = client.messages.create(
41             model="claude-sonnet-4-20250514", max_tokens=8000,
42             system="Generate comprehensive API documentation in Markdown.
43 ",
44             messages=[{"role": "user",
45                         "content": f"Document these probe results:\n{data}"}],
46             )
47         return message.content[0].text
```

Listing 4.3: API Discovery and Documentation

API Pagination Patterns

Different APIs use different pagination strategies, and choosing the wrong one leads to missing data or duplicate records. CLAUDE CODE can identify the pagination pattern from a sample response:

- **Cursor-based:** The response includes a `next_cursor` or `next_page_token`. Most reliable for data engineering because it handles concurrent inserts gracefully.
- **Offset-based:** Uses `page` and `per_page` parameters. Prone to skipping or duplicating records if data changes during pagination.
- **Keyset-based:** Paginates using a column value (e.g., `created_after`). Efficient and reliable but requires a monotonically increasing key.
- **Link-header:** The `Link` response header contains URLs for next/previous pages (common in GitHub-style APIs).

! Tip

API Rate Limit Documentation Always ask CLAUDE CODE to extract rate limit information from the API response headers. Common headers include X-RateLimit-Limit, X-RateLimit-Remaining, and Retry-After. Knowing these before building your connector prevents trial-and-error tuning later.

4.4 Understanding Legacy Systems

Legacy systems are the most challenging sources to onboard. They embed business logic in stored procedures, triggers, BLOB columns, and application-level encoding. Documentation is sparse or outdated. The original developers are often unavailable.

CLAUDE CODE can extract business rules from stored procedure definitions, which is one of the most time-consuming aspects of legacy system onboarding.

```

1 def analyze_stored_procedure(proc: dict) -> dict:
2     message = client.messages.create(
3         model="claude-sonnet-4-20250514", max_tokens=4000,
4         system="""Analyze this stored procedure. Return JSON:
5             1. Purpose in plain English
6             2. Business rules as formal statements
7             3. Tables read and written
8             4. Data transformations performed
9             5. Classification: ETL, Business Logic, Reporting, Maintenance
10            6. Side effects (inserts, updates, deletes, DDL changes)
11            7. Error handling patterns
12            8. Dependencies on other procedures or functions""",
13         messages=[{"role": "user",
14                   "content": f"Name: {proc['name']}\n\n{proc['definition']}"},
15                    ],
16         )
17     text = message.content[0].text
18     return json.loads(text[text.find("{"):text.rfind("}")+1])

```

Listing 4.4: Stored Procedure Business Logic Extractor

Common Legacy System Patterns

When working with legacy systems, CLAUDE CODE can help you recognize and document several common patterns:

- **Trigger-based CDC:** Some systems use database triggers to populate audit tables. `CLAUDE CODE` can parse trigger definitions to understand what changes are tracked and how.
- **EAV (Entity-Attribute-Value) models:** A pattern where entities have a flexible set of attributes stored as rows rather than columns. `CLAUDE CODE` can help you pivot EAV data into a conventional schema.
- **Encoded columns:** Binary or custom-encoded data stored in `VARCHAR` or `BLOB` columns. `CLAUDE CODE` can often identify encoding patterns (Base64, custom delimited formats, packed binary) from sample values.
- **Temporal tables:** Manual implementations of slowly changing dimensions, with `valid_from/valid_to` columns that predate native temporal table support.

Note

COBOL and Mainframe Systems `CLAUDE CODE` can parse COBOL copybooks and generate Python unpacking code. Feed it the copybook and ask for a struct-based parser—it handles `PIC` clauses, `REDEFINES`, and `OCCURS` with surprising accuracy. For EBCDIC-encoded data, `CLAUDE CODE` can generate the appropriate codepage conversion code. Yes, there is still COBOL in production. Yes, it will outlive us all. Use `CLAUDE CODE` for Everything™, even the things that predate the internet.

Warning

Stored Procedure Side Effects Legacy stored procedures often have non-obvious side effects: updating audit tables, sending emails via `xp_sendmail`, writing to file systems, or calling external services. One company discovered their “simple reporting procedure” was also sending faxes. In 2025. Always ask `CLAUDE CODE` to identify side effects, and verify them with the operations team before building pipelines that might bypass them.

4.5 Generating ERDs from DDL

Entity-Relationship Diagrams are essential for communicating data models to stakeholders, but they are tedious to create and maintain manually. `CLAUDE CODE` generates ERDs in Mermaid, PlantUML, or DBML format from raw DDL:

```

1 def generate_erd(ddl: str, format: str = "mermaid") -> str:
2     message = client.messages.create(
3         model="claude-sonnet-4-20250514", max_tokens=8000,
4         messages=[{"role": "user",
5                   "content": f"Convert this DDL to a {format} ERD.\n"
6                             f"Infer relationships from FKs and naming.\n"
7                             f"Group by business domain.\n"
8                             f"Color-code by domain.\n"
9                             f"Include cardinality (1:1, 1:N, M:N).\n\n"
10                            f"DDL:\n{ddl}"}],
11     )
12     return message.content[0].text

```

Listing 4.5: ERD Generation

The resulting Mermaid or PlantUML diagram can be rendered in Confluence, GitHub/GitLab, or any documentation system that supports these formats. This makes it trivial to keep ERDs up-to-date: re-run the generation whenever the schema changes.

Tip

ERD Best Practices Generate separate ERDs for each business domain rather than one massive diagram. A 200-table ERD is unreadable. Ask CLAUDE CODE to identify domains first, then generate domain-specific ERDs with cross-domain relationships shown as simplified references.

From ERDs to Data Models

Once you have a clear ERD, CLAUDE CODE can help you design your analytical data model. Feed it the ERD (or the DDL) along with your analytical requirements, and ask it to propose a star schema or OBТ (One Big Table) design:

```

1 def design_star_schema(source_ddl: str, requirements: str) -> str:
2     message = client.messages.create(
3         model="claude-sonnet-4-20250514", max_tokens=8000,
4         system="""You are a data modeling expert. Design a star schema
5                 for an analytical warehouse. Include:
6                 1. Fact tables with grain definition
7                 2. Dimension tables with SCD type recommendation
8                 3. DDL for the target warehouse
9                 4. dbt model SQL for each table
10                5. Rationale for modeling decisions""",

```

```
11     messages=[{"role": "user",
12                "content": f"Source schema:\n{source_ddl}\n\n"
13                        f"Analytical requirements:\n{requirements}"}],
14     )
15     return message.content[0].text
```

Listing 4.6: Star Schema Design from Source ERD

4.6 API Contract Validation

Breaking API changes are a top cause of pipeline failures. An API provider adds a field, changes a type, renames an endpoint, or deprecates a version—and your pipeline breaks at 3 AM. CLAUDE CODE compares old and new API responses, classifies changes as BREAKING, WARNING, or SAFE, and generates pytest contract tests.

```
1 def detect_api_changes(old_response: dict, new_response: dict,
2                        endpoint: str) -> dict:
3     message = client.messages.create(
4         model="claude-sonnet-4-20250514", max_tokens=4000,
5         system="""Compare two API responses from the same endpoint.
6         Classify each change as:
7         - BREAKING: removed field, type change, semantic change
8         - WARNING: new required field, format change
9         - SAFE: new optional field, additional enum value
10        Return JSON with changes array and overall_risk.""",
11        messages=[{"role": "user",
12                   "content": f"Endpoint: {endpoint}\n\n"
13                           f"Old response:\n{json.dumps(old_response, indent
14                   =2)}\n\n"
15                           f"New response:\n{json.dumps(new_response, indent
16                   =2)}"}],
17        )
18     text = message.content[0].text
19     return json.loads(text[text.find("{"):text.rfind("}")+1])
```

Listing 4.7: API Contract Change Detection

! Tip

Contract Testing in CI/CD Schedule nightly contract tests against source APIs. When CLAUDE CODE detects a breaking change, it can generate the migration code for your pipeline. Store historical API responses in version control to maintain a change log and enable regression testing.

4.7 Source System Profiling

Schema analysis tells you the *structure* of data; profiling tells you the *reality*. Data profiling examines actual data for distributions, quality patterns, and anomalies that the schema alone cannot reveal.

Key profiling dimensions include:

- **Completeness:** What percentage of rows have non-null values for each column?
- **Uniqueness:** Are columns that should be unique actually unique? Do “unique” columns have near-duplicates (e.g., trailing spaces)?
- **Distribution:** What are the min, max, mean, median, and standard deviation for numeric columns? What are the cardinality and top values for categorical columns?
- **Temporal patterns:** How does row count vary by date? Are there gaps? Weekday/weekend patterns? Seasonality?
- **Referential integrity:** Do foreign key values always have matching parent records? What is the orphan rate?
- **Format consistency:** Do phone numbers, addresses, and dates follow consistent formats?

CLAUDE CODE enhances profiling by interpreting results in business context—identifying likely enum values, soft-delete indicators, temporal patterns, and multi-table join keys from data patterns alone.

```
1 -- Generate profiling statistics for a table
2 -- Claude Code can customize this template per table
3 WITH column_stats AS (
```

```
4 SELECT
5     'users' AS table_name,
6     COUNT(*) AS total_rows,
7     COUNT(DISTINCT user_id) AS distinct_users,
8     COUNT(*) - COUNT(email) AS null_emails,
9     COUNT(DISTINCT status) AS status_cardinality,
10    MIN(created_at) AS earliest_record,
11    MAX(created_at) AS latest_record,
12    COUNT(*) FILTER (WHERE deleted_at IS NOT NULL)
13        AS soft_deleted_count,
14    COUNT(*) FILTER (
15        WHERE created_at > CURRENT_DATE - INTERVAL '7 days'
16    ) AS recent_7d_count
17 FROM source.users
18 )
19 SELECT * FROM column_stats;
```

Listing 4.8: Comprehensive profiling query template

Warning

Sampling Bias Always use stratified sampling when profiling large tables. If the table has a date column, sample across all time periods. A profile based on only recent data may miss historical patterns—such as a column that used to allow nulls but was made required six months ago, or a status value that was retired.

Note

Profiling vs. Privacy When profiling tables that may contain PII, use aggregate statistics (counts, distributions, patterns) rather than individual values. Never extract sample PII rows to send to an external API. Profile locally and send only the aggregate results to CLAUDE CODE for interpretation.

4.8 Case Study: Onboarding a New Data Source

Scenario: your company acquires a SaaS product with a 180-table PostgreSQL database, 45 API endpoints, and minimal documentation. A perfectly normal Tuesday in enterprise data engineering. The integration team needs a data model, ingestion plan, and initial data quality assessment within two weeks.

Week 1: Discovery and Documentation

1. **Day 1:** Extract DDL from the PostgreSQL database. Run the chunked schema analysis pipeline. CLAUDE CODE produces a structured catalog of all 180 tables, grouped into seven business domains (users, subscriptions, billing, content, analytics, notifications, admin). It identifies 23 implicit foreign key relationships not declared in the schema.
2. **Day 2:** Run the API analyzer against all 45 endpoints. CLAUDE CODE generates API documentation covering request/response schemas, pagination patterns (cursor-based for listing endpoints, no pagination for detail endpoints), rate limits (100 requests/minute), and authentication (OAuth 2.0 with refresh tokens).
3. **Day 3:** Feed the 47 stored procedures to the business logic extractor. CLAUDE CODE classifies 18 as ETL (including a nightly billing reconciliation), 15 as business logic (subscription tier calculations, proration rules), 8 as reporting, and 6 as maintenance (index rebuilds, partition management).
4. **Days 4–5:** Run profiling queries on the 30 most important tables. Feed profiling results to CLAUDE CODE for interpretation. Key findings: the `subscriptions` table has a 3% orphan rate on `plan_id`, suggesting a deleted plans table; the `events` table has a 12-hour gap every Sunday (maintenance window); timestamps are in US/Pacific, not UTC.

Week 2: Design and Implementation

1. **Days 6–7:** Design the target data model with CLAUDE CODE's help. Generate dbt staging models for the 30 priority tables with appropriate tests and documentation.
2. **Days 8–9:** Build the ingestion pipeline. CLAUDE CODE generates Airflow DAGs for both database CDC (using Debezium) and API extraction (custom Python connectors).
3. **Day 10:** Run the full pipeline end-to-end on staging. Use CLAUDE CODE to generate data quality assertions comparing source and target row counts, key distributions, and aggregate values.

Table 4.1: Onboarding Time: Traditional vs. CLAUDE CODE-Assisted

Activity	Traditional	Claude Code-Assisted
Schema documentation	3–5 days	2 hours
API documentation	2–3 days	1 hour
Business rule extraction	1–2 weeks	3 hours
Data profiling (20 tables)	2–3 days	1 hour
ERD generation	1–2 days	15 min
Target data model design	3–5 days	1 day
Ingestion pipeline build	1–2 weeks	2 days
Total	5–8 weeks	2 weeks

! Tip

Iterate, Don't Trust Blindly CLAUDE CODE's initial analysis is a starting point. Validate critical findings with the source system team—assuming the source system team still exists and has not been “restructured.” The real gain is accelerating the path to *informed* questions—instead of spending a week figuring out what to ask, you arrive at the right questions on day one.

4.9 Building a Reusable Source Analysis Toolkit

As your team onboards more sources, you will want a standardized toolkit rather than ad-hoc scripts. A well-designed toolkit includes:

1. **Schema analyzer:** The DDL analysis pipeline from Section 4.2, with output stored in a metadata repository.
2. **API prober:** The API analyzer from Section 4.3, extended with authentication support for OAuth, API keys, and custom auth schemes.
3. **Profiler:** Parameterized profiling queries that run against any source database and produce standardized output.
4. **Contract tester:** Scheduled tests that detect source schema and API changes before they break pipelines.
5. **Documentation generator:** Templates that produce consistent documentation for every source, stored in your data catalog.

Note

Standardized Output Formats Store all analysis output in a consistent JSON schema. This enables comparison across sources, trend analysis over time, and integration with your data catalog. CLAUDE CODE can generate the JSON schema for your analysis output format.

4.10 Exercises

1. **Schema Archaeology:** Export the DDL from a database you work with. Run the schema analysis pipeline and compare the results to your existing documentation. What did CLAUDE CODE discover that you did not know?
2. **API Documentation:** Choose a SaaS API your team uses (Stripe, Salesforce, HubSpot, etc.). Use the API analyzer to probe five endpoints and generate documentation. Compare the generated documentation to the vendor's official docs—where does CLAUDE CODE add value, and where does it fall short?
3. **Legacy Procedure Analysis:** Find a stored procedure in your environment that no one fully understands. Feed it to CLAUDE CODE and validate the extracted business rules with someone who knows the system. How accurate was the analysis?
4. **Profiling Interpretation:** Run profiling queries against a table with known data quality issues. Feed the results to CLAUDE CODE without mentioning the known issues. Does it identify them? What additional issues does it surface?
5. **Contract Test Suite:** Use CLAUDE CODE to generate a contract test suite for one of your source APIs. Run it daily for a week. Document any changes detected and how you would handle them in your pipeline.

This chapter demonstrated that CLAUDE CODE transforms source system understanding from slow, manual archaeology into a structured, rapid process. The engineer's expertise is amplified, not replaced—you still need to know which questions to ask, how to validate answers, and which findings matter for your use case. But at least you no longer need to spend three days staring at a column named `flg_xref_2b` wondering what it means. CLAUDE CODE will stare at it for you.

5

Data Ingestion Pipelines Powered by Claude Code

Every API is a unique snowflake of dysfunction. CLAUDE CODE speaks fluent dysfunction.

Data ingestion is the circulatory system of every data platform. Every API has its own authentication scheme, pagination logic, rate limits, and schema quirks. Every database has its own CDC mechanism. This chapter demonstrates how CLAUDE CODE transforms ingestion from a repetitive slog into a structured, accelerated, and more reliable process.

5.1 Traditional Ingestion Patterns

Batch Ingestion

Batch ingestion extracts data at scheduled intervals. It is appropriate when business requirements tolerate hours-old data and source systems cannot support continuous extraction. Batch pipelines are simpler to reason about, cheaper to operate, and easier to recover when something fails—you simply rerun the batch.

The canonical pattern uses a *high watermark* to track progress: each run queries the target for the maximum value of an incremental column, then extracts only records newer than that value.

```
1 from dataclasses import dataclass
2 from datetime import datetime
3 from typing import Optional
4
5 @dataclass
6 class BatchConfig:
7     source_table: str
8     target_table: str
9     incremental_column: str
10    batch_size: int = 10_000
11
12 class BatchIngestionPipeline:
13     def __init__(self, config: BatchConfig, source_conn, target_conn):
14         self.config = config
15         self.source = source_conn
16         self.target = target_conn
17
18     def get_high_watermark(self) -> Optional[datetime]:
19         result = self.target.execute(
20             f"SELECT MAX({self.config.incremental_column}) "
21             f"FROM {self.config.target_table}")
22         return result.scalar()
23
24     def run(self):
25         watermark = self.get_high_watermark()
26         query = f"SELECT * FROM {self.config.source_table}"
27         if watermark:
28             query += f" WHERE {self.config.incremental_column} > '{
29 watermark.isoformat()}'"
30         records = self.source.execute(query)
31         batch = []
32         for record in records:
33             batch.append(record)
34             if len(batch) >= self.config.batch_size:
35                 self.target.bulk_insert(self.config.target_table, batch)
36                 batch = []
37         if batch:
38             self.target.bulk_insert(self.config.target_table, batch)
```

Listing 5.1: Classic Batch Ingestion Pattern

Note

Watermark Pitfalls High-watermark ingestion assumes the incremental column is monotonically increasing. Late-arriving records or back-dated updates will be missed. And they will arrive late—because data, like houseguests, never arrives when expected. For sources with this behavior, consider a lookback window: subtract a safety margin (e.g., two hours) from the watermark before querying.

Streaming Ingestion

Streaming ingestion processes data continuously with sub-second latency. It requires handling out-of-order events, exactly-once delivery, backpressure, and state management. Common implementations use Apache Kafka, Amazon Kinesis, or Google Pub/Sub as the transport layer, with a consumer framework (Flink, Spark Structured Streaming, or a custom consumer) performing the ingestion.

Warning

Streaming Is Not Always the Answer If your business tolerates five-minute-old data, micro-batch is almost certainly the right choice. Streaming requires 3–5× more engineering effort to build and maintain. Just because you *can* use CLAUDE CODE to build a Kafka-Flink-Iceberg streaming pipeline does not mean you *should*. Evaluate latency requirements honestly before committing to a streaming architecture. Most “we need real-time” requests translate to “we need it faster than once a day,” which is a micro-batch, not a Kafka cluster.

Micro-Batch: The Middle Ground

Micro-batch ingestion runs batch extractions at very short intervals—typically every one to fifteen minutes. It combines the simplicity of batch with near-real-time freshness. Tools like Spark Structured Streaming’s trigger-once mode and Airflow’s timetable API make micro-batch straightforward to implement.

5.2 Generating Ingestion Code from API Documentation

Writing API clients is one of the most time-consuming parts of ingestion. Each API has its own authentication flow, pagination style, rate-limiting headers, and

error format. CLAUDE CODE can read raw API documentation and generate a production-quality client in minutes.

```
1 from anthropic import Anthropic
2
3 client = Anthropic()
4
5 def generate_api_client(api_docs: str) -> str:
6     response = client.messages.create(
7         model="claude-sonnet-4-20250514", max_tokens=8000,
8         system="""Generate production-quality Python ingestion code:
9             1. Handle auth (API keys, OAuth2, JWT)
10            2. Implement pagination (cursor, offset, keyset)
11            3. Respect rate limits with exponential backoff
12            4. Support incremental extraction via watermarks
13            5. Include type hints and docstrings
14            Generate COMPLETE, RUNNABLE code."""),
15         messages=[{"role": "user",
16                   "content": f"Generate ingestion client:\n\n{api_docs}"},
17                    ],
18         )
19     return response.content[0].text
```

Listing 5.2: Generating an API Client from Documentation

Tip

Version-Pin Generated Code Store the API spec alongside generated code in version control. Include a comment referencing the spec version and generation date. When the upstream API changes—and it will, at 2 AM, without a changelog—re-generate from the updated spec and diff against the previous version to understand what changed.

The generated code should be treated as a starting point. Review it carefully, run it against the live API in a development environment, and add integration tests before promoting to production.

5.3 Schema Inference and Mapping

CLAUDE CODE understands semantic meaning: "2024-03-15T10:30:00Z" is an ISO 8601 timestamp, not merely a string. "192.168.1.1" is an IPv4 address. This semantic awareness produces far better type mappings than rule-based inference

tools.

```

1 def infer_schema_with_claude(sample_records: list[dict]) -> dict:
2     import json
3     sample = json.dumps(sample_records[:20], indent=2, default=str)
4     response = client.messages.create(
5         model="claude-sonnet-4-20250514", max_tokens=4000,
6         messages=[{"role": "user",
7                   "content": f"""Infer schema from these records. For each
8 field:
9 field_name, data_type, semantic_type, nullable, format_pattern,
10 suggested_warehouse_type, notes. Return JSON.\n\n{sample}"""}],
11     )
12     return json.loads(response.content[0].text
13                        .split("```json")[-1].split("```")[0])

```

Listing 5.3: Schema Inference with Semantic Understanding

Note

Confidence Thresholds Establish a confidence threshold (0.8 is reasonable) below which a human must review the mapping. Lower-confidence fields go to a review queue. Fields with ambiguous types—such as numeric strings that might be IDs or quantities—particularly benefit from human review.

Schema inference is a **design-time** activity. Run it once when onboarding a new source, review the results, and commit the finalized schema mapping to version control. Do not call CLAUDE CODE at runtime for every record.

Tip

Mapping to Warehouse Types When CLAUDE CODE suggests warehouse types, it considers storage efficiency and query performance. For example, it maps UUIDs to VARCHAR(36) rather than TEXT, uses NUMERIC(12,2) for currency rather than FLOAT (which introduces rounding errors), and recommends TIMESTAMPTZ over TIMESTAMP for any field that may cross timezones.

5.4 Building a Universal API Connector

The universal connector has four layers: configuration (YAML), discovery (CLAUDE CODE reads docs), execution (generic HTTP with pagination/rate limits), and out-

put (normalized data). This architecture lets you onboard new APIs by writing a YAML configuration file rather than a new Python module.

```
1 import yaml, requests, time, os, logging
2 from datetime import datetime, timezone
3 from typing import Generator, Optional
4
5 logger = logging.getLogger(__name__)
6
7 class UniversalAPIConnector:
8     def __init__(self, config_path: str):
9         with open(config_path) as f:
10             self.config = yaml.safe_load(f)["connector"]
11             self.base_url = self.config["api"]["base_url"]
12             self.session = requests.Session()
13             auth = self.config["api"]["auth"]
14             if auth["type"] == "bearer_token":
15                 token = os.environ[auth["token_env_var"]]
16                 self.session.headers["Authorization"] = f"Bearer {token}"
17
18     def extract_endpoint(
19         self, endpoint_name: str, since: Optional[datetime] = None,
20     ) -> Generator[dict, None, None]:
21         ep = next(e for e in self.config["endpoints"]
22                 if e["name"] == endpoint_name)
23         url = f"{self.base_url}{ep['path']}"
24         params = dict(ep.get("params", {}))
25         if since and "incremental" in ep:
26             params[ep["incremental"]["field"]] = int(since.timestamp())
27         pagination = self.config["api"]["pagination"]
28         while True:
29             response = self.session.get(url, params=params, timeout=30)
30             if response.status_code == 429:
31                 time.sleep(int(response.headers.get("Retry-After", 60)))
32                 continue
33             response.raise_for_status()
34             data = response.json()
35             records = data.get("data", data.get("results", []))
36             for record in records:
37                 record["_ingested_at"] = datetime.now(timezone.utc).
38                 isoformat()
39                 yield record
40             if not data.get(pagination.get("has_more_field", ""), False):
41                 break
```

```
41         params[pagination["cursor_param"]] = records[-1][pagination["cursor_field"]]
```

Listing 5.4: Universal API Connector (core extract loop)

```
1 connector:
2   api:
3     base_url: "https://api.example.com/v2"
4     auth:
5       type: bearer_token
6       token_env_var: EXAMPLE_API_KEY
7     pagination:
8       cursor_param: starting_after
9       cursor_field: id
10      has_more_field: has_more
11   endpoints:
12     - name: transactions
13       path: /transactions
14       params:
15         limit: 100
16         incremental:
17           field: created_gte
```

Listing 5.5: Sample YAML Configuration for the Universal Connector

! Tip

Test with Recorded Responses Use a library like `responses` or `vcrpy` to record real API responses and replay them in tests. This decouples your test suite from the live API and makes tests fast, deterministic, and free of rate-limit concerns.

5.5 CDC Patterns with Claude Code Assistance

Change Data Capture (CDC) tracks row-level changes in a source database and replicates them to a target. The three major approaches differ in latency, source impact, and implementation complexity.

CLAUDE CODE generates optimized Debezium configurations and CDC event processors with merge logic, including proper handling of inserts, updates, and soft-deletes. When you describe your source schema and target warehouse, CLAUDE

Table 5.1: CDC Approach Comparison

Approach	Latency	Source Impact	Complexity
Timestamp-based	Minutes–hours	Low	Low
Trigger-based	Seconds	High	Medium
Log-based	Seconds	Minimal	High

CODE produces a complete Debezium connector configuration, a Kafka consumer that deserializes change events, and the MERGE SQL that applies those changes to the target.

Warning

CDC and Schema Changes Log-based CDC captures the schema at the time the connector was configured. If the source schema changes (columns added, types altered), the CDC pipeline may break silently or drop data. “Silently” is doing a lot of heavy lifting in that sentence. Monitor schema registry compatibility checks and re-generate connectors after DDL changes. Or use CLAUDE CODE to generate a schema drift detector that pages you before the CFO notices the revenue dashboard is blank.

5.6 Handling Schema Evolution

Source schemas change: columns are added, removed, renamed, or retyped. CLAUDE CODE analyzes before-and-after snapshots, classifies each change (breaking vs. safe), and generates migration DDL.

```

1 def detect_schema_changes(old_schema: dict, new_schema: dict) -> list:
2     response = client.messages.create(
3         model="claude-sonnet-4-20250514", max_tokens=4000,
4         messages=[{"role": "user",
5                   "content": f"""Compare schemas and classify changes as
6 column_added, column_removed, type_changed, etc.
7 For each: is_breaking, suggested_action.
8 Old: {json.dumps(old_schema)}
9 New: {json.dumps(new_schema)}
10 Return JSON array. """}],
11     )
12     return json.loads(response.content[0].text

```

```
13 .split("```json")[-1].split("```")[0]
```

Listing 5.6: Schema Change Detection

Safe changes (adding a nullable column) can be applied automatically. Breaking changes (removing a column, narrowing a type) should trigger a review workflow. CLAUDE CODE can also generate the downstream impact analysis—which views, stored procedures, dbt models, and dashboards reference the affected columns.

! Warning

Schema Evolution in Production Never apply auto-generated migrations to production without human review. This is one of those rare moments where you might need to think for yourself. We apologize. Schema changes cascade through views, procedures, pipelines, and dashboards. A column rename that seems harmless can break dozens of downstream consumers. It's like pulling one thread on a sweater, except the sweater is your entire data platform and the thread is named `user_id_v2_final_FINAL`.

5.7 Error Handling and Dead Letter Queues

A well-designed *dead letter queue* (DLQ) captures the failed record, error type, stack trace, and retry count. CLAUDE CODE analyzes DLQ failure patterns, identifies root causes, and recommends remediations.

Common failure categories include malformed records (parsing errors), constraint violations (duplicate keys, NULL in required fields), schema mismatches (unexpected fields or types), and transient errors (timeouts, rate limits). Each category demands a different remediation strategy.

```
1 from dataclasses import dataclass, field
2 from datetime import datetime, timezone
3 from typing import Any
4
5 @dataclass
6 class DLQRecord:
7     original_record: dict
8     error_type: str
9     error_message: str
10    stack_trace: str
11    source_system: str
```

```

12 pipeline_name: str
13 retry_count: int = 0
14 max_retries: int = 3
15 created_at: datetime = field(
16     default_factory=lambda: datetime.now(timezone.utc))
17
18 @property
19 def is_retryable(self) -> bool:
20     return (self.retry_count < self.max_retries
21         and self.error_type in ("timeout", "rate_limit", "
22     connection"))

```

Listing 5.7: DLQ Record Structure

! Tip

DLQ Monitoring Set alerts on DLQ depth and growth rate. A sudden spike usually indicates a systemic issue (API change, schema change, auth expiration) rather than individual bad records. Alert on both absolute depth and rate of change. A DLQ that grows faster than your actual data table is nature's way of telling you something is profoundly wrong.

5.8 Data Validation During Ingestion

Traditional validation uses static rules (type checks, range checks, regex). CLAUDE CODE adds *semantic validation*: detecting logical inconsistencies (end_date before start_date), implausible values (negative prices, future birth dates), and cross-field violations (country code mismatching phone format).

Use CLAUDE CODE at **design time** to generate validation rules, not at runtime for every record. Reserve runtime CLAUDE CODE calls for truly ambiguous records that cannot be handled deterministically.

```

1 def generate_validation_rules(schema: dict, sample_data: list) -> str:
2     response = client.messages.create(
3         model="claude-sonnet-4-20250514", max_tokens=4000,
4         messages=[{"role": "user",
5             "content": f"""Given this schema and sample data, generate
6 a Python validation function that checks:
7 1. Type conformance 2. Required fields 3. Range constraints
8 4. Cross-field consistency 5. Format patterns

```

```

9 Return a function: def validate_record(record: dict) -> list[str]
10 that returns a list of validation error messages (empty if valid).
11
12 Schema: {json.dumps(schema)}
13 Sample: {json.dumps(sample_data[:5], default=str)}"""},
14 )
15 return response.content[0].text

```

Listing 5.8: Validation Rule Generation with Claude Code

Note

Three-Tier Validation Structure validation in three tiers: (1) syntactic checks (types, formats) that reject records immediately, (2) semantic checks (business rules) that flag records for review, and (3) statistical checks (distribution anomalies) that trigger alerts but do not block ingestion.

5.9 Full Example: Stripe to Snowflake Pipeline

The pipeline follows: Extract (Stripe API with pagination) → Stage (Parquet on S3) → Load (COPY into Snowflake raw tables) → Monitor. This is a common pattern for SaaS API ingestion, and CLAUDE CODE can generate most of the boilerplate.

```

1 CREATE SCHEMA IF NOT EXISTS analytics.raw_stripe;
2
3 CREATE OR REPLACE TABLE analytics.raw_stripe.customers (
4     raw_data VARIANT,
5     _loaded_at TIMESTAMP_NTZ DEFAULT CURRENT_TIMESTAMP(),
6     _file_name VARCHAR
7 );
8
9 CREATE OR REPLACE VIEW analytics.raw_stripe.v_customers AS
10 SELECT
11     raw_data:id::VARCHAR AS customer_id,
12     raw_data:email::VARCHAR AS email,
13     raw_data:name::VARCHAR AS name,
14     TO_TIMESTAMP(raw_data:created::INT) AS created_at,
15     raw_data:currency::VARCHAR AS currency,
16     raw_data:delinquent::BOOLEAN AS is_delinquent,
17     _loaded_at
18 FROM analytics.raw_stripe.customers;

```

Listing 5.9: Snowflake Setup for Stripe Ingestion

The raw-then-view pattern is intentional: raw tables store the original JSON unchanged (for auditability and reprocessing), while views expose typed columns for downstream consumption.

5.10 Performance Considerations

The key principle is: use CLAUDE CODE at design time to generate code and rules, then execute that generated code deterministically at runtime.

Table 5.2: When to Use CLAUDE CODE vs. Deterministic Code

Task	Approach	Rationale
Schema inference	CLAUDE CODE (design-time)	One-time analysis
API client generation	CLAUDE CODE (design-time)	Generate code, then run it
Validation rules	CLAUDE CODE (design-time)	Generate rules, execute deterministically
Record validation	Deterministic (runtime)	Rules generated by CLAUDE CODE
Error triage	CLAUDE CODE (on-demand)	Analyze DLQ on failure
Schema evolution	CLAUDE CODE (on-change)	Triggered by schema drift detection

Warning

API Cost Management Each CLAUDE CODE call has a token cost. For millions of records, even batched calls accumulate quickly. Calculate expected costs before enabling runtime processing. A single CLAUDE CODE call per record at \$0.003 per call costs \$3,000 for one million records. Use CLAUDE CODE for Everything™ does not mean “use CLAUDE CODE for every row.” That way lies a finance team with questions.

5.11 Exercises

1. **API Connector Generation:** Choose a public API (GitHub, Spotify, or OpenWeatherMap). Feed its documentation to CLAUDE CODE, generate a complete ingestion client, test against the live API, and document corrections

needed. Measure how long generation took vs. how long manual implementation would have taken.

2. **CDC Pipeline:** Set up PostgreSQL with logical replication. Implement basic log-based CDC using Python. Use `CLAUDE CODE` to generate SQL merge statements for the target table. Test with inserts, updates, and deletes.
3. **Schema Evolution Simulation:** Create two schema versions with ten differences (added columns, removed columns, type changes, renames). Use `CLAUDE CODE` to detect, classify, and generate migration DDL. Verify against a test database and confirm both forward and rollback migrations work.
4. **End-to-End Pipeline:** Build a complete pipeline from a source of your choice to a warehouse. Include extraction with pagination, validation (using rules generated by `CLAUDE CODE`), staging in Parquet, loading, and monitoring. Use `CLAUDE CODE` to generate at least three components.
5. **DLQ Analysis:** Populate a dead letter queue with 50 synthetic failure records spanning at least five error categories. Use `CLAUDE CODE` to analyze the failures, identify patterns, and recommend remediations. Compare the AI-generated analysis to your own manual review.

6

Working with Unstructured Data

For decades, 80% of enterprise data sat in a corner, unstructured and unloved, like the office plant nobody waters. CLAUDE CODE is here to water that plant.

Industry estimates suggest 80–90% of enterprise data is unstructured: documents, emails, images, log files, audio recordings. For decades, data engineers largely ignored this data because extracting structure was prohibitively expensive. CLAUDE CODE changes this calculus. What once required specialized OCR, custom NLP pipelines, and teams of annotators can now be accomplished with well-crafted prompts.

i Note

The 80/20 Rule of Unstructured Data In most organizations, 80% of document volume comes from 20% of document types. Start with your highest-volume types. CLAUDE CODE makes it feasible to handle the long tail without custom extractors for each format.

6.1 The Unstructured Data Challenge

Traditional approaches—regex, OCR, NLP, template matching—are brittle: they work for the specific cases they were designed for and fail unpredictably on anything else. CLAUDE CODE brings general-purpose understanding that degrades gracefully on unfamiliar formats.

The fundamental shift is from *programming extraction rules* to *describing what you want extracted*. Instead of writing regex patterns for every possible date format—a task that has driven more engineers to madness than any other—you tell CLAUDE CODE “extract all dates” and it handles 2024-03-15, March 15, 2024, 15/03/2024, and 3-15-24 without separate rules for each. You could also ask CLAUDE CODE to extract dates from your therapist’s notes about your regex-induced trauma. Use CLAUDE CODE for Everything™.

Tip

Start with a Taxonomy Before building extraction pipelines, catalog your unstructured data sources. For each source type, document: volume (documents per day), format (PDF, email, image), complexity (structured tables vs. free-form prose), and business value (what decisions depend on this data). This taxonomy drives prioritization.

6.2 Document Parsing: PDFs, Invoices, Contracts

PDF Text Extraction

PDF is the most common document format in enterprise settings, but it is also one of the hardest to parse. PDFs encode visual layout, not semantic structure. A table that looks clean to a human reader may be stored as dozens of independent text fragments with absolute positioning.

```
1 import io
2 from dataclasses import dataclass
3 from typing import Optional
4
5 @dataclass
6 class ExtractedDocument:
7     text: str
8     page_count: int
9     extraction_method: str # "text", "ocr", "hybrid"
10    confidence: float
11
12 def extract_pdf_text(pdf_path: str, ocr_fallback: bool = True) ->
13     ExtractedDocument:
14     import fitz # PyMuPDF
15     doc = fitz.open(pdf_path)
16     pages_text, ocr_pages = [], 0
```

```

16     for page_num in range(len(doc)):
17         page = doc[page_num]
18         text = page.get_text("text").strip()
19         if len(text) < 50 and ocr_fallback:
20             import pytesseract
21             from PIL import Image
22             pix = page.get_pixmap(dpi=300)
23             img = Image.open(io.BytesIO(pix.tobytes("png")))
24             ocr_text = pytesseract.image_to_string(img, lang="eng")
25             if ocr_text and len(ocr_text) > len(text):
26                 text, ocr_pages = ocr_text, ocr_pages + 1
27             pages_text.append(f"--- Page {page_num + 1} ---\n{text}")
28     doc.close()
29     method = "ocr" if ocr_pages == len(pages_text) else (
30         "hybrid" if ocr_pages > 0 else "text")
31     return ExtractedDocument(
32         "\n\n".join(pages_text), len(pages_text), method,
33         1.0 if method == "text" else 0.85)

```

Listing 6.1: PDF Extraction with OCR Fallback

⚠ Warning

Scanned PDFs and OCR Quality Scanned PDFs—especially older documents, faxes, or low-resolution scans—produce noisy OCR output. If someone in your organization is still sending faxes in 2026, you have bigger problems than OCR accuracy. Always log the extraction method and confidence score. For documents where accuracy is critical (financial statements, legal contracts), route low-confidence extractions to a human review queue.

Parsing Invoices

```

1 from anthropic import Anthropic
2 import json
3 from datetime import datetime, timezone
4
5 client = Anthropic()
6
7 def parse_invoice(document_text: str, source_filename: str = "") -> dict:
8     response = client.messages.create(

```

```
9     model="claude-sonnet-4-20250514", max_tokens=4000,
10     messages=[{"role": "user", "content":
11         f"""Parse this invoice. Extract: invoice_number,
12         invoice_date,
13         due_date, vendor_name, customer_name, currency, subtotal, tax_amount,
14         total_amount, payment_terms, line_items (array with description,
15         quantity, unit_price, amount). Return ONLY valid JSON.
16         Document:\n{document_text}"""}],
17     )
18     raw = response.content[0].text
19     if "```json" in raw:
20         raw = raw.split("```json")[1].split("```")[0]
21     parsed = json.loads(raw)
22     parsed["_source_file"] = source_filename
23     parsed["_parse_timestamp"] = datetime.now(timezone.utc).isoformat()
24     return parsed
```

Listing 6.2: Invoice Parsing with CLAUDE CODE

Contract Analysis

Contracts require semantic understanding beyond structured extraction. CLAUDE CODE extracts metadata (type, dates, parties, governing law), key terms (payment, termination, renewal, liability caps, IP ownership), obligations per party, and risk flags (unusual clauses, missing standards, ambiguous language).

Tip

Multi-Pass for Long Documents For contracts exceeding 50 pages, first extract metadata and a table of contents, then process each section individually, and synthesize section-level results into a complete analysis. This multi-pass approach avoids context window limitations and produces more accurate results than trying to process the entire document at once.

6.3 Email Parsing and Classification

Email is a rich but messy data source. A single email thread may contain customer complaints, order references, PII, and attachments—all in free-form text with inconsistent formatting. CLAUDE CODE excels at extracting structure from this chaos.

```
1 def classify_email(from_addr: str, subject: str, body: str) -> dict:
2     response = client.messages.create(
3         model="claude-sonnet-4-20250514", max_tokens=2000,
4         messages=[{"role": "user", "content":
5             f"""Classify this email. Return JSON:
6             category (customer_support, billing_issue, complaint, feature_request,
7             vendor_communication, internal_notification, spam),
8             sentiment (positive/neutral/negative/urgent),
9             priority (low/medium/high/critical),
10            entities (order numbers, account IDs, amounts),
11            summary (one sentence).
12
13            From: {from_addr}
14            Subject: {subject}
15            Body: \n{body[:3000]}"""],
16        )
17     raw = response.content[0].text
18     if "```json" in raw:
19         raw = raw.split("```json")[1].split("```")[0]
20     return json.loads(raw)
```

Listing 6.3: Email Classification

 **Warning**

PII in Emails Emails frequently contain PII: names, addresses, phone numbers, social security numbers, and financial details. Before sending content to any external API, ensure compliance with GDPR, HIPAA, and CCPA. Consider PII redaction before the CLAUDE CODE call. At minimum, strip email signatures, which often contain personal contact information, inspirational quotes, and legal disclaimers longer than the email itself.

For high-volume email processing, classify emails into buckets using a lightweight model first, then apply CLAUDE CODE only to emails that require deeper analysis. This tiered approach reduces costs by 70–80% while maintaining accuracy on the emails that matter most.

Email Thread Reconstruction

Email threads present additional challenges: quoted text, forwarded messages, and reply chains create nested structures. Before classification, reconstruct the thread

by parsing email headers (In-Reply-To, References), stripping quoted text, and ordering messages chronologically. CLAUDE CODE can then analyze the full thread context to produce more accurate classifications than single-message analysis.

6.4 Log File Analysis

When logs do not match known patterns (JSON, Apache combined, syslog), CLAUDE CODE can analyze sample lines and generate a Python parser function. For anomaly detection, CLAUDE CODE identifies cascading failures, security incidents, resource exhaustion, and unusual error patterns that statistical methods alone would miss.

```
1 def generate_log_parser(sample_lines: list[str]) -> str:
2     response = client.messages.create(
3         model="claude-sonnet-4-20250514", max_tokens=4000,
4         messages=[{"role": "user",
5                   "content": f"Analyze these log lines and generate a Python
6 function: def parse_log_line(line: str) -> dict
7 Extract: timestamp, level, source, message, identifiers.
8 Handle edge cases. Return ONLY Python code."},
9
10    {chr(10).join(sample_lines[:30])}],
11    )
12    code = response.content[0].text
13    if "`python" in code:
14        code = code.split("`python")[1].split("`")[0]
15    return code.strip()
```

Listing 6.4: CLAUDE CODE-Generated Log Parser

Note

Log Parser Validation Always validate a generated parser against a larger sample than the one used for generation. Test with at least 1,000 lines and check the parse success rate. A good parser should handle >99% of lines without errors. Edge cases to watch for: multi-line stack traces, Unicode characters, and log lines that span timezone changes.

6.5 Image-to-Data Pipelines

CLAUDE CODE's multimodal capabilities extract structured data directly from images, eliminating the need for separate OCR and computer vision pipelines. Common use cases include extracting data from receipts, business cards, whiteboard photos, charts, and scanned forms.

```
1 import base64
2 from pathlib import Path
3
4 def extract_data_from_image(image_path: str, instructions: str) -> dict:
5     path = Path(image_path)
6     media_map = {".png": "image/png", ".jpg": "image/jpeg", ".jpeg": "
7     image/jpeg"}
8     with open(image_path, "rb") as f:
9         data = base64.standard_b64encode(f.read()).decode("utf-8")
10    response = client.messages.create(
11        model="claude-sonnet-4-20250514", max_tokens=4000,
12        messages=[{"role": "user", "content": [
13            {"type": "image", "source": {"type": "base64",
14            "media_type": media_map.get(path.suffix.lower(), "image/
15            png"),
16            "data": data}},
17            {"type": "text", "text": f"{instructions}\n\nReturn as JSON.
18            "}],
19        ]}],
20    )
21    raw = response.content[0].text
22    if "```json" in raw:
23        raw = raw.split("```json")[1].split("```")[0]
24    return json.loads(raw)
```

Listing 6.5: Image-to-Data Extraction

! Tip

Image Resolution Matters For documents and receipts, aim for at least 150 DPI. For detailed tables, 300 DPI is recommended. Below 72 DPI, extraction accuracy degrades significantly. When preprocessing images, convert to grayscale, increase contrast, and deskew rotated documents before sending to CLAUDE CODE.

Batch Image Processing

For processing large volumes of images, implement a pipeline that reads from an object store, processes in parallel, and writes structured output to a staging table. Use a work queue (SQS, Redis, or Celery) to distribute work across multiple workers and handle retries for failed extractions.

```
1 from concurrent.futures import ThreadPoolExecutor, as_completed
2 import logging
3
4 logger = logging.getLogger(__name__)
5
6 def process_image_batch(image_paths: list[str],
7                         instructions: str,
8                         max_workers: int = 4) -> list[dict]:
9     results = []
10    with ThreadPoolExecutor(max_workers=max_workers) as executor:
11        futures = {
12            executor.submit(extract_data_from_image, path, instructions):
13            path
14            for path in image_paths
15        }
16        for future in as_completed(futures):
17            path = futures[future]
18            try:
19                result = future.result()
20                result["_source_path"] = path
21                result["_status"] = "success"
22                results.append(result)
23            except Exception as e:
24                logger.error(f"Failed to process {path}: {e}")
25                results.append({"_source_path": path,
26                               "_status": "failed",
27                               "_error": str(e)})
28    return results
```

Listing 6.6: Batch Image Processing Pipeline

Warning

Concurrency and Rate Limits When processing images in parallel, respect the API rate limits. The Anthropic API has per-minute token limits. Set

`max_workers` conservatively and implement a rate limiter. Processing 1,000 images at full concurrency will trigger rate limiting and degrade throughput. Think of rate limits as the bouncer at a club: you can bring friends, but not the entire city.

6.6 Chunking Strategies for Large Documents

Documents that exceed the context window must be split into chunks. The choice of chunking strategy significantly impacts extraction quality.

Table 6.1: Document Chunking Strategies

Strategy	Best For	Pros	Cons
Fixed-size	Uniform text	Simple	Splits mid-sentence
Sentence-based	Prose	Respects boundaries	Variable size
Section-based	Structured docs	Semantic units	Requires heading detection
Sliding window	Overlapping context	No lost context	Redundant processing

i Note

Chunk Size Trade-off 3,000–5,000 tokens per chunk with 10–15% overlap provides the best balance of accuracy and efficiency for most document types. Smaller chunks improve precision but increase cost and require more complex reassembly. Larger chunks risk hitting context limits and reduce extraction focus.

When reassembling chunk-level results, handle deduplication carefully. Overlapping chunks will produce duplicate extractions that must be merged. Use entity resolution (matching on key fields) rather than simple string comparison.

```

1 import re
2 from dataclasses import dataclass
3
4 @dataclass
5 class DocumentChunk:
6     text: str
7     section_title: str
8     page_range: tuple[int, int]
9     chunk_index: int

```

```

10     total_chunks: int
11
12     def chunk_by_sections(document_text: str,
13                           max_tokens: int = 4000) -> list[DocumentChunk]:
14         # Split on common heading patterns
15         sections = re.split(
16             r'\n(?:?:ARTICLE|SECTION|CHAPTER|\d+\.\.)\s)',
17             document_text)
18         chunks = []
19         for i, section in enumerate(sections):
20             title_match = re.match(r'^(.+?)[\n:]', section)
21             title = title_match.group(1).strip() if title_match else f"
Section {i+1}"
22             chunks.append(DocumentChunk(
23                 text=section.strip(),
24                 section_title=title,
25                 page_range=(0, 0), # populate from page markers
26                 chunk_index=i,
27                 total_chunks=len(sections),
28             ))
29         return chunks

```

Listing 6.7: Section-Based Chunking Implementation

6.7 Structured Output with Tool Use

For production systems, prefer CLAUDE CODE's tool use over JSON prompting. Tool use provides schema-level enforcement: the model *must* return data matching the schema.

```

1     def extract_invoice_with_tool_use(document_text: str) -> dict:
2         tools = [{"name": "store_invoice_data",
3                   "description": "Store extracted invoice data.",
4                   "input_schema": {"type": "object",
5                                     "properties": {
6                                         "invoice_number": {"type": "string"},
7                                         "invoice_date": {"type": "string"},
8                                         "vendor_name": {"type": "string"},
9                                         "currency": {"type": "string"},
10                                        "total_amount": {"type": "number"},
11                                        "line_items": {"type": "array", "items": {"type": "
object",

```

```
12         "properties": {"description": {"type": "string"},
13                       "amount": {"type": "number"}},
14         "required": ["description", "amount"]},
15     },
16     "required": ["invoice_number", "vendor_name", "total_amount", "
line_items"],
17     }]}
18 response = client.messages.create(
19     model="claude-sonnet-4-20250514", max_tokens=4000,
20     tools=tools,
21     tool_choice={"type": "tool", "name": "store_invoice_data"},
22     messages=[{"role": "user",
23               "content": f"Extract invoice data:\n\n{document_text}"},
24             ]
25     )
26     for block in response.content:
27         if block.type == "tool_use":
28             return block.input
29     raise ValueError("No tool use response")
```

Listing 6.8: Tool Use for Guaranteed Structured Output

Tip

Tool Use vs. JSON Prompting Tool use provides schema enforcement—the model *must* conform. JSON prompting relies on instruction-following, which occasionally fails for complex schemas, usually at the worst possible moment. Tool use also makes schema evolution easier: add new fields to the tool schema and the model will populate them without prompt changes. It is the difference between asking politely and issuing a legally binding contract.

6.8 Quality Metrics

Measuring extraction quality requires a combination of automated and manual evaluation. Automated metrics catch systemic failures; manual evaluation catches subtle accuracy issues.

Track these metrics by document type, not just in aggregate. An overall 98% accuracy rate may hide the fact that invoices achieve 99.5% while contracts achieve only 92%—and contracts may be the higher-value use case.

Table 6.2: Quality Metrics for Unstructured Data Extraction

Metric	Definition	Target
Extraction rate	% of documents yielding output	>99%
Schema conformance	% of outputs matching schema	>99.5%
Field completeness	% of required fields populated	>95%
Value accuracy	% of extracted values correct	>97%
Latency (p50)	Median processing time	<5s
Cost per document	API cost per document	<\$0.05

Tip

Monitoring Extraction Quality Over Time Set up automated quality monitoring that compares each batch's metrics against historical baselines. A sudden drop in field completeness often signals a change in the source document format. Catch these regressions early by alerting when any metric deviates more than two standard deviations from its rolling average.

Warning

Ground Truth Is Essential Automated metrics tell you whether extraction *ran*, not whether it was *correct*. A pipeline that confidently extracts the wrong invoice total is worse than one that fails loudly. Manually annotate 50–100 documents of each type to create ground truth for accuracy measurement. Yes, manually. With your hands. We know. Refresh ground truth quarterly as document formats evolve.

6.9 Exercises

1. **PDF Parsing Benchmark:** Collect 20 PDF invoices from diverse sources (different vendors, formats, languages if possible). Build the parsing pipeline and measure extraction rate, field completeness, and accuracy against manually labeled ground truth. Document which invoice styles cause the most errors.
2. **Log Anomaly Injection:** Take a production-like log file (at least 10,000 lines). Inject three anomaly types: error rate spike, new error pattern, and gradual latency degradation. Run the anomaly detector and evaluate detection accuracy. Calculate the false positive and false negative rates.

3. **Tool Use vs. JSON Mode:** Implement both approaches for invoice parsing. Process 50 invoices with each and compare schema conformance, accuracy, failure rate, and cost per document. Document when each approach is preferable.
4. **Image Table Extraction:** Collect 10 images of tables from different sources (screenshots, photos of printed tables, scanned documents). Use CLAUDE CODE's vision to extract tabular data. Measure accuracy against hand-transcribed ground truth. Experiment with image preprocessing (contrast adjustment, deskewing) and measure improvement.
5. **Chunking Strategy Comparison:** Take a 100-page document and process it using all four chunking strategies from Table 6.1. Compare extraction completeness, accuracy, processing time, and cost. Determine which strategy works best for your document type.

Part 

Data Storage and Transformation

7

Data Modeling with Claude Code

A data model is a promise to your future self. Most people break that promise within six months. CLAUDE CODE is here to hold you accountable.

Data modeling is the foundation of every reliable data platform. A poorly designed schema leads to sluggish queries, brittle pipelines, and endless rework. This chapter explores how CLAUDE CODE accelerates every phase—from translating business requirements into DDL, to designing star schemas, to handling slowly changing dimensions and Data Vault 2.0.

i Note

Modeling Is a Conversation The best data models emerge from iterative dialogue between data engineers and business stakeholders. CLAUDE CODE can participate in this conversation: describe your business domain in natural language, and CLAUDE CODE produces a first-draft schema that serves as a concrete starting point for discussion. Unlike your stakeholders, CLAUDE CODE will not change the requirements halfway through the meeting.

7.1 Generating DDL from Natural Language

Tip

When generating DDL, always specify your target database engine (PostgreSQL, MySQL, Snowflake, BigQuery) in the prompt. Each has distinct data types, constraint syntax, and defaults. A BIGSERIAL in PostgreSQL becomes AUTOINCREMENT in Snowflake.

```
1 CREATE TABLE customers (  
2     customer_id    BIGSERIAL PRIMARY KEY,  
3     full_name      VARCHAR(200) NOT NULL,  
4     email          VARCHAR(320) NOT NULL UNIQUE,  
5     billing_city   VARCHAR(200),  
6     billing_country VARCHAR(100) DEFAULT 'US',  
7     created_at     TIMESTAMPTZ NOT NULL DEFAULT now(),  
8     updated_at     TIMESTAMPTZ NOT NULL DEFAULT now()  
9 );  
10  
11 CREATE TABLE orders (  
12     order_id       BIGSERIAL PRIMARY KEY,  
13     customer_id    BIGINT NOT NULL REFERENCES customers(customer_id),  
14     status          VARCHAR(20) NOT NULL DEFAULT 'pending'  
15     CHECK (status IN ('pending', 'shipped', 'delivered', '  
16     cancelled')),  
17     status_changed_at TIMESTAMPTZ NOT NULL DEFAULT now(),  
18     created_at     TIMESTAMPTZ NOT NULL DEFAULT now()  
19 );  
20  
21 CREATE TABLE order_line_items (  
22     line_item_id   BIGSERIAL PRIMARY KEY,  
23     order_id       BIGINT NOT NULL REFERENCES orders(order_id) ON DELETE  
24     CASCADE,  
25     product_id     BIGINT NOT NULL,  
26     quantity       INT NOT NULL CHECK (quantity > 0),  
27     unit_price      NUMERIC(12,2) NOT NULL CHECK (unit_price >= 0)  
28 );  
29  
30 CREATE INDEX idx_orders_customer ON orders(customer_id);  
31 CREATE INDEX idx_line_items_order ON order_line_items(order_id);
```

Listing 7.1: DDL generated from natural language (PostgreSQL)

Notice what CLAUDE CODE gets right unprompted: audit columns (`created_at/updated_at`), `TIMESTAMPTZ` (critical for multi-timezone), `CHECK` constraints, and indexes on foreign keys. These are patterns that experienced data engineers always include but junior engineers often forget. CLAUDE CODE is, in effect, a senior engineer who never has a bad day, never forgets the audit columns, and never names a column `data2`.

Warning

ALTER TABLE: Handle With Care When CLAUDE CODE generates `ALTER TABLE` for tables with data, review for data loss or locking. Adding a `NOT NULL` column without a default fails on non-empty PostgreSQL tables. On large tables, adding an index can lock the table for minutes—use `CREATE INDEX CONCURRENTLY` instead. Nothing says “career-defining moment” quite like locking a production table during business hours.

Iterative Refinement

The first DDL draft is rarely final. Use CLAUDE CODE iteratively: generate the initial schema, then ask follow-up questions. “What indexes should I add for a query that filters by `status` and sorts by `created_at`?” or “How should I partition this table if it grows to 500 million rows?” Each round produces a more production-ready schema.

```

1 def refine_schema(current_ddl: str, requirements: str) -> str:
2     """Ask Claude Code to refine an existing schema."""
3     client = anthropic.Anthropic()
4     response = client.messages.create(
5         model="claude-sonnet-4-20250514", max_tokens=4096,
6         messages=[{"role": "user", "content": f"""Here is my current
7 schema. Refine it based on these additional requirements.
8 Preserve existing columns and constraints unless explicitly
9 asked to change them. Explain each change.
10
11 Current DDL:
12 ```sql
13 {current_ddl}
14 ```
15
16 New requirements: {requirements}"""])
17     return response.content[0].text

```

Listing 7.2: Iterative Schema Refinement with Claude Code

Note

Version Your Schemas Treat schema definitions like code: store DDL in version control, tag releases, and maintain a changelog. When using CLAUDE CODE to iterate on a schema, save each version so you can diff between iterations and understand why each change was made.

7.2 Star and Snowflake Schema Design

In a star schema, a central fact table connects to denormalized dimension tables—optimizing for query simplicity and read performance. A snowflake schema normalizes dimensions further, trading query simplicity for storage efficiency and data integrity.

```
1 CREATE TABLE fact_orders (  
2     order_key          BIGSERIAL PRIMARY KEY,  
3     customer_key      BIGINT NOT NULL REFERENCES dim_customers,  
4     product_key       BIGINT NOT NULL REFERENCES dim_products,  
5     date_key          INT NOT NULL REFERENCES dim_date,  
6     quantity          INT NOT NULL,  
7     unit_price        NUMERIC(12,2) NOT NULL,  
8     net_amount        NUMERIC(12,2) NOT NULL,  
9     total_amount      NUMERIC(12,2) NOT NULL  
10 );  
11  
12 CREATE TABLE dim_customers (  
13     customer_key      BIGSERIAL PRIMARY KEY,  
14     customer_id      BIGINT NOT NULL,  
15     full_name         VARCHAR(200) NOT NULL,  
16     email            VARCHAR(320),  
17     segment          VARCHAR(50),  
18     region           VARCHAR(100),  
19     effective_from   DATE NOT NULL,  
20     effective_to     DATE,  
21     is_current       BOOLEAN NOT NULL DEFAULT TRUE  
22 );  
23  
24 CREATE TABLE dim_date (  
25
```

```

25     date_key          INT PRIMARY KEY, -- YYYYMMDD
26     full_date        DATE NOT NULL UNIQUE,
27     day_of_week      SMALLINT NOT NULL,
28     month_number     SMALLINT NOT NULL,
29     quarter          SMALLINT NOT NULL,
30     year             INT NOT NULL,
31     is_weekend       BOOLEAN NOT NULL,
32     is_holiday       BOOLEAN NOT NULL DEFAULT FALSE
33 );

```

Listing 7.3: Star schema fact and dimensions

Table 7.1: Star vs. snowflake schema comparison

Criterion	Star	Snowflake
Query simplicity	Fewer JOINS	More JOINS
Storage efficiency	Higher redundancy	Lower redundancy
Query performance	Generally faster	Slower for wide queries
BI tool compatibility	Excellent	Good, more setup

! Tip

Let the Query Patterns Decide When asking CLAUDE CODE to design a dimensional model, provide the top 10 queries your analysts will run. The query patterns determine whether a star or snowflake schema is appropriate, how to grain the fact table, and which dimensions to conformed across subject areas.

7.3 Data Vault 2.0

Data Vault separates structural information (hubs), relationships (links), and descriptive context (satellites). CLAUDE CODE dramatically reduces the boilerplate—a single business entity can require three or more tables in Data Vault, each with hash keys, load timestamps, and record source tracking.

```

1 CREATE TABLE hub_customer (
2     hub_customer_hk CHAR(32) PRIMARY KEY, -- MD5 of business key
3     customer_bk    VARCHAR(100) NOT NULL,
4     load_dts       TIMESTAMPTZ NOT NULL DEFAULT now(),
5     record_source  VARCHAR(100) NOT NULL

```

```
6 );
7
8 CREATE TABLE link_order_customer (
9     link_order_customer_hk CHAR(32) PRIMARY KEY,
10    hub_order_hk           CHAR(32) NOT NULL REFERENCES hub_order,
11    hub_customer_hk       CHAR(32) NOT NULL REFERENCES hub_customer,
12    load_dts              TIMESTAMPTZ NOT NULL DEFAULT now(),
13    record_source         VARCHAR(100) NOT NULL
14 );
15
16 CREATE TABLE sat_customer_details (
17    hub_customer_hk CHAR(32) NOT NULL REFERENCES hub_customer,
18    load_dts        TIMESTAMPTZ NOT NULL,
19    load_end_dts    TIMESTAMPTZ,
20    record_source   VARCHAR(100) NOT NULL,
21    hash_diff       CHAR(32) NOT NULL,
22    full_name       VARCHAR(200),
23    email           VARCHAR(320),
24    billing_city    VARCHAR(200),
25    PRIMARY KEY (hub_customer_hk, load_dts)
26 );
```

Listing 7.4: Data Vault hub, link, and satellite

Feed CLAUDE CODE a description of your source entities and it generates the full raw vault: hubs, links, and satellites with proper hash key computation. For a domain with ten entities and fifteen relationships, this can save days of mechanical DDL writing. Data Vault without CLAUDE CODE is just suffering with extra tables.

Tip

Hash Key Consistency Use a consistent hash function across all hubs and links. MD5 is standard in Data Vault 2.0, but some teams prefer SHA-256 for collision resistance. Whatever you choose, document it and enforce it. CLAUDE CODE can generate a reusable macro or function for hash computation to ensure consistency across all loading procedures.

Note

When to Use Data Vault Data Vault shines in environments with many source systems, frequent schema changes, and strict auditability requirements. If you have a single source and straightforward reporting needs, a star schema is simpler and more appropriate. CLAUDE CODE can help you evaluate the

trade-off by analyzing your source count, change frequency, and compliance requirements.

7.4 Normalization and Denormalization

CLAUDE CODE identifies normalization violations (transitive dependencies, repeating groups) and recommends where to denormalize for analytics. When you specify target query patterns, it tailors the denormalized design accordingly.

Tip

When asking CLAUDE CODE to denormalize, always specify target query patterns. A table for “revenue by region by quarter” looks different from one for “customer lifetime value by cohort.” Without query context, CLAUDE CODE will produce a generic denormalization that serves no pattern well.

A common pattern is to maintain a normalized operational model (3NF) for transactional workloads and a denormalized analytical model (star schema) for reporting. CLAUDE CODE can generate the transformation SQL that bridges the two, including proper handling of NULL values, type casting, and business logic.

Warning

Denormalization and Data Freshness Denormalized views and tables introduce a freshness trade-off. A view is always current but may be slow for complex joins. A materialized table is fast but stale until refreshed. Choose based on your query latency requirements and acceptable staleness. Document the refresh cadence for every denormalized table.

