



Design cybernetics; self-adaptive systems; socio-cognitive design; feedback loops; systems design methodology; adaptive product development; management cybernetics

Aarav R. Thompson

Global Institute of Data Systems

ABSTRACT

Background: The rapid shift from on-premises, traditional data warehousing to cloud-native platforms has introduced a complex interplay between architectural paradigms, operational elasticity, and cost models. Cloud data warehouses such as Snowflake and Amazon Redshift represent two dominant approaches—one emphasizing decoupled storage and compute with fully managed elasticity, and the other offering a blend of tightly integrated cluster-based compute with cloud-native extensions. Comparative evaluation of performance and cost across realistic decision-support workloads remains critical for practitioners and researchers seeking to align technical design with business constraints. This study synthesizes theoretical foundations, empirical tuning methods, and benchmark-driven evaluation strategies into an integrated framework for adaptive performance and cost optimization. (Gupta & Sharma, 2021; Malik & Bhatia, 2022; Rao & Karim, 2018)

Objective: The work aims to (1) present a structured, reproducible methodology for comparing Snowflake and Redshift across typical analytical workloads; (2) explicate system-level tuning levers, query optimization practices, and storage/caching strategies; and (3) provide a layered decision framework that maps workload characteristics to optimal configuration patterns and cost-control techniques. (Patil & Zade, 2023; Li & Martin, 2020; Tanaka & Foster, 2021)

Methods: We synthesize insights from benchmark literature (TPC-DS), vendor documentation, and peer-reviewed studies to propose a five-layer analytical framework: Workload Characterization, Storage & Compression, Compute Resource Management, Query & Index Optimization, and Cost Governance. Each layer is described with prescriptive tactics and expected trade-offs. The framework incorporates query profiling, adaptive resource scaling strategies, materialization and caching policies, and schema/encoding choices, grounded in both Snowflake and Redshift operational models. (TPC Council, 2023; Snowflake, 2024; AWS, 2024)

Results: Descriptive analysis highlights that decoupled storage-compute architectures—epitomized by Snowflake—tend to offer superior elasticity and concurrency handling for bursty, multi-tenant analytic workloads, while cluster-based Redshift with Spectrum and concurrency scaling can achieve cost advantages under predictable, sustained throughput when tuned carefully. Storage encodings and micro-partitioning choices materially affect I/O and query latency across both platforms. Query profiling and predicate-pushdown strategies consistently yield order-of-magnitude improvements for complex star-schema TPC-DS style queries when combined with tailored sort and distribution strategies (Redshift) or clustering/caching (Snowflake). (Malik & Bhatia, 2022; Rao & Karim, 2018; Tanaka & Foster, 2021; Narayanan & Chu, 2020)

Conclusions: There is no singular "best" platform; rather, workload signatures and organizational priorities determine optimal choices. The proposed five-layer framework enables systematic selection and tuning, balancing performance and cost across common analytical scenarios. Adoption of this framework supports transparency in architectural trade-offs and produces predictable operational outcomes when applied with rigorous profiling and continuous cost monitoring. Future work should empirically validate the framework across diverse industry datasets and explore automated tuning systems that map telemetry to configuration recommendations. (Singh &

Das, 2021; Malviya, 2025)

KEYWORDS

Cloud data warehousing, Snowflake, Amazon Redshift, query optimization, cost governance, TPC-DS, workload tuning

INTRODUCTION

The proliferation of cloud platforms for data warehousing has fundamentally altered how organizations design extract-transform-load (ETL) pipelines, model analytical schemas, and budget for large-scale decision-support systems. Historically, traditional data warehouses followed tightly coupled architectures—monolithic appliances with fixed compute and storage resources—requiring careful capacity planning and frequent over-provisioning to handle peak loads (Gupta & Sharma, 2021). With the advent of cloud-native services, providers introduced new paradigms: separation of storage from compute, elastic provision of resources, serverless abstractions, and usage-based billing models. These shifts have produced substantial opportunities for cost optimization and operational agility, yet they have also amplified complexity: the same workload may behave differently under different platform-specific optimization levers and cost schemes, and practitioners often lack rigorous methods to match workload characteristics to the most economical configurations (Rao & Karim, 2018; Malik & Bhatia, 2022).

