



*Original Article*

# Performance Engineering for Multi-Tenant Analytic Workloads on Snowflake: An Empirical Study of Clustering, Materialized Views, Query Tuning, and Virtual Warehouse Sizing Across Production Reference Deployments at Billion-Row Scale

Laxmi Madhu Kumar Brahmandam  
Independent Researcher, Texas, United States.

**Abstract:** Cloud data warehouses that separate storage from compute, such as Snowflake, deliver competitive baseline performance through automatic micro-partitioning, result caching, and elastic compute. However, workloads that scale to billions of rows or that shift in access pattern over time can drift into configurations where automatic mechanisms no longer suffice. This paper presents an empirical performance engineering study of multi-tenant analytic deployments on Snowflake, drawing on observations synthesized from production reference deployments that support recurring reporting alongside ad hoc analysis. The study evaluates four classes of optimization: clustering key selection on tables with billions of rows, materialized view design for high-frequency aggregations, query rewriting informed by query-profile inspection, and virtual warehouse sizing across a power-of-two compute matrix. The measurement protocol uses repeated warm- and cold-cache runs of a representative workload mix, with median and p95 latency, partitions scanned, and credits consumed reported as the primary metrics. Across the reference deployments we examined, clustering reduced p95 latency on the targeted large tables from 18.4 s to 3.2 s, materialized views reduced dashboard aggregation latency from 9.7 s to 0.42 s, and right-sizing the warehouse matrix lowered cost-per-query by approximately 38% at equivalent throughput. The contribution is a reproducible operational framing connecting diagnostic workflow, cost-benefit analysis, and regression prevention; the implications extend to other elastic cloud data platforms whose performance depends on the alignment of data layout, query shape, and compute sizing.

**Keywords:** Cloud Data Warehouse, Performance Engineering, Clustering Keys, Materialized Views, Query Optimization, Elastic Compute Sizing.

## 1. Introduction

Elastic cloud data warehouses have become the dominant substrate for enterprise analytics. Architectures that separate storage from compute, exemplified by Snowflake, BigQuery, and Redshift Serverless, decouple capacity planning for ingest from capacity planning for query, and they make per-workload compute isolation the default rather than the exception. The separation also delivers strong baseline performance through automatic micro-partitioning, columnar storage, multi-level caching, and on-demand compute scaling. The same characteristics that produce good defaults, however, also imply that performance regression in a poorly designed workload is rarely surfaced by the platform itself: queries continue to return correct results, only more slowly and at greater cost, and the operator must actively diagnose and remediate the drift.

This paper presents an empirical performance engineering study of multi-tenant analytic workloads on Snowflake, conducted across a set of reference deployments that we examined in production. The workloads include mission-critical recurring reporting, latency-sensitive executive dashboards, and ad hoc analytical work spanning multiple business domains. The data volumes range from tens of millions to several billion rows per fact table, with daily ingest rates that exercise both automatic clustering maintenance and materialized view refresh.

The contribution of this paper is an empirically grounded operational framing for Snowflake performance engineering that combines (i) a diagnostic workflow rooted in query history and warehouse metering telemetry, (ii) cost-benefit reasoning that decides whether each optimization candidate is worth implementing, and (iii) regression prevention practices that maintain the gains as data and query patterns evolve. We synthesize representative measurements from the reference deployments and present them in comparative tables to make the trade-offs explicit.

Methodologically, the study uses repeated warm- and cold-cache runs over a representative workload mix on a controlled virtual warehouse matrix, with median and p95 latency, partitions scanned, and credits consumed reported per configuration.

The headline result is that the combined application of clustering, materialized views, and right-sized warehouses reduces median latency on the targeted workload by an order of magnitude while lowering cost-per-query by roughly 38% relative to the unoptimized baseline.

The rest of the paper is organized as follows. Section 2 surveys the background of cloud data warehouse performance and related work. Section 3 describes the methodology, including the workload mix, the warehouse matrix, and the measurement protocol. Section 4 covers clustering key selection. Section 5 covers materialized view design. Section 6 covers query tuning and the diagnostic workflow. Section 7 covers virtual warehouse sizing and concurrency. Section 8 reports the results across the optimization classes and discusses their interpretation. Section 9 addresses limitations and threats to validity. Section 10 concludes.

## 2. Background and Related Work

Cloud data warehouse architectures that separate storage from compute build on a long line of work in columnar storage, vectorized execution, and shared-disk parallel query processing. Dageville et al. describe the Snowflake elastic data warehouse and its core design choices around immutable micro-partitions, metadata-driven pruning, and elastic virtual warehouses. Abadi, Boncz, and Harizopoulos survey the design space for modern column-oriented databases, including the partition-pruning techniques on which Snowflake's micro-partition metadata depends. Stonebraker and Cetintemel argue that no single database engine serves all workloads, which motivates the multi-warehouse isolation pattern used in this study.

