Introduction

# Why Data Engineering?

"The purpose of computing is insight, not numbers."

– Richard Hamming [1]

In the modern world, nearly every interaction generates data. When you stream music, order food through an app, or even tap your credit card at the store, those small actions produce digital traces. Yet raw data on its own does not drive decisions or build products. The systems that collect, organize, and transform it into something useful are designed by data engineers. They are the professionals who create the pipelines and platforms that turn chaotic information into business value.

Data engineering is no longer a niche technical role. It's one of the fastest growing careers in technology and business. While the U.S. Bureau of Labor Statistics (BLS) doesn't publish a separate category for data engineers, related roles provide a strong indicator of growth. For example, employment of data scientists is projected to grow 34 percent from 2024 to 2034, much faster than the average for all occupations. Database administrators and architects are expected to grow at about 4 percent, and software developers at about 15–16 percent over the same period [2], [3]. Industry reports also highlight that demand for data engineers is expanding rapidly, with many organizations creating new roles as they recognize the need for clean, reliable, and well-structured data to compete in a digital economy.

The career outlook is equally attractive when it comes to compensation. The average salary for data engineers in the United States is about 130,000 dollars per year, with entry-level engineers starting around 100,000 dollars and senior engineers earning more than 150,000 dollars, according to Glassdoor data [4]. In technology hubs such as the San Francisco Bay Area, salaries can rise above 200,000 dollars, especially at major companies like Google, Meta, and Dropbox. This level of compensation reflects the critical role data engineers play in modern organizations, where a broken pipeline can mean millions in lost revenue or missed insights.

What sets data engineering apart is its centrality. Without data engineers, data science models cannot be trained, dashboards cannot be built, and executives cannot make informed decisions. The entire ecosystem of analytics, artificial intelligence (AI), and

digital transformation depends on the reliability of the infrastructure that data engineers design. They are often described as the backbone of successful organizations because they enable everyone else to succeed with data.

For those considering a career path, data engineering offers more than stability and salary. It provides the opportunity to work at the intersection of technology, strategy, and innovation. Whether in finance, healthcare, e-commerce, or government, every sector is leaning on data engineering to power its future. That makes this a field worth mastering.

In this book, I'll walk you through the world of data engineering, from its foundational concepts to real-world applications. We'll explore how raw data is extracted, transformed, and loaded into systems where it can deliver value. Along the way, I'll break down the modern data stack, highlight the tools that power today's most data-driven companies, and walk through the principles that separate good pipelines from great ones. Whether you're just starting out or looking to deepen your expertise, this book aims to make the complex systems behind data engineering both accessible and actionable.

# Chapter 1
# History of Data Engineering

"The best way to predict the future is to invent it."

– Alan Kay [5]

# Historic Data Engineering

Every company today claims it wants to "become more data-driven." But the truth is, people have been data-driven for thousands of years.

We observed, measured, remembered, and acted on information out of necessity. Spotting animal tracks, tracking seasons, and learning from patterns in nature weren't optional. Rather, it was necessary for survival. In many ways, data engineering is just a modern extension of those ancient instincts, only now the datasets are bigger, and the stakes are measured in dollars instead of dinner.

For most of human history, data was physical: etched on clay tablets, scribbled in ledgers, stored in human memory. Ancient Egyptians recorded grain inventories and cultivable land [6]. Babylonians tracked celestial movements to predict seasons [7]. The Ancient Greeks kept detailed recorded voting outcomes from early democratic assemblies [8]. Early data "engineering" was done with pen and parchment and required discipline, logic, and long hours. These systems are archaic to us, sure, but they had structure: columns, rows, dates, categories. The ideas behind data schemas existed long before relational databases did.

While the term "data engineering" didn't exist until much later, the functions associated with it like moving data between systems, defining data structures, and ensuring consistency were already present in the era of hierarchical and network databases in the mid 20th century.