Cloud data warehouses such as Snowflake and Amazon Redshift represent distinct design choices within the cloud-native continuum. Snowflake popularized an architecture with a cloud-resident, centrally managed storage layer and multiple independent compute warehouses that can scale independently; its micro-partitioning, automatic clustering advisory, and time-travel features are often touted for simplifying administration (Snowflake, 2024; Tanaka & Foster, 2021). Amazon Redshift, built on a cluster-based model but increasingly cloud-native with features like Redshift Spectrum and concurrency scaling, provides strong performance for large-scale, sustained analytic throughput when table distribution keys, sort keys, and compression encodings are meticulously chosen (AWS, 2024; AWS Redshift Compression Guide, 2024). The differences in operational models drive divergent best practices for tuning and cost control; understanding these is essential for architects, data engineers, and business stakeholders responsible for analytics SLAs and budgets (Patil & Zade, 2023; Narayanan & Chu, 2020).

This paper addresses a central, pragmatic problem: how can organizations consistently select, tune, and govern cloud data warehouse deployments to achieve required performance objectives while minimizing total cost of ownership? Prior work has examined specific optimization strategies—indexing techniques in relational databases (James & Ortega, 2019), caching and storage efficiency in Snowflake (Tanaka & Foster, 2021), or the economics of cloud data warehouses more broadly (Rao & Karim, 2018)—but there is a gap in holistic frameworks that integrate workload characterization, platform-specific tuning levers, benchmark-driven evaluation, and cost governance mechanisms into a single prescriptive methodology. Additionally, contemporary studies that contrast Snowflake and Redshift are often limited in scope (single benchmark, narrow workload types) or do not fully articulate the mapping from query telemetry to configuration changes (Patil & Zade, 2023; Malik & Bhatia, 2022).

We present a comprehensive, layered framework that synthesizes the existing literature, vendor guidance, and benchmark norms (TPC-DS) to produce actionable, theoretically grounded recommendations for platform selection, tuning, and operational governance. The framework is designed to be vendor-agnostic at the higher layers—centering on workload classification and cost policy—and vendor-specific at the tuning layer where Redshift and Snowflake diverge sharply. This work aims to provide both researchers and practitioners with a structured approach to navigating the complex trade-offs inherent in cloud data warehousing, particularly in the context of contemporary, decision-support workloads common to insurance, finance, retail, and manufacturing.

METHODOLOGY

This study constructs a theoretical and prescriptive framework through systematic synthesis of peer-reviewed literature, vendor technical documentation, and established industry benchmarks. The methodological approach comprises three interlocking components: a conceptual framework design, mapping of platform-specific tuning levers, and a benchmark-aligned evaluation methodology for comparative analysis.

Conceptual Framework Design

We adopt a layered decomposition to manage the multi-dimensional optimization problem. The decomposition supports modular reasoning and traceable mapping from workload characteristics to tuning actions. The five layers are:

