# Optimize Warehouse Configuration in Snowflake

## Real-World Problem: Unpredictable Performance or Cost Overruns in Batch & Real-Time Processing

### Industry Case Study: Amazon's Real-Time Logistics Data Processing

Amazon ingests and processes logistics, delivery, and supply chain updates every few seconds. These updates trigger tasks that push delivery status changes, generate metrics for dashboards, and archive raw data for audit purposes.

**Challenge:** During surge periods (like Prime Day), massive parallel data comes in from multiple regions. If warehouses are undersized or misconfigured:

- Processing gets delayed
- Queues grow
- Cost skyrockets due to inefficient scaling

### Snowflake's Solution: Granular, Elastic, Configurable Warehouses

Snowflake decouples storage and compute. You can assign the right compute power (warehouse) to your workload, automatically or manually scale them, suspend when idle, and avoid over-provisioning.

A **Virtual Warehouse** is a cluster of compute resources in Snowflake that processes queries and tasks.

| Component | Description |
|-----------|-------------|
| Size | XS, S, M, L, XL, up to 6XL |
| Scaling | Multi-cluster (min/max clusters) |
| Billing | Per-second (1-minute minimum) |
| Suspension | Auto-suspend when idle |
| Auto-Resume | Instantly resumes on demand |

## Key Warehouse Configuration Parameters

### 1. Size (XS to 6XL)

Choose warehouse size based on:

- **Query complexity** (e.g., joins, aggregations)
- **Data volume**
- **Concurrency needs**

Example: For a real-time dashboard that refreshes every minute using small result sets, an **XS warehouse** is optimal.

In contrast, **Google Cloud's Billing Team** uses **L or XL** for hourly aggregations across terabytes of usage data.

## 2. Multi-Cluster Warehouses

Enable **multi-cluster** for high concurrency workloads (e.g., BI dashboards with 100s of users).

- · Automatically spins up multiple clusters
- · Avoids query queuing
- · Charges per active cluster only

```
CREATE WAREHOUSE BI_CLUSTERED_WH
  WITH WAREHOUSE_SIZE = 'MEDIUM'
  AUTO_SUSPEND = 60
  AUTO_RESUME = TRUE
  MIN_CLUSTER_COUNT = 1
  MAX_CLUSTER_COUNT = 3
  SCALING_POLICY = 'ECONOMY';
```

**Netflix** uses multi-cluster warehouses to power peak-time content view dashboards for internal content analysts.

## 3. Auto-Suspend and Auto-Resume

- · **Auto-suspend** minimizes cost by shutting down idle warehouses.
- · **Auto-resume** ensures they spin up automatically when queried.
- · Recommended idle timeout: 60–300 seconds

```
ALTER WAREHOUSE TRANSFORM_WH SET AUTO_SUSPEND = 120;
```

Meta's ML Feature Engineering Pipelines use AUTO_SUSPEND = 60 to avoid billing for idle time between batch jobs.

## 4. Scaling Policy: STANDARD vs ECONOMY

- · **STANDARD**: Spins up new clusters immediately if needed.
- · **ECONOMY**: Delays scaling slightly to save cost, ideal for non-urgent workloads.

Use SCALING_POLICY = 'ECONOMY' for scheduled tasks like nightly rollups.
Caution: ECONOMY may cause queueing if concurrency suddenly spikes.

# Best Practices

| Scenario | Best Warehouse Setup |
|---|---|
| Real-time streaming | XS or S + auto-resume + single cluster |
| High concurrency dashboards | MED or L + multi-cluster (2–4) |
| Nightly batch jobs | M or L + auto-suspend + economy policy |
| Development/testing | XS + suspend in 60s |

# Tech example: Optimizing for YouTube Analytics

YouTube's data team runs hourly user engagement aggregations and ad click reports.

| Table | Volume | Concurrency | Recommended WH Config |
|---|---|---|---|
| video_watch_log | 1B+/day | Medium | L + auto-suspend + 2-cluster |
| ad_click_events | 500M/day | High | M + multi-cluster (min 2, max 5) |
| metadata_staging | Low | Low | XS, suspend in 60s |

## Summary

Snowflake offers warehouse configurations that adapt to diverse data workloads:
- · Right-sizing ensures optimal performance without waste
- · Multi-cluster enables concurrency without delay
- · Suspend/resume features keep cost in check

Treat your Snowflake warehouses like tuned engines. Not every job needs a Ferrari—sometimes a hybrid gets the job done faster and cheaper.

Happy Learning
Regards
Saransh Jain