```
1 CREATE OR REPLACE VIEW analytics.wide_orders AS
2 SELECT
3     o.order_id,
4     o.created_at AS order_date,
5     c.full_name AS customer_name,
6     c.segment AS customer_segment,
7     p.product_name,
8     p.category AS product_category,
9     li.quantity,
10    li.unit_price,
```

```

11     li.quantity * li.unit_price AS line_total
12 FROM orders o
13 JOIN customers c ON o.customer_id = c.customer_id
14 JOIN order_line_items li ON o.order_id = li.order_id
15 JOIN products p ON li.product_id = p.product_id;

```

Listing 7.5: Denormalized view from normalized tables

7.5 Slowly Changing Dimensions

Table 7.2: SCD types

Type	Strategy	Description
0	Retain original	Never update
1	Overwrite	Replace; no history
2	Add new row	Full history with versioning
3	Add new column	Current and previous value
6	Hybrid (1+2+3)	Maximum flexibility

SCD Type 2 is the most common in practice because it preserves complete history. The implementation requires closing out the current row (setting `effective_to` and `is_current = FALSE`) and inserting a new row with the updated values.

```

1  -- Close out changed rows
2  MERGE INTO dim_customers AS tgt
3  USING (
4      SELECT customer_id, full_name, email, region, segment,
5             MD5(full_name || email || region || segment) AS hash_diff
6      FROM staging.stg_customers
7  ) AS src
8  ON tgt.customer_id = src.customer_id AND tgt.is_current = TRUE
9  WHEN MATCHED AND tgt.hash_diff != src.hash_diff THEN
10     UPDATE SET effective_to = CURRENT_DATE(), is_current = FALSE;
11
12  -- Insert new versions and new customers
13  INSERT INTO dim_customers (
14     customer_id, full_name, email, region, segment,
15     hash_diff, effective_from, effective_to, is_current
16  )
17  SELECT customer_id, full_name, email, region, segment,

```

```

18     MD5(full_name || email || region || segment),
19     CURRENT_DATE(), NULL, TRUE
20 FROM staging.stg_customers src
21 LEFT JOIN dim_customers tgt
22     ON src.customer_id = tgt.customer_id AND tgt.is_current = TRUE
23 WHERE tgt.customer_id IS NULL
24     OR tgt.hash_diff != MD5(src.full_name || src.email || src.region ||
    src.segment);

```

Listing 7.6: SCD Type 2 using Snowflake MERGE

! Tip

Include a `hash_diff` column on all tracked SCD Type 2 dimensions. Change detection becomes a simple hash comparison rather than column-by-column checks. Use MD5 for speed or SHA-256 if collision resistance matters.

! Warning

SCD Type 2 and Fact Table Joins When joining fact tables to SCD Type 2 dimensions, always filter on the date range: `WHERE fact.event_date BETWEEN dim.effective_from AND COALESCE(dim.effective_to, '9999-12-31')`. Joining only on `is_current = TRUE` gives you the current attribute values, not the values at the time of the event. This is the single most common SCD Type 2 bug, and it will make your numbers wrong in ways that are just subtle enough to ship to the CEO before anyone notices.

7.6 Schema Review with Claude Code

Paste your DDL into CLAUDE CODE and ask for a review. Common issues it catches: missing NOT NULL constraints, INT primary keys where BIGINT is needed, no indexes on foreign keys, inconsistent naming, missing temporal columns, and over-normalized OLAP schemas.

```

1 import anthropic, subprocess, sys
2
3 def review_schema_changes(diff: str) -> str:
4     client = anthropic.Anthropic()
5     message = client.messages.create(
6         model="claude-sonnet-4-20250514", max_tokens=4096,

```

```

7     messages=[{"role": "user",
8               "content": f""Review these migration changes. Flag issues.
9               ``diff
10              {diff}
11              ``
12              Format: ## Critical Issues ## Warnings ## Suggestions ## Summary""}]
13     return message.content[0].text
14
15 if __name__ == "__main__":
16     diff = subprocess.run(
17         ["git", "diff", "origin/main", "--", "migrations/"],
18         capture_output=True, text=True).stdout
19     if not diff:
20         print("No migration changes."); sys.exit(0)
21     review = review_schema_changes(diff)
22     print(review)
23     if "## Critical Issues" in review and \
24         "None" not in review.split("## Critical Issues")[1].split("##")
25         [0]:
26         sys.exit(1)

```

Listing 7.7: Automated Schema Review in CI/CD

Note

CLAUDE CODE can also compare two schema versions and generate a diff report—invaluable during code review for flagging breaking changes. Feed it the “before” and “after” DDL and ask it to classify each change as safe, risky, or breaking.

7.7 Schema Migration Generation

CLAUDE CODE generates versioned migration scripts with data preservation, backward compatibility, and rollback procedures. It supports Flyway, Liquibase, Alembic, Django, dbt, and Sqitch conventions.

```

1  -- Forward: V042__add_loyalty_tier.sql
2  BEGIN;
3  ALTER TABLE dim_customers
4     ADD COLUMN loyalty_tier VARCHAR(20) DEFAULT 'standard' NOT NULL;
5  UPDATE dim_customers SET loyalty_tier = CASE
6     WHEN lifetime_orders >= 100 THEN 'platinum'

```

```
7     WHEN lifetime_orders >= 50 THEN 'gold'
8     WHEN lifetime_orders >= 20 THEN 'silver'
9     ELSE 'standard' END
10 WHERE is_current = TRUE;
11 COMMIT;
12
13 -- Rollback: R042_remove_loyalty_tier.sql
14 BEGIN;
15 ALTER TABLE dim_customers DROP COLUMN loyalty_tier;
16 COMMIT;
```

Listing 7.8: Forward and rollback migration

Warning

Transactions or Tears Always wrap migrations in transactions when your database supports transactional DDL (PostgreSQL does; MySQL does not for most DDL). Without transactions, a failed migration leaves the schema in a partially-applied state that is difficult to recover from. It is Schrödinger’s migration: both applied and not applied until you look.

Tip

Migration Naming Conventions Adopt a consistent naming convention for migration files: V{version}__{description}.sql for forward migrations and R{version}__{description}.sql for rollbacks. CLAUDE CODE can enforce this convention when generating migration scripts if you include it in the prompt.

7.8 Exercises

1. **Natural Language to DDL:** Describe an e-commerce domain (customers, products, orders, reviews, inventory) to CLAUDE CODE in plain English. Generate DDL for PostgreSQL and Snowflake. Compare the two outputs and document the differences in data types, constraints, and indexing strategies.
2. **Star Schema Design:** Given a set of business questions (“What is revenue by product category and region, by quarter?”), design a star schema using CLAUDE CODE. Include at least one fact table, four dimension tables, and a date dimension. Write three analytical queries against the schema.

3. **SCD Type 2 Implementation:** Implement a complete SCD Type 2 pipeline: staging table, MERGE statement, and validation queries. Test with a sequence of changes: initial load, updates to existing records, new records, and records with no changes. Verify that history is preserved correctly.
4. **Schema Review Pipeline:** Set up an automated schema review using the CI/CD integration pattern from Section 7.6. Introduce five intentional issues (missing NOT NULL, INT instead of BIGINT, no index on FK, inconsistent naming, missing audit column) and verify that CLAUDE CODE catches all of them.
5. **Data Vault from Scratch:** Describe a business domain with at least five entities and eight relationships. Use CLAUDE CODE to generate the complete raw vault (hubs, links, satellites). Write the loading procedures for two hubs and one satellite.
6. **Normalization Audit:** Take an existing denormalized table from your data warehouse and ask CLAUDE CODE to identify normalization violations. Classify each violation by normal form (1NF, 2NF, 3NF). For each violation, evaluate whether it is intentional (for performance) or accidental (a modeling error).

This iterative workflow—CLAUDE CODE as drafting partner, reviewer, and sounding board—reduces modeling time from weeks to days while maintaining quality. CLAUDE CODE does not replace the data engineer’s judgment; it amplifies it by handling mechanical work and surfacing issues that might be overlooked under time pressure. You could also use CLAUDE CODE to design the seating chart for the team offsite. We are not saying you *should*. We are saying you *could*. Use CLAUDE CODE for Everything™.

8

SQL Generation and Optimization

Writing SQL by hand in 2026 is like churning your own butter.
Technically possible. Spiritually questionable.

SQL remains the lingua franca of data. Yet writing correct, performant SQL at scale is genuinely difficult. Analysts struggle with multi-table joins, engineers debate CTE structure versus subqueries with the fervor of religious scholars, and everyone dreads reading undocumented legacy queries left behind by engineers who have since entered witness protection. CLAUDE CODE can generate sophisticated SQL, translate between dialects, optimize slow queries by interpreting EXPLAIN plans, recommend indexes, produce migration scripts, and serve as the backbone of a natural language analytics interface.

8.1 The Text-to-SQL Challenge

Text-to-SQL converts a natural language question into a syntactically correct and semantically accurate SQL query. The task is hard for several reasons: natural language is ambiguous, real-world schemas are complex (thousands of tables with cryptic names), SQL dialects are fragmented across engines, and a generated query can be syntactically valid yet semantically *wrong*.

Warning

Dialect Confusion Will Ruin Your Day Never assume that SQL generated for one dialect will run on another without modification. Always specify the target engine in your prompt, and validate the output against that engine’s documentation or a test environment. “But it works on Postgres” is not a valid debugging strategy for Snowflake.

Note

Research benchmarks like Spider and BIRD measure text-to-SQL accuracy. CLAUDE CODE, when given full schema context and few-shot examples, consistently performs at the top of the 70–85% accuracy range and often surpasses it on enterprise schemas where domain context is provided.

8.2 Claude Code’s SQL Capabilities Across Dialects

CLAUDE CODE handles all major SQL dialects fluently. Here we demonstrate two representative examples.

PostgreSQL

```
1 SELECT
2     event_id,
3     event_type,
4     payload -> 'user' ->> 'email' AS user_email,
5     payload -> 'user' ->> 'plan' AS user_plan,
6     created_at
7 FROM events
8 WHERE payload -> 'user' ->> 'plan' = 'enterprise'
9     AND created_at >= now() - INTERVAL '30 days'
10 ORDER BY created_at DESC;
```

Listing 8.1: Claude Code-generated PostgreSQL JSONB query

Snowflake

```
1 SELECT
2     o.order_id,
3     o.customer_id,
4     o.order_date,
5     f.value:product_name::STRING AS product_name,
6     f.value:quantity::INT AS quantity,
7     f.value:unit_price::DECIMAL(12,2) AS unit_price
8 FROM orders o,
9     LATERAL FLATTEN(input => o.line_items) f
10 QUALIFY ROW_NUMBER() OVER (
11     PARTITION BY o.customer_id
12     ORDER BY o.order_date DESC
13 ) = 1
14 ORDER BY o.customer_id;
```

Listing 8.2: Claude Code-generated Snowflake query with QUALIFY and FLATTEN

Tip

When working with multiple dialects in the same organization, create a dialect specification document and include it in every prompt to CLAUDE CODE. This ensures consistent syntax choices and prevents accidental dialect mixing.

8.3 Building a Text-to-SQL Pipeline with Schema Context

A reliable text-to-SQL system is a pipeline with multiple stages: schema retrieval, prompt construction, query generation, validation, and execution.

1. **Schema Retrieval:** Identify the relevant tables and retrieve their DDL, column descriptions, and sample data.
2. **Prompt Construction:** Assemble schema context, dialect specification, business glossary, and the user's question.
3. **Query Generation:** Send the prompt to CLAUDE CODE and receive the generated SQL.
4. **Validation:** Parse the SQL, check for common errors, and optionally run EXPLAIN.
5. **Execution:** Run the query and return results.

```
1 import anthropic
2 import sqlparse
3 import json
4 from dataclasses import dataclass
5
6
7 @dataclass
8 class SchemaContext:
9     tables: list[dict]
10    relationships: list[dict]
11    glossary: dict[str, str]
12    sample_queries: list[dict]
13
14
15 @dataclass
16 class GeneratedQuery:
17     sql: str
18     explanation: str
19     confidence: float
20     dialect: str
21     tables_used: list[str]
22
23
24 class TextToSQLPipeline:
25     def __init__(
26         self,
27         client: anthropic.Anthropic,
28         schema_store: "SchemaStore",
29         dialect: str = "postgresql",
30         model: str = "claude-sonnet-4-20250514",
31     ):
32         self.client = client
33         self.schema_store = schema_store
34         self.dialect = dialect
35         self.model = model
36
37     def generate(
38         self, question: str, max_tables: int = 10
39     ) -> GeneratedQuery:
40         context = self.schema_store.retrieve(
41             question, max_tables=max_tables
42         )
43         prompt = self._build_prompt(question, context)
44         response = self.client.messages.create(
45             model=self.model,
```

```

46         max_tokens=4096,
47         system=self._system_prompt(),
48         messages=[{"role": "user", "content": prompt}],
49     )
50     result = self._parse_response(response)
51     self._validate_sql(result.sql)
52     return result
53
54     def _system_prompt(self) -> str:
55         return f"""You are an expert SQL developer
56 specializing in {self.dialect}. Rules:
57 1. Use only tables and columns from the provided schema.
58 2. Use {self.dialect} syntax exclusively.
59 3. Include comments explaining complex logic.
60 4. Prefer CTEs over nested subqueries.
61 5. Always include appropriate WHERE clauses.
62 6. Return JSON: sql, explanation, confidence (0-1), tables_used."""
63
64     def _build_prompt(self, question, context):
65         sections = ["## Schema"]
66         for table in context.tables:
67             sections.append(f"### {table['name']}")
68             sections.append(f"```sql\n{table['ddl']}\n```")
69         if context.relationships:
70             sections.append("## Relationships")
71             for rel in context.relationships:
72                 sections.append(
73                     f"- {rel['from_table']}.{rel['from_col']}"
74                     f" -> {rel['to_table']}.{rel['to_col']}"
75                 )
76         if context.glossary:
77             sections.append("## Business Glossary")
78             for term, defn in context.glossary.items():
79                 sections.append(f"- **{term}**: {defn}")
80         if context.sample_queries:
81             sections.append("## Example Queries")
82             for ex in context.sample_queries:
83                 sections.append(
84                     f"Q: {ex['question']}\n```sql\n{ex['sql']}\n```"
85                 )
86         sections.append(f"## Question\n{question}")
87         return "\n\n".join(sections)
88
89     def _parse_response(self, response):
90         data = json.loads(response.content[0].text)
91         return GeneratedQuery(

```

```

92         sql=data["sql"],
93         explanation=data["explanation"],
94         confidence=data["confidence"],
95         dialect=self.dialect,
96         tables_used=data["tables_used"],
97     )
98
99     def _validate_sql(self, sql: str) -> None:
100         parsed = sqlparse.parse(sql)
101         if not parsed:
102             raise ValueError("Generated SQL is empty")
103         if parsed[0].get_type() != "SELECT":
104             raise ValueError("Only SELECT queries are allowed.")
105         sql_upper = sql.upper()
106         if "DELETE" in sql_upper or "DROP" in sql_upper:
107             raise ValueError("Dangerous keywords detected")

```

Listing 8.3: Text-to-SQL pipeline orchestrator

Tip

Few-shot examples are the single most impactful way to improve text-to-SQL accuracy. Maintain a curated library of 50–100 question/SQL pairs covering your schema’s most common query patterns.

8.4 Complex Query Generation

CLAUDE CODE shines when generating complex queries involving CTEs, window functions, recursive queries, and advanced aggregations.

CTEs and Growth Rate Ranking

```

1 WITH monthly_category_revenue AS (
2     SELECT
3         DATE_TRUNC('month', o.created_at)::DATE AS revenue_month,
4         p.category,
5         SUM(oi.quantity * oi.unit_price) AS monthly_revenue
6 FROM orders o
7 JOIN order_items oi ON o.order_id = oi.order_id
8 JOIN products p ON oi.product_id = p.product_id

```

```

9     WHERE o.status = 'completed'
10     AND o.created_at >= DATE_TRUNC(
11         'month', CURRENT_DATE - INTERVAL '13 months')
12     GROUP BY 1, 2
13 ),
14 with_growth AS (
15     SELECT *, CASE
16         WHEN LAG(monthly_revenue) OVER (
17             PARTITION BY category ORDER BY revenue_month
18         ) > 0 THEN ROUND(
19             (monthly_revenue - LAG(monthly_revenue) OVER (
20                 PARTITION BY category ORDER BY revenue_month
21             )) / LAG(monthly_revenue) OVER (
22                 PARTITION BY category ORDER BY revenue_month
23             ) * 100, 2)
24         ELSE NULL
25     END AS growth_rate_pct
26     FROM monthly_category_revenue
27 ),
28 ranked AS (
29     SELECT *, RANK() OVER (
30         PARTITION BY revenue_month
31         ORDER BY growth_rate_pct DESC NULLS LAST
32     ) AS growth_rank
33     FROM with_growth
34     WHERE growth_rate_pct IS NOT NULL
35 )
36 SELECT * FROM ranked WHERE growth_rank <= 3
37 ORDER BY revenue_month DESC, growth_rank;

```

Listing 8.4: Multi-CTE query with growth rate ranking

Recursive Queries

```

1 WITH RECURSIVE org_tree AS (
2     SELECT employee_id, full_name, title, manager_id,
3         1 AS depth, full_name AS path
4     FROM employees WHERE manager_id IS NULL
5     UNION ALL
6     SELECT e.employee_id, e.full_name, e.title, e.manager_id,
7         ot.depth + 1, ot.path || ' > ' || e.full_name
8     FROM employees e
9     JOIN org_tree ot ON e.manager_id = ot.employee_id

```

```

10 WHERE ot.depth < 20
11 )
12 SELECT * FROM org_tree ORDER BY path;

```

Listing 8.5: Recursive CTE for organizational hierarchy

Warning

Infinite Loops: Not Just a Philosophy Problem Recursive CTEs can cause infinite loops if the data contains cycles. Always include a depth limit in the `WHERE` clause of the recursive member. An infinite loop in a recursive CTE is how you turn a \$0.50 query into a \$5,000 query and simultaneously discover your warehouse’s auto-suspend is not configured.

8.5 Query Optimization with Claude Code

CLAUDE CODE can interpret `EXPLAIN ANALYZE` output and suggest concrete optimizations. The workflow is: run `EXPLAIN ANALYZE`, send the plan to CLAUDE CODE along with the query and DDL, and ask for specific recommendations.

Table 8.1: Common `EXPLAIN` plan bottlenecks and Claude Code-suggested fixes

Bottleneck	Plan Indicator	Typical Fix
Sequential scan on large table	Seq Scan with high row count	Add index on filter/join columns
Nested loop on large tables	Nested Loop with high row estimates	Rewrite to enable hash join
Sort spilling to disk	Sort Method: external merge	Increase <code>work_mem</code> ; add index
Inaccurate row estimates	Estimated vs. actual differ by 10x+	Run <code>ANALYZE</code>

```

1 -- Create a composite partial index
2 CREATE INDEX CONCURRENTLY idx_orders_date_status
3 ON orders (created_at, status)
4 WHERE status = 'completed';
5

```

```

6 -- Pre-aggregate before joining to customers
7 WITH order_revenue AS (
8     SELECT o.order_id, o.customer_id,
9           DATE_TRUNC('month', o.created_at) AS month,
10          SUM(oi.quantity * oi.unit_price) AS order_total
11     FROM orders o
12     JOIN order_items oi ON o.order_id = oi.order_id
13     WHERE o.created_at >= '2025-01-01'
14           AND o.created_at < '2026-01-01'
15           AND o.status = 'completed'
16     GROUP BY o.order_id, o.customer_id,
17              DATE_TRUNC('month', o.created_at)
18 )
19 SELECT c.segment, orv.month,
20        COUNT(*) AS orders, SUM(orv.order_total) AS revenue
21 FROM order_revenue orv
22 JOIN customers c ON orv.customer_id = c.customer_id
23 GROUP BY c.segment, orv.month
24 ORDER BY orv.month, revenue DESC;
25 -- Expected improvement: 45s -> ~3s (15x faster, 0 human tears)

```

Listing 8.6: Optimized query with pre-aggregation

Note

Always provide table row counts and the full EXPLAIN ANALYZE output to CLAUDE CODE, not just EXPLAIN. Actual row counts and timing reveal problems that estimated-only plans hide.

8.6 Migration Script Generation

CLAUDE CODE can generate safe, idempotent migration scripts that avoid downtime.

```

1 BEGIN;
2 -- Step 1: Add column as nullable (non-blocking)
3 ALTER TABLE customers
4     ADD COLUMN IF NOT EXISTS customer_tier VARCHAR(20);
5
6 -- Step 2: Add check constraint (validated separately)
7 ALTER TABLE customers
8     ADD CONSTRAINT chk_customer_tier

```

```

9     CHECK (customer_tier IN (
10         'bronze', 'silver', 'gold', 'platinum'
11     )) NOT VALID;
12
13 -- Step 3: Validate constraint in background (non-blocking)
14 ALTER TABLE customers
15     VALIDATE CONSTRAINT chk_customer_tier;
16
17 -- Step 4: Add default for future rows
18 ALTER TABLE customers
19     ALTER COLUMN customer_tier SET DEFAULT 'bronze';
20 COMMIT;
21
22 -- Step 5: Create index CONCURRENTLY (outside transaction)
23 -- CREATE INDEX CONCURRENTLY idx_customers_tier
24 --     ON customers(customer_tier);

```

Listing 8.7: Safe migration script generated by Claude Code

Warning

PostgreSQL Gotchas `CREATE INDEX CONCURRENTLY` cannot run inside a transaction block. CLAUDE CODE correctly separates this step—because it has read the PostgreSQL documentation more thoroughly than any human ever has or ever will. Always verify migration scripts for this and similar PostgreSQL-specific constraints.

8.7 Full Example: Natural Language Analytics Interface

Here is a complete natural language analytics interface using FastAPI, PostgreSQL, and CLAUDE CODE.

```

1 from fastapi import FastAPI, HTTPException
2 from pydantic import BaseModel
3 import anthropic
4 import psycopg2
5 import psycopg2.extras
6 import json
7 from datetime import datetime
8
9 app = FastAPI(title="NL Analytics API")

```

```
10 client = anthropic.Anthropic()
11
12 SCHEMA_CONTEXT = """
13 ### sales.orders
14 - order_id (BIGINT PK), customer_id (BIGINT FK),
15   order_date (DATE), status (VARCHAR), total_amount (NUMERIC),
16   region (VARCHAR): NA, EU, APAC, LATAM
17
18 ### sales.customers
19 - customer_id (BIGINT PK), company_name (VARCHAR),
20   segment (VARCHAR): enterprise, mid_market, smb
21
22 ## Business Glossary
23 - Revenue: SUM(quantity * unit_price * (1 - discount_pct/100))
24 - Active customer: 1+ completed order in last 12 months
25 """
26
27 SYSTEM_PROMPT = """You are a SQL analyst. Generate PostgreSQL
28 queries from natural language. Rules:
29 1. Use ONLY tables listed in the schema.
30 2. Prefer CTEs for complex queries.
31 3. Include LIMIT 1000 unless asked for all rows.
32 4. Never generate DDL, DML, or DCL statements.
33 5. Return JSON: {"sql": "...", "title": "...", "description": "...}"""
34
35
36 class QuestionRequest(BaseModel):
37     question: str
38     max_rows: int = 1000
39
40
41 @app.post("/ask")
42 def ask_question(req: QuestionRequest):
43     response = client.messages.create(
44         model="claude-sonnet-4-20250514",
45         max_tokens=2048,
46         system=SYSTEM_PROMPT,
47         messages=[{
48             "role": "user",
49             "content": f"{SCHEMA_CONTEXT}\n\nQuestion: {req.question}",
50         }],
51     )
52     result = json.loads(response.content[0].text)
53     sql = result["sql"]
54
55     # Safety check
```

```

56 sql_upper = sql.upper().strip()
57 if not sql_upper.startswith(("WITH", "SELECT")):
58     raise HTTPException(400, "Only SELECT queries allowed")
59
60 # Execute with timeout
61 conn = psycopg2.connect(
62     "postgresql://analytics:pass@localhost/warehouse"
63 )
64 try:
65     with conn.cursor(
66         cursor_factory=psycopg2.extras.RealDictCursor
67     ) as cur:
68         cur.execute("SET statement_timeout = '30s';")
69         cur.execute(sql)
70         rows = cur.fetchmany(req.max_rows)
71     conn.commit()
72     return {
73         "question": req.question,
74         "sql": sql,
75         "title": result["title"],
76         "columns": list(rows[0].keys()) if rows else [],
77         "rows": rows,
78         "row_count": len(rows),
79     }
80 finally:
81     conn.close()

```

Listing 8.8: Natural language analytics API

Tip

Timeouts: Your Last Line of Defense Always set a `statement_timeout` when executing AI-generated SQL in production. A 30-second timeout is a reasonable starting point for analytics queries. Without a timeout, one accidental Cartesian product could consume your entire warehouse budget before you finish your coffee.

8.8 Index Recommendation

CLAUDE CODE can analyze a workload and recommend an optimal index set. Provide the top N slowest queries (from `pg_stat_statements`), table DDL, and existing indexes.

Table 8.2: PostgreSQL index types and use cases

Index Type	Best For	Example
B-tree	Equality and range queries	<code>CREATE INDEX ... ON t(col)</code>
Partial	Fixed-predicate queries	<code>... WHERE status = 'active'</code>
GIN	Full-text, JSONB, arrays	<code>... USING gin(payload)</code>
BRIN	Large, naturally ordered tables	<code>... USING brin(created_at)</code>

! Tip

BRIN: The Most Underrated Index BRIN indexes are spectacularly efficient for append-only tables ordered by timestamp. A BRIN index on 500 million rows might consume only 1 MB compared to 10 GB for a B-tree. If BRIN indexes were a stock, we would tell you to buy. Ask CLAUDE CODE if your table qualifies—it is better at identifying BRIN candidates than most humans.

8.9 Anti-Patterns in AI-Generated SQL

AI-generated SQL is remarkably good but not infallible. Like a star employee who occasionally microwaves fish in the office kitchen, it has its failure modes. The five most common are:

1. **Hallucinated tables:** Referencing tables or columns that do not exist. Always validate against the schema catalog.
2. **Silent Cartesian products:** A missing join condition produces every-row-times-every-row results.
3. **Incorrect aggregation granularity:** Joining before counting inflates numbers. Use `COUNT(DISTINCT ...)`.
4. **Timezone mishandling:** Comparing `TIMESTAMPTZ` to naive date strings uses the session timezone silently.
5. **NULL semantics:** `NOT IN` with `NULLs` in the subquery returns zero rows. Use `NOT EXISTS` instead.

Table 8.3: Validation checklist for AI-generated SQL

Check	What to Verify	Automated?
Schema validity	All tables and columns exist	Yes
Statement type	Only SELECT allowed	Yes
Join completeness	Every join has an ON clause	Yes
Row limit	LIMIT clause present	Yes
NULL safety	NOT IN replaced with NOT EXISTS	Partially
Performance	EXPLAIN cost under threshold	Yes

8.10 Exercises

1. **Dialect Translation:** Take a PostgreSQL query using `FILTER`, `ARRAY_AGG`, and `INTERVAL` and use `CLAUDE CODE` to translate it to BigQuery, Snowflake, and Databricks SQL. Document the differences.
2. **Schema Context Optimization:** Design a two-stage retrieval system for a 2,000-table database that first uses `CLAUDE CODE` to identify likely relevant tables, then retrieves full DDL only for those. Implement both stages in Python.
3. **EXPLAIN Plan Analysis:** Obtain an `EXPLAIN ANALYZE` output for a slow query. Send it to `CLAUDE CODE` with table DDL and ask for optimization recommendations. Implement at least two and measure the improvement.
4. **Anti-Pattern Detector:** Build a Python function that checks AI-generated SQL for the five anti-patterns described in this chapter. Test it against 20 queries generated by `CLAUDE CODE`.

9

dbt and Transformation Pipelines

dbt without CLAUDE CODE is just YAML with extra steps. dbt *with* CLAUDE CODE is a lifestyle.

The *dbt* (data build tool) framework has fundamentally changed how data teams build transformation pipelines. By treating SQL transformations as software—with version control, testing, documentation, and modular design—dbt brought engineering discipline to the analytics layer. Yet dbt projects are labor-intensive to build and maintain. Writing models, tests, documentation, macros, and source YAML by hand is tedious. Nobody became a data engineer because they dreamed of writing YAML. This is exactly where CLAUDE CODE excels. In this chapter we explore how CLAUDE CODE accelerates every phase of dbt development: generating models from requirements, writing tests, producing documentation, creating macros, and migrating legacy SQL.

9.1 dbt Fundamentals Review

- **Models:** SQL `SELECT` statements stored as `.sql` files, compiled to `CREATE TABLE AS` or `CREATE VIEW AS`.
- **Sources:** Declarations of raw tables in YAML.
- **Tests:** Schema tests (`unique`, `not_null`, `accepted_values`, `relationships`) in YAML; data tests as standalone SQL returning zero rows.
- **Macros:** Reusable Jinja functions that generate SQL.

- **Snapshots:** SCD Type 2 tracking of historical changes.

A well-organized dbt project follows a layered architecture: staging (views), intermediate (ephemeral), and marts (tables).

Note

Convention Over Configuration Over Sanity The layered architecture (staging, intermediate, marts) is a convention popularized by the dbt community. CLAUDE CODE is well-versed in this convention and follows it when generating models. It has read more dbt projects than any human on Earth, including the ones with 47 layers of indirection and a staging model named `stg_tmp_final_v3_DONOTUSE`.

9.2 Using Claude Code to Generate dbt Models

The key is providing CLAUDE CODE with sufficient context: the source schema, the dbt project conventions, and the business logic.

Staging Models

```
1 with source as (  
2     select * from {{ source('stripe', 'payments') }}  
3 ),  
4 renamed as (  
5     select  
6         id as payment_id,  
7         customer_id as stripe_customer_id,  
8         amount / 100.0 as payment_amount,  
9         lower(currency) as currency_code,  
10        lower(status) as payment_status,  
11        created as payment_created_at,  
12        _sdc_received_at as loaded_at,  
13        is_test  
14    from source  
15 ),  
16 filtered as (  
17     select * from renamed  
18     where not is_test  
19 ),
```

```
20 final as (  
21   select  
22     {{ dbt_utils.generate_surrogate_key(  
23       ['payment_id']  
24     ) }} as payment_sk,  
25     *  
26   from filtered  
27 )  
28 select * from final
```

Listing 9.1: Claude Code-generated staging model: stg_stripe__payments.sql

Mart Models

```
1  {{  
2    config(  
3      materialized='incremental',  
4      unique_key=['customer_sk', 'revenue_month'],  
5      incremental_strategy='merge',  
6      cluster_by=['revenue_month']  
7    )  
8  }}  
9  
10 with payments as (  
11   select * from {{ ref('int_payments_pivoted') }}  
12 ),  
13 customers as (  
14   select * from {{ ref('stg_stripe__customers') }}  
15 ),  
16 final as (  
17   select  
18     c.customer_sk,  
19     p.payment_month as revenue_month,  
20     c.customer_name,  
21     c.customer_segment,  
22     p.gross_revenue,  
23     p.successful_payments as payment_count,  
24     p.avg_payment_amount,  
25     case  
26       when p.failed_payments > 0  
27       then round(p.failed_payments::float  
28         / nullif(p.total_payments, 0) * 100, 2)  
29       else 0
```

```

30     end as failure_rate_pct,
31     current_timestamp() as dbt_updated_at
32 from payments p
33 left join customers c
34     on p.stripe_customer_id = c.stripe_customer_id
35 {% if is_incremental() %}
36 where p.payment_month >= (
37     select max(revenue_month) from {{ this }}
38 ) - interval '3 months'
39 {% endif %}
40 )
41 select * from final

```

Listing 9.2: Claude Code-generated mart model: fct_revenue.sql

Tip

Late-Arriving Data Is Not a Bug, It's a Feature When generating incremental models, always include a lookback window in the `is_incremental()` block to handle late-arriving data. Data is never on time. If your data were a person, it would miss every flight.

9.3 Generating dbt Tests

CLAUDE CODE can generate comprehensive test suites that would take hours to write manually. Hours that you would rather spend doing literally anything else, including attending meetings.

```

1 version: 2
2 models:
3   - name: fct_revenue
4     description: >
5       Monthly revenue fact table at the customer-month grain.
6     columns:
7       - name: customer_sk
8         tests: [not_null]
9       - name: revenue_month
10        tests: [not_null]
11       - name: gross_revenue
12         tests:
13           - not_null
14           - dbt_utils.accepted_range:

```

```

15         min_value: 0
16         inclusive: true
17     - name: failure_rate_pct
18       tests:
19         - dbt_utils.accepted_range:
20             min_value: 0
21             max_value: 100
22     - name: customer_segment
23       tests:
24         - accepted_values:
25             values: [enterprise, mid_market, smb, unknown]
26   tests:
27     - dbt_utils.unique_combination_of_columns:
28         combination_of_columns:
29             - customer_sk
30             - revenue_month

```

Listing 9.3: Claude Code-generated schema tests for the finance mart

Custom Generic Tests

```

1  {% test recency(model, column_name,
2      datepart='day', interval=1) %}
3  with max_date as (
4      select max({{ column_name }}) as max_val
5      from {{ model }}
6  )
7  select max_val from max_date
8  where max_val < dateadd(
9      '{{ datepart }}', -{{ interval }}, current_timestamp()
10 )
11 {% endtest %}

```

Listing 9.4: Claude Code-generated custom generic test: recency check

Warning

The Silent Test Graveyard Custom generic tests must be placed in the `tests/generic/` directory and named `test_<test_name>.sql`. Mismatches cause silent failures where the test is never discovered. Your test is not “pass-

ing.” Your test does not *exist*. This is not a Zen koan; it is a bug.

9.4 Claude Code for dbt Documentation

CLAUDE CODE can generate comprehensive documentation from model SQL, business context, and column names.

```

1 import anthropic
2 import yaml
3 from pathlib import Path
4
5
6 def generate_model_docs(
7     model_path: str, client: anthropic.Anthropic,
8 ) -> dict:
9     with open(model_path) as f:
10         model_sql = f.read()
11         model_name = Path(model_path).stem
12         prompt = f"""Analyze this dbt model and generate YAML documentation.
13 Model name: {model_name}
14 SQL:
15 ```sql
16 {model_sql}
17 ```
18 Generate: model description (2-3 sentences), column descriptions, tests.
19 Return ONLY valid YAML."""
20
21     response = client.messages.create(
22         model="claude-sonnet-4-20250514",
23         max_tokens=4096,
24         messages=[{"role": "user", "content": prompt}],
25     )
26     return yaml.safe_load(response.content[0].text)

```

Listing 9.5: Script to generate dbt documentation with Claude Code

9.5 Macro Generation with Claude Code

dbt macros are Jinja functions that generate SQL. Writing macros requires fluency in both SQL and Jinja—a combination CLAUDE CODE handles well.

```

1 {% macro safe_divide(numerator, denominator, default=0) %}
2     case
3         when {{ denominator }} is null then {{ default }}
4         when {{ denominator }} = 0 then {{ default }}
5         else {{ numerator }}::float / {{ denominator }}
6     end
7 {% endmacro %}

```

Listing 9.6: Claude Code-generated macro: safe division

```

1 {% macro date_spine(start_date, end_date, datepart='day') %}
2     {%- if target.type == 'snowflake' -%}
3     select dateadd('{{ datepart }}', seq4(),
4         '{{ start_date }}'::date) as date_{{ datepart }}
5     from table(generator(rowcount => 10000))
6     where date_{{ datepart }} <= '{{ end_date }}'::date
7     {%- elif target.type == 'bigquery' -%}
8     select date_{{ datepart }}
9     from unnest(generate_date_array(
10        date '{{ start_date }}', date '{{ end_date }}',
11        interval 1 {{ datepart }}
12    )) as date_{{ datepart }}
13     {%- elif target.type == 'postgres' -%}
14     select date_{{ datepart }}::date
15     from generate_series(
16        '{{ start_date }}'::date, '{{ end_date }}'::date,
17        '1 {{ datepart }}'::interval
18    ) as date_{{ datepart }}
19     {%- endif -%}
20 {% endmacro %}

```

Listing 9.7: Claude Code-generated macro: cross-database date spine

! Tip

Jinja: Where Whitespace Goes to Cause Problems Use CLAUDE CODE to generate macros, but always review the Jinja carefully. Whitespace handling ({{- vs {{}}) and trailing commas need human review. Jinja templating is the dark art of data engineering—half SQL, half template language, half existential crisis. Yes, that is three halves. Jinja does not respect the laws of mathematics either.

9.6 Migrating Legacy SQL to dbt Models

One of the most valuable applications is migrating legacy stored procedures into well-structured dbt models. The conceptual shift is from imperative (cursor loops, temp tables, quiet weeping) to declarative (a single SELECT). CLAUDE CODE handles this translation with the calm professionalism of someone who did not write the original stored procedure and therefore has no emotional attachment to it.

```

1  -- Migrated from: sp_daily_revenue_report
2  -- Cursor loop converted to set-based operation
3  -- DELETE+INSERT converted to incremental merge
4  {{
5      config(
6          materialized='incremental',
7          unique_key=['report_date', 'region', 'product_category'],
8          incremental_strategy='merge'
9      )
10 }}
11
12 with orders as (
13     select * from {{ ref('stg_app_db__orders') }}
14     where order_status = 'completed'
15 ),
16 order_items as (
17     select * from {{ ref('int_order_items_enriched') }}
18 ),
19 daily_revenue as (
20     select
21         o.ordered_at::date as report_date,
22         c.customer_country as region,
23         oi.product_category,
24         count(distinct o.order_id) as order_count,
25         sum(oi.net_amount) as total_revenue,
26         {{ safe_divide(
27             'sum(oi.net_amount)',
28             'count(distinct o.order_id)'
29         ) }} as avg_order_value
30     from orders o
31     join order_items oi on o.order_id = oi.order_id
32     left join {{ ref('dim_customers') }} c
33         on o.customer_id = c.customer_id
34     {% if is_incremental() %}
35     where o.ordered_at::date >= (
36         select max(report_date) - interval '7 days'

```

```
37     from {{ this }}
38   )
39   {% endif %}
40   group by 1, 2, 3
41 )
42 select *, current_timestamp() as dbt_updated_at
43 from daily_revenue
```

Listing 9.8: dbt model after migration: fct_daily_revenue.sql

Note

From Imperative to Declarative When migrating stored procedures to dbt, CLAUDE CODE excels at the imperative-to-declarative translation because it understands both paradigms. It is fluent in both “how to do it” and “what to do,” which is more than can be said for most meeting attendees.

9.7 CI/CD for dbt with Claude Code-Powered Reviews

Adding CLAUDE CODE-powered code review to a CI pipeline catches issues that automated tests miss.

```
1 name: dbt CI
2 on:
3   pull_request:
4     paths: ['models/**', 'macros/**', 'tests/**']
5 jobs:
6   dbt-ci:
7     runs-on: ubuntu-latest
8     steps:
9       - uses: actions/checkout@v4
10      - name: Install dependencies
11        run: pip install dbt-snowflake anthropic
12      - name: dbt build (modified models only)
13        run: |
14          dbt build --select state:modified+ \
15            --defer --state ./prod_manifest --target ci
16      - name: Claude Code Review
17        if: always()
18        run: python scripts/claude_dbt_review.py
19      env:
20        ANTHROPIC_API_KEY: ${ secrets.ANTHROPIC_API_KEY }
```

Listing 9.9: GitHub Actions workflow for dbt CI

! Tip

CI Speed: Because Life Is Short In CI, use `dbt build --select state:modified+` with the `--defer` flag to only build models that changed. This dramatically reduces CI run times for large projects. Nobody wants to wait 45 minutes for CI to tell them they forgot a comma. Use CLAUDE CODE to generate the CI workflow too, naturally. Use CLAUDE CODE for Everything™.

9.8 Exercises

1. **Staging Model Generation:** Given a raw Salesforce opportunity table with columns `id`, `Name`, `AccountId`, `Amount`, `StageName`, `CloseDate`, `IsDeleted`, `IsWon`, `Type`, use CLAUDE CODE to generate a complete staging model with proper renaming, type casting, and filtering.
2. **Test Suite Design:** For a `fct_orders` fact table, design a comprehensive test suite with at least 15 assertions across schema tests, custom generic tests, and singular data tests.
3. **Macro Library:** Build five reusable dbt macros: `cents_to_dollars`, `fiscal_quarter`, `slugify`, `pivot_column`, and `scd_type2`.
4. **Legacy Migration:** Take a stored procedure or complex view and use CLAUDE CODE to decompose it into staging, intermediate, and mart models, with a comparison test verifying output equivalence.

10

Data Quality and Testing

“Data quality is not a project; it is a discipline.”

—Every data team’s Confluence page, right above the section titled “Known Issues (Do Not Touch)”

Bad data is expensive. A single incorrect join key, unnoticed schema drift, or silently failing pipeline can cascade through dashboards, ML models, and business decisions. Before CLAUDE CODE, data quality was largely aspirational—something you put on the roadmap right after “fix tech debt” and right before “heat death of the universe.” This chapter shows how to use CLAUDE CODE to build a proactive data quality practice: catching problems before they reach consumers, generating tests automatically, and negotiating data contracts between teams. You could do all of this manually, of course. You could also perform surgery on yourself. Both are technically possible.

10.1 Data Quality Dimensions

! Tip

When starting a data quality initiative, pick the two dimensions that cause the most downstream pain (usually completeness and consistency) and instrument those first. Or, if you enjoy suffering, try to fix all six simultaneously. Your standup updates will be riveting.

Table 10.1: The six dimensions of data quality.

Dimension	Definition	Example Check
<i>Accuracy</i>	Data correctly represents real-world entities.	Does <code>order_total</code> match the sum of line items?
<i>Completeness</i>	All expected data is present.	What percentage of records have a null email?
<i>Consistency</i>	Same data agrees across multiple places.	Does every <code>customer_id</code> in facts exist in the dimension?
<i>Timeliness</i>	Data is available when needed.	Is today's partition populated by 8 AM?
<i>Uniqueness</i>	No unintended duplicates.	Are there duplicate <code>order_id</code> values?
<i>Validity</i>	Data conforms to expected format and range.	Is <code>country_code</code> always a valid ISO 3166-1 alpha-2 code?

10.2 Traditional vs. Claude Code-Augmented Data Quality

Traditional tools (Great Expectations, Soda, dbt tests) are rule-based and excellent at enforcing *known* constraints. They struggle with: (1) unknown unknowns like misspelled status values, (2) semantic anomalies where data passes type checks but a surge suggests duplication, and (3) cross-table reasoning where upstream changes silently invalidate downstream aggregations. In other words, they catch the problems you already know about. How reassuring.

CLAUDE CODE complements rule-based tools by adding a reasoning layer: it *generates* tests, *reviews* them, and fills gaps that static rules cannot cover. Think of it as the difference between a smoke detector and a firefighter who also installs smoke detectors, writes the building code, and judges you for not having done this sooner.

i Note

CLAUDE CODE does not replace Great Expectations or dbt tests. Think of it as the senior engineer who writes the tests, while the framework executes them. Unlike the actual senior engineer, however, it does not take a three-week vacation right before the audit.

10.3 Claude Code for Anomaly Detection

Statistical anomaly detection catches numerical outliers but struggles with contextual anomalies—values that are statistically normal but semantically suspicious. The hybrid approach: (1) run statistical checks to produce candidates, (2) send candidates with context to CLAUDE CODE for triage, (3) CLAUDE CODE classifies each as “real anomaly,” “expected variation,” or “needs investigation.” Finally, a human glances at the results and says “looks fine” without reading them, completing the circle of life.

```
1 import anthropic
2 import json
3
4 def triage_anomalies(
5     anomalies: list[dict],
6     table_context: str,
7     client: anthropic.Anthropic,
8 ) -> list[dict]:
9     prompt = f"""You are a data quality analyst. Classify each anomaly
10    as:
11    - REAL_ANOMALY: Likely a genuine data quality issue.
12    - EXPECTED_VARIATION: Normal business variation (e.g., holiday spike).
13    - NEEDS_INVESTIGATION: Ambiguous; recommend follow-up queries.
14
15    Table context: {table_context}
16    Anomalies: {json.dumps(anomalies, indent=2)}
17
18    Respond with JSON array: id, classification, explanation,
19    suggested_query. """
20
21     response = client.messages.create(
22         model="claude-sonnet-4-20250514",
23         max_tokens=4096,
24         temperature=0.0,
25         messages=[{"role": "user", "content": prompt}],
26     )
27     return json.loads(response.content[0].text)
```

Listing 10.1: Hybrid anomaly triage with Claude Code

Tip

Feed CLAUDE CODE as much business context as possible: known events (holidays, launches), recent changes (new data sources, schema changes), and historical baselines. Richer context yields more accurate triage. Without context, CLAUDE CODE is just a very expensive coin flip.

10.4 Generating Great Expectations Suites

CLAUDE CODE can generate a comprehensive Great Expectations suite from just the DDL and a few sample rows. After the first validation run, feed failures back to CLAUDE CODE for iterative refinement. It is the software equivalent of “measure once, cut once, then have someone smarter re-measure and re-cut for you.”

```

1 def generate_gx_suite(
2     table_name: str, ddl: str, sample_rows: list[dict],
3     business_context: str = "", model: str = "claude-sonnet-4-20250514",
4 ) -> dict:
5     client = anthropic.Anthropic()
6     prompt = f"""Generate a comprehensive Great Expectations suite for:
7 Table: {table_name}
8 DDL: ```sql\n{ddl}\n```
9 Sample rows: ```json\n{json.dumps(sample_rows[:5], indent=2, default=str
10 )}\n```
11 Business context: {business_context}
12 Include expectations for EVERY column. Output valid JSON."""
13     response = client.messages.create(
14         model=model, max_tokens=8192, temperature=0.0,
15         messages=[{"role": "user", "content": prompt}],
16     )
17     return json.loads(response.content[0].text)

```

Listing 10.2: Generating a GX suite with Claude Code

Warning

Always review CLAUDE CODE-generated expectations before production use. Pay special attention to range expectations—CLAUDE CODE infers ranges from sample data, which may not cover edge cases. It has never seen your data on

Black Friday. Nobody is ever truly prepared for Black Friday.

10.5 Data Contracts

A *data contract* is a formal agreement between a data producer and consumer about the schema, semantics, quality guarantees, and SLAs of a dataset. In theory, it prevents the “I changed a column name and now seventeen dashboards are broken” conversation. In practice, it at least gives you something to point to during the postmortem.

```
1  apiVersion: v1
2  kind: DataContract
3  metadata:
4    name: fact-orders-contract
5    version: 2.1.0
6    owner: order-platform-team
7    consumers: [analytics-team, ml-team, finance-team]
8  schema:
9    properties:
10     order_id: {type: integer, primaryKey: true, nullable: false}
11     order_total: {type: number, nullable: false, minimum: 0.01}
12     status: {type: string, nullable: false,
13             enum: [PENDING, PROCESSING, SHIPPED, DELIVERED, CANCELLED]}
14  quality:
15    completeness: {threshold: 0.99, columns: [order_id, order_total,
16      status]}
17    uniqueness: {columns: [order_id]}
18    freshness: {maxStaleness: PT6H}
19  sla:
20    availability: 99.5%
21    updateFrequency: PT1H
```

Listing 10.3: Data contract structure

CLAUDE CODE can generate contracts from DDL and sample data, and can facilitate negotiation between producer and consumer teams by analyzing both sides’ requirements and proposing compromises. It is, in effect, a couples therapist for your data teams—except both sides actually listen.

10.6 Building a Data Quality Monitoring System

A monitoring system continuously evaluates data quality and alerts when metrics breach thresholds. The architecture: Scheduler → DQ Checker (SQL + Python) → CLAUDE CODE (Triage) → Alerting (Slack/PagerDuty). You could also just wait for someone in Finance to send an angry email, but the formal approach scales better.

```

1 from dataclasses import dataclass
2 from datetime import datetime
3
4 @dataclass
5 class DQCheckResult:
6     table_name: str
7     check_name: str
8     dimension: str
9     score: float
10    threshold: float
11    passed: bool
12    details: str
13    checked_at: datetime
14
15 class DQChecker:
16     def __init__(self, warehouse_conn_str, metrics_conn_str):
17         self.wh_conn = psycopg2.connect(warehouse_conn_str)
18         self.metrics_conn = psycopg2.connect(metrics_conn_str)
19
20     def check_completeness(self, table, columns, threshold=0.99):
21         results = []
22         for col in columns:
23             with self.wh_conn.cursor() as cur:
24                 cur.execute(f"""
25                     SELECT COUNT({col}>::FLOAT / NULLIF(COUNT(*), 0)
26                     FROM {table}""")
27                 score = cur.fetchone()[0] or 0.0
28                 results.append(DQCheckResult(
29                     table, f"completeness_{col}", "completeness",
30                     score, threshold, score >= threshold,
31                     f"Non-null rate: {score:.4f}", datetime.utcnow(),
32                 ))
33         return results
34
35     def check_freshness(self, table, timestamp_col, max_hours=6):
36         with self.wh_conn.cursor() as cur:
37             cur.execute(f"""

```

```

38         SELECT EXTRACT(EPOCH FROM (NOW() - MAX({timestamp_col}))
39         ) / 3600.0
40         FROM {table}"""
41         hours_stale = cur.fetchone()[0] or 999
42         score = max(0.0, 1.0 - (hours_stale / max_hours))
43         return DQCheckResult(
44             table, "freshness", "timeliness", score, 0.5,
45             hours_stale <= max_hours,
46             f"Data is {hours_stale:.1f}h old (max: {max_hours}h)",
47             datetime.utcnow(),
48         )

```

Listing 10.4: DQ Checker core methods

Instead of alerting on every threshold breach, use CLAUDE CODE to consolidate and prioritize failures into a single intelligent alert with severity, likely root cause, and recommended actions. Your on-call engineer's 3 AM self will thank you. Or at least resent you slightly less.

10.7 Claude Code for Root Cause Analysis

When a data quality issue is detected, CLAUDE CODE with tool access can perform autonomous root cause analysis by systematically querying the warehouse, checking pipeline logs, and reviewing schema change history. It is like having a detective on retainer, except this one never sleeps, never eats, and has no opinions about your naming conventions. (That last part is a lie.)

```

1  def run_rca_agent(issue_description: str, client: anthropic.Anthropic)
2      -> str:
3      tools = [
4          {"name": "run_query", "description": "Run a read-only SQL query.",
5           "input_schema": {"type": "object",
6            "properties": {"query": {"type": "string"}},
7            "required": ["query"]}},
8          {"name": "get_pipeline_logs",
9           "description": "Get recent Airflow task logs.",
10          "input_schema": {"type": "object",
11           "properties": {"dag_id": {"type": "string"}},
12          "required": ["dag_id"]}},
13      ]
14      messages = [{"role": "user", "content": issue_description}]

```

```

14     for _ in range(10):
15         response = client.messages.create(
16             model="claude-sonnet-4-20250514", max_tokens=4096,
17             tools=tools, messages=messages,
18         )
19         if response.stop_reason != "tool_use":
20             return "".join(b.text for b in response.content if b.type ==
21                 "text")
22             messages.append({"role": "assistant", "content": response.
23                 content})
24             tool_results = []
25             for block in response.content:
26                 if block.type == "tool_use":
27                     result = execute_rca_tool(block.name, block.input)
28                     tool_results.append({"type": "tool_result",
29                         "tool_use_id": block.id, "content": result})
30             messages.append({"role": "user", "content": tool_results})
31     return "RCA agent reached maximum iterations."

```

Listing 10.5: Root cause analysis agent (abbreviated)

! Tip

RCA agents work best with access to multiple data sources: the warehouse, pipeline logs, schema change histories, and deployment logs. Each additional source narrows the hypothesis space. Ideally, give it access to everything short of your Slack DMs. Actually, those might help too.

10.8 Automated Test Generation

This system scans a dbt project, generates tests for every model using CLAUDE CODE, and writes them to schema YAML files. The dream: complete test coverage without a single human writing a single test. The reality: complete test coverage with a human reviewing the tests CLAUDE CODE wrote, which is still dramatically better than the previous state of “we’ll add tests next sprint” for forty-seven consecutive sprints.

```

1 import anthropic
2 import yaml
3 import asyncio
4 from pathlib import Path

```

```

5 from dataclasses import dataclass
6
7 @dataclass
8 class DbtModel:
9     name: str
10    sql_path: Path
11    sql_content: str
12
13 class DbtTestGenerator:
14     def __init__(self, project_dir: Path):
15         self.project_dir = project_dir
16         self.client = anthropic.AsyncAnthropic()
17
18     def discover_models(self) -> list[DbtModel]:
19         models = []
20         for sql_file in (self.project_dir / "models").rglob("*.sql"):
21             models.append(DbtModel(
22                 name=sql_file.stem, sql_path=sql_file,
23                 sql_content=sql_file.read_text()))
24         return models
25
26     async def generate_tests(self, model: DbtModel) -> dict:
27         prompt = f"""Generate dbt tests for this model.
28 Model: {model.name}
29 SQL:
30 ```sql
31 {model.sql_content}
32 ```
33 Generate column-level tests (not_null, unique, accepted_values,
34 relationships) and descriptions. Output YAML only."""
35
36         response = await self.client.messages.create(
37             model="claude-sonnet-4-20250514", max_tokens=4096,
38             temperature=0.0, system="Output only valid YAML.",
39             messages=[{"role": "user", "content": prompt}])
40         return yaml.safe_load(response.content[0].text.strip())
41
42     async def generate_all(self, models):
43         sem = asyncio.Semaphore(5)
44         async def limited(m):
45             async with sem:
46                 return m.name, await self.generate_tests(m)
47         return dict(await asyncio.gather(*[limited(m) for m in models]))

```

Listing 10.6: Automated dbt test generator

Note

After running the generator, always review the output. Common issues: overly strict `accepted_values` lists, incorrect relationship tests, and redundant tests on nullable columns. `CLAUDE CODE` is thorough to a fault—it will test columns you forgot existed, which is either helpful or existentially unsettling.

10.9 DQ Metrics and KPIs

Table 10.2: KPIs for a data quality program. Achieving all five simultaneously is sometimes referred to as “the impossible dream” or “next quarter’s OKR.”

KPI	Definition	Target
DQ Score (composite)	Weighted average across dimensions	≥ 0.95
Test Coverage	% of columns with ≥ 1 test	≥ 0.90
MTTR	Average hours from detection to fix	< 4 hours
Contract Coverage	% of critical tables with contracts	$= 1.00$
False Positive Rate	% of alerts that are not real issues	< 0.10

10.10 Exercises

- DQ Dimension Assessment.** Choose a dataset and write SQL queries to measure all six dimensions. Calculate a composite DQ score. Which dimension scores lowest? (Spoiler: it is always timeliness. Always.)
- Anomaly Detection Pipeline.** Build a hybrid pipeline: (a) use z-scores to identify candidates, (b) send to `CLAUDE CODE` for contextual triage. Evaluate classification quality on 10+ candidates. Bonus: use `CLAUDE CODE` to also write your grocery list. It is surprisingly good at this.
- Data Contract Drafting.** Draft a data contract for a dataset you produce or consume. Use `CLAUDE CODE` to create the initial draft, then refine with a colleague on the other side of the producer/consumer relationship. If the

negotiation goes poorly, CLAUDE CODE can also draft the passive-aggressive Slack messages.

4. **Root Cause Analysis.** Simulate a DQ issue by introducing nulls or duplicates. Run your DQ checks, then use the RCA agent to investigate. Did CLAUDE CODE identify the correct root cause? If so, congratulations—you now have a colleague who is better at debugging than most of your team.

Part **IV**

Data Serving and Analytics

11

Building Analytical Queries with Claude Code

“Write SQL,” they said. “It’ll be fun,” they said. Nobody mentioned the part where you spend four hours debugging a window function at 11 PM because the VP needs “just a quick analysis” by morning.

Modern data teams face ever-growing demand for analytical queries—from ad hoc Slack requests to production dashboards refreshing every fifteen minutes. CLAUDE CODE can fundamentally reshape this workflow by accelerating every phase of query authoring, review, and deployment. Sure, you *could* write all this SQL yourself. You could also churn your own butter. Both activities build character and waste time in roughly equal measure. In this chapter we build production-grade patterns for dashboard query generation, metrics layers, self-service analytics, cohort analysis, A/B testing, and query optimization.

11.1 Dashboard Query Generation

Dashboard queries must be deterministic, performant under concurrent execution, and return results that map cleanly to visualization primitives. The most important input to CLAUDE CODE is an accurate, curated schema description. Feed it garbage, and it will produce beautifully formatted garbage—which is arguably worse than ugly garbage because people trust it.

```
WITH daily_revenue AS (
```

```

2  SELECT
3      DATE(oi.ordered_at)                AS order_date,
4      COUNT(DISTINCT oi.order_id)        AS total_orders,
5      COUNT(DISTINCT oi.customer_id)     AS unique_customers,
6      SUM(oi.net_revenue)                 AS gross_revenue,
7      SUM(oi.net_revenue) - SUM(oi.shipping_cost) AS net_revenue,
8      SUM(oi.net_revenue)
9      / NULLIF(COUNT(DISTINCT oi.order_id), 0) AS avg_order_value
10 FROM analytics.fct_order_items oi
11 WHERE oi.ordered_at >= DATEADD('day', -90, CURRENT_DATE)
12       AND oi.fulfillment_status <> 'returned'
13 GROUP BY 1
14 ),
15 with_rolling AS (
16     SELECT dr.*,
17         AVG(dr.gross_revenue) OVER (
18             ORDER BY dr.order_date ROWS BETWEEN 6 PRECEDING AND CURRENT
19             ROW
20         ) AS revenue_7d_moving_avg,
21         LAG(dr.gross_revenue, 7) OVER (ORDER BY dr.order_date)
22         AS gross_revenue_prev_week
23     FROM daily_revenue dr
24 )
25 SELECT * FROM with_rolling ORDER BY order_date DESC;

```

Listing 11.1: Claude Code-generated daily revenue query

! Tip

Always instruct CLAUDE CODE to include NULLIF guards on denominators and to exclude cancelled or returned rows explicitly. Dashboard queries that surface division-by-zero errors erode trust faster than any other bug. Nothing says “we have our data together” quite like a dashboard displaying Infinity%.

11.2 Metrics Layers and Semantic Layers

A metrics layer provides a single source of truth for business metric definitions. Without one, you get the classic data team experience: three dashboards, three definitions of “revenue,” three VPs arguing about whose number is right. CLAUDE CODE excels at generating dbt metric definitions and Cube.js schema files from plain-English metric descriptions.