1. **Workload Characterization:** Defining workload signatures using dimensions such as query complexity (number of joins, aggregations), concurrency profile (peak concurrent users/queries), data freshness (batch vs. streaming), data distribution (skew, cardinalities), and access patterns (point lookups vs. full table scans). This layer emphasizes the necessity of instrumenting query telemetry and ETL cadence observations to classify workload types quantitatively. Foundational work on query profiling and characterization informs this layer (Lee & Martin, 2020; TPC Council, 2023).
2. **Storage & Compression:** Strategies for physical data layout, encodings, micro-partitioning/clustering, and cold/warm data segregation. This layer synthesizes guidance on compression encodings and storage types, acknowledging that encodings reduce I/O and storage cost but require careful selection based on column cardinality and query predicates (AWS Redshift Compression Guide, 2024; Tanaka & Foster, 2021).
3. **Compute Resource Management:** Elasticity, concurrency control, and autoscaling strategies. The layer contrasts Snowflake's independent warehouses and auto-suspend/auto-resume behaviors with Redshift's cluster management, elastic resize, and concurrency scaling. The focus is on matching compute supply with workload demand to avoid both performance degradation and cost overrun (Singh & Das, 2021; AWS Elastic Resize Documentation, 2023; Snowflake Warehouse Management, 2024).
4. **Query & Index Optimization:** Query rewrite rules, materialized views, result set caching, distribution/sort key design (Redshift), and clustering keys/micro-partitions (Snowflake). Here we apply canonical indexing and optimization literature to cloud contexts and specify platform-specific tactics (James & Ortega, 2019; Narayanan & Chu, 2020).
5. **Cost Governance and Policy:** Billing-aware policies, reservation vs. on-demand decisions, workload prioritization, resource tagging for chargeback, and automated meta-policies for downscaling idle resources. This layer draws on cost-analysis literature and the latest platform billing models to provide governance tactics and monitoring requirements (Rao & Karim, 2018; Malviya, 2025).

Mapping Platform-Specific Tuning Levers

Within the layered design, we enumerate specific operational levers and articulate expected impacts under different workload signatures. Examples include:

- **Snowflake:** warehouse sizing policies, multi-cluster warehouses for concurrency, clustering keys, materialized views, result cache utilization, and time-travel retention policy adjustments (Snowflake, 2024; Tanaka & Foster, 2021).

- Redshift: choosing sort and distribution keys, column encodings, spectrum usage to offload colder data to S3, concurrency scaling, elastic resize scheduling, and analyzing query plans with the Redshift Query Analyzer (AWS, 2024; AWS Redshift Compression Guide, 2024; AWS Concurrency Scaling, 2023).

For each lever, we specify triggers (telemetry signals) that should prompt action, e.g., rising average queue time, frequent spill-to-disk events, high micro-partition pruning misses, or disproportionate egress charges.

Benchmark-Aligned Evaluation Methodology

Although this article synthesizes existing findings rather than presenting new experimental data, it aligns recommendations with the TPC-DS benchmark and real-world workload analogues to ensure practical relevance. The TPC benchmark provides canonical workload templates, dataset scales, and query mixes that are widely recognized for decision-support systems; aligning with it enables comparative clarity (TPC Council, 2023). We describe how to translate TPC-DS query types into workload signatures and illustrate how each layer's levers would be applied to those signatures. This alignment also defines the evaluation metrics practitioners should monitor—query latency distributions, 95th percentile latency, cost per query, concurrency-limited throughput, and storage cost per terabyte.

Synthesis Process and Citation Discipline

To ensure each claim is traceable, we systematically cite the relevant literature for theoretical claims (e.g., query profiling techniques from Lee & Martin, 2020), vendor documentation for operational behavior (Snowflake, 2024; AWS, 2024), and economic analyses for cost-impact claims (Rao & Karim, 2018; Malviya, 2025). The methodology explicitly avoids new performance testing in order to preserve reproducibility and to focus on producing a broadly applicable prescriptive framework based on the best available secondary sources.

RESULTS

The results section provides a descriptive analysis of framework application across canonical workload signatures and articulates expected outcomes from recommended tunings. Because this work synthesizes literature and vendor guidance rather than reporting primary experimental data, results are framed as theoretically and empirically grounded expectations and trade-offs.

