

# Building Scalable Data Lakes with Databricks and Snowflake: An Enterprise Approach

Santosh Pashikanti\*

**Citation:** Pashikanti S. Building Scalable Data Lakes with Databricks and Snowflake: An Enterprise Approach. *J Artif Intell Mach Learn & Data Sci* 2024, 2(2), 2064-2067. DOI: doi.org/10.51219/JAIMLD/Santosh-pashikanti/453

**Received:** 03 April, 2024; **Accepted:** 28 April, 2024; **Published:** 30 April, 2024

\*Corresponding author: Santosh Pashikanti, Independent Researcher, USA

**Copyright:** © 2024 Pashikanti S., This is an open-access article distributed under the terms of the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

## ABSTRACT

Recent advancements in cloud-based data platforms have facilitated unprecedented opportunities for enterprises to store, process and analyze massive volumes of data. Among the leading solutions in this sphere, Databricks and Snowflake offer robust, scalable and secure infrastructures for building modern data lakes and enabling advanced analytics. This white paper presents an enterprise-oriented approach to designing and implementing scalable data lakes using Databricks and Snowflake. This paper delves into deep technical architecture, discuss associated challenges, propose solutions and provide methodological guidance. Furthermore, this paper illustrate practical case studies and use cases that demonstrate the synergy between Databricks and Snowflake and outline best practices for successful implementation.

**Keywords:** Data lakes, Databricks, Snowflake, Big data architecture, Data Lake house, Cloud computing, Enterprise analytics, Scalability

## 1. Introduction

Enterprises increasingly rely on robust data-driven strategies to gain competitive advantages. The growing abundance of real-time, semi-structured and unstructured data has led to the evolution of modern data lake architectures. Traditional data warehouses have struggled to meet the needs of large-scale analytics, machine learning (ML) and advanced analytics workflows due to high complexity, limited scalability and rigid schema constraints<sup>1</sup>.

Databricks and Snowflake have emerged as two prominent cloud-based platforms designed to address these limitations:

- **Databricks** provides a unified analytics platform, incorporating Apache Spark's distributed computing capabilities and Delta Lake's reliability layer for transactional consistency and schema management<sup>2</sup>. It enables data engineering, data science, machine learning and business intelligence (BI) at scale.

- **Snowflake** offers a cloud-native data warehouse and analytics engine that supports multi-cluster shared data architecture, seamless scaling of compute and storage and robust security and governance features<sup>3</sup>.

By combining Databricks and Snowflake into a synergistic data solution, enterprises can build scalable data lakes for complex workflows. This white paper provides a holistic technical approach to conceptualizing, architecting and implementing an enterprise-level data lake environment using these platforms.

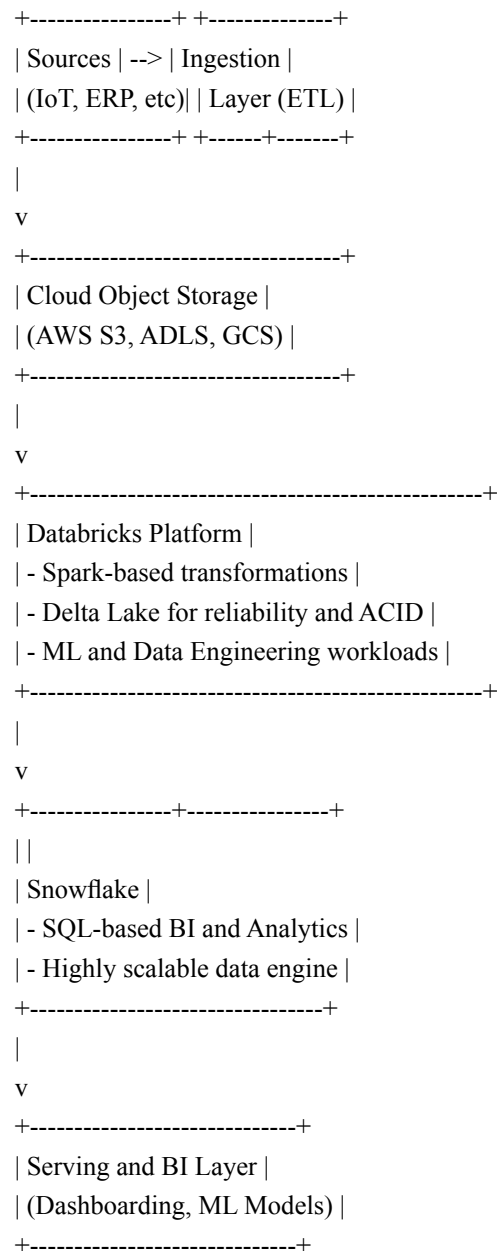
## 2. Architecture Overview

### 2.1. Architectural pillars

- **Ingestion layer:** Responsible for ingesting various data types (structured, semi-structured, unstructured) from diverse sources (IoT devices, transactional systems, third-party APIs) into cloud storage.

