

# Serverless SQL: Querying Cloud Data Lakes with On-Demand Engines

**Authors:** <sup>1</sup> Areej Mustafa, <sup>2</sup> Zillay Huma

**Corresponding Author:** [areejmustafa703@gmail.com](mailto:areejmustafa703@gmail.com)

## Abstract

Serverless SQL engines are revolutionizing the way organizations query massive datasets stored in cloud data lakes by eliminating the need for persistent infrastructure and allowing on-demand, scalable data access. This paradigm enables users to run SQL queries directly on cloud-based storage systems like Amazon S3, Azure Data Lake, or Google Cloud Storage without provisioning servers or managing clusters. As businesses increasingly adopt cloud-native architectures, serverless SQL offers a compelling alternative to traditional data warehouses, particularly for exploratory analytics, ad-hoc reporting, and cost-efficient big data processing. This paper explores the principles, architecture, and use cases of serverless SQL engines, with a focus on their integration with data lake storage. It further examines performance considerations, data governance, cost models, and security strategies. By analyzing key technologies such as AWS Athena, Google BigQuery, and Azure Synapse Serverless, this paper highlights the opportunities and challenges associated with this modern approach to querying data at scale.

**Keywords:** Serverless SQL, data lakes, cloud analytics, on-demand query engines, AWS Athena, Google BigQuery, Azure Synapse Serverless, cost optimization, data governance, SQL-as-a-service

## Introduction

As data volumes grow exponentially, traditional data processing paradigms are increasingly challenged by scalability, cost, and complexity[1]. Cloud computing has enabled organizations to offload infrastructure concerns, shifting the focus toward elasticity and usage-based pricing.

<sup>1</sup> University of Gujrat, Pakistan

<sup>2</sup> University of Gujrat, Pakistan

In this context, data lakes have emerged as scalable repositories for raw and semi-structured data, offering flexibility over rigid schemas found in data warehouses[2]. However, effectively extracting insights from these vast stores has necessitated innovation in querying technologies. Serverless SQL is one such innovation that redefines the way analytics is conducted in the cloud era[3].

Unlike conventional data warehouses that require provisioning of compute clusters and management of storage, serverless SQL engines allow users to execute queries directly against data stored in cloud-native formats such as Parquet, ORC, JSON, or CSV[4]. This is achieved without persistent servers, reducing operational overhead and allowing for pay-per-query or pay-per-scan models. This mode of operation is particularly well-suited to environments with variable workloads, exploratory data analysis needs, and unpredictable query patterns[5].

The architecture of serverless SQL engines is fundamentally different from that of traditional relational databases. It typically involves a distributed query execution layer that spins up automatically in response to a user's query[6]. The query engine fetches metadata, plans the execution path, parallelizes the read operation across distributed storage systems, and returns the result to the user. Examples include Amazon Athena, which operates over Amazon S3; Google BigQuery, which abstracts compute from storage and charges per query processed; and Azure Synapse Serverless, which offers on-demand querying over Azure Data Lake Storage[7].

The serverless model brings a host of benefits. It democratizes access to big data by allowing data scientists, analysts, and even business users to use familiar SQL syntax on large datasets without having to understand the underlying infrastructure[8]. It offers financial efficiency by charging only for what is used, eliminating idle compute time. Furthermore, it supports a decoupled architecture, enabling independent scaling of compute and storage[9].

However, serverless SQL is not without challenges. Performance optimization is a key concern, particularly since users pay based on data scanned. Poorly optimized queries, unpartitioned datasets, and inefficient data formats can significantly inflate costs and query times[10]. Data governance also becomes complex when multiple teams and users access shared data lakes with

varying levels of sensitivity. Ensuring compliance, security, and lineage in such an environment requires robust integration with identity management and access control tools[11].

Additionally, serverless SQL engines often come with limitations around advanced SQL capabilities, concurrency, and support for procedural logic. These limitations must be considered when designing data architectures that rely heavily on SQL-based transformations or require transactional consistency[12].

Despite these concerns, the value proposition of serverless SQL is compelling. As enterprises move towards real-time insights and operational intelligence, the ability to query large volumes of semi-structured data without standing infrastructure represents a significant advancement[10]. Serverless SQL is especially powerful in hybrid architectures where structured data from transactional systems is combined with unstructured or semi-structured data in data lakes to derive richer insights[13].