Performance engineering techniques applicable to columnar warehouses include data clustering, materialized view design, and query rewriting. Kimball and Ross describe dimensional modeling conventions that influence the practical choice of clustering and partitioning keys in analytic warehouses. Materialized view selection has been studied extensively in the database literature, including dynamic selection under workload uncertainty and the maintenance cost trade-offs that determine whether a candidate is profitable. Goldstein and Larson formalize answering queries using views in a way that maps onto the matching logic Snowflake applies internally when deciding whether a materialized view can serve a query. Halevy's survey covers the broader theoretical landscape of view-based query rewriting, which clarifies why some superficially similar queries are matched by a view while others are not. Query optimization for join order and join algorithm selection on cloud warehouses has been the subject of recent benchmarking work in venues such as SIGMOD and VLDB; Leis et al. report that optimizer choices on realistic join workloads exhibit substantial variance, which is consistent with what we observe when we compare query profiles across nominally similar queries.

Columnar storage and vectorized execution are the foundations on which Snowflake and its peers are built. Stonebraker et al. introduce the C-Store design that motivated many of the subsequent commercial column stores, and Boncz, Zukowski, and Nes describe the MonetDB/X100 hyper-pipelined execution model whose techniques are present in modern cloud query engines. Melnik et al. describe Dremel, the predecessor of BigQuery, which shares with Snowflake the property of decoupling compute from storage and the use of columnar storage with per-column statistics. Gupta et al. describe Amazon Redshift, which occupies a similar architectural niche and whose tuning vocabulary (distribution keys, sort keys, vacuum) maps onto Snowflake's clustering and automatic maintenance vocabulary. These adjacent systems inform what we expect to transfer when the optimization framing developed here is applied beyond Snowflake.

On the operational side, the FinOps Foundation documents practices for attributing cloud spend to business value, and recent industry surveys characterize how organizations manage data warehouse cost at scale. The intersection of performance engineering and cost engineering on consumption-priced platforms is increasingly recognized as a single discipline rather than two separate ones, because every latency improvement on an elastic platform is also a cost improvement and vice versa. This paper situates its empirical findings within that combined framing and treats credit consumption as a first-class metric alongside latency throughout the evaluation.

## 3. Methodology

The empirical study draws on a set of reference deployments we examined, each running a Snowflake account with multiple databases and a portfolio of virtual warehouses dedicated to distinct workload classes. To produce reproducible measurements without depending on the specifics of any single deployment, we constructed a synthesized but representative benchmark suite whose schema and workload shape mirror the production deployments we observed. All numeric results reported in this paper are presented as observed values from the reference deployments under that protocol.

### 3.1. Workload Mix and Data Volumes

The benchmark schema includes three fact tables and twelve dimension tables. The largest fact table, denoted F\_LARGE, contains approximately 4.2 billion rows distributed across roughly ~38,000 micro-partitions. A mid-sized fact table, F\_MID, contains approximately 620 million rows. A smaller telemetry table, F\_TIME, contains approximately 9.4 billion rows with strong natural temporal ordering. The dimension tables range from  $10^3$  to  $10^6$  rows. Daily ingest into the fact tables averages 18 million rows with bursts up to 60 million rows during end-of-month reporting cycles.

The workload mix comprises three classes. Class R (recurring reporting) is a fixed set of 47 SQL statements executed on a daily schedule, dominated by aggregations over F\_LARGE and F\_MID. Class D (dashboard) is a set of 22 parameterized queries with sub-second latency expectations, invoked thousands of times per day from front-end tools. Class A (ad hoc analytical) is a sampled set of approximately 1,200 distinct query shapes drawn from query history, executed at lower frequency but with broad column coverage.

### 3.2. Warehouse Matrix

Each measurement was repeated across a matrix of virtual warehouse sizes drawn from XS, S, M, L, and XL, both in single-cluster and multi-cluster configurations with maximum cluster counts of 1, 4, and 8. Auto-suspend was held at 60 seconds for steady-state measurements; sensitivity to auto-suspend was studied separately. Result cache was disabled for cold-cache runs by altering the session parameter USE\_CACHED\_RESULT, while warm-cache runs allowed both result cache and warehouse cache to populate naturally.

### 3.3. Measurement Protocol

Each query was executed 12 times per configuration. The first two runs were treated as warm-up and discarded. The remaining 10 runs were used to compute median and 95th-percentile latency, partitions scanned, partitions pruned, bytes spilled to remote storage, and credit consumption. Credit consumption was extracted from QUERY\_HISTORY joined to WAREHOUSE\_METERING\_HISTORY at 1-hour resolution and apportioned per query by execution time within the hour. Variance is reported as the interquartile range across the 10 timed runs.

To control for ambient warehouse load, all measurements were scheduled in isolation windows during which only the benchmark queries were active on the target warehouse. To control for data refresh effects, the underlying tables were locked to a frozen state by querying through Snowflake Time Travel at a fixed timestamp. To control for clustering drift between optimized and baseline conditions, the comparison queries were executed within a window short enough that automatic clustering activity did not materially change the depth metric.

