# Why Learn Data Modelling and SQL in AI Systems?

A Research Report on the Role of Structured Data in AI Workflows

## Introduction

In the evolving world of data science and artificial intelligence (AI), the ability to model, organize, and retrieve data efficiently is crucial. Data modelling and SQL (Structured Query Language) are not just technical skills; they are foundational tools that ensure AI systems function reliably, ethically, and efficiently. This report explores how structured data design impacts machine learning (ML) systems and why mastering these tools is essential for any aspiring AI professional.

## 1. Data Retrieval Directly Affects ML Model Performance

The speed and accuracy of machine learning model training depend heavily on how data is stored and accessed. Large datasets must be queried, cleaned, and loaded efficiently—delays here slow down experimentation cycles and model updates. Structured SQL databases with proper indexing dramatically reduce this lag.

### Impact on Data Accessibility and Speed

Efficient data storage and retrieval directly impact the **speed** at which AI/ML models can access training data. Training a machine learning model requires feeding it large volumes of structured or unstructured data. If the underlying **data storage system** (e.g., SQL databases, data warehouses, or data lakes) is poorly designed or slow to retrieve data, it leads to **bottlenecks** in the training pipeline.

### Enabling Scalable Workflows

Modern AI systems require **scalable data retrieval methods**. Well-modeled databases allow for **partitioning**, **indexing**, and **caching**, which are crucial for handling big data in training environments. SQL queries with optimized indexes ensure fast filtering and aggregation, reducing the time needed to prepare training batches.

**Real-World Example**: Netflix uses **Apache Spark** with **SQL-based transformations** to process viewing history and user behaviour from massive datasets, enabling real-time model updates for personalized recommendations.

*Insight:* Fast data access leads to faster model iteration, improving time-to-insight and operational efficiency.

## 2. Well-Modeled Data Minimizes Technical Debt

Technical debt in ML systems arises from messy, unstructured, or redundant data. When databases are poorly designed, developers spend time fixing issues rather than building features or improving models. Clean schemas (like star or snowflake models) reduce errors, standardize formats, and enhance data consistency across systems.

*Insight:* Clean data structures reduce rework, lower operational risk, and support long-term scalability.

## 3. Governance and Monitoring Depend on Structured Databases

In regulated industries, AI systems must be explainable, auditable, and secure. This requires maintaining clear data lineage, logging changes, and enforcing access control. SQL-based systems allow for comprehensive monitoring, logging, and auditing—capabilities critical for compliance with laws like GDPR and HIPAA.

*Insight:* Structured databases enable transparency and accountability in AI decision-making.

**Real-World Examples**

- **Netflix** uses SQL-based tools like Presto and Hive to power ML models for user recommendations. Data engineers ensure structured formats and partitioning strategies to make querying petabytes of data efficient
- **Capital One** utilizes SQL-based data lakes for real-time fraud detection and auditing. The structured format allows their teams to track data usage and ensure regulatory compliance.

## Reflection: Connection to My Learning

As a participant in the **Omantel Data Leaders Program**, my learning journey has gone beyond just theoretical knowledge. Through hands-on projects and training in tools like **Power BI**, **SQL**, and **Python**, I've experienced firsthand how data modeling and querying directly influence the quality and success of analytics and AI solutions.

Working with real datasets, I've learned to identify key data structures, define relationships using **primary and foreign keys**, and apply **data normalization** to eliminate redundancy and improve consistency. These skills have enabled me to build more **accurate dashboards**, **clean data pipelines**, and **insightful reports**. What seemed at first like backend tasks are now clearly critical for building AI solutions that are **scalable**, **transparent**, and **reliable**.

This research reinforced the idea that **AI is not just about algorithms—it's about data**. If the foundation of data is weak, even the most advanced machine learning models will produce flawed results. Poor data quality leads to biased predictions, training errors, and operational delays. That's why structured, well-modeled data isn't just a technical preference—it's a **strategic asset**.

## References

- Breck, E., Cai, S., Nielsen, E., Salib, M., & Sculley, D. (2017). *The hidden technical debt in machine learning systems*. In *Proceedings of the 30th International Conference on Neural Information Processing Systems* (pp. 2503–2511). NeurIPS. https://papers.nips.cc/paper_files/paper/2015/file/86df7dcfd896fcaf2674f757a2463eba-Paper.pdf
- Google Cloud. (n.d.). *BigQuery architecture overview*. Google Cloud. https://cloud.google.com/bigquery/docs/introduction
- LearnSQL.com. (2025, April 1). *Why learning SQL beats relying on AI: understanding data still matters*. LearnSQL.com. https://learnsql.com/blog/understanding-data-matters
- DZone. (2025, February 11). *SQL as the backbone of big data and AI powerhouses*. DZone. https://dzone.com/articles/sql-the-backbone-of-big-data-and-ai-powerhouses
- WhereScape. (n.d.). *What is a data model?* WhereScape. https://www.wherescape.com/blog/what-is-a-data-model/
- Alex The Analyst. (2021, January 20). *What is a data model?* [Video]. YouTube. https://www.youtube.com/watch?v=u750dq3Undo
- TechWorld with Nana. (2022, May 8). *What is SQL? Why use SQL?* [Video]. YouTube. https://www.youtube.com/watch?v=PX8qprV1txo