These early systems, like IBM's Information Management System (IMS) conceptualized in the 1960s, which stored data in a tree-like structure, and the CODASYL network model, which used a more flexible graph of predefined relationships, required engineers to manually define relationships, optimize access paths, and write complex low-level procedural code to retrieve data [9], [10]. Working with these databases demanded deep technical knowledge of the schema and system architecture, which meant that in practice, data access was limited to a small group of trained specialists.

| Customer |
| --- |
| Customer_ID |
| Name |

Places

| Order |
| --- |
| Order_ID |
| Date |

Contains

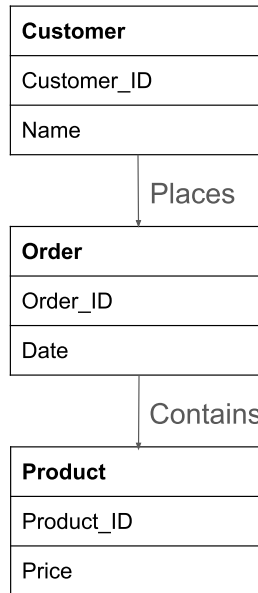| Product |
| --- |
| Product_ID |
| Price |

Figure 1: Simple CODAYSL Model

The shift from this rigid, expert-only approach to the more accessible and modular form of data engineering we recognize today didn't truly begin until the invention of the relational database by Edgar F. Codd, a British computer scientist at IBM, in the 1970s [11].

In his revolutionary paper, *A Relational Model of Data for Large Shared Data Banks*, Codd proposed a new model that directly addressed the limitations of the hierarchical and network databases that came before it. Specifically, he wrote:

"The relational view (or model) of data ... appears to be superior in several respects to the graph or network model ... It provides a means of describing data with its natural structure only—that is, without superimposing any additional structure for machine representation purposes [12]."

Primary Key

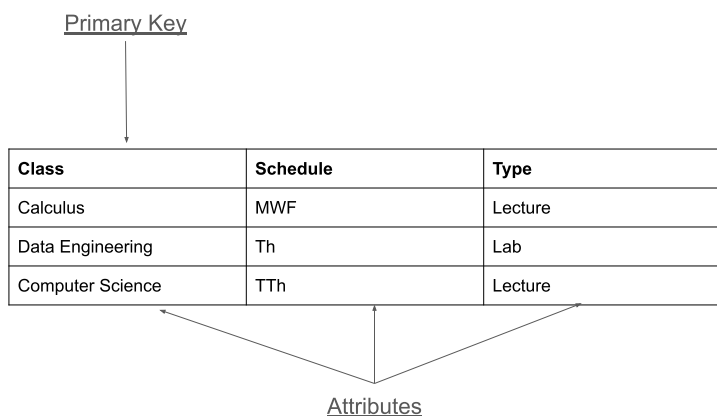| Class | Schedule | Type |
|---|---|---|
| Calculus | MWF | Lecture |
| Data Engineering | Th | Lab |
| Computer Science | TTh | Lecture |

Attributes

Figure 2: The Simple Relational Model

What Codd highlighted here was simple but profound. He argued that data should be modeled based on its real-world meaning, not the internal logic of how machines store or retrieve it. In doing so, he separated the logical structure of data (the relationship it represents) from its physical storage (how it's actually saved on disk).

This principle changed everything. It meant that engineers no longer had to manually define traversal paths, worry about pointer chains, or embed access logic directly into their application code. Instead, they could describe what data they wanted, and let the database engine figure out how to get it.

As a continuation of Codd's work, a team of IBM researchers began working on an experimental relational database system called System R in 1974 [13], [14]. The project was a success. Not only did it prove that relational databases could achieve competitive performance, it also laid the groundwork for what would eventually become the Structured Query Language or SQL, now the industry standard for interacting with relational data today.

Originally dubbed the Structured English Querying Language or SEQUEL, the researchers designed a language that allowed users to interact with data in a human-friendly way [15]. By using phras-

es like SELECT, WHERE, JOIN, etc. which resembled plain English, they aimed to make querying accessible to non-technical users [16]. This eliminated the need for writing complex procedural code and allowed people to simply describe the data they wanted, rather than detailing how to retrieve it.

Due to a trademark discrepancy with a UK-based aerospace company, IBM dropped the vowels from SEQUEL and renamed the language SQL to avoid legal complications [15]. Despite the change in spelling, many still continue to pronounce it as "sequel."

The success of System R and SQL didn't stay confined to IBM's labs. By the early 1980s, companies like Oracle, Ingres, and IBM itself had begun commercializing relational databases. SQL became the de facto standard, adopted not just by technologists but by entire organizations. It became the foundation for countless enterprise systems, from payroll to inventory management and helped businesses turn operational data into strategic insight.