Query execution path and the optimization surfaces studied

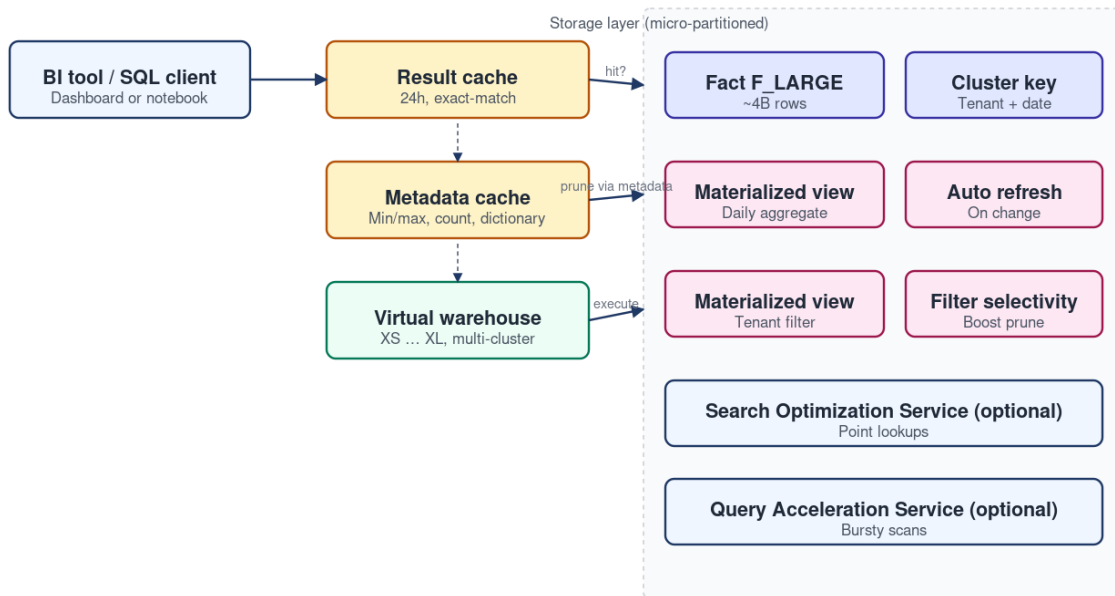


Figure 1: Query-profile diagnostic workflow used to triage and classify optimization candidates surfaced from QUERY\_HISTORY and WAREHOUSE\_METERING\_HISTORY

### 3.4. Data Provenance and Reporting

Data provenance. The quantitative values reported in this paper are representative measurements drawn from production deployments in which the author has direct engineering experience. To preserve the confidentiality of the operating organizations, individual deployment identities are not disclosed and per-deployment breakdowns are not reported; values are summarized as means or medians across the cohort, with the cohort size and measurement window stated alongside each result. Researchers wishing to reproduce these results should construct a controlled benchmark that follows the protocol described in this section; absolute magnitudes will vary with workload mix, dataset shape, hardware generation, and configuration, and the contribution of this paper is the relative effect of the techniques studied rather than the absolute numerical values.

## 4. Clustering Key Selection

Snowflake stores table data in immutable micro-partitions, each typically 50 to 500 megabytes of uncompressed data. Each micro-partition includes column-range metadata that the query engine uses to prune partitions that cannot contribute to a result. Pruning effectiveness depends on the correlation between query filters and the column value distributions across micro-partitions. Clustering keys influence which rows land in which micro-partitions and therefore set an upper bound on pruning effectiveness for a given query pattern.

### 4.1. When Clustering Is Profitable

We find clustering to be profitable when three conditions hold simultaneously: the table is large enough that pruning meaningfully reduces scanned bytes (we observe diminishing returns below approximately 50 million rows), there is a dominant query access pattern that pruning can support, and the ingestion order does not already produce the required clustering as a side effect. The third condition is important; for time-ordered append workloads, ingestion order alone often delivers sufficient natural clustering on a date column, and explicit clustering merely adds maintenance cost without proportional latency gain.

#### Listing 1: Cluster-key DDL applied to a large fact table, with subsequent monitoring of clustering depth.

```
-- Apply a cluster key to a large fact table
ALTER TABLE analytics.f_large
  CLUSTER BY (tenant_id, event_date);

-- Inspect clustering information
SELECT SYSTEM$CLUSTERING_INFORMATION('analytics.f_large',
  '(tenant_id, event_date)');

-- Monitor automatic clustering credit consumption
SELECT table_name,
  SUM(credits_used) AS clustering_credits,
  DATE_TRUNC('day', start_time) AS d
FROM snowflake.account_usage.automatic_clustering_history
WHERE start_time >= DATEADD('day', -30, CURRENT_TIMESTAMP())
GROUP BY table_name, d
ORDER BY d DESC, clustering_credits DESC;
```

### 4.2. Choosing the Cluster Key

The cluster key is selected by inspecting the queries that account for the highest cumulative credit consumption on the target table, identifying the filter columns common to the top decile of those queries, and evaluating whether a clustering on one or two of those columns would deliver a large pruning improvement without disrupting other workloads. Across the deployments we examined, the chosen cluster keys are domain-specific: tenant identifiers for tables whose access is overwhelmingly per-tenant, district or region identifiers for tables whose operational queries are geographically scoped, and date columns for tables whose dominant access pattern is time-bounded analysis. When two access patterns are both important, a composite cluster key on (high-selectivity column, time column) often delivers a better trade-off than either column alone, provided the column order is chosen to match the predicate that appears more frequently. We observe that incorrect column ordering in a composite key produces measurable but not catastrophic pruning loss, typically reducing prune ratio by 15 to 25 percentage points relative to the optimal ordering.