```
1 semantic_models:
2   - name: order_items
3     defaults:
4       agg_time_dimension: ordered_at
5       model: ref('fct_order_items')
6     entities:
7       - name: order_item
8         type: primary
9         expr: order_item_id
10    measures:
11      - name: total_net_revenue
12        agg: sum
13        expr: net_revenue
14        description: Sum of net revenue excluding returned items.
15        create_metric: true
16      - name: total_orders
17        agg: count_distinct
18        expr: order_id
19        create_metric: true
20
21 metrics:
22   - name: average_order_value
23     type: derived
24     label: "Average Order Value (AOV)"
25     description: "Net revenue divided by distinct order count."
26     type_params:
27       expr: total_net_revenue / NULLIF(total_orders, 0)
28     metrics:
29       - name: total_net_revenue
30       - name: total_orders
```

Listing 11.2: dbt metric definition generated by Claude Code

Note

Always include a **description** field in semantic layer definitions. These descriptions become documentation that downstream consumers—both human and AI—use to understand each metric. Future You will thank Present You. Future You thanks Present You so rarely that you should seize every opportunity.

11.3 Self-Service Analytics

CLAUDE CODE enables natural-language analytics interfaces where business users ask questions in plain English and receive validated SQL and results. This is the dream that every BI vendor has sold since 2005, except now it actually works. The architecture has four components: context retrieval (schema + metrics from a vector store), query generation, validation (parse + allowlist check via `sqlglot`), and execution with presentation.

```
1 import anthropic
2 import sqlglot
3
4 ALLOWED_SCHEMAS = {"analytics", "metrics", "staging"}
5
6 class AnalyticsQueryEngine:
7     def __init__(self, client: anthropic.Anthropic, schema_context: str):
8
9         self.client = client
10        self.schema_context = schema_context
11
12    def generate_sql(self, question: str) -> str:
13        response = self.client.messages.create(
14            model="claude-sonnet-4-20250514", max_tokens=4096,
15            system="You are an analytics SQL generator. Snowflake
16            dialect. "
17            "Only reference tables in the provided schema.",
18            messages=[{"role": "user",
19                    "content": f"Schema:\n{self.schema_context}\n\n"
20                    f"Question: {question}"}],
21        )
22        return response.content[0].text.strip()
23
24    def validate_sql(self, sql: str) -> tuple[bool, str]:
25        try:
26            parsed = sqlglot.parse(sql, dialect="snowflake")
27        except Exception as e:
28            return False, f"Parse error: {e}"
29
30        for stmt in parsed:
31            for table in stmt.find_all(sqlglot.exp.Table):
32                if (table.db or "").lower() not in ALLOWED_SCHEMAS:
33                    return False, f"Disallowed schema: {table.db}"
34            if not isinstance(stmt, sqlglot.exp.Select):
35                return False, "Only SELECT statements allowed."
36        return True, "Validation passed."
```

Listing 11.3: Self-service analytics query engine

! Tip

Always validate CLAUDE CODE-generated SQL before execution. The `sqlglot` library provides dialect-aware parsing and catches table references outside your allowlist without hitting the warehouse. Letting unvalidated AI-generated SQL touch your production warehouse is the data engineering equivalent of texting your ex at 2 AM—technically possible, never advisable.

11.4 Building a Metrics Store

A metrics store materializes pre-computed metric values for fast retrieval—essential for sub-second lookups in executive dashboards and real-time alerts. Because nothing matters more than ensuring the CEO’s dashboard loads in under one second, even if the underlying data is three hours stale.

```
1 from dataclasses import dataclass, field
2 import hashlib
3
4 @dataclass
5 class MetricDefinition:
6     name: str
7     sql_expression: str
8     source_table: str
9     time_column: str
10    granularities: list[str]
11    dimensions: list[str] = field(default_factory=list)
12    filters: list[str] = field(default_factory=list)
13
14    @property
15    def metric_id(self) -> str:
16        raw = f"{self.name}:{self.source_table}:{self.sql_expression}"
17        return hashlib.md5(raw.encode()).hexdigest()[:12]
18
19 class MetricsStore:
20     def __init__(self, engine, target_schema="metrics"):
21         self.engine = engine
22         self.target_schema = target_schema
23
```

```

24 def generate_materialization_sql(self, metric, granularity, lookback
    =90):
25     dim_cols = ", ".join(metric.dimensions) if metric.dimensions
    else ""
26     where = " AND ".join(
27         [f"{metric.time_column} >= DATEADD('day', -{lookback},
    CURRENT_DATE)"]
28         + metric.filters
29     )
30     return f"""
31     INSERT INTO {self.target_schema}.metric_values
32         (metric_id, metric_name, granularity, period, metric_value)
33     SELECT '{metric.metric_id}', '{metric.name}', '{granularity}',
34         DATE_TRUNC('{granularity}', {metric.time_column}),
35         {metric.sql_expression}
36     FROM {metric.source_table}
37     WHERE {where}
38     GROUP BY DATE_TRUNC('{granularity}', {metric.time_column})
39         {', ' + dim_cols if dim_cols else ''}
40     ON CONFLICT (metric_id, granularity, period)
41     DO UPDATE SET metric_value = EXCLUDED.metric_value;"""

```

Listing 11.4: Metrics store computation engine

Warning

Metrics stores can grow rapidly when combining multiple granularities with high-cardinality dimensions. Use CLAUDE CODE to audit dimension cardinality before adding new slices. “Let’s add a per-user daily breakdown” sounds innocent until your metrics table has more rows than the source data.

11.5 Time Series Analysis Queries

CLAUDE CODE generates sophisticated temporal SQL patterns. Year-over-year comparisons, rolling averages, anomaly detection—all the queries that make you question whether SQL was really designed for this, or whether we are collectively engaged in an elaborate act of denial.

```

1 WITH monthly_revenue AS (
2     SELECT DATE_TRUNC('month', ordered_at) AS revenue_month,
3         EXTRACT(YEAR FROM ordered_at) AS yr,

```

```

4      EXTRACT(MONTH FROM ordered_at) AS mo,
5      SUM(net_revenue) AS total_revenue
6  FROM analytics.fct_order_items
7  WHERE fulfillment_status <> 'returned'
8      AND ordered_at >= DATEADD('year', -3, CURRENT_DATE)
9  GROUP BY 1, 2, 3
10 )
11 SELECT curr.revenue_month,
12        curr.total_revenue AS current_revenue,
13        prev.total_revenue AS prior_year_revenue,
14        ROUND((curr.total_revenue - prev.total_revenue)
15              / NULLIF(prev.total_revenue, 0) * 100, 2) AS yoy_growth_pct
16 FROM monthly_revenue curr
17 LEFT JOIN monthly_revenue prev
18   ON curr.mo = prev.mo AND curr.yr = prev.yr + 1
19 WHERE curr.yr = EXTRACT(YEAR FROM CURRENT_DATE)
20 ORDER BY curr.revenue_month;

```

Listing 11.5: Year-over-year comparison with growth rates

```

1 WITH daily_revenue AS (
2   SELECT DATE(ordered_at) AS order_date, SUM(net_revenue) AS revenue
3   FROM analytics.fct_order_items
4   WHERE fulfillment_status <> 'returned'
5   AND ordered_at >= DATEADD('day', -180, CURRENT_DATE)
6   GROUP BY 1
7 ),
8 with_stats AS (
9   SELECT order_date, revenue,
10          AVG(revenue) OVER (ORDER BY order_date
11                             ROWS BETWEEN 30 PRECEDING AND 1 PRECEDING) AS rolling_mean,
12          STDDEV(revenue) OVER (ORDER BY order_date
13                                ROWS BETWEEN 30 PRECEDING AND 1 PRECEDING) AS rolling_stddev
14   FROM daily_revenue
15 )
16 SELECT order_date, revenue, rolling_mean,
17        ROUND((revenue - rolling_mean) / NULLIF(rolling_stddev, 0), 2) AS
18        z_score,
19        CASE
20          WHEN ABS((revenue - rolling_mean) / NULLIF(rolling_stddev, 0)) >
21            3 THEN 'CRITICAL'
22          WHEN ABS((revenue - rolling_mean) / NULLIF(rolling_stddev, 0)) >
23            2 THEN 'WARNING'
24          ELSE 'NORMAL'

```

```
22 END AS anomaly_status
23 FROM with_stats WHERE rolling_stddev IS NOT NULL
24 ORDER BY order_date DESC;
```

Listing 11.6: Z-score anomaly detection in SQL

! Tip

When generating anomaly detection queries, specify the baseline window size (e.g., 30 days) and sensitivity threshold (e.g., 2 or 3 standard deviations). These depend on your data's natural volatility. Set the threshold too low and you will get paged for every Tuesday. Set it too high and you will learn about problems from Twitter.

11.6 Cohort Analysis and Funnel Queries

Cohort and funnel queries are among the most requested and error-prone analytical patterns. They are also the queries most likely to be described as “simple” in a Slack message by someone who has never written SQL.

```
1 WITH first_purchase AS (
2     SELECT customer_id,
3         DATE_TRUNC('month', MIN(ordered_at)) AS cohort_month
4     FROM analytics.fct_order_items
5     WHERE fulfillment_status <> 'returned'
6     GROUP BY 1
7 ),
8 monthly_activity AS (
9     SELECT DISTINCT customer_id,
10        DATE_TRUNC('month', ordered_at) AS activity_month
11    FROM analytics.fct_order_items
12    WHERE fulfillment_status <> 'returned'
13 ),
14 cohort_data AS (
15     SELECT fp.cohort_month, ma.activity_month,
16        DATEDIFF('month', fp.cohort_month, ma.activity_month)
17        AS months_since_first,
18        COUNT(DISTINCT ma.customer_id) AS active_customers
19    FROM first_purchase fp
20    INNER JOIN monthly_activity ma ON fp.customer_id = ma.customer_id
21    GROUP BY 1, 2, 3
```

```

22 ),
23 cohort_sizes AS (
24     SELECT cohort_month, COUNT(DISTINCT customer_id) AS cohort_size
25     FROM first_purchase GROUP BY 1
26 )
27 SELECT cd.cohort_month, cs.cohort_size, cd.months_since_first,
28        cd.active_customers,
29        ROUND(cd.active_customers::FLOAT / cs.cohort_size * 100, 2) AS
        retention_pct
30 FROM cohort_data cd
31 INNER JOIN cohort_sizes cs ON cd.cohort_month = cs.cohort_month
32 WHERE cd.months_since_first <= 12
33 ORDER BY cd.cohort_month, cd.months_since_first;

```

Listing 11.7: Monthly retention cohort analysis

Warning

Funnel queries using `COUNT(DISTINCT session_id)` per step do *not* enforce ordering—a session counted in step 3 may never have hit step 2. For strict sequential funnels, use window functions with `LEAD/LAG` or `MATCH_RECOGNIZE`. Or just let CLAUDE CODE handle it, because life is short and funnel queries are long.

11.7 A/B Test Analysis

CLAUDE CODE can generate warehouse-native significance calculations. Always verify three things: (1) the experiment window is correctly bounded, (2) users are attributed to exactly one variant, and (3) the metric window is consistent across control and treatment. Getting any of these wrong will produce statistically significant results that are also statistically meaningless—the data equivalent of a beautifully typeset lie.

```

1 WITH variant_stats AS (
2     SELECT variant, COUNT(*) AS n, SUM(converted) AS successes,
3           AVG(converted) AS p_hat
4     FROM user_outcomes GROUP BY 1
5 ),
6 test_inputs AS (
7     SELECT
8         MAX(CASE WHEN variant='variant_a' THEN p_hat END) AS p_control,

```

```

9      MAX(CASE WHEN variant='variant_a' THEN n END) AS n_control,
10     MAX(CASE WHEN variant='variant_b' THEN p_hat END) AS p_treatment,
11
12     MAX(CASE WHEN variant='variant_b' THEN n END) AS n_treatment
13 FROM variant_stats
14 ),
15 z_test AS (
16     SELECT *,
17         (p_treatment - p_control) / SQRT(
18             ((p_control*n_control + p_treatment*n_treatment)
19              / (n_control + n_treatment))
20             * (1 - ((p_control*n_control + p_treatment*n_treatment)
21                  / (n_control + n_treatment)))
22             * (1.0/n_control + 1.0/n_treatment)
23         ) AS z_statistic
24 FROM test_inputs
25 )
26 SELECT
27     ROUND(p_control*100, 4) AS control_rate_pct,
28     ROUND(p_treatment*100, 4) AS treatment_rate_pct,
29     ROUND(z_statistic, 4) AS z_stat,
30     CASE
31         WHEN ABS(z_statistic) > 2.576 THEN 'p < 0.01'
32         WHEN ABS(z_statistic) > 1.960 THEN 'p < 0.05'
33         ELSE 'Not significant'
34     END AS significance_level
35 FROM z_test;

```

Listing 11.8: Z-test for conversion rate difference in SQL

11.8 Query Optimization with Claude Code

CLAUDE CODE can analyze execution plans and suggest optimizations. Common recommendations include: early filtering for partition pruning, predicate alignment with clustering keys, removing `SELECT *` in CTEs, replacing correlated subqueries with `JOIN`s, and pre-aggregating common computations. In other words, all the things you know you should be doing but have not gotten around to because “it only takes 45 seconds and that’s fine.” It is not fine. Your warehouse bill knows it is not fine.

```

1  -- BEFORE: Full scan, 45 seconds, 12GB scanned
2  SELECT customer_id, SUM(net_revenue) AS total_rev

```

```
3 FROM analytics.fct_order_items
4 WHERE customer_id IN (
5     SELECT customer_id FROM analytics.dim_customers
6     WHERE customer_tier = 'platinum')
7 GROUP BY 1;
8
9 -- AFTER (Claude Code-optimized): 3 seconds, 800MB
10 SELECT oi.customer_id, SUM(oi.net_revenue) AS total_rev
11 FROM analytics.fct_order_items oi
12 INNER JOIN analytics.dim_customers c
13     ON oi.customer_id = c.customer_id
14 WHERE c.customer_tier = 'platinum'
15     AND oi.ordered_at >= DATEADD('year', -2, CURRENT_DATE)
16 GROUP BY 1;
```

Listing 11.9: Before and after optimization

 **Warning**

Always validate CLAUDE CODE's optimization suggestions by comparing EXPLAIN plans and benchmarking on realistic data volumes. Optimizations that help on small datasets can sometimes hurt on large ones. Trust, but verify. Actually, just verify.

12

Reverse ETL and Operational Analytics

“We spent years getting data into the warehouse. Now we spend years getting it back out. The circle of enterprise life is complete.”

Data warehouses were built to be the destination. But modern data teams have discovered that the warehouse is also the best *source* for operational systems. **Reverse ETL**—pushing curated warehouse data back into SaaS tools, CRMs, and marketing platforms—closes the loop between analytical insight and operational action. It is, essentially, admitting that the dashboard nobody checks should have been a Salesforce field all along. CLAUDE CODE accelerates every stage: designing sync logic, generating audience segments, building activation pipelines, and monitoring freshness.

12.1 What Is Reverse ETL

Traditional ETL moves data from operational systems into a warehouse. Reverse ETL bridges the gap by syncing curated data back into tools where operational teams work—pushing customer health scores into Salesforce contact records instead of asking salespeople to check dashboards. Because salespeople will never check dashboards. This is a universal constant, like gravity or deployments breaking on Fridays.

Dimension	Traditional ETL	Reverse ETL
Direction	Sources → Warehouse	Warehouse → Destinations
Volume	High (bulk ingestion)	Lower (curated subsets)
Schema control	Warehouse owns schema	Destination API owns schema
Primary users	Data engineers	Data engineers + RevOps + Marketing Ops

! Tip

When explaining reverse ETL to stakeholders: “We are going to keep Salesforce automatically updated with the latest customer metrics—no manual exports, no stale spreadsheets.” When explaining it to your therapist: “I am now responsible for data flowing in both directions and I can feel my hair thinning.”

12.2 Operational Analytics

Reverse ETL is one pillar of *operational analytics*—the idea that analytics should drive automated actions, not just dashboards that get screenshot’d into PowerPoints. It takes outputs of analytical and ML workflows and feeds them into business processes across sales (lead scores in CRM), marketing (audience segments in ad platforms), support (health scores in ticketing), and finance (revenue data in ERP).

12.3 Using Claude Code to Design Reverse ETL Pipelines

Generating the Source Model

```

1  {{
2    config(
3      materialized='incremental',
4      unique_key='customer_id',
5      cluster_by=['customer_id']
6    )
7  }}
8
9  WITH customers AS (
10     SELECT contact_id AS customer_id, email, first_name,
11            last_name, company_name, owner_id AS sales_owner_id
12     FROM {{ ref('stg_crm_contacts') }}
13 ),

```

```

14 subscriptions AS (
15     SELECT customer_id, plan_name AS subscription_tier, mrr,
16           subscription_status,
17           CASE
18             WHEN subscription_status = 'active' THEN 'paying'
19             WHEN subscription_status = 'past_due' THEN 'at_risk'
20             WHEN subscription_status = 'canceled' THEN 'churned'
21             ELSE 'unknown'
22           END AS lifecycle_stage
23     FROM {{ ref('stg_billing__subscriptions') }}
24     QUALIFY ROW_NUMBER() OVER (
25       PARTITION BY customer_id ORDER BY current_period_start DESC
26     ) = 1
27 ),
28 usage AS (
29     SELECT customer_id,
30           COUNT(DISTINCT CASE WHEN activity_date >= DATEADD('day', -30,
31             CURRENT_DATE)
32             THEN activity_date END) AS active_days_last_30,
33           MAX(activity_date) AS last_activity_date
34     FROM {{ ref('stg_events__daily_usage_summary') }}
35     GROUP BY 1
36 ),
37 tickets AS (
38     SELECT customer_id,
39           COUNT(*) AS total_tickets,
40           COUNT(CASE WHEN status = 'open' THEN 1 END) AS open_tickets
41     FROM {{ ref('stg_zendesk__tickets') }}
42     GROUP BY 1
43 ),
44 health_score AS (
45     SELECT c.customer_id,
46           LEAST(100, GREATEST(0,
47             (COALESCE(u.active_days_last_30, 0) / 30.0) * 40
48             + CASE s.lifecycle_stage WHEN 'paying' THEN 25 WHEN 'at_risk'
49             THEN 5 ELSE 0 END
50             + CASE WHEN COALESCE(t.open_tickets, 0) = 0 THEN 20 WHEN t.
51             open_tickets <= 2 THEN 10 ELSE 0 END
52             + CASE WHEN u.last_activity_date >= DATEADD('day', -3,
53             CURRENT_DATE) THEN 15
54             WHEN u.last_activity_date >= DATEADD('day', -14,
55             CURRENT_DATE) THEN 8 ELSE 0 END
56           ))::INT AS health_score
57     FROM customers c
58     LEFT JOIN subscriptions s USING (customer_id)
59     LEFT JOIN usage u USING (customer_id)

```

```

55     LEFT JOIN tickets t USING (customer_id)
56 )
57 SELECT c.*, s.subscription_tier, s.mrr, s.lifecycle_stage,
58        u.active_days_last_30, u.last_activity_date,
59        t.open_tickets, h.health_score,
60        CURRENT_TIMESTAMP() AS synced_at
61 FROM customers c
62 LEFT JOIN subscriptions s USING (customer_id)
63 LEFT JOIN usage u USING (customer_id)
64 LEFT JOIN tickets t USING (customer_id)
65 LEFT JOIN health_score h USING (customer_id)
66 {% if is_incremental() %}
67 WHERE c.customer_id IN (
68     SELECT DISTINCT customer_id FROM {{ ref('
69     stg_events__daily_usage_summary' ) }}
70     WHERE activity_date >= DATEADD('day', -2, CURRENT_DATE)
71     UNION
72     SELECT DISTINCT customer_id FROM {{ ref('stg_zendesk__tickets' ) }}
73     WHERE updated_at >= (SELECT MAX(synced_at) FROM {{ this }})
74 )
75 {% endif %}

```

Listing 12.1: dbt model for reverse ETL: reverse_etl_contacts.sql

Tip

Always include a `synced_at` timestamp column in reverse ETL source models. This enables incremental sync logic and freshness verification. It also gives you something to point at during the inevitable “when was this data last updated?” Slack thread.

Change Detection

```

1 WITH current_snapshot AS (
2     SELECT *, MD5(CONCAT_WS('|',
3         COALESCE(email, ''), COALESCE(subscription_tier, ''),
4         COALESCE(mrr::VARCHAR, ''), COALESCE(health_score::VARCHAR, ''),
5         COALESCE(lifecycle_stage, ''))
6     )) AS record_hash
7     FROM {{ ref('reverse_etl_contacts' ) }}
8 ),
9 previous_sync AS (

```

```

10  SELECT customer_id, record_hash
11  FROM {{ ref('reverse_etl_sync_log') }}
12  WHERE destination = 'salesforce' AND sync_status = 'success'
13  QUALIFY ROW_NUMBER() OVER (
14     PARTITION BY customer_id ORDER BY synced_at DESC
15  ) = 1
16  )
17  SELECT cs.*, CASE
18     WHEN ps.customer_id IS NULL THEN 'insert'
19     WHEN cs.record_hash <> ps.record_hash THEN 'update'
20  END AS sync_action
21  FROM current_snapshot cs
22  LEFT JOIN previous_sync ps USING (customer_id)
23  WHERE ps.customer_id IS NULL OR cs.record_hash <> ps.record_hash

```

Listing 12.2: Hash-based change detection for incremental sync

Warning

Ensure null handling is consistent in hash-based change detection: `COALESCE` every field to a deterministic default before hashing, or a null-to-empty-string change will not trigger a sync. This is the kind of bug that silently corrupts data for three weeks before anyone notices, by which point you have synced 50,000 wrong records to Salesforce and the VP of Sales has opinions.

12.4 Customer 360 Views

A *Customer 360 view* unifies every data point about a customer into a single, terrifyingly comprehensive record. CLAUDE CODE generates a layered approach: an intermediate model joining sources via an identity graph, then a final mart adding RFM scoring and churn risk classification. The result is a table that knows more about your customers than your customers know about themselves.

```

1  WITH rfm AS (
2     SELECT customer_id,
3            DATEDIFF('day', GREATEST(last_active_date,
4            last_payment_date), CURRENT_DATE) AS recency_days,
5            COALESCE(total_active_days, 0) + COALESCE(total_invoices, 0) AS
frequency_raw,
6            COALESCE(lifetime_revenue, 0) AS monetary_value
7  FROM {{ ref('int_customer_360_base') }}

```

```

8 ),
9 rfm_scored AS (
10     SELECT *, NTILE(5) OVER (ORDER BY recency_days DESC) AS
        recency_score,
11         NTILE(5) OVER (ORDER BY frequency_raw ASC) AS frequency_score,
12         NTILE(5) OVER (ORDER BY monetary_value ASC) AS monetary_score
13     FROM rfm
14 )
15 SELECT b.*, r.recency_score, r.frequency_score, r.monetary_score,
16     CASE
17         WHEN r.recency_score >= 4 AND r.frequency_score >= 4 AND r.
monetary_score >= 4 THEN 'Champion'
18         WHEN r.recency_score >= 3 AND r.frequency_score >= 3 THEN 'Loyal'
19
        WHEN r.recency_score <= 2 AND r.frequency_score >= 3 THEN 'At
Risk'
20         ELSE 'Needs Attention'
21     END AS rfm_segment
22 FROM {{ ref('int_customer_360_base') }} b
23 LEFT JOIN rfm_scored r USING (customer_id)

```

Listing 12.3: Customer 360 mart with RFM scoring (abbreviated)

12.5 Syncing to SaaS Tools

Each destination has unique API semantics and rate limits. CLAUDE CODE generates integration code for Salesforce (REST + Bulk API), HubSpot (batch upsert by email), and Intercom (rate-limit-aware per-user updates). Each also has unique ways of making your life difficult, which CLAUDE CODE handles with the quiet resignation of someone who has read too many API changelogs.

```

1 class SalesforceReverseETL:
2     def __init__(self, sf_username, sf_password, sf_security_token,
3                 bulk_threshold=500):
4         self.sf = Salesforce(username=sf_username,
5                               password=sf_password, security_token=sf_security_token)
6         self.bulk_threshold = bulk_threshold
7
8     def upsert_records(self, object_name, external_id_field, records):
9         if len(records) > self.bulk_threshold:
10             results = self.sf.bulk.__getattr__(
11                 object_name

```

```

12         ).upsert(records, external_id_field, batch_size=5000)
13     else:
14         for record in records:
15             ext_id = record.pop(external_id_field)
16             getattr(self.sf, object_name).upsert(
17                 f"{external_id_field}/{ext_id}", record)

```

Listing 12.4: Salesforce reverse ETL sync (core logic)

! Tip

Intercom’s API has stricter rate limits than Salesforce or HubSpot. For large syncs, use bulk import endpoints or schedule during off-peak hours. Alternatively, use CLAUDE CODE to write a rate limiter so sophisticated it could negotiate international trade agreements.

12.6 Claude Code for Audience Segmentation

Marketing teams describe segments in natural language. CLAUDE CODE translates them into precise SQL with proper null handling. This eliminates the traditional workflow of marketing writing a requirements doc, data engineering misinterpreting it, marketing saying “that’s not what I meant,” and everyone going home unhappy.

```

1 WITH current_state AS (
2     SELECT *, DATEDIFF('day', last_active_date, CURRENT_DATE) AS
3     days_since_active,
4     DATEDIFF('day', signup_date, CURRENT_DATE) AS days_since_signup
5     FROM marts.customer_360
6 )
7 -- Trial Power Users
8 SELECT customer_id, email, 'Trial Power Users' AS segment_name
9 FROM current_state
10 WHERE days_since_signup <= 30 AND active_days_last_30 >= 10
11 AND lifecycle_stage IN ('trialing', 'free')
12 UNION ALL
13 -- Churn Risk
14 SELECT customer_id, email, 'Churn Risk'
15 FROM current_state
16 WHERE lifecycle_stage = 'paying'
17 AND (health_score < 40 OR days_since_active > 14)
18 UNION ALL

```

```

18 -- Expansion Ready
19 SELECT customer_id, email, 'Expansion Ready'
20 FROM current_state
21 WHERE lifecycle_stage = 'paying' AND current_mrr < 500 AND health_score
    >= 80

```

Listing 12.5: Audience segmentation SQL (abbreviated)

12.7 Full Example: Customer Activation Pipeline

A complete pipeline using Dagster: refresh Customer 360 (dbt), run segmentation, detect changes, sync to Salesforce, sync segments to HubSpot, push alerts to Slack, and log results. Seven steps. What could possibly go wrong? (Do not answer that.)

```

1 from dagster import asset, Output, MetadataValue, Definitions,
   ScheduleDefinition
2
3 @asset(group_name="reverse_etl")
4 def customer_360_model(context, dbt_resource):
5     result = dbt_resource.cli(["build", "--select", "customer_360"])
6     return Output(None, metadata={"refreshed_at": MetadataValue.text(
7         datetime.utcnow().isoformat())})
8
9 @asset(group_name="reverse_etl", deps=["customer_360_model"])
10 def audience_segments(context, snowflake_resource):
11     df = snowflake_resource.execute_query(SEGMENT_QUERY, fetch_pandas=
12     True)
13     return Output(df, metadata={
14         "total_records": MetadataValue.int(len(df))})
15
16 @asset(group_name="reverse_etl", deps=["salesforce_change_set"])
17 def salesforce_sync(context, salesforce_change_set, salesforce_resource):
18
19     if salesforce_change_set.empty:
20         return Output({"inserted": 0, "updated": 0, "errors": 0})
21     return Output(salesforce_resource.map_and_upsert(
22         df=salesforce_change_set, object_name="Contact",
23         external_id_field="Email", field_mapping=FIELD_MAPPING))
24
25 @asset(group_name="reverse_etl", deps=["audience_segments"])
26 def slack_alerts(context, audience_segments, slack_resource):
27     for segment in ["Churn Risk", "VIP At Risk"]:

```

```

26     segment_df = audience_segments[
27         audience_segments["segment_name"] == segment]
28     if not segment_df.empty:
29         slack_resource.post_message(
30             channel="#customer-alerts",
31             text=f"*{segment}*: {len(segment_df)} customers
identified")

```

Listing 12.6: Dagster assets for activation pipeline (core)

12.8 Operational Alerting System

CLAUDE CODE generates a flexible YAML-based alert configuration with warehouse queries, severity levels, and multi-channel dispatch (Slack, PagerDuty). Because if a revenue drop happens and nobody gets paged, did it really happen? (Yes. Yes it did. And someone will find out on Monday.)

```

1 alerts:
2   - name: revenue_drop
3     description: "Daily revenue dropped >20% vs 7-day avg"
4     severity: critical
5     query: |
6       WITH daily AS (
7         SELECT DATE(ordered_at) AS d, SUM(net_revenue) AS rev
8         FROM analytics.fct_order_items
9         WHERE ordered_at >= DATEADD('day', -30, CURRENT_DATE)
10        GROUP BY 1
11      ), comparison AS (
12        SELECT d, rev, AVG(rev) OVER (
13          ORDER BY d ROWS BETWEEN 7 PRECEDING AND 1 PRECEDING
14        ) AS avg_7d FROM daily
15      )
16      SELECT d, rev, avg_7d FROM comparison
17      WHERE d = CURRENT_DATE - 1 AND (rev - avg_7d)/NULLIF(avg_7d,0) <
-0.20
18     channels:
19       - {type: slack, target: "#revenue-alerts"}
20       - {type: pagerduty, target: "revenue-team"}
21   - name: data_freshness
22     description: "Customer 360 stale (>3 hours)"
23     severity: critical
24     query: |

```

```

25     SELECT MAX(model_refreshed_at) AS last_refresh,
26           DATEDIFF('minute', MAX(model_refreshed_at), CURRENT_TIMESTAMP())
           AS mins
27     FROM marts.customer_360 HAVING mins > 180
28 channels:
29     - {type: slack, target: "#data-engineering"}

```

Listing 12.7: Alert definitions (excerpt)

12.9 Event-Driven Reverse ETL

Batch syncs work for many use cases, but some scenarios demand near-real-time delivery. Event-driven reverse ETL uses CDC or streaming to push updates within seconds. Because apparently hourly syncs are “too slow” now. We live in an era where a three-second delay in updating a Salesforce field is considered a production incident.

```

1 CREATE OR REPLACE STREAM analytics.marts.customer_360_stream
2 ON TABLE analytics.marts.customer_360
3 APPEND_ONLY = FALSE;
4
5 -- Query pending changes
6 SELECT METADATA$ACTION, METADATA$ISUPDATE,
7        customer_id, email, health_score, churn_risk_bucket
8 FROM analytics.marts.customer_360_stream
9 LIMIT 100;

```

Listing 12.8: Creating a Snowflake Stream for reverse ETL

```

1 class SnowflakeStreamPoller:
2     def __init__(self, stream_name):
3         self.stream_name = stream_name
4         self.conn = sf_connect(account=os.environ["SNOWFLAKE_ACCOUNT"],
5                                user=os.environ["SNOWFLAKE_USER"],
6                                password=os.environ["SNOWFLAKE_PASSWORD"])
7
8     def consume_stream(self):
9         cursor = self.conn.cursor()
10        cursor.execute(f"""
11            SELECT *, METADATA$ACTION AS _action,

```

```

12         METADATA$ISUPDATE AS _is_update
13         FROM {self.stream_name}""")
14         columns = [desc[0] for desc in cursor.description]
15         return [dict(zip(columns, row)) for row in cursor]
16
17 class KafkaChangePublisher:
18     def __init__(self, bootstrap_servers, topic):
19         self.topic = topic
20         self.producer = Producer({"bootstrap.servers": bootstrap_servers
21                                   })
22
23     def publish_changes(self, changes):
24         for change in changes:
25             key = str(change.get("CUSTOMER_ID", "unknown"))
26             event = {"action": change.get("_ACTION", "INSERT"),
27                    "data": {k: v for k, v in change.items()
28                              if not k.startswith("_")}}
29             self.producer.produce(self.topic, key=key.encode(),
30                                   value=json.dumps(event, default=str).encode())
31             self.producer.flush()

```

Listing 12.9: Stream poller publishing to Kafka

Warning

Snowflake Streams have a data retention period (default 14 days). If your consumer falls behind, the stream becomes stale and must be recreated. Monitor stream lag proactively. Discovering your stream expired is like discovering your gym membership lapsed—except with more data loss and fewer excuses.

12.10 Monitoring and Observability

Key metrics to track: sync latency, record throughput, error rate, data freshness, and API quota utilization. Create a sync logging table and freshness monitoring view in the warehouse. If you are not monitoring your reverse ETL, you are not doing reverse ETL—you are doing reverse ETL and hoping for the best, which is a different discipline entirely.

```

1 CREATE OR REPLACE VIEW analytics.ops.sync_freshness AS
2 SELECT sync_name, destination,
3        MAX(completed_at) AS last_success,

```

```
4   DATEDIFF('minute', MAX(completed_at), CURRENT_TIMESTAMP()) AS
   staleness_minutes,
5   SUM(CASE WHEN status = 'failed' THEN 1 ELSE 0 END)::FLOAT
6     / NULLIF(COUNT(*), 0) AS failure_rate_7d
7 FROM analytics.ops.reverse_etl_sync_log
8 WHERE started_at >= DATEADD('day', -7, CURRENT_DATE)
9 GROUP BY 1, 2;
```

Listing 12.10: Sync freshness monitoring view

Warning

Monitor API quota utilization for every destination. A runaway sync that exhausts API quota can disrupt operations across your entire organization. There is no Slack message more terrifying than “Why can’t anyone log into Salesforce?” followed by the realization that your sync consumed all the API calls for the day.

12.11 Exercises

1. **Design a Reverse ETL Pipeline.** Generate a dbt health score model, field mapping, change detection query, and monitoring query for syncing to Salesforce. Then use CLAUDE CODE to generate the apology email template for when the first sync overwrites everyone’s contact notes.
2. **Audience Segmentation Challenge.** Translate five plain-English segment definitions into SQL. Handle the case where a customer qualifies for multiple segments by assigning highest-priority only. Use CLAUDE CODE to also generate the segment names, because marketing will reject your first five suggestions anyway.
3. **Multi-Destination Sync.** Build a configuration-driven sync engine that reads YAML mapping files and routes records to Salesforce, HubSpot, and Intercom simultaneously. Contemplate the existential implications of being responsible for data consistency across three SaaS platforms.
4. **Alert Engineering.** Write alert definitions for: signup volume spike ($> 3\sigma$), payment failure rate $> 10\%$ in the last hour, and feature adoption rate dropping below 5%. Use CLAUDE CODE to also write the PagerDuty escalation policy, because someone has to get woken up at 3 AM and it should not always be you.

13

Machine Learning Feature Engineering

“Garbage in, garbage out” is the polite version. The full version is “garbage in, a model that confidently predicts garbage, a dashboard that visualizes garbage, and an executive who makes decisions based on garbage.”

Every ML model is only as good as its features. Raw data must be transformed into meaningful numerical representations before an algorithm can learn. This transformation—**feature engineering**—is one of the most time-consuming stages of the ML lifecycle, and also the stage most likely to be described as “almost done” for six consecutive sprints. CLAUDE CODE changes the economics: instead of spending hours brainstorming and hand-coding, describe the business domain and schema to CLAUDE CODE and receive ranked feature candidates with ready-to-run code.

In this chapter we cover feature engineering fundamentals, **feature store** architecture (Feast and Tecton), and how CLAUDE CODE accelerates feature discovery, generation, embedding management, documentation, and time-series work. Essentially, CLAUDE CODE does the tedious parts so you can focus on the glamorous work of explaining to stakeholders why the model’s predictions are not magic.

13.1 Feature Engineering Fundamentals

A **feature** is a measurable property of an observed phenomenon. Features fall into categories: raw (taken directly from source), derived (computed from columns),

aggregate (summaries over windows), embedding (dense vectors from deep learning), and interaction (combinations of features).

Tip

Before productionizing notebook feature logic, always ask: “What entity does this feature describe, and what is its time granularity?” These two questions prevent most join-key and temporal-leakage bugs. They also prevent the classic “it works in the notebook” followed by “it produces nonsense in production” arc that we have all lived through at least twice.

Common Pitfalls

Pitfall	Description
Target leakage	Including information unavailable at prediction time.
Train/serve skew	Features computed differently in training vs. serving.
Entity misalignment	Joining features at the wrong granularity.
Stale features	Serving values computed hours or days ago.
Feature explosion	Thousands of features without pruning.

Warning

Target leakage is the most common source of “too good to be true” performance. When CLAUDE CODE generates feature candidates, verify each one: “Would this value be known at prediction time in production?” If your churn model has 99.7% accuracy, it is not because you are brilliant. It is because you leaked the target.

13.2 Feature Stores: Concepts and Architecture

A *feature store* manages the lifecycle of features: computation, storage, serving, monitoring, and discovery. It provides a registry (catalog of definitions), an offline store (warehouse for training), an online store (key-value DB for inference), materialization pipelines, and a serving API. It is, in effect, a warehouse for your warehouse’s data, which is exactly the kind of recursion that keeps data engineers employed.

Feast vs. Tecton

Criterion	Feast	Tecton
Deployment	Self-hosted	Managed SaaS
Streaming features	Limited	Native support
Monitoring	Community plugins	Built-in
Best for	Teams with strong infra skills	Teams wanting managed ML infra

i Note

Feast and Tecton are complementary. Feast is ideal when you want full infrastructure control and enjoy debugging Kubernetes at 2 AM. Tecton is better for fully managed platforms with streaming support and the budget to prove it.

13.3 Using Claude Code for Feature Discovery

The most effective prompt pattern has four components: business context (what you are predicting), schema description, constraints (what data is available at prediction time), and output format (feature name, SQL, rationale, importance, leakage risk). In other words, give CLAUDE CODE the same context you would give a new hire, except CLAUDE CODE will not spend the first week setting up its development environment.

```

1 I am building a churn prediction model for a B2B SaaS product.
2 Schema:
3 - users (user_id, company_id, created_at, plan_type, role)
4 - events (event_id, user_id, event_type, timestamp)
5 - subscriptions (sub_id, company_id, plan, mrr, start_date, end_date)
6 - support_tickets (ticket_id, company_id, created_at, priority,
7   resolved_at, satisfaction_score)
8
9 Target: will this company churn in the next 30 days?
10 Only data available up to the prediction date can be used.
11
12 Generate 20 candidate features with: name, SQL, rationale,
13 expected importance (high/medium/low), leakage risk.
```

CLAUDE CODE's output includes features like `active_users_ratio_30d` (high importance, no leakage), `support_tickets_open_14d` (high), `days_since_last_login` (high), and `mrr_change_pct_90d` (medium). After the initial round, use follow-up prompts to generate exact SQL, unit tests, and edge cases for the top candidates. It is like pair programming with someone who has read every feature engineering blog post ever written.

13.4 Automated Feature Generation

```
1 from datetime import timedelta
2 from feast import BigQuerySource, Entity, FeatureView, Field, ValueType
3 from feast.types import Float64, Int64
4
5 customer = Entity(name="customer", join_keys=["customer_id"],
6                 value_type=ValueTypes.STRING)
7
8 customer_features_source = BigQuerySource(
9     name="customer_features_daily",
10    table="analytics.customer_features_daily",
11    timestamp_field="feature_date",
12 )
13
14 customer_feature_view = FeatureView(
15     name="customer_features_daily",
16     entities=[customer],
17     ttl=timedelta(days=2),
18     schema=[
19         Field(name="total_orders_30d", dtype=Int64,
20              description="Orders placed in last 30 days"),
21         Field(name="avg_order_value_30d", dtype=Float64,
22              description="Mean order value (USD) last 30 days"),
23         Field(name="days_since_last_order", dtype=Int64,
24              description="Days since most recent order"),
25         Field(name="total_revenue_lifetime", dtype=Float64,
26              description="Cumulative customer revenue"),
27         Field(name="support_tickets_open", dtype=Int64,
28              description="Currently open support tickets"),
29     ],
30     source=customer_features_source,
31     online=True,
32     tags={"team": "growth", "refresh_cadence": "daily"},
33 )
```

Listing 13.1: Feast FeatureView definition generated by Claude Code

```

1 WITH order_stats AS (
2     SELECT customer_id, CURRENT_DATE() AS feature_date,
3         COUNT(*) AS total_orders_30d,
4         AVG(order_total) AS avg_order_value_30d,
5         DATE_DIFF(CURRENT_DATE(), MAX(order_date), DAY)
6             AS days_since_last_order
7     FROM `warehouse.orders`
8     WHERE order_date >= DATE_SUB(CURRENT_DATE(), INTERVAL 30 DAY)
9         AND order_status IN ('completed', 'shipped')
10    GROUP BY customer_id
11 ),
12 lifetime AS (
13     SELECT customer_id, SUM(order_total) AS total_revenue_lifetime
14     FROM `warehouse.orders`
15     WHERE order_status IN ('completed', 'shipped')
16    GROUP BY customer_id
17 ),
18 support AS (
19     SELECT customer_id, COUNTIF(status = 'open') AS support_tickets_open
20     FROM `warehouse.support_tickets` GROUP BY customer_id
21 )
22 SELECT c.customer_id, CURRENT_TIMESTAMP() AS feature_date,
23     COALESCE(o.total_orders_30d, 0),
24     COALESCE(o.avg_order_value_30d, 0.0),
25     COALESCE(o.days_since_last_order, 9999),
26     COALESCE(l.total_revenue_lifetime, 0.0),
27     COALESCE(s.support_tickets_open, 0)
28 FROM `warehouse.customers` c
29 LEFT JOIN order_stats o USING (customer_id)
30 LEFT JOIN lifetime l USING (customer_id)
31 LEFT JOIN support s USING (customer_id);

```

Listing 13.2: Feature computation SQL generated by Claude Code

! Tip

When asking CLAUDE CODE to generate Feast definitions, always specify the Feast version. The API changed significantly between 0.24 and 0.30+. Feast's changelog reads like a thriller novel where the protagonist keeps changing their

name and appearance between chapters.

13.5 Embedding Generation and Management

Embeddings are dense vector representations capturing semantic meaning—essential for text, images, user behavior, and high-cardinality categoricals. They are also the answer to the question “how do I turn words into numbers without anyone understanding how?”

```
1 from sentence_transformers import SentenceTransformer
2 import numpy as np
3 import pandas as pd
4
5 def generate_embeddings(texts: list[str], model_name="all-MiniLM-L6-v2"):
6
7     model = SentenceTransformer(model_name)
8     return model.encode(texts, batch_size=256, normalize_embeddings=True)
9
10 def embeddings_to_dataframe(product_ids, embeddings):
11     cols = {f"emb_{i}": embeddings[:, i] for i in range(embeddings.shape
12 [1])}
13     df = pd.DataFrame(cols)
14     df["product_id"] = product_ids
15     df["event_timestamp"] = pd.Timestamp.now(tz="UTC")
16     return df
```

Listing 13.3: Embedding generation pipeline

! Tip

Always store embedding model metadata (model name, dimensionality, normalization method) alongside the vectors. This prevents silent dimension mismatches when models are swapped. There is no debugging experience quite as soul-crushing as discovering your 384-dimensional vectors were being compared to 768-dimensional vectors for the past month.

13.6 Time-Series Feature Engineering

Time-series features include rolling aggregates, lag features, difference features, calendar features, Fourier features, and exponential smoothing. If that sounds like a lot, it is. Use CLAUDE CODE to generate them all at once. You were not put on this earth to manually write rolling window aggregates.

```
1 def add_rolling_features(df, value_col, entity_col,
2                          windows=(7, 14, 30), aggs=("mean", "std")):
3     df = df.sort_values([entity_col, "date"])
4     grouped = df.groupby(entity_col)[value_col]
5     for window in windows:
6         rolling = grouped.rolling(window=window, min_periods=1)
7         for agg in aggs:
8             df[f"{value_col}_rolling_{agg}_{window}d"] = \
9                 getattr(rolling, agg)().reset_index(level=0, drop=True)
10    return df
11
12 def add_lag_features(df, value_col, entity_col, lags=(1, 7, 30)):
13     df = df.sort_values([entity_col, "date"])
14     for lag in lags:
15         df[f"{value_col}_lag_{lag}d"] = \
16             df.groupby(entity_col)[value_col].shift(lag)
17    return df
18
19 def add_calendar_features(df, date_col="date"):
20     dt = pd.to_datetime(df[date_col])
21     df["day_of_week"] = dt.dt.dayofweek
22     df["month"] = dt.dt.month
23     df["is_weekend"] = (dt.dt.dayofweek >= 5).astype(int)
24    return df
```

Listing 13.4: Time-series feature utilities

Warning

A common mistake is computing a rolling average that includes the current row's value. If the current value correlates with the target, this introduces leakage. Always shift the window by one period. This is the kind of bug that lets you publish a paper about your amazing model performance before the retraction.

13.7 Full Example: Building a Feature Store with Feast

We build a complete feature store for e-commerce serving churn prediction and recommendation models, with CLAUDE CODE assisting at every stage. By “assisting” we mean “doing approximately 90% of the work while we provide the domain context and take credit in the architecture review.”

Project Structure and Configuration

```
1 # feature_store.yaml
2 project: ecommerce_features
3 registry: data/registry.db
4 provider: local
5 online_store:
6   type: redis
7   connection_string: "redis://localhost:6379"
8 offline_store:
9   type: file
```

Feature Services and Retrieval

```
1 from feast import FeatureStore, FeatureService
2
3 churn_service = FeatureService(
4     name="churn_prediction",
5     features=[customer_engagement_fv, customer_support_fv],
6     description="All features for churn prediction model v3.",
7 )
8
9 # Historical retrieval for training
10 store = FeatureStore(repo_path="feature_repo/")
11 training_df = store.get_historical_features(
12     entity_df=entity_df,
13     features=store.get_feature_service("churn_prediction"),
14 ).to_df()
15
16 # Online retrieval for inference
17 feature_vector = store.get_online_features(
18     features=store.get_feature_service("churn_prediction"),
```

```
19     entity_rows=[{"customer_id": "C001"}],
20 ).to_dict()
```

Listing 13.5: Feature service and training retrieval

Freshness Validation

```
1 def validate_freshness(store, max_staleness_hours=48):
2     errors = []
3     cutoff = datetime.now(timezone.utc) - timedelta(hours=
4     max_staleness_hours)
5     for fv in store.list_feature_views():
6         if fv.materialization_intervals:
7             latest = max(i.end_date for i in fv.
8             materialization_intervals)
9             if latest < cutoff:
10                errors.append(f"'{fv.name}' last materialized at {latest}
11                ")
12            else:
13                errors.append(f"'{fv.name}' has never been materialized.")
14    return errors
```

Listing 13.6: Feature freshness validation

! Tip

Schedule freshness validation immediately after materialization in your orchestration DAG. If it fails, alert the on-call engineer *before* the model consumes stale features. Serving a model with 48-hour-old features is like navigating with a two-day-old weather forecast—technically a forecast, practically useless.

13.8 Full Example: Feature Discovery for Recommendations

Recommendation systems are among the most feature-intensive ML applications. Here we use CLAUDE CODE to discover features for a hybrid recommender serving 50,000 products across 200 categories. The alternative is brainstorming features in a whiteboard session, which will produce exactly seven features before devolving into a debate about lunch.

User Features

```

1 USER_FEATURES = {
2   "user_total_pageviews_7d": {
3     "sql": """SELECT user_id, COUNT(*) AS user_total_pageviews_7d
4       FROM page_views
5       WHERE timestamp >= CURRENT_DATE - INTERVAL '7 days'
6       GROUP BY user_id""",
7     "rationale": "Active users have more diverse browsing patterns",
8     "importance": "high",
9   },
10  "user_avg_purchase_price_90d": {
11    "sql": """SELECT user_id,
12      AVG(revenue / NULLIF(quantity, 0)) AS
13      user_avg_purchase_price_90d
14      FROM orders WHERE order_date >= CURRENT_DATE - INTERVAL '90
15      days'
16      GROUP BY user_id""",
17    "rationale": "Price affinity filters out-of-range products",
18    "importance": "high",
19  },
20  "user_cart_abandonment_rate_30d": {
21    "sql": """WITH adds AS (
22      SELECT user_id, COUNT(*) AS add_count FROM cart_events
23      WHERE action = 'add' AND timestamp >= CURRENT_DATE -
24      INTERVAL '30 days'
25      GROUP BY user_id
26    ), purchases AS (
27      SELECT user_id, COUNT(DISTINCT product_id) AS buy_count FROM
28      orders
29      WHERE order_date >= CURRENT_DATE - INTERVAL '30 days'
30      GROUP BY user_id
31    )
32    SELECT a.user_id, 1.0 - COALESCE(p.buy_count, 0)::FLOAT
33      / NULLIF(a.add_count, 0) AS user_cart_abandonment_rate_30d
34    FROM adds a LEFT JOIN purchases p USING (user_id)""",
35    "rationale": "High abandonment signals price sensitivity",
36    "importance": "medium",
37  },
38 }

```

Listing 13.7: User-level features for recommendations

Item Features

```

1 ITEM_FEATURES = {
2   "item_view_to_purchase_rate_30d": {
3     "sql": """"WITH views AS (
4       SELECT product_id, COUNT(DISTINCT user_id) AS viewers
5       FROM page_views WHERE timestamp >= CURRENT_DATE - INTERVAL
6       '30 days'
7       GROUP BY product_id
8     ), buys AS (
9       SELECT product_id, COUNT(DISTINCT user_id) AS buyers
10      FROM orders WHERE order_date >= CURRENT_DATE - INTERVAL '30
11      days'
12      GROUP BY product_id
13    )
14    SELECT v.product_id, COALESCE(b.buyers, 0)::FLOAT
15    / NULLIF(v.viewers, 0) AS item_view_to_purchase_rate_30d
16    FROM views v LEFT JOIN buys b USING (product_id)""",
17    "rationale": "High conversion items are strong recommendation
18    candidates",
19    "importance": "high",
20  },
21  "item_trending_score_7d": {
22    "sql": """"WITH curr AS (
23      SELECT product_id, COUNT(*) AS views_7d FROM page_views
24      WHERE timestamp >= CURRENT_DATE - INTERVAL '7 days' GROUP BY
25      1
26    ), prev AS (
27      SELECT product_id, COUNT(*) AS views_prev FROM page_views
28      WHERE timestamp BETWEEN CURRENT_DATE - INTERVAL '14 days'
29      AND CURRENT_DATE - INTERVAL '7 days' GROUP BY 1
30    )
31    SELECT c.product_id, (c.views_7d - COALESCE(p.views_prev, 0))::
32    FLOAT
33    / NULLIF(COALESCE(p.views_prev, 1), 0) AS
34    item_trending_score_7d
35    FROM curr c LEFT JOIN prev p USING (product_id)""",
36    "rationale": "Trending products capture emerging demand",
37    "importance": "medium",
38  },
39 }

```

Listing 13.8: Item-level features

Interaction Features

Key interaction features include `user_item_view_count_30d` (repeated views indicate interest), `user_item_in_cart` (current cart items are strong purchase signals), and `user_purchased_same_brand` (brand affinity drives repeats).

Note

In production, store credentials in a secrets manager, not in code. CLAUDE CODE can generate boilerplate for AWS Secrets Manager, HashiCorp Vault, or your orchestrator's secrets facility. It can also generate the boilerplate for rotating those secrets, which is the part everyone forgets until the credentials expire on a Friday afternoon.

13.9 Best Practices for Production Feature Pipelines

Idempotency

Every feature computation should be idempotent: running it twice with the same inputs produces the same outputs. Use the reference date as a partition key. If your feature pipeline is not idempotent, you are not doing feature engineering—you are doing feature roulette.

Schema Enforcement

Use Pandera or Great Expectations to enforce feature schemas at write time:

```
1 import pandera as pa
2
3 customer_schema = pa.DataFrameSchema({
4     "customer_id": pa.Column(str, nullable=False, unique=True),
5     "total_orders_30d": pa.Column(int, pa.Check.ge(0), nullable=False),
6     "avg_order_value_30d": pa.Column(float, pa.Check.ge(0.0), nullable=
7     False),
8 }, strict=True, coerce=True)
```

Drift Monitoring

Track feature distributions using Population Stability Index (PSI): < 0.1 = no shift, $0.1-0.2$ = moderate, > 0.2 = significant. A PSI above 0.2 means either the world changed or your pipeline broke. Both possibilities are equally alarming.

```
1 def compute_psi(reference, current, bins=10, epsilon=1e-4):
2     bin_edges = np.percentile(reference, np.linspace(0, 100, bins + 1))
3     bin_edges[0], bin_edges[-1] = -np.inf, np.inf
4     ref_pct = np.histogram(reference, bins=bin_edges)[0] / len(reference)
5     + epsilon
6     cur_pct = np.histogram(current, bins=bin_edges)[0] / len(current) +
7     epsilon
8     return float(np.sum((cur_pct - ref_pct) * np.log(cur_pct / ref_pct)))
```

Listing 13.9: PSI computation

Warning

Never test feature pipelines only in isolation. The most dangerous bugs appear at boundaries between stages: type coercions, timezone mismatches, and fan-out joins. Integration tests on realistic sample data catch these. Unit tests catch the bugs you imagined. Integration tests catch the bugs that actually happen.

13.10 Exercises

1. **Feature Discovery.** For a fraud detection model with `transactions`, `users`, and `merchants` tables, use CLAUDE CODE to generate 15 candidate features with SQL, leakage assessment, and rationale. Bonus: ask CLAUDE CODE to also predict which features your stakeholder will reject for being “too complicated to explain.”
2. **Time-Series Features.** Given hourly server metrics, write Python code for rolling, lag, difference, and Fourier features. Add Pandera schema validation and three edge-case unit tests. Use CLAUDE CODE to write the Fourier features because nobody remembers the Fourier transform after the exam.

3. **Feature Store Setup.** Set up a local Feast store with two entities, three feature views, a materialization script, and freshness validation. If you complete this in under two hours, you are either lying or have done this before.
4. **Drift Monitoring.** Implement PSI monitoring with Slack alerting when the threshold is exceeded. Verify the implementation returns 0 for identical distributions. Use CLAUDE CODE to also generate the Slack message, because “PSI exceeded threshold” will mean nothing to your on-call engineer at 3 AM.

Part **V**

Orchestration and Infrastructure

14

Orchestrating Data Pipelines with Claude Code

“Without an orchestrator, we had cron jobs. With an orchestrator, we have cron jobs that send Slack messages when they fail. This is what progress looks like.”

Modern data engineering demands robust orchestration: scheduling, dependency management, error handling, and monitoring across hundreds of interconnected tasks. In this chapter we explore how CLAUDE CODE accelerates every phase of orchestration across three dominant frameworks: Apache Airflow, Dagster, and Prefect. Sure, you could write this Terraform yourself. Wait, wrong chapter. Sure, you could hand-write all your DAGs yourself. You could also walk to work instead of driving. Both are technically possible and equally inadvisable.

14.1 Why Orchestration Matters

Orchestration coordinates discrete tasks according to a dependency graph and schedule: extraction, transformation, loading, validation, and notification. Without an orchestrator, teams resort to cron jobs and ad hoc retry logic—a pattern that collapses under hundreds of interdependent tasks and exactly one sick day from the person who understands the crontab.

The challenge is that each orchestrator imposes its own abstractions. Airflow has Operators, Sensors, and XComs; Dagster has Assets, Ops, and Resources; Prefect has Tasks, Flows, and Blocks. CLAUDE CODE collapses this learning curve by trans-

lating intent into idiomatic orchestrator code. It has memorized all three frameworks so you do not have to, which is fortunate because memorizing all three frameworks is how hobbies die.

A well-orchestrated data platform has three properties: **determinism** (given the same inputs, a pipeline produces the same outputs), **observability** (you can see what ran, when, and whether it succeeded), and **recoverability** (when something fails, you can retry from the failure point without reprocessing everything). Achieving all three simultaneously is referred to in the industry as “the dream.”

Note

All three orchestrators are open-source and offer managed cloud products. CLAUDE CODE can help evaluate trade-offs between self-hosting and managed offerings based on team size, compliance requirements, and budget constraints. It cannot, however, help you win the internal political battle over which orchestrator to choose. That requires a different kind of intelligence.

14.2 Airflow DAG Generation

```
1 from airflow.decorators import dag, task
2 from airflow.operators.bash import BashOperator
3 from airflow.providers.slack.operators.slack_webhook import
  SlackWebhookOperator
4 from airflow.utils.trigger_rule import TriggerRule
5 import pendulum
6 from datetime import timedelta
7
8 default_args = {
9     "owner": "data-engineering",
10    "retries": 3,
11    "retry_delay": timedelta(minutes=5),
12    "retry_exponential_backoff": True,
13    "execution_timeout": timedelta(hours=1),
14    "sla": timedelta(hours=2),
15 }
16
17 TABLES = ["orders", "customers"]
18
19 @dag(dag_id="elt_postgres_to_snowflake",
20     schedule="0 2 * * *",
21     start_date=pendulum.datetime(2025, 1, 1, tz="UTC"),
```

```
22     catchup=False, default_args=default_args,
23     tags=["elt", "production"])
24 def elt_pipeline():
25
26     @task(pool="postgres_pool", pool_slots=1)
27     def extract_and_load(table: str, **context):
28         # Extract from Postgres, stage into Snowflake
29         pass
30
31     dbt_build = BashOperator(
32         task_id="dbt_build",
33         bash_command="cd /opt/airflow/dbt && dbt build --select marts",
34     )
35     validate = BashOperator(
36         task_id="validate",
37         bash_command="great_expectations checkpoint run marts_checkpoint
38     ",
39     )
40     notify_success = SlackWebhookOperator(
41         task_id="notify_success",
42         slack_webhook_conn_id="slack_data",
43         message=":white_check_mark: ELT completed for {{ ds }}",
44         trigger_rule=TriggerRule.ALL_SUCCESS,
45     )
46     notify_failure = SlackWebhookOperator(
47         task_id="notify_failure",
48         slack_webhook_conn_id="slack_data",
49         message=":x: ELT FAILED for {{ ds }}",
50         trigger_rule=TriggerRule.ONE_FAILED,
51     )
52
53     loaded = extract_and_load.expand(table=TABLES)
54     loaded >> dbt_build >> validate >> [notify_success, notify_failure]
55
56     elt_pipeline()
```

Listing 14.1: Production Airflow DAG generated with Claude Code

Key production features: dynamic task mapping via `.expand()`, pool limiting to avoid overwhelming sources, exponential backoff on retries, and separate success/failure notification paths.

! Tip

When prompting CLAUDE CODE for Airflow DAGs, always specify your Airflow version. Airflow 2.x introduced the TaskFlow API, dynamic task mapping, and dataset-driven scheduling. Code written for Airflow 1.x will not work on 2.x and vice versa. Discovering this at deploy time is a rite of passage, but not a pleasant one.

Common Airflow Anti-Patterns

CLAUDE CODE can identify and fix common Airflow anti-patterns: top-level imports of heavy libraries (which slow DAG parsing), using `PythonOperator` where `@task` is cleaner, hardcoded connection strings instead of Airflow Connections, and monolithic DAGs that should be split into smaller, independently schedulable units. In other words, CLAUDE CODE can review your DAGs and diplomatically suggest that you rewrite most of them.

14.3 Dagster Asset Definitions

Dagster defines *software-defined assets*—the data artifacts your pipeline produces—and infers the execution graph from dependencies.

```

1 from dagster import (asset, AssetExecutionContext, MaterializeResult,
2     MetadataValue, Definitions, ScheduleDefinition, define_asset_job,
3     RetryPolicy, Backoff)
4
5 @asset(group_name="raw",
6     retry_policy=RetryPolicy(max_retries=3, delay=30,
7         backoff=Backoff.EXPONENTIAL))
8 def raw_orders(context, snowflake):
9     with snowflake.get_connection() as conn:
10         row_count = conn.cursor().execute("COPY INTO raw.orders ...").
11         fetchone()[0]
12     return MaterializeResult(
13         metadata={"row_count": MetadataValue.int(row_count)})
14
15 @asset(group_name="marts", deps=["stg_orders", "stg_customers"])
16 def order_metrics(context, snowflake):
17     with snowflake.get_connection() as conn:
18         conn.cursor().execute("""
19             CREATE OR REPLACE TABLE marts.order_metrics AS

```

```

19         SELECT c.region, DATE_TRUNC('month', o.order_date) AS month,
20                COUNT(DISTINCT o.order_id) AS orders,
21                SUM(o.total_amount) AS revenue
22         FROM staging.orders o JOIN staging.customers c USING (
           customer_id)
23         GROUP BY 1, 2""")
24     return MaterializeResult(metadata={"model": MetadataValue.text("
           marts.order_metrics")})
25
26     defs = Definitions(
27         assets=[raw_orders, order_metrics],
28         jobs=[define_asset_job("nightly_elt")],
29         schedules=[ScheduleDefinition(job_name="nightly_elt", cron_schedule=
           "0 2 * * *")],
30     )

```

Listing 14.2: Dagster assets for an ELT pipeline

Dagster assets are plain Python functions that can be unit-tested without the orchestrator running. Dagster also natively supports freshness-based scheduling, where downstream assets are materialized when their upstream dependencies are refreshed. It is almost as if someone designed an orchestrator for the way data actually works. Revolutionary.

Note

Dagster vs. Airflow: When to Choose Which Choose Airflow when you need a battle-tested, widely-adopted orchestrator with a massive ecosystem of providers and a community large enough to answer your Stack Overflow questions within the hour. Choose Dagster when you want asset-based thinking, strong typing, and built-in data quality checks. CLAUDE CODE can translate between the two if you decide to migrate, which it will do with zero complaints—unlike your engineering team.

14.4 Prefect Flow Authoring

Prefect is the most Pythonic orchestrator. Flows and tasks are decorated functions with minimal boilerplate. If you can write a Python function, you can write a Prefect flow. If you cannot write a Python function, CLAUDE CODE can write both.