The remainder of this paper explores two dimensions of this emerging paradigm. First, it analyzes the architectural and functional characteristics of serverless SQL engines, including integration with cloud data lakes, execution models, and optimization techniques[14]. Second, it examines practical use cases, operational considerations, and best practices that organizations can adopt to harness the power of serverless SQL while maintaining cost efficiency and data governance[15].

## **Architectural Foundations and Optimization in Serverless SQL Engines:**

Serverless SQL engines are designed around the principle of ephemeral compute. Rather than relying on continuously running servers or clusters, they instantiate query processing resources on demand. This architectural approach enables a high degree of elasticity and cost efficiency. Typically, these engines consist of three major components: the query interface (often SQL-based), a distributed execution engine, and an integration layer with object storage systems where the data resides[16].

Amazon Athena, for instance, relies on Presto (now Trino) as its query execution engine and integrates directly with Amazon S3. When a query is submitted, Athena retrieves metadata from the AWS Glue Data Catalog, parses the SQL, and launches parallel workers to read data from S3[17]. Similarly, Google BigQuery separates compute and storage by maintaining a multi-tenant execution fabric that processes SQL jobs against data stored in Colossus, Google's distributed file system. Azure Synapse Serverless uses the T-SQL language to query data in Azure Data Lake Storage, relying on a distributed SQL engine that spins up for each query session[18].

The optimization of queries in serverless SQL environments is fundamentally linked to data layout and storage format. Columnar storage formats like Parquet and ORC are highly recommended because they enable selective reading of columns, reducing the amount of data scanned and improving query performance. These formats also support compression and predicate pushdown, which further reduce resource consumption[19].

Partitioning is another crucial technique. Data in cloud data lakes should be partitioned on frequently queried columns such as date or region. Query engines can then prune irrelevant partitions during query planning, reducing the scan footprint. However, excessive partitioning can lead to metadata overhead and should be balanced carefully[20].

Another consideration is metadata management. Serverless engines rely heavily on metadata catalogs to interpret the schema of external data. Tools such as AWS Glue, Google Data Catalog, and Azure Purview provide schema registries, data classification, and lineage tracking. Accurate and up-to-date metadata ensures efficient query planning and minimizes runtime errors[21].

Cost control is a critical aspect of operating in a serverless SQL model. Since most engines charge by data scanned or processed, query efficiency directly impacts cost. Organizations can implement guardrails such as query cost estimators, user quotas, or even automated query rewriting to limit inefficient operations. Additionally, policies can be configured to restrict full-table scans or to limit access to high-cost datasets[22].

Security is implemented through a combination of IAM (Identity and Access Management), data encryption, and audit logging. Serverless SQL engines often integrate with native cloud security tools, enabling fine-grained access control down to row and column levels. Data at rest is encrypted using storage-level encryption, while query logs are often captured in systems like AWS CloudTrail or Azure Monitor for auditing purposes[23].

Concurrency and latency are operational factors that vary across platforms. While BigQuery is designed for massively parallel query execution with high concurrency, Athena and Synapse Serverless may face throttling under heavy loads. Organizations with mission-critical SLAs should benchmark platforms under peak conditions and consider hybrid setups where hot data is cached in fast-access stores[24].

Despite these intricacies, serverless SQL engines are constantly evolving, with enhancements in caching, federated query capabilities, and machine learning integration. For instance, BigQuery ML allows users to train and apply machine learning models using SQL directly on their datasets. This broadens the scope of serverless SQL beyond simple analytics, making it a central tool in the modern data stack[25].

## **Operational Use Cases, Governance, and Best Practices:**

The versatility of serverless SQL engines makes them suitable for a wide array of operational use cases across industries. One of the most prominent applications is in ad-hoc analytics and data exploration. Data scientists and analysts can query vast data lakes without having to wait for infrastructure provisioning, accelerating the pace of insight generation. For example, marketing teams can analyze web traffic logs stored in Parquet format on S3 using Athena to identify customer behavior patterns in real-time[26].

Another major use case is in federated querying. Serverless SQL engines can access multiple data sources across cloud-native and external systems. Google BigQuery's federated queries allow integration with data in Google Sheets, Cloud SQL, and even third-party databases like

PostgreSQL or MySQL. This flexibility supports complex business reporting scenarios where information resides in silos[25].