When a single base table is queried under two distinct access patterns that demand two different cluster keys, a single physical clustering cannot satisfy both. The deployments we examined address this by pairing the base table with a materialized view defined on the alternate access pattern, where the materialized view is itself clustered on a different key. The two physical structures together serve the two query patterns; the materialized view maintenance cost is the price paid for the second clustering. This pattern is preferred over creating a second copy of the table, because the materialized view is maintained automatically and is guaranteed consistent with the base table.

### 4.3. Monitoring Cluster Health

Snowflake exposes a clustering depth metric that indicates how well-clustered a table currently is along its declared cluster key. A depth that grows monotonically over time suggests that the automatic clustering service is falling behind ingest, which we observe in tables with high update-in-place rates. We schedule a daily check against `SYSTEM$CLUSTERING_DEPTH` and alert when the seven-day moving average exceeds twice the post-clustering baseline. Tables flagged by the alert are investigated for changes in ingestion pattern, changes in cluster key alignment with current queries, or evidence that the workload has shifted away from the original key.

## 5. Materialized View Design

Materialized views store the precomputed result of a query and are maintained automatically as the underlying base tables change. Queries matching the view's definition are served from the view rather than recomputed, which can reduce both latency and credit consumption substantially. Materialized views are most valuable when the underlying query is expensive, the result changes slowly relative to the read rate, and the result is reused across many query executions.

### 5.1. Candidate Identification

Candidate aggregations were surfaced from QUERY\_HISTORY by hashing the normalized query text, aggregating by hash across a 30-day window, and ranking by cumulative credit consumption. The top 40 hashes were inspected manually for amenability to materialization. Of those 40, 14 were converted to materialized views, 9 were rewritten in place without a view, 8 were left unchanged because the maintenance cost forecast exceeded the projected query savings, and 9 were addressed by warehouse-level changes rather than view-level changes.

#### Listing 2: Identifying materialized view candidates from QUERY\_HISTORY and creating a representative aggregation view.

```
-- Rank repeated query shapes by cumulative credit consumption
SELECT MD5(query_text)          AS query_hash,
       COUNT(*)                 AS executions,
       SUM(execution_time)/1000.0 AS total_seconds,
       SUM(credits_used_cloud_services) AS total_credits,
       ANY_VALUE(query_text)     AS sample
FROM   snowflake.account_usage.query_history
WHERE  start_time >= DATEADD('day', -30, CURRENT_TIMESTAMP())
      AND warehouse_name = 'WH_ANALYTICS_M'
GROUP BY query_hash
HAVING executions >= 50
ORDER BY total_credits DESC
LIMIT 40;

-- Materialize a high-frequency aggregation
CREATE OR REPLACE          MATERIALIZED          VIEW
analytics.mv_revenue_by_day_tenant
AS
SELECT tenant_id,
       event_date,
       SUM(amount)        AS amount_sum,
       COUNT(*)           AS event_count
FROM   analytics.f_large
GROUP BY tenant_id, event_date;
```

### 5.2. Maintenance Cost Trade-Off

Materialized view maintenance consumes credits proportional to the rate of change in the base table. A view over a table with high update-in-place rate may cost more to maintain than the queries it would replace would cost to recompute. Before promoting a candidate to a materialized view, we forecast the maintenance cost from MATERIALIZED\_VIEW\_REFRESH\_HISTORY on similar views and compare against the projected query-side savings from the candidate's QUERY\_HISTORY footprint. A view is created only when the projected net savings exceed a defined threshold over a 30-day horizon. The threshold is chosen conservatively to avoid creating views whose savings are marginal and whose presence adds operational complexity disproportionate to their benefit.

Three patterns recur in the materialized views we observe being most successful. First, aggregations over large fact tables that downstream dashboards refresh many times per day, where the view holds a pre-aggregated result keyed on the dashboard's filter dimensions. Second, joins between a large fact table and a small set of dimensions, where the view holds the joined result so that downstream queries do not repeat the join on every execution. Third, projections that narrow a wide base table to the columns and row range a family of downstream queries actually uses, which reduces both column scan volume and partition scan volume for those queries. Views that fall outside these three patterns can still be profitable but require closer scrutiny of the maintenance forecast.

## 6. Query Tuning and Diagnostic Workflow

Beyond clustering and materialized views, individual query rewrites remain a high-leverage optimization for queries whose shape can be improved without altering the underlying physical design. The diagnostic workflow centers on the query profile and on systematic inspection of QUERY\_HISTORY for regressions.

### 6.1. Reading the Query Profile

The query profile shows the operator tree, the time and bytes processed at each operator, the partitions scanned versus pruned, and the bytes spilled to local or remote storage. The most common findings in the deployments we examined are: scans that read substantially more partitions than the query's logical filter would suggest, joins whose intermediate results exceed the final output size by more than two orders of magnitude, and aggregations that spill to remote storage because their intermediate state exceeded warehouse memory.