This was the era when the data engineer role started to take shape. While it wasn't always called that, professionals were now tasked with designing schemas, building ETL pipelines, maintaining data integrity, and making data accessible across departments.

As organizations began collecting more data than ever before, a new challenge emerged: how do you store and analyze massive volumes of historical data across different systems, departments, and use cases? Traditional relational databases were optimized for transactions, that being things like updating a record, retrieving a customer profile, or processing an order. But they weren't designed to support large-scale analytics over time.

# Modern Data Engineering

The solution came in the concept of data warehousing, formalized in the 1990s by computer scientist Bill Inmon, who is widely regarded as the father of the data warehouse. According to Inmon, a data warehouse is not a specific technology but an architectural approach to managing data to support analysis and decision-making. As he defines it in *Building the Data Warehouse*, "A data warehouse is

a subject-oriented, integrated, nonvolatile, and time-variant collection of data in support of management's decisions."

Inmon emphasized that this distinction between architecture and technology is essential to understand. As he elegantly puts it: "There is the architecture, then there is the underlying technology, and they are two very different things. Unquestionably, there is a relationship between data warehousing and database technology, but they are most certainly not the same. Data warehousing requires the support of many different kinds of technology [17]."

Inmon was a proponent of a top-down approach to data warehousing, where an enterprise-wide data model is designed first, and individual data marts are created afterward to serve specific business needs. Instead of focusing on just one department at a time, the company would take a big-picture view and design the structure of the warehouse to serve the entire organization. Once that foundation is in place, smaller pieces of the system (data marts) can be built to serve the needs of individual teams like sales, marketing, or customer service.

While a good idea in theory, in practice it faced many limitations. For example, when companies attempted to implement this, they often ran into the problem of long development timelines. Because the data warehouse had to be fully designed and built to serve
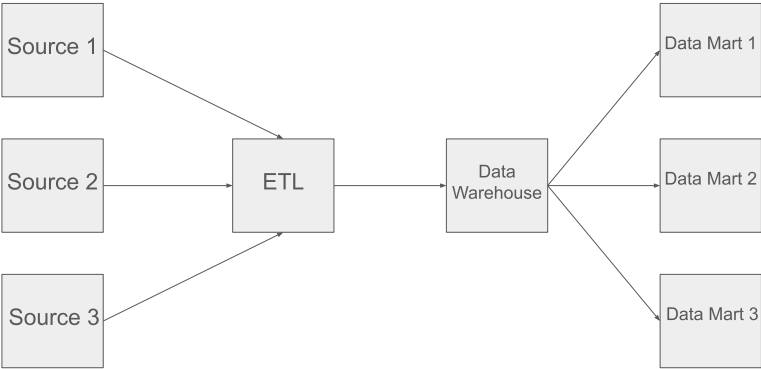


Figure 3: Inmon's Top-down Approach

the entire organization before any one team could use it, projects could take months or even years before delivering any real business value.

In addition to the delays, the initial cost of building such a large, centralized system was substantial. It required significant investment in infrastructure, tools, and specialized personnel. The approach also lacked flexibility. Business needs would often shift during development, but the rigid structure made it difficult to adapt quickly. Teams were frequently left waiting for updates or changes to be modeled and approved across the entire system before they could move forward [18].

These challenges left many organizations searching for a more practical and responsive alternative, something that could show results sooner and adjust as the business evolved.

Inmon's colleague Ralph Kimball proposed a different strategy, one that addressed many of the limitations found in the top-down model [19]. Instead of starting with a centralized data warehouse, Kimball suggested building individual data marts first. Each data mart would focus on a specific business area, such as sales, marketing, or inventory. These marts could then be connected using shared dimensions and rolled up into a broader data warehouse over time [20].
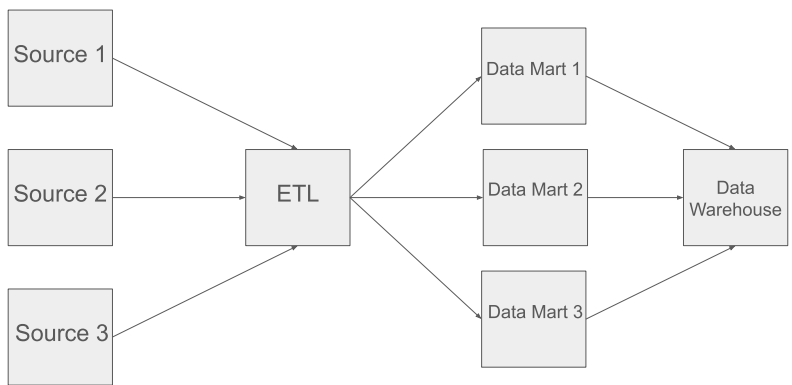