Workload Characterization Outcomes

Workload classification yields discrete archetypes, each requiring different optimization emphases:

- Ad-hoc, high-concurrency BI workloads: Characterized by many concurrent, medium-complexity queries, often with user-driven ad-hoc filters and frequent small aggregations. For such workloads, Snowflake's multi-cluster warehouses and result caching are especially beneficial because they decouple user concurrency from compute contention and maintain query-level result caches that frequently satisfy repetitive BI queries (Tanaka & Foster, 2021; Snowflake, 2024). Redshift can accommodate this pattern using concurrency scaling and WLM queues, but careful queue configuration and possible over-provisioning of compute may be needed to avoid queue wait times (AWS Concurrency Scaling, 2023; Singh & Das, 2021).
- Large-scale ETL and batch analytics: Ingest-heavy, long-running ETL jobs that produce wide, transformed tables are well-suited to Redshift's cluster approach when the workload is relatively stable and predictable, allowing the operator to optimize distribution styles and sort keys for sustained throughput. Redshift Spectrum can offload historical data to S3, reducing storage footprint and cost while preserving queryability (AWS Redshift Spectrum Guide, 2023). Snowflake can also handle this workload effectively, particularly when storage growth and compute isolation for ETL pipelines are desired; Snowflake's automatic micro-partitioning reduces the need for manual sort

key tuning but still benefits from well-chosen clustering keys when selective predicates are common (Snowflake, 2024; Tanaka & Foster, 2021).

- Mixed workloads with highly selective queries and point lookups: These favor architectures that minimize scan overhead and take advantage of predicate pushdown and appropriate data encodings. Redshift's column encodings and sort keys can optimize these patterns; Snowflake's clustering key strategy and pruning of micro-partitions are key levers to reduce scanned data volumes (AWS Redshift Compression Guide, 2024; Tanaka & Foster, 2021).

These archetypes align with the TPC-DS delineation of query mixes, mapping TPC-DS query complexities to real-world workload patterns to help practitioners choose tuning paths (TPC Council, 2023; Patil & Zade, 2023).

Storage & Compression Outcomes

Selecting compression encodings and physical layouts provides immediate returns in reduced storage costs and I/O. Literature and vendor guidance converge on the principle that column-specific compression yields twofold benefits: reduced storage footprint and reduced I/O when queries read a subset of columns. The AWS Redshift Compression Guide enumerates encodings (DELTA, LZ0, ZSTD, etc.) that suit different cardinalities; choosing the wrong encoding can increase decode overhead and reduce performance (AWS Redshift Compression Guide, 2024). Snowflake's automatic micro-partitioning and columnar compression aim to relieve the user of encoding selection, but operators still influence performance through clustering keys and by minimizing time-travel retention if historical versions are not required (Snowflake, 2024; Tanaka & Foster, 2021).

Expected outcomes when storage tuning is applied correctly include substantial reductions in scan volume for selective queries (often >70% reduction in scanned bytes), decreased per-query latency, and predictable storage costs. Conversely, improper encoding or overly aggressive time-travel retention can increase storage costs materially.

Compute Resource Management Outcomes

The compute layer is the primary lever for aligning performance with cost in cloud billing models. Snowflake's per-second billing for warehouses with auto-suspend and auto-resume features allows for fine-grained control: micro-bursts can be absorbed without long-term cost commitments, and multi-cluster warehouses can be provisioned to handle concurrency spikes. For workloads with unpredictable peaks and many small queries, this often results in better cost-performance profiles relative to static clusters (Snowflake, 2024; Tanaka & Foster, 2021).

Redshift's elastic resize and concurrency scaling permit scaling up for heavy workloads and scaling down to conserve cost, but resizing operations and cluster reconfiguration can impose temporal overhead and require planning. Redshift Reserved Instances or Savings Plans can lower baseline costs for predictable workloads, but they reduce flexibility (AWS Elastic Resize Documentation, 2023; AWS Concurrency Scaling, 2023). The literature emphasizes the governance trade-off: predictable workloads that can commit to reservations in Redshift may achieve lower cost per query, while volatile or multi-tenant workloads typically benefit from Snowflake-style elasticity (Rao & Karim, 2018; Singh & Das, 2021).