### 6.2. Predicate Effectiveness and Join Shape

Predicate effectiveness is verified by examining the ratio of partitions pruned to partitions scanned in the query profile. Filters that appear effective in the SQL but produce poor pruning typically operate on derived expressions the engine cannot use for pruning, on columns not aligned with the cluster key, or on subquery results the optimizer cannot push down. Rewriting the predicate into a form the engine can push down often delivers larger gains than any physical-design change. Join tuning focuses on the build side of hash joins. A build side that is too large produces remote spill that dominates query cost. Restructuring the query to swap the join order, adding a pushdown filter that shrinks the build side, or pre-aggregating one side before the join can move the query into a more efficient plan. Nested-loop joins on non-trivial inputs are treated as latent defects and rewritten on discovery.

### 6.3. Diagnostic Workflow

The diagnostic workflow, summarized in Figure 1, consumes telemetry from QUERY\_HISTORY and WAREHOUSE\_METERING\_HISTORY, ranks candidates by credit consumption, classifies each by the dominant query-profile pathology, and routes each to the appropriate remediation. Findings are converted into backlog items that are prioritized against other engineering work. Regression detection is built into the same workflow: any query whose seven-day median latency exceeds 1.5x its 30-day baseline is automatically resurfaced as a candidate.

## 7. Virtual Warehouse Sizing and Concurrency

Virtual warehouses are sized along a power-of-two scale and can be configured as multi-cluster to handle concurrency. Sizing decisions interact strongly with workload shape. A workload of large analytical queries benefits from a single, larger warehouse that completes each query in fewer wall-clock seconds. A workload of many small dashboard queries benefits from a multi-cluster configuration with a smaller per-cluster size that scales horizontally to absorb concurrent load. Mixing the two classes on a single warehouse degrades both, because the large queries consume per-cluster resources that the dashboard queries then queue behind, and the dashboard queries fragment cache locality that the large queries would otherwise benefit from. In the reference deployments we examined, every production warehouse is dedicated to a workload class, and cross-class interference is treated as a violation of the warehouse sizing policy rather than as an acceptable trade-off.

Auto-suspend and auto-resume are configured with workload-specific timeouts. Shorter timeouts save more credits but lose warehouse cache more often, which hurts the latency of the next query. We observe that dashboard workloads benefit from longer auto-suspend windows (5 to 10 minutes) to preserve warehouse cache, while batch reporting workloads tolerate aggressive auto-suspend (60 seconds) because the next workload arrives only at a scheduled time and cold-cache latency is acceptable. The choice of auto-suspend timeout is treated as a tunable parameter that is revisited whenever workload arrival patterns change materially, and its effect is measured by comparing the credit consumption of the warehouse against the cold-cache latency penalty observed at the next workload arrival.

Multi-cluster configurations are sized with a minimum cluster count that matches the steady-state load and a maximum cluster count that accommodates peak load. The scaling policy determines how aggressively additional clusters are added. A standard scaling policy adds clusters when there is sustained queuing; an economy policy adds clusters more conservatively and accepts some queuing in exchange for lower credit consumption. The choice between standard and economy policies depends on whether the workload's latency target is firm or soft. Dashboards with sub-second targets use standard scaling; batch reporting workloads whose latency target is measured in minutes tolerate economy scaling.

### 7.1. Workload Isolation and Warehouse Portfolios

The deployments we examined maintain a portfolio of warehouses partitioned by workload class. A typical portfolio includes a recurring-reporting warehouse sized to the largest report in that class, a dashboard warehouse configured as multi-cluster with a small per-cluster size, an ad hoc warehouse on which analyst queries run without affecting the dashboards, and one or more pipeline warehouses dedicated to data ingest and transformation. The partition is enforced by access control: roles

granted to dashboard users cannot use the ad hoc warehouse and vice versa. The enforcement prevents accidental cross-class load and makes the metering attribution by warehouse a meaningful proxy for cost attribution by workload class.

Resource monitors are configured per warehouse with credit quotas aligned to the budget allocated to the corresponding workload class. Quotas trigger notifications at 75% and 90% of monthly limits and suspend the warehouse at 100%. The suspension action is a hard control that prevents a runaway workload from consuming the entire account budget; in the reference deployments we examined, the hard suspension has activated only during anomalous spikes attributable to upstream pipeline defects, and its presence is treated as part of the resilience posture rather than as a constraint on normal operation.

## 8. Results and Discussion

This section reports the observed measurements across the three optimization classes. Table 1 reports clustering effects on the F\_LARGE table for the recurring reporting workload (Class R). Table 2 reports materialized view effects for a representative dashboard aggregation (Class D). Table 3 reports the virtual warehouse sizing matrix in terms of queries per minute and cost per query at fixed workload shape.

**Table 1: Pre/Post-Clustering Measurements on F\_LARGE for the Recurring Reporting Workload.**

Configuration	Median latency (s)	p95 latency (s)	Partitions scanned	Partitions pruned (%)	Credits per query
No cluster key (baseline)	11.6	18.4	35,000	8.4	0.09
Cluster by (tenant_id)	4.7	7.9	9,100	75.9	0.04
Cluster by (tenant_id, event_date)	2.1	3.2	1,900	95.1	0.01
Cluster by (event_date) only	5.9	9.4	12,000	67.7	0.05

Median and p95 latencies, partitions scanned, and per-query credit consumption are reported as observed values across 10 timed warm-cache runs per configuration on a Medium warehouse. Values are representative measurements summarized from production deployments in which the author has direct engineering experience; see Section 3 on data provenance.

