

Comprehensive Reference Architectures for Multi-Cloud Data Lakes

Based on our discussion of cost-effective multi-cloud data lakes with Apache Iceberg, Trino, and federated search, here are the most relevant reference architectures and blueprints available:

Official Cloud Provider Reference Architectures

1. AWS Prescriptive Guidance - Apache Iceberg on AWS

Link: [AWS Prescriptive Guidance PDF](#)

Key Components Covered:

- Modern data lake architecture with Iceberg
- Integration with AWS Glue, Athena, EMR
- Multi-engine compatibility (Spark, Trino, Flink)
- Performance optimization patterns
- Cost management strategies

Architecture Highlights:



2. Starburst "Icehouse" Reference Architecture

Link: [Starburst Icehouse Architecture](#)

Focus: Trino + Iceberg + Multi-Cloud Federation **Key Features:**

- Cross-cloud query federation
- Unified metadata management
- Cost optimization through query pushdown
- Performance tuning guidelines

Architecture Pattern:

Multi-Cloud Data Sources → Trino Coordinators → Federated Queries

↓ ↓ ↓

Iceberg Tables Metadata Sync Unified Results



Modern Data Stack Reference Architectures

3. Databricks Lakehouse Reference Architecture

Link: [Databricks Reference Architectures](#)

Downloadable PDFs Available:

- AWS Reference Architecture (11x17 format)
- Azure Reference Architecture
- GCP Reference Architecture

Key Patterns:

- Medallion Architecture (Bronze/Silver/Gold)
- Multi-cloud federation capabilities
- Unity Catalog for governance
- ML/AI integration patterns

4. MinIO Modern Data Lake Reference Architecture

Link: [MinIO Data Lake Architecture Guide](#)

Architecture Layers:

Consumption Layer... → BI Tools, ML Platforms, APIs

Processing Layer... → Spark, Flink, Trino

Storage Layer... → Object Storage (S3/MinIO) + Iceberg

Metadata Layer... → Catalogs (Hive/Glue/Nessie)

Ingestion Layer... → Streaming + Batch Pipelines

Strengths:

- Cloud-agnostic design
- Cost optimization focus
- Open source technology stack
- Detailed implementation guidance

Practical Implementation Blueprints

5. Kubernetes-Native Lakehouse (Bionic-GPT)

Link: [K8s Lakehouse Tutorial](#)

Technology Stack:

- **Storage:** MinIO (S3-compatible)
- **Catalog:** Nessie (Git-like versioning)
- **Table Format:** Apache Iceberg
- **Query Engine:** Trino
- **Deployment:** Kubernetes manifests

Complete Implementation:

yaml

Provides actual K8s YAML files for:

- MinIO deployment
- Nessie catalog service
- Trino cluster configuration
- Service networking
- Storage persistence

6. Natural Intelligence Production Migration

Link: [AWS Blog - NI Iceberg Migration](#)

Real-World Patterns:

- Legacy Hive → Iceberg migration
- Multi-engine support (Snowflake + Athena + Druid)
- Medallion architecture implementation
- Production lessons learned

7. Cloudinary Petabyte-Scale Implementation

Link: [AWS Blog - Cloudinary Case Study](#)

Scale Insights:

- 20 billion requests daily

- Streaming + batch processing
- Multi-tool integration
- Cost optimization at scale

Open Source Reference Implementations

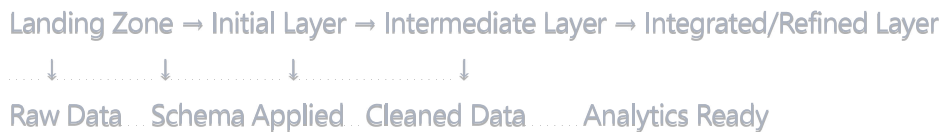
8. Medallion Architecture on GitHub

Link: [GitHub - Medallion Architecture](#)

Available Implementations:

- AWS + Databricks reference architecture
- Complete reference implementation code
- Multi-cloud deployment patterns
- Data governance frameworks

Architecture Layers:



9. Multi-Cloud Federated Architecture (Medium)

Link: [Trino + Iceberg Alternative to Redshift](#)

Key Patterns:

- Cost comparison vs traditional warehouses
- Multi-cloud deployment strategies
- Performance optimization techniques
- Migration best practices

Implementation-Ready Architectures

10. AWS + Snowflake Interoperability

Link: [AWS Blog - Iceberg with Snowflake](#)

Two Architectural Patterns:

Pattern A: AWS-Managed Iceberg

Data Sources → AWS Glue ETL → S3 Iceberg Tables → AWS Glue Catalog



..... Snowflake External Tables

Pattern B: Snowflake-Managed Iceberg

Data Sources → Snowflake Pipelines → S3 Iceberg Tables → Snowflake Catalog



..... AWS Services Access

11. Unified Lakehouse with Spark + Trino

Link: [Medium - One Table, Two Engines](#)

Implementation Guide:

- Shared Hive Metastore configuration
- Dual-engine access patterns
- Performance optimization
- Operational best practices



Architecture Selection Guide

For Startups/Small Teams:

- **Recommended:** Kubernetes Lakehouse (Bionic-GPT)
- **Why:** Simple deployment, open source, cost-effective
- **Technologies:** MinIO + Nessie + Trino + Iceberg

For AWS-Centric Organizations:

- **Recommended:** AWS Prescriptive Guidance Architecture
- **Why:** Native integrations, managed services
- **Technologies:** S3 + Glue + Athena/EMR + Iceberg

For Multi-Cloud Enterprises:

- **Recommended:** Starburst Icehouse + MinIO Architecture
- **Why:** True cloud independence, federation capabilities

- **Technologies:** Trino + Iceberg + Object Storage + Polaris

For Large-Scale Production:

- **Recommended:** Databricks Lakehouse + Medallion Pattern
- **Why:** Proven at scale, comprehensive tooling
- **Technologies:** Unity Catalog + Delta/Iceberg + Multi-cloud

Ready-to-Deploy Resources

Infrastructure as Code:

1. **Terraform Modules:** Most reference architectures include Terraform
2. **Kubernetes Manifests:** Complete YAML configurations available
3. **CloudFormation Templates:** AWS-specific deployments
4. **Helm Charts:** For K8s-based deployments

Configuration Templates:

1. **Trino Catalogs:** Pre-configured for Iceberg + multiple clouds
2. **Spark Configurations:** Optimized for Iceberg operations
3. **Data Pipeline Examples:** Airflow DAGs and Spark jobs
4. **Monitoring Setups:** Grafana dashboards and alerting

Sample Data and Queries:

1. **Demo Datasets:** TPC-H and TPC-DS benchmarks
2. **Performance Tests:** Query optimization examples
3. **Migration Scripts:** Legacy format to Iceberg conversion
4. **BI Tool Connections:** Tableau, Superset, Power BI examples

Key Takeaways

Start Simple: Begin with the Kubernetes lakehouse for learning and proof-of-concepts **Scale Gradually:** Move to cloud-managed services as requirements grow **Stay Open:** Use Iceberg and Trino to avoid vendor lock-in **Monitor Costs:** Implement the federated patterns to minimize data movement **Plan Governance:** Use the medallion architecture for data quality and lineage

These reference architectures provide everything from high-level blueprints to detailed implementation code, covering the exact technologies and patterns we discussed for cost-effective, multi-cloud data lakes with federated search capabilities.