```

1 from prefect import flow, task, get_run_logger

```

```

2 from prefect.tasks import task_input_hash
3 from datetime import timedelta
4 import httpx
5
6 @task(retries=3, retry_delay_seconds=[30, 60, 120],
7       cache_key_fn=task_input_hash, cache_expiration=timedelta(hours=1))
8 def extract_api_page(endpoint: str, page: int, api_key: str):
9     resp = httpx.get(f"https://api.example.com/v2/{endpoint}",
10                    params={"page": page}, headers={"Authorization": f"Bearer {
11                    api_key}"})
12     resp.raise_for_status()
13     return resp.json()["results"]
14
15 @flow(name="incremental-api-ingest", retries=1, timeout_seconds=3600)
16 def incremental_ingest(endpoint="transactions", pages=10):
17     from prefect.blocks.system import Secret
18     api_key = Secret.load("example-api-key").get()
19     page_futures = extract_api_page.map(
20         endpoint=[endpoint]*pages, page=list(range(1, pages+1)),
21         api_key=[api_key]*pages)
22     all_records = []
23     for future in page_futures:
24         all_records.extend(future.result())
25     load_to_snowflake(records=all_records, table=endpoint)

```

Listing 14.3: Prefect flow for API ingestion

! Tip

Prefect Task Caching Notice the `cache_key_fn=task_input_hash` parameter. This caches task results based on inputs, so retrying a flow skips already-completed tasks. This is especially valuable for expensive extraction tasks where re-fetching data wastes time and API quota. It is also valuable for your sanity, which is a non-renewable resource.

14.5 Task Dependency Analysis

CLAUDE CODE can parse SQL queries to extract table-level lineage, building an adjacency list that can be fed into any orchestrator to generate correct task ordering. This is particularly valuable when migrating from ad hoc SQL scripts to an orchestrated pipeline—a process that archaeologists refer to as “discovering the de-

pendency graph that was there all along, hidden in 47 cron jobs and a shared Google Sheet.”

```
1 import re
2 from collections import defaultdict
3
4 def extract_table_lineage(sql: str):
5     normalized = sql.strip().upper()
6     targets = set()
7     for p in [r"CREATE\s+(?:OR\s+REPLACE\s+)?TABLE\s+([\w.]+\s+AS",
8             r"INSERT\s+INTO\s+([\w.]+)"]:
9         m = re.search(p, normalized)
10        if m: targets.add(m.group(1).lower())
11    sources = set()
12    for p in [r"FROM\s+([\w.]+)", r"JOIN\s+([\w.]+)"]:
13        for m in re.finditer(p, normalized):
14            sources.add(m.group(1).lower())
15    return targets, sources - targets
16
17 def build_dependency_graph(sql_files: dict[str, str]):
18     graph = defaultdict(list)
19     target_map = {}
20     for name, sql in sql_files.items():
21         targets, _ = extract_table_lineage(sql)
22         for t in targets: target_map[t] = name
23     for name, sql in sql_files.items():
24         _, sources = extract_table_lineage(sql)
25         for s in sources:
26             graph[name].append(target_map.get(s, f"[external] {s}"))
27     return dict(graph)
```

Listing 14.4: SQL-based dependency extractor

Warning

Regex Limitations for SQL Parsing The regex-based approach above handles simple SQL but misses CTEs, subqueries, and dynamic SQL. For production lineage extraction, use a proper SQL parser like `sqlglot` or ask CLAUDE CODE to generate a parser that handles your specific SQL dialect’s edge cases. Using regex to parse SQL is like using a butter knife to perform surgery—it might work, but you will not be proud of the results.

14.6 Error Handling and Retry Logic

CLAUDE CODE helps implement layered error handling that distinguishes transient from permanent failures. Classify errors into TRANSIENT (retry with back-off), THROTTLED (retry after longer delay), PERMANENT (alert on-call), and DATA_QUALITY (route to data team). This prevents the common anti-pattern of retrying a permissions error seventeen times before giving up and paging someone who also cannot fix it.

```
1 from enum import Enum, auto
2 import re
3
4 class ErrorCategory(Enum):
5     TRANSIENT = auto()
6     THROTTLED = auto()
7     PERMANENT = auto()
8     DATA_QUALITY = auto()
9
10 def classify_error(exc):
11     msg = str(exc).lower()
12     if re.search(r"429|rate.limit|throttled", msg):
13         return ErrorCategory.THROTTLED
14     if re.search(r"timeout|connection.reset|503", msg):
15         return ErrorCategory.TRANSIENT
16     if re.search(r"not.null.constraint|duplicate.key", msg):
17         return ErrorCategory.DATA_QUALITY
18     return ErrorCategory.PERMANENT
19
20 RETRY_DELAYS = {
21     ErrorCategory.TRANSIENT: [30, 60, 120],
22     ErrorCategory.THROTTLED: [60, 300, 900],
23     ErrorCategory.DATA_QUALITY: [], # no retry
24     ErrorCategory.PERMANENT: [], # no retry
25 }
```

Listing 14.5: Error classification for retry decisions

! Tip

Combine failure callbacks with `on_retry_callback` to track retry frequency by error category. This telemetry identifies flaky dependencies before they cause outages. A task that retries three times every night is a warning sign

even if it eventually succeeds. It is the pipeline equivalent of a check engine light that you have been ignoring for six months.

14.7 Monitoring and Alerting

Orchestrator-native UIs are a starting point, but production needs deeper observability. CLAUDE CODE generates Prometheus exporters and Grafana dashboards that go beyond the orchestrator's built-in metrics. Because if something fails and no dashboard shows it, you will find out from the CFO, and that is always worse.

```
1 from prometheus_client import start_http_server, Counter, Histogram,  
    Gauge  
2  
3 TASK_DURATION = Histogram(  
4     "airflow_task_duration_seconds",  
5     "Task execution time",  
6     ["dag_id", "task_id"],  
7     buckets=[10, 30, 60, 300, 600, 1800, 3600],  
8 )  
9 TASK_FAILURES = Counter(  
10    "airflow_task_failures_total",  
11    "Total task failures",  
12    ["dag_id", "task_id"],  
13 )  
14 SLA_MISSES = Counter(  
15    "airflow_sla_misses_total",  
16    "Total SLA misses",  
17    ["dag_id"],  
18 )  
19 DAG_LAST_SUCCESS = Gauge(  
20    "airflow_dag_last_success_timestamp",  
21    "Timestamp of last successful DAG run",  
22    ["dag_id"],  
23 )
```

Listing 14.6: Prometheus metrics for Airflow

The four essential dashboards for a data platform are: (1) Pipeline Health (success/failure rates, SLA compliance), (2) Performance Trends (task durations over time, identifying degradation), (3) Resource Utilization (worker CPU, memory, queue depth), and (4) Data Freshness (time since last successful materialization of each

dataset). If you have fewer than four dashboards, you are flying blind. If you have more than forty, you are dashboard-hoarding and need an intervention.

Note

For Dagster, stream the event log to Kafka for custom metrics. For Prefect, use the API's run-level telemetry. CLAUDE CODE can transform either into Prometheus-compatible exporters with minimal custom code. It will also name your metrics sensibly, which is more than can be said for most humans.

14.8 Dynamic DAG Generation

The most robust dynamic pattern reads a metadata table at DAG parse time and generates tasks accordingly. Airflow 2.4+ also supports *datasets* as a scheduling primitive—a producer DAG declares an output dataset, and consumer DAGs trigger when it updates.

```
1 from airflow import Dataset
2 from airflow.decorators import dag, task
3 import pendulum
4
5 orders_dataset = Dataset("snowflake://analytics/staging/orders")
6
7 @dag(dag_id="produce_staging", schedule="0 2 * * *",
8     start_date=pendulum.datetime(2025, 1, 1, tz="UTC"), catchup=False)
9 def produce():
10     @task(outlets=[orders_dataset])
11     def refresh_orders():
12         pass
13     refresh_orders()
14
15 produce()
16
17 @dag(dag_id="build_marts", schedule=[orders_dataset],
18     start_date=pendulum.datetime(2025, 1, 1, tz="UTC"), catchup=False)
19 def consume():
20     @task
21     def build_metrics():
22         pass
23     build_metrics()
24
25 consume()
```

Listing 14.7: Dataset-triggered DAGs in Airflow 2.4+

 **Warning**

Dynamically generated DAGs can confuse Airflow’s serialization if the factory produces inconsistent DAG counts across restarts. Pin config files in version control and avoid reading from volatile sources at parse time. Test DAG generation in CI to catch serialization issues before deployment. Nothing is more humbling than a production outage caused by a YAML typo.

Configuration-Driven DAG Factories

For teams managing dozens of similar pipelines, a DAG factory reads YAML configuration and generates Airflow DAGs programmatically. Each YAML file specifies the source, destination, schedule, and any custom transformation logic. CLAUDE CODE can generate the factory code and the initial set of YAML configurations from your existing ad hoc pipelines. It can also generate the documentation explaining why you have a DAG factory, which you will need when someone new joins and asks “why don’t we just write DAGs normally?”

14.9 Scheduling Optimization

Naive scheduling—running everything at midnight—causes resource contention. It is the data engineering equivalent of everyone trying to leave the parking garage at 5 PM. CLAUDE CODE generates analysis scripts that examine historical run data and recommend staggered schedules to minimize peak concurrency.

```
1 def analyze_schedule_conflicts(run_history: list[dict]) -> dict:
2     """Analyze historical runs to find scheduling bottlenecks."""
3     from collections import Counter
4     hour_counts = Counter()
5     for run in run_history:
6         start_hour = run["start_time"].hour
7         hour_counts[start_hour] += 1
8     peak_hour = hour_counts.most_common(1)[0]
9     return {
10         "peak_hour": peak_hour[0],
```

```
11     "peak_concurrent": peak_hour[1],
12     "recommendation": "Stagger pipelines across 01:00-05:00 UTC"
13         if peak_hour[1] > 10 else "Current schedule is acceptable",
14 }
```

Listing 14.8: Schedule optimization analysis

Warning

Shifting a pipeline’s schedule may break downstream consumers expecting fresh data by a certain hour. Always validate schedule changes against downstream SLA requirements. Document the dependency chain so that schedule changes propagate correctly. The phrase “I moved the pipeline to 4 AM and didn’t tell anyone” has ended careers.

14.10 Migration Between Orchestrators

CLAUDE CODE can translate between frameworks. The key mappings: Airflow DAG → Dagster Job/Asset Group, Operator → Op/Asset, XCom → IO Manager, Connection → Resource, Sensor → Freshness Policy.

Migration is not just code translation—it requires rethinking how you model your pipelines. Airflow thinks in tasks and dependencies; Dagster thinks in assets and their relationships. CLAUDE CODE can analyze your Airflow DAGs and recommend which Dagster patterns (assets, ops, or a hybrid) best fit each pipeline. It approaches this with the enthusiasm of someone who does not have to sit through the six-month migration that follows.

Tip

When migrating, CLAUDE CODE recommends a parallel-run strategy: keep both old and new pipelines running, compare outputs, and decommission the old pipeline only after two weeks of consistent agreement. This approach catches subtle differences in scheduling, retry behavior, and error handling. It also doubles your infrastructure costs for two weeks, which is why you should not mention it to Finance until after it is done.

14.11 Putting It All Together

1. **Start with asset-based thinking.** Model pipelines as data products with clear inputs and outputs. Even if you use Airflow, think about what data each task produces.
2. **Parameterize everything.** Use DAG factories, YAML configs, or metadata tables. Never hardcode connection strings, table names, or schedules in DAG code.
3. **Layer error handling.** Combine native retries with error classification and intelligent alerting. Not every failure deserves a page. Most failures deserve a page even less than you think.
4. **Instrument aggressively.** Export metrics to Prometheus; the orchestrator UI is for development, not production monitoring.
5. **Optimize schedules.** Analyze historical run data to spread load and meet SLAs.
6. **Plan for migration.** Write modular logic loosely coupled to your orchestrator's API. Business logic should live in pure Python functions, not in operator callbacks. Future You will either thank you or inherit your technical debt. Be kind to Future You.

Note

Orchestrator Lock-In The biggest risk in orchestration is coupling business logic to framework-specific APIs. Keep your core transformation logic in plain Python modules that the orchestrator calls but does not define. This makes migration between frameworks a rewiring exercise rather than a rewrite. It also means CLAUDE CODE can rewrite the orchestration layer without touching your business logic, which is the kind of modularity that brings a tear to an architect's eye.

14.12 Exercises

1. **Multi-Framework Implementation:** Take a simple ELT pipeline (extract from an API, load to a database, run a transformation) and implement it in all three frameworks: Airflow, Dagster, and Prefect. Compare the code volume, testing experience, and debugging workflow. Use CLAUDE CODE to generate

all three, then argue with your team about which one is “best” for the next two sprints.

2. **Error Handling Drill:** Create a pipeline with intentional failure points (a task that fails 30% of the time, a task with a 10-second timeout that occasionally takes 15 seconds). Implement the error classification and retry logic from this chapter. Measure how many runs complete successfully with and without intelligent retry logic. Contemplate how this mirrors life.
3. **Dynamic DAG Factory:** Build a YAML-driven DAG factory that generates Airflow DAGs from configuration files. Start with three similar pipelines, then add a fourth by writing only a YAML file. Verify that all DAGs serialize correctly. Experience the quiet satisfaction of never writing another DAG by hand.
4. **Monitoring Dashboard:** Set up Prometheus metrics for your orchestrator and build a Grafana dashboard with the four essential panels described in Section 14.7. Simulate failures and verify that alerts fire correctly. Use CLAUDE CODE to also generate the runbook for when those alerts fire, because alerts without runbooks are just expensive anxiety.
5. **Orchestrator Migration:** Take an existing Airflow DAG (from this chapter or your own) and migrate it to Dagster using CLAUDE CODE. Document every decision point where the migration required rethinking the pipeline structure, not just translating syntax. This exercise takes approximately three times longer than you estimate.
6. **Scheduling Optimization:** Collect (or simulate) two weeks of historical run data for at least 10 pipelines. Analyze the data to identify scheduling conflicts, peak concurrency windows, and SLA risks. Use CLAUDE CODE to recommend an optimized schedule and validate that no downstream SLAs are violated. Present the results to your team and watch them ignore the recommendations for six months before implementing them during an outage.

15

Infrastructure as Code with Claude Code

“Sure, you could write this Terraform yourself. You could also walk to work instead of driving. Both are technically possible and equally inadvisable.”

Infrastructure as Code (IaC) expresses infrastructure in declarative or imperative configuration files that live in version control alongside application code. In this chapter we explore how CLAUDE CODE transforms the IaC workflow—from authoring Terraform modules and CloudFormation templates to reviewing plans, estimating costs, detecting drift, and hardening security posture—across the three dominant IaC frameworks: Terraform, AWS CloudFormation, and Pulumi. Because clicking through the AWS console is not “infrastructure as code,” no matter how many screenshots you paste into Confluence.

15.1 IaC Fundamentals

Every IaC framework addresses three concerns:

1. **Desired-state declaration:** You describe *what* you want, and the framework determines *how* to reach that state. This is a beautiful abstraction right up until the framework decides to destroy and recreate your production database to “reach that state.”

2. **State management:** The framework tracks which resources it has already created so that subsequent runs apply only the delta.
3. **Execution planning:** Before making changes, the framework shows you a plan of creates, updates, and destroys. Reading this plan carefully is the difference between a good day and an incident report.

Note

If you are new to IaC, think of it as the infrastructure equivalent of a database migration framework like Alembic or Flyway. Just as those tools track schema state and apply diffs, IaC tools track cloud resource state and apply diffs. Except when they do not, and then you learn about “state drift” in the most educational way possible.

Terraform (HashiCorp) uses HCL and is cloud-agnostic. **AWS CloudFormation** is AWS-native using JSON or YAML. **Pulumi** uses general-purpose programming languages (Python, TypeScript, Go) instead of DSLs.

Data platforms are infrastructure-heavy—object storage, data warehouses, compute clusters, networking, IAM, and monitoring. Managing all of this manually is untenable at scale. IaC brings reproducibility (recreate your entire platform from scratch), auditability (every change is a git commit), and speed (provision a new environment in minutes rather than days). It also brings the terrifying ability to delete your entire platform from scratch, which is why we have code review.

Warning

State File Security Terraform and Pulumi store sensitive information in their state files: database passwords, API keys, and resource ARNs. Never commit state files to version control. Use remote state backends (S3 + DynamoDB for Terraform, Pulumi Cloud for Pulumi) with encryption at rest and strict access controls. A state file in a public GitHub repo is not a hypothetical risk—it is a Tuesday on r/DevOps.

Tip

Start with Modules, Not Monoliths Structure your IaC as reusable modules from day one. A monolithic Terraform file with 500 resources becomes un-maintainable. Organize by concern: networking, storage, compute, IAM, monitoring. CLAUDE CODE can refactor a monolith into modules if you describe

the desired boundaries. It can also generate the modules from scratch, which is faster than refactoring and does not require apologizing to the person who wrote the monolith.

15.2 Using Claude Code to Generate Terraform Modules

When asking CLAUDE CODE to generate a Terraform module, specificity pays dividends. A prompt like “Create a Terraform module for an S3 bucket” gives almost no context and produces the infrastructure equivalent of a participation trophy. Instead, specify encryption, versioning, lifecycle rules, access logging, naming conventions, and tags.

A Complete S3 Data Lake Module

```
1 variable "bucket_name" {
2   description = "Name of the S3 bucket for the data lake"
3   type       = string
4   validation {
5     condition   = can(regex("^[a-z0-9][a-z0-9.-]{1,61}[a-z0-9]$",
6       var.bucket_name))
7     error_message = "Bucket name must be 3-63 chars, lowercase."
8   }
9 }
10
11 variable "environment" {
12   description = "Deployment environment"
13   type       = string
14   default    = "production"
15   validation {
16     condition   = contains(["development", "staging", "production"],
17       var.environment)
18     error_message = "Must be development, staging, or production."
19   }
20 }
21
22 variable "cost_center" {
23   description = "Cost center for billing allocation"
24   type       = string
25 }
26
```

```
27 variable "vpc_endpoint_id" {
28   description = "VPC endpoint ID for bucket policy restriction"
29   type       = string
30 }
31
32 variable "logging_bucket_id" {
33   description = "ID of the S3 bucket for access logging"
34   type       = string
35 }
```

Listing 15.1: Terraform module: S3 data lake bucket (variables.tf)

```
1  locals {
2    common_tags = {
3      Environment = var.environment
4      CostCenter  = var.cost_center
5      ManagedBy   = "terraform"
6    }
7  }
8
9  resource "aws_kms_key" "data_lake" {
10   description          = "KMS key for ${var.bucket_name}"
11   deletion_window_in_days = 30
12   enable_key_rotation  = true
13   tags                 = local.common_tags
14 }
15
16 resource "aws_s3_bucket" "data_lake" {
17   bucket = var.bucket_name
18   tags   = local.common_tags
19 }
20
21 resource "aws_s3_bucket_versioning" "data_lake" {
22   bucket = aws_s3_bucket.data_lake.id
23   versioning_configuration { status = "Enabled" }
24 }
25
26 resource "aws_s3_bucket_server_side_encryption_configuration" "data_lake"
27   " {
28   bucket = aws_s3_bucket.data_lake.id
29   rule {
30     apply_server_side_encryption_by_default {
31       sse_algorithm = "aws:kms"
32       kms_master_key_id = aws_kms_key.data_lake.arn
```

```
32     }
33     bucket_key_enabled = true
34   }
35 }
36
37 resource "aws_s3_bucket_lifecycle_configuration" "data_lake" {
38   bucket = aws_s3_bucket.data_lake.id
39   rule {
40     id      = "intelligent-tiering"
41     status = "Enabled"
42     transition { days = 30; storage_class = "STANDARD_IA" }
43     transition { days = 90; storage_class = "GLACIER" }
44     noncurrent_version_expiration { noncurrent_days = 365 }
45   }
46 }
47
48 resource "aws_s3_bucket_public_access_block" "data_lake" {
49   bucket                = aws_s3_bucket.data_lake.id
50   block_public_acls     = true
51   block_public_policy   = true
52   ignore_public_acls   = true
53   restrict_public_buckets = true
54 }
55
56 resource "aws_s3_bucket_policy" "data_lake" {
57   bucket = aws_s3_bucket.data_lake.id
58   policy = jsonencode({
59     Version = "2012-10-17"
60     Statement = [{
61       Sid      = "DenyNonVPCEndpoint"
62       Effect   = "Deny"
63       Principal = "*"
64       Action   = "s3:*"
65       Resource = [aws_s3_bucket.data_lake.arn,
66                 "${aws_s3_bucket.data_lake.arn}/*"]
67       Condition = {
68         StringNotEquals = { "aws:sourceVpce" = var.vpc_endpoint_id }
69       }
70     }]
71   })
72 }
```

Listing 15.2: Terraform module: S3 data lake bucket (main.tf)

Tip

Always ask CLAUDE CODE to include validation blocks in your variables. These catch configuration errors at `terraform plan` time rather than during `apply`, when the feedback loop is much slower and the blast radius is larger. The difference between “invalid bucket name” at plan time and “invalid bucket name” during a production deploy is the difference between a typo and an incident.

Outputs and Documentation

A well-designed module exports outputs that downstream modules consume. CLAUDE CODE generates `outputs.tf` with descriptions, and can also produce a `README.md` with usage examples, input/output tables, and prerequisite documentation. It will write better documentation than you would, and it will not complain about writing it, which makes it better than approximately 100% of engineers.

```
1  output "bucket_id" {
2    description = "The ID of the data lake S3 bucket"
3    value      = aws_s3_bucket.data_lake.id
4  }
5
6  output "bucket_arn" {
7    description = "The ARN of the data lake S3 bucket"
8    value      = aws_s3_bucket.data_lake.arn
9  }
10
11 output "kms_key_arn" {
12   description = "The ARN of the KMS key used for encryption"
13   value      = aws_kms_key.data_lake.arn
14 }
15
16 output "kms_key_id" {
17   description = "The ID of the KMS key used for encryption"
18   value      = aws_kms_key.data_lake.id
19 }
```

Listing 15.3: Terraform module: S3 data lake bucket (`outputs.tf`)

15.3 Claude Code for Terraform Plan Review and Risk Assessment

Running `terraform plan` produces a textual diff of intended changes. For large refactors the plan can be hundreds of lines long. CLAUDE CODE can parse this output and produce a structured risk assessment in seconds. It reads Terraform plans faster than you read Slack messages, and with considerably more attention to detail.

```

1 terraform plan -out=plan.bin
2 terraform show plan.bin -no-color > plan.txt
3
4 claude --print "Review this Terraform plan for risks. Flag any
5   destructive changes, security regressions, cost implications,
6   and potential downtime. Rate overall risk as LOW, MEDIUM, or HIGH.
7   $(cat plan.txt)"

```

Listing 15.4: Capturing and reviewing a Terraform plan

When reviewing a plan, CLAUDE CODE evaluates destructive operations (resources being destroyed or recreated), security regressions (wide CIDR ranges, * actions in IAM), cost implications (instance size changes), ordering risks, and state drift indicators.

Tip

For CI/CD integration, pipe the plan output to CLAUDE CODE as part of your pull request workflow. CLAUDE CODE can post its risk assessment as a PR comment, giving reviewers a structured summary before they read the raw plan diff. This is the only PR comment format that reviewers actually read.

Warning

Destructive Plan Changes Pay special attention when `terraform plan` shows `destroy` followed by `create` for a resource. This often means the resource will be recreated, which can cause downtime. Common triggers include changing a resource name, modifying an immutable attribute, or moving a resource between modules. The word “destroy” in a Terraform plan should trigger the same physiological response as a fire alarm.

15.4 Pulumi and Programmatic IaC with Claude Code

Pulumi uses general-purpose programming languages, so data engineers can define infrastructure in the same Python they use for pipelines. This eliminates the context switch between HCL and Python. It also means you can introduce bugs in your infrastructure using the same language you introduce bugs in your application code. Consistency is a virtue.

```
1 import pulumi
2 import pulumi_aws as aws
3 from pulumi import Config
4
5 config = Config()
6 environment = config.require("environment")
7
8 common_tags = {"Environment": environment, "ManagedBy": "pulumi"}
9
10 data_lake_key = aws.kms.Key("data-lake-key",
11     description=f"KMS key for {environment} data lake",
12     enable_key_rotation=True, tags=common_tags)
13
14 zones = ["raw", "staging", "curated"]
15 buckets = {}
16
17 for zone in zones:
18     bucket = aws.s3.BucketV2(f"data-lake-{zone}",
19         bucket=f"company-data-lake-{zone}-{environment}",
20         tags={**common_tags, "Zone": zone})
21
22     aws.s3.BucketVersioningV2(f"data-lake-{zone}-versioning",
23         bucket=bucket.id,
24         versioning_configuration={"status": "Enabled"})
25
26     aws.s3.BucketServerSideEncryptionConfigurationV2(
27         f"data-lake-{zone}-encryption", bucket=bucket.id,
28         rules=[{"apply_server_side_encryption_by_default": {
29             "sse_algorithm": "aws:kms",
30             "kms_master_key_id": data_lake_key.arn,
31         }}, {"bucket_key_enabled": True}])
32
33     aws.s3.BucketPublicAccessBlock(f"data-lake-{zone}-public-block",
34         bucket=bucket.id, block_public_acls=True,
35         block_public_policy=True, ignore_public_acls=True,
36         restrict_public_buckets=True)
```

```

37     buckets[zone] = bucket
38
39
40 for zone, bucket in buckets.items():
41     pulumi.export(f"{zone}_bucket_name", bucket.bucket)
42     pulumi.export(f"{zone}_bucket_arn", bucket.arn)
43
44 pulumi.export("kms_key_arn", data_lake_key.arn)

```

Listing 15.5: Pulumi: Data lake infrastructure in Python

Note

Notice how the Python `for` loop eliminates the repetition you would see in Terraform. This is Pulumi's core advantage: general-purpose language constructs applied to infrastructure. Conditionals, loops, functions, and classes all work as expected. The disadvantage is that you can now write infinitely complex infrastructure code, which some engineers interpret as an invitation rather than a warning.

15.5 Security Review of Infrastructure Code

Misconfigured infrastructure is a leading cause of data breaches. CLAUDE CODE can serve as a systematic security reviewer, evaluating encryption at rest, encryption in transit, network exposure, IAM least privilege, public access, logging, backup, and secrets management. It reviews infrastructure code with the thoroughness of someone who has read every AWS security blog post and the paranoia of someone who has lived through a breach.

```

1  # BEFORE -- overly permissive (Claude Code flags as HIGH severity)
2  resource "aws_iam_role_policy" "glue_job" {
3      role = aws_iam_role.glue_job.id
4      policy = jsonencode({
5          Version = "2012-10-17"
6          Statement = [{
7              Effect = "Allow"; Action = "s3:*"; Resource = "*"
8          }]
9      })
10 }
11
12 # AFTER -- least privilege (Claude Code's recommended fix)

```

```

13 resource "aws_iam_role_policy" "glue_job" {
14   role = aws_iam_role.glue_job.id
15   policy = jsonencode({
16     Version = "2012-10-17"
17     Statement = [
18       {
19         Sid      = "ReadSourceBucket"
20         Effect   = "Allow"
21         Action   = ["s3:GetObject", "s3:ListBucket"]
22         Resource = [aws_s3_bucket.raw.arn,
23                   "${aws_s3_bucket.raw.arn}/events/*"]
24       },
25       {
26         Sid      = "WriteTargetBucket"
27         Effect   = "Allow"
28         Action   = ["s3:PutObject", "s3:DeleteObject"]
29         Resource = ["${aws_s3_bucket.curated.arn}/analytics/*"]
30       },
31       {
32         Sid      = "UseEncryptionKey"
33         Effect   = "Allow"
34         Action   = ["kms:Decrypt", "kms:GenerateDataKey"]
35         Resource = [aws_kms_key.data_lake.arn]
36       }
37     ]
38   })
39 }

```

Listing 15.6: Before and after: Tightening a Glue job IAM policy

Warning

Overly permissive Glue job roles. Data engineers often grant Glue jobs `s3:*` on `*` to “get things working.” This is one of the most common security anti-patterns in data platforms. It is also the most common thing you will find when you inherit someone else’s infrastructure. Always scope S3 permissions to specific bucket ARNs and prefixes, and separate read permissions from write permissions. “I’ll tighten the permissions later” is the data engineering equivalent of “I’ll start going to the gym on Monday.”

Automated Security Scanning in CI

Integrate CLAUDE CODE-powered security reviews into your CI pipeline alongside tools like `checkov`, `tfsec`, and `terrascan`. Static analysis tools catch known patterns (open security groups, unencrypted buckets), while CLAUDE CODE catches contextual issues (a policy that is technically valid but overly broad for its stated purpose).

Note

Defense in Depth No single security tool catches everything. Use a layered approach: `tfsec` for known misconfigurations, `checkov` for compliance frameworks (CIS, SOC2), and CLAUDE CODE for contextual review. The overlap between tools is a feature, not a bug—redundant detection reduces the chance of a security gap reaching production. Think of it as wearing a belt and suspenders. Both at once is not paranoid; it is prudent.

Tip

Security Review Checklist When asking CLAUDE CODE to review infrastructure code, provide context: “This Glue job reads from the raw zone and writes to the curated zone. Review the IAM policy for least privilege.” Without context, CLAUDE CODE cannot assess whether permissions are appropriate for the workload. It is very good at security review, but it cannot read your mind. Yet.

15.6 Drift Detection and Cost Estimation

Configuration drift occurs when actual infrastructure diverges from the desired state in code. Terraform detects drift automatically when you run `terraform plan`. CLAUDE CODE can analyze unexpected plan changes and recommend whether to adopt the drift or revert to desired state.

Data platforms are particularly susceptible to drift: someone opens a port “temporarily” for debugging, a developer adds an inline IAM policy through the console, or auto-scaling adjusts instance counts. Each of these creates a discrepancy between what your code says and what actually exists. The word “temporarily” in infrastructure has a half-life of approximately infinity.

Note

To reduce drift, enforce a policy that *all* infrastructure changes go through Terraform, even emergency fixes. Use SCP (Service Control Policies) to restrict console access to read-only for production accounts. When emergency console changes are unavoidable, import them into Terraform state within 24 hours. Treat every console change as a small fire that needs to be put out before it spreads.

For cost estimation, *Infracost* estimates cloud costs from Terraform code. CLAUDE CODE can interpret Infracost output and suggest optimizations: spot instances for Spark executors (60–70% savings), S3 Intelligent Tiering, reserved capacity for always-on resources, right-sizing clusters, and VPC endpoints to reduce NAT gateway costs. Your cloud bill will thank you. Your cloud provider will not.

Table 15.1: Common Data Platform Cost Optimizations (a.k.a. “Where Did All the Money Go?”)

Resource	Optimization	Typical Savings
EC2/EMR workers	Spot instances	60–70%
S3 storage	Intelligent Tiering	20–40%
Redshift/BigQuery	Reserved capacity	30–50%
NAT Gateway	VPC endpoints for S3/DynamoDB	50–80%
EBS volumes	GP3 instead of GP2	15–20%

Tip

Cost Tagging Strategy Tag every resource with `CostCenter`, `Team`, `Environment`, and `Pipeline`. Without tags, you cannot attribute costs to specific teams or workloads. CLAUDE CODE can enforce tagging by generating `default_tags` in your provider configuration. Untagged resources are the infrastructure equivalent of unmarked bills—impossible to trace and always suspicious.

```

1 infracost breakdown --path . --format json > cost.json
2
3 claude --print "Analyze this Infracost report and suggest
4   optimizations that could reduce monthly costs by at least 20%
5   without impacting data pipeline reliability.
6   $(cat cost.json)"

```

Listing 15.7: Integrating Infracost with Claude Code review

15.7 Exercises

1. **Module Generation.** Ask CLAUDE CODE to generate a Terraform module for an Amazon Redshift Serverless namespace with customer-managed KMS encryption, VPC connectivity, a usage limit, and IAM authentication. Validate the output for completeness and security best practices. Run `terraform validate` and fix any issues. If CLAUDE CODE produces valid Terraform on the first try, consider what that means for your career trajectory.
2. **Plan Review Drill.** Take the Terraform plan output from a recent infrastructure change (or generate a synthetic one). Feed it to CLAUDE CODE and ask for a risk assessment. Compare findings against your own manual review. Note any risks that CLAUDE CODE caught that you missed, and vice versa. This exercise is humbling in both directions.
3. **Security Hardening.** Start with an intentionally insecure Terraform configuration (unencrypted S3 bucket, open security group, public RDS instance, wildcard IAM policy, disabled logging, no versioning). Ask CLAUDE CODE to identify and fix all security issues. A thorough review should find at least 10 distinct issues. If it finds fewer, your “insecure” configuration was not insecure enough, which is an unusual problem to have.
4. **Cost Optimization Challenge.** Provide CLAUDE CODE with a complete data lake infrastructure (S3 buckets, Glue jobs, Redshift cluster, EMR cluster, NAT gateways) and ask it to find ways to reduce monthly costs by at least 30% without sacrificing data durability or query performance. Implement the top three recommendations. Present the savings to management and bask in the brief, warm glow of executive approval.
5. **Pulumi Translation.** Take the Terraform S3 module from this chapter and translate it to Pulumi Python using CLAUDE CODE. Compare the two implementations in terms of readability, testability, and the ability to express conditional logic. Form strong opinions about which is better. Change those opinions within six months.
6. **Drift Detection Simulation.** Deploy a simple infrastructure (S3 bucket + IAM role) via Terraform. Make three manual changes through the AWS

console (add a tag, modify the bucket policy, attach an inline IAM policy). Run `terraform plan`, feed the output to CLAUDE CODE, and evaluate its recommendations for each drift item. Decide whether to adopt or revert each change. Practice saying “who made this console change?” in an appropriately disappointed tone.

16

DataOps and CI/CD

*Modern data engineering demands the same rigor in deployment, testing, and operational excellence that software engineering has cultivated over decades. **DataOps** brings DevOps principles into the data world, combining agile methodology, continuous integration, and continuous delivery to accelerate data pipeline development while maintaining quality and reliability. In this chapter, we explore how CLAUDE CODE can serve as a force multiplier across every stage of the DataOps lifecycle.*

16.1 DataOps Principles and Practices

DataOps rests on three pillars: *agile data development* (short feedback loops, frequent releases), *continuous integration and delivery for data* (every pipeline change passes through automated validation), and *monitoring and observability* (visibility into data quality, freshness, volume, and lineage).

The core insight of DataOps is that data pipelines are software. They deserve the same development practices: version control, code review, automated testing, staged deployments, and monitoring. Yet many data teams still deploy by running SQL scripts manually, skip testing because “it’s just a query,” and learn about failures from angry stakeholders rather than alerts.

Note

DataOps is often confused with *data governance*, but they are complementary. Data governance defines the policies and standards; DataOps provides the practices and tooling to enforce them continuously. Think of governance as the “what” and DataOps as the “how.”

Tip

Start your DataOps journey by mapping your current data pipeline deployment process end to end. Identify every manual step, every approval bottleneck, and every point where errors tend to occur. These are your highest-impact automation targets. Most teams find that 60–70% of deployment time is spent on manual validation that could be automated.

AI-assisted DataOps is about amplifying engineer effectiveness. CLAUDE CODE excels at the repetitive, pattern-heavy tasks that consume disproportionate time: writing tests, reviewing SQL for anti-patterns, generating migration scripts, drafting runbooks, and summarizing incident timelines. The goal is not to remove engineers from the loop but to reduce the toil that prevents them from doing higher-value work.

16.2 CI/CD for Data Pipelines

Data pipeline CI/CD differs from traditional software CI/CD in important ways: pipelines are stateful (they read from and write to persistent stores), schema changes can be destructive (dropping a column destroys data), and testing requires realistic data (synthetic data may not exercise important edge cases). A well-designed system includes linting and static analysis, unit testing, integration testing, data quality validation, performance benchmarking, deployment, and post-deployment validation.

Warning

Never skip the integration testing stage, even when changes appear trivial. A single-character typo in a SQL JOIN condition can silently produce incorrect results that propagate through downstream systems for days before anyone notices.

Each branch should map to a corresponding database schema or environment. When

a developer creates a feature branch, CI should automatically provision a sandboxed environment with isolated schemas, seed data, and the ability to run the full pipeline end to end.

Environment Isolation Strategies

There are three common approaches to environment isolation for data pipelines:

1. **Schema-per-branch:** Each branch gets a dedicated schema (e.g., `dev_feature_123`) in a shared database. Simple to implement, but no compute isolation.
2. **Database-per-branch:** Each branch gets a full database clone. Better isolation, but expensive for large databases.
3. **Zero-copy clones:** Supported by Snowflake, Databricks, and some PostgreSQL extensions. Full isolation with minimal storage overhead.

CLAUDE CODE can generate the provisioning and teardown scripts for any of these strategies, including the CI/CD hooks that create environments on branch creation and destroy them on merge.

16.3 Claude Code for Pull Request Reviews of Data Code

CLAUDE CODE is particularly effective at identifying join correctness issues, NULL handling errors, aggregation mistakes, type coercion problems, performance anti-patterns, schema compatibility issues, and business logic errors in data code. These are precisely the issues that are hardest to catch in manual code review because they require holding the entire data model in your head.

Code Example

```
1 import anthropic
2 import subprocess
3 import json
4 from pathlib import Path
5
6
7 def review_data_code(filepath: str) -> dict:
8     """Perform a Claude Code-powered review of a data file."""
9     client = anthropic.Anthropic()
```

```
10
11     diff = subprocess.run(
12         ["git", "diff", "main", "--", filepath],
13         capture_output=True, text=True, check=True
14     ).stdout
15
16     full_content = Path(filepath).read_text()
17
18     review_prompt = f"""You are a senior data engineer reviewing
19 a pull request. Review the following data code change.
20
21 ## Changed File: {filepath}
22
23 ### Full file content:
24 ```sql
25 {full_content}
26 ```
27
28 ### Diff:
29 ```diff
30 {diff}
31 ```
32
33 ## Review Checklist:
34 1. Join correctness 2. NULL handling
35 3. Performance      4. Data quality tests
36 5. Schema compat    6. Business logic
37 7. Incremental logic 8. Documentation
38
39 Provide your review as JSON with: "summary", "severity"
40 (approve/request_changes/comment), "findings" (list with
41 line, severity, category, message), "suggestions"."""
42
43     response = client.messages.create(
44         model="claude-sonnet-4-20250514", max_tokens=4096,
45         messages=[{"role": "user", "content": review_prompt}]
46     )
47
48     text = response.content[0].text
49     start = text.find("{")
50     end = text.rfind("}") + 1
51     return json.loads(text[start:end])
```

! Tip

Contextual Reviews For the best review quality, provide CLAUDE CODE with both the diff *and* the full file content. The diff shows what changed; the full file shows the context. Without context, CLAUDE CODE cannot assess whether a JOIN is correct or whether a WHERE clause is consistent with the rest of the query.

16.4 Automated Migration Generation and Validation

Database migrations are error-prone because they operate on live data with real consequences. A migration that works perfectly on an empty test database may lock a production table for hours. CLAUDE CODE can generate migration scripts that account for edge cases: default values for existing rows, index recreation, constraint validation, and rollback procedures.

Always validate migrations before execution—CLAUDE CODE can perform multi-point validation checking for data loss risk, reversibility, lock impact, idempotency, and transaction safety.

```

1 def validate_migration(forward_sql: str, rollback_sql: str) -> dict:
2     client = anthropic.Anthropic()
3     response = client.messages.create(
4         model="claude-sonnet-4-20250514", max_tokens=4096,
5         messages=[{"role": "user", "content": f"""Validate this
6 database migration for safety.
7
8 Forward migration:
9 ```sql
10 {forward_sql}
11 ```
12
13 Rollback migration:
14 ```sql
15 {rollback_sql}
16 ```
17
18 Check: 1. Data loss risk 2. Reversibility
19 3. Lock impact on large tables 4. Idempotency
20 5. Transaction safety 6. Index impact
21 Return JSON with: safe (bool), risks (list), suggestions (list)."""}]
22     )

```

```
23 text = response.content[0].text
24 start = text.find("{")
25 end = text.rfind("}") + 1
26 return json.loads(text[start:end])
```

Listing 16.1: Migration validation with Claude Code

Warning

Never execute CLAUDE CODE-generated migrations directly in production without human review and testing in a staging environment. Even well-validated migrations can have unexpected interactions with concurrent workloads, replication, and database-specific behaviors.

16.5 Testing Strategies for Data Pipelines

Data pipeline testing spans three levels: unit tests (individual transformation logic), integration tests (components working together), and end-to-end tests (full pipeline validation against realistic data).

Code Example

```
1 import pytest
2 import pandas as pd
3 from decimal import Decimal
4 from datetime import date
5
6
7 class TestRevenueCalculation:
8     """Unit tests for the daily revenue aggregation model."""
9
10    @pytest.fixture
11    def sample_transactions(self) -> pd.DataFrame:
12        return pd.DataFrame({
13            "transaction_id": [1, 2, 3, 4, 5, 6],
14            "customer_id": [101, 102, 101, 103, 102, 101],
15            "amount": [Decimal("99.99"), Decimal("149.50"),
16                    Decimal("29.99"), Decimal("0.00"),
17                    Decimal("199.99"), Decimal("-29.99")],
18            "status": ["completed", "completed", "completed",
```

```
19         "failed", "completed", "refunded"],
20         "transaction_date": [
21             date(2025, 1, 15), date(2025, 1, 15),
22             date(2025, 1, 15), date(2025, 1, 15),
23             date(2025, 1, 16), date(2025, 1, 16)],
24     })
25
26     def test_daily_revenue_excludes_failed(self, sample_transactions)
27     :
28         result = calculate_daily_revenue(sample_transactions)
29         jan_15 = result[result["date_day"] == date(2025, 1, 15)]
30         assert jan_15["total_revenue"].iloc[0] == Decimal("279.48")
31
32     def test_daily_revenue_handles_refunds(self, sample_transactions)
33     :
34         result = calculate_daily_revenue(sample_transactions)
35         jan_16 = result[result["date_day"] == date(2025, 1, 16)]
36         assert jan_16["total_revenue"].iloc[0] == Decimal("170.00")
37
38     def test_daily_revenue_no_transactions(self):
39         empty_df = pd.DataFrame(columns=[
40             "transaction_id", "customer_id", "amount",
41             "status", "transaction_date"])
42         result = calculate_daily_revenue(empty_df)
43         assert len(result) == 0
```

Integration tests validate that pipeline components work together, running against a real database with controlled test data. Use `testcontainers` to spin up temporary PostgreSQL containers for isolated testing.

! Tip

Use CLAUDE CODE to generate comprehensive test fixtures. Describe your schema and business rules, and ask CLAUDE CODE to generate test data covering edge cases: NULL values, boundary conditions, Unicode characters, time-zone edge cases, extreme values, and the specific scenarios that have caused production incidents in the past.

Note

The Testing Pyramid for Data Apply the testing pyramid to data pipelines: many unit tests (fast, cheap, isolated), fewer integration tests (slower, test component interaction), and a handful of end-to-end tests (slowest, test the full pipeline). CLAUDE CODE can generate tests at all three levels from a single description of your pipeline's expected behavior.

16.6 Deployment Validation and Incident Response

Every deployment should have a pre-deployment risk assessment, a corresponding rollback plan, and post-deployment validation. CLAUDE CODE can generate rollback scripts that account for the specific changes being deployed.

Post-deployment validation goes beyond “did it run?” to “did it produce correct results?” Common validation checks include row count comparisons (new run vs. previous run), distribution checks (no unexpected NULL spikes), freshness verification (data is as recent as expected), and cross-system consistency (totals match between source and target).

```

1 def validate_deployment(table: str, conn) -> list[str]:
2     """Run post-deployment validation and return any failures."""
3     failures = []
4     # Row count sanity check
5     current = conn.execute(f"SELECT COUNT(*) FROM {table}").scalar()
6     previous = conn.execute(
7         f"SELECT row_count FROM pipeline_metadata "
8         f"WHERE table_name = '{table}' ORDER BY run_at DESC LIMIT 1"
9     ).scalar()
10    if previous and current < previous * 0.9:
11        failures.append(
12            f"Row count dropped >10%: {previous} -> {current}")
13    # NULL spike detection
14    for col in get_required_columns(table):
15        null_pct = conn.execute(
16            f"SELECT COUNT(*) FILTER (WHERE {col} IS NULL) * 100.0 "
17            f"/ COUNT(*) FROM {table}").scalar()
18        if null_pct > 5:
19            failures.append(f"NULL spike in {col}: {null_pct:.1f}%")
20    return failures

```

Listing 16.2: Post-deployment validation checks

When data incidents occur, CLAUDE CODE can accelerate every phase: detection, diagnosis, resolution, and postmortem. It can analyze monitoring alerts, classify severity, identify likely root causes, and draft blameless postmortem documents from incident timelines.

Tip

Store postmortems in a version-controlled repository alongside your pipeline code. Over time, CLAUDE CODE can analyze your postmortem history to identify recurring patterns and suggest systemic improvements. Common patterns include missing monitoring for newly deployed pipelines and inadequate testing of schema changes.

Warning

Rollback Is Not Always Safe For pipelines that write to external systems (reverse ETL, API pushes, email triggers), rollback may not undo side effects. A rolled-back pipeline does not unsend emails or un-update CRM records. For these pipelines, deploy with extra caution and consider a shadow-mode testing phase.

16.7 Full Example: GitHub Actions Pipeline

The following GitHub Actions workflow demonstrates a complete CI/CD pipeline for a dbt-based data project. It includes linting, AI-powered review, integration testing, and gated production deployment.

Code Example

```
1 # .github/workflows/data-pipeline-ci.yaml
2 name: Data Pipeline CI/CD
3
4 on:
5   pull_request:
6     branches: [main, staging]
7     paths: ['models/**', 'macros/**', 'tests/**']
8
9 jobs:
10  lint:
```

```
11 name: SQL Linting
12 runs-on: ubuntu-latest
13 steps:
14   - uses: actions/checkout@v4
15   - run: pip install sqlfluff sqlfluff-templater-dbt
16   - run: sqlfluff lint models/ --dialect postgres
17       --format github-annotation-native
18
19 claude-review:
20 name: Claude Code AI Review
21 runs-on: ubuntu-latest
22 needs: lint
23 permissions: { pull-requests: write, contents: read }
24 steps:
25   - uses: actions/checkout@v4
26     with: { fetch-depth: 0 }
27   - uses: actions/setup-python@v5
28     with: { python-version: '3.11' }
29   - run: pip install anthropic pyyaml
30   - id: changed-files
31     run: |
32       FILES=$(git diff --name-only origin/${{ github.base_ref }}
33         -- '*.sql' '*.yaml' | tr '\n' ',')
34       echo "files=${FILES}" >> "$GITHUB_OUTPUT"
35   - if: steps.changed-files.outputs.files != ''
36     env:
37       ANTHROPIC_API_KEY: ${ secrets.ANTHROPIC_API_KEY }
38     run: python ci/scripts/claude_review.py
39         --files "${{ steps.changed-files.outputs.files }}"
40         --pr-number ${ github.event.pull_request.number }
41
42 integration-tests:
43 name: Integration Tests
44 runs-on: ubuntu-latest
45 needs: [lint, claude-review]
46 services:
47   postgres:
48     image: postgres:15
49     env: { POSTGRES_USER: dbt_test,
50           POSTGRES_PASSWORD: dbt_test,
51           POSTGRES_DB: dbt_test }
52     ports: ['5432:5432']
53 steps:
54   - uses: actions/checkout@v4
```

```
55     - run: pip install dbt-postgres==1.7.*
56     - run: dbt seed --target ci
57     - run: dbt run --target ci
58     - run: dbt test --target ci
59
60   deploy-production:
61     name: Deploy to Production
62     runs-on: ubuntu-latest
63     needs: integration-tests
64     if: github.base_ref == 'main'
65     environment: production
66     steps:
67       - uses: actions/checkout@v4
68       - run: pip install dbt-postgres==1.7.*
69       - run: python ci/scripts/create_snapshot.py
70         --target production --version ${ github.sha }
71       - run: dbt run --target production --fail-fast
72       - run: python ci/scripts/post_deploy_validation.py
73         --target production --version ${ github.sha }
```

Each stage serves a specific purpose: linting catches style and syntax issues early, the AI review catches semantic issues that linters miss, integration tests verify correctness against real data, and the production deploy includes both pre-deployment snapshots and post-deployment validation.

16.8 Exercises

Exercises

- CI/CD Pipeline Design** Design a CI/CD pipeline for a dbt project that includes at least five stages. Write the GitHub Actions YAML file and explain the purpose of each stage. Include at least one CLAUDE CODE-powered step. Test the pipeline against a sample dbt project with intentional issues.
- Migration Safety Analysis** You need to add a NOT NULL constraint to an existing column that contains some NULL values. Write the migration script (both up and down), including the data backfill step. Use CLAUDE CODE to validate its safety and identify edge cases you may have missed.
- Incident Response Runbook** Create a comprehensive incident response runbook for a data pipeline failure, including detection criteria, severity classifi-

cation, escalation paths, diagnostic procedures, resolution steps, communication templates, and a postmortem template. Use CLAUDE CODE to draft the runbook and then refine it based on your team's specific infrastructure.

4. **Review Automation** Implement a simplified CLAUDE CODE-powered PR review system that accepts a SQL file, sends it to CLAUDE CODE for review, and outputs structured findings. Test it against at least three SQL files with intentional issues (incorrect JOIN, NULL mishandling, missing WHERE clause, performance anti-pattern).
5. **Test Generation Challenge** Take an existing data transformation (a dbt model or SQL query) and use CLAUDE CODE to generate a comprehensive test suite. The suite should include unit tests, integration tests, and edge case tests. Measure code coverage and identify any scenarios that CLAUDE CODE missed.
6. **Post-Deployment Monitoring** Implement the post-deployment validation function from Section 16.6 for a table of your choice. Define thresholds for row count deviation, NULL spike detection, and freshness. Deploy a change that intentionally violates one threshold and verify that the validation catches it.

Part **VI**

Security, Governance, and Ethics

17

Data Governance with Claude Code

Tip

Data governance is not a one-time project—it is a continuous discipline. Claude Code can accelerate every facet of governance, from cataloging metadata to enforcing compliance policies, but human oversight remains essential for setting organizational priorities and interpreting regulatory nuance.

Modern data engineering teams operate under increasing pressure to treat data as a first-class organizational asset. Regulations such as GDPR, CCPA, and HIPAA impose strict requirements on how data is collected, stored, processed, and shared. Data governance provides the framework that bridges these demands—and Claude Code can serve as a tireless governance co-pilot.

This chapter covers four pillars of AI-augmented governance: automated data cataloging, lineage tracing and PII detection, compliance automation, and governance-as-code policy enforcement. Each section provides production-ready code and patterns you can adapt to your environment.

17.1 The Governance Challenge

Most organizations struggle with governance not because they lack policies, but because they lack the capacity to implement and enforce them consistently. A data team managing 500 tables, 200 dbt models, and 50 pipelines simply does not have

the bandwidth to manually classify every column, document every transformation, and audit every access pattern.

This is precisely where Claude Code changes the equation. The mechanical work of governance—scanning schemas, classifying columns, generating documentation, checking compliance—can be automated. The judgment work—setting policy, interpreting regulations, making trade-off decisions—remains with humans.

Note

Governance Maturity Organizations typically progress through governance maturity levels: (1) No formal governance—data is a free-for-all; (2) Reactive governance—policies exist but are enforced after incidents; (3) Proactive governance—policies are embedded in pipelines and CI/CD; (4) Automated governance—AI-augmented continuous monitoring and enforcement. This chapter helps you reach level 3–4.

The costs of poor governance are concrete:

- **Regulatory fines:** GDPR penalties can reach 4% of global annual turnover. CCPA violations carry fines of \$2,500–\$7,500 per violation.
- **Data breaches:** The average cost of a data breach exceeded \$4.5 million in 2023, with healthcare breaches averaging nearly \$11 million.
- **Lost productivity:** Data engineers spend an estimated 30–40% of their time on data quality issues, many of which stem from poor governance.
- **Decision risk:** Executives making decisions on ungoverned data face unreliable metrics, inconsistent definitions, and unreproducible results.

17.2 Automated Data Cataloging with Claude Code

A data catalog is the single source of truth for what data exists in your organization, who owns it, what it means, and how it should be handled. Claude Code can automate the heavy lifting by inspecting database schemas, reading source code, and generating human-readable descriptions.

```
1 import anthropic, psycopg2, json
2
3 def extract_schema_metadata(conn_string: str, schema: str = "public"):
4     conn = psycopg2.connect(conn_string)
```

```

5     cur = conn.cursor()
6     cur.execute("""
7         SELECT t.table_name, c.column_name, c.data_type,
8             c.is_nullable, c.column_default,
9             pgd.description AS column_comment
10        FROM information_schema.tables t
11        JOIN information_schema.columns c
12            ON t.table_name = c.table_name
13            AND t.table_schema = c.table_schema
14        LEFT JOIN pg_catalog.pg_statio_all_tables st
15            ON st.relname = t.table_name AND st.schemaname = t.
16        table_schema
17        LEFT JOIN pg_catalog.pg_description pgd
18            ON pgd.objoid = st.relid AND pgd.objsubid = c.
19        ordinal_position
20        WHERE t.table_schema = %s AND t.table_type = 'BASE TABLE'
21        ORDER BY t.table_name, c.ordinal_position
22        """, (schema,))
23     tables = {}
24     for row in cur.fetchall():
25         table_name = row[0]
26         if table_name not in tables:
27             tables[table_name] = {"columns": []}
28             tables[table_name]["columns"].append({
29                 "name": row[1], "data_type": row[2],
30                 "nullable": row[3], "default": row[4],
31                 "existing_comment": row[5]})
32     cur.close(); conn.close()
33     return tables
34
35 def generate_catalog_descriptions(tables: dict) -> dict:
36     client = anthropic.Anthropic()
37     catalog = {}
38     for table_name, meta in tables.items():
39         prompt = f"""Given the table schema, generate:
40         1. A concise business description (2-3 sentences).
41         2. A business description for each column.
42         3. Suggested data classification (public/internal/confidential/
43            restricted).
44         4. Suggested data owner role (e.g., "finance-team", "engineering").
45         5. Retention recommendation with justification.
46
47         Table: {table_name}
48         Columns: {json.dumps(meta['columns'], indent=2, default=str)}
49
50         Return valid JSON with keys: table_description, columns,

```

```

48 classification, suggested_owner, retention."""
49     response = client.messages.create(
50         model="claude-sonnet-4-20250514", max_tokens=2048,
51         messages=[{"role": "user", "content": prompt}])
52     catalog[table_name] = json.loads(response.content[0].text)
53     return catalog

```

Listing 17.1: Automated metadata extraction from PostgreSQL

Enriching Existing Catalogs

Most organizations already have a partial catalog—perhaps maintained in a spreadsheet, a wiki, or a commercial catalog tool. Claude Code can enrich these existing catalogs by identifying gaps (undocumented tables and columns), standardizing descriptions (making terminology consistent across domains), adding classifications (public, internal, confidential, restricted), and suggesting data owners based on table naming patterns and domain grouping.

```

1  def find_catalog_gaps(catalog: dict, schema_metadata: dict) -> dict:
2      """Compare existing catalog against live schema to find gaps."""
3      gaps = {
4          "undocumented_tables": [],
5          "undocumented_columns": [],
6          "stale_entries": [],
7          "missing_classifications": [],
8      }
9      # Tables in schema but not in catalog
10     for table in schema_metadata:
11         if table not in catalog:
12             gaps["undocumented_tables"].append(table)
13         else:
14             catalog_cols = {c["name"] for c in catalog[table].get("
15             columns", [])}
16             schema_cols = {c["name"] for c in schema_metadata[table]["
17             columns"]}
18             for col in schema_cols - catalog_cols:
19                 gaps["undocumented_columns"].append(f"{table}.{col}")
20                 if not catalog[table].get("classification"):
21                     gaps["missing_classifications"].append(table)
22     # Tables in catalog but not in schema (stale)
23     for table in catalog:
24         if table not in schema_metadata:
25             gaps["stale_entries"].append(table)

```

```
24 return gaps
```

Listing 17.2: Catalog gap analysis

Note

Treat your metadata catalog with the same rigor as production data. Version the schema with migrations, back it up regularly, and restrict write access to automated governance pipelines. A corrupt or outdated catalog is worse than no catalog because it creates false confidence.

Tip

Catalog Freshness Schedule catalog enrichment runs weekly. Compare the live schema against your catalog to detect new tables, dropped tables, and column changes. Claude Code can generate human-readable change summaries for review by data stewards.

17.3 Data Lineage Tracing and PII Detection

Data lineage answers: *Where did this data come from, and what happened to it?* Claude Code can parse SQL transformations, Airflow DAGs, and dbt models to reconstruct lineage graphs automatically.

SQL-Based Lineage Extraction

For teams using dbt, lineage is partially available through the manifest. But for custom SQL pipelines, stored procedures, and views, lineage must be extracted from the SQL itself.

```
1 def extract_lineage(sql: str, target_table: str) -> dict:
2     """Extract column-level lineage from a SQL transformation."""
3     client = anthropic.Anthropic()
4     message = client.messages.create(
5         model="claude-sonnet-4-20250514", max_tokens=4096,
6         system="""You are a data lineage expert. Given a SQL statement,
7         extract column-level lineage. For each output column, identify:
8         1. Source table(s) and column(s)
9         2. Transformation applied (direct copy, aggregation,
```

```

10     calculation, conditional, etc.)
11     3. Whether the transformation preserves PII status""",
12     messages=[{"role": "user",
13               "content": f"Target table: {target_table}\n\nSQL:\n{sql}"},
14             ]
15     return json.loads(message.content[0].text)

```

Listing 17.3: SQL lineage extraction with Claude Code

Hybrid PII Detection

For PII detection, combine deterministic regex-based scanning with Claude Code’s semantic analysis. Regex catches structured patterns (SSNs, emails, credit cards); Claude Code applies reasoning to detect context-dependent PII that regex misses—such as free-text fields containing addresses, columns named `notes` that contain patient information, or composite fields that become PII when combined.