Serverless SQL is also widely used in ETL and data transformation pipelines. While traditional ETL tools extract and load data into staging areas, serverless engines can process and filter raw data directly within the data lake. This approach minimizes data movement and ensures scalability. Scheduled queries can be orchestrated through workflow tools like AWS Step Functions or Google Cloud Composer, enabling fully automated pipelines[27].

In regulated industries such as finance and healthcare, data governance is essential. Serverless SQL engines must integrate with centralized governance frameworks to ensure data compliance and accountability. For example, access to financial records stored in a data lake can be controlled using AWS Lake Formation, which provides attribute-based access control and audit logs. Similarly, tagging sensitive columns in the data catalog can enforce masking or obfuscation when accessed via SQL[28].

Monitoring and observability are equally important. Cloud-native platforms provide detailed query execution logs, performance metrics, and error reports. These logs can be ingested into observability platforms like Datadog, Splunk, or native services such as AWS CloudWatch and Google Cloud Logging. By monitoring query patterns, administrators can identify performance bottlenecks, unused datasets, or potentially risky user behavior[29].

Best practices for serverless SQL adoption include educating users on query efficiency, especially in cost-sensitive environments. Query optimization guidelines, such as using filters early in the WHERE clause, selecting only necessary columns, and avoiding wildcard scans, should be part of developer onboarding. Usage dashboards and real-time cost monitors can also promote accountability and encourage efficient usage[30].

Another best practice is to maintain versioned schemas and data contracts. This ensures that downstream applications or queries do not break when data formats change. Schema evolution is supported by most data lake formats, but requires careful coordination when used in production systems[31].

To support multi-tenant use cases, organizations can employ strategies such as logical data separation using tenant IDs, dynamic query rewriting, or even isolated metadata catalogs per business unit. These strategies help maintain performance isolation and prevent data leakage[32].

Finally, the future of serverless SQL is poised to intersect with artificial intelligence. With the rise of intelligent query assistants and natural language interfaces, serverless SQL engines could serve as backends to conversational analytics systems, expanding access to non-technical users. Integration with real-time stream processing and support for advanced analytics will further blur the line between batch and interactive workloads[33].

## Conclusion

In conclusion, serverless SQL offers a robust, flexible, and cost-effective framework for querying cloud data lakes. With the right architectural design, governance practices, and user education, it can unlock the full potential of cloud-native data analytics. Serverless SQL has emerged as a transformative approach for querying data in cloud environments, combining scalability, cost efficiency, and ease of access. By decoupling compute from storage and leveraging on-demand execution, it empowers organizations to perform large-scale analytics without the burden of infrastructure management. While challenges remain in areas such as query optimization and data governance, the continued evolution of serverless technologies and best practices positions serverless SQL as a cornerstone of modern data architectures.

## References:

- [1] A. S. Shethiya, "Rise of LLM-Driven Systems: Architecting Adaptive Software with Generative AI," *Spectrum of Research*, vol. 3, no. 2, 2023.
- [2] I. Salehin *et al.*, "AutoML: A systematic review on automated machine learning with neural architecture search," *Journal of Information and Intelligence*, vol. 2, no. 1, pp. 52-81, 2024.
- [3] A. S. Shethiya, "Redefining Software Architecture: Challenges and Strategies for Integrating Generative AI and LLMs," *Spectrum of Research*, vol. 3, no. 1, 2023.
- [4] M. Noman, "Safe Efficient Sustainable Infrastructure in Built Environment," 2023.
- [5] A. S. Shethiya, "Next-Gen Cloud Optimization: Unifying Serverless, Microservices, and Edge Paradigms for Performance and Scalability," *Academia Nexus Journal*, vol. 2, no. 3, 2023.
- [6] M. Noman, "Precision Pricing: Harnessing AI for Electronic Shelf Labels," 2023.