Query & Index Optimization Outcomes

Query profiling and structural optimization consistently demonstrate outsized impact. Several case studies and empirical analyses indicate that careful query rewrites, predicate pushdown, and the use of materialized views can reduce execution time by orders of magnitude for complex, join-heavy queries (Lee & Martin, 2020; Narayanan & Chu, 2020). For Redshift, appropriate distribution keys reduce network shuffle during joins substantially; correct sort key selection reduces I/O for range scans. Snowflake lacks user-controlled distribution keys but compensates

with micro-partitioning and clustering; the consequence is a trade-off between administrative simplicity and the degree of granular control available for micro-optimizations (Snowflake, 2024; James & Ortega, 2019).

For both platforms, the implementation of materialized views or result caching must be governed by workload access patterns and data freshness requirements. When query results are reusable across sessions and data change is limited, materialized views or caching reduce repeated computation costs and improve responsiveness. However, maintenance cost and staleness risk must be carefully balanced (Tanaka & Foster, 2021; Narayanan & Chu, 2020).

Cost Governance Outcomes

Informed cost governance combines predictive budgeting with reactive controls. Many organizations adopt hybrid strategies: reserving baseline capacity for predictable pipelines and relying on on-demand elasticity or multi-cluster provisioning for ad-hoc analytical loads. Chargeback and tagging systems improve transparency, enabling cost attribution by team, department, or workload (Rao & Karim, 2018; Malviya, 2025). The governance layer also includes automated rules—e.g., auto-suspend below utilization thresholds, retention window limits, and policies to move cold data to cheaper storage tiers or to Redshift Spectrum/S3—driven by telemetry.

The expected result of robust governance is a reduction in “wasted” spend—idle compute and excessive time-travel retention—that often constitutes a non-trivial fraction of cloud spend in poorly managed warehouses. Malviya (2025) and Rao & Karim (2018) both report actionable frameworks for identifying and capturing these savings.

DISCUSSION

The preceding results synthesize a complex body of knowledge into a prescriptive framework for cloud data warehousing. This section provides deeper interpretation, explores theoretical implications, articulates caveats and limitations, and suggests directions for empirical validation and automation.

Interpretation of Key Trade-offs

The most important high-level conclusion is that architectural choice and operational discipline jointly determine cost-performance outcomes. Snowflake’s decoupled architecture simplifies operational management and provides powerful elasticity, which is particularly advantageous for organizations with unpredictable query patterns or varied teams running independent workloads (Snowflake, 2024; Tanaka & Foster, 2021). The decoupling reduces the risk of noisy neighbors because compute isolation is built into warehouse provisioning. The consequence is lower administrative overhead and more predictable per-query response-time behavior for multi-tenant environments.

Conversely, Redshift’s cluster-based model—with carefully chosen sort and distribution keys, compression encodings, and the use of Spectrum for colder data—permits operators to squeeze additional efficiency for predictable, sustained workloads (AWS, 2024; AWS Redshift Compression Guide, 2024). The ability to reserve capacity further reduces per-query cost when utilization is high and consistent. Thus, cost-performance is not a property solely of the platform, but of the combination of workload predictability, operators’ ability to tune platform-specific levers, and willingness to make economic commitments (Rao & Karim, 2018; Singh & Das, 2021).

The role of query profiling emerges as foundational: without rigorous telemetry and workload classification, even the best-laid tuning plans will be mismatched to actual demand. Techniques and tools for profiling—capturing distribution of query latencies, identifying heavy queries, and mapping query shapes to resource consumption—form the necessary diagnostic layer that triggers subsequent tuning actions (Lee & Martin, 2020; Narayanan & Chu, 2020).

Theoretical Implications

From an architectural theory perspective, this analysis underscores the continuing relevance of classical database optimization principles—even in cloud contexts where storage-compute separation and managed services abstract

many low-level choices. Principles such as locality of reference, minimizing shuffles in distributed joins, and choosing physical layouts to align with predicate selectivity remain central. What changes in cloud platforms is the set of levers and their associated cost semantics: decisions now must consider not only latency and throughput but also per-second billing, storage egress, and reservation economics. This integration of economic modeling into system optimization is an important theoretical shift in the design space of modern data management systems (Li & Martin, 2020; Rao & Karim, 2018).