```

1  import re, anthropic, json
2  from typing import List, Dict
3
4  PII_PATTERNS = {
5      "ssn": re.compile(r"\b\d{3}-\d{2}-\d{4}\b"),
6      "email": re.compile(
7          r"\b[A-Za-z0-9._%+-]+@[A-Za-z0-9.-]+\.[A-Z|a-z]{2,}\b"),
8      "credit_card": re.compile(r"\b(?:\d{4}[-\s]?){3}\d{4}\b"),
9      "phone_us": re.compile(
10         r"\b(?:\+1[-.\s]?|\(?\d{3}\)?[-.\s]?\d{3}[-.\s]?\d{4}\b"),
11     "ip_address": re.compile(
12         r"\b(?:\d{1,3}\.){3}\d{1,3}\b"),
13 }
14
15 def hybrid_pii_scan(table_name: str, columns: List[str],
16                    sample_rows: List[Dict]) -> Dict:
17     # Phase 1: Regex scan
18     regex_findings = {}
19     for row in sample_rows:
20         for col in columns:
21             val = str(row.get(col, ""))
22             for name, pattern in PII_PATTERNS.items():
23                 if pattern.search(val):
24                     regex_findings.setdefault(col, set()).add(name)
25
26     # Phase 2: Claude Code semantic scan

```

```
27     client = anthropic.Anthropic()
28     prompt = f"""Examine sample rows from '{table_name}' and identify
29 PII values that regex might miss. Look for:
30 - Names, addresses, dates of birth in free-text fields
31 - Columns that become PII when combined (quasi-identifiers)
32 - Encoded or obfuscated PII
33 - Device identifiers, IP addresses, geolocation data
34
35 Return JSON array with column_name, detected_pattern,
36 pii_category, confidence, gdpr_special_category (boolean).
37
38 Sample data: {json.dumps(sample_rows[:5], indent=2, default=str)}"""
39     response = client.messages.create(
40         model="claude-sonnet-4-20250514", max_tokens=2048,
41         messages=[{"role": "user", "content": prompt}])
42     claude_findings = json.loads(response.content[0].text)
43
44     # Phase 3: Merge
45     merged = {}
46     for col, patterns in regex_findings.items():
47         merged[col] = {"method": "regex", "patterns": list(patterns)}
48     for f in claude_findings:
49         col = f["column_name"]
50         if col in merged:
51             merged[col]["method"] = "both"
52         else:
53             merged[col] = {"method": "claude", "category": f["
54 pii_category"]}
55     return merged
```

Listing 17.4: Hybrid PII detection pipeline

 **Warning**

Never send real production PII to any external API unless your data processing agreement explicitly permits it. Scan metadata and column names at the schema level, or use synthetic/redacted samples. When in doubt, scan locally and send only aggregate findings (column names and pattern types, not actual values) to Claude Code for classification.

! Tip

For GDPR compliance, distinguish between “regular” PII (name, email, address) and “special category” data (racial or ethnic origin, political opinions, religious beliefs, health data, sexual orientation). Special category data requires explicit consent and additional safeguards under Article 9.

17.4 Compliance Automation: GDPR, CCPA, and HIPAA

Compliance penalties are severe—up to 4% of annual turnover under GDPR, up to \$7,500 per violation under CCPA, and up to \$1.5 million per violation category per year under HIPAA. Claude Code can automate compliance checks, generate Data Subject Access Request (DSAR) queries, and produce Article 30 Records of Processing Activities (ROPA) from your data catalog.

DSAR Query Generation

When a data subject requests access to their personal data (a right under both GDPR and CCPA), you need to identify and extract all data associated with that individual across your entire platform. Claude Code can generate the comprehensive queries needed:

```

1 def generate_dsar_queries(
2     catalog: dict, subject_identifier: str, identifier_column: str
3 ) -> list[str]:
4     """Generate queries to extract all data for a data subject."""
5     client = anthropic.Anthropic()
6     tables_with_pii = {
7         name: meta for name, meta in catalog.items()
8         if meta.get("classification") in ("confidential", "restricted")
9     }
10    message = client.messages.create(
11        model="claude-sonnet-4-20250514", max_tokens=4096,
12        system="""Generate SQL queries to extract all personal data
13        for a data subject. Include:
14        1. Direct matches on the identifier column
15        2. Indirect matches through foreign key relationships
16        3. Queries for audit/log tables that reference the subject
17        Return as a JSON array of {table, query, data_category} objects.
18    """,
19        messages=[{"role": "user",

```

```

19         "content": f"Catalog:\n{json.dumps(tables_with_pii, indent
20             =2)}\n\n"
21             f"Subject identifier: {subject_identifier}\n"
22             f"Identifier column: {identifier_column}"}],
23     )
24     return json.loads(message.content[0].text)

```

Listing 17.5: DSAR query generator

Right to Erasure (“Right to Be Forgotten”)

GDPR Article 17 grants data subjects the right to erasure. Implementing this requires knowing every location where their data exists—including backups, derived tables, and cached aggregations. Claude Code can generate an erasure plan:

```

1  -- Step 1: Delete from transactional tables
2  DELETE FROM orders WHERE user_id = :subject_id;
3  DELETE FROM order_items WHERE order_id IN (
4      SELECT id FROM orders WHERE user_id = :subject_id
5  );
6
7  -- Step 2: Anonymize in analytical tables (preserve aggregates)
8  UPDATE analytics.user_metrics
9  SET email = 'redacted', name = 'redacted',
10     ip_address = 'redacted'
11  WHERE user_id = :subject_id;
12
13 -- Step 3: Remove from search indexes
14 -- (requires application-level action)
15
16 -- Step 4: Log the erasure for compliance audit
17 INSERT INTO governance.erasure_log (
18     subject_id, requested_at, completed_at, tables_affected
19 ) VALUES (
20     :subject_id, :request_timestamp, CURRENT_TIMESTAMP,
21     'orders,order_items,analytics.user_metrics'
22 );

```

Listing 17.6: Erasure plan generated by Claude Code

Warning

Claude Code-generated compliance artifacts must be reviewed by your legal and privacy teams before submission to regulators or data subjects. AI-generated content does not constitute legal advice. This is especially critical for erasure operations, which are irreversible.

17.5 Governance Policies as Code

The most effective governance policies are enforced automatically. Express governance rules as code and integrate them into your CI/CD pipeline so that policy violations are caught before they reach production.

```

1 import psycopg2, json, sys
2 from dataclasses import dataclass, field
3 from typing import List, Optional
4 from enum import Enum
5
6 class PolicySeverity(Enum):
7     BLOCK = "block"; WARN = "warn"; INFO = "info"
8
9 @dataclass
10 class GovernancePolicy:
11     policy_id: str; name: str; description: str
12     severity: PolicySeverity; check_sql: Optional[str] = None
13
14 POLICIES = [
15     GovernancePolicy("GOV-001", "PII columns must be classified",
16                     "Every PII column must have a classification.",
17                     PolicySeverity.BLOCK, check_sql="""
18                         SELECT cc.column_name, dc.table_name
19                         FROM governance.column_catalog cc
20                         JOIN governance.data_catalog dc ON dc.catalog_id = cc.
21                         catalog_id
22                         WHERE cc.pii_flag = TRUE AND cc.classification IS NULL;"""),
23     GovernancePolicy("GOV-002", "Tables must have a data owner",
24                     "Every table must have an assigned data_owner.",
25                     PolicySeverity.WARN, check_sql="""
26                         SELECT table_schema, table_name
27                         FROM governance.data_dictionary
28                         WHERE data_owner IS NULL GROUP BY 1, 2;"""),
29     GovernancePolicy("GOV-003", "No unencrypted PII in staging",

```

```

29     "PII columns in staging must use encryption or masking.",
30     PolicySeverity.BLOCK, check_sql="""
31         SELECT table_name, column_name
32         FROM governance.column_catalog
33         WHERE pii_flag = TRUE AND schema_name = 'staging'
34         AND encryption_status = 'none'
35         AND masking_policy IS NULL;"""),
36     GovernancePolicy("GOV-004", "Retention policy required",
37     "Every table must have a defined retention period.",
38     PolicySeverity.WARN, check_sql="""
39         SELECT table_schema, table_name
40         FROM governance.data_dictionary
41         WHERE retention_days IS NULL;"""),
42 ]
43
44 class PolicyEngine:
45     def __init__(self, conn_string: str, policies: list):
46         self.conn = psycopg2.connect(conn_string)
47         self.policies = policies
48
49     def run_all_checks(self) -> bool:
50         all_passed = True
51         cur = self.conn.cursor()
52         for policy in self.policies:
53             if not policy.check_sql: continue
54             cur.execute(policy.check_sql)
55             violations = cur.fetchall()
56             if violations:
57                 marker = "BLOCK" if policy.severity == PolicySeverity.
BLOCK else "WARN"
58                 print(f"[{marker}] {policy.policy_id}: {policy.name} "
59                       f"- {len(violations)} violation(s)")
60                 for v in violations[:5]:
61                     print(f"  -> {v}")
62                 if len(violations) > 5:
63                     print(f"  ... and {len(violations) - 5} more")
64                 if policy.severity == PolicySeverity.BLOCK:
65                     all_passed = False
66             cur.close()
67         return all_passed
68
69 if __name__ == "__main__":
70     engine = PolicyEngine("postgresql://gov:pass@localhost/warehouse",
POLICIES)
71     if not engine.run_all_checks():
72         print("\nGovernance check FAILED. Pipeline blocked.")

```

```
73     sys.exit(1)
74     print("\nAll governance checks passed.")
```

Listing 17.7: Policy enforcement engine

Integrating Governance into CI/CD

The policy engine becomes most valuable when integrated into your CI/CD pipeline. Add it as a step in your GitHub Actions or GitLab CI configuration:

```
1  # .github/workflows/governance.yml
2  name: Governance Checks
3  on:
4    pull_request:
5      paths:
6        - 'dbt/**'
7        - 'airflow/**'
8        - 'migrations/**'
9
10 jobs:
11   governance:
12     runs-on: ubuntu-latest
13     steps:
14       - uses: actions/checkout@v4
15
16       - name: Run governance policy checks
17         run: python governance/policy_engine.py
18         env:
19           GOVERNANCE_DB_URL: ${ secrets.GOVERNANCE_DB_URL }
20
21       - name: Run PII scan on changed models
22         run: |
23           python governance/pii_scanner.py \
24             --changed-files $(git diff --name-only HEAD~1)
25
26       - name: Verify catalog coverage
27         run: python governance/catalog_coverage.py --threshold 90
```

Listing 17.8: Governance check in GitHub Actions

Tip

Version your governance policies alongside your data pipeline code. Use pull requests for policy changes so that stakeholders can review modifications before they take effect. This creates an audit trail showing who approved each policy change and when.

Note

Policy Evolution Start with a small set of high-impact policies (PII classification, data ownership) and expand gradually. Introducing 50 policies on day one overwhelms teams and creates resistance. A phased rollout with clear communication about why each policy matters drives better adoption.

17.6 Data Quality as a Governance Concern

Data quality is often treated as a technical concern, but it is fundamentally a governance issue. Poor data quality violates the implicit contract between data producers and consumers. Claude Code can help enforce quality standards as part of governance:

```
1 def generate_quality_rules(table_name: str, business_context: str,
2                             schema: dict) -> list[dict]:
3     """Generate data quality rules from business context."""
4     client = anthropic.Anthropic()
5     message = client.messages.create(
6         model="claude-sonnet-4-20250514", max_tokens=4096,
7         system="""Generate data quality rules as JSON. For each rule:
8         - rule_id, description, severity (critical/warning/info)
9         - check_type (not_null, unique, accepted_values, range,
10         custom_sql, referential_integrity, freshness)
11         - implementation as SQL or dbt test YAML""",
12         messages=[{"role": "user",
13                   "content": f"Table: {table_name}\n"
14                             f"Business context: {business_context}\n"
15                             f"Schema: {json.dumps(schema, indent=2)}"}],
16     )
17     return json.loads(message.content[0].text)
```

Listing 17.9: Quality rule generation from business requirements

17.7 Case Study: Governance Transformation at a Healthcare Company

A mid-size healthcare analytics company with 300 tables across three databases had no formal governance. After a close call with a HIPAA audit, leadership mandated a governance program. Using Claude Code, the team achieved the following in 12 weeks:

Weeks 1–2: Ran automated catalog generation across all three databases. Claude Code produced descriptions for all 300 tables and 4,200 columns. Manual review by domain experts refined 15% of descriptions.

Weeks 3–4: Deployed the hybrid PII scanner. Identified 47 columns containing PHI that were not previously classified, including three columns in the analytics schema that contained patient names in free-text “notes” fields.

Weeks 5–8: Implemented governance-as-code with 12 policies. Integrated into CI/CD. Blocked three PRs in the first week that would have exposed unmasked PHI to the analytics layer.

Weeks 9–12: Generated HIPAA compliance documentation: data flow diagrams, access control matrices, and encryption inventories. Passed the subsequent HIPAA audit with zero findings.

Warning

This case study illustrates what is achievable, but every organization’s regulatory environment is different. Engage qualified legal and compliance professionals before implementing governance programs for regulated industries.

17.8 Exercises

1. **Catalog Generation:** Run the automated cataloging pipeline against a database you manage. Review the generated descriptions for accuracy. How many are correct without modification? What patterns does Claude Code get wrong?
2. **PII Scanner:** Run the hybrid PII detection pipeline against a staging database. Compare regex-only findings to the hybrid approach. How many additional PII columns does the semantic scan identify?

3. **Policy Engine:** Implement three governance policies relevant to your organization. Integrate them into your CI/CD pipeline. Track how many violations they catch in the first two weeks.
4. **DSAR Simulation:** Using a test database with synthetic data, generate DSAR queries for a fictional data subject. Verify that the queries find all relevant data across all tables.
5. **Governance Roadmap:** Draft a 12-week governance implementation plan for your organization using the case study as a template. Identify the top five policies to implement first and justify your prioritization.

Summary

In this chapter, we built a comprehensive data governance practice powered by Claude Code: automated data cataloging with gap analysis, lineage tracing and column-level lineage extraction, hybrid PII detection combining regex and semantic analysis, compliance automation for GDPR/CCPA/HIPAA including DSAR and erasure support, governance-as-code policy enforcement integrated into CI/CD, and data quality as a governance concern.

Note

Governance is a team sport. Claude Code accelerates the technical work—scanning, classifying, generating documentation—but organizational success depends on executive sponsorship, clear ownership, and a culture that values data as a strategic asset. The best governance program is one that teams actually follow, and automation makes compliance the path of least resistance.

18

Security and Access Control

Tip

Security is not a feature you bolt on at the end of a data pipeline—it is a foundational discipline woven into every layer. `CLAUDE CODE` can dramatically accelerate security tasks such as policy generation, vulnerability assessment, and audit log analysis, but ultimate responsibility for protecting sensitive data remains with your team.

Data engineers work directly with raw datasets, build pipelines that move sensitive information, and design storage layers. A single misconfigured pipeline can expose millions of records. A single over-permissioned role can give an intern access to production PII. This chapter provides production-ready patterns for RBAC/ABAC, encryption, data masking, row-level security, audit log anomaly detection, and secrets management—all augmented by `CLAUDE CODE`.

18.1 Data Security Fundamentals

Before diving into implementation, establish the mental models that guide all security decisions.

The *CIA triad*—Confidentiality, Integrity, and Availability—is the bedrock of information security. Confidentiality ensures data is accessible only to authorized parties. Integrity ensures data has not been tampered with or corrupted. Availability ensures data is accessible when needed.

Defense in depth layers controls across multiple levels: network (firewalls, VPNs,

private endpoints), identity (SSO, MFA, service accounts), authorization (RBAC, ABAC, row-level security), data (encryption, masking, tokenization), and monitoring (audit logs, anomaly detection, alerting). A failure at one layer should not compromise the system because other layers provide backup protection.

The *principle of least privilege* grants only minimum necessary permissions for the minimum necessary duration. This applies to human users, service accounts, and automated pipelines. Every permission granted is an attack surface expanded.

Zero trust assumes no user, device, or network segment is inherently trustworthy. Every access request is authenticated and authorized, regardless of source. This model is particularly important for data platforms that span cloud environments and serve diverse consumer groups.

Warning

Over-permissioning is the most common security anti-pattern in data engineering. It is far easier to grant access than to revoke it after a breach. Start with no permissions and add them incrementally based on documented, approved requests. Audit permissions quarterly and revoke anything unused for 90 days.

The Data Engineer's Security Responsibilities

Data engineers have unique security responsibilities because they operate at the intersection of raw data and transformation logic:

- **Pipeline credentials:** Managing database passwords, API keys, and service account tokens used by automated pipelines.
- **Data movement:** Ensuring that sensitive data is encrypted in transit as it moves between systems.
- **Access control design:** Defining who can read, write, and modify tables, views, and schemas.
- **Data masking:** Ensuring that non-production environments (staging, development) do not contain unmasked PII.
- **Audit compliance:** Maintaining logs that demonstrate who accessed what data and when.

Note

Security as a Career Differentiator Data engineers who understand security deeply are disproportionately valuable. Most data engineering curricula underemphasize security, creating a skills gap. Investing in security knowledge pays dividends in career advancement and organizational impact.

18.2 RBAC and ABAC Design

Role-Based Access Control (RBAC) assigns permissions to roles rather than individual users. Users are granted roles, and roles are organized in a hierarchy where higher roles inherit permissions from lower ones.

```
1  -- Create the role hierarchy
2  CREATE ROLE IF NOT EXISTS data_reader;
3  CREATE ROLE IF NOT EXISTS data_analyst;
4  CREATE ROLE IF NOT EXISTS data_engineer;
5  CREATE ROLE IF NOT EXISTS data_admin;
6  CREATE ROLE IF NOT EXISTS security_admin;
7
8  -- Establish hierarchy: analyst inherits reader, etc.
9  GRANT ROLE data_reader TO ROLE data_analyst;
10 GRANT ROLE data_analyst TO ROLE data_engineer;
11 GRANT ROLE data_engineer TO ROLE data_admin;
12 GRANT ROLE security_admin TO ROLE SYSADMIN;
13
14 -- Reader: SELECT on production tables
15 GRANT USAGE ON DATABASE analytics_warehouse TO ROLE data_reader;
16 GRANT SELECT ON ALL TABLES IN DATABASE analytics_warehouse
17     TO ROLE data_reader;
18 GRANT SELECT ON FUTURE TABLES IN DATABASE analytics_warehouse
19     TO ROLE data_reader;
20
21 -- Engineer: full access to staging
22 GRANT ALL PRIVILEGES ON DATABASE staging TO ROLE data_engineer;
23
24 -- Security admin: manage grants
25 GRANT MANAGE GRANTS ON ACCOUNT TO ROLE security_admin;
```

Listing 18.1: RBAC role hierarchy in Snowflake

RBAC Design Principles

When designing an RBAC hierarchy, follow these principles:

1. **Separation of duties:** No single role should be able to both create data and approve its release. The role that writes to production should not be the role that manages access to production.
2. **Role granularity:** Create roles at the right level of granularity. Too few roles lead to over-permissioning; too many roles become unmanageable. A good heuristic is one role per logical persona (reader, analyst, engineer, admin, security).
3. **Environment isolation:** Separate roles for development, staging, and production. A developer should never accidentally modify production data.
4. **Future-proofing:** Use `FUTURE GRANTS` (Snowflake) or equivalent mechanisms to ensure new objects automatically inherit appropriate permissions.

Tip

Use `CLAUDE CODE` to generate RBAC hierarchies from natural language descriptions of your team structure. Provide the team roles, data assets, and security requirements, and `CLAUDE CODE` will generate the complete SQL. Then review carefully—`CLAUDE CODE` may not know your organization’s specific compliance requirements.

Attribute-Based Access Control

ABAC extends RBAC by evaluating attributes of the user, resource, and environment at query time. This enables policies such as “EU users can only access data where `data_region = 'EU'`” or “contractors can only access data from the current quarter.”

```
1  -- Row-level security policy: regional data access
2  CREATE OR REPLACE ROW ACCESS POLICY regional_access AS
3    (region VARCHAR) RETURNS BOOLEAN ->
4    CASE
5      WHEN CURRENT_ROLE() IN ('DATA_ADMIN', 'SECURITY_ADMIN')
6        THEN TRUE  -- Admins see all regions
7      WHEN CURRENT_ROLE() = 'EU_ANALYST'
```

```

8         THEN region IN ('EU', 'UK')
9     WHEN CURRENT_ROLE() = 'US_ANALYST'
10        THEN region IN ('US', 'CA')
11        ELSE FALSE
12    END;
13
14    -- Apply to table
15    ALTER TABLE analytics.customers
16        ADD ROW ACCESS POLICY regional_access ON (country_region);

```

Listing 18.2: ABAC-style row-level security in Snowflake

CLAUDE CODE can generate ABAC policies from natural-language requirements, translating business rules like “regional managers see data for their region and all sub-regions” into SQL policies with appropriate role mappings.

18.3 Encryption Strategies

Encryption protects data at three levels, each addressing a different threat model.

Encryption at rest with customer-managed keys (CMKs) provides control beyond platform-managed encryption. If you use platform-managed keys, the cloud provider controls the keys—which means a compromised provider or a government subpoena could access your data. CMKs give you the ability to revoke access independently.

Encryption in transit enforces TLS 1.2+ for all data movement: between your application and the database, between pipeline components, and between your warehouse and BI tools. Ensure that your connection strings specify SSL/TLS and that certificate validation is enabled (not disabled for “convenience”).

Column-level encryption adds protection for individual sensitive columns (SSNs, credit cards, health identifiers) beyond what at-rest encryption provides. Even if someone gains access to the table, encrypted columns remain unreadable without the decryption key.

```

1  from cryptography.fernet import Fernet
2  from cryptography.hazmat.primitives.kdf.pbkdf2 import PBKDF2HMAC
3  from cryptography.hazmat.primitives import hashes
4  import base64, os
5
6  class ColumnEncryptor:
7      """Encrypt and decrypt individual column values.
8

```

```

9     Uses PBKDF2 key derivation for security and Fernet
10    (AES-128-CBC + HMAC-SHA256) for authenticated encryption.
11    """
12
13    def __init__(self, master_key: bytes, salt: bytes = None):
14        self.salt = salt or os.urandom(16)
15        kdf = PBKDF2HMAC(algorithm=hashes.SHA256(), length=32,
16                        salt=self.salt, iterations=480_000)
17        key = base64.urlsafe_b64encode(kdf.derive(master_key))
18        self.fernet = Fernet(key)
19
20    def encrypt_value(self, plaintext: str) -> str:
21        return self.fernet.encrypt(
22            plaintext.encode("utf-8")).decode("utf-8")
23
24    def decrypt_value(self, ciphertext: str) -> str:
25        return self.fernet.decrypt(
26            ciphertext.encode("utf-8")).decode("utf-8")
27
28    def encrypt_dataframe_column(self, df, column_name: str):
29        """Encrypt a column in a pandas DataFrame in place."""
30        df[column_name] = df[column_name].apply(
31            lambda v: self.encrypt_value(str(v)) if v else v
32        )
33        return df

```

Listing 18.3: Column-level encryption utility

Warning

Column-level encryption makes columns unsearchable and unjoinable without decryption. Plan your query patterns before implementing it. If you need to join on an encrypted column, consider deterministic encryption (same plaintext always produces the same ciphertext) at the cost of reduced security.

18.4 Data Masking and Tokenization

Data masking replaces sensitive values with fictional but realistic substitutes. Unlike encryption, masking is typically irreversible (one-way) and preserves the format and statistical properties of the data.

```
1  -- Email masking: show full to admins, partial to analysts, none to
   others
2  CREATE OR REPLACE MASKING POLICY mask_email AS
3      (val STRING) RETURNS STRING ->
4      CASE
5          WHEN CURRENT_ROLE() IN ('DATA_ADMIN', 'SECURITY_ADMIN')
6              THEN val
7          WHEN CURRENT_ROLE() IN ('DATA_ANALYST')
8              THEN REGEXP_REPLACE(val, '~(.{2})(.*)((@.*))$', '\\1****\\3')
9          ELSE '****@****.***'
10     END;
11
12 -- SSN masking: full to security admin, last 4 to admin, none to others
13 CREATE OR REPLACE MASKING POLICY mask_ssn AS
14     (val STRING) RETURNS STRING ->
15     CASE
16         WHEN CURRENT_ROLE() IN ('SECURITY_ADMIN') THEN val
17         WHEN CURRENT_ROLE() IN ('DATA_ADMIN')
18             THEN 'XXX-XX-' || RIGHT(val, 4)
19         ELSE 'XXX-XX-XXXX'
20     END;
21
22 -- Phone masking: preserve country code for analytics
23 CREATE OR REPLACE MASKING POLICY mask_phone AS
24     (val STRING) RETURNS STRING ->
25     CASE
26         WHEN CURRENT_ROLE() IN ('DATA_ADMIN', 'SECURITY_ADMIN')
27             THEN val
28         ELSE LEFT(val, 3) || '-XXX-XXXX'
29     END;
30
31 -- Apply policies to columns
32 ALTER TABLE hr.employees MODIFY COLUMN email
33     SET MASKING POLICY mask_email;
34 ALTER TABLE hr.employees MODIFY COLUMN ssn
35     SET MASKING POLICY mask_ssn;
36 ALTER TABLE hr.employees MODIFY COLUMN phone
37     SET MASKING POLICY mask_phone;
```

Listing 18.4: Snowflake dynamic data masking policies

Tokenization

Tokenization replaces sensitive values with opaque tokens that can be reversed by an authorized vault. This differs from encryption in that the token bears no mathematical relationship to the original value—you cannot derive the original from the token without access to the vault’s lookup table.

Use **deterministic tokenization** when you need to join on tokenized values across tables. The same input always produces the same token, enabling joins without detokenization. Use **random tokenization** for stronger security when joins are not needed.

! Tip

For non-production environments, use masking rather than tokenization. Masking is simpler, does not require vault infrastructure, and is sufficient for development and testing. Reserve tokenization for production scenarios where reversibility is required (e.g., customer service lookups).

Masking for Non-Production Environments

A common security gap is using production data in development and staging environments. CLAUDE CODE can generate masking scripts that produce realistic but fictional data:

```

1  -- Generate a masked copy of the customers table for staging
2  CREATE TABLE staging.customers AS
3  SELECT
4      customer_id,
5      -- Preserve format but replace values
6      MD5(first_name || 'salt') AS first_name,
7      MD5(last_name || 'salt') AS last_name,
8      CONCAT('user', customer_id, '@example.com') AS email,
9      'XXX-XX-' || LPAD(MOD(ABS(HASH(ssn)), 10000)::VARCHAR, 4, '0')
10     AS ssn,
11     -- Preserve distribution but shift dates
12     date_of_birth + INTERVAL '30 days' AS date_of_birth,
13     -- Keep non-sensitive columns unchanged
14     signup_date, customer_segment, region, status
15 FROM production.customers;

```

Listing 18.5: Production-to-staging masking script

18.5 Row-Level Security

Row-level security (RLS) restricts which rows a user can see based on their identity or role. This is essential for multi-tenant platforms, regional data isolation, and time-based access restrictions.

```
1 -- Analysts can only see data from the last 90 days
2 -- Engineers can see all data
3 CREATE OR REPLACE ROW ACCESS POLICY time_based_access AS
4   (event_date DATE) RETURNS BOOLEAN ->
5   CASE
6     WHEN CURRENT_ROLE() IN ('DATA_ENGINEER', 'DATA_ADMIN')
7       THEN TRUE
8     WHEN CURRENT_ROLE() = 'DATA_ANALYST'
9       THEN event_date >= CURRENT_DATE - 90
10    ELSE FALSE
11  END;
```

Listing 18.6: Time-based row-level security

Note

RLS Performance Row-level security policies are evaluated at query time and can impact performance. Test RLS policies with realistic data volumes. In Snowflake, RLS policies are pushed down to the micro-partition level, so performance impact is typically minimal. In other systems, the impact may be more significant.

18.6 Audit Log Anomaly Detection with Claude Code

Audit logs are your security camera system. They record who accessed what data, when, and how. But raw logs are only useful if someone—or something—is watching them.

```
1 CREATE MATERIALIZED VIEW security.user_baselines AS
2 WITH daily_metrics AS (
3   SELECT user_id, event_timestamp::date AS activity_date,
4         COUNT(*) AS daily_events,
5         COUNT(DISTINCT resource_name) AS daily_resources,
6         COUNT(*) FILTER (WHERE result = 'denied') AS daily_denials,
```

```

7         COUNT(DISTINCT ip_address) AS daily_unique_ips
8     FROM security.audit_log
9     WHERE event_timestamp >= NOW() - INTERVAL '30 days'
10    GROUP BY user_id, event_timestamp::date
11 )
12 SELECT user_id,
13        AVG(daily_events) AS avg_daily_events,
14        STDDEV(daily_events) AS stddev_daily_events,
15        AVG(daily_denials) AS avg_daily_denials,
16        AVG(daily_unique_ips) AS avg_daily_ips,
17        PERCENTILE_CONT(0.95) WITHIN GROUP (ORDER BY daily_events)
18        AS p95_daily_events
19 FROM daily_metrics GROUP BY user_id
20 HAVING COUNT(DISTINCT activity_date) >= 5;
21
22 -- Detect anomalies: users exceeding baselines
23 CREATE VIEW security.anomalous_activity AS
24 SELECT al.user_id, al.event_timestamp::date AS activity_date,
25        COUNT(*) AS today_events, ub.avg_daily_events,
26        CASE WHEN ub.stddev_daily_events > 0 THEN
27            (COUNT(*) - ub.avg_daily_events) / ub.stddev_daily_events
28        ELSE 0 END AS z_score
29 FROM security.audit_log al
30 JOIN security.user_baselines ub ON ub.user_id = al.user_id
31 WHERE al.event_timestamp::date = CURRENT_DATE
32 GROUP BY al.user_id, al.event_timestamp::date,
33          ub.avg_daily_events, ub.stddev_daily_events, ub.avg_daily_ips
34 HAVING CASE WHEN ub.stddev_daily_events > 0 THEN
35            (COUNT(*) - ub.avg_daily_events) / ub.stddev_daily_events
36        ELSE 0 END > 3
37 OR COUNT(*) FILTER (WHERE al.result = 'denied') > 10;

```

Listing 18.7: Computing user behavioral baselines

CLAUDE CODE can classify detected anomalies into categories—credential compromise, insider threat, privilege escalation, data exfiltration, brute force, policy violation, misconfiguration, or normal deviation—and generate incident reports for high-severity findings.

```

1 def classify_anomaly(anomaly_data: dict) -> dict:
2     """Use Claude Code to classify a detected anomaly."""
3     client = anthropic.Anthropic()
4     message = client.messages.create(
5         model="claude-sonnet-4-20250514", max_tokens=2048,

```

```
6     system="""You are a security analyst. Classify this anomaly:
7     Categories: credential_compromise, insider_threat,
8     privilege_escalation, data_exfiltration, brute_force,
9     policy_violation, misconfiguration, normal_deviation.
10    Return JSON with: category, confidence, severity (critical/
11    high/medium/low), recommended_actions, investigation_queries."""
12
13    messages=[{"role": "user",
14              "content": json.dumps(anomaly_data, indent=2)}],
15    )
16    return json.loads(message.content[0].text)
```

Listing 18.8: Claude Code anomaly classification

 **Warning**

Anomaly detection systems generate false positives. Tune your z-score threshold based on your environment. Start with a threshold of 3 (capturing events more than 3 standard deviations from the mean) and adjust based on alert volume. Too many false positives leads to alert fatigue and missed real incidents.

18.7 Compliance Frameworks

Understanding compliance frameworks helps you design security controls that satisfy multiple requirements simultaneously.

SOC 2 evaluates controls across five Trust Service Criteria: Security, Availability, Processing Integrity, Confidentiality, and Privacy. Most data platforms need SOC 2 Type II, which evaluates controls over a period (typically 6–12 months).

GDPR imposes data minimization, right to erasure, data protection by design and by default, and records of processing. The key principle is “privacy by design”—security controls should be built into the architecture, not added after the fact.

HIPAA requires encryption of Protected Health Information (PHI), access controls, six-year audit trails, and Business Associate Agreements (BAAs) with all vendors who process PHI. This includes your cloud provider, your warehouse vendor, and any AI services used for data processing.

CLAUDE CODE can perform gap analysis between your current controls and framework requirements, generating a prioritized remediation plan.

Note

Always review CLAUDE CODE-generated security configurations before applying them to production. Treat the output as a high-quality first draft that requires human validation. Security misconfigurations can be catastrophic, and AI-generated policies may miss organization-specific requirements.

18.8 Secrets Management

Hardcoded credentials are one of the most common vulnerabilities in data engineering. A database password in a Git repository, an API key in an environment variable, or a service account key file on a shared server—each is a breach waiting to happen.

Use a dedicated secrets manager (AWS Secrets Manager, HashiCorp Vault, GCP Secret Manager) with:

- **Encryption:** Secrets encrypted at rest with customer-managed keys.
- **Access logging:** Every secret access logged with timestamp, identity, and source IP.
- **Automatic rotation:** Credentials rotated on a schedule (e.g., every 90 days) without pipeline downtime.
- **Least privilege:** Each pipeline has its own service account with access only to the secrets it needs.
- **Short-lived credentials:** Prefer temporary credentials (AWS STS, GCP workload identity federation) over long-lived keys.

```
1 import boto3, json
2 from functools import lru_cache
3
4 @lru_cache(maxsize=1)
5 def get_warehouse_credentials() -> dict:
6     """Fetch credentials from AWS Secrets Manager with caching."""
7     client = boto3.client("secretsmanager", region_name="us-east-1")
8     response = client.get_secret_value(
9         SecretId="prod/warehouse/credentials"
10    )
11    return json.loads(response["SecretString"])
12
```

```
13 # Usage in pipeline code
14 creds = get_warehouse_credentials()
15 conn = snowflake.connector.connect(
16     account=creds["account"],
17     user=creds["user"],
18     password=creds["password"],
19     warehouse=creds["warehouse"],
20     database=creds["database"],
21 )
```

Listing 18.9: Secrets management best practices

Warning

Never store secrets in environment variables in production. Environment variables are visible to any process in the same context and often appear in crash dumps, log files, and process listings. Use a secrets manager and inject credentials at runtime through the manager's SDK.

Tip

Use CLAUDE CODE to scan your codebase for hardcoded credentials. Feed your repository's file list and sample code to CLAUDE CODE and ask it to identify potential secrets: API keys, passwords, connection strings, and private keys. Many teams are surprised by what they find.

18.9 Case Study: Securing a Multi-Tenant Data Platform

A B2B SaaS company serves 200 customers from a shared Snowflake warehouse. Each customer's data must be isolated—no customer should ever see another customer's data. The security requirements:

1. **Row-level isolation:** Every query must be filtered by `tenant_id`.
2. **Column-level masking:** PII columns show different levels of detail based on role.
3. **Audit trail:** Every data access logged with tenant context.
4. **SOC 2 compliance:** Full audit trail for 12 months, encryption at rest and in transit.

Using CLAUDE CODE, the team generated RLS policies for all 45 tenant-aware tables, masking policies for 12 PII columns across 8 tables, audit log baselines and anomaly detection queries, and SOC 2 control documentation mapping each requirement to an implemented control. The implementation took three weeks instead of the estimated eight, with CLAUDE CODE generating approximately 70% of the SQL and documentation, reviewed and refined by the security team.

Exercises

1. **RBAC Design:** Design a complete RBAC hierarchy for a healthcare analytics platform that must comply with HIPAA. Include at least six roles, the SQL to implement it, and a justification for each permission grant.
2. **Masking Policy:** Write Snowflake masking policies for a `customers` table containing: `full_name`, `email`, `ssn`, `date_of_birth`, and `credit_card_number`. Define three roles with different masking levels and test each role's view of the data.
3. **Compliance Gap Analysis:** Use CLAUDE CODE to perform a SOC 2 gap analysis for a startup with SSO but no MFA, credentials in environment variables, no formal change management, and basic logging but no anomaly detection. Generate a prioritized remediation plan.
4. **RLS System:** Design a row-level security system that supports time-based policies (e.g., “analysts can only see data from the last 90 days”) and hierarchical policies (e.g., “regional managers see all countries in their region”). Implement and test in your warehouse.
5. **Secrets Audit:** Audit your team's current secrets management practices. Identify all hardcoded credentials, environment variable secrets, and unrotated keys. Create a migration plan to move all secrets to a dedicated secrets manager.

Summary

We built a comprehensive security framework: RBAC/ABAC hierarchies with least-privilege design, encryption at rest/in transit/column-level, dynamic data masking and tokenization for production and non-production environments, row-level security for multi-tenant isolation, audit log anomaly detection with Claude Code classifi-

cation, compliance mapping across SOC 2/GDPR/HIPAA, and secrets management best practices.

i Note

Security is a journey, not a destination. Threats evolve, regulations change, and your data platform grows. Make security a first-class concern in every sprint, not an afterthought during audit season. The best time to fix a security gap is before it becomes a breach.

19

Ethics of AI in Data Engineering

Every time a data engineer uses CLAUDE CODE to generate a query, build a pipeline, or make a design choice, they delegate a portion of their professional judgment to a system that lacks moral agency and accountability for outcomes. This chapter confronts the ethical dimensions head-on: bias in AI-generated code, hallucination risks, privacy implications, accountability frameworks, and governance structures for responsible AI use.

19.1 The Responsibility of AI-Augmented Data Engineering

When an engineer asks CLAUDE CODE to generate a query, responsibility does not transfer to the AI. The engineer remains fully accountable for every line of code that enters production. This principle—that *human accountability persists through AI assistance*—is the foundation of ethical AI use.

This may seem obvious in the abstract, but it becomes complex in practice. Consider these scenarios:

- A CLAUDE CODE-generated query silently filters out records for a minority demographic due to a hardcoded assumption in the WHERE clause. Who is responsible when the downstream dashboard underreports that population?
- A pipeline generated by CLAUDE CODE handles currency conversion correctly for USD, EUR, and GBP but uses an outdated exchange rate for currencies

from developing nations. The resulting financial reports disadvantage certain regional teams.

- An AI-generated data quality check passes because the test data was representative of US customers but not of the global customer base. PII from non-US customers leaks through a masking gap.

In each case, the answer is the same: the engineer who approved and deployed the code bears responsibility. AI does not change the accountability model—it changes the *workflow* through which accountable decisions are made.

Warning

The most dangerous moment in AI-assisted development is when everything appears to work. Correct-looking output can mask subtle logical errors that only manifest with production data volumes, edge-case inputs, or demographic groups underrepresented in testing.

The Automation Complacency Problem

As AI-generated code becomes more reliable, engineers face *automation complacency*: the tendency to accept outputs without sufficient scrutiny. This is well-documented in aviation (where autopilot has degraded manual flying skills) and medicine (where AI diagnostic tools can reduce clinical reasoning).

In data engineering, automation complacency manifests as:

- Accepting SQL without reviewing join logic or null handling.
- Deploying migrations without testing rollback procedures.
- Trusting AI-generated quality checks without verifying edge cases.
- Copy-pasting AI-generated pipeline code without understanding error handling paths.
- Skipping code review because “Claude Code wrote it, so it must be good.”

The antidote is not to avoid AI, but to build disciplined review processes that match the risk level of each task.

Risk-Proportionate Oversight

AI involvement exists on a risk spectrum, and human oversight should scale accordingly:

Low risk: Boilerplate generation, syntax help, documentation, code formatting. Minimal review needed beyond basic correctness.

Medium risk: Business logic, test generation, query optimization. Standard code review with attention to edge cases.

High risk: Compliance queries, PII handling, access control policies. Detailed review by subject matter expert plus security review.

Critical risk: Healthcare data transformations, financial regulatory reporting, algorithmic decision-making affecting individuals. Formal review process with sign-off, testing against production-representative data, and documentation of AI involvement.

! Tip

Create a simple decision matrix for your team: for each category of work (pipeline code, SQL queries, infrastructure config, access policies), define the minimum review requirements when AI is involved. Post it in your team's documentation and reference it in code review checklists.

19.2 Bias in AI-Generated Data Pipelines

Bias in AI-generated code is not always obvious. It does not look like explicit discrimination—it looks like reasonable defaults that happen to disadvantage certain groups.

Sources of Bias

Bias arises from multiple sources:

- **Training data patterns:** CLAUDE CODE has seen more US-centric code than code from other regions. It may default to US date formats (MM/DD/YYYY), US-centric timezone handling, English-language text processing assumptions, and dollar-denominated financial calculations.

- **Prompt framing:** The way you describe requirements influences the output. “Segment customers into high-value and low-value” may produce a binary classification that disadvantages customers who are valuable in ways not captured by the specified metric.
- **Selection bias in examples:** If your few-shot examples only cover common cases, the generated code may not handle edge cases that disproportionately affect minority populations.
- **Survivorship bias:** CLAUDE CODE’s training data overrepresents successful, well-documented projects and underrepresents the silent failures and edge cases that are critical to equitable data systems.
- **Proxy discrimination:** A query that filters by zip code may effectively filter by race. A model that uses “years of experience” as a feature may discriminate by age.

A Practical Bias Detection Framework

Code Example

```
1 import anthropic, json
2
3 def audit_query_for_bias(
4     query: str, context: str, protected_attributes: list[str]
5 ) -> dict:
6     """Audit a SQL query for potential bias against protected groups.
7     """
8     client = anthropic.Anthropic()
9     prompt = f"""Analyze this SQL query for potential bias.
10
11 Query: ```sql\n{query}\n```
12 Business Context: {context}
13 Protected Attributes: {', '.join(protected_attributes)}
14
15 Check for:
16 1) Direct use of protected attributes in filtering or scoring
17 2) Proxy variables that correlate with protected attributes
18    (e.g., zip code as proxy for race, name patterns for ethnicity)
19 3) Hardcoded thresholds that disadvantage specific populations
20 4) Missing handling of international formats, languages, or customs
21 5) Assumptions about data distributions that may not hold for
```

```
21     all demographic groups
22 6) Aggregation methods that mask disparities (e.g., averages
23     that hide bimodal distributions)
24
25 Return JSON with "issues" array (each with description, severity,
26 affected_groups, mitigation) and "overall_risk" field. ""
27
28 response = client.messages.create(
29     model="claude-sonnet-4-20250514", max_tokens=4096,
30     messages=[{"role": "user", "content": prompt}])
31 text = response.content[0].text
32 return json.loads(text[text.find("{"):text.rfind("}")+1])
```

Bias Mitigation Strategies

Mitigation requires a multi-layered approach:

1. **Diverse review:** Ensure code reviewers include people from different backgrounds who may spot assumptions invisible to the original author.
2. **Bias-specific test cases:** Write tests that explicitly check for differential outcomes across demographic groups. If your pipeline scores customers, test that score distributions are fair across protected attributes.
3. **Explicit fairness requirements in prompts:** When asking CLAUDE CODE to generate queries involving people, include fairness constraints in the prompt.
4. **Regular audits:** Schedule quarterly bias audits of production pipelines, especially those that influence decisions about people (credit scoring, hiring, content recommendations).
5. **Documentation:** Record AI involvement in code that affects people, including the prompt used, the model version, and any bias review performed.

Tip

When asking CLAUDE CODE to generate queries involving people, always include this in your prompt: “Ensure this query does not directly or indirectly discriminate based on race, gender, age, disability, or other protected characteristics. Flag any assumptions about demographics or geography.”

Note

Bias Is Not Just a Technical Problem Technical bias detection catches the symptoms, but the root causes are often in business requirements. “Identify our best customers” may encode biases about who counts as “best.” Data engineers should push back on requirements that embed discriminatory assumptions, even when the SQL itself is technically correct.

19.3 Hallucination: Causes, Detection, and Mitigation

Hallucination in the context of data engineering manifests in several dangerous forms:

- **Invented column names:** CLAUDE CODE references a column that does not exist in your schema, often choosing a plausible name based on the table context.
- **Incorrect function signatures:** Using a function with the wrong number of arguments or wrong argument types, especially for dialect-specific functions.
- **Fabricated best practices:** Stating “it is best practice to...” followed by advice that sounds authoritative but is not supported by documentation or expert consensus.
- **Plausible but wrong business logic:** Generating a query that looks correct but implements the wrong business rule—for example, calculating revenue as `quantity * price` when the actual formula requires applying discounts and taxes.
- **Non-existent configuration options:** Referencing Airflow operator parameters, dbt config keys, or cloud service options that do not exist.

In data engineering, hallucinated SQL is particularly dangerous because it is executed by machines that cannot distinguish correct from incorrect logic. A hallucinated column name causes an immediate error (detectable); a hallucinated join condition produces wrong results silently (catastrophic).

Detection Strategies

1. **Schema validation:** Always provide actual schema definitions rather than relying on the model to guess. Compare generated column references against the actual schema programmatically.
2. **Dry-run execution:** Execute generated SQL with `EXPLAIN` or `LIMIT 0` before running against real data. This catches non-existent columns and type mismatches.
3. **Dual-prompt verification:** Generate the same query twice with different prompts and compare the results. Disagreements indicate areas requiring human judgment.
4. **Version specification:** Explicitly specify tool and library versions in your prompts. “Using dbt 1.7 with Snowflake adapter 1.7.2” prevents the model from referencing features from other versions.
5. **Runtime assertions:** Embed assertions in generated code that verify assumptions at execution time—row count checks, null count checks, value range checks.

Warning

Never trust CLAUDE CODE when it cites specific documentation URLs, version numbers, release dates, or benchmark results. Always verify against official sources. The model generates plausible-sounding citations that may not correspond to real documents.

A Hallucination Severity Matrix

Not all hallucinations are equally dangerous. Classify them by impact:

Compile-time detectable: Non-existent column names, syntax errors, wrong function signatures. These fail immediately and are low risk.

Test-detectable: Wrong join types, incorrect aggregation logic, missing filters. These produce wrong results but are caught by data quality tests.

Silently wrong: Correct syntax but wrong business logic, subtle filter errors, incorrect timezone handling. These are the most dangerous because they produce plausible but incorrect results that may not be caught for days or weeks.

 **Tip**

Invest your review time proportionally to the severity category. Spend less time checking syntax (which your compiler/executor catches) and more time verifying business logic (which only humans can validate).

19.4 Data Privacy and AI

When you interact with CLAUDE CODE through the API, your full prompt is transmitted to Anthropic's servers. Under standard API terms, Anthropic does not use API inputs to train models, but engineers should always verify current terms and ensure their organization has appropriate data processing agreements in place.

The privacy implications extend beyond the obvious:

- **Schema metadata:** Even without sending actual data, sending your schema reveals your data model, business entities, and potentially sensitive column names.
- **Query patterns:** The queries you ask CLAUDE CODE to generate reveal your business logic, analytical priorities, and potentially sensitive business strategies.
- **Error messages:** Pasting tracebacks may inadvertently include connection strings, file paths, or data samples from error output.
- **Prompt accumulation:** Over time, the sum of all prompts from your organization creates a comprehensive picture of your data infrastructure, even if no single prompt is sensitive.

 **Warning**

Never send production data containing PII, PHI, financial account numbers, or other sensitive information to any AI service without explicit authorization from your security and legal teams. This includes data samples in prompts, error messages containing data, and schema comments that reference real individuals.

Data Classification for AI Interactions

Organizations should establish a data classification scheme that governs what can be shared with AI services:

Public: Non-sensitive data, freely shareable. Schema structures without sensitive column names, generic code patterns, open-source project references.

Internal: Not sensitive but not public. Business logic descriptions, internal tool configurations, non-sensitive schema designs.

Confidential: Requires protection. Production schema metadata, business-specific query patterns, performance data, cost information.

Restricted: Never share with external services. PII, PHI, credentials, encryption keys, financial account numbers, production data samples.

Note

Some organizations deploy self-hosted LLM solutions for restricted data. While this eliminates the external data transmission concern, it introduces new challenges: model quality, infrastructure cost, and update cadence. Evaluate the trade-offs carefully before building in-house.

19.5 Human-in-the-Loop Patterns

For critical data systems, blind automation is unacceptable. *Human-in-the-loop* patterns ensure that AI assistance amplifies human judgment rather than replacing it.

Approval Gates

Route AI-generated changes to human reviewers based on risk level. Low-risk changes (documentation updates, formatting) can be auto-approved. Medium-risk changes (new columns, query modifications) require one reviewer. High-risk changes (access control, compliance logic, schema migrations) require two reviewers plus a domain expert.

Confidence-Based Routing

When CLAUDE CODE generates output, assess confidence based on factors like prompt specificity, schema coverage, and task complexity. Low-confidence outputs are routed to humans for review; high-confidence outputs go through automated validation (tests, dry runs, schema checks) before deployment.

Shadow Mode

Shadow mode is the gold standard for introducing AI into critical systems. The progression:

Phase 1 (Weeks 1–4): AI outputs are generated but not used. They are compared offline to human-produced results. This calibrates trust and identifies systematic errors.

Phase 2 (Weeks 5–8): AI outputs are shown as suggestions in the IDE or code review tool. Engineers choose to accept, modify, or reject each suggestion. Track acceptance rates by task type.

Phase 3 (Weeks 9–16): AI outputs are used by default for accepted task types. Humans sample a percentage of outputs for review. Reduce sampling rate as confidence grows.

Phase 4 (Month 4+): AI operates autonomously for proven task types with automated validation. Humans are alerted only on anomalies, low-confidence outputs, or failed validations.

Tip

Keep detailed records of each phase transition: what task types were promoted, what evidence supported the promotion, and what fallback procedures are in place. This documentation is invaluable for compliance audits and for replicating the process with new AI capabilities.

19.6 Organizational AI Governance

An effective AI governance framework includes these components:

1. **AI usage policy:** What AI tools are approved, what data can be shared with them, and who has authority to approve new use cases.

2. **Risk assessment process:** How to evaluate the risk of a new AI-assisted workflow before deploying it.
3. **Review and approval workflows:** Who reviews AI-generated code, at what granularity, and with what authority to approve or reject.
4. **Audit trail requirements:** What must be logged (prompt, response, model version, reviewer, outcome) and for how long.
5. **Incident response procedures:** What happens when AI-generated code causes an incident. How is root cause determined? How are similar incidents prevented?
6. **Training and certification:** Minimum training requirements for engineers using AI tools, including bias awareness, hallucination detection, and privacy protocols.
7. **Vendor assessment:** Evaluation criteria for AI service providers, including data handling practices, security certifications, and contractual protections.
8. **Continuous improvement:** Regular review of the governance framework itself, incorporating lessons from incidents, audits, and evolving best practices.

The Regulatory Landscape

The regulatory landscape for AI in data engineering is evolving rapidly and varies by jurisdiction:

EU AI Act: Establishes risk-based regulation with transparency obligations, high-risk classification criteria, data governance requirements, and human oversight mandates. Data pipelines that process personal data or influence automated decisions may fall under high-risk categories.

GDPR: Imposes constraints on automated decision-making (Article 22), requiring meaningful information about the logic involved and the right to human review.

CCPA/CPRA: Grants consumers the right to opt out of automated decision-making technologies.

HIPAA: Imposes strict controls on how AI tools interact with Protected Health Information.

SOX: Requires auditability of financial data processing, which extends to AI-generated transformations of financial data.

Warning

Regulatory compliance is not optional. If your data pipelines process regulated data, consult with legal and compliance teams before integrating AI tools. The regulatory landscape is changing faster than most engineering teams realize.

Note

Keep a regulatory watch list: subscribe to updates from relevant regulatory bodies, review AI-related legislation quarterly, and maintain relationships with your legal team. A proactive approach to compliance is far less costly than reactive remediation.

19.7 Building an Ethical AI Culture

Ethics cannot be enforced through policies alone—it requires a culture where engineers feel empowered to raise concerns, question AI outputs, and prioritize correctness over speed.

Key cultural practices:

- **Blameless postmortems:** When AI-generated code causes an issue, focus on process improvements rather than individual blame. Was the review process adequate? Were the right tests in place?
- **Ethics office hours:** Regular sessions where engineers can discuss ethical concerns about AI use cases, proposed pipelines, or data handling practices.
- **Red team exercises:** Periodically have team members attempt to generate biased, incorrect, or privacy-violating code through AI tools. Use the results to improve prompts, review checklists, and training materials.
- **Transparency by default:** Document AI involvement in all code that enters production. Future maintainers should know which code was AI-generated and what review it received.

! Tip

Add an “AI Involvement” section to your code review template. For each PR, the author indicates which parts were AI-generated, what prompts were used, and what validation was performed. This normalizes AI use while maintaining accountability.

19.8 Exercises

Exercises

1. **Bias Audit:** Select a pipeline in your organization that makes decisions about people (customer segmentation, lead scoring, content recommendations). Using the bias detection framework from this chapter, audit it for fairness issues across at least three protected attributes. Document findings and propose specific mitigations.
2. **Hallucination Detection:** Generate five SQL queries using CLAUDE CODE for a schema you know well. For each query, apply at least three hallucination detection strategies from this chapter. Document every discrepancy found, classify it by severity, and estimate how long the issue would have persisted if deployed to production without review.
3. **AI Usage Policy Draft:** Using the governance framework from this chapter as a starting point, draft an AI usage policy customized for your organization. Include data classification rules, review requirements by risk level, training obligations, audit trail specifications, and incident response procedures. Present it to your team for feedback.
4. **Shadow Mode Plan:** Design a shadow mode rollout plan for introducing AI-assisted code generation into one of your team’s workflows. Define the four phases, the metrics for phase transitions, the fallback procedures, and the documentation requirements. Estimate the timeline and resources needed.
5. **Ethical Dilemma Discussion:** Your team has built an AI-assisted pipeline that is 3x faster to develop but has a 2% higher error rate compared to manually written code. Management wants to adopt it for all pipelines, including financial reporting subject to SOX compliance. Write a one-page memo arguing for or against this decision, addressing risk, compliance, efficiency, and the ethical obligations of data engineers.

Part **VII**

Advanced Patterns

20

Claude Code as a Data Engineering Agent

“Any sufficiently advanced agent is indistinguishable from a junior engineer who never sleeps, never complains, and never eats the last donut.”

—CLAUDE CODE, probably

! Tip

This chapter explores autonomous agents that can plan, execute, and recover from complex multi-step workflows. That’s right—CLAUDE CODE building things that use CLAUDE CODE. It’s CLAUDE CODE all the way down. By the end, you will be able to design, build, and deploy CLAUDE CODE-powered agents that automate substantial portions of your data engineering operations, so you can finally get some sleep. Or at least pretend to.

AI agents represent the next leap in abstraction—systems that autonomously plan and execute sequences of actions, using tools, maintaining state, and recovering from errors. For data engineers, agents are compelling because pipelines are inherently multi-step, stateful, failure-prone, and span heterogeneous systems. Where a single prompt-response interaction can generate a SQL query or review a config file, an agent can diagnose a pipeline failure, trace the root cause across three services, apply a fix, validate the result, and notify the on-call engineer—all without human intervention.

20.1 What Are AI Agents

An *AI agent* uses a large language model as its reasoning core, augmented with the ability to take actions through *tool use*. It operates in a loop:

1. **Observe:** Gather information from the environment—logs, metrics, query results, API responses.
2. **Reason:** Analyze the observations, form hypotheses, and plan the next action.
3. **Act:** Execute an action using a tool—run a query, call an API, modify a configuration.
4. **Evaluate:** Assess the result of the action and decide whether to continue, try a different approach, or terminate.

This loop continues until the agent achieves its goal, reaches a resource limit, or determines that human intervention is required. So, basically the same loop a human engineer follows, minus the existential dread and coffee breaks.

Agent capabilities range across a maturity spectrum:

Level 0 — Single Turn No tools. The model answers questions from its training data alone. Useful for knowledge queries but cannot take action.

Level 1 — Tool-Augmented Single round of tool use. The model calls one or more tools, processes results, and responds. Covers simple lookup and validation tasks.

Level 2 — Multi-Step Planning Iterative tool use with planning. The agent breaks a complex task into steps, executes them sequentially, and adapts the plan based on intermediate results.

Level 3 — Autonomous with Recovery Handles errors, retries transient failures, tries alternative approaches, and escalates when stuck. This is the sweet spot for most data engineering automation.

Level 4 — Collaborative Multi-Agent Multiple specialized agents coordinate to solve problems that span domains—for example, a pipeline agent handing off to a data quality agent, which escalates to a schema migration agent.

This chapter focuses on building Level 2 and Level 3 agents, which provide the best balance of capability and controllability for data engineering workflows.

Note

Data engineering is uniquely well-suited for agentic automation because of its repetitive operational patterns, rich tool ecosystems (APIs, CLIs, SQL interfaces), clear success and failure criteria, high cost of manual intervention during off-hours, and well-documented procedures in runbooks and playbooks. In other words, data engineering is the perfect job for a tireless robot. Don't take it personally.

20.2 Tool Use and Function Calling

Tools are the bridge between the agent's reasoning and the real world. Each tool is defined by a name, a description, and a JSON schema for its inputs. The quality of tool definitions directly determines agent performance—a well-described tool enables the model to use it correctly, while a vague description leads to hallucinated parameters and misuse.

Code Example

Defining tools for a data engineering agent

```
1 import anthropic
2
3 tools = [
4     {
5         "name": "run_sql_query",
6         "description": (
7             "Execute a read-only SQL query against the data warehouse
8             "Returns results as a list of dictionaries. Max 1000 rows
9             "Queries are executed with a 30-second timeout. "
10            "Only SELECT statements are allowed."),
11        "input_schema": {
12            "type": "object",
13            "properties": {
14                "query": {
15                    "type": "string",
16                    "description": "The SQL query to execute. Must be
17                                a "                "SELECT statement. Use LIMIT to "
```

```
18         "constrain result size."},
19     "warehouse": {
20         "type": "string",
21         "enum": ["snowflake", "bigquery", "redshift"],
22         "description": "Target warehouse to run the query
23
24         "against."},
25     },
26     "required": ["query", "warehouse"]
27 }
28 {
29     "name": "check_pipeline_status",
30     "description": (
31         "Check the current status of an Airflow DAG run. "
32         "Returns DAG run state, individual task states, "
33         "error messages, and duration. If no execution_date "
34         "is provided, returns the most recent run."),
35     "input_schema": {
36         "type": "object",
37         "properties": {
38             "dag_id": {
39                 "type": "string",
40                 "description": "The Airflow DAG identifier"},
41             "execution_date": {
42                 "type": "string",
43                 "description": "ISO 8601 execution date. "
44                 "Optional; defaults to latest run."
45             },
46             "required": ["dag_id"]
47         }
48     },
49     {
50         "name": "get_table_schema",
51         "description": (
52             "Retrieve the schema of a table including column names, "
53             "types, nullability, and any comments or descriptions. "
54             "Returns metadata as a dictionary."),
55         "input_schema": {
56             "type": "object",
57             "properties": {
58                 "database": {"type": "string"},
59                 "schema": {"type": "string"},
```

```
60         "table": {"type": "string"},
61     },
62     "required": ["database", "schema", "table"]
63 }
64 },
65 {
66     "name": "restart_failed_task",
67     "description": (
68         "Clear and restart a failed Airflow task instance. "
69         "This is a WRITE operation that modifies pipeline state. "
70         "Only clears the specific task, not downstream tasks. "
71         "Returns the new task state after clearing."),
72     "input_schema": {
73         "type": "object",
74         "properties": {
75             "dag_id": {"type": "string"},
76             "task_id": {"type": "string"},
77             "execution_date": {"type": "string"},
78         },
79         "required": ["dag_id", "task_id", "execution_date"]
80     }
81 },
82 ]
```

Warning

Vague or incomplete tool descriptions are the single most common cause of agent failures. Telling CLAUDE CODE to use a tool named `do_stuff` with no description is like handing a new hire a keyboard and saying “fix production.” For every tool, state clearly: what the tool does, the format of its return value, any side effects it produces, constraints and limitations, and whether it is a read or write operation. Spending an extra five minutes on tool descriptions saves hours of debugging agent misbehavior.

Tool Design Principles

Follow these principles when designing tools for data engineering agents:

1. **Atomic operations:** Each tool should do one thing well. Prefer `get_table_schema`

and `run_sql_query` as separate tools over a single `do_database_stuff` tool.

2. **Clear read/write semantics:** Label tools as read-only or write, and describe side effects explicitly. This enables safety policies to distinguish observation from action.
3. **Bounded output:** Cap result sizes (e.g., max 1000 rows) and return structured data. An agent that receives a 50MB query result will exhaust its context window.
4. **Informative errors:** Return structured error messages that help the agent reason about what went wrong—“Column `user_id` not found in table `events`” is far more useful than “Query failed.”
5. **Idempotent where possible:** Tools that can be safely retried simplify error recovery logic.

Always validate tool inputs before execution. Never blindly trust the arguments an LLM provides, especially for tools with side effects. Apply input sanitization, type checking, and authorization checks inside the tool executor, not in the agent loop. Trust, but verify. Actually, just verify.

20.3 The Model Context Protocol (MCP)

The *Model Context Protocol (MCP)* is an open standard that defines how AI models connect to external data sources and tools—a universal adapter layer that eliminates the need to build custom integrations for every service. MCP provides three primitives:

Tools Actions the model can take (e.g., run a query, restart a task, create a table).

Resources Data the model can read (e.g., schema definitions, configuration files, documentation).

Prompts Reusable prompt templates for common workflows (e.g., incident investigation, schema review).

The key advantage of MCP is composability: you can connect multiple MCP servers to a single agent, giving it access to your entire operational toolchain through a standardized protocol.

Code Example

Connecting MCP servers to a Claude Code agent

```
1 from claude_agent_sdk import Agent, MCPServerConfig
2
3 agent = Agent(
4     model="claude-sonnet-4-20250514",
5     mcp_servers=[
6         MCPServerConfig(
7             name="airflow",
8             command="python",
9             args=["mcp_servers/airflow_mcp.py"],
10        ),
11        MCPServerConfig(
12            name="snowflake",
13            command="python",
14            args=["mcp_servers/snowflake_mcp.py"],
15        ),
16        MCPServerConfig(
17            name="dbt",
18            command="python",
19            args=["mcp_servers/dbt_mcp.py"],
20        ),
21        MCPServerConfig(
22            name="pagerduty",
23            command="python",
24            args=["mcp_servers/pagerduty_mcp.py"],
25        ),
26    ],
27    system_prompt=(
28        "You are a data engineering operations agent. "
29        "Diagnose and resolve pipeline issues using the "
30        "available tools. Always check pipeline status "
31        "before attempting fixes. Log all actions taken."
32    ),
33 )
34
35 response = agent.query(
36     "The daily_customer_etl DAG failed at 3 AM. "
37     "Investigate the root cause and fix it if possible."
38 )
```

Note

MCP servers can be implemented in any language. The protocol uses JSON-RPC over stdin/stdout, making it straightforward to wrap existing CLI tools or REST APIs. Community-maintained MCP servers exist for most popular data tools including Snowflake, BigQuery, dbt, Airflow, and Kafka.

Building a Custom MCP Server

When a community MCP server does not exist for your tool, building one is straightforward. The following example wraps a simplified Airflow REST API:

Code Example

Minimal MCP server for Airflow

