

**Snowflake SnowPro Certification - Dustin Liu**

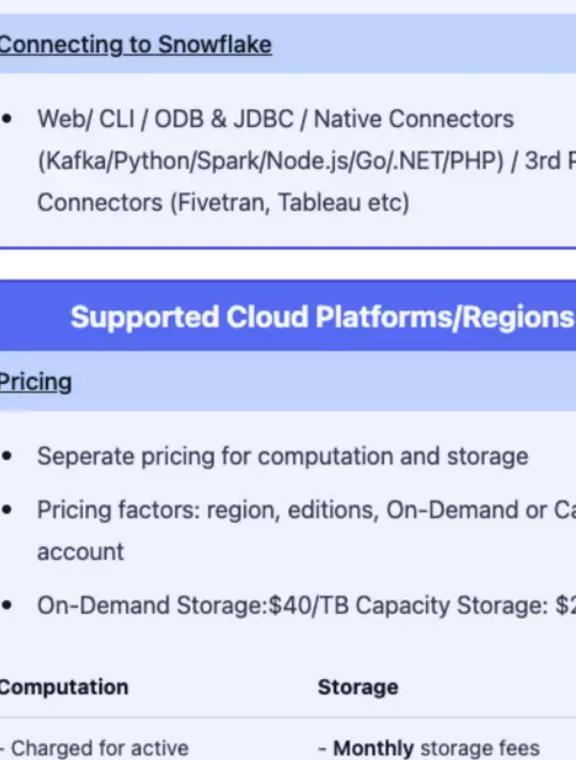
By lucas.benam@gmail.com

**Key Concepts & Architecture****Data Platform as a Cloud Service**

- OLAPaaS
- New SQL query based on public cloud with no-premises option
- Managed Service
- decoupled storage and computation

**Architecture**

A hybrid of traditional shared-disk and shared-nothing database architectures

**Storage**

• Stored managed by Snowflake

• Data objects are not directly visible (but through SQL)

**Query Processing**

• Using "virtual warehouses" to process queries

• Each VWH is an independent MPP/Parallel Processing compute cluster

**Cloud Services**

• Services that coordinate activities across Snowflake

• Authentication, Infrastructure management, Metadata management etc.

from the cloud provider(NWS/Azure/GCP)

• Query planning and optimization

**Connecting to Snowflake**

• Web (GL) JDBC / ODBC / Native Connectors

(Kafka/Python/Park/Node.js/SQL/NET/HP) / 3rd Party Connectors (Fivetran, Tableau etc)

**Supported Cloud Platforms/Regions****Pricing**

• Separate pricing for computation and storage

• Pricing factors: regions, editions, On-Demand or Capacity account

• On-Demand Storage:\$40/TB Capacity Storage: \$23/TB

**Computation****Storage**

- Charged for active

warehouses per hour

- Based on size of the

warehouse

- Cost calculated after

compression

- Cloud Providers

- Billed by second (minimum of 1 min)

- Charged in Snowflake credits

**Regions**

• North America, Europe, Asia Pacific

• Supported platforms: AWS/GCP

• Regions(Data geographically)Compute resources geographically provisioned;

• AWS (North America - EU-4, EU-5, Asia Pacific-5)

• Supported platforms: AWS/GCP

• Regions(Data geographically)Compute resources geographically provisioned;

• AWS (North America - EU-4, EU-5, Asia Pacific-5)

• Azure (North America - EU-2, Asia Pacific-3)

• GCP (North America - EU-2, Asia Pacific-3)

• Excluding Government Regions

• Across region account is not supported, single account for each region.

**Organizations**

An organization is a first-class Snowflake object that links the accounts owned by your business entity

• By default, the maximum number of accounts in an organization cannot exceed 25.

• Each account in your organization can have its own set of users, roles, databases, and warehouses - billed for usage at all of your accounts on a single bill.

**Virtual Warehouses****Overview of Warehouses**

• Resources provision for performing SELECT and DMA(DELETE, INSERT, UPDATE, COPY INTO)

• Started and stopped at anytime, resize at anytime, running queries not affected but only new queries

• Auto-suspension and Auto-resumption are enabled by default and configurable

10 t-shirt sizes with X-large as default:

• Impact on credit usage and billing

**Impact on Loading**

larger size doesn't always

improve data loading performance

• Impact on Query Processing

The approach to choose size: starting with a smaller size -&gt; monitor cost/performance -&gt; resize if needed

The number of queries a warehouse can concurrently process is determined by the size and complexity of each query

**Multi-cluster Warehouses - Scale Out**

Only available starting with Enterprise Edition

• Up to 10 warehouses

• Support all the same properties and actions as single warehouses

**Mode:**

• Maximized - same value for both max and min # of warehouses (value must be larger than 1)

• Effective for statically controlling the available compute resources

• Auto-scale: different max and min # of warehouses

• Based on scaling policy

**Scaling Policy**

Strategy Cluster Starts Cluster Shuts Down

Standard/default: - Favoring starting strategy

- Immediately when a query is submitted

- Successive queries are queued or delayed

- Only fast scaling is currently running

- faster than the query

- longer warehouse waiting time

- busier for at least 6 minutes

**Scale-up vs Scale-out**

Scale-up: Resizing to a larger warehouse is to accommodate more complex querying workload

Scale-out: To accommodate more concurrent queries/consumption. Size and complexity of query determines the concurrency (also # of queries is recommended)

**Caching****Snowflake Cache Layers**

Snowflake has three types of caching to optimize performance. Select the three types of caches from the list:

Metadata, Warehouse and Results

**Snowflake Caching**

Result Cache: Which holds the results of every query

executed in the past 24 hours and available across virtual warehouses. (result cache can last for 31 days maximum)

Metadata Cache: Improve cold time for queries against commonly used tables.