The framework also invites formal modeling of workload elasticity—characterizing the stochastic process of query arrivals and resource utilization and optimizing provisioning decisions under uncertain demand and non-linear cost functions. Such modeling would extend traditional capacity planning into stochastic control problems that account for cloud pricing schemes, reservation discounts, and the latency-sensitive nature of interactive analytics. While beyond this paper’s scope, these formal models are promising areas for rigorous research, combining queuing theory, optimization, and cost-aware systems design (Singh & Das, 2021; Malviya, 2025).

Limitations and Caveats

A first limitation is that this study synthesizes existing results and vendor guidance rather than conducting primary experimental benchmarking. While aligning the framework with TPC-DS and vendor documentation ensures practical relevance, practitioners should validate recommendations in their unique context using representative datasets and realistic concurrency. The exact gains from an optimization—e.g., percent reduction in query latency due to clustering or encoding changes—vary with data distributions, query mixes, and cloud region pricing.

A second caveat pertains to vendor feature evolution. Cloud platforms continuously introduce new features, pricing models, and managed optimizations that can materially change the cost-performance calculus. For instance, new caching layers, federation capabilities, or changes in billing granularity could alter the relative benefits of the levers discussed. Therefore, the framework must be treated as a living set of guidelines that requires periodic re-evaluation against vendor updates and fresh telemetry (Snowflake, 2024; AWS, 2024).

A related limitation is that not all enterprises have the operational capacity to implement certain layers fully. Large enterprises with dedicated data platform teams might execute complex distribution key strategies and automated governance policies; smaller organizations may prefer Snowflake’s managed model to avoid heavy operational overhead. The framework is intended to be adaptable along these organizational capability dimensions, but the burden of implementation should be acknowledged.

Future Work and Automation Opportunities

There are compelling avenues for empirical validation and automation. A natural extension is to implement the framework in a decision-support tool that consumes query telemetry, identifies workload archetypes, and outputs recommended configuration changes such as warehouse sizes, clustering keys, materialized views to create, or compression encodings to apply. Such a tool could use supervised learning from historical telemetry to predict the impact of tuning actions, or employ reinforcement learning to iteratively apply low-risk actions and measure outcomes. However, automation must be conservative; platform changes can produce cascading effects and require human oversight, especially when reservations or irreversible maintenance windows are involved.

Empirical work could validate the framework across industries (insurance, retail, finance) and dataset scales, ideally publishing open datasets and reproducible experiments. Additionally, research into formal modeling of the economic decision problem—combining workload stochastic models with cloud pricing structures—would provide theoretical guarantees on cost-optimal provisioning strategies under uncertainty (Patil & Zade, 2023; Li & Martin, 2020).

CONCLUSION

Cloud data warehousing presents a rich, multidimensional optimization problem bridging systems engineering, query optimization theory, and cost governance. This paper contributes a comprehensive five-layer framework that integrates workload characterization, storage and compression strategies, compute resource management, query-level optimizations, and cost governance policy. By mapping concrete, vendor-specific levers—such as Snowflake’s multi-cluster warehouses and clustering keys or Redshift’s sort/distribution keys and spectrum offload—to workload archetypes, the framework enables practitioners to systematically align platform choice and configuration with performance and cost objectives.

Key takeaways include: (1) workload profiling is indispensable; without it, tuning decisions will be misaligned with actual demand (Lee & Martin, 2020); (2) Snowflake’s decoupled architecture simplifies concurrency and elasticity concerns and often benefits unpredictable, multi-tenant workloads (Snowflake, 2024; Tanaka & Foster, 2021); (3) Redshift’s cluster-centric model can achieve cost efficiencies for predictable, sustained workloads when operators apply careful distribution, sort, and compression strategies and leverage spectrum for cold data (AWS, 2024; AWS Redshift Compression Guide, 2024); and (4) rigorous cost governance—combining reservations, auto-suspend policies, and chargeback—reduces wasteful spend (Rao & Karim, 2018; Malviya, 2025).