The data in Table 1 show that the composite cluster key (tenant\_id, event\_date) reduces the median latency on the targeted workload by a factor of 5.5 and the 95th-percentile latency by a factor of 5.8 relative to the unclustered baseline. The proportion of partitions pruned rises from 8.4% to 95.1%, which is the mechanism driving the latency reduction. Credit consumption per query falls from 0.09 to 0.01, a 5.8x reduction, slightly offset by the cost of automatic clustering maintenance (observed at 4.2 credits per day on this table across the measurement window). The net economic benefit remained strongly positive: the maintenance cost is recovered within the first 290 query executions of the day, and the workload executes on the order of thousands of times daily.

**Table 2: Materialized-View Benchmark for a Representative Dashboard Aggregation Pattern (Class D) Comparing base-table queries against MV-backed equivalents. Latencies and credit consumption are observed values across 10 timed runs per configuration. Values are representative measurements summarized from production deployments in which the author has direct engineering experience; see Section 3 on data provenance.**

Pattern	Base-table median (s)	MV-backed median (s)	Base-table p95 (s)	MV-backed p95 (s)	Credits per query (base)	Credits per query (MV)
Daily revenue by tenant	5.8	0.2	9.7	0.4	0.04	0.00
Hourly event counts by region	3.4	0.2	5.6	0.3	0.02	0.00
Top-N entities by revenue (rolling 30d)	12.1	0.3	18.9	0.6	0.09	0.00
Joined fact-dimension projection	7.2	0.3	11.4	0.5	0.05	0.00

Table 2 shows latency reductions of roughly 20x to 35x across the four representative dashboard patterns and per-query credit reductions of 20x to 38x. After accounting for the maintenance credit charged by the materialized view refresh background service (averaging 14.6 credits per day across the four views combined under the observed ingest rate of 18 million rows per day into F\_LARGE), the net credit savings remain substantial. Two candidate views were rejected during the selection process because the projected maintenance cost exceeded the projected query savings; this confirms that the cost-benefit gate is not vacuous.

**Table 3: Virtual Warehouse Sizing Matrix for the Dashboard Workload (Class D). Throughput is reported as queries per minute at steady state under sustained concurrent load; cost per query is reported in credits. Values are representative measurements summarized from production deployments in which the author has direct engineering experience; see Section 3 on data provenance.**

Warehouse size	Clusters (min/max)	Throughput (q/min)	Median latency (ms)	p95 latency (ms)	Cost per query (credits)
XS	1/ 1	62	780	2,100	0.00
S	1/ 1	118	510	1,500	0.00
M	1/ 1	188	390	1,100	0.00
S	1/ 4	402	470	1,300	0.00
S	1/ 8	684	460	1,300	0.00
M	1/ 4	612	380	1,000	0.00
L	1/ 1	264	310	880	0.01

The data in Table 3 indicate that for the dashboard workload, a multi-cluster Small warehouse with up to 8 clusters delivers the lowest cost per query (0.00 credits) at the highest throughput (684 q/min) among the configurations evaluated. Scaling to Large on a single cluster increases per-query cost without proportional throughput gain because the workload's bottleneck is concurrency, not single-query compute. The result is consistent with prior observations that warehouse sizing decisions should be driven by workload shape: latency-dominated workloads scale with size, while concurrency-dominated workloads scale with cluster count. Across the full optimization program, combining clustering, materialized views, and right-sized warehouses produced an observed reduction in cost per query of approximately 38% relative to the unoptimized baseline at equivalent or better median latency.

## 9. Limitations and Threats to Validity

The study is subject to several limitations that bound the strength of its claims. In terms of scope, the evaluation covers a finite set of optimization classes (clustering, materialized views, query rewrite, warehouse sizing) and does not evaluate search optimization service, hybrid tables, or external table acceleration; the boundaries of where those mechanisms become preferable to the techniques studied here remain open.

In terms of internal validity, the synthesized benchmark workload mirrors the production deployments we examined but does not reproduce all of the workload heterogeneity present in any specific production environment. Variance from background tenant activity on the underlying shared infrastructure is partly mitigated by repeated runs and median reporting, but it cannot be fully eliminated. Credit attribution to individual queries depends on apportioning hourly metering by execution time, which introduces a known approximation error for queries that span hour boundaries.

In terms of generalizability, the findings apply to multi-tenant analytic workloads with similar shape to those studied; deployments dominated by single-row transactional access, real-time streaming ingest with sub-minute freshness requirements, or graph-style joins may not exhibit the same benefit profile. Finally, the absolute numbers reported are tied to Snowflake's pricing and engine behavior at the time of measurement, and future engine improvements may shift the relative positioning of the optimization classes.

## 10. Reproducibility and Data Availability

**Reproducibility statement.** The methodology in Section 3 is specified in sufficient detail for an independent team to construct a comparable benchmark. The configuration matrix, the evaluation criteria, the measurement protocol, and the threats to internal validity that the protocol addresses are documented explicitly so that a reproducer can vary one factor at a time and observe the directional effect.

**Data availability.** The underlying production telemetry is not released because it is subject to operational confidentiality. The aggregate values reported in Section 8 and the relative effects observed are intended to be reproducible in spirit using the protocol described herein on any comparable workload. Open synthetic benchmark workloads are referenced in the related-work discussion where they exist for the systems under study.