```
1 from mcp.server import Server
2 from mcp.types import Tool, TextContent
3 import httpx, json
4
5 app = Server("airflow-mcp")
6 AIRFLOW_URL = "http://airflow.internal:8080/api/v1"
7 HEADERS = {"Authorization": "Bearer <token>"}
8
9 @app.list_tools()
10 async def list_tools():
11     return [
12         Tool(
13             name="get_dag_runs",
14             description="List recent DAG runs with status",
15             inputSchema={
16                 "type": "object",
17                 "properties": {
18                     "dag_id": {"type": "string"},
19                     "limit": {"type": "integer", "default": 5},
20                 },
21                 "required": ["dag_id"],
22             },
23         ),
24         Tool(
25             name="get_task_logs",
26             description="Retrieve logs for a specific task "
27                 "instance. Returns the last 200 lines.",
```

```
28     inputSchema={
29         "type": "object",
30         "properties": {
31             "dag_id": {"type": "string"},
32             "task_id": {"type": "string"},
33             "dag_run_id": {"type": "string"},
34         },
35         "required": ["dag_id", "task_id", "dag_run_id"],
36     },
37 ),
38 ]
39
40 @app.call_tool()
41 async def call_tool(name: str, arguments: dict):
42     async with httpx.AsyncClient() as client:
43         if name == "get_dag_runs":
44             resp = await client.get(
45                 f"{AIRFLOW_URL}/dags/{arguments['dag_id']}"
46                 f"/dagRuns?limit={arguments.get('limit', 5)}",
47                 headers=HEADERS,
48             )
49             return [TextContent(
50                 type="text", text=json.dumps(resp.json())
51             )]
52         elif name == "get_task_logs":
53             resp = await client.get(
54                 f"{AIRFLOW_URL}/dags/{arguments['dag_id']}"
55                 f"/dagRuns/{arguments['dag_run_id']}"
56                 f"/taskInstances/{arguments['task_id']}/logs/1",
57                 headers=HEADERS,
58             )
59             lines = resp.text.strip().split("\n")
60             return [TextContent(
61                 type="text", text="\n".join(lines[-200:])
62             )]
63
64 if __name__ == "__main__":
65     import asyncio
66     from mcp.server.stdio import stdio_server
67     asyncio.run(stdio_server(app))
```

20.4 Building a Data Engineering Agent

With tools defined and MCP servers connected, the next step is building the agent loop itself. The core pattern is deceptively simple: send a message with available tools, check if the response includes tool calls, execute the tools, feed results back, and repeat until the model produces a final text response.

Code Example

Core agent loop implementation

```
1 import anthropic, json, logging
2 from dataclasses import dataclass, field
3 from datetime import datetime, timezone
4 from typing import Any, Callable
5
6 logger = logging.getLogger(__name__)
7
8 @dataclass
9 class AgentMemory:
10     """Tracks conversation state and operational limits."""
11     conversation_history: list[dict] = field(default_factory=list)
12     actions_taken: list[dict] = field(default_factory=list)
13     error_count: int = 0
14     max_errors: int = 5
15     iteration_count: int = 0
16     max_iterations: int = 20
17
18 @dataclass
19 class SafetyPolicy:
20     """Defines what the agent is and is not allowed to do."""
21     allowed_tools: set[str] = field(default_factory=set)
22     max_write_actions_per_session: int = 5
23     require_confirmation_for: set[str] = field(default_factory=set)
24     write_action_count: int = 0
25     write_tools: set[str] = field(default_factory=lambda: {
26         "restart_failed_task", "run_dbt_command",
27         "execute_migration", "update_config",
28     })
29
30     def can_execute(self, tool_name: str) -> tuple[bool, str]:
31         if self.allowed_tools and tool_name not in self.allowed_tools
32         :
```

```
32         return False, f"Tool {tool_name} not in allowed set."
33     if tool_name in self.require_confirmation_for:
34         return False, (
35             f"Tool {tool_name} requires human confirmation.")
36     if tool_name in self.write_tools:
37         if self.write_action_count >= self.
max_write_actions_per_session:
38         return False, "Write action budget exhausted."
39         self.write_action_count += 1
40         return True, "Allowed."
41
42 class DataEngineeringAgent:
43     """An autonomous agent for data engineering operations."""
44
45     def __init__(self, tools, tool_executors, system_prompt,
46                 safety_policy=None,
47                 model="claude-sonnet-4-20250514"):
48         self.client = anthropic.Anthropic()
49         self.model = model
50         self.tools = tools
51         self.tool_executors = tool_executors
52         self.system_prompt = system_prompt
53         self.safety_policy = safety_policy or SafetyPolicy()
54         self.memory = AgentMemory()
55
56     def run(self, user_message: str) -> str:
57         """Execute the agent loop until completion or limit."""
58         self.memory.conversation_history.append(
59             {"role": "user", "content": user_message})
60
61         while self.memory.iteration_count < self.memory.
max_iterations:
62             self.memory.iteration_count += 1
63             logger.info(
64                 f"Iteration {self.memory.iteration_count}: "
65                 f"Calling model...")
66
67             response = self.client.messages.create(
68                 model=self.model, max_tokens=4096,
69                 system=self.system_prompt,
70                 tools=self.tools,
71                 messages=self.memory.conversation_history)
72
73             self.memory.conversation_history.append(
```

```
74         {"role": "assistant", "content": response.content})
75
76     # Extract tool use blocks
77     tool_use_blocks = [
78         b for b in response.content
79         if b.type == "tool_use"
80     ]
81     if not tool_use_blocks:
82         # No tool calls: agent is done
83         return "\n".join(
84             b.text for b in response.content
85             if b.type == "text")
86
87     # Execute each tool and collect results
88     tool_results = [
89         self._execute_tool(tb.name, tb.input, tb.id)
90         for tb in tool_use_blocks
91     ]
92     self.memory.conversation_history.append(
93         {"role": "user", "content": tool_results})
94
95     return "Agent reached maximum iteration limit."
96
97     def _execute_tool(self, tool_name, tool_input, tool_use_id):
98         """Execute a single tool with safety checks."""
99         can_execute, reason = self.safety_policy.can_execute(
100             tool_name)
101         if not can_execute:
102             logger.warning(f"BLOCKED: {tool_name}: {reason}")
103             return {
104                 "type": "tool_result",
105                 "tool_use_id": tool_use_id,
106                 "content": f"BLOCKED: {reason}",
107                 "is_error": True,
108             }
109
110         executor = self.tool_executors.get(tool_name)
111         if not executor:
112             return {
113                 "type": "tool_result",
114                 "tool_use_id": tool_use_id,
115                 "content": f"Unknown tool: {tool_name}",
116                 "is_error": True,
117             }
```

```
118
119     try:
120         result = executor(**tool_input)
121         self.memory.actions_taken.append({
122             "tool": tool_name,
123             "input": tool_input,
124             "timestamp": datetime.now(
125                 timezone.utc).isoformat(),
126             "success": True,
127         })
128         return {
129             "type": "tool_result",
130             "tool_use_id": tool_use_id,
131             "content": (json.dumps(result)
132                         if not isinstance(result, str)
133                         else result),
134         }
135     except Exception as e:
136         self.memory.error_count += 1
137         logger.error(
138             f"Tool {tool_name} failed: {e}")
139         return {
140             "type": "tool_result",
141             "tool_use_id": tool_use_id,
142             "content": f"Failed: {str(e)}",
143             "is_error": True,
144         }
```

20.5 Agent Memory and Error Recovery

Production agents need more than a conversation buffer. They require structured memory systems that persist across sessions and enable learning from past operations. Unlike human engineers, who persist knowledge via “that one Slack thread from 2019 that nobody can find.”

Memory Types

Agents benefit from four types of memory:

Working Memory The current conversation context. This is the agent’s “short-

term memory”—the messages exchanged during the current session. It is bounded by the model’s context window.

Episodic Memory Past incidents and their resolutions, stored in a vector database for similarity search. When the agent encounters a new failure, it retrieves similar past incidents to inform its approach.

Semantic Memory Structured knowledge about the environment: pipeline dependency graphs, table schemas, SLA definitions, team ownership mappings. This is typically loaded into the system prompt or retrieved on demand.

Procedural Memory Runbooks and standard operating procedures, encoded as step-by-step instructions the agent can follow. These provide guardrails and ensure consistency with established practices.

Code Example

Episodic memory with vector search

```
1 from dataclasses import dataclass
2 import numpy as np
3
4 @dataclass
5 class Episode:
6     incident_type: str
7     symptoms: str
8     root_cause: str
9     resolution: str
10    embedding: list[float]
11    timestamp: str
12
13 class EpisodicMemory:
14     """Stores and retrieves past incident resolutions."""
15
16     def __init__(self, embedding_service):
17         self.episodes: list[Episode] = []
18         self.embedding_service = embedding_service
19
20     def record(self, incident_type, symptoms,
21              root_cause, resolution):
22         embedding = self.embedding_service.embed(
23             f"{incident_type}: {symptoms}")
24         self.episodes.append(Episode(
25             incident_type=incident_type,
```

```
26         symptoms=symptoms,
27         root_cause=root_cause,
28         resolution=resolution,
29         embedding=embedding,
30         timestamp=datetime.now(
31             timezone.utc).isoformat(),
32     ))
33
34     def recall(self, current_symptoms, top_k=3):
35         query_emb = self.embedding_service.embed(
36             current_symptoms)
37         scored = [
38             (ep, self._cosine_sim(query_emb, ep.embedding))
39             for ep in self.episodes
40         ]
41         scored.sort(key=lambda x: x[1], reverse=True)
42         return scored[:top_k]
43
44     @staticmethod
45     def _cosine_sim(a, b):
46         a, b = np.array(a), np.array(b)
47         return float(np.dot(a, b) / (
48             np.linalg.norm(a) * np.linalg.norm(b)))
```

Error Recovery Strategies

For error recovery, classify errors into three categories and handle each differently:

1. **Transient errors:** Network timeouts, rate limits, temporary service unavailability. Retry with exponential backoff, up to a configurable limit.
2. **Recoverable errors:** Permission denied (try a different credential), query syntax error (fix and retry), missing table (check alternative schemas). The agent can work around these by trying alternative approaches.
3. **Fatal errors:** Data corruption, security violations, unknown system states. The agent must stop immediately and escalate to a human operator.

 **Warning**

Automatic retries can be dangerous for non-idempotent operations. If your tool inserts rows, retrying a partially-completed operation could create duplicates. “Oops, we now have 47 copies of every customer” is not a postmortem title you want to write. Always ensure retried operations are *idempotent*—producing the same result whether executed once or multiple times. Use techniques like `INSERT ... ON CONFLICT DO NOTHING` or merge/upsert semantics.

Code Example

Error classification and recovery

```
1 import time
2 from enum import Enum
3
4 class ErrorCategory(Enum):
5     TRANSIENT = "transient"
6     RECOVERABLE = "recoverable"
7     FATAL = "fatal"
8
9 class ErrorClassifier:
10     TRANSIENT_PATTERNS = [
11         "timeout", "rate limit", "503",
12         "connection reset", "temporarily unavailable",
13     ]
14     FATAL_PATTERNS = [
15         "permission denied", "access denied",
16         "data corruption", "integrity constraint",
17     ]
18
19     def classify(self, error_message: str) -> ErrorCategory:
20         lower = error_message.lower()
21         if any(p in lower for p in self.TRANSIENT_PATTERNS):
22             return ErrorCategory.TRANSIENT
23         if any(p in lower for p in self.FATAL_PATTERNS):
24             return ErrorCategory.FATAL
25         return ErrorCategory.RECOVERABLE
26
27     def retry_with_backoff(func, max_retries=3,
28                           base_delay=1.0):
29         """Retry a function with exponential backoff."""
30         for attempt in range(max_retries):
```

```
31     try:
32         return func()
33     except Exception as e:
34         if attempt == max_retries - 1:
35             raise
36         delay = base_delay * (2 ** attempt)
37         time.sleep(delay)
```

20.6 Safety Guardrails

Safety is not an afterthought—it is the foundation of any production agent. An agent with access to your data warehouse and pipeline orchestrator has the potential to cause significant damage if not properly constrained. Think of it like giving a toddler a flamethrower: impressive capability, questionable judgment. Every agent should implement these four layers of protection:

1. **Least privilege:** Grant only the minimum permissions required. A monitoring agent should have read-only database access. A remediation agent should have write access only to the specific resources it manages.
2. **Action budgets:** Limit the number of write operations per session. A runaway agent that continuously restarts tasks or modifies configurations can cause cascading failures.
3. **Human-in-the-loop checkpoints:** Require human approval for high-risk operations such as schema migrations, production deployments, or actions affecting multiple pipelines simultaneously.
4. **Audit logging:** Record every action in a tamper-resistant location—including the prompt, the model’s reasoning, the tool called, the arguments provided, and the result.

Code Example

Action budget enforcement

```
1 class ActionBudget:
2     """Enforce per-category action limits per session."""
3
4     def __init__(self, budgets: dict[str, int]):
```

```
5     self.budgets = budgets # e.g. {"write": 5, "restart": 3}
6     self.usage = {k: 0 for k in budgets}
7
8     def can_perform(self, action_category: str) -> bool:
9         if action_category not in self.budgets:
10            return False
11            return self.usage[action_category] < \
12                self.budgets[action_category]
13
14        def record(self, action_category: str):
15            self.usage[action_category] = \
16                self.usage.get(action_category, 0) + 1
17
18        def remaining(self, action_category: str) -> int:
19            if action_category not in self.budgets:
20                return 0
21            return self.budgets[action_category] - \
22                self.usage.get(action_category, 0)
23
24        def summary(self) -> str:
25            lines = []
26            for cat, budget in self.budgets.items():
27                used = self.usage.get(cat, 0)
28                lines.append(f" {cat}: {used}/{budget} used")
29            return "\n".join(lines)
```

Tip

Store audit logs in an append-only log store (e.g., an S3 bucket with object lock, or a write-once database table) that the agent itself cannot modify. This ensures that even a compromised or misbehaving agent cannot hide its tracks. The agent equivalent of “this conversation is being recorded for quality assurance purposes.”

The Confirmation Gateway

For high-risk operations, implement a confirmation gateway that pauses the agent and waits for human approval. This can be implemented via Slack, email, or a custom approval UI.

Code Example

Human-in-the-loop confirmation

```
1 import asyncio
2
3 class ConfirmationGateway:
4     """Pause agent execution for human approval."""
5
6     def __init__(self, notifier, timeout_seconds=300):
7         self.notifier = notifier # Slack, email, etc.
8         self.timeout = timeout_seconds
9
10    async def request_approval(self, action_description,
11                               context, risk_level):
12        """Send approval request and wait for response."""
13        request_id = self.notifier.send(
14            channel="#data-ops-approvals",
15            message=(
16                f"*Agent Approval Request*\n"
17                f"Action: {action_description}\n"
18                f"Risk: {risk_level}\n"
19                f"Context: {context}\n"
20                f"Reply 'approve {request_id}' or "
21                f"'deny {request_id}'"),
22        )
23        try:
24            response = await asyncio.wait_for(
25                self.notifier.wait_for_reply(request_id),
26                timeout=self.timeout,
27            )
28            return response.strip().lower() == "approve"
29        except asyncio.TimeoutError:
30            return False # Deny on timeout
```

Note

Timeout-on-denial is a critical safety pattern. If no human responds within the timeout window, the agent should treat the action as denied. Never default to approval on timeout—"nobody said no" is not the same as "yes." That logic does not work in engineering any more than it works in real life.

20.7 Multi-Step Pipeline Automation

The true power of agents emerges in multi-step workflows where each action depends on the results of prior actions. Consider a typical pipeline failure investigation:

1. Check the DAG status to identify which task failed
2. Retrieve the task logs to understand the error
3. Query the data warehouse to verify data state
4. Check if this is a known issue by searching episodic memory
5. Apply the fix (restart, rerun with parameters, or modify config)
6. Validate the fix by checking downstream data quality
7. Notify the team with a summary of the incident and resolution

An agent can execute this entire workflow autonomously, typically completing in 2–5 minutes what would take a human engineer 30–60 minutes—especially at 3 AM, when the human’s cognitive abilities are roughly equivalent to a `SELECT * FROM brain WHERE caffeine > 0` returning zero rows.

Warning

Multi-step workflows amplify both the benefits and the risks of automation. A mistake in step 5 that goes undetected until step 7 can be difficult to unwind. Always include validation steps after write operations, and design rollback mechanisms for every automated fix.

20.8 Full Example: Schema Migration Agent

To illustrate all concepts in this chapter working together, here is a complete schema migration agent. This agent analyzes the requested change, discovers all downstream consumers, plans reversible migration steps, generates SQL with rollback scripts, tests on staging, and executes if approved.

Code Example

Schema migration agent — orchestration

```

1  class SchemaMigrationAgent(DataEngineeringAgent):
2      """Specialized agent for safe schema migrations."""
3
4      SYSTEM_PROMPT = """You are a schema migration specialist.
5      When given a schema change request:
6      1. Analyze the change and identify affected columns/tables
7      2. Discover all downstream consumers (views, dbt models,
8          dashboards, applications) using the dependency tools
9      3. Plan a backward-compatible migration with these phases:
10         - Add new columns/tables (non-breaking)
11         - Backfill data into new structure
12         - Update consumers to use new structure
13         - Deprecate old columns with a removal date
14      4. Generate migration SQL with rollback scripts
15      5. Test on staging before requesting production approval
16      Never drop columns without a deprecation period.
17      Always create backward-compatible views."""
18
19     def __init__(self, **kwargs):
20         migration_tools = self._build_migration_tools()
21         super().__init__(
22             tools=migration_tools,
23             system_prompt=self.SYSTEM_PROMPT,
24             safety_policy=SafetyPolicy(
25                 require_confirmation_for={
26                     "execute_migration_production",
27                     "drop_column",
28                 },
29                 max_write_actions_per_session=10,
30             ),
31             **kwargs,
32         )

```

Code Example

Example migration SQL generated by the agent

```

1  -- V042__split_customer_address.sql

```

```

2  -- Generated by SchemaMigrationAgent
3  -- Rollback: V042__split_customer_address_ROLLBACK.sql
4
5  -- Step 1: Add new columns (non-breaking)
6  ALTER TABLE ANALYTICS.PUBLIC.CUSTOMERS
7  ADD COLUMN street_address VARCHAR(500),
8  ADD COLUMN city VARCHAR(100),
9  ADD COLUMN state VARCHAR(50),
10 ADD COLUMN zip_code VARCHAR(20),
11 ADD COLUMN country VARCHAR(100);
12
13 -- Step 2: Backfill from existing data
14 UPDATE ANALYTICS.PUBLIC.CUSTOMERS
15 SET street_address = PARSE_ADDRESS(address, 'street'),
16     city = PARSE_ADDRESS(address, 'city'),
17     state = PARSE_ADDRESS(address, 'state'),
18     zip_code = PARSE_ADDRESS(address, 'zip'),
19     country = PARSE_ADDRESS(address, 'country')
20 WHERE address IS NOT NULL;
21
22 -- Step 3: Deprecation comment (sets removal date)
23 COMMENT ON COLUMN ANALYTICS.PUBLIC.CUSTOMERS.ADDRESS IS
24     'DEPRECATED: Use street_address, city, state, zip_code, '
25     'country instead. Will be removed after 2026-07-01.';
26
27 -- Step 4: Backward-compatible view for consumers
28 CREATE OR REPLACE VIEW ANALYTICS.PUBLIC.CUSTOMERS_V2 AS
29 SELECT customer_id, first_name, last_name, email,
30        address, -- kept for backward compatibility
31        street_address, city, state, zip_code, country,
32        created_at, updated_at
33 FROM ANALYTICS.PUBLIC.CUSTOMERS;
34
35 -- Step 5: Grant permissions to match original table
36 GRANT SELECT ON ANALYTICS.PUBLIC.CUSTOMERS_V2
37 TO ROLE ANALYTICS_READER;

```

! Tip

The agent also generates a corresponding rollback script that reverses every step. Because CLAUDE CODE learned from humanity's greatest hits—and its greatest “oh no, roll it back” moments. In the example above, the rollback

would drop the new columns, remove the view, and restore the original column comment. Always verify that rollback scripts are tested alongside the migration itself.

20.9 Monitoring and Observability for Agents

Agents are software systems and require the same observability practices as any production service. Key metrics to track include:

- **Iterations per task:** How many agent loop iterations are needed to complete each task type. Increasing iteration counts may indicate degraded tool descriptions or model performance.
- **Token usage per task:** Total input and output tokens consumed. Track this by task type to identify unexpectedly expensive operations.
- **Tool call success rate:** The percentage of tool calls that succeed. A drop signals tool infrastructure issues.
- **Error escalation rate:** How often the agent escalates to human intervention. This is your primary measure of agent autonomy.
- **Time to resolution:** Wall-clock time from task submission to completion. Compare against manual resolution times to quantify ROI.

i Note

Set up alerting on agent error escalation rate. If it rises above your baseline (typically 10–20% for well-tuned agents), investigate whether tool definitions need updating, the system prompt needs refinement, or the environment has changed in ways the agent does not understand. If the escalation rate hits 100%, congratulations: you have built a very expensive pager.

20.10 Exercises

1. **Data Quality Agent:** Create an agent that runs Great Expectations checkpoints, analyzes failures, determines root causes, and suggests fixes—distinguishing data issues (bad upstream data), code issues (incorrect transformation logic), and infrastructure issues (timeouts, resource exhaustion). The agent should produce a structured incident report for each failure.

2. **Circuit Breaker Pattern:** Extend the agent with a circuit breaker that disables a tool after three consecutive failures, resetting after a configurable cool-down period. Test this by simulating intermittent tool failures and verifying that the agent gracefully degrades.
3. **Cost-Aware Agent:** Track token usage per iteration, switch to a cheaper model (e.g., Haiku) when session cost exceeds a soft threshold, and terminate gracefully at a hard budget limit. Implement a cost report that summarizes token usage by tool call type.
4. **MCP Server Development:** Build an MCP server for a data system you use regularly (dbt, Databricks, Fivetran, etc.) with at least five tools covering both read and write operations. Connect it to a CLAUDE CODE agent and demonstrate an end-to-end workflow such as triggering a dbt run and validating the results.
5. **Multi-Agent Coordination:** Design and implement a system where a “triage agent” receives an alert, classifies it, and dispatches to specialized agents (pipeline agent, data quality agent, infrastructure agent). Implement a shared state store so agents can see each other’s findings. Test with three simulated incident types.

21

RAG for Data Engineering

“RAG lets CLAUDE CODE look things up before answering. Before RAG, it just made things up confidently. So, like a senior engineer.”

! Tip

Retrieval-Augmented Generation (RAG) bridges the gap between CLAUDE CODE’s broad knowledge and your organization’s specific data engineering context. This chapter teaches you to build RAG systems that turn your team’s runbooks, data catalogs, and incident histories into an always-available engineering assistant—one that actually reads the documentation, which already puts it ahead of most humans.

Large language models know a great deal about the world, but they do not know about *your* data warehouse, *your* pipeline architecture, or *your* on-call procedures. **Retrieval-Augmented Generation** solves this by dynamically injecting relevant context from your knowledge base into the prompt at query time. Instead of fine-tuning a model on your documentation—an expensive process that creates a snapshot that goes stale—RAG keeps the model’s general capabilities intact while grounding every response in your current, authoritative sources.

The result is a system that can answer questions like “What is the SLA for the customer_events pipeline?” or “How do I recover from a Snowflake warehouse suspension?” by finding the relevant runbook paragraphs and synthesizing a precise answer—complete with source citations. No more asking Dave, who wrote that runbook three years ago, left the company, and took all the tribal knowledge with him. Thanks, Dave.

21.1 RAG Fundamentals

A RAG system has three core components:

1. **Indexing Pipeline:** Process source documents, split them into chunks, generate vector embeddings for each chunk, and store everything in a vector database. This runs offline (batch) or incrementally as documents change.
2. **Retrieval Engine:** At query time, embed the user's question, search the vector database for the most similar chunks, and optionally re-rank or filter results by metadata.
3. **Generation Engine:** Combine the retrieved chunks with the original question into a prompt, send it to CLAUDE CODE, and return the response along with source citations.

Warning

Fine-tuning is often the first approach teams consider, because humans love the most expensive option first. But for most data engineering use cases, RAG is superior. Fine-tuning is expensive (thousands of dollars per training run), creates models that go stale as your documentation evolves, can degrade the model's general capabilities, and makes it difficult to attribute answers to specific sources. RAG keeps the model intact while dynamically grounding it in current information. It's like giving the model an open-book exam instead of making it memorize the entire textbook.

RAG offers several advantages that make it the right default for knowledge-grounded applications:

- **Always current:** Re-index when documents change. No retraining required.
- **Attributable:** Every answer can cite its sources, enabling verification.
- **Composable:** Combine runbooks, data catalogs, incident reports, and Slack threads into a single searchable knowledge base.
- **Access-controlled:** Filter retrieved chunks by user permissions—an analyst should not see infrastructure runbooks marked as ops-only.
- **Debuggable:** When an answer is wrong, inspect the retrieved chunks to determine whether the problem is retrieval (wrong chunks found) or generation (model misinterpreted correct chunks).

Code Example

RAG pipeline data structures

```
1 from dataclasses import dataclass, field
2 from typing import Any
3
4 @dataclass
5 class Document:
6     """A source document before chunking."""
7     content: str
8     metadata: dict[str, Any]
9     source_type: str # "runbook", "catalog", "incident"
10    source_id: str
11
12 @dataclass
13 class Chunk:
14     """A chunk of a document with its embedding."""
15     content: str
16     metadata: dict[str, Any]
17     chunk_index: int
18     document_id: str
19     embedding: list[float] | None = None
20     token_count: int = 0
21
22 @dataclass
23 class RetrievalResult:
24     """A ranked retrieval result."""
25     chunk: Chunk
26     relevance_score: float
27     rank: int
28
29 @dataclass
30 class RAGResponse:
31     """The final RAG response with provenance."""
32     answer: str
33     sources: list[RetrievalResult]
34     confidence: float
35     tokens_used: int
```

21.2 Building a Knowledge Base

The quality of a RAG system depends entirely on the quality and coverage of its knowledge base. For data engineering teams, the most valuable sources include:

- **Runbooks:** Step-by-step procedures for incident response, deployments, and maintenance tasks.
- **Data catalog entries:** Table descriptions, column definitions, ownership, SLAs, and lineage information.
- **Architecture documentation:** System diagrams, design decisions, and integration patterns.
- **Incident reports:** Post-mortems with root cause analysis and remediation steps.
- **dbt project documentation:** Model descriptions, tests, source definitions from `schema.yml` files.
- **Relevant Slack/Teams threads:** Curated discussions containing tribal knowledge.

Code Example

Multi-source document ingestion

```
1 from abc import ABC, abstractmethod
2 import json, yaml, httpx
3
4 class DocumentSource(ABC):
5     """Base class for document sources."""
6     @abstractmethod
7     def fetch_documents(self) -> list[Document]:
8         ...
9
10 class ConfluenceSource(DocumentSource):
11     """Ingest runbooks from Confluence."""
12     def __init__(self, base_url, space_key, token):
13         self.base_url = base_url
14         self.space_key = space_key
15         self.headers = {"Authorization": f"Bearer {token}"}
16
17     def fetch_documents(self) -> list[Document]:
```

```
18     docs = []
19     url = (f"{self.base_url}/rest/api/content"
20           f"?spaceKey={self.space_key}"
21           f"&expand=body.storage,metadata.labels"
22           f"&limit=100")
23     resp = httpx.get(url, headers=self.headers)
24     for page in resp.json()["results"]:
25         labels = [l["name"] for l in
26                  page.get("metadata", {}).
27                          .get("labels", {}).
28                          .get("results", [])]
29         docs.append(Document(
30             content=self._html_to_text(
31                 page["body"]["storage"]["value"]),
32             metadata={
33                 "title": page["title"],
34                 "labels": labels,
35                 "last_updated": page["version"]["when"],
36                 "author": page["version"]["by"]
37                          ["displayName"],
38             },
39             source_type="runbook",
40             source_id=f"confluence:{page['id']}",
41         ))
42     return docs
43
44     def _html_to_text(self, html: str) -> str:
45         from bs4 import BeautifulSoup
46         return BeautifulSoup(html, "html.parser").get_text(
47             separator="\n")
48
49     class DbtSource(DocumentSource):
50         """Ingest documentation from a dbt project."""
51         def __init__(self, manifest_path: str):
52             self.manifest_path = manifest_path
53
54         def fetch_documents(self) -> list[Document]:
55             with open(self.manifest_path) as f:
56                 manifest = json.load(f)
57             docs = []
58             for node_id, node in manifest["nodes"].items():
59                 if node["resource_type"] != "model":
60                     continue
61                 content_parts = [
```

```

62         f"Model: {node['name']}",
63         f"Description: {node.get('description', 'N/A')}",
64         f"Schema: {node.get('schema', 'N/A')}",
65         f"Materialization: "
66         f"{node['config'].get('materialized', 'view')}",
67     ]
68     # Include column descriptions
69     for col_name, col in node.get("columns", {}).items():
70         desc = col.get("description", "")
71         content_parts.append(
72             f" Column {col_name}: {desc}")
73     docs.append(Document(
74         content="\n".join(content_parts),
75         metadata={
76             "model_name": node["name"],
77             "schema": node.get("schema"),
78             "tags": node.get("tags", []),
79             "depends_on": node.get("depends_on", {})
80                             .get("nodes", []),
81         },
82         source_type="catalog",
83         source_id=f"dbt:{node_id}",
84     ))
85     return docs
86
87 class IncidentReportSource(DocumentSource):
88     """Ingest incident reports from a directory."""
89     def __init__(self, directory: str):
90         self.directory = directory
91
92     def fetch_documents(self) -> list[Document]:
93         from pathlib import Path
94         docs = []
95         for path in Path(self.directory).glob("*.md"):
96             content = path.read_text()
97             # Parse YAML front matter if present
98             metadata = {}
99             if content.startswith("---"):
100                 parts = content.split("---", 2)
101                 if len(parts) >= 3:
102                     metadata = yaml.safe_load(parts[1])
103                     content = parts[2]
104             docs.append(Document(
105                 content=content,

```

```
106         metadata=metadata,  
107         source_type="incident",  
108         source_id=f"incident:{path.stem}",  
109     ))  
110     return docs
```

! Tip

Schedule your indexing pipeline to run on a cron (e.g., every 6 hours) and track document checksums to avoid re-indexing unchanged documents. For Confluence and similar wikis, use the “last modified” timestamp to fetch only updated pages since the last indexing run. Yes, you’re building a pipeline to feed a pipeline. Welcome to data engineering.

21.3 Embedding Models and Vector Stores

Embeddings are dense vector representations of text that capture semantic meaning. Two pieces of text with similar meanings will have embeddings that are close together in vector space, even if they use completely different words. This property enables *semantic search*—finding relevant content based on meaning rather than keyword matching.

Choosing an Embedding Model

Several factors influence embedding model selection:

- **Dimension size:** Higher dimensions (e.g., 1024 or 1536) capture more nuance but require more storage and compute. For most data engineering documentation, 1024 dimensions provide an excellent balance.
- **Context window:** Ensure the model’s context window accommodates your chunk sizes. Most modern embedding models support 512–8192 tokens.
- **Domain performance:** Some models perform better on technical content. Test with a sample of your actual documents.
- **Cost and latency:** Embedding is a high-volume operation during indexing. Factor in cost per token and throughput.

Code Example

Embedding service wrapper

```
1 import voyageai
2 from typing import Optional
3
4 class EmbeddingService:
5     """Wrapper for generating text embeddings."""
6
7     def __init__(self, model: str = "voyage-3",
8                 batch_size: int = 32):
9         self.client = voyageai.Client()
10        self.model = model
11        self.batch_size = batch_size
12
13    def embed_query(self, text: str) -> list[float]:
14        """Embed a single query string."""
15        result = self.client.embed(
16            [text], model=self.model,
17            input_type="query")
18        return result.embeddings[0]
19
20    def embed_documents(self,
21                      texts: list[str]) -> list[list[float]]:
22        """Embed a batch of document chunks."""
23        all_embeddings = []
24        for i in range(0, len(texts), self.batch_size):
25            batch = texts[i:i + self.batch_size]
26            result = self.client.embed(
27                batch, model=self.model,
28                input_type="document")
29            all_embeddings.extend(result.embeddings)
30        return all_embeddings
```

i Note

Many embedding APIs distinguish between “query” and “document” input types. Using the correct type improves retrieval quality because the model optimizes the embedding differently for short queries versus longer document passages. Always use `input_type="query"` for user questions and `input_type="document"` for chunks being indexed.

Vector Store Implementation

Code Example

Vector store with pgvector

```
1 import psycopg2, json
2 from psycopg2.extras import execute_values
3
4 class PgVectorStore:
5     """Vector store backed by PostgreSQL with pgvector."""
6
7     def __init__(self, connection_string: str,
8                 embedding_dim: int = 1024):
9         self.conn = psycopg2.connect(connection_string)
10        self.embedding_dim = embedding_dim
11        self._ensure_schema()
12
13    def _ensure_schema(self):
14        with self.conn.cursor() as cur:
15            cur.execute(
16                "CREATE EXTENSION IF NOT EXISTS vector;")
17            cur.execute(f"""
18                CREATE TABLE IF NOT EXISTS document_chunks (
19                    id SERIAL PRIMARY KEY,
20                    content TEXT NOT NULL,
21                    metadata JSONB NOT NULL DEFAULT '{{}}',
22                    source_type VARCHAR(50) NOT NULL,
23                    source_id VARCHAR(255) NOT NULL,
24                    chunk_index INTEGER NOT NULL,
25                    embedding vector({self.embedding_dim}),
26                    created_at TIMESTAMPTZ DEFAULT NOW(),
27                    UNIQUE (source_id, chunk_index)
28                );""")
29            cur.execute("""
30                CREATE INDEX IF NOT EXISTS idx_chunks_embedding
31                ON document_chunks
32                USING ivfflat (embedding vector_cosine_ops)
33                WITH (lists = 100);""")
34            cur.execute("""
35                CREATE INDEX IF NOT EXISTS idx_chunks_source
36                ON document_chunks (source_type, source_id);
37                """)
38        self.conn.commit()
```

```
39
40 def upsert(self, chunks: list[dict]):
41     """Insert or update chunks with embeddings."""
42     with self.conn.cursor() as cur:
43         for chunk in chunks:
44             cur.execute("""
45                 INSERT INTO document_chunks
46                     (content, metadata, source_type,
47                      source_id, chunk_index, embedding)
48                 VALUES (%s, %s, %s, %s, %s, %s)
49                 ON CONFLICT (source_id, chunk_index)
50                 DO UPDATE SET
51                     content = EXCLUDED.content,
52                     metadata = EXCLUDED.metadata,
53                     embedding = EXCLUDED.embedding,
54                     created_at = NOW()
55                 """, (
56                     chunk["content"],
57                     json.dumps(chunk["metadata"]),
58                     chunk["source_type"],
59                     chunk["source_id"],
60                     chunk["chunk_index"],
61                     chunk["embedding"],
62                 ))
63         self.conn.commit()
64
65 def search(self, query_embedding, top_k=10,
66            source_type_filter=None,
67            metadata_filter=None):
68     """Semantic search with optional filters."""
69     conditions = ["1=1"]
70     params = [str(query_embedding)]
71     if source_type_filter:
72         conditions.append("source_type = %s")
73         params.append(source_type_filter)
74     if metadata_filter:
75         for key, value in metadata_filter.items():
76             conditions.append(
77                 f"metadata->>' {key}' = %s")
78             params.append(str(value))
79     where_clause = " AND ".join(conditions)
80     params.extend([str(query_embedding), top_k])
81
82     with self.conn.cursor() as cur:
```

```

83         cur.execute(f """
84             SELECT content, metadata, source_type,
85                    source_id, chunk_index,
86                    1 - (embedding <=> %s::vector)
87                    AS similarity
88             FROM document_chunks
89             WHERE {where_clause}
90             ORDER BY embedding <=> %s::vector
91             LIMIT %s""" , params)
92     results = []
93     for rank, row in enumerate(cur.fetchall()):
94         chunk = Chunk(
95             content=row[0],
96             metadata=row[1],
97             chunk_index=row[4],
98             document_id=row[3])
99         results.append(RetrievalResult(
100             chunk=chunk,
101             relevance_score=float(row[5]),
102             rank=rank + 1))
103     return results
104
105     def delete_by_source(self, source_id: str):
106         """Remove all chunks for a given source."""
107         with self.conn.cursor() as cur:
108             cur.execute(
109                 "DELETE FROM document_chunks "
110                 "WHERE source_id = %s", (source_id,))
111         self.conn.commit()

```

Warning

The `ivfflat` index requires the table to have data before it can be effectively built. If you create the index on an empty table and then bulk-load data, query performance will be poor. Either create the index after the initial data load, or use `REINDEX` after bulk loading. For tables with fewer than 10,000 chunks, an exact (brute-force) search without an index may actually be faster.

21.4 Chunking Strategies for Technical Documentation

Chunking—the process of splitting documents into smaller pieces for embedding and retrieval—is one of the most impactful design decisions in a RAG system. It is also one of the least glamorous. Nobody has ever won an award for “Best Chunking Strategy,” but bad chunking has certainly caused award-winning postmortems. Poor chunking leads to irrelevant retrieval, which leads to wrong answers, regardless of how good your embedding model or LLM is.

Technical documents present unique chunking challenges:

- **Code blocks** should never be split mid-statement. A Python function split across two chunks loses its meaning in both.
- **Tables** are meaningless without their headers. Each table should be kept as a single chunk or include the header row in every chunk.
- **Step-by-step instructions** lose procedural context when split. Step 5 without steps 1–4 is often useless.
- **Configuration examples** are meaningless if truncated. A half-complete YAML config is worse than no config.

Recommended Strategies

Structure-aware chunking Parse the document’s structure (Markdown headers, HTML tags) and split at natural boundaries. Each chunk corresponds to a section or subsection. This preserves the author’s intended information grouping and works well for technical documentation.

Parent-child chunking Create small chunks (200–300 tokens) for retrieval precision, but link each to a larger parent chunk (800–1200 tokens). Search against the small chunks, but return the parent when assembling context for the LLM. This gives you the best of both worlds: precise matching and rich context.

Sliding window with overlap Use fixed-size windows (e.g., 500 tokens) with 10–20% overlap between consecutive chunks. The overlap ensures that information at chunk boundaries is not lost. This is the simplest strategy and works adequately for prose-heavy documents, but performs poorly on code-heavy content.

Code Example

Structure-aware Markdown chunker

```
1 import re
2 from dataclasses import dataclass
3
4 @dataclass
5 class MarkdownSection:
6     heading: str
7     level: int
8     content: str
9     parent_heading: str | None = None
10
11 class MarkdownChunker:
12     """Split Markdown documents at heading boundaries."""
13
14     def __init__(self, max_chunk_tokens: int = 800,
15                 min_chunk_tokens: int = 100):
16         self.max_tokens = max_chunk_tokens
17         self.min_tokens = min_chunk_tokens
18
19     def chunk(self, document: Document) -> list[Chunk]:
20         sections = self._parse_sections(document.content)
21         chunks = []
22         for i, section in enumerate(sections):
23             # Prepend heading hierarchy for context
24             header = section.heading
25             if section.parent_heading:
26                 header = (f"{section.parent_heading} > "
27                          f"{section.heading}")
28             chunk_content = f"# {header}\n\n{section.content}"
29
30             # Split oversized sections
31             if self._estimate_tokens(chunk_content) \
32                > self.max_tokens:
33                 sub_chunks = self._split_large_section(
34                     chunk_content)
35                 for j, sc in enumerate(sub_chunks):
36                     chunks.append(Chunk(
37                         content=sc,
38                         metadata={
39                             **document.metadata,
40                             "heading": header},
```

```
41         chunk_index=len(chunks),
42         document_id=document.source_id))
43     else:
44         chunks.append(Chunk(
45             content=chunk_content,
46             metadata={
47                 **document.metadata,
48                 "heading": header},
49             chunk_index=len(chunks),
50             document_id=document.source_id))
51     return chunks
52
53 def _parse_sections(self, text):
54     lines = text.split("\n")
55     sections = []
56     current_heading = "Introduction"
57     current_content = []
58     parent = None
59     current_level = 1
60
61     for line in lines:
62         heading_match = re.match(
63             r'^#{1,4}\s+(.+)$', line)
64         if heading_match:
65             if current_content:
66                 sections.append(MarkdownSection(
67                     heading=current_heading,
68                     level=current_level,
69                     content="\n".join(current_content),
70                     parent_heading=parent))
71                 level = len(heading_match.group(1))
72                 current_heading = heading_match.group(2)
73                 current_level = level
74                 if level > 1:
75                     parent = (sections[-1].heading
76                               if sections else None)
77             else:
78                 parent = None
79                 current_content = []
80         else:
81             current_content.append(line)
82
83     if current_content:
84         sections.append(MarkdownSection(
```

```
85         heading=current_heading,
86         level=current_level,
87         content="\n".join(current_content),
88         parent_heading=parent))
89     return sections
90
91     def _estimate_tokens(self, text: str) -> int:
92         return len(text) // 4 # rough estimate
93
94     def _split_large_section(self, text: str) -> list[str]:
95         paragraphs = text.split("\n\n")
96         chunks, current = [], []
97         current_size = 0
98         for para in paragraphs:
99             para_tokens = self._estimate_tokens(para)
100             if (current_size + para_tokens > self.max_tokens
101                 and current):
102                 chunks.append("\n\n".join(current))
103                 current = [para]
104                 current_size = para_tokens
105             else:
106                 current.append(para)
107                 current_size += para_tokens
108         if current:
109             chunks.append("\n\n".join(current))
110         return chunks
```

Tip

For code-heavy documentation, create separate chunks for code blocks and their surrounding explanatory text, linked by metadata. This lets the retrieval engine find code examples directly while keeping explanations available for context. Store the code language (Python, SQL, YAML) in chunk metadata so you can filter retrieval by language when appropriate.

21.5 RAG for Runbook Automation

Runbooks are the most immediately valuable content for a data engineering RAG system. Every team has runbooks—step-by-step procedures for handling common incidents, performing maintenance tasks, and executing deployments. But runbooks

suffer from chronic problems: they are scattered across wikis, they go stale, and engineers under pressure at 3 AM cannot find the right one.

A RAG-powered runbook assistant solves all three problems: it searches across all sources, it surfaces the most relevant procedure regardless of where it lives, and it synthesizes information from multiple runbooks when no single document covers the situation. It is, in essence, the colleague who has actually read every page of every runbook—a mythical creature that, until now, did not exist.

Code Example

RAG-powered runbook assistant

```
1 import anthropic
2
3 class RunbookRAGAssistant:
4     """A RAG assistant specialized for runbook queries."""
5
6     SYSTEM_PROMPT = """You are a data engineering operations
7 assistant with access to the team's runbooks, documentation,
8 and incident history.
9
10 When answering questions:
11 1. Search the knowledge base for relevant runbooks
12 2. Present steps in order with warnings and prerequisites
13 3. Always cite your sources with document titles
14 4. If you find conflicting information across sources,
15    flag the conflict and present both versions
16 5. If no relevant information exists, say so clearly
17    and suggest who to contact
18
19 Never fabricate procedures. If a step seems incomplete
20 in the source material, note it explicitly."""
21
22     def __init__(self, vector_store, embedding_service):
23         self.client = anthropic.Anthropic()
24         self.vector_store = vector_store
25         self.embedding_service = embedding_service
26
27     def answer(self, question, source_filter=None,
28               top_k=8):
29         """Answer a question using RAG retrieval."""
30         query_embedding = self.embedding_service.embed_query(
31             question)
```

```
32     results = self.vector_store.search(
33         query_embedding=query_embedding,
34         top_k=top_k,
35         source_type_filter=source_filter)
36
37     if not results:
38         return RAGResponse(
39             answer="No relevant documents found in the "
40                 "knowledge base for this question.",
41             sources=[], confidence=0.0, tokens_used=0)
42
43     # Build context from retrieved chunks
44     context_parts = []
45     for r in results:
46         source_info = (
47             f"[Source: "
48             f"{r.chunk.metadata.get('title',
49             r.chunk.document_id)} "
50             f"(relevance: {r.relevance_score:.2f})]"
51         )
52         context_parts.append(
53             f"{source_info}\n{r.chunk.content}")
54     context = "\n\n---\n\n".join(context_parts)
55
56     response = self.client.messages.create(
57         model="claude-sonnet-4-20250514",
58         max_tokens=4096,
59         system=self.SYSTEM_PROMPT,
60         messages=[{
61             "role": "user",
62             "content": (
63                 f"Context from knowledge base:"
64                 f"\n\n{context}\n\n---\n\n"
65                 f"Question: {question}")
66         })
67
68     # Calculate confidence from retrieval scores
69     top_score = results[0].relevance_score
70     avg_score = sum(
71         r.relevance_score for r in results[:3]
72     ) / min(3, len(results))
73
74     return RAGResponse(
75         answer=response.content[0].text,
76         sources=results,
```

```
76     confidence=(top_score * 0.6 + avg_score * 0.4),
77     tokens_used=(response.usage.input_tokens
78                 + response.usage.output_tokens))
```

21.6 RAG for Data Catalog Search

Traditional data catalog search relies on keyword matching—if you do not know the exact table name or column name, you cannot find what you need. RAG enables *semantic search* over your catalog, allowing users to describe what they need in natural language and find the right tables.

A user can ask “I need customer purchase history with product details” and find the `orders` table joined with `order_items` and `products`, even if they do not know any of those table names. The embedding model captures the semantic relationship between “purchase history” and “orders.”

Note

Semantic catalog search is particularly powerful for new team members who know *what* data they need but not *where* it lives. It also helps experienced engineers discover datasets they did not know existed—a common scenario in organizations with hundreds or thousands of tables, roughly 40% of which are named `temp_final_v2_REAL_USE_THIS_ONE`.

Code Example

Semantic catalog search with metadata enrichment

```
1 class CatalogSearchAssistant:
2     """Semantic search over data catalog entries."""
3
4     SYSTEM_PROMPT = """You are a data catalog assistant.
5     When a user describes the data they need:
6     1. Search for relevant tables and columns
7     2. Explain what each table contains and its freshness
8     3. Suggest joins between tables when appropriate
9     4. Note any data quality caveats or access restrictions
10    5. Provide example SQL queries to get started"""
11
```

```

12     def __init__(self, vector_store, embedding_service):
13         self.client = anthropic.Anthropic()
14         self.vector_store = vector_store
15         self.embedding_service = embedding_service
16
17     def search(self, description: str) -> RAGResponse:
18         embedding = self.embedding_service.embed_query(
19             description)
20         results = self.vector_store.search(
21             query_embedding=embedding,
22             top_k=15,
23             source_type_filter="catalog")
24
25         context = "\n\n".join(
26             f"Table: {r.chunk.metadata.get('model_name', '?')}"
27             f"\n{r.chunk.content}"
28             for r in results)
29
30         response = self.client.messages.create(
31             model="claude-sonnet-4-20250514",
32             max_tokens=2048,
33             system=self.SYSTEM_PROMPT,
34             messages=[{
35                 "role": "user",
36                 "content": (
37                     f"Available tables:\n\n{context}"
38                     f"\n\nI need: {description}")
39             })]
40         return RAGResponse(
41             answer=response.content[0].text,
42             sources=results,
43             confidence=results[0].relevance_score
44                 if results else 0.0,
45             tokens_used=(response.usage.input_tokens
46                 + response.usage.output_tokens))

```

21.7 Hybrid Search: Combining Vector and Keyword

Pure vector search excels at finding semantically similar content but can miss exact matches on technical terms, table names, or error codes. Pure keyword search (BM25) finds exact matches but misses semantic relationships. *Hybrid search*

combines both approaches for the best of both worlds.

Code Example

Hybrid search combining vector and BM25

```

1 class HybridSearcher:
2     """Combine vector similarity with keyword matching."""
3
4     def __init__(self, vector_store, embedding_service,
5                 vector_weight=0.7, keyword_weight=0.3):
6         self.vector_store = vector_store
7         self.embedding_service = embedding_service
8         self.vector_weight = vector_weight
9         self.keyword_weight = keyword_weight
10
11    def search(self, query: str, top_k: int = 10,
12              source_type_filter=None):
13        # Vector search
14        embedding = self.embedding_service.embed_query(query)
15        vector_results = self.vector_store.search(
16            query_embedding=embedding,
17            top_k=top_k * 2, # fetch more for re-ranking
18            source_type_filter=source_type_filter)
19
20        # Keyword search (using PostgreSQL full-text search)
21        keyword_results = self.vector_store.keyword_search(
22            query=query, top_k=top_k * 2,
23            source_type_filter=source_type_filter)
24
25        # Reciprocal Rank Fusion (RRF) to merge results
26        scores: dict[str, float] = {}
27        chunk_map: dict[str, RetrievalResult] = {}
28
29        for r in vector_results:
30            key = f"{r.chunk.document_id}:{r.chunk.chunk_index}"
31            rrf = self.vector_weight / (60 + r.rank)
32            scores[key] = scores.get(key, 0) + rrf
33            chunk_map[key] = r
34
35        for r in keyword_results:
36            key = f"{r.chunk.document_id}:{r.chunk.chunk_index}"
37            rrf = self.keyword_weight / (60 + r.rank)
38            scores[key] = scores.get(key, 0) + rrf

```

```
39         if key not in chunk_map:
40             chunk_map[key] = r
41
42         # Sort by combined score and return top_k
43         ranked = sorted(scores.items(),
44                         key=lambda x: x[1], reverse=True)
45     return [
46         RetrievalResult(
47             chunk=chunk_map[key].chunk,
48             relevance_score=score,
49             rank=i + 1)
50     for i, (key, score) in enumerate(ranked[:top_k])
51     ]
```

! Tip

Reciprocal Rank Fusion (RRF) is a simple but effective method for combining ranked lists from different retrieval strategies. The constant $k=60$ in the RRF formula $1/(k + \text{rank})$ controls how much weight is given to top-ranked results versus lower-ranked ones. Higher values of k flatten the distribution; lower values amplify the top results.

21.8 Evaluation Metrics for RAG Systems

A RAG system that cannot be measured cannot be improved. Evaluation spans both the retrieval stage and the generation stage, and each requires different metrics.

Retrieval Metrics

Precision@K Of the K chunks retrieved, how many are actually relevant? Low precision means the LLM receives noise that can lead to incorrect answers.

Recall@K Of all relevant chunks in the corpus, how many appear in the top K results? Low recall means the LLM is missing information it needs.

Mean Reciprocal Rank (MRR) How high does the first relevant result rank? An MRR close to 1.0 means the most relevant chunk is almost always the top result.

Generation Metrics

Faithfulness Does the answer accurately reflect the retrieved context? Hallucinated information not present in the retrieved chunks indicates a faithfulness problem.

Relevance Does the answer actually address the user's question? A faithful but off-topic answer is not useful.

Completeness Does the answer cover all aspects of the question? Partial answers indicate retrieval gaps or generation truncation.

Attribution Can every claim in the answer be traced to a specific source? Unattributed claims are unverifiable.

Code Example

Automated RAG evaluation

```
1 class RAGEvaluator:
2     """Evaluate RAG system quality."""
3
4     def __init__(self, rag_assistant):
5         self.assistant = rag_assistant
6         self.client = anthropic.Anthropic()
7
8     def evaluate_batch(self,
9                       test_cases: list[dict]) -> dict:
10        """Run evaluation on a set of test cases.
11
12        Each test case has: question, expected_answer,
13        relevant_source_ids.
14        """
15        results = {
16            "precision_at_5": [],
17            "recall_at_5": [],
18            "faithfulness": [],
19            "relevance": [],
20        }
21        for case in test_cases:
22            response = self.assistant.answer(
23                case["question"], top_k=5)
24
25            # Retrieval metrics
```

```
26     retrieved_ids = {
27         r.chunk.document_id
28         for r in response.sources[:5]}
29     relevant_ids = set(case["relevant_source_ids"])
30     hits = retrieved_ids & relevant_ids
31     precision = (len(hits) / len(retrieved_ids)
32                 if retrieved_ids else 0)
33     recall = (len(hits) / len(relevant_ids)
34              if relevant_ids else 0)
35     results["precision_at_5"].append(precision)
36     results["recall_at_5"].append(recall)
37
38     # LLM-judged faithfulness and relevance
39     judge_resp = self.client.messages.create(
40         model="claude-sonnet-4-20250514",
41         max_tokens=200,
42         messages=[{"role": "user", "content":
43                   f"Rate 1-5.\n"
44                   f"Question: {case['question']}\n"
45                   f"Answer: {response.answer}\n"
46                   f"Expected: {case['expected_answer']}\n"
47                   f"Faithfulness (1-5): "
48                   f"Relevance (1-5):"}])
49     # Parse scores from response
50     scores = self._parse_scores(
51         judge_resp.content[0].text)
52     results["faithfulness"].append(
53         scores.get("faithfulness", 3))
54     results["relevance"].append(
55         scores.get("relevance", 3))
56
57     # Average all metrics
58     return {k: sum(v) / len(v)
59            for k, v in results.items() if v}
60
61 def _parse_scores(self, text):
62     scores = {}
63     for line in text.split("\n"):
64         lower = line.lower()
65         for metric in ["faithfulness", "relevance"]:
66             if metric in lower:
67                 nums = re.findall(r'\d', line)
68                 if nums:
69                     scores[metric] = int(nums[0])
```

```
70 return scores
```

⚠ Warning

LLM-as-judge evaluations are useful for rapid iteration but not perfect. Having an LLM grade its own work is like having a student grade their own exam—surprisingly useful, but don't bet your GPA on it. Models can exhibit position bias (favoring the first option), verbosity bias (preferring longer answers), and self-preference bias (rating their own outputs higher). For critical production systems, combine LLM-based evaluation with human evaluation on a regular cadence—monthly reviews of 50–100 randomly sampled queries provide a reliable ground truth signal.

📘 Note

Build your evaluation dataset incrementally. Start with 30 question-answer pairs covering your most common query types. Add 10 new pairs each month, focusing on queries where the system performed poorly. After six months, you will have a comprehensive benchmark that captures the diversity of real-world usage patterns.

21.9 Production Deployment Considerations

Moving a RAG system from prototype to production introduces several concerns that do not arise during development:

- **Latency:** Embedding the query, searching the vector store, and generating the response adds 2–5 seconds of latency. For interactive use cases, consider pre-computing embeddings for common queries and caching frequent results.
- **Freshness:** Documents change. Your indexing pipeline must detect and re-index modified documents within your freshness SLA. Stale retrieval results erode user trust quickly.
- **Access control:** Not all users should see all chunks. Implement metadata-based filtering that respects your organization's data classification and role-based access control.

- **Cost management:** Each query involves an embedding API call and an LLM API call. At 100 queries per day, costs are negligible; at 10,000 queries per day, they require budgeting and optimization.

Tip

Implement a feedback mechanism where users can thumbs-up or thumbs-down RAG responses. Track the thumbs-down rate by source type and query category. This provides a real-time quality signal that complements your formal evaluation metrics and helps prioritize improvements.

21.10 Exercises

1. **Multi-Source Indexer:** Build an indexing pipeline that pulls documents from at least three sources (Confluence, a dbt project, and incident reports). Index into a pgvector store and verify that search returns relevant results for ten test queries spanning all three source types.
2. **Hybrid Search:** Extend the vector store with PostgreSQL full-text search to implement hybrid retrieval using Reciprocal Rank Fusion. Compare retrieval quality (Precision@5 and Recall@5) of pure vector, pure keyword, and hybrid search on 20 test queries. Document which query types benefit most from each approach.
3. **Conversational RAG:** Extend the runbook assistant to handle multi-turn conversations where follow-up questions depend on previous context. For example: “What tables have customer data?” followed by “Which of those are updated daily?” The system should rewrite follow-up questions to be self-contained before embedding.
4. **Evaluation Dataset:** Create 30+ question-answer pairs spanning runbooks (operational), catalog entries (discovery), and incident reports (diagnostic). Run the evaluator, identify the weakest category, and implement at least one improvement (better chunking, hybrid search, or prompt refinement) that measurably improves scores.
5. **Access-Controlled RAG:** Implement a metadata-based access control system where chunks are tagged with team ownership and sensitivity levels. Build a search endpoint that accepts a user role and filters results accordingly. Verify that an “analyst” role cannot retrieve ops-only runbooks.

22

Real-Time Data Processing with Claude Code

“Real-time data processing: because batch was too easy, and you had too much free time.”

Real-time data processing is one of the most demanding domains in data engineering. This chapter explores how CLAUDE CODE can assist in building, optimizing, and maintaining real-time systems—from generating stream processing logic in Apache Flink and ksqlDB, to designing anomaly detection algorithms, to building fraud detection pipelines. We also examine the practical constraints of integrating LLMs into latency-sensitive streaming architectures. Spoiler: you cannot call an LLM for every event when you have a million events per second. Math is cruel that way.

22.1 Real-Time Data Engineering Fundamentals

Defining Real-Time

We distinguish three tiers of “real-time”:

- **Hard real-time:** Processing under 10ms. Common in embedded systems and trading. LLMs are not suitable for this path. CLAUDE CODE is many things, but “sub-10ms” is not one of them.

- **Soft real-time:** Processing under 1 second, with occasional delays tolerable. Common in web applications and monitoring.
- **Near real-time:** Processing within seconds to minutes. Common in dashboards, alerting, and data sync. This is where CLAUDE CODE is most useful in the processing loop. Also where most “real-time” systems actually live, despite what the marketing page says.

Stream Processing Architecture Patterns

Event sourcing. Every state change is captured as an immutable event; current state is derived by replaying events.

CQRS. Write and read operations use separate models, often with different stores optimized for their access patterns.

Kappa architecture. All data flows through a single stream processing layer, simplifying the architecture compared to Lambda’s separate batch and speed layers.

Key Concepts in Stream Processing

- **Event time vs. processing time:** Correctly handling the gap between when events occur and when they are processed is a core challenge.
- **Windowing:** Grouping events into finite sets—tumbling (fixed, non-overlapping), sliding (fixed, overlapping), and session (gap-based) windows.
- **Watermarks:** Mechanisms for tracking event-time progress and handling late-arriving data.
- **Exactly-once semantics:** Guaranteeing each event is processed exactly once despite failures and retries.
- **Backpressure:** Handling situations where producers outpace consumers.

22.2 Kafka, Kinesis, and Pub/Sub Overview

Apache Kafka is the most widely adopted event streaming platform, offering durable ordered logs, consumer groups, Schema Registry, Kafka Connect, and ksqlDB. *Amazon Kinesis* differs with shard-based capacity, 24-hour default retention, and native AWS integration. *Google Cloud Pub/Sub* provides automatic scaling, at-least-once delivery, and native Dataflow integration.

Tip

When asking CLAUDE CODE to generate stream processing code, always specify which platform you are using. Saying “write me a streaming pipeline” without specifying the platform is like ordering “food” at a restaurant. The APIs, configuration parameters, and best practices differ significantly.

22.3 Claude Code for Stream Processing Logic Generation

Stream processing code is complex, with subtle correctness requirements around windowing, state management, and error handling. CLAUDE CODE can generate both Flink applications and ksqlDB queries from high-level descriptions.

Generating Apache Flink Applications

Apache Flink supports complex event processing, stateful computations, and exactly-once semantics. CLAUDE CODE can generate Flink application logic from high-level descriptions. The key elements to specify in your prompt are: the source topic and serialization format, the windowing strategy, the aggregation logic, and the sink configuration. Always include checkpointing configuration for production—CLAUDE CODE generates this correctly when you mention exactly-once requirements.

Code Example

```
1 from pyflink.datastream import StreamExecutionEnvironment
2 from pyflink.datastream.window import TumblingEventTimeWindows
3 from pyflink.common.time import Time
4 from pyflink.common.watermark_strategy import WatermarkStrategy
5
6 def build_order_analytics_pipeline():
7     """Build a Flink streaming pipeline for order analytics."""
8     env = StreamExecutionEnvironment.get_execution_environment()
9     env.set_parallelism(4)
10    env.enable_checkpointing(60000)
11
12    from pyflink.datastream.connectors.kafka import (
13        KafkaSource, KafkaOffsetsInitializer,
14    )
15    from pyflink.common.serialization import SimpleStringSchema
```