**Virtual Local Disk Cache**

Which is used to cache data used by SQL queries, data retrieved from the Remote Disk storage, and cached in SSD memory.

**Storage-layer:**

Remote Disk: Which holds the long term storage. This level is responsible for data resilience

**Release****About\_Release**

• Release frequency: Weekly

• During a release, new customer requests/questions/connections transparently move over to the newer version

• Types: full, patch and behavior change release

• Staged Release Process: Three stages - early access for Association members &amp; preview for selected users

**Edition/Service Plan Comparison**

The Snowflake Edition that your organization chooses determines the unit costs for the credits and the data storage you use.

Standard Enterprise Business Critical Virtual Private

\$1.70/Credit \$4/Credit \$5-A/Credit Contact Snowflake

Note: The cost here is an estimation only

**Monitoring****Only visible to ACCOUNTADMIN role by default**

(Website UI -&gt; Account-Resource Monitors), can be set on either account or warehouse level.

• Monitor credit usage by user-managed VW and VW used by cloud services

• Create by ACCOUNTADMIN only but can enable other roles to view/modify

• Limit can be set for a specified interval or date range:

• Effectively for statically controlling the available compute resources

• Auto-scale: different max and min # of warehouses

• Based on scaling policy

**Approach****Continuous Loading****Continuous Loading**

• Most frequent method - Designed to load small volumes of data

- Loading from stages

- COPY command

- Transformations possible

- Snippet (Serverless feature)

**Files:****External Stage Internal Stage**

- External cloud provider - Local storage maintained by Snowflake

- GCP - Snowflake

- Azure - Snowflake

- Database object created in

**File types:**

• Structured: CSV (the most performant for loading), TSV, ETC

• Semi-structured: JSON, Avro, ORC, Parquet, XML, Compressed:

• ungzipped files - automatically compressed by gzip/parquet

• bzp2, deflate, raw\_deflate

• Snippet: automatically detect the following compression format: Broth, Zstandard

**Encryption:**

• Encrypted File: Auto encryptions using 128-bit keys

or 256 bit keys (additional configuration)

• Encrypted File: key needs to be provided to SF

**Planning a Data Load**

Loading can affect query performance.

Dedicating separate warehouses for loading and querying operations to optimize performance.

• Normally, a smaller warehouse (Small, Medium, Large) would be sufficient.

• To optimise the parallel loading, produce data file roughly 100 - 250MB, compressed is recommended

• Loading very large files (e.g. 100 GB or larger) is not recommended.

**Staging Data**

Both internal and external references can include a path.

SF recommendations:

• Partitioning the data into logical path (such as geographical location, source identifiers, date)

• Organising data files by path will take advantage of parallel operations.

**Loading Data**

• By path (internal stages) / prefix (Amazon S3 bucket)

• Specifying a list of specific files to load (fastest approach)

• max: 1000 files

• The max file size default: 16MB

• Using pattern matching to identify specific files by pattern (slower approach)

• Applied differently for bulk loads vs Snowpipe loads

• copyOptions parameter:

• Validation Mode: RETURN\_N\_ROWS, RETURN\_ERRORS, RETURN\_ALL\_ERRORS - validate errors without actually loading into table

• ON\_ERROR\_N - actions to follow

• Removing data from stage once loading is completed (REMOVE command)

• Loading semi-structured files

• as type VARIANT

• compounded values can be flattened into rows using FLATTEN function

• Transformation along loading is supported but NOT for all functions

**Comparison**

The Snowflake Edition that your organization chooses

determines the unit costs for the credits and the data storage you use.

Standard Enterprise Business Critical Virtual Private

\$1.70/Credit \$4/Credit \$5-A/Credit Contact Snowflake

Note: The cost here is an estimation only

**Monitoring****Only visible to ACCOUNTADMIN role by default**

(Website UI -&gt; Account-Resource Monitors), can be set on either account or warehouse level.

• Monitor credit usage by user-managed VW and VW used by cloud services

• Create by ACCOUNTADMIN only but can enable other roles to view/modify

• Limit can be set for a specified interval or date range:

• Effectively for statically controlling the available compute resources

• Auto-scale: different max and min # of warehouses

• Based on scaling policy

**Approach****Continuous Loading**

• Most frequent method - Designed to load small volumes of data

- Loading from stages

- COPY