- **Storage layer:** Provides cheap, durable and infinitely scalable storage. Common options include AWS S3, Azure Data Lake Storage (ADLS) or Google Cloud Storage (GCS).
- **Processing/compute layer:** Utilizes Databricks (Spark-based) clusters for running transformations, batch processing, streaming and advanced analytics workloads.
- **Data warehouse and analysis layer:** Snowflake provides a cloud data warehouse engine for high-performance SQL-based analytics and enterprise BI workloads.
- **Serving layer:** Data is served to end-users, ML pipelines or external applications through APIs, BI tools or direct query access.

## 2.2. Conceptual diagram



**Figure 1:** High-level architecture for a Databricks–Snowflake data lake solution.

## 3. Detailed Technical Architecture

### 3.1. Ingestion methodologies

Data ingestion into the data lake can be performed using the

following methods:

- **Batch ingestion:** Large-volume data loads at defined intervals using Spark jobs or Snowflake’s bulk loading methods.
  - Tools such as **Databricks Auto Loader**<sup>2</sup> for incremental file-based ingestion.
  - **Snowflake COPY** commands for efficient data loading from cloud storage<sup>3</sup>.
- **Streaming ingestion:** Real-time data capture from sources generating continuous data streams (e.g., Apache Kafka, Azure Event Hubs).
  - Databricks Structured Streaming for continuous data ingestion into Delta tables.
  - Snowflake Snowpipe for near-real-time ingestion of micro-batches.
- **API-based Ingestion:** Ingestion from REST APIs or external systems.
  - Integration with cloud-based data integration services (Azure Data Factory, AWS Glue, etc.).
  - Managed or custom connectors in Databricks.

### 3.2. Storage and data lakehouse

**Delta Lake** on Databricks introduces ACID transactions, schema enforcement and time travel on top of raw data files in cloud storage<sup>2</sup>. This serves as a Lakehouse paradigm—merging the benefits of data lakes (scalable, low-cost storage) and data warehouses (transactional reliability, data governance).

#### Key considerations:

- **Partitioning and clustering:** Partition data by date or region to optimize query performance.
- **Delta transaction logs:** Maintain consistency for concurrent reads/writes.
- **Data versioning:** Allows rollback or consistent ML feature generation.

### 3.3. Processing and transformation

Databricks clusters utilize **Apache Spark** to perform transformations, augmenting raw data with analytics-friendly schemas. Workloads can be classified as:

- **Batch workloads:** Large-scale transformations and aggregations using Spark SQL, PySpark or Scala.
- **Streaming workloads:** Real-time analytics and event processing using Databricks’ Structured Streaming or Spark Streaming.
- **Machine learning:** Training sophisticated ML/DL models via Spark MLlib, TensorFlow or other frameworks.

Databricks notebooks and jobs orchestrate these tasks, integrating seamlessly with CI/CD pipelines for continuous integration and delivery<sup>2</sup>.

### 3.4. Snowflake data warehouse

Snowflake complements Databricks by providing a high-performance SQL engine for enterprise BI consumption<sup>3</sup>. The key architectural elements of Snowflake include:

- **Separation of compute and storage:** Compute clusters (virtual warehouses) scale independently of data storage, enabling cost-effective elasticity.
- **Micro-partitioning:** Data is automatically split into micro-partitions for efficient query performance without manual partitioning overhead.
- **Cloning and time travel:** Allows zero-copy clones for quick environment replication and historical data access for compliance.
- **Security and governance:** Role-based access control (RBAC), data masking policies, encryption and network security configurations for enterprise-grade compliance.

### 3.5 Integration between Databricks and Snowflake

Although Databricks can store and process data in the data lake and Snowflake can manage and analyze structured data, synergy often emerges when both are integrated:

- **Data exchange:** Use **Delta Sharing** or external tables to read/write data between Databricks and Snowflake.
- **JDBC/ODBC connectivity:** Spark can read from or write to Snowflake using connectors such as the Snowflake Spark Connector<sup>4</sup>.
- **Orchestration:** Tools like **Airflow**, **Azure Data Factory** or **Databricks Workflows** schedule end-to-end pipelines across both platforms.

## 4. Implementation Methodology

### 4.1. Step-by-step process

- **Environment setup**
  - Provision Databricks workspace and Snowflake account.
  - Configure security credentials (AWS IAM roles, Azure service principals or GCP service accounts) to ensure proper access to cloud storage.
- **Data lake creation**
  - Set up a cloud storage bucket/container as the centralized data lake.
  - Configure data lake settings for versioning, encryption and lifecycle policies.

### Ingestion pipelines

- Implement automated pipelines (batch, streaming or both) using Databricks notebooks, Spark jobs or Snowflake COPY into staging tables.
- Validate data correctness and completeness.
- **Data Transformation with Databricks**
  - Curate raw data into Silver (cleansed) and Gold (aggregated) layers.
  - Apply business logic, transformations and data quality checks.
  - Store curated data in Delta tables.
- **Snowflake Data Warehouse:**
  - Use Snowflake external tables or ingestion pipelines to load curated data into Snowflake.
  - Design star or snowflake schemas for BI consumption if needed.
  - Implement security controls (RBAC, row-level security).