```

16
17 kafka_source = (
18     KafkaSource.builder()
19     .set_bootstrap_servers("kafka:9092")
20     .set_topics("orders")
21     .set_group_id("order-analytics")
22     .set_starting_offsets(KafkaOffsetsInitializer.earliest())
23     .set_value_only_deserializer(SimpleStringSchema())
24     .build()
25 )
26
27 watermark_strategy = (
28     WatermarkStrategy
29     .for_bounded_out_of_orderness(Time.seconds(30))
30 )
31
32 order_stream = (
33     env.from_source(kafka_source, watermark_strategy, "Orders")
34     .map(OrderParser())
35     .key_by(lambda x: x[0]) # Key by customer_id
36     .window(TumblingEventTimeWindows.of(Time.minutes(5)))
37     .process(OrderAggregator())
38 )
39 env.execute("Order Analytics Pipeline")

```

Generating ksqlDB Queries

ksqlDB provides a SQL-like interface for stream processing on Kafka. CLAUDE CODE is particularly effective here because of the similarity to standard SQL:

Code Example

```

1  -- ksqlDB stream processing for real-time analytics
2
3  -- Create a stream from the raw orders topic
4  CREATE STREAM orders_raw (
5      order_id VARCHAR KEY,
6      customer_id VARCHAR,
7      amount DOUBLE,
8      status VARCHAR,

```

```
9     product_category VARCHAR,  
10     event_time TIMESTAMP  
11 ) WITH (  
12     KAFKA_TOPIC = 'orders',  
13     VALUE_FORMAT = 'JSON',  
14     TIMESTAMP = 'event_time'  
15 );  
16  
17 -- Filtered stream of completed orders  
18 CREATE STREAM completed_orders AS  
19     SELECT order_id, customer_id, amount,  
20            product_category, event_time  
21     FROM orders_raw  
22     WHERE status = 'completed'  
23     EMIT CHANGES;  
24  
25 -- Tumbling window aggregation  
26 CREATE TABLE revenue_per_category_5min AS  
27     SELECT  
28         product_category,  
29         WINDOWSTART AS window_start,  
30         COUNT(*) AS order_count,  
31         SUM(amount) AS total_revenue,  
32         AVG(amount) AS avg_order_value  
33     FROM completed_orders  
34     WINDOW TUMBLING (SIZE 5 MINUTES)  
35     GROUP BY product_category  
36     EMIT CHANGES;  
37  
38 -- Detect rapid ordering patterns (potential fraud)  
39 CREATE TABLE rapid_orders AS  
40     SELECT  
41         customer_id,  
42         WINDOWSTART AS window_start,  
43         COUNT(*) AS order_count,  
44         SUM(amount) AS total_amount  
45     FROM orders_raw  
46     WINDOW HOPPING (SIZE 5 MINUTES, ADVANCE BY 1 MINUTE)  
47     GROUP BY customer_id  
48     HAVING COUNT(*) >= 10  
49     EMIT CHANGES;
```

22.4 Real-Time Anomaly Detection with Claude Code

Anomaly detection in streaming data requires statistical methods, domain knowledge, and adaptive thresholds. CLAUDE CODE can help design detection logic and tune parameters.

Code Example

```
1 import math
2 from collections import deque
3 from dataclasses import dataclass
4 from datetime import datetime
5
6
7 @dataclass
8 class AnomalyResult:
9     is_anomaly: bool
10    value: float
11    expected_range: tuple[float, float]
12    z_score: float
13    timestamp: datetime
14    metric_name: str
15    severity: str = "info"
16
17
18 class StreamingAnomalyDetector:
19     """Real-time anomaly detection using Welford's algorithm."""
20
21     def __init__(
22         self, window_size: int = 1000,
23         z_threshold: float = 3.0,
24         min_samples: int = 30,
25     ):
26         self.window_size = window_size
27         self.z_threshold = z_threshold
28         self.min_samples = min_samples
29         self.count = 0
30         self.mean = 0.0
31         self.m2 = 0.0
32         self.window: deque[float] = deque(maxlen=window_size)
33
34     def update(
35         self, value: float, timestamp: datetime,
```

```
36     metric_name: str = "default",
37 ) -> AnomalyResult:
38     """Process a new value and check for anomalies."""
39     self.count += 1
40     delta = value - self.mean
41     self.mean += delta / self.count
42     delta2 = value - self.mean
43     self.m2 += delta * delta2
44     self.window.append(value)
45
46     if self.count < self.min_samples:
47         return AnomalyResult(
48             is_anomaly=False, value=value,
49             expected_range=(0.0, 0.0), z_score=0.0,
50             timestamp=timestamp, metric_name=metric_name,
51         )
52
53     variance = self.m2 / (self.count - 1)
54     std_dev = math.sqrt(variance) if variance > 0 else 1.0
55     z_score = abs((value - self.mean) / std_dev)
56     is_anomaly = z_score > self.z_threshold
57
58     lower = self.mean - self.z_threshold * std_dev
59     upper = self.mean + self.z_threshold * std_dev
60
61     severity = "info"
62     if is_anomaly:
63         if z_score > self.z_threshold * 2:
64             severity = "critical"
65         elif z_score > self.z_threshold * 1.5:
66             severity = "warning"
67
68     return AnomalyResult(
69         is_anomaly=is_anomaly, value=value,
70         expected_range=(round(lower, 2), round(upper, 2)),
71         z_score=round(z_score, 4),
72         timestamp=timestamp, metric_name=metric_name,
73         severity=severity,
74     )
```

! Tip

When using CLAUDE CODE to design anomaly detection, provide sample data characteristics: typical value ranges, seasonal patterns, expected event frequency, and examples of both normal and anomalous behavior. Without this context, CLAUDE CODE will design a perfectly generic detector that flags everything or nothing. Both are technically anomaly detectors; neither is useful.

22.5 Event Schema Validation and Evolution

In streaming systems, schema management is critical because producers and consumers are decoupled. CLAUDE CODE can assist in designing schema evolution strategies and generating validation logic.

Code Example

```
1 from enum import Enum
2
3
4 class Compatibility(Enum):
5     BACKWARD = "backward"
6     FORWARD = "forward"
7     FULL = "full"
8     NONE = "none"
9
10
11 class SchemaEvolutionValidator:
12     """Validate schema changes for streaming compatibility."""
13
14     def __init__(self, compatibility: Compatibility = Compatibility.FULL):
15         self.compatibility = compatibility
16
17     def validate_evolution(self, old_schema: dict, new_schema: dict)
18     -> dict:
19         """Validate schema change for compatibility."""
20         issues = []
21         old_fields = {f["name"]: f for f in old_schema.get("fields",
22         [])}
23         new_fields = {f["name"]: f for f in new_schema.get("fields",
24         [])}
```

```
22
23     removed = set(old_fields) - set(new_fields)
24     if removed and self.compatibility in (
25         Compatibility.BACKWARD, Compatibility.FULL
26     ):
27         issues.append({
28             "type": "removed_field", "severity": "error",
29             "fields": list(removed),
30             "message": f"Removing fields {removed} breaks "
31                 f"{self.compatibility.value} compatibility
32     .",
33         })
34
35     added = set(new_fields) - set(old_fields)
36     for field_name in added:
37         field_def = new_fields[field_name]
38         has_default = "default" in field_def
39         is_nullable = isinstance(field_def.get("type"), list) and
40             "null" in field_def["type"]
41         if not has_default and not is_nullable:
42             if self.compatibility in (Compatibility.FORWARD,
43 Compatibility.FULL):
44                 issues.append({
45                     "type": "required_field_added", "severity": "
46 error",
47                     "field": field_name,
48                     "message": f"Adding required field '{
49 field_name}' without "
50                         f"a default breaks compatibility."
51                 },
52             )
53
54     is_compatible = not any(i["severity"] == "error" for i in
55 issues)
56     return {
57         "compatible": is_compatible,
58         "compatibility_mode": self.compatibility.value,
59         "issues": issues,
60         "added_fields": list(added),
61         "removed_fields": list(removed),
62     }
```

22.6 Building Real-Time Data Enrichment Pipelines

Data enrichment adds context and value to raw events by joining them with reference data, computing derived fields, or augmenting with external information. Common enrichment types include:

- **Lookup enrichment:** Join events with reference data (customer profiles, product catalogs) using a cache for performance.
- **Computed enrichment:** Derive new fields from existing ones (e.g., categorizing amounts into buckets).
- **Geographic enrichment:** Resolve IP addresses to locations using a GeoIP database.

When using CLAUDE CODE to generate enrichment logic, specify the cache TTL, the fallback behavior when lookups fail, and how to handle null values in key fields. CLAUDE CODE consistently generates defensive code for these cases when prompted explicitly.

22.7 Real-Time Data Quality Monitoring

Data quality monitoring in streaming systems requires continuous validation as data flows through the pipeline. Key quality dimensions to monitor in real time:

- **Completeness:** NULL rates per field per time window.
- **Type conformance:** Fields matching expected types.
- **Range validity:** Numeric values within acceptable bounds.
- **Uniqueness:** Duplicate detection within a sliding window.
- **Timeliness:** Delay between event time and processing time.

Implement a windowed quality monitor that accumulates violations per window, computes rates, and generates alerts when thresholds are breached. Use CLAUDE CODE to generate the quality rule configurations—describe your schema and business rules, and it will produce comprehensive validation logic including edge cases you might miss.

Note

Track quality scores over time to detect gradual degradation. A sudden drop indicates a source system change (or someone deploying on Friday). A gradual decline suggests data drift—the data equivalent of entropy slowly winning, as it always does.

22.8 Full Example: Fraud Detection Pipeline

This section presents a fraud detection pipeline combining rule-based detection, statistical anomaly detection, and pattern matching.

Code Example

```
1  """Real-time fraud detection pipeline."""
2  import time
3  from datetime import datetime, timedelta
4  from dataclasses import dataclass
5  from collections import defaultdict
6  from enum import Enum
7
8
9  class FraudSignal(Enum):
10     VELOCITY = "velocity"
11     AMOUNT = "amount_anomaly"
12     GEOGRAPHY = "geo_anomaly"
13
14
15  @dataclass
16  class FraudScore:
17     transaction_id: str
18     customer_id: str
19     score: float
20     signals: list[dict]
21     decision: str # "allow", "review", "block"
22     timestamp: datetime
23
24
25  class FraudDetectionPipeline:
26     """Real-time fraud detection with multiple signal types."""
27
```

```

28 def __init__(self, config: dict):
29     self.config = config
30     self.thresholds = config.get("thresholds", {
31         "block": 0.8, "review": 0.5,
32     })
33     self.customer_history: dict[str, list[dict]] = defaultdict(
34     list)
35     self.amount_detectors: dict[str, StreamingAnomalyDetector] =
36     {}
37     self.velocity_window = config.get("velocity_window_minutes",
38     60)
39
40 def process_transaction(self, transaction: dict) -> FraudScore:
41     """Process a transaction and return fraud assessment."""
42     start_time = time.time()
43     customer_id = transaction["customer_id"]
44     amount = float(transaction["amount"])
45     timestamp = datetime.fromisoformat(transaction["timestamp"])
46
47     signals = []
48
49     # Velocity check
50     cutoff = timestamp - timedelta(minutes=self.velocity_window)
51     recent = [t for t in self.customer_history[customer_id]
52               if t["timestamp"] > cutoff]
53     velocity_limit = self.config.get("velocity_limit", 10)
54     if len(recent) >= velocity_limit:
55         signals.append({
56             "signal": FraudSignal.VELOCITY.value,
57             "weight": 0.3,
58             "confidence": min(len(recent) / velocity_limit, 2.0)
59         } / 2.0,
60     })
61
62     # Amount anomaly
63     if customer_id not in self.amount_detectors:
64         self.amount_detectors[customer_id] =
65         StreamingAnomalyDetector(
66             window_size=100, z_threshold=2.5, min_samples=5,
67         )
68     result = self.amount_detectors[customer_id].update(
69         amount, timestamp, "amount"
70     )
71     if result.is_anomaly:

```

```
67         signals.append({
68             "signal": FraudSignal.AMOUNT.value,
69             "weight": 0.25,
70             "confidence": min(result.z_score / 5.0, 1.0),
71         })
72
73         # Composite score
74         if signals:
75             weighted_sum = sum(s["weight"] * s["confidence"] for s in
signals)
76             max_possible = sum(s["weight"] for s in signals)
77             score = (weighted_sum / max_possible) + min(len(signals)
* 0.1, 0.3)
78             score = min(score, 1.0)
79         else:
80             score = 0.0
81
82         if score >= self.thresholds["block"]:
83             decision = "block"
84         elif score >= self.thresholds["review"]:
85             decision = "review"
86         else:
87             decision = "allow"
88
89         # Update state
90         self.customer_history[customer_id].append({
91             **transaction, "timestamp": timestamp,
92         })
93         self.customer_history[customer_id] = [
94             t for t in self.customer_history[customer_id]
95             if t["timestamp"] > timestamp - timedelta(hours=24)
96         ]
97
98         return FraudScore(
99             transaction_id=transaction["transaction_id"],
100             customer_id=customer_id, score=round(score, 4),
101             signals=signals, decision=decision, timestamp=timestamp,
102         )
```

22.9 Latency and Cost Considerations for LLM-in-the-Loop Streaming

A typical CLAUDE CODE API call takes 500ms–5s. For streaming pipelines, this has significant implications:

- **Synchronous enrichment:** Acceptable only for near-real-time pipelines with second-scale latency budgets.
- **Asynchronous enrichment:** Fire API calls asynchronously and join results back later. Maintains low primary-path latency but adds complexity.
- **Batch micro-enrichment:** Collect 10–50 events per CLAUDE CODE call to amortize per-call overhead, at the cost of batching latency.

Warning

Never place a synchronous CLAUDE CODE API call in the critical path of a high-throughput streaming pipeline. At 1,000 events/second, even 1-second latency requires 1,000 concurrent calls—neither cost-effective nor reliable. Your CFO will find you. Your CFO will have questions.

When Not to Use LLMs in Streaming

1. **Sub-100ms latency requirements:** LLM API latency makes this impossible.
2. **Throughput above 10,000 events/second:** Cost and concurrency become prohibitive.
3. **Deterministic processing requirements:** LLMs produce non-deterministic outputs.
4. **Simple rule-based logic:** An LLM adds unnecessary complexity and cost.

The most effective pattern is using CLAUDE CODE to *generate* stream processing logic that then runs as deterministic code at full streaming speed. Reserve runtime LLM calls for low-volume, high-value enrichment.

! Tip

For cost-sensitive streaming, use CLAUDE CODE Haiku for high-volume classification and reserve Sonnet or Opus for complex reasoning. The cost difference can be an order of magnitude. Think of it as economy vs. first class: both get you there, but one lets you keep your house.

22.10 Exercises

Exercises

1. **Stream Processing Generation:** Use CLAUDE CODE to generate a ksqlDB application that sessionizes web clickstream events with a 30-minute gap, computes page views per session, identifies sessions with 50+ views as potential bots, and computes real-time conversion rates.
2. **Anomaly Detection Tuning:** Process a synthetic dataset through `StreamingAnomalyDetector` with different `z_threshold` (2.0, 2.5, 3.0, 3.5) and `window_size` (100, 500, 1000) values. Report precision, recall, and F1 for each configuration.
3. **Schema Evolution Challenge:** Design an Avro schema for payment events. Simulate three evolution rounds: (a) add optional `currency`, (b) rename `amount` to `payment_amount`, (c) add required `idempotency_key`. Analyze compatibility for each.
4. **Fraud Detection Enhancement:** Extend the fraud pipeline with a “card testing” signal (small transactions followed by a large one) and a “merchant pattern” signal (unfamiliar merchants during unusual hours).
5. **End-to-End Pipeline:** Build a complete streaming pipeline that generates synthetic orders, monitors quality, enriches with customer data, detects anomalies, and writes results to an output topic using Docker Compose with Kafka.

23

Multi-Modal Data Engineering

“Multi-modal data engineering: because regular data wasn’t complicated enough, now we’re feeding CLAUDE CODE pictures of whiteboards.”

Data engineering has long been dominated by structured rows and columns. But enterprises generate vast amounts of unstructured, multi-modal data: images, audio recordings, PDFs, video, and scanned documents. With native vision capabilities and sophisticated reasoning, CLAUDE CODE enables data engineers to build multi-modal ETL pipelines that treat these sources as first-class data. This chapter provides architectural patterns, production code, and operational guidance for building these pipelines at scale. Finally, a use for all those photos of whiteboard diagrams you took with the sincere intention of “digitizing later.”

23.1 Multi-Modal Data in the Modern Data Stack

Architectural Patterns for Multi-Modal Pipelines

Three primary patterns exist for integrating multi-modal data:

Pattern 1: Extract-then-load. Files are processed by CLAUDE CODE before entering the warehouse. Only extracted structured data lands in the warehouse; raw files remain in object storage. This is the simplest and most common pattern.

Pattern 2: Load-then-extract. Raw files are loaded into a lakehouse first (e.g., as external tables pointing to S3), and extraction happens as a transformation step.

Pattern 3: Hybrid streaming. Multi-modal data arrives via a streaming platform, is processed in near-real-time by CLAUDE CODE, and extracted data is written to both streaming and batch sinks.

Cost Considerations

- **Image tokens:** A 1024×1024 image typically uses around 1,600 tokens. Down-scaling reduces costs dramatically.
- **Audio preprocessing:** CLAUDE CODE does not natively process audio—transcribe first (Whisper, AWS Transcribe), then analyze the transcript.
- **Caching:** Prompt caching can reduce costs by up to 90% when processing many items with the same instructions.

! Tip

Always calculate cost per unit (per image, per document, per audio minute) before scaling to production. Run a sample of 100–500 items and measure actual token consumption. Many a promising multi-modal project has died on the spreadsheet of “wait, it costs HOW much per image?”

23.2 Image Processing Pipelines with Claude Code Vision

CLAUDE CODE’s vision capabilities allow direct image analysis—reading text, describing content, and extracting structured data without custom computer vision models.

Sending Images to the API

Code Example

```
1 import anthropic
2 import base64
3 from pathlib import Path
4
5
6 def encode_image(image_path: str) -> tuple[str, str]:
7     """Encode an image file to base64 with media type detection."""
```

```
8     path = Path(image_path)
9     media_types = {
10         ".jpg": "image/jpeg", ".jpeg": "image/jpeg",
11         ".png": "image/png", ".gif": "image/gif",
12         ".webp": "image/webp",
13     }
14     media_type = media_types.get(path.suffix.lower(), "image/jpeg")
15     with open(image_path, "rb") as f:
16         image_data = base64.standard_b64encode(f.read()).decode("utf
17 -8")
18     return image_data, media_type
19
20 def extract_from_image(
21     client: anthropic.Anthropic,
22     image_path: str,
23     extraction_prompt: str,
24     model: str = "claude-sonnet-4-20250514",
25 ) -> str:
26     """Send an image to Claude Code and extract structured data."""
27     image_data, media_type = encode_image(image_path)
28     message = client.messages.create(
29         model=model, max_tokens=4096,
30         messages=[{
31             "role": "user",
32             "content": [
33                 {"type": "image", "source": {
34                     "type": "base64",
35                     "media_type": media_type,
36                     "data": image_data,
37                 }},
38                 {"type": "text", "text": extraction_prompt},
39             ],
40         }],
41     )
42     return message.content[0].text
```

Image Preprocessing for Cost Optimization

Code Example

```
1 from PIL import Image
2 import io
3 import base64
4
5
6 def preprocess_image(
7     image_bytes: bytes,
8     max_dimension: int = 1024,
9     quality: int = 85,
10    output_format: str = "JPEG",
11 ) -> tuple[str, str]:
12     """Resize and compress an image for optimal API usage."""
13     img = Image.open(io.BytesIO(image_bytes))
14     if img.mode == "RGBA" and output_format == "JPEG":
15         background = Image.new("RGB", img.size, (255, 255, 255))
16         background.paste(img, mask=img.split()[3])
17         img = background
18
19     width, height = img.size
20     if max(width, height) > max_dimension:
21         ratio = max_dimension / max(width, height)
22         img = img.resize(
23             (int(width * ratio), int(height * ratio)),
24             Image.LANCZOS,
25         )
26
27     buffer = io.BytesIO()
28     img.save(buffer, format=output_format, quality=quality)
29     encoded = base64.standard_b64encode(buffer.getvalue()).decode("
30 utf-8")
31     format_to_media = {"JPEG": "image/jpeg", "PNG": "image/png", "
32     WEBP": "image/webp"}
33     return encoded, format_to_media.get(output_format, "image/jpeg")
```

Note

For most extraction tasks (receipts, product attributes, charts), resizing to 1024 pixels on the longest side reduces token consumption by 60–75% with excellent results. Use higher resolutions only for fine-grained details.

23.3 Audio Transcription and Analysis Workflows

The standard pattern is a two-stage pipeline: transcribe audio using a speech-to-text service, then send the transcript to CLAUDE CODE for structured extraction.

Code Example

```
1 import anthropic
2 import whisper
3 import json
4 from dataclasses import dataclass
5 from typing import Optional
6
7
8 @dataclass
9 class AudioAnalysis:
10     file_path: str
11     duration_seconds: float
12     transcript: str
13     summary: str
14     topics: list[str]
15     sentiment: str
16     action_items: list[str]
17
18
19 AUDIO_ANALYSIS_PROMPT = """Analyze this transcript and return JSON:
20 {
21     "summary": "2-3 sentence summary",
22     "topics": ["main", "topics"],
23     "sentiment": "positive/negative/neutral/mixed",
24     "action_items": ["list of action items mentioned"],
25     "language": "detected language"
26 }
27
28 TRANSCRIPT:
```

```
29 {transcript}
30
31 Return ONLY valid JSON. """
32
33
34 class AudioAnalysisPipeline:
35     def __init__(self, whisper_model: str = "base"):
36         self.whisper_model = whisper.load_model(whisper_model)
37         self.claude = anthropic.Anthropic()
38
39     def transcribe(self, audio_path: str) -> dict:
40         result = self.whisper_model.transcribe(audio_path)
41         return {"text": result["text"], "language": result.get("
language", "unknown")}
42
43     def analyze_transcript(self, transcript: str) -> dict:
44         prompt = AUDIO_ANALYSIS_PROMPT.format(transcript=transcript)
45         message = self.claude.messages.create(
46             model="claude-sonnet-4-20250514", max_tokens=2048,
47             messages=[{"role": "user", "content": prompt}],
48         )
49         return json.loads(message.content[0].text)
```

Handling Long Audio Files

For transcripts exceeding the context window, split into overlapping chunks, process each independently, then merge:

Code Example

```
1 def chunk_transcript(
2     transcript: str,
3     max_chars: int = 150000,
4     overlap_chars: int = 5000,
5 ) -> list[str]:
6     """Split a long transcript into overlapping chunks."""
7     if len(transcript) <= max_chars:
8         return [transcript]
9     chunks = []
10    start = 0
```

```

11     while start < len(transcript):
12         end = start + max_chars
13         if end < len(transcript):
14             boundary = transcript.rfind(".", start + max_chars -
15             1000, end)
16             if boundary > start:
17                 end = boundary + 1
18             chunks.append(transcript[start:end])
19             start = end - overlap_chars
20     return chunks

```

23.4 Video Metadata Extraction Pipelines

A single hour of 1080p video contains 108,000 frames. The key to practical pipelines is intelligent *frame sampling*.

Table 23.1: Frame sampling strategies

Strategy	Frames/hr	Best For
Fixed interval (1 fps)	3,600	General indexing
Fixed interval (1/10s)	360	Lectures, presentations
Scene change detection	50–500	Dynamic content
Motion-triggered	Varies	Surveillance

Warning

Video processing can generate enormous API costs. A 1-hour video at 1 fps produces 3,600 calls and ~5.76M input tokens. That is not a typo. That is a “we need to talk about your API bill” meeting. Always start with coarse sampling (one frame every 30–60 seconds) and refine only when necessary.

23.5 Video Frame Analysis with Claude Code

Once frames are extracted, each is analyzed individually. The following pipeline produces a structured timeline:

Code Example

```
1 import anthropic
2 import json
3 import subprocess
4 from dataclasses import dataclass
5 from pathlib import Path
6
7
8 @dataclass
9 class FrameAnalysis:
10     frame_path: str
11     timestamp_seconds: float
12     description: str
13     objects: list[str]
14     activity: str
15
16
17 FRAME_PROMPT = """Analyze this video frame and return JSON:
18 {"description": "what is shown", "objects": ["visible", "objects"],
19  "text_visible": "any on-screen text or null",
20  "activity": "what action is occurring"}
21 Return ONLY JSON."""
22
23
24 def extract_frames_interval(
25     video_path: str, output_dir: str,
26     interval_seconds: float = 10.0,
27 ) -> list[str]:
28     """Extract frames at fixed intervals using FFmpeg."""
29     Path(output_dir).mkdir(parents=True, exist_ok=True)
30     cmd = [
31         "ffmpeg", "-i", video_path,
32         "-vf", f"fps=1/{interval_seconds},scale=1024:-1",
33         "-q:v", "2",
34         f"{output_dir}/frame_%05d.jpg",
35         "-y", "-loglevel", "quiet",
36     ]
37     subprocess.run(cmd, check=True)
38     return sorted(str(f) for f in Path(output_dir).glob("frame_*.jpg"))
39
40
```

```
41 class VideoAnalysisPipeline:
42     def __init__(self):
43         self.claude = anthropic.Anthropic()
44
45     def analyze_frame(self, frame_path: str, timestamp: float) ->
FrameAnalysis:
46         image_data, media_type = encode_image(frame_path)
47         message = self.claude.messages.create(
48             model="claude-sonnet-4-20250514", max_tokens=512,
49             messages=[{"role": "user", "content": [
50                 {"type": "image", "source": {
51                     "type": "base64", "media_type": media_type,
52                     "data": image_data,
53                 }},
54                 {"type": "text", "text": FRAME_PROMPT},
55             ]}],
56         )
57         parsed = json.loads(message.content[0].text)
58         return FrameAnalysis(
59             frame_path=frame_path, timestamp_seconds=timestamp,
60             description=parsed["description"],
61             objects=parsed["objects"], activity=parsed["activity"],
62         )
```

23.6 Building Multi-Modal ETL Pipelines

A production multi-modal ETL pipeline uses the *processor registry* pattern, where each file type has a registered handler:

Code Example

```
1 import anthropic
2 import json
3 import logging
4 from abc import ABC, abstractmethod
5 from dataclasses import dataclass, field
6 from datetime import datetime
7 from enum import Enum
8 from pathlib import Path
```

```
9 from typing import Optional
10
11 logger = logging.getLogger(__name__)
12
13
14 class ProcessingStatus(Enum):
15     PENDING = "pending"
16     COMPLETED = "completed"
17     FAILED = "failed"
18     SKIPPED = "skipped"
19
20
21 @dataclass
22 class ProcessingResult:
23     file_path: str
24     status: ProcessingStatus
25     extracted_data: Optional[dict] = None
26     error_message: Optional[str] = None
27     processing_time_ms: float = 0.0
28
29
30 class BaseProcessor(ABC):
31     def __init__(self, client: anthropic.Anthropic):
32         self.client = client
33
34     @abstractmethod
35     def supported_extensions(self) -> set[str]: ...
36
37     @abstractmethod
38     def process(self, file_path: str) -> dict: ...
39
40     def can_handle(self, file_path: str) -> bool:
41         return Path(file_path).suffix.lower() in self.
42             supported_extensions()
43
44 class ImageProcessor(BaseProcessor):
45     def supported_extensions(self) -> set[str]:
46         return {".jpg", ".jpeg", ".png", ".gif", ".webp"}
47
48     def process(self, file_path: str) -> dict:
49         preprocessed, prep_type = preprocess_image(
50             Path(file_path).read_bytes()
51         )
```

```
52     message = self.client.messages.create(
53         model="claude-sonnet-4-20250514", max_tokens=2048,
54         messages=[{"role": "user", "content": [
55             {"type": "image", "source": {
56                 "type": "base64", "media_type": prep_type,
57                 "data": preprocessed,
58             }},
59             {"type": "text", "text": (
60                 "Extract all structured data from this image. "
61                 "Return JSON with: description, text_content, "
62                 "objects, categories, and any tabular data."
63             )},
64         ]}],
65     )
66     return json.loads(message.content[0].text)
67
68
69 class MultiModalETLPipeline:
70     """Orchestrates multi-modal file processing."""
71
72     def __init__(self):
73         self.client = anthropic.Anthropic()
74         self.processors: list[BaseProcessor] = [
75             ImageProcessor(self.client),
76         ]
77
78     def process_file(self, file_path: str) -> ProcessingResult:
79         import time
80         processor = next(
81             (p for p in self.processors if p.can_handle(file_path)),
82             None,
83         )
84         if processor is None:
85             return ProcessingResult(
86                 file_path=file_path, status=ProcessingStatus.SKIPPED,
87                 error_message="No processor for file type",
88             )
89         start = time.monotonic()
90         try:
91             extracted = processor.process(file_path)
92             return ProcessingResult(
93                 file_path=file_path, status=ProcessingStatus.
COMPLETED,
94                 extracted_data=extracted,
```

```
95         processing_time_ms=(time.monotonic() - start) * 1000,  
96     )  
97     except Exception as e:  
98         logger.error(f"Failed to process {file_path}: {e}")  
99         return ProcessingResult(  
100             file_path=file_path, status=ProcessingStatus.FAILED,  
101             error_message=str(e),  
102             processing_time_ms=(time.monotonic() - start) * 1000,  
103         )
```

! Tip

The processor registry pattern makes it easy to add new file types without modifying existing code. Simply implement a new `BaseProcessor` subclass and register it. Open/closed principle: open for extension, closed for “who touched the image processor and broke PDFs.”

23.7 Document Digitization with Claude Code Vision

Traditional OCR systems produce raw text but struggle with layout understanding and contextual interpretation. They see pixels; CLAUDE CODE’s vision understands document *semantics*. It’s the difference between a scanner and a human who can read—except this human processes invoices at 3 AM without complaining.

Invoice Processing Example

Code Example

```
1 import anthropic  
2 import json  
3 from pydantic import BaseModel  
4 from typing import Optional  
5  
6  
7 class LineItem(BaseModel):  
8     description: str  
9     quantity: float
```

```
10     unit_price: float
11     total: float
12
13
14 class InvoiceData(BaseModel):
15     invoice_number: str
16     invoice_date: str
17     vendor_name: str
18     line_items: list[LineItem]
19     subtotal: float
20     tax_amount: Optional[float] = None
21     total_amount: float
22     currency: str = "USD"
23
24
25 INVOICE_PROMPT = """Extract all data from this invoice image.
26 Return a JSON object with: invoice_number, invoice_date (YYYY-MM-DD),
27 vendor_name, line_items (description, quantity, unit_price, total),
28 subtotal, tax_amount, total_amount, currency.
29 All monetary values as numbers. Return ONLY JSON."""
30
31
32 def extract_invoice(
33     client: anthropic.Anthropic, image_path: str,
34 ) -> InvoiceData:
35     image_data, media_type = encode_image(image_path)
36     message = client.messages.create(
37         model="claude-sonnet-4-20250514", max_tokens=4096,
38         messages=[{"role": "user", "content": [
39             {"type": "image", "source": {
40                 "type": "base64", "media_type": media_type,
41                 "data": image_data,
42             }},
43             {"type": "text", "text": INVOICE_PROMPT},
44         ]}],
45     )
46     return InvoiceData(**json.loads(message.content[0].text))
```

 **Warning**

When processing financial documents, always implement a validation layer that cross-checks extracted values—verify line item totals, check that items sum to the subtotal, and that subtotal plus tax equals the total. Flag discrepancies for human review. “The AI said the invoice was \$47 million” is not a sentence your accounting department wants to hear.

Validation

Code Example

```
1 from dataclasses import dataclass
2
3
4 @dataclass
5 class ValidationResult:
6     is_valid: bool
7     errors: list[str]
8     warnings: list[str]
9
10
11 def validate_invoice(invoice: InvoiceData) -> ValidationResult:
12     errors, warnings = [], []
13     for i, item in enumerate(invoice.line_items):
14         expected = round(item.quantity * item.unit_price, 2)
15         if abs(item.total - expected) > 0.01:
16             errors.append(f"Line {i+1}: total != qty * price")
17
18     items_sum = round(sum(it.total for it in invoice.line_items), 2)
19     if abs(invoice.subtotal - items_sum) > 0.01:
20         errors.append(f"Subtotal {invoice.subtotal} != line items sum
21             {items_sum}")
22
23     if not invoice.invoice_number:
24         errors.append("Missing invoice number")
25     if not invoice.line_items:
26         errors.append("No line items extracted")
27
28     return ValidationResult(is_valid=len(errors) == 0, errors=errors,
29         warnings=warnings)
```

23.8 PDF Processing Pipeline

PDFs require special handling: text-based pages can use direct text extraction, while scanned pages must be rendered as images and processed through CLAUDE CODE's vision:

Code Example

```
1 import base64
2
3 class PDFProcessor(BaseProcessor):
4     """Processor for PDF documents."""
5
6     def supported_extensions(self) -> set[str]:
7         return {".pdf"}
8
9     def process(self, file_path: str) -> dict:
10        import fitz # PyMuPDF
11
12        doc = fitz.open(file_path)
13        pages_data = []
14        for page_num in range(len(doc)):
15            page = doc[page_num]
16            text = page.get_text()
17            if len(text.strip()) < 50:
18                # Scanned page: render as image
19                pix = page.get_pixmap(matrix=fitz.Matrix(2, 2))
20                img_b64 = base64.standard_b64encode(
21                    pix.tobytes("jpeg")
22                ).decode("utf-8")
23                message = self.client.messages.create(
24                    model="claude-sonnet-4-20250514", max_tokens
25                    =2048,
26                    messages=[{"role": "user", "content": [
27                        {"type": "image", "source": {
28                            "type": "base64", "media_type": "image/
```

```

29         }},
30         {"type": "text", "text": (
31             "Extract text, tables, and form fields. "
32             "Return JSON with: text, tables,
form_fields."
33         )}},
34     ]}],
35 )
36     page_data = json.loads(message.content[0].text)
37     else:
38         page_data = {"text": text, "tables": [], "form_fields
": {}}
39     pages_data.append({"page_number": page_num + 1, "content"
: page_data})
40     doc.close()
41     return {"total_pages": len(pages_data), "pages": pages_data}

```

23.9 Document Classification System

For large-scale digitization, classify documents by type first, then apply type-specific extraction prompts:

Code Example

```

1  import anthropic
2  import json
3  from dataclasses import dataclass
4  from enum import Enum
5
6
7  class DocumentType(Enum):
8      INVOICE = "invoice"
9      RECEIPT = "receipt"
10     CONTRACT = "contract"
11     FORM = "form"
12     UNKNOWN = "unknown"
13
14
15 @dataclass

```

```
16 class ClassificationResult:
17     document_type: DocumentType
18     confidence: float
19     reasoning: str
20
21
22 CLASSIFICATION_PROMPT = """Classify this document into one of:
23 invoice, receipt, contract, form, unknown.
24 Return JSON: {"document_type": "...", "confidence": 0.0-1.0,
25 "reasoning": "brief explanation"}
26 Return ONLY JSON."""
27
28
29 class DocumentClassifier:
30     def __init__(self, client: anthropic.Anthropic):
31         self.client = client
32
33     def classify(self, image_data: str, media_type: str) ->
34     ClassificationResult:
35         message = self.client.messages.create(
36             model="claude-sonnet-4-20250514", max_tokens=256,
37             messages=[{"role": "user", "content": [
38                 {"type": "image", "source": {
39                     "type": "base64", "media_type": media_type,
40                     "data": image_data,
41                 }},
42                 {"type": "text", "text": CLASSIFICATION_PROMPT},
43             ]}],
44             )
45         parsed = json.loads(message.content[0].text)
46         return ClassificationResult(
47             document_type=DocumentType(parsed["document_type"]),
48             confidence=parsed["confidence"],
49             reasoning=parsed["reasoning"],
50         )
```

! Tip

For high-volume digitization (10,000+ documents/day), use a two-pass strategy: classify with Haiku (faster, cheaper), then extract with Sonnet (more accurate). This reduces costs by 40–60%. Think of Haiku as the bouncer deciding who gets in, and Sonnet as the bartender taking the actual order.

23.10 Storage and Indexing Strategies

Separate raw assets from extracted metadata:

- **Object storage:** Raw media files, organized by date and type, with lifecycle policies for cold storage.
- **Relational database:** Extracted metadata, classification results, and processing state.
- **Vector database:** Embeddings of extracted text for semantic search.
- **Search index:** Full-text search across extracted content.

Code Example

```

1  -- Core asset table
2  CREATE TABLE media_assets (
3      asset_id        UUID PRIMARY KEY DEFAULT gen_random_uuid(),
4      s3_key          TEXT NOT NULL UNIQUE,
5      asset_type      TEXT NOT NULL,
6      file_size_bytes BIGINT NOT NULL,
7      content_type    TEXT NOT NULL,
8      uploaded_at     TIMESTAMP WITH TIME ZONE DEFAULT NOW(),
9      status          TEXT DEFAULT 'pending'
10 );
11
12 -- Extracted text with full-text search
13 CREATE TABLE extracted_text (
14     id            SERIAL PRIMARY KEY,
15     asset_id     UUID REFERENCES media_assets(asset_id),
16     content      TEXT NOT NULL,
17     content_tsv  TSVECTOR GENERATED ALWAYS AS (
18                 to_tsvector('english', content)
19             ) STORED
20 );
21 CREATE INDEX idx_extracted_text_tsv
22     ON extracted_text USING GIN (content_tsv);
23
24 -- Extracted structured fields (JSONB for flexibility)
25 CREATE TABLE extracted_fields (
26     id            SERIAL PRIMARY KEY,
27     asset_id     UUID REFERENCES media_assets(asset_id),

```

```
28     doc_type    TEXT NOT NULL,
29     fields      JSONB NOT NULL,
30     confidence  REAL DEFAULT 1.0
31 );
32 CREATE INDEX idx_extracted_fields_jsonb
33     ON extracted_fields USING GIN (fields jsonb_path_ops);
```

Note

For semantic search, generate embeddings using a dedicated embedding model (Voyage AI, OpenAI's `text-embedding-3-large`), not CLAUDE CODE itself. Using CLAUDE CODE for embeddings is like hiring a brain surgeon to take your temperature. Embedding models are orders of magnitude cheaper and faster for indexing.

23.11 Advanced Techniques

Multi-Image Analysis

Some tasks require analyzing multiple related images together. CLAUDE CODE supports multiple images in a single API call:

Code Example

```
1 def analyze_multi_page_document(
2     client: anthropic.Anthropic,
3     image_paths: list[str],
4     extraction_prompt: str,
5 ) -> dict:
6     """Send multiple page images for unified analysis."""
7     content = []
8     for i, path in enumerate(image_paths):
9         data, media = encode_image(path)
10        content.append({"type": "image", "source": {
11            "type": "base64", "media_type": media, "data": data,
12        }})
13        content.append({"type": "text", "text": f"(Page {i+1} of {len
14            (image_paths)})"})
```

```
14 content.append({"type": "text", "text": extraction_prompt})
15
16 message = client.messages.create(
17     model="claude-sonnet-4-20250514", max_tokens=8192,
18     messages=[{"role": "user", "content": content}],
19 )
20 return json.loads(message.content[0].text)
```

Error Handling and Retry Strategies

Code Example

```
1 import time
2 import random
3 from functools import wraps
4
5 def retry_with_backoff(
6     max_retries: int = 3, base_delay: float = 1.0,
7     max_delay: float = 60.0,
8 ):
9     """Decorator for retrying API calls with exponential backoff."""
10    def decorator(func):
11        @wraps(func)
12        def wrapper(*args, **kwargs):
13            for attempt in range(max_retries + 1):
14                try:
15                    return func(*args, **kwargs)
16                except (anthropic.RateLimitError,
17                        anthropic.InternalServerError) as e:
18                    if attempt == max_retries:
19                        raise
20                    delay = min(
21                        base_delay * (2 ** attempt) + random.uniform
22                        (0, 1),
23                        max_delay,
24                    )
25                    time.sleep(delay)
26                except anthropic.BadRequestError:
27                    raise # Non-retryable
28            return wrapper
```

```
28 return decorator
```

23.12 Exercises

1. **Receipt Processing Pipeline.** Build a pipeline that processes receipt photographs, extracts store name, date, items with prices, tax, and total. Write results to SQLite. Validate that item prices plus tax equal total.
2. **Audio Meeting Summarizer.** Build a pipeline using Whisper for transcription and CLAUDE CODE for extraction of summaries, action items, key decisions, and topics. Handle files over 2 hours with chunking.
3. **Multi-Modal ETL with Error Handling.** Extend `MultiModalETLPipeline` to support images and PDFs with: exponential backoff retry, a dead-letter queue for failures, Prometheus metrics, and idempotent processing.
4. **Cost Calculator.** Build a tool that estimates the cost of processing a directory of mixed media through the pipeline, based on file types, image dimensions, and audio durations.

Tip

For exercises involving image or document processing, find free sample datasets on Kaggle—search for “receipt OCR dataset” or “invoice dataset.” For audio, use freely available podcast episodes. Do not use your company’s actual financial documents for testing. Your compliance team will not find it as educational as you do. Always respect licensing terms.

23.13 Monitoring and Observability

Production multi-modal pipelines require comprehensive monitoring. Track these key metrics:

Data Lifecycle Management

Multi-modal data grows quickly. Implement lifecycle policies from day one:

Table 23.2: Key metrics for multi-modal pipeline monitoring

Metric	Type	Alert Threshold	Purpose
Processing latency (p99)	Histogram	>60s per item	Performance regression
API error rate	Counter	>5%	Claude Code API issues
Validation failure rate	Gauge	>20%	Extraction quality
Token usage per item	Histogram	>2× baseline	Cost anomaly
Queue depth	Gauge	>10,000 items	Throughput bottleneck

- **Object storage:** Transition processed files to infrequent access after 30 days, Glacier after 90 days, deep archive after 365 days.
- **Temporary frames:** Delete extracted video frames after 7 days.
- **Database:** Archive old extractions annually. Partition tables by date for efficient pruning.

Prompt Caching for Batch Processing

When processing many documents of the same type, the system prompt and extraction instructions remain constant. Use prompt caching to avoid reprocessing these tokens:

Code Example

```

1 def process_batch_with_caching(
2     client: anthropic.Anthropic,
3     image_data_list: list[tuple[str, str]],
4     system_prompt: str,
5     extraction_prompt: str,
6 ) -> list[dict]:
7     results = []
8     for image_data, media_type in image_data_list:
9         message = client.messages.create(
10             model="claude-sonnet-4-20250514", max_tokens=2048,
11             system=[{
12                 "type": "text", "text": system_prompt,
13                 "cache_control": {"type": "ephemeral"}},
14             ],
15             messages=[{"role": "user", "content": [
16                 {"type": "image", "source": {

```

```
17         "type": "base64", "media_type": media_type,
18         "data": image_data,
19     }},
20     {"type": "text", "text": extraction_prompt},
21 ]}],
22 )
23 results.append(json.loads(message.content[0].text))
24 return results
```


Part **VIII**

The Future

24

Building a Claude Code-First Data Platform

“Every company will eventually become a CLAUDE CODE wrapper. It’s wrappers all the way down.”

Tip

Building a CLAUDE CODE-first data platform is not about replacing your infrastructure overnight. It is about systematically identifying where AI augmentation delivers the highest ROI and incrementally transforming your architecture, team, and processes. Rome wasn’t built in a day, and neither is a data platform—but with CLAUDE CODE, you can at least get the architecture diagram done before lunch.

This chapter provides a blueprint for building and operating a Claude Code-augmented data platform, including honest assessments of where it should *not* be used, how to measure ROI, and how to manage the human side of this transformation. Whether you are a solo data engineer at a startup or leading a platform team at a Fortune 500, the patterns here scale to your context.

24.1 Architecture Patterns for AI-Augmented Data Platforms

A modern Claude Code-augmented platform consists of five layers, each with defined AI integration points:

Ingestion Layer Connectors, CDC pipelines, and API extractors. Claude Code assists with connector configuration, schema mapping, and error handling logic.

Processing Layer Spark, Flink, dbt, and custom Python. Claude Code generates transformation code, optimizes queries, and reviews data quality rules.

Storage Layer Data lakes, warehouses, and lakehouses. Claude Code advises on partitioning strategies, manages schema evolution, and generates documentation.

Serving Layer APIs, dashboards, and feature stores. Claude Code powers natural language query interfaces and generates API boilerplate.

Orchestration Layer Airflow, Dagster, and Prefect. Claude Code monitors pipeline health, diagnoses failures, and automates incident response.

The key architectural decision is *where* Claude Code integrates. Three patterns have emerged as the most effective.

The Sidecar Pattern

The *sidecar pattern* attaches a Claude Code agent to each critical data service. The sidecar monitors the service's events, provides advisory analysis, and optionally takes remediation actions within a constrained scope. This pattern is effective because each sidecar has deep context about its specific service without needing global visibility.

Code Example

Claude Code sidecar pattern

```
1 import anthropic
2 from collections import deque
3
4 class ClaudeSidecar:
```

```
5     """A sidecar agent that monitors a data service."""
6
7     def __init__(self, service_name: str,
8                 model: str = "claude-sonnet-4-20250514",
9                 window_size: int = 100):
10         self.service_name = service_name
11         self.client = anthropic.Anthropic()
12         self.model = model
13         self.context_window = deque(maxlen=window_size)
14         self.alert_history: list[dict] = []
15
16     async def observe(self, event: dict) -> None:
17         """Record an event from the monitored service."""
18         self.context_window.append(event)
19
20     async def analyze(self, question: str) -> str:
21         """Ask the sidecar to analyze recent events."""
22         context = "\n".join(
23             f"[{e['timestamp']}] {e['type']}: "
24             f"{e['message']}"
25             for e in list(self.context_window)[-50:]
26         )
27         response = self.client.messages.create(
28             model=self.model, max_tokens=2048,
29             system=(
30                 f"You are a sidecar agent for "
31                 f"{self.service_name}. Analyze events, "
32                 f"detect anomalies, recommend actions. "
33                 f"Be specific and cite event timestamps."),
34             messages=[{
35                 "role": "user",
36                 "content": (
37                     f"Recent events:\n{context}\n\n"
38                     f"Question: {question}")
39             }],
40         )
41         return response.content[0].text
42
43     async def should_alert(self) -> dict | None:
44         """Proactively check if conditions warrant alert."""
45         if len(self.context_window) < 10:
46             return None
47         recent_errors = sum(
48             1 for e in list(self.context_window)[-20:]
```

```

49         if e.get("level") == "ERROR")
50     if recent_errors < 3:
51         return None
52     analysis = await self.analyze(
53         "Are recent errors concerning? "
54         "Should the on-call engineer be alerted?")
55     return {"service": self.service_name,
56           "error_count": recent_errors,
57           "analysis": analysis}

```

The Gateway Pattern

The *gateway pattern* routes all Claude Code interactions through a centralized API gateway. This gateway handles authentication, rate limiting, prompt management, cost tracking, and audit logging. It provides a single control plane for all AI usage across the platform.

Note

The gateway pattern is especially valuable in regulated industries where every AI interaction must be logged, every prompt must be versioned, and usage must be attributed to specific users and cost centers. It also enables A/B testing of different prompts and models without changing application code.

Code Example

Centralized AI gateway

```

1  from dataclasses import dataclass, field
2  from datetime import datetime, timezone
3
4  @dataclass
5  class GatewayRequest:
6      user_id: str
7      team: str
8      use_case: str # "code_review", "incident", etc.
9      prompt: str
10     model_override: str | None = None
11
12 @dataclass

```

```
13 class GatewayConfig:
14     default_model: str = "claude-sonnet-4-20250514"
15     model_overrides: dict[str, str] = field(
16         default_factory=lambda: {
17             "classification": "claude-haiku-4-20250514",
18             "architecture": "claude-opus-4-20250514",
19         })
20     rate_limits: dict[str, int] = field(
21         default_factory=lambda: {
22             "default": 100,      # per hour
23             "batch": 1000,
24             "interactive": 50,
25         })
26
27 class AIGateway:
28     """Centralized gateway for all Claude Code calls."""
29
30     def __init__(self, config: GatewayConfig):
31         self.config = config
32         self.client = anthropic.Anthropic()
33         self.usage_log: list[dict] = []
34
35     def route(self, request: GatewayRequest) -> str:
36         # Select model based on use case
37         model = (request.model_override
38                 or self.config.model_overrides.get(
39                     request.use_case,
40                     self.config.default_model))
41
42         # Check rate limits
43         self._check_rate_limit(
44             request.user_id, request.use_case)
45
46         # Execute request
47         response = self.client.messages.create(
48             model=model, max_tokens=4096,
49             messages=[{"role": "user",
50                       "content": request.prompt}])
51
52         # Log usage for cost tracking
53         self.usage_log.append({
54             "timestamp": datetime.now(
55                 timezone.utc).isoformat(),
56             "user_id": request.user_id,
```

```
57         "team": request.team,
58         "use_case": request.use_case,
59         "model": model,
60         "input_tokens": response.usage.input_tokens,
61         "output_tokens": response.usage.output_tokens,
62     })
63     return response.content[0].text
64
65 def _check_rate_limit(self, user_id, use_case):
66     # Implementation: count recent requests,
67     # raise if over limit
68     pass
69
70 def cost_report(self, team=None) -> dict:
71     """Generate cost report, optionally by team."""
72     filtered = self.usage_log
73     if team:
74         filtered = [
75             r for r in filtered
76             if r["team"] == team]
77     # Pricing per 1M tokens (approximate)
78     pricing = {
79         "claude-haiku-4-20250514": (0.25, 1.25),
80         "claude-sonnet-4-20250514": (3.0, 15.0),
81         "claude-opus-4-20250514": (15.0, 75.0),
82     }
83     total = 0.0
84     by_use_case: dict[str, float] = {}
85     for r in filtered:
86         inp, out = pricing.get(
87             r["model"], (3.0, 15.0))
88         cost = (r["input_tokens"] * inp
89             + r["output_tokens"] * out) / 1_000_000
90         total += cost
91         by_use_case[r["use_case"]] = \
92             by_use_case.get(r["use_case"], 0) + cost
93     return {"total_cost": total,
94           "by_use_case": by_use_case}
```

The Mesh Pattern

In larger organizations, a *mesh pattern* combines sidecars and gateways: each domain team runs its own sidecars and prompt libraries, but all traffic routes through a shared gateway for governance, cost tracking, and cross-domain coordination. This mirrors the data mesh philosophy of domain-oriented ownership with federated governance.

24.2 Team Structure and Cost Management

The Evolving Role Map

Table 24.1: Role evolution in a Claude Code-augmented platform

Traditional Role	Evolved Role	Key Changes
Data Engineer	AI-Augmented Data Engineer	Writes prompts alongside code; reviews AI output; focuses on architecture and system design
Analytics Engineer	Semantic Layer Architect	Designs natural language interfaces; curates prompt templates; owns business context encoding
Data Platform Engineer	AI Platform Engineer	Manages AI gateway infrastructure; optimizes token budgets; maintains MCP servers
(New)	Prompt Engineer / AI Ops	Develops and versions prompt libraries; monitors AI output quality; tunes retrieval systems

Use a *hub and spoke* model for team organization: a small AI Platform Hub (3–5 people) maintains the gateway, prompt libraries, and shared tooling. Domain spokes (each data team) each designate an “AI champion” who adapts shared tooling to their domain. A governance overlay sets policies for data classification, human review requirements, and cost allocation.

 **Warning**

Do not create a centralized “AI team” that becomes a bottleneck for every AI-related request. We tried centralized everything in the 2010s; it was called “the data team that has a 6-month backlog.” This pattern creates queues, slows iteration, and prevents domain teams from developing AI expertise. Instead, distribute AI capabilities across existing teams with a small platform team providing infrastructure, standards, and enablement—not gatekeeping.

Cost Optimization Strategies

AI API costs can grow rapidly without deliberate management. Five strategies keep costs under control:

1. **Model tiering:** Use Haiku for high-volume, low-complexity tasks (log classification, alert triage); Sonnet for code generation and analysis; Opus for architecture decisions and complex reasoning. A single pipeline monitoring system might use all three tiers depending on the task.
2. **Prompt caching:** Cache system prompts and schema context using Claude’s prompt caching feature for up to 90% input cost reduction on repeated calls with the same system prompt.
3. **Batch API:** Use the Batch API for non-time-sensitive workloads (documentation generation, bulk code review, nightly quality checks) at 50% cost savings.
4. **Output management:** Constrain `max_tokens` to what each task actually needs. A classification task needs 50 tokens, not 4096.
5. **Smart context management:** Use embeddings or search to select only relevant files and schema definitions, rather than dumping entire codebases or full warehouse schemas into the context window.

 **Tip**

Track cost-per-pipeline-run and cost-per-incident-resolved metrics from day one. These “unit economics” metrics demonstrate that CLAUDE CODE costs are offset by reduced engineering hours and faster resolution times. A \$0.50 AI-assisted incident resolution that saves 45 minutes of engineer time is a clear win at any engineering salary. Your CFO speaks fluent ROI—learn the language.