There is no universal best platform—the optimal choice depends on workload shape, organizational capacity, and cost flexibility. The proposed framework offers a structured path to making those choices defensible, repeatable, and traceable. Future empirical validation and automation of the framework could further reduce operational overhead and improve responsiveness to changing workload patterns, thereby advancing both practical cloud operations and the theoretical basis for cost-aware data management.

REFERENCES

1. Gupta, H., & Sharma, M. (2021). A Comparative Study of Traditional and Cloud-Native Data Warehousing Platforms. *International Journal of Computer Applications*, 183(35), 1–6.
2. Malik, A., & Bhatia, R. (2022). Cloud Data Warehousing: Performance Optimization in Snowflake. *Journal of Cloud Computing*, 11(1), 55–73.
3. Patil, A., & Zade, A. (2023). Performance Tuning in Big Data Environments: RDBMS vs. Snowflake. *Data Engineering Journal*, 9(4), 204–219.
4. Li, X., & Martin, P. (2020). Query Optimization in Hybrid Cloud Databases. *ACM Transactions on Database Systems*, 45(2), 17–36.
5. TPC Council. (2023). TPC Benchmark™ DS Overview. <https://www.tpc.org/tpcds/>
6. Rao, D., & Karim, A. (2018). Cost and Performance Analysis of Cloud Data Warehouses. *International Journal of Cloud Applications*, 6(3), 155–172.
7. James, T., & Ortega, M. (2019). Indexing Techniques in Relational Databases: A Review. *Database Technology Review*, 13(2), 45–63.
8. Narayanan, V., & Chu, H. (2020). Optimizing SQL Performance in Cloud-based Platforms. *Journal of Data Engineering*, 22(1), 89–104.
9. Lee, J. H., & Martin, P. (2020). The Role of Query Profiling in Performance Tuning. *ACM Transactions on Database Systems*, 45(4), 29–52.
10. Singh, R., & Das, S. (2021). Elastic Compute in Modern Data Warehousing. *Cloud Systems Review*, 17(1), 134–150.

11. Malviya, S. (2025). A Five-Layer Framework for Cost Optimization in Snowflake: Applied to P&C Insurance Workloads. *The American Journal of Interdisciplinary Innovations and Research*, 7(07), 28–43. <https://doi.org/10.37547/tajir/Volume07Issue07-04>
12. Tanaka, M., & Foster, B. (2021). Caching and Storage Optimization in Snowflake. *Journal of Cloud Optimization*, 10(3), 200–218.
13. AWS. (2024). AWS Redshift Documentation. <https://aws.amazon.com/redshift/>
14. Snowflake. (2024). Snowflake Documentation. <https://docs.snowflake.com>
15. AWS Redshift Compression Guide. (2024). https://docs.aws.amazon.com/redshift/latest/dg/c_Compression_encodings.html
16. AWS Glue Documentation. (2024). <https://docs.aws.amazon.com/glue/latest/dg/whatis-glue.html>
17. AWS Redshift Query Analyzer Guide. (2024). <https://aws.amazon.com/redshift/features/query-analyzer/>
18. AWS Elastic Resize Documentation. (2023). <https://docs.aws.amazon.com/redshift/latest/mgmt/elastic-resize.html>
19. AWS Concurrency Scaling. (2023). <https://docs.aws.amazon.com/redshift/latest/dg/concurrency-scaling.html>
20. AWS Redshift Spectrum Guide. (2023). <https://aws.amazon.com/redshift/features/spectrum/>
21. Snowflake Warehouse Management. (2024). <https://docs.snowflake.com/en/userguide/warehouses-overview.html>