- 
- [7] A. S. Shethiya, "Machine Learning in Motion: Real-World Implementations and Future Possibilities," *Academia Nexus Journal*, vol. 2, no. 2, 2023.
  - [8] M. Noman, "Potential Research Challenges in the Area of Plethysmography and Deep Learning," 2023.
  - [9] A. S. Shethiya, "LLM-Powered Architectures: Designing the Next Generation of Intelligent Software Systems," *Academia Nexus Journal*, vol. 2, no. 1, 2023.
  - [10] M. Noman, "Machine Learning at the Shelf Edge Advancing Retail with Electronic Labels," 2023.
  - [11] A. S. Shethiya, "Learning to Learn: Advancements and Challenges in Modern Machine Learning Systems," *Annals of Applied Sciences*, vol. 4, no. 1, 2023.
  - [12] A. S. Shethiya, "Scalability and Performance Optimization in Web Application Development," *Integrated Journal of Science and Technology*, vol. 2, no. 1, 2025.
  - [13] A. S. Shethiya, "Load Balancing and Database Sharding Strategies in SQL Server for Large-Scale Web Applications," *Journal of Selected Topics in Academic Research*, vol. 1, no. 1, 2025.
  - [14] N. Mazher and I. Ashraf, "A Systematic Mapping Study on Cloud Computing Security," *International Journal of Computer Applications*, vol. 89, no. 16, pp. 6-9, 2014.
  - [15] A. S. Shethiya, "Deploying AI Models in .NET Web Applications Using Azure Kubernetes Service (AKS)," *Spectrum of Research*, vol. 5, no. 1, 2025.
  - [16] A. S. Shethiya, "Building Scalable and Secure Web Applications Using .NET and Microservices," *Academia Nexus Journal*, vol. 4, no. 1, 2025.
  - [17] N. Mazher, I. Ashraf, and A. Altaf, "Which web browser work best for detecting phishing," in *2013 5th International Conference on Information and Communication Technologies*, 2013: IEEE, pp. 1-5.
  - [18] A. S. Shethiya, "AI-Assisted Code Generation and Optimization in .NET Web Development," *Annals of Applied Sciences*, vol. 6, no. 1, 2025.
  - [19] N. Mazher and I. Ashraf, "A Survey on data security models in cloud computing," *International Journal of Engineering Research and Applications (IJERA)*, vol. 3, no. 6, pp. 413-417, 2013.
  - [20] A. S. Shethiya, "Adaptive Learning Machines: A Framework for Dynamic and Real-Time ML Applications," *Annals of Applied Sciences*, vol. 5, no. 1, 2024.
  - [21] A. S. Shethiya, "AI-Enhanced Biometric Authentication: Improving Network Security with Deep Learning," *Academia Nexus Journal*, vol. 3, no. 1, 2024.
  - [22] A. S. Shethiya, "Architecting Intelligent Systems: Opportunities and Challenges of Generative AI and LLM Integration," *Academia Nexus Journal*, vol. 3, no. 2, 2024.
  - [23] A. S. Shethiya, "Decoding Intelligence: A Comprehensive Study on Machine Learning Algorithms and Applications," *Academia Nexus Journal*, vol. 3, no. 3, 2024.
  - [24] A. S. Shethiya, "Engineering with Intelligence: How Generative AI and LLMs Are Shaping the Next Era of Software Systems," *Spectrum of Research*, vol. 4, no. 1, 2024.
  - [25] I. Ashraf and N. Mazher, "An Approach to Implement Matchmaking in Condor-G," in *International Conference on Information and Communication Technology Trends*, 2013, pp. 200-202.
  - [26] A. S. Shethiya, "Ensuring Optimal Performance in Secure Multi-Tenant Cloud Deployments," *Spectrum of Research*, vol. 4, no. 2, 2024.
  - [27] A. S. Shethiya, "From Code to Cognition: Engineering Software Systems with Generative AI and Large Language Models," *Integrated Journal of Science and Technology*, vol. 1, no. 4, 2024.
  - [28] A. Nishat, "Towards Next-Generation Supercomputing: A Reconfigurable Architecture Leveraging Wireless Networks," 2020.



- [29] A. S. Shethiya, "Smarter Systems: Applying Machine Learning to Complex, Real-Time Problem Solving," *Integrated Journal of Science and Technology*, vol. 1, no. 1, 2024.
- [30] A. Nishat, "Future-Proof Supercomputing with RAW: A Wireless Reconfigurable Architecture for Scalability and Performance," 2022.
- [31] A. Nishat, "The Role of IoT in Building Smarter Cities and Sustainable Infrastructure," *International Journal of Digital Innovation*, vol. 3, no. 1, 2022.
- [32] A. Nishat, "AI Meets Transfer Pricing: Navigating Compliance, Efficiency, and Ethical Concerns," *Aitoz Multidisciplinary Review*, vol. 2, no. 1, pp. 51-56, 2023.
- [33] A. Nishat, "Artificial Intelligence in Transfer Pricing: Unlocking Opportunities for Tax Authorities and Multinational Enterprises," *Aitoz Multidisciplinary Review*, vol. 2, no. 1, pp. 32-37, 2023.