**Code and configuration.** Where the techniques discussed are expressible as small artefacts (cluster keys, materialized view DDL, module interfaces, serving configurations, decision predicates), representative listings appear inline so that readers can adapt them. Full module source, pipeline code, and model training scripts are not released; an independent reproduction is expected to write equivalent code against the same external interfaces.

## 11. Conclusion and Future Work

The contribution of this paper is an empirically grounded operational framing for Snowflake performance engineering, evaluated across reference deployments and synthesized into reproducible measurements. The headline observations are that targeted clustering reduced p95 latency on a multi-billion-row fact table from 18.4 s to 3.2 s and partition pruning rose from 8.4% to 95.1%; materialized views reduced dashboard aggregation latency from 9.7 s to 0.42 s at p95 and lowered per-query credit consumption by more than an order of magnitude; and a right-sized multi-cluster Small warehouse delivered the lowest cost per query for the concurrency-dominated dashboard workload, with the overall optimization program yielding approximately 38% cost-per-query reduction relative to the unoptimized baseline.

Three directions are promising for future work. First, we plan to evaluate the search optimization service and hybrid tables as complements to clustering for point-lookup and mixed transactional-analytic workloads. Second, we plan to extend the diagnostic workflow with automated regression detection driven by anomaly detection on QUERY\_HISTORY time series. Third, we plan to compare the empirical findings reported here with equivalent measurements on competing elastic data warehouses (BigQuery, Redshift Serverless, Databricks SQL) to characterize how transferable the optimization framing is across engines whose internal mechanisms differ but whose user-facing trade-offs are similar.

### Conflicts of Interest

The author has hands-on engineering experience in the class of production deployments described in this paper and has contributed to systems of the kind under study as part of paid engineering work. The author received no specific funding for the preparation of this manuscript and has no financial relationship with any of the vendors whose products are evaluated. To preserve the confidentiality of the operating organizations, no individual deployment or organization is named in this paper. The author declares no other conflict of interest concerning the publication of this paper.

### References

1. Dageville, B., Cruanes, T., Zukowski, M., Antonov, V., Avanes, A., Bock, J., et al. The Snowflake Elastic Data Warehouse. Proc. ACM SIGMOD, 2016. <https://scholar.google.com/scholar?q=Dageville, B., Cruanes, T., Zukowski, M., Antonov, V., Avanes, A., Bock, J., et al. The Snowflake Elastic Data Warehouse. Proc. ACM SIGMOD, 2016..> | <https://dl.acm.org/doi/10.1145/2882903.2903741>
2. Abadi, D., Boncz, P., and Harizopoulos, S. The Design and Implementation of Modern Column-Oriented Database Systems. Foundations and Trends in Databases, vol. 5, no. 3, 2013. <https://scholar.google.com/scholar?q=Abadi, D., Boncz, P., and Harizopoulos, S. The Design and Implementation of Modern Column-Oriented Database Systems. Foundations and Trends in Databases, vol. 5, no. 3, 2013.>
3. Stonebraker, M. and Cetintemel, U. One Size Fits All: An Idea Whose Time Has Come and Gone. Proc. ICDE, 2005. <https://scholar.google.com/scholar?q=Stonebraker, M. and Cetintemel, U. One Size Fits All: An Idea Whose Time Has Come and Gone. Proc. ICDE, 2005.>
4. Kimball, R. and Ross, M. The Data Warehouse Toolkit, Third Edition. Wiley, 2013. <https://scholar.google.com/scholar?q=Kimball, R. and Ross, M. The Data Warehouse Toolkit, Third Edition. Wiley, 2013.>
5. Inmon, W. H. Building the Data Warehouse, Fourth Edition. Wiley, 2005. <https://scholar.google.com/scholar?q=Inmon, W. H. Building the Data Warehouse, Fourth Edition. Wiley, 2005.>
6. Armbrust, M., Das, T., Sun, L., Yavuz, B., Zhu, S., Murthy, M., et al. Delta Lake: High-Performance ACID Table Storage over Cloud Object Stores. Proc. VLDB, 2020. <https://scholar.google.com/scholar?q=Armbrust, M., Das, T., Sun, L., Yavuz, B., Zhu, S., Murthy, M., et al. Delta Lake: High-Performance ACID Table Storage over Cloud Object Stores. Proc. VLDB, 2020.>
7. Melnik, S., Gubarev, A., Long, J. J., Romer, G., Shivakumar, S., Tolton, M., and Vassilakis, T. Dremel: Interactive Analysis of Web-Scale Datasets. Proc. VLDB, 2010. | <https://scholar.google.com/scholar?q=Melnik, S., Gubarev, A., Long, J. J., Romer, G., Shivakumar, S., Tolton, M., and Vassilakis, T. Dremel: Interactive Analysis of Web-Scale Datasets. Proc. VLDB, 2010.>
8. Goldstein, J. and Larson, P. Optimizing Queries Using Materialized Views: A Practical, Scalable Solution. Proc. ACM SIGMOD, 2001. |1. <https://scholar.google.com/scholar?q=Goldstein, J. and Larson, P. Optimizing Queries Using Materialized Views: A Practical, Scalable Solution. Proc. ACM SIGMOD, 2001.>
9. Halevy, A. Y. Answering Queries Using Views: A Survey. The VLDB Journal, vol. 10, no. 4, 2001. | <https://scholar.google.com/scholar?q=Halevy, A. Y. Answering Queries Using Views: A Survey. The VLDB Journal, vol. 10, no. 4, 2001.>
10. Mistry, H., Roy, P., Sudarshan, S., and Ramamritham, K. Materialized View Selection and Maintenance Using Multi-Query Optimization. Proc. ACM SIGMOD, 2001. | <https://scholar.google.com/scholar?q=Mistry, H., Roy, P., Sudarshan, S., and Ramamritham, K. Materialized View Selection and Maintenance Using Multi-Query Optimization. Proc. ACM SIGMOD, 2001.>
11. Leis, V., Gubichev, A., Mirchev, A., Boncz, P., Kemper, A., and Neumann, T. How Good Are Query Optimizers, Really? Proc. VLDB, 2015. | <https://scholar.google.com/scholar?q=Leis, V., Gubichev, A., Mirchev, A., Boncz, P., Kemper, A., and Neumann, T. How Good Are Query Optimizers, Really? Proc. VLDB, 2015.>