Figure 4: Kimball's Bottom-up Approach

This approach allowed for a much more agile way of building data infrastructure. Making changes no longer required reworking the entire system. Instead, teams could update or expand just the relevant data mart. The initial setup was also faster, making it easier for organizations to deliver value to the business early in the process, rather than waiting months or years for a fully centralized warehouse to be completed [21].

As data volumes continued to grow in the 2000s, traditional relational data warehouses began to show their limitations. Web-scale companies were now dealing with data that arrived continuously, in massive quantities, and in a variety of formats. This included user activity logs, clickstream data, product images, and free-form text which was far more complex than the clean, structured tables that warehouses were built to handle. These systems were expensive to scale, slow to adapt, and poorly suited for real-time or semi-structured data.

In response, engineers began turning to distributed computing systems, which distribute work across multiple machines to handle large workloads. Instead of relying on a single powerful server, these systems break a large task into smaller parts and run them in parallel on a cluster of computers. This is what gives the approach its name. It distributes both data and processing across a system.

This made it possible to store and analyze far larger datasets than traditional systems could handle. However, it also introduced new challenges. Engineers now had to figure out how to split the work effectively, manage coordination between machines, and design systems that were fault-tolerant (handle failures without losing progress or data). Despite the complexity, this shift enabled a level of scalability and speed that centralized warehouses simply could not achieve.

One of the most influential breakthroughs came from Google in 2004 with the introduction of the MapReduce abstraction. In the widely influential paper, *MapReduce: Simplified Data Processing on Large Clusters*, Google engineers Jeffrey Dean and Sanjay Ghemawat

outlined a programming model that made it possible to process vast amounts of data in parallel across distributed systems [22].

The key idea was to split work into two simple operations: map, which applies a function to each record in a dataset and produces intermediate key–value pairs, and reduce, which aggregates or summarizes those pairs by key. For example, mapping could count the occurrence of each word in a document, while reducing would sum those counts across all documents. This approach abstracted away many of the complexities of distributed computing, allowing engineers to focus on writing simple "map" and "reduce" functions while the underlying system handled fault tolerance, data distribution, and execution.
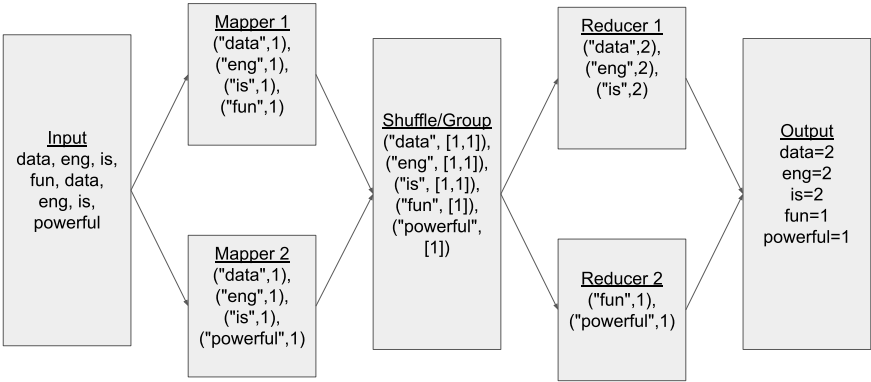


Figure 5: Simple MapReduce Example

The authors described a key challenge of the time: although the computations themselves were often conceptually simple, they became difficult to scale. As they put it, "The issues of how to parallelize the computation, distribute the data, and handle failures conspire to obscure the original simple computation with large amounts of complex code." MapReduce offered a way to simplify this, allowing developers to express their logic cleanly while the system took care of the difficult distributed systems engineering behind the scenes.

Inspired by the MapReduce paper, engineer Doug Cutting, along with Mike Cafarella, began building an open-source implementation. When the system outgrew its original purpose, they spun it out into a new Apache project, which Cutting named Hadoop, after his son's toy elephant [23].