### Analytics and Visualization:

- Connect BI tools (Tableau, Power BI, etc.) or ML applications to either Databricks or Snowflake for analytics.
- Use Snowflake's built-in features (e.g., Snowflake SQL) for quick ad-hoc analyses.

### Monitoring and Optimization:

- Track cluster performance, query optimization and storage usage.
- Implement autoscaling to handle peak loads and reduce cost.
- Monitor data pipeline health and data governance metrics.

## 5. Challenges and Solutions

### 5.1. Data governance and security

- **Challenge:** Large-scale data environments complicate enforcement of security policies and compliance requirements.
- **Solution:** Centralize identity and access management using RBAC in Snowflake, Azure Active Directory or AWS IAM. Enforce column-level or row-level security, data masking and encryption at rest and in transit.

### 5.2. Performance optimization

- **Challenge:** Inefficient queries and unoptimized cluster sizing can lead to high costs and slow performance.
- **Solution:**
  - For Databricks, configure autoscaling clusters, optimize Spark jobs with partition pruning and leverage Delta Lake's Z-order clustering.
  - For Snowflake, optimize micro-partition usage, set appropriate warehouse sizes for concurrency and regularly monitor query profiles.

### 5.3. Data consistency and reliability

- **Challenge:** Simultaneous reads/writes or schema drifts can compromise data correctness in a large data lake.
- **Solution:** Databricks Delta Lake's ACID transactions mitigate concurrency conflicts. Automated schema evolution features ensure consistent schema handling<sup>2</sup>.

### 5.4. Integration Complexity

- **Challenge:** Integrating multiple cloud services (Databricks, Snowflake, Kafka, etc.) can be complex and error-prone.
- **Solution:** Simplify orchestration with tools like Airflow or Databricks Workflows. Leverage official Snowflake and Databricks connectors for secure data exchange<sup>4</sup>.

### 5.5. Cost management

- **Challenge:** Running large clusters and high storage volumes can accumulate substantial costs.
- **Solution:**
  - Utilize autoscaling and spot instances in Databricks to reduce compute spend.
  - Monitor Snowflake credits consumed by queries and scale down or pause warehouses during low demand.
  - Optimize data retention policies and tiered storage.

## 6. Case Studies and Use Cases

### 6.1. Case study: global retail company

A global retail company needed to unify data from point-of-sale (POS) systems, e-commerce sites and supply chain logs. By creating a data lake on **Azure Data Lake Storage** and processing data through **Databricks** (batch for historical data and streaming for real-time orders), the company curated analytics-ready data in Delta tables. Then, curated data was loaded into **Snowflake** for detailed sales analytics and enterprise BI<sup>2</sup>. The result was a 40% reduction in data processing time and near real-time insights for supply chain optimization.

### 6.2. Use case: Financial services compliance

A financial institution leveraged **Databricks** to ingest large volumes of transaction logs in real time, applying advanced ML-based anomaly detection. **Snowflake** served as a secure data warehouse for compliance reporting, allowing auditing teams to run complex SQL queries on an as-needed basis. The integrated approach provided transparent data lineage, robust governance and agile analytics, crucial for meeting stringent regulatory requirements<sup>5</sup>.

### 6.3 Use Case: Healthcare Analytics

Healthcare providers handling claims data used **Databricks** for large-scale data processing-cleansing medical records and generating feature sets for predictive modeling. **Snowflake** acted as the central query engine for clinical dashboards and patient population analytics, benefiting from minimal overhead, auto-scaling and HIPAA-compliant security measures<sup>6</sup>.

## 7. Conclusion

Building an enterprise-scale data lake requires addressing data ingestion, storage, processing, governance and analytics, all while remaining cost-effective and secure. Databricks and Snowflake offer complementary capabilities: Databricks excels in large-scale data processing and machine learning with Spark and Delta Lake, whereas Snowflake delivers robust SQL-based analytics with a cloud-native data warehouse engine. By combining these platforms organizations can achieve a versatile and scalable data architecture that accommodates diverse workloads, from real-time streaming to advanced analytics and BI reporting.

The methodology outlined in this white paper-covering ingestion, transformation, data lake house creation and integration-provides a starting blueprint for enterprises aiming to harness the full power of their data. As organizations continue to expand data-driven initiatives, implementing best practices in data governance, security, cost management and performance optimization will remain critical for success.

## 8. References

1. <https://www.oreilly.com/library/view/hadoop-the-definitive/9781491901687/>
2. <https://docs.databricks.com/>
3. <https://docs.snowflake.com/>
4. <https://docs.snowflake.com/en/user-guide/spark-connector>
5. <https://www.databricks.com/solutions/industries/financial-services>
6. <https://www.databricks.com/solutions/industries/healthcare-and-life-sciences>