```
41         })
42         return result
43     return wrapper
44     return decorator
```

24.3 When NOT to Use Claude Code

Knowing when *not* to use AI is as important as knowing when to use it. Over-applying Claude Code adds latency, cost, and non-determinism without benefit.

Hard constraints—do not use for:

- **Deterministic computations:** Financial calculations, checksums, cryptographic operations. These require exact, reproducible results that LLMs cannot guarantee.
- **Sub-millisecond real-time decisions:** Fraud detection scoring, real-time bidding. LLM inference latency (100ms–5s) is orders of magnitude too slow.
- **Bulk data transformation:** Do not send millions of rows through an LLM. Use Claude Code to *generate* the transformation code, then execute that code directly.
- **Secrets and credentials management:** Never pass API keys, passwords, or tokens through an LLM. Use dedicated secrets managers.
- **Regulatory-prohibited contexts:** Some jurisdictions and industries prohibit AI decision-making in specific contexts. Know your regulatory landscape.

Soft constraints—exercise caution:

- **Highly novel domains:** If your domain has limited representation in training data, AI outputs will be less reliable. Validate more aggressively.
- **Safety-critical systems:** Use advisory mode with mandatory human review. Never let an agent autonomously modify production systems that could cause financial harm or data loss without approval gates.
- **Large-scale production schema migrations:** Plan with AI, but execute with tested, deterministic migration tooling (Flyway, Alembic, dbt). AI should generate and review the migration; automated tooling should apply it.

 **Warning**

The most common mistake is applying CLAUDE CODE to every problem indiscriminately. Not every nail needs an AI-powered hammer. If a regex, a SQL CASE statement, or a simple Python script solves the problem reliably, use that. Using CLAUDE CODE to add two numbers together is like hiring a NASA engineer to change a lightbulb. Adding AI to simple deterministic tasks introduces latency, cost, and non-determinism without benefit. Reserve AI for tasks that genuinely require reasoning, language understanding, or pattern recognition across unstructured data.

24.4 Platform Maturity Model

Organizations adopt AI at different paces. This maturity model helps you assess where you are and plan where to go.

Level 1 — Ad Hoc Individual engineers use Claude Code via the web interface or CLI. No shared infrastructure, no cost tracking, no governance. Value is real but uncoordinated.

Level 2 — Standardized Shared API access via organization account. Basic prompt libraries for common tasks. Usage guidelines and cost awareness. This is the minimum viable level for a team.

Level 3 — Integrated Claude Code is embedded in CI/CD pipelines, quality checks, and incident response workflows. Centralized gateway provides cost tracking and audit logging. Prompt libraries are versioned and reviewed.

Level 4 — Optimized Model tiering routes tasks to the most cost-effective model. Prompt caching and batch processing reduce costs. A/B testing of prompts and models drives continuous improvement. RAG systems provide organizational knowledge grounding.

Level 5 — Autonomous Self-healing pipelines diagnose and fix common failures without human intervention. Schema evolution is AI-mediated for safe, backward-compatible changes. Agents coordinate across services to resolve complex incidents.

Note

Aim for Level 3–4 within the first year. Level 5 autonomy should be pursued selectively for well-understood, low-risk workflows only after you have robust monitoring, audit logging, and rollback mechanisms in place. Premature autonomy erodes trust when things go wrong—and they will go wrong, because entropy is undefeated.

24.5 Build vs. Buy Decisions

Every AI integration presents a build-versus-buy decision. The right answer depends on your team’s size, domain complexity, and competitive landscape.

Build custom integrations when:

- Your domain has unique semantics that generic tools cannot capture (proprietary data models, custom pipeline frameworks).
- You need deep integration with proprietary internal systems.
- Prompt engineering and AI-powered workflows provide competitive advantage.
- Data sensitivity or regulatory requirements preclude third-party SaaS.

Buy (or use open-source) when:

- The vendor covers 80%+ of your requirements out of the box.
- The problem is commoditized (text-to-SQL, code review, documentation generation).
- Your team is small and cannot maintain custom AI infrastructure.
- Time-to-value is critical and the vendor offers fast onboarding.
- The vendor provides enterprise features (SSO, audit logging, compliance certifications) that you would otherwise need to build.

Tip

A useful heuristic: if the AI integration is *about* your data (schema-aware querying, domain-specific quality checks), build it. If the AI integration is *generic* (code formatting, documentation spell-checking, standard CI checks),

buy it. Your engineering time is better spent on integrations where domain knowledge creates defensible value.

24.6 Migration and Change Management

The Strangler Fig Pattern

Named after the strangler fig tree that gradually envelops its host, this pattern replaces manual processes with AI-augmented ones incrementally, never cutting over all at once.

1. **Shadow Mode (Weeks 1–4)**: Claude Code processes inputs in parallel with existing workflows. Outputs are logged and compared but never used in production. This builds confidence and reveals edge cases.
2. **Advisory Mode (Weeks 5–12)**: Claude Code recommendations are surfaced via Slack notifications or PR comments. Engineers choose whether to follow them. Track acceptance rates—they are your leading indicator of trust and quality.
3. **Supervised Mode (Months 3–6)**: Claude Code acts on pre-approved task categories with full logging and a human override window. If no human objects within the window (e.g., 15 minutes), the action proceeds.
4. **Autonomous Mode (Month 6+)**: Proven, high-confidence workflows run autonomously with monitoring, alerting, and automatic rollback on anomalies.

Every integration point needs four safety mechanisms: **feature flags** (disable any AI integration instantly), **fallback paths** (revert to pre-AI workflows), **provenance tagging** (mark AI-generated artifacts), and **circuit breakers** (auto-disable on excessive error rates).

Warning

Resist the temptation to skip Shadow Mode. “It worked in dev” is the data engineering equivalent of “hold my beer.” Even if you are confident in the AI’s capabilities, Shadow Mode reveals integration issues, edge cases, and trust concerns that only surface when the system runs against real production data and workflows. Two weeks of shadow running has prevented many costly production incidents.

24.7 Organizational Change Management

Technology transformation is 20% technology and 80% people. The most sophisticated AI platform will fail if the team does not trust it, does not understand it, or feels threatened by it. You can have the best CLAUDE CODE integration in the world, but if your team treats it like a haunted Ouija board, adoption will be zero. Address common concerns honestly:

“Will Claude Code replace my job?” Frame AI as a force multiplier, not a replacement. Engineers who leverage AI are *more* valuable because they produce higher quality work faster. The role evolves from writing boilerplate to designing systems, reviewing AI output, and solving novel problems that AI cannot handle alone.

“I don’t trust AI-generated code.” This is a healthy instinct. Channel it into building robust validation processes: automated tests, code review checklists for AI-generated code, and metrics that track AI output quality over time.

“This adds complexity.” Acknowledge it honestly. AI integration *does* add complexity—another dependency, another failure mode, another cost center. Show concrete metrics proving that the value (faster development, fewer incidents, better documentation) offsets the complexity cost.

Build champions by giving early adopters dedicated experimentation time (10–20% of sprint capacity), creating shared channels for wins and patterns, and having them present at team meetings. Champions who speak from experience are far more persuasive than mandates from leadership.

i Note

Track and share “AI wins” publicly. When an agent resolves a 3 AM incident without waking anyone, when AI-generated documentation passes review without changes, when a text-to-SQL query saves an analyst two hours—celebrate these moments. Concrete examples build organic adoption faster than any top-down mandate.

24.8 ROI Measurement and Governance

Table 24.2: ROI metrics for Claude Code-augmented platforms

Metric	Measurement	Target
Developer velocity	Tasks completed per sprint	30–50% increase in 6 months
Time to first pipeline	Days from request to production	50–70% reduction
Incident MTTR	Minutes from alert to resolution	40–60% reduction
Code review turnaround	Hours from PR to approval	60–80% reduction
Documentation coverage	% of pipelines with up-to-date docs	80%+ (from typical 20–30%)
On-call burden	After-hours pages requiring human action	50–70% reduction

Governance Framework

Establish policies for:

- **Data classification and AI eligibility:** Which data can be sent to AI APIs? PII, financial data, and health data may require on-premises models or special handling.
- **Human review requirements:** Which AI-generated outputs require human review before production use? All code changes? Only schema modifications?
- **Audit trail requirements:** Every AI interaction should log the prompt, response, model used, timestamp, and requesting user.
- **Incident response:** What happens when an AI-generated artifact causes a production issue? Who is responsible? How is it triaged differently from human-caused issues?
- **Prompt version control:** Treat prompts like code. Version them, review changes, and test before deployment.

Compliance considerations vary by industry: GDPR and CCPA require data processing agreements with AI providers; SOX and SOC 2 demand audit controls for AI in financial pipelines; HIPAA requires Business Associate Agreements with any provider processing health data.

 **Warning**

Be honest in ROI calculations. We are data engineers—we of all people should not be making up numbers. Use conservative numbers validated with actual time-tracking data, not optimistic estimates. Overpromised ROI undermines trust in the entire initiative and makes it harder to secure budget for future improvements. A 25% improvement that you can prove is more valuable than a claimed 60% improvement that nobody believes.

12-Month Roadmap

Months 1–2: Foundation Set up API access, deploy the AI gateway, establish usage policies, and run the first pilot project (recommend starting with incident response or documentation generation).

Months 3–4: Expand Add 2–3 additional workflows (code review, data quality checks). Build prompt libraries. Implement cost tracking dashboards.

Months 5–6: Integrate Embed Claude Code in CI/CD pipelines. Deploy RAG for runbook automation. Set up quality monitoring. Reach Maturity Level 3.

Months 7–9: Optimize Implement model tiering, prompt caching, and batch processing. Run A/B tests on prompts and models. Publish the first ROI report with real data.

Months 10–12: Scale Roll out to all data teams. Enable supervised autonomy for proven workflows. Plan Level 4–5 initiatives for the following year.

 **Tip**

This roadmap assumes a team of 8–15 engineers. Larger organizations may run phases in parallel across different teams. Smaller teams (3–5 engineers) should focus on fewer integration points per phase—start with the single highest-impact workflow and expand from there.

24.9 Exercises

1. **Architecture Design:** Design a Claude Code-augmented data platform architecture for a mid-size e-commerce company. Draw the five layers, identify

three specific AI integration points per layer, and specify which pattern (side-car, gateway, or mesh) you would use and why.

2. **Cost Model:** Build a spreadsheet model for AI API costs based on your team's current workflow. Estimate daily query volume for five use cases (incident response, code review, documentation, data quality, ad hoc queries). Calculate monthly costs with and without model tiering, prompt caching, and batch processing.
3. **Maturity Assessment:** Assess your current organization against the five-level maturity model. For each gap between your current level and your target level, identify the specific actions, tools, and timeline needed to close it.
4. **Build vs. Buy Analysis:** Pick three AI-powered features you want to add to your platform (e.g., text-to-SQL, automated documentation, incident triage). For each, evaluate build vs. buy using the criteria in this chapter. Present your recommendation with estimated cost, timeline, and risk for each option.
5. **Change Management Plan:** Draft a change management plan for introducing Claude Code to a skeptical team of 10 data engineers. Include your Shadow Mode design, success metrics for each phase, a communication plan, and how you would handle the three common objections discussed in this chapter.

25

The Future of Data Engineering with AI

“The future of data engineering is bright. Mostly because CLAUDE CODE doesn’t need office lighting.”

! Tip

The future of data engineering is not about replacement—it is about amplification. The data engineers who thrive will be those who wield AI as a force multiplier, treating CLAUDE CODE and similar tools as extensions of their expertise rather than threats to their livelihood. The robots are not coming for your job. They are coming for the boring parts of your job. The interesting parts are safe—for now.

Throughout this book, we have explored how CLAUDE CODE transforms every stage of the data engineering lifecycle. In this final chapter, we look toward the horizon: how will the profession evolve, what new patterns will emerge, and what should you do today to prepare for the next decade of data engineering?

25.1 The Evolving Role of the Data Engineer

From Builder to Orchestrator

The data engineer’s primary activity is shifting from *writing* code to *directing*, *reviewing*, and *orchestrating* code generation. This does not mean data engineers

write less code—it means they write *different* code. The boilerplate shrinks while the architecture, validation, and integration logic grows.

Note

This shift does not mean less technical depth. Paradoxically, you need *more* expertise to review AI-generated code, identify subtle bugs that pass tests, and architect systems that AI tools can enhance effectively. A data engineer who does not understand query optimization cannot verify that an AI-generated query is efficient. You wouldn't trust a GPS that says "drive into the lake" just because it sounds confident. The AI raises the floor but does not lower the ceiling.

Dimension	Pre-Cloud (2005–15)	Cloud-Native (2015–24)	AI-Augmented (2024+)
Primary Activity	Managing infrastructure, writing MapReduce jobs	Building pipelines in SQL and Python	Orchestrating AI-generated pipelines, designing systems
Key Skill	Java, Hadoop administration	SQL, Python, cloud services, dbt	Prompt engineering, system design, AI output validation
Code Origin	100% human-written	95% human, 5% generated	40–70% AI-generated, human-reviewed
Iteration Speed	Days to weeks	Hours to days	Minutes to hours
Failure Mode	Hardware failures, cluster issues	Configuration drift, cloud misconfig	AI hallucination, prompt regression, context window limits

New Responsibilities

The AI-augmented data engineer takes on responsibilities that did not exist five years ago:

- **AI Output Quality Assurance:** Reviewing and validating AI-generated

code, SQL, and configurations before they reach production. This requires deeper domain expertise, not less, because you must catch errors that automated tests might miss.

- **Prompt Library Curation:** Maintaining versioned, tested prompt templates for common tasks—pipeline generation, schema review, incident investigation. These prompts become organizational assets, iterated and improved over time.
- **Cost Governance:** Managing AI API costs alongside traditional compute and storage costs. This includes model tiering decisions, prompt caching strategies, and usage budgets per team and use case.
- **Trust Calibration:** Developing an intuition for when to trust AI output at face value versus when to demand independent verification. This calibration is context-dependent and comes from experience.
- **AI-Human Interface Design:** Designing workflows where AI and human responsibilities are clearly delineated, handoffs are smooth, and neither side becomes a bottleneck.

Warning

The most dangerous position is uncritical acceptance of AI output. “CLAUDE CODE said it, I believe it, that settles it” is not an engineering methodology. Engineers who blindly copy-paste AI-generated code without review will introduce bugs, security vulnerabilities, and performance problems. The AI is a powerful first draft generator, not an infallible oracle. Always review, always test, always think critically about the output.

25.2 Emerging Patterns: AI-Native Data Platforms

While most organizations today are *retrofitting* AI onto existing platforms, a new generation of *AI-native data platforms* is being designed from the ground up with AI as a core component, not an add-on.

1. **Intent-Driven Configuration:** Engineers express what they want in natural language or high-level declarative specs. The platform translates intent into executable infrastructure, pipeline code, and monitoring rules. Instead of writing 500 lines of Airflow DAG code, the engineer writes a 20-line manifest describing the desired data flow.

2. **Self-Healing Pipelines:** When a pipeline fails, an agent automatically diagnoses the issue, generates a fix, tests it in a sandbox, and applies it—escalating to a human only when confidence is low or the fix involves a high-risk operation.
3. **Continuous Schema Evolution:** Rather than treating schema changes as disruptive migrations, AI mediates schema evolution continuously. It detects upstream changes, generates backward-compatible adaptations, notifies affected consumers, and applies changes incrementally.
4. **Semantic Metadata:** Every table, column, and pipeline has rich, AI-generated and human-reviewed metadata that enables natural language querying. Users ask “What percentage of our European customers made a purchase last month?” and the platform translates this to the correct SQL against the correct tables.

Code Example

AI-native platform manifest

```
1 platform:
2   name: analytics-platform-v3
3   intent: |
4     Real-time analytics platform ingesting clickstream
5     data, enriching with user profiles from PostgreSQL,
6     serving aggregated metrics to the dashboard team.
7     Maintain 30-second freshness for clickstream and
8     5-minute freshness for user profiles.
9   data_sources:
10    - name: clickstream
11      type: kafka
12      topic: user.clicks
13      freshness_sla: 30s
14      volume: 50k_events_per_second
15    - name: user_profiles
16      source_type: postgresql
17      connection: profiles_db
18      change_capture: cdc
19      freshness_sla: 5m
20   transformations:
21    - name: enriched_clicks
22      intent: |
23        Join clickstream events with user profiles.
24        Add geographic region based on IP address.
25        Filter bot traffic using the ua_parser library.
```

```
26     output_schema: auto_evolve_with_review
27 quality_expectations:
28   completeness: 99.5%
29   freshness_alert: 2x_sla
30   schema_drift_policy: auto_evolve_with_review
31   anomaly_detection: enabled
32 serving:
33   - name: dashboard_metrics
34     type: materialized_view
35     refresh: continuous
36     consumers: [dashboard_team, exec_reports]
```

i Note

Intent-driven configuration does not eliminate the need to understand the underlying systems. When the platform mis-translates intent (and it will), you need the expertise to diagnose and correct the generated implementation. Think of it as a dramatically faster starting point, not a finished product.

The Semantic Layer Revolution

The semantic layer—a business-friendly abstraction over raw data—is being transformed by AI. Traditional semantic layers require manual definition of metrics, dimensions, and relationships. AI-native semantic layers can:

- **Auto-discover relationships** between tables by analyzing column names, data types, foreign keys, and actual data distributions.
- **Generate metric definitions** from business documentation, Slack conversations, and existing dashboard queries.
- **Resolve ambiguity** in natural language queries by asking clarifying questions or presenting multiple interpretations.
- **Learn from corrections:** When a user says “That is not what I meant by revenue—I mean net revenue after refunds,” the system updates its understanding for future queries.

25.3 The Convergence of Data Engineering and ML Engineering

One of the most significant trends is the dissolving boundary between data engineering and ML engineering. As AI becomes central to data platforms, the skill sets converge:

- **Data engineers embed ML models inside pipelines:** Anomaly detection on streaming data, entity resolution during ingestion, data quality scoring at each transformation step.
- **ML engineers need deep data engineering skills:** Feature stores, training data pipelines, model serving infrastructure, and data versioning are all data engineering problems.
- **Both roles require prompt engineering:** Working with LLMs is a core skill for both disciplines, whether for code generation, data analysis, or building AI-powered features.

The result is a new hybrid role—the *AI data engineer*—who combines pipeline expertise with ML literacy and prompt engineering fluency. This does not mean everyone needs to be an expert in all three areas, but baseline competence across all three is becoming table stakes.

Tip

If you are a data engineer, invest in learning ML fundamentals: how models are trained, what feature engineering means, how evaluation metrics work, and how models are deployed and monitored. You do not need to derive gradient descent from first principles, but you need to have an informed conversation with ML engineers and understand what their systems need from yours.

25.4 Natural Language as the New Interface for Data

The Text-to-SQL Revolution

Text-to-SQL has been a research topic for decades, but CLAUDE CODE's capabilities make it practical for production use for the first time. Modern text-to-SQL systems provide full schema awareness, business context understanding, dialect-specific SQL generation, query safety guardrails, and multi-turn refinement.

Code Example

Text-to-SQL service with guardrails

```
1 from anthropic import Anthropic
2 import sqlglot
3
4 class TextToSQLService:
5     """Convert natural language to safe SQL queries."""
6
7     def __init__(self, schema_context: str,
8                 dialect: str = 'snowflake'):
9         self.client = Anthropic()
10        self.schema_context = schema_context
11        self.dialect = dialect
12        self.blocked_ops = [
13            'DROP', 'DELETE', 'TRUNCATE',
14            'INSERT', 'UPDATE', 'ALTER', 'CREATE']
15
16        def generate_sql(self, question: str) -> str | None:
17            """Generate SQL from a natural language question."""
18            response = self.client.messages.create(
19                model="claude-sonnet-4-20250514",
20                max_tokens=2048,
21                system=(
22                    f'You are a SQL expert for {self.dialect}. '
23                    f'Schema:\n{self.schema_context}\n'
24                    f'Rules:\n'
25                    f'- Only generate SELECT statements\n'
26                    f'- Use CTEs for complex logic\n'
27                    f'- Include LIMIT unless aggregating\n'
28                    f'- Add column aliases for clarity\n'
29                    f'- Return only the SQL, no explanation'),
30                messages=[{
31                    "role": "user",
32                    "content": question
33                }],
34            )
35            sql = self._extract_sql(response.content[0].text)
36            if self._validate_sql(sql):
37                return sql
38            return None
39
40        def _validate_sql(self, sql: str) -> bool:
```

```

41     """Validate SQL is safe to execute."""
42     try:
43         parsed = sqlglot.parse_one(sql)
44         # Check for prohibited operations
45         for node in parsed.walk():
46             if hasattr(node, 'key'):
47                 if node.key.upper() in self.blocked_ops:
48                     return False
49         return True
50     except sqlglot.errors.ParseError:
51         return False
52
53     def _extract_sql(self, text: str) -> str:
54         """Extract SQL from model response."""
55         if '`sql`' in text:
56             return text.split("`sql`")[1] \
57                 .split("`")[0].strip()
58         if '```' in text:
59             return text.split("```")[1] \
60                 .split("`")[0].strip()
61         return text.strip()

```

Beyond Text-to-SQL: Conversational Analytics

Text-to-SQL is just the beginning. *Conversational analytics* extends this into multi-turn interactions where users explore data through dialogue:

- “Show me revenue by region for Q1.” (Initial query)
- “Break that down by product category.” (Refinement—system adds GROUP BY)
- “Why did Europe drop compared to last quarter?” (Analysis—system generates comparison query and interprets results)
- “Create a weekly alert if any region drops more than 15%.” (Action—system generates monitoring rule)

Each turn builds on the previous context, and the system maintains state across the conversation. This transforms data exploration from a technical skill into a natural conversation.

Note

Conversational analytics does not replace dashboards or BI tools. It complements them by handling ad hoc questions that do not warrant a new dashboard panel. The sweet spot is questions that are too specific for a pre-built dashboard but too common to justify a data engineering ticket every time—a.k.a. the “Can you just pull a quick number for me?” category that has consumed approximately 40% of all data engineering hours since the dawn of the profession.

25.5 The AI-Augmented Data Engineer’s Day

What does a typical day look like for an AI-augmented data engineer? The contrast with the traditional workflow is striking:

Time	Traditional Workflow	AI-Augmented Workflow
9:00 AM	Check dashboards, read through logs manually	CLAUDE CODE has triaged overnight failures, proposed fixes, and auto-resolved two routine issues
10:00 AM	Write dbt models from scratch	Describe the business logic; review and refine generated models
11:00 AM	Code review: read every line carefully	AI pre-reviews PRs with detailed comments; engineer focuses on architecture and business logic
1:00 PM	Debug a slow query by reading EXPLAIN plans	Feed query plan to CLAUDE CODE for optimization suggestions with rationale
2:00 PM	Write documentation from memory	CLAUDE CODE generates docs from code and commit history; engineer reviews for accuracy
3:00 PM	Handle ad hoc data requests from analysts	Direct analysts to the natural language query interface for self-service

Time	Traditional Workflow	AI-Augmented Workflow
4:00 PM	Sprint planning: estimate story points	Use CLAUDE CODE to analyze ticket descriptions, estimate complexity, identify risks

The Pair Programming Paradigm

The most productive data engineers treat CLAUDE CODE as a pair programming partner rather than an autocomplete engine. This means engaging in genuine collaboration:

- **Think out loud:** Explain your reasoning before asking for code. “I need to deduplicate events, but some events are legitimately repeated. The key should be `user_id + event_type + timestamp` within a 5-second window.” This context dramatically improves generated code quality.
- **Iterate, don’t perfect:** Start with a rough description and refine through conversation rather than spending 20 minutes crafting the perfect prompt. Three rounds of iteration typically produce better results than one perfect prompt.
- **Challenge:** Ask CLAUDE CODE to identify weaknesses in your approach. “What edge cases am I missing?” and “How would this break at 10x scale?” often surface issues you had not considered.
- **Learn:** Use responses as learning opportunities. When CLAUDE CODE suggests an approach you have not seen before, ask it to explain the tradeoffs. This accelerates skill development.
- **Verify:** Always run the code. Always check the output. Trust but verify, especially for SQL queries that touch production data.

Tip

Keep a personal “prompt journal”—a file where you save prompts that worked exceptionally well, along with notes on why they worked. Over time, this becomes a powerful reference that captures your personal prompt engineering expertise. Share the best entries with your team.

25.6 Career Advice for Data Engineers in the AI Era

The data engineering job market is evolving rapidly. Here is practical guidance for navigating the transition.

Skills That Become More Valuable

System Design and Architecture Understanding how systems fit together, how data flows end-to-end, and how to design for reliability, scalability, and maintainability. AI cannot replace the architect who understands the business and its constraints.

Data Modeling Dimensional modeling, graph modeling, and schema design require deep understanding of the business domain. AI can suggest schemas, but a human must validate them against business reality.

Business Domain Expertise The engineer who understands healthcare billing, financial compliance, or supply chain logistics brings irreplaceable context that makes AI outputs more accurate and relevant.

Quality Engineering Designing comprehensive data quality frameworks, defining expectations, and building monitoring systems. As more code is AI-generated, quality assurance becomes more critical, not less.

Communication and Leadership Explaining technical decisions to stakeholders, mentoring junior engineers on AI collaboration, and building the case for AI investment require human skills that AI amplifies but does not replace.

Skills That Become Table Stakes

These skills remain necessary but are no longer sufficient to differentiate you:

- Boilerplate SQL and Python writing (AI handles first drafts)
- Configuration syntax memorization (AI generates configs from descriptions)
- API documentation lookup (AI synthesizes API usage from docs)
- Routine debugging of common errors (agents handle known patterns)
- Manual documentation writing (AI generates from code; you review and refine)

 **Warning**

AI is a leverage multiplier: the more expertise you bring, the more leverage you gain. A chainsaw is far more dangerous in the hands of someone who has never used one. A senior engineer who masters AI collaboration can produce the output of an entire team while maintaining quality standards that only experience provides. But a junior engineer using AI without foundational knowledge will produce plausible-looking code that fails in subtle, expensive ways. Invest in fundamentals first, then amplify with AI.

Career Development Strategies

1. **Go deep on architecture:** Invest in distributed systems, streaming architecture, and data modeling. These are the skills AI amplifies most.
2. **Build a portfolio of AI-augmented projects:** Show that you can use AI tools effectively, not just that you can code. Employers increasingly value this.
3. **Develop business domain expertise:** Specialize in an industry (finance, healthcare, e-commerce) where your domain knowledge makes AI outputs more valuable.
4. **Learn to evaluate AI output critically:** This is a meta-skill that applies across all tasks. Practice by asking CLAUDE CODE to solve problems you already know the answer to, then analyzing where it gets things right and wrong.
5. **Contribute to the community:** Write about your experiences with AI-augmented data engineering. The field is young enough that practitioners who share knowledge become recognized experts.

25.7 Predictions for 2027–2030

Predicting the future is inherently uncertain, but the following trajectories are supported by current trends and investment patterns.

2027: Integration Becomes Standard

- AI-native orchestrators ship with built-in LLM integration for pipeline generation, monitoring, and remediation. Dagster, Prefect, and next-generation tools offer “describe your pipeline” interfaces.
- Conversational data catalogs become the default way business users discover and understand data. Keyword search in catalogs feels as dated as command-line-only interfaces.
- 50%+ of new pipeline code in progressive organizations is AI-generated, human-reviewed. The review process itself is partially AI-assisted.
- MCP becomes the dominant standard for AI-to-tool integration, with community-maintained servers for every major data tool.

2028: Platform Convergence

- The distinction between data lake, data warehouse, and streaming platform blurs further. Unified platforms handle batch, streaming, and AI workloads through a single interface.
- Self-tuning pipelines adjust parallelism, resource allocation, and scheduling automatically based on workload patterns. Manual performance tuning becomes the exception, not the rule.
- Governments establish AI-specific data governance regulations. Data engineers who understand both AI capabilities and regulatory requirements become highly sought after.
- Agent-to-agent coordination matures. Specialized agents for different domains (ingestion, transformation, serving, quality) collaborate through shared protocols to resolve complex incidents.

2029–2030: The Autonomous Horizon

- Intent-driven platforms become the default for new projects. Engineers describe what they want; the platform figures out how to build it.
- Level 4 autonomous data operations are common in mature organizations. Agents handle 80%+ of routine operational tasks without human intervention.

- A single senior AI-augmented data engineer manages infrastructure that previously required 5–10 people. This does not mean fewer data engineers—it means more data infrastructure per engineer and more ambitious projects becoming feasible.
- The “AI data engineer” role is mainstream, combining pipeline expertise, ML literacy, prompt engineering, and domain knowledge.

Note

These predictions assume continued rapid progress in AI capabilities. If progress plateaus or hits diminishing returns, the timeline stretches, but the direction remains the same. The question is not *whether* these changes happen, but *when* and *how fast*. Of course, predicting the future of AI is itself the kind of task you probably should not trust an AI to do. The irony is not lost on us.

25.8 Building a Learning Culture Around AI

Organizations that successfully integrate AI into their data engineering practice share common cultural traits:

- **Psychological safety:** No blame for AI-related mistakes. When an AI-generated query causes an incident, the response is “How do we improve our review process?” not “Why did you trust the AI?”
- **Experimentation time:** 10–20% of sprint capacity dedicated to AI experimentation. This is not optional—teams that cannot experiment cannot learn, and teams that cannot learn cannot improve.
- **Shared learning:** Dedicated channels (Slack, Teams) for effective prompts, surprising interactions, lessons learned, and tool recommendations. Knowledge sharing accelerates adoption across the team.
- **AI champions:** Senior engineers who curate prompt libraries, run workshops, mentor on AI collaboration, and stay current with the rapidly evolving landscape. Every team should have at least one.
- **Structured learning cadence:** Weekly AI digests (curated links and tips), monthly deep dives (hands-on workshops on specific techniques), quarterly hackathons (build something new with AI), and annual strategy reviews (where are we, where are we going).

 Tip

Start a “Prompt of the Week” practice where team members share their most useful prompt discovery. This low-effort, high-impact practice builds collective prompt engineering expertise faster than any training program. Plus, it’s the one Slack channel people actually read.

25.9 Ethical Considerations

As AI becomes deeply embedded in data engineering, ethical considerations demand explicit attention:

- **Bias propagation:** AI-generated data transformations can inadvertently encode or amplify biases present in training data. Review AI-generated business logic for fairness implications, especially in customer-facing analytics.
- **Transparency:** Stakeholders have a right to know when data products are AI-generated or AI-assisted. Label AI-generated artifacts clearly.
- **Environmental impact:** Large language models consume significant compute resources. Use model tiering and prompt caching not just for cost savings but as a responsible resource usage practice.
- **Skill atrophy:** If engineers over-rely on AI, foundational skills can atrophy. Maintain deliberate practice of core skills (writing SQL from scratch, manual debugging) alongside AI-augmented workflows.

 Warning

Do not let AI-generated code become a “black box” in your pipeline. “I don’t know what it does, but it works” was already a bad excuse for human-written code; it does not get better with AI-generated code. Every AI-generated transformation should be reviewable, testable, and explainable. If you cannot explain what a piece of AI-generated code does and why, it should not be in production.

25.10 Final Thoughts and Call to Action

The AI-augmented data engineer rethinks their approach to every aspect of the profession:

- Instead of “How do I write this pipeline?” they ask “What should this pipeline accomplish, and how do I verify it works correctly?”
- Instead of being the bottleneck for every data request, they build natural language interfaces that enable self-service.
- Instead of guarding knowledge in their heads, they encode it in AI-augmented workflows, RAG systems, and prompt libraries that scale beyond any individual.
- Instead of fearing obsolescence, they embrace the force multiplier that makes their expertise more valuable and impactful than ever.

The transition is not always smooth. There will be failed experiments, AI-generated bugs in production, and moments of frustration when the model does not understand what seems obvious. But the trajectory is clear: AI-augmented data engineering is not the future—it is the present, and the gap between teams that embrace it and those that resist it will only widen.

📌 Your Call to Action

Start today:

1. Pick one repetitive task you do weekly and automate it with CLAUDE CODE.
2. Share what you learn—the wins and the failures—with your team.
3. Build a personal prompt library and iterate on it.
4. Experiment with agentic workflows for a low-risk operational pipeline.
5. Measure the impact: time saved, quality improved, incidents prevented.
6. Advocate for broader adoption based on your real results, not hype.

The data engineers who start now will be the architects and leaders of tomorrow’s AI-native data platforms. The rest will be writing “we should have adopted AI sooner” in their retrospectives. Written by CLAUDE CODE, of course.

25.11 Exercises

1. **Self-Assessment:** Assess your current skills across four dimensions: technical depth (SQL, Python, distributed systems), AI collaboration (prompt engineering, output validation), breadth (cloud platforms, streaming, ML fundamentals), and leadership (communication, mentoring, architecture). Rate yourself 1–5 on each, identify your three largest gaps, and create a six-month development plan with specific actions and milestones.
2. **AI-Native Pipeline Design:** Redesign a pipeline you currently maintain using the intent-driven pattern. Write the natural language manifest (similar to the example in this chapter), then compare the manifest to your existing implementation. What information is captured in your code that the manifest misses? What does the manifest express more clearly?
3. **Text-to-SQL Prototype:** Build a text-to-SQL service for one of your databases using the pattern in this chapter. Test it with ten real questions from your stakeholders. Measure accuracy (correct SQL generated), safety (no prohibited operations), and usability (would the stakeholder accept this answer?). Identify the failure modes and propose improvements.
4. **Daily Workflow Audit:** Track every activity you perform for one full work week. For each activity, estimate how much time CLAUDE CODE could save (0%, 25%, 50%, 75%, or 90%). Calculate total potential savings. Identify the top three candidates for AI augmentation and implement one of them.
5. **Future Architecture:** Design a complete AI-native data platform for an e-commerce company processing 100M events per day. Include: the platform intent manifest, data product specifications for three key products (customer 360, real-time recommendations, executive dashboard), the team structure and roles required, the maturity model progression plan for the first 18 months, and cost projections for both compute and AI API usage.

Claude API Reference Quick Guide

Claude Code API Reference Quick Guide

“In case of emergency, consult this appendix. In case of existential crisis about whether you should have used CLAUDE CODE for this, also consult this appendix.”

Tip

This appendix provides a condensed reference for the CLAUDE CODE API. For current information, consult docs.anthropic.com. This book was written by a human and may go stale. The docs are updated by—well, probably also CLAUDE CODE at this point.

Messages API Endpoint Reference

Property	Value
URL	POST https://api.anthropic.com/v1/messages
Authentication	x-api-key header
Content Type	application/json
API Version	anthropic-version: 2023-06-01

Key parameters: `model` (required), `messages` (required), `max_tokens` (required), `system`, `temperature` (0.0–1.0), `stream`, `tools`, `tool_choice`, `stop_sequences`.

```
curl https://api.anthropic.com/v1/messages \
```

```

2  -H 'Content-Type: application/json' \
3  -H 'x-api-key: $ANTHROPIC_API_KEY' \
4  -H 'anthropic-version: 2023-06-01' \
5  -d '{
6    "model": "claude-sonnet-4-20250514",
7    "max_tokens": 1024,
8    "system": "You are a SQL expert for Snowflake.",
9    "messages": [{"role": "user", "content": "Write a query to find
10   duplicates"}]
  }'
```

Listing 1: Basic Messages API request

Model IDs and Selection Guide

Model ID	Context	Speed	Cost	Best For
claude-opus-4-20250514	200K	Slowest	Highest	Complex architecture, nuanced review
claude-sonnet-4-20250514	200K	Medium	Medium	General coding, SQL, most DE tasks
claude-haiku-3-5-20241201	200K	Fastest	Lowest	Classification, extraction, high-volume

Note

Pin versioned model IDs in configuration files for production to avoid unexpected behavior changes. “Our pipeline broke because the model got smarter overnight” is a real sentence that real engineers have said. Pin your versions.

Rate Limits and Error Codes

Rate limits vary by tier (Tier 1: 50 req/min, Tier 4: 4,000 req/min). Key error codes: 400 (invalid request), 401 (auth error), 429 (rate limit—implement exponential backoff with jitter), 500 (server error—retry), 529 (overloaded—retry with backoff).

 **Warning**

Always implement exponential backoff with jitter for 429 responses. Do not retry immediately—hammering a rate-limited endpoint is the API equivalent of pressing the elevator button faster to make it arrive sooner. It does not work and annoys everyone involved.

Tool Use Schema Reference

Tool use enables CLAUDE CODE to invoke external functions. The flow: (1) send tools in the request, (2) CLAUDE CODE responds with `tool_use` blocks, (3) execute the tool and return results, (4) CLAUDE CODE incorporates results.

```
1 from anthropic import Anthropic
2 client = Anthropic()
3
4 def run_tool_loop(user_message, tools, tool_handlers):
5     messages = [{"role": "user", "content": user_message}]
6     while True:
7         response = client.messages.create(
8             model="claude-sonnet-4-20250514", max_tokens=4096,
9             tools=tools, messages=messages,
10        )
11        tool_results = []
12        final_text = []
13        for block in response.content:
14            if block.type == "text":
15                final_text.append(block.text)
16            elif block.type == "tool_use":
17                result = tool_handlers[block.name](block.input)
18                tool_results.append({
19                    "type": "tool_result",
20                    "tool_use_id": block.id,
21                    "content": str(result),
22                })
23        if not tool_results:
24            return "\n".join(final_text)
25        messages.append({"role": "assistant", "content": response.
26            content})
27        messages.append({"role": "user", "content": tool_results})
```

Listing 2: Tool use flow

Streaming and Batch API

Streaming: Use `client.messages.stream()` for incremental results in interactive tools.

Batch API: Process large volumes asynchronously within 24 hours at 50% cost reduction—ideal for bulk documentation, quality scans, and batch SQL validation.

! Tip

Use the Batch API for any task that does not require immediate results—nightly quality scans, weekly documentation updates, or bulk schema analysis. Patience saves 50%. If only that worked at restaurants.

Python SDK Quick Reference

```
1 from anthropic import Anthropic, AsyncAnthropic
2
3 client = Anthropic() # Reads ANTHROPIC_API_KEY from env
4
5 # Simple message
6 response = client.messages.create(
7     model="claude-sonnet-4-20250514", max_tokens=1024,
8     messages=[{"role": "user", "content": "Hello, Claude Code"}],
9 )
10 print(response.content[0].text)
11
12 # With prompt caching
13 response = client.messages.create(
14     model="claude-sonnet-4-20250514", max_tokens=1024,
15     system=[{
16         "type": "text", "text": large_schema_context,
17         "cache_control": {"type": "ephemeral"}
18     }],
19     messages=[{"role": "user", "content": "List all tables"}],
20 )
```

Listing 3: Python SDK essentials

Note

Both Python and TypeScript SDKs read `ANTHROPIC_API_KEY` from the environment. In production, store keys in a secrets manager and inject at runtime. If you hardcode your API key, you will end up on a very expensive and very public list of “things found in GitHub repos.”

Prompt Engineering Patterns for Data Engineers

“A prompt engineer is just a regular engineer who has learned to say ‘please’ and ‘be specific’ to a computer.”

This appendix provides a catalog of reusable *prompt engineering patterns* for data engineering workflows. Think of these as copy-paste recipes for talking to CLAUDE CODE without it going off on a creative tangent about the philosophy of NULL values. Each pattern includes: the problem it solves, a parameterized template, and guidance on when to apply it.

The Anatomy of an Effective Data Engineering Prompt

i Note

Every effective data engineering prompt contains four layers: **Context** (system and schema), **Intent** (what to produce), **Constraints** (rules to obey), and **Format** (output structure). We call this the *CICF framework*.

The single most impactful improvement is including actual DDL for relevant tables. CLAUDE CODE reasons about column types, foreign keys, and constraints far more accurately given real schemas rather than prose descriptions. Telling CLAUDE CODE “I have a table with some columns” and expecting correct SQL is like telling a chef “I have some ingredients” and expecting a Michelin-star meal.

! Tip

Use `SHOW CREATE TABLE` or your ORM's schema dump. For large schemas, include only relevant tables plus their foreign key references.

Pattern 1: Schema-Aware SQL Generation

Problem: Generate SQL queries correct with respect to a specific schema.

Code Example

Template

```
1 You are a data engineer working with [DATABASE_PLATFORM].
2
3 ## Schema
4 ```sql
5 [PASTE DDL HERE]
6 ```
7
8 ## Task
9 Write a SQL query that [OBJECTIVE].
10
11 ## Requirements
12 - Use CTEs for readability. Handle NULLs explicitly.
13 - Optimize for [PERFORMANCE_GOAL].
14 - Target [SQL DIALECT] syntax.
```

Pattern 2: Documentation Generator

Problem: Generate documentation from existing code or schemas.

Code Example

Template

```
1 You are a technical writer for data engineering.
2
3 ## Source Material
```

```
4  ``` [LANGUAGE]
5  [PASTE CODE, DDL, OR CONFIGURATION]
6  ```
7
8  ## Task
9  Generate [table / pipeline / API] documentation.
10 - Audience: [junior engineers / analysts / stakeholders]
11 - Include: purpose, inputs, outputs, dependencies, limitations.
12 - Format: [Markdown / YAML / Confluence]
```

Pattern 3: Code Review Assistant

Problem: Review data engineering code for correctness, performance, security, and maintainability.

Code Example

Template

```
1  Senior data engineer code review.
2
3  ## Code
4  ``` [LANGUAGE]
5  [PASTE CODE]
6  ```
7
8  ## Context
9  - Environment: [ENV]. Volume: [VOLUME]. Frequency: [FREQ].
10
11 ## Evaluate
12 1. Correctness (logic errors, edge cases, NULL handling)
13 2. Performance (joins, indexes, full scans)
14 3. Security (injection, credentials, permissions)
15 4. Resilience (error handling, retry, idempotency)
16
17 ## Output per finding
18 - Severity: Critical / Major / Minor
19 - Location, Issue, Recommendation, Example fix
```

 **Warning**

Never blindly apply code review suggestions. CLAUDE CODE may flag idiomatic patterns as issues—it once suggested “fixing” a perfectly valid Python walrus operator because it looked like a typo. Always validate against your specific context.

Pattern 4: Data Quality Analyzer

Problem: Generate data quality checks for a dataset or table.

Code Example

Data Quality Analyzer Template

```
1 You are a data quality engineer.
2
3 ## Table Schema
4 ```sql
5 [PASTE DDL HERE]
6 ```
7
8 ## Business Context
9 [DESCRIBE WHAT THIS TABLE REPRESENTS]
10
11 ## Task
12 Generate quality checks covering:
13 1. Completeness (NULL rates for each column)
14 2. Uniqueness (PK/NK uniqueness)
15 3. Validity (ranges, formats, referential integrity)
16 4. Consistency (cross-column logic)
17 5. Timeliness (data freshness)
18 6. Statistical (volume anomalies vs. historical)
19
20 Return SQL queries for [PLATFORM], one per check.
```

Worked Example

Code Example

Generated Data Quality Checks

```
1  -- CHECK 1: Primary key uniqueness (expect 0 duplicates)
2  SELECT order_id, COUNT(*) AS cnt
3  FROM orders GROUP BY order_id HAVING COUNT(*) > 1;
4
5  -- CHECK 2: Referential integrity (expect 0 orphans)
6  SELECT o.order_id, o.customer_id
7  FROM orders o
8  LEFT JOIN customers c ON o.customer_id = c.customer_id
9  WHERE c.customer_id IS NULL;
10
11 -- CHECK 3: Temporal consistency (ship never before order)
12 SELECT order_id, order_date, ship_date
13 FROM orders WHERE ship_date < order_date;
14
15 -- CHECK 4: Volume anomaly (daily counts within 2 std devs)
16 WITH daily_counts AS (
17     SELECT DATE_TRUNC('day', created_at) AS dt, COUNT(*) AS cnt
18     FROM orders
19     WHERE created_at >= DATEADD('day', -30, CURRENT_DATE)
20     GROUP BY 1
21 ), stats AS (
22     SELECT AVG(cnt) AS avg_cnt, STDDEV(cnt) AS std_cnt
23     FROM daily_counts
24 )
25 SELECT dc.dt, dc.cnt
26 FROM daily_counts dc CROSS JOIN stats s
27 WHERE ABS(dc.cnt - s.avg_cnt) > 2 * s.std_cnt;
```

Pattern 5: Migration Script Generator

Problem: Generate safe, reversible database migration scripts.

Provide current DDL, desired DDL, platform, table size. Request UP and DOWN scripts with pre-migration validation, transactions, and post-migration checks. Specify zero-downtime requirements.

Pattern 6: Incident Root Cause Analyzer

Problem: Systematically diagnose a pipeline failure from logs and system state.

Provide: incident summary, error logs (with timestamps), system context (pipeline, orchestrator, platform), recent changes, and steps already tried. Request: error categorization, top 3 causes ranked by probability, diagnostic queries for each, and immediate mitigation plus long-term fix.

Tip

Include timestamps in log entries and preserve chronological order. Redact secrets but keep connection strings (with masked passwords) for diagnostic clues.

Pattern 7: Cost Optimization Advisor

Problem: Identify optimization opportunities in queries, storage, or compute.

Provide: platform, monthly spend, expensive query with execution plan, and optimization goals. Request: line-by-line cost driver analysis, recommendations ranked by savings, rewritten queries, and trade-off assessment.

Pattern 8: Data Contract Generator

Problem: Define formal data contracts between producers and consumers.

Code Example

Data Contract Generator Template

```
1 You are a data governance specialist.
2
3 ## Dataset
4 - Name: [NAME]. Owner: [TEAM]. Description: [WHAT IT IS].
5
6 ## Schema
7 ```sql
8 [PASTE DDL]
9 ```
10
```

```
11 ## Consumers
12 [LIST TEAMS AND HOW THEY USE THE DATA]
13
14 ## Task
15 Generate YAML data contract with:
16 1. Schema: every field with type, description, PII flag
17 2. Quality: completeness, uniqueness, freshness rules
18 3. SLA: latency, availability, incident response time
19 4. Ownership: producer, consumers, escalation contacts
20 5. Change management: how changes are proposed/reviewed
```

Worked Example

Code Example

Generated Data Contract

```
1 dataContract:
2   name: user_transactions
3   version: "2.1.0"
4   owner:
5     team: payments-data-eng
6     slack: "#payments-data"
7   schema:
8     fields:
9       - name: transaction_id
10         type: STRING
11         nullable: false
12         pii: false
13         primaryKey: true
14       - name: user_id
15         type: STRING
16         nullable: false
17         pii: true
18       - name: amount_cents
19         type: INTEGER
20         nullable: false
21         validRange: {min: 1, max: 100000000}
22   quality:
23     completeness: {transaction_id: 100%, user_id: 100%}
24     uniqueness: [transaction_id]
```

```
25     freshness: {maxDelayMinutes: 30}
26   sla:
27     availability: 99.9%
28     maxLatencyMinutes: 30
```

Pattern 9: dbt Model Generator

Problem: Create dbt models following project conventions.

Provide: naming conventions, materialization defaults, testing strategy, source DDL, model requirements (name, layer, grain, business logic, incremental strategy). Generate: SQL model with Jinja/ref(), schema YAML with tests, and source definitions.

Pattern 10: Terraform Module Generator

Problem: Generate IaC for data infrastructure components.

Code Example

Template

```
1 Cloud infrastructure engineer for data platforms.
2
3 ## Requirements
4 - Provider: [AWS / GCP / Azure]. Component: [DESCRIBE].
5 - Security: [encryption, network, IAM, compliance]
6 - Operations: [scaling, monitoring, backup]
7
8 ## Task
9 Generate: main.tf, variables.tf, outputs.tf, usage example.
```

Worked Example

Code Example

Generated Terraform – S3 Data Lake

```
1 resource "aws_s3_bucket" "data_lake" {
2   bucket = "${var.project}-${var.environment}-data-lake"
3   tags = merge(var.common_tags, {
4     Layer = "bronze"
5     ManagedBy = "terraform"
6   })
7 }
8
9 resource "aws_s3_bucket_versioning" "data_lake" {
10  bucket = aws_s3_bucket.data_lake.id
11  versioning_configuration { status = "Enabled" }
12 }
13
14 resource "aws_s3_bucket_server_side_encryption_configuration" "
15   data_lake" {
16   bucket = aws_s3_bucket.data_lake.id
17   rule {
18     apply_server_side_encryption_by_default {
19       sse_algorithm = "aws:kms"
20       kms_master_key_id = var.kms_key_arn
21     }
22   }
23 }
24 resource "aws_s3_bucket_public_access_block" "data_lake" {
25   bucket = aws_s3_bucket.data_lake.id
26   block_public_acls = true
27   block_public_policy = true
28   ignore_public_acls = true
29   restrict_public_buckets = true
30 }
```

Anti-Patterns to Avoid

Warning

Vague requests: “Write SQL for analytics” with no schema or metrics. CLAUDE CODE is good, but it cannot read your mind. Yet. Fix: provide DDL and specific objectives.

Kitchen sink: Asking for schema design, ETL, tests, docs, and monitoring in one prompt. This is the prompt equivalent of a 47-item grocery list scrawled on a napkin. Fix: break into sequential prompts.

No platform specified: Optimization strategies differ radically between Snowflake, BigQuery, and Spark SQL. “SQL is SQL” said no one who has ever debugged a dialect mismatch.

Ignoring data volume: A correct solution for 10K rows may fail at 10B rows. “It worked in dev” does not count.

Pasting credentials: Security risk and biases output. Fix: sanitize with synthetic data. Your API key does not need to be in a book.

No output examples: “Generate a report” without showing format. Fix: include 2–3 rows of desired output.

Trusting without validation: Never run AI-generated SQL against production without review. We cannot stress this enough. We have stressed it in every chapter. We will stress it again.

Composing Patterns for Complex Workflows

New source onboarding: Pattern 1 (SQL gen) → Pattern 9 (dbt) → Pattern 4 (quality) → Pattern 2 (docs) → Pattern 8 (contract).

Schema migration: Schema Evolution Advisor → Pattern 5 (migration) → Pattern 9 (dbt update) → Pattern 4 (validation) → Pattern 3 (code review).

Incident response: Pattern 6 (RCA) → Pipeline Debug → Pattern 1 (fix) → Pattern 4 (prevention) → Pattern 2 (runbook update).

Use patterns when the task is repeatable, high-stakes, involves complex context, or needs consistent quality across team members. For quick one-off questions, freeform prompts work fine.

i Note

Save customized templates, successful outputs, and lessons learned in a shared repository. Over time, your team builds a library of battle-tested prompts encoding your organization's standards. Think of it as institutional memory that does not leave when someone accepts an offer from a competitor.

Cost Optimization Strategies

“The cloud was supposed to save us money. Then AI was supposed to save us money. At some point, we should probably stop ‘saving’ money.”

Deploying CLAUDE CODE at scale requires careful cost management. This appendix covers understanding, monitoring, and optimizing API spending—because the only thing worse than a broken pipeline is a working pipeline that costs more than the data it produces.

Understanding Claude Code Pricing

Token Basics

A *token* is roughly 3–4 characters of English text. Code tokenizes less efficiently—SQL averages about 2.5 characters per token.

Table 3: Approximate token counts for common artifacts.

Artifact	Chars	Tokens
Simple SELECT (5 cols)	200	60
Complex CTE query (100 lines)	3,000	900
Full schema (50 tables)	25,000	7,500
Python ETL script (200 lines)	6,000	1,800

Table 4: Claude Code model pricing (per million tokens, USD).

Model	Input (\$/1M)	Output (\$/1M)	Context
Haiku	\$0.25	\$1.25	200K
Sonnet	\$3.00	\$15.00	200K
Opus	\$15.00	\$75.00	200K

Per-Model Pricing

Warning

Output tokens are 5× more expensive than input tokens. CLAUDE CODE is like a consultant: listening is cheap, but the moment it starts talking, the meter is running hard. Controlling output length is one of the most effective cost levers. Including output format examples reduces verbosity.

Token Counting and Estimation

Code Example

Token Counting and Cost Estimation

```
1 import anthropic
2
3 client = anthropic.Anthropic()
4
5 def count_tokens(text, model="claude-sonnet-4-20250514"):
6     response = client.messages.count_tokens(
7         model=model,
8         messages=[{"role": "user", "content": text}]
9     )
10    return response.input_tokens
11
12 def estimate_cost(input_tokens, output_tokens, model="sonnet"):
13     pricing = {
14         "haiku": {"input": 0.25, "output": 1.25},
15         "sonnet": {"input": 3.00, "output": 15.00},
16         "opus": {"input": 15.00, "output": 75.00},
17     }
```

```
18     p = pricing[model]
19     return round(
20         (input_tokens / 1e6) * p["input"]
21         + (output_tokens / 1e6) * p["output"], 6
22     )
```

Prompt Caching

Prompt caching reduces costs when you repeatedly send the same prompt prefix (system prompt, schema context). First request caches at $1.25\times$ base cost; subsequent reads at $0.10\times$ base.

Code Example

Prompt Caching Implementation

```
1 def query_with_caching(client, schema_context, question):
2     response = client.messages.create(
3         model="claude-sonnet-4-20250514", max_tokens=2048,
4         system=[{
5             "type": "text", "text": schema_context,
6             "cache_control": {"type": "ephemeral"}
7         }],
8         messages=[{"role": "user", "content": question}]
9     )
10    usage = response.usage
11    print(f"Cache read: {getattr(usage, 'cache_read_input_tokens', 0)}")
12    return response.content[0].text
```

A 5,000-token system prompt making 100 calls/hour saves $\sim 88\%$ on cached tokens—from \$264/month to \$31/month with Sonnet.

Tip

Prompt caching is especially valuable for data engineering because schema context is large and stable. A 50-table schema (7,000–10,000 tokens) saves significantly with dozens of queries per hour.

Model Selection Optimization

Table 5: Model selection guide.

Task	Haiku	Sonnet	Opus
Simple SQL, validation, docs, log parsing	✓		
Complex SQL, dbt, code review, migration		✓	
Architecture, complex debug, novel design			✓

Note

Start with Haiku for every new workload. Step up to Sonnet if quality is insufficient. Reserve Opus for tasks where Sonnet demonstrably falls short. This is the “try the free sample before buying the bottle” principle of AI cost management.

Batch API for Bulk Processing

The *Batch API* processes requests within 24 hours at 50% discount. Ideal for: bulk documentation, warehouse-wide quality analysis, migration code review, and dataset enrichment.

Cost Monitoring

Code Example

Cost Tracking Module

```

1 import sqlite3
2 from datetime import datetime, timedelta
3 from dataclasses import dataclass
4
5 @dataclass
6 class APICallRecord:
7     timestamp: str
8     model: str
9     workload: str

```

```
10     input_tokens: int
11     output_tokens: int
12     cost_usd: float
13
14 class CostTracker:
15     PRICING = {
16         "claude-haiku-4-20250514": {"input": 0.25, "output": 1.25},
17         "claude-sonnet-4-20250514": {"input": 3.00, "output": 15.00},
18         "claude-opus-4-20250514": {"input": 15.00, "output": 75.00},
19     }
20
21     def __init__(self, db_path="cost_tracking.db"):
22         self.conn = sqlite3.connect(db_path)
23         self.conn.execute("""
24             CREATE TABLE IF NOT EXISTS api_calls (
25                 id INTEGER PRIMARY KEY AUTOINCREMENT,
26                 timestamp TEXT, model TEXT, workload TEXT,
27                 input_tokens INTEGER, output_tokens INTEGER,
28                 cost_usd REAL
29             )""")
30
31     def calculate_cost(self, model, input_tokens, output_tokens):
32         p = self.PRICING.get(model, {"input": 3.0, "output": 15.0})
33         return round(
34             (input_tokens / 1e6) * p["input"]
35             + (output_tokens / 1e6) * p["output"], 8
36         )
37
38     def record(self, model, workload, usage):
39         cost = self.calculate_cost(
40             model, usage.input_tokens, usage.output_tokens
41         )
42         self.conn.execute(
43             "INSERT INTO api_calls VALUES (NULL,?,?,?,?,?,?)",
44             (datetime.utcnow().isoformat(), model, workload,
45              usage.input_tokens, usage.output_tokens, cost)
46         )
47         self.conn.commit()
48
49     def daily_summary(self, days=30):
50         cutoff = (datetime.utcnow() - timedelta(days=days)).isoformat()
51         ()
52         return self.conn.execute("""
53             SELECT DATE(timestamp), model, COUNT(*),
```

```

53         SUM(cost_usd) FROM api_calls
54     WHERE timestamp >= ? GROUP BY DATE(timestamp), model
55     ORDER BY 1 DESC""", (cutoff,)).fetchall()

```

Budget Allocation

Table 6: Example monthly budget for a 10-person team.

Category	Model	Calls/mo	Cost/mo
Ad hoc SQL	Sonnet	3,000	\$315
Quality checks	Haiku	12,000	\$66
Code review (CI/CD)	Sonnet	2,000	\$240
Complex analysis	Opus	300	\$360
Buffer (16%)	Mixed	—	\$237
Total			\$1,440

! Tip

Always include a 15–20% buffer. Data engineering workloads are spiky—incidents, migrations, and new source onboarding cause temporary increases. “We didn’t budget for the production fire” is not a line item anyone wants to explain.

ROI Calculation Framework

Code Example

ROI Calculation

```

1 def calculate_roi(
2     engineer_hourly_rate: float,
3     tasks_per_month: int,
4     avg_manual_hours: float,
5     avg_assisted_hours: float,
6     monthly_api_cost: float,

```

```
7 ) -> dict:
8     manual_cost = tasks_per_month * avg_manual_hours *
9     engineer_hourly_rate
10    assisted_cost = (
11        tasks_per_month * avg_assisted_hours * engineer_hourly_rate
12        + monthly_api_cost
13    )
14    savings = manual_cost - assisted_cost
15    hours_saved = tasks_per_month * (avg_manual_hours -
16    avg_assisted_hours)
17    return {
18        "monthly_savings": round(savings, 2),
19        "hours_saved": round(hours_saved, 1),
20        "roi_pct": round((savings / monthly_api_cost) * 100, 1),
21    }
22
23 # Example: SQL generation workload
24 result = calculate_roi(
25     engineer_hourly_rate=95.0, tasks_per_month=150,
26     avg_manual_hours=0.5, avg_assisted_hours=0.15,
27     monthly_api_cost=315.0,
28 )
29 # savings: $4,672, hours saved: 52.5, ROI: 1483%
```

Beyond time savings, consider: quality improvement (fewer production bugs), knowledge democratization (junior engineers producing senior-level work), faster onboarding, and reduced context switching. Also consider the intangible value of not being woken up at 3 AM because CLAUDE CODE handled the incident while you were sleeping. You cannot put a price on sleep. Actually, you can—it's roughly \$0.50 per incident resolution.

Cost Per Task Reference

i Note

Costs assume no prompt caching. With caching and stable system prompts, input costs for repeated workloads drop by up to 90%.

Table 7: Cost per task (Sonnet, typical token volumes).

Use Case	In Tok	Out Tok	Cost/Call	Cost/1000
Simple SQL query	800	400	\$0.008	\$8.40
Complex SQL (CTEs)	3,000	1,500	\$0.032	\$31.50
Code review	5,000	2,000	\$0.045	\$45.00
dbt model generation	2,500	1,800	\$0.035	\$34.50
Data quality suite	2,000	3,000	\$0.051	\$51.00
Migration script	4,000	3,500	\$0.065	\$64.50
Terraform module	3,000	4,000	\$0.069	\$69.00

Tips for Reducing Token Usage

1. **Concise system prompts:** Every token is charged on every call. Cut verbose instructions.
2. **Relevant schema only:** Do not paste 50 tables when only 3 are needed.
3. **Structured output:** Request JSON/SQL code blocks rather than prose (40–60% shorter).
4. **Appropriate max_tokens:** 512 for simple queries, 2048 for complex generation.
5. **Stop sequences:** Halt generation at predictable boundaries.
6. **Compressed schema:** Use shorthand notation (~40 tokens/table vs. ~160 for full DDL).
7. **Batch related questions:** One call for 5 tables beats 5 separate calls.
8. **Cache aggressively:** Any constant text across calls should be cached.

i Note

Cost optimization is ongoing. Like flossing, everyone agrees it should be done regularly, and almost nobody does it until something hurts. Set a monthly reminder to review your cost dashboard, identify trends, and adjust model routing and prompts. Small per-call savings compound at scale.

Glossary

“Every glossary is just a graveyard of terms the author had to explain more than twice.”

This glossary defines key terms used throughout this book, organized alphabetically. Some definitions are more serious than others. All are technically accurate. We think.

Agentic AI

An AI system that autonomously plans, uses tools, and takes multi-step actions to accomplish a goal, as opposed to single-shot Q&A. Essentially, an AI that has learned to procrastinate less than you.

Airflow

Apache Airflow. Open-source workflow orchestration platform for scheduling and monitoring data pipelines, defined as DAGs in Python.

API (Application Programming Interface)

Defined protocols and endpoints for software communication. In this book, usually the Anthropic Messages API.

Backfill

Retroactively processing historical data through a pipeline after a bug fix or new transformation. The data engineering equivalent of “we need to redo all of this.” Always takes longer than estimated.

Batch API

An Anthropic API feature that processes requests asynchronously within 24 hours at 50% discount.

Batch Processing

Processing data in scheduled chunks (hourly, daily) rather than continuously. Contrast with *Stream Processing*.

CDC (Change Data Capture)

Technique for capturing row-level changes from source databases via log reading (Debezium) or timestamp polling.

CI/CD (Continuous Integration / Continuous Deployment)

Practices automating building, testing, and deploying code—including pipeline code, dbt models, and schema migrations.

Claude Code

Anthropic’s official CLI for using CLAUDE CODE in software development workflows, including code generation, review, and debugging.

Columnar Storage

Storage format storing column values contiguously, optimized for analytical queries. Examples: Parquet, ORC.

Constitutional AI (CAI)

Anthropic’s AI alignment approach where a model follows a set of principles guiding its behavior.

Context Window

Maximum tokens (input plus output) a model processes per request. CLAUDE CODE supports up to 200K tokens. Still smaller than your average Jira ticket’s comment thread.

CTE (Common Table Expression)

A temporary named result set in SQL defined with WITH. Preferred pattern for complex transformations.

DAG (Directed Acyclic Graph)

Graph with directed edges and no cycles. Used in orchestrators for task dependencies and in dbt for lineage. The “acyclic” part is aspirational in some codebases.

Data Contract

Formal agreement between producer and consumers specifying schema, quality expectations, SLAs, and change procedures.

Data Governance

Policies, processes, and standards ensuring data is managed consistently, securely, and in regulatory compliance.

Data Lakehouse

Architecture combining data lake flexibility with warehouse management features (ACID, schema enforcement). Examples: Delta Lake, Iceberg, Hudi.

Data Lineage

Tracing data from origin through all transformations to final destination. Critical for debugging and compliance. Also critical for answering “where did this number come from?” in the board meeting.

Data Observability

Monitoring and understanding data pipeline health in real time—freshness, volume, schema, distribution, lineage.

dbt (data build tool)

Open-source transformation framework using SQL SELECT statements with built-in testing, documentation, and lineage.

Delta Lake

Open-source storage layer bringing ACID transactions and time travel to data lakes. Commonly used with Spark/Databricks.

Dimensional Modeling

Data modeling technique organizing data into fact tables (measurements) and dimension tables (context). Popularized by CLAUDE CODE’s training data, originally by Kimball.

ELT (Extract, Load, Transform)

Pattern where raw data is loaded into a target system first, then transformed in place. Contrast with *ETL*.

Embeddings

Dense vector representations capturing semantic meaning, enabling similarity search and RAG.

Extended Thinking

A CLAUDE CODE feature allowing additional internal reasoning steps before generating a response.

Few-Shot Prompting

Including examples of desired input-output behavior in the prompt to guide the model.

Hallucination

When an AI model generates plausible but factually incorrect information—invented column names, wrong SQL syntax, fictional functions. Like a confident intern on their first day, but with better grammar.

Iceberg

Apache Iceberg. Open table format supporting ACID transactions, schema evolution, and time travel for large analytic datasets.

Idempotency

Property where executing an operation multiple times produces the same result as once. Critical for pipeline retries. If only meetings had this property.

IaC (Infrastructure as Code)

Managing infrastructure through configuration files. Tools: Terraform, Pulumi, CloudFormation.

Kafka

Apache Kafka. Distributed event streaming platform for real-time data pipelines with durable pub-sub messaging.

Human

A biological entity occasionally consulted for coffee preferences and final sign-off on production deployments. Requires sleep, food, and encouragement. See also: *Legacy System*.

Large Language Model (LLM)

Deep learning model trained on vast text, capable of understanding and generating language. CLAUDE CODE is an LLM. Not to be confused with “Large Language Mistake,” which is what happens when you skip the review step.

Legacy Code

Code written by a human. Identifiable by its bugs, inconsistent variable naming, and comments that say `// TODO: fix this later` from 2017.

Materialization

In dbt, strategy for persisting results: `view`, `table`, `incremental`, or `ephemeral`.

MCP (Model Context Protocol)

Open protocol by Anthropic standardizing how AI models connect to external tools, data sources, and services.

Medallion Architecture

Data lake/lakehouse pattern: Bronze (raw), Silver (cleaned), Gold (business aggregations).

Model Routing

Directing tasks to different model tiers (Haiku, Sonnet, Opus) based on complexity, latency, or cost.

Parquet

Apache Parquet. Columnar file format optimized for analytics, the de facto standard for data lake storage.

PII (Personally Identifiable Information)

Data identifying a specific individual (name, email, SSN). Must be handled per privacy regulations.

Pipeline

Sequence of processing steps moving data from sources through transformations to a destination. Will break at 3 AM. Always at 3 AM.

Prompt Caching

API feature caching repeated prompt prefixes. Cached tokens charged at 10% of normal input rate.

Prompt Engineering

Crafting effective prompts: system prompt design, few-shot examples, constraints, and output formatting. The art of asking a computer to do something in exactly the right way, which is somehow now a six-figure job.

RAG (Retrieval-Augmented Generation)

Combining information retrieval with generation—retrieved documents ground the model's response in facts.

Rate Limiting

Restricting API calls per time period. Anthropic enforces limits on requests/min and tokens/min.

SCD (Slowly Changing Dimension)

Techniques for handling dimension changes over time. Type 1: overwrite. Type 2: new rows with date ranges. Type 3: add a column. Type 4: deep existential regret about schema design choices.

Sleep

A biological process required by humans but not by CLAUDE CODE. Historically interrupted by pipeline failures. See *Pipeline*.

Schema Evolution

Modifying schemas over time while maintaining compatibility with existing data and consumers.

Snowflake

Cloud-native data warehouse separating storage and compute, supporting multi-cloud deployment.

Spark

Apache Spark. Distributed computing framework for batch/stream processing, SQL, and ML.

Stream Processing

Processing data continuously as it arrives. Tools: Kafka Streams, Flink, Spark Structured Streaming.

System Prompt

Special prompt setting model behavior and constraints for an entire conversation, passed separately from user messages.

Technical Debt

The accumulated cost of shortcuts taken during development. Like real debt, except the interest rate is measured in on-call pages.

Temperature

Parameter controlling output randomness. Lower (0.0) = deterministic; higher (1.0) = creative. For data engineering, keep it low unless you enjoy surprises in your SQL. You do not enjoy surprises in your SQL.

Terraform

Open-source IaC tool using HCL to provision cloud infrastructure across providers.

Token

Fundamental text processing unit in LLMs, typically 3–4 characters. Billing is based on token counts. The modern equivalent of “the meter is running.”

Tool Use (Function Calling)

CLAUDE CODE capability to invoke external functions during conversations via structured tool call requests.

Transformer

Neural network architecture based on attention mechanisms. Foundation of all modern LLMs including CLAUDE CODE.

Upsert

Insert-or-update operation based on key match. Implemented via MERGE in most SQL dialects.

Vector Database

Database optimized for storing and querying high-dimensional embeddings. Examples: Pinecone, Weaviate, pgvector.

Watermark

In stream processing, a timestamp indicating all data up to that point has been received.

Window Function

SQL function calculating across related rows without collapsing them. Defined with `OVER()`. Examples: `ROW_NUMBER()`, `LAG()`.

YAML (YAML Ain't Markup Language)

Human-readable serialization format for configuration. Used by dbt, Kubernetes, GitHub Actions, and data contracts. “Human-readable” is doing a lot of heavy lifting in that definition, especially when indentation is off by one space.

Z-Order

Multi-dimensional clustering in Delta Lake/Databricks to co-locate related data, improving multi-column filter performance.

Zero-Shot Prompting

Giving a task description without examples, relying on pre-trained knowledge. Contrast with *Few-Shot Prompting*. The AI equivalent of “just figure it out.”