With support from Yahoo!, Hadoop rapidly matured. As Cutting later reflected on his website, "After one year, Hadoop was used daily by many research groups within Yahoo!. After two years it generated Yahoo!'s web search index, achieving web-scale. Now, after three years, Hadoop holds the big-data sort record and the project has become a de-facto industry standard for big-data computing, used by scores of companies [24]."

The success of Hadoop was driven by Google's ideas and the engineering talent at Yahoo!. Their work made it possible for organizations outside of Big Tech to process data at scale using affordable, distributed infrastructure.

Despite its massive impact, Hadoop came with real challenges. It was powerful, but complex. Running a Hadoop cluster required deep technical expertise and significant operational overhead. Jobs were written in Java, making it inaccessible for analysts and business users. It was also fundamentally built for batch processing, which made it less suited for real-time analytics. As data grew faster and more diverse, and as companies sought more agile, cloud-friendly solutions, new tools began to emerge that prioritized simplicity, speed, and accessibility [25].

One of the most significant developments during this period was the emergence of Apache Spark. Initially developed at UC Berkeley's AMPLab in 2009, deployed in production by 2010, and later donated to the Apache Software Foundation in 2013, Spark was created to overcome many of the limitations of Hadoop and MapReduce.

Unlike previous systems, which required writing intermediate results to disk between steps, Spark supported running multiple operations over the same dataset in memory. This made it significantly more efficient for many workloads, especially those

involving iterative processing like machine learning or data exploration. Spark introduced a flexible programming model through a concept called resilient distributed datasets (RDDs), which allowed for fault-tolerant, parallel computation.

Perhaps even more importantly, Spark offered unified APIs in accessible languages like Python and Scala, expanding its usability beyond Java developers. As of today, Spark supports Python, SQL, Scala, Java, and R, making it one of the most versatile platforms in the data engineering ecosystem.

Spark could handle batch, streaming, and interactive big data workloads, making it far more adaptable than MapReduce. It quickly became a popular choice for building ETL pipelines, powering streaming machine learning workflows, and scaling analytics efforts. While Spark retained compatibility with the Hadoop Distributed File System (HDFS), its emphasis on speed, flexibility, and developer productivity marked a major shift in how modern data systems were built [26].

Among the earliest breakthroughs in cloud-native analytics was Google BigQuery, launched in 2011. Based on Google's Dremel technology, BigQuery introduced a fully managed, serverless approach to querying massive datasets with SQL [27]. Around the same time, Amazon Redshift, launched in 2012, brought data warehousing to the Amazon Web Services ecosystem, allowing teams to migrate their analytical workloads into the cloud [28].

A few years later, Snowflake redefined expectations with a platform designed from the ground up for cloud computing. By separating storage from compute, Snowflake made it possible to scale workloads up or down instantly and support multiple users or teams without performance bottlenecks.

As these cloud data warehouses gained adoption, they were often paired with modern business intelligence platforms like Tableau, Microsoft Power BI, and Google Looker. These tools gave analysts and stakeholders the ability to explore data visually, build dashboards, and generate reports directly from the warehouse. This convergence of scalable storage, familiar query languages, and in-

tuitive visualization marked the beginning of the modern analytics era.

And that brings us to where we are today. Data engineering is no longer just about moving data from one place to another. It's about designing systems that are scalable, reliable, and accessible across entire organizations. Engineers are expected to work across cloud platforms, orchestrate ETL/ELT workflows, maintain data quality, and empower analysts and stakeholders through well-modeled, accessible datasets.

The history of data engineering is the story of how we've learned to tame complexity at ever greater scales. What began as marks on clay tablets and ink in ledgers has evolved into massively parallel systems, real-time pipelines, and cloud-native platforms that power decisions across every modern industry. Each era, from relational databases to Hadoop to Spark and the modern data stack, emerged to solve the limitations of the last.

Today, data engineers sit at the center of this evolution. They build enterprise-grade systems that move and transform data, while also ensuring its reliability, accessibility, and trustworthiness across the organization. They are connectors between raw numbers and real-world insight.

Understanding how we got here provides context and direction for the future of our profession. The challenges we face today, whether that be scaling machine learning systems or maintaining data quality in real time, echo the same themes: how to handle more data, faster, with fewer errors, and greater impact.