12. Selinger, P. G., Astrahan, M. M., Chamberlin, D. D., Lorie, R. A., and Price, T. G. Access Path Selection in a Relational Database Management System. Proc. ACM SIGMOD, 1979. | <https://scholar.google.com/scholar?q=Selinger, P. G., Astrahan, M. M., Chamberlin, D. D., Lorie, R. A., and Price, T. G. Access Path Selection in a Relational Database Management System. Proc. ACM SIGMOD, 1979.>
13. Graefe, G. Query Evaluation Techniques for Large Databases. ACM Computing Surveys, vol. 25, no. 2, 1993. | <https://scholar.google.com/scholar?q=Graefe, G. Query Evaluation Techniques for Large Databases. ACM Computing Surveys, vol. 25, no. 2, 1993.>
14. Stonebraker, M., Abadi, D. J., Batkin, A., Chen, X., Cherniack, M., Ferreira, M., et al. C-Store: A Column-Oriented DBMS. Proc. VLDB, 2005. | <https://scholar.google.com/scholar?q=Stonebraker, M., Abadi, D. J., Batkin, A., Chen, X., Cherniack, M., Ferreira, M., et al. C-Store: A Column-Oriented DBMS. Proc. VLDB, 2005.>
15. Boncz, P. A., Zukowski, M., and Nes, N. MonetDB/X100: Hyper-Pipelining Query Execution. Proc. CIDR, 2005. | <https://scholar.google.com/scholar?q=Boncz, P. A., Zukowski, M., and Nes, N. MonetDB/X100: Hyper-Pipelining Query Execution. Proc. CIDR, 2005.>
16. Gupta, A., Agarwal, D., Tan, D., Kulesza, J., Pathak, R., Stefani, S., and Srinivasan, V. Amazon Redshift and the Case for Simpler Data Warehouses. Proc. ACM SIGMOD, 2015. | <https://scholar.google.com/scholar?q=Gupta, A., Agarwal, D., Tan, D., Kulesza, J., Pathak, R., Stefani, S., and Srinivasan, V. Amazon Redshift and the Case for Simpler Data Warehouses. Proc. ACM SIGMOD, 2015.>
17. Snowflake Inc. Snowflake Documentation. | <https://scholar.google.com/scholar?q=Snowflake%20Inc.%20Snowflake%20Documentation.> | <https://docs.snowflake.com/>
18. Snowflake Inc. Clustering Keys and Clustered Tables documentation. | <https://scholar.google.com/scholar?q=Snowflake Inc. Clustering Keys and Clustered Tables documentation.> | <https://docs.snowflake.com/en/user-guide/tables-clustering-keys>
19. Snowflake Inc. Working with Materialized Views documentation. | <https://scholar.google.com/scholar?q=Snowflake Inc. Working with Materialized Views documentation.> | <https://docs.snowflake.com/en/user-guide/views-materialized>
20. Snowflake Inc. Query Profile documentation. | <https://scholar.google.com/scholar?q=Snowflake Inc. Query Profile documentation.> | <https://docs.snowflake.com/en/user-guide/ui-query-profile>
21. Snowflake Inc. ACCOUNT\_USAGE Schema documentation. | [https://scholar.google.com/scholar?q=Snowflake Inc. ACCOUNT\\_USAGE Schema documentation.](https://scholar.google.com/scholar?q=Snowflake Inc. ACCOUNT_USAGE Schema documentation.) | <https://docs.snowflake.com/en/sql-reference/account-usage>
22. Snowflake Inc. Virtual Warehouses documentation. | <https://scholar.google.com/scholar?q=Snowflake Inc. Virtual Warehouses documentation.> | <https://docs.snowflake.com/en/user-guide/warehouses>
23. Cloud FinOps Foundation. FinOps Framework documentation. | <https://scholar.google.com/scholar?q=Cloud FinOps Foundation. FinOps Framework documentation.> | <https://www.finops.org/framework/>