

Hairvana

Cloud Data Pipeline

Shuhan Chang
Chemseddine Nadour
Cyprien Cazenave
Mathieu Chabirand
Vanessa Gunathasan



Agenda

- ❑ Introduction
- ❑ Cloud Data Pipeline
- ❑ Business Impact
- ❑ Conclusion

[Presenter's
First Name]

Data Ingestion and Storage

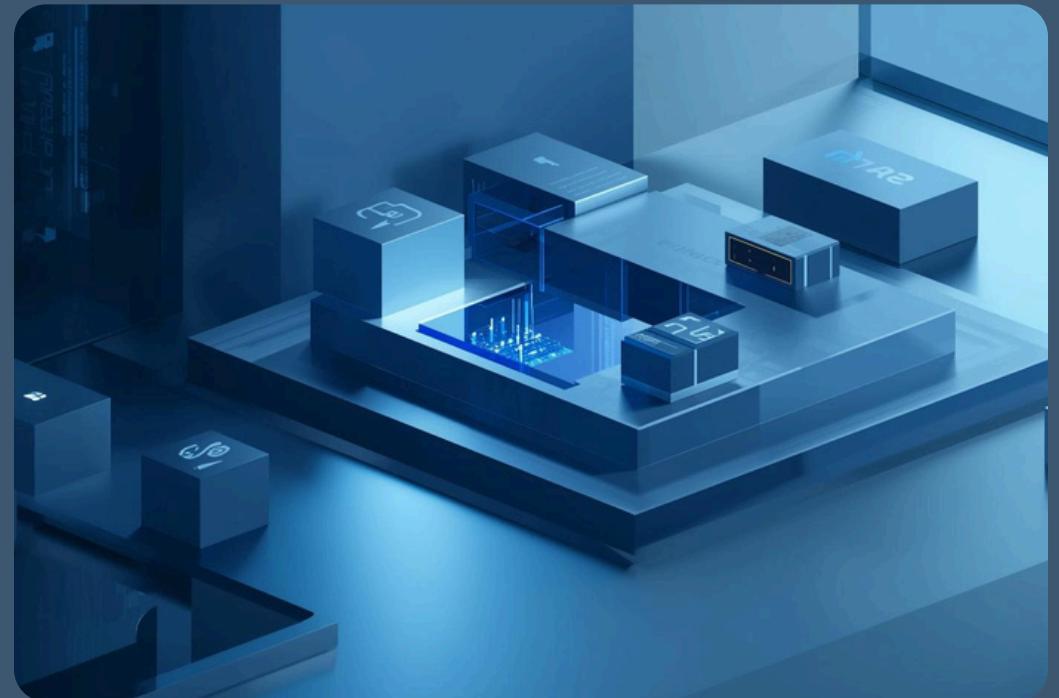
Sources

Integrating diverse data sources enhances insights for Hairvana's haircare solutions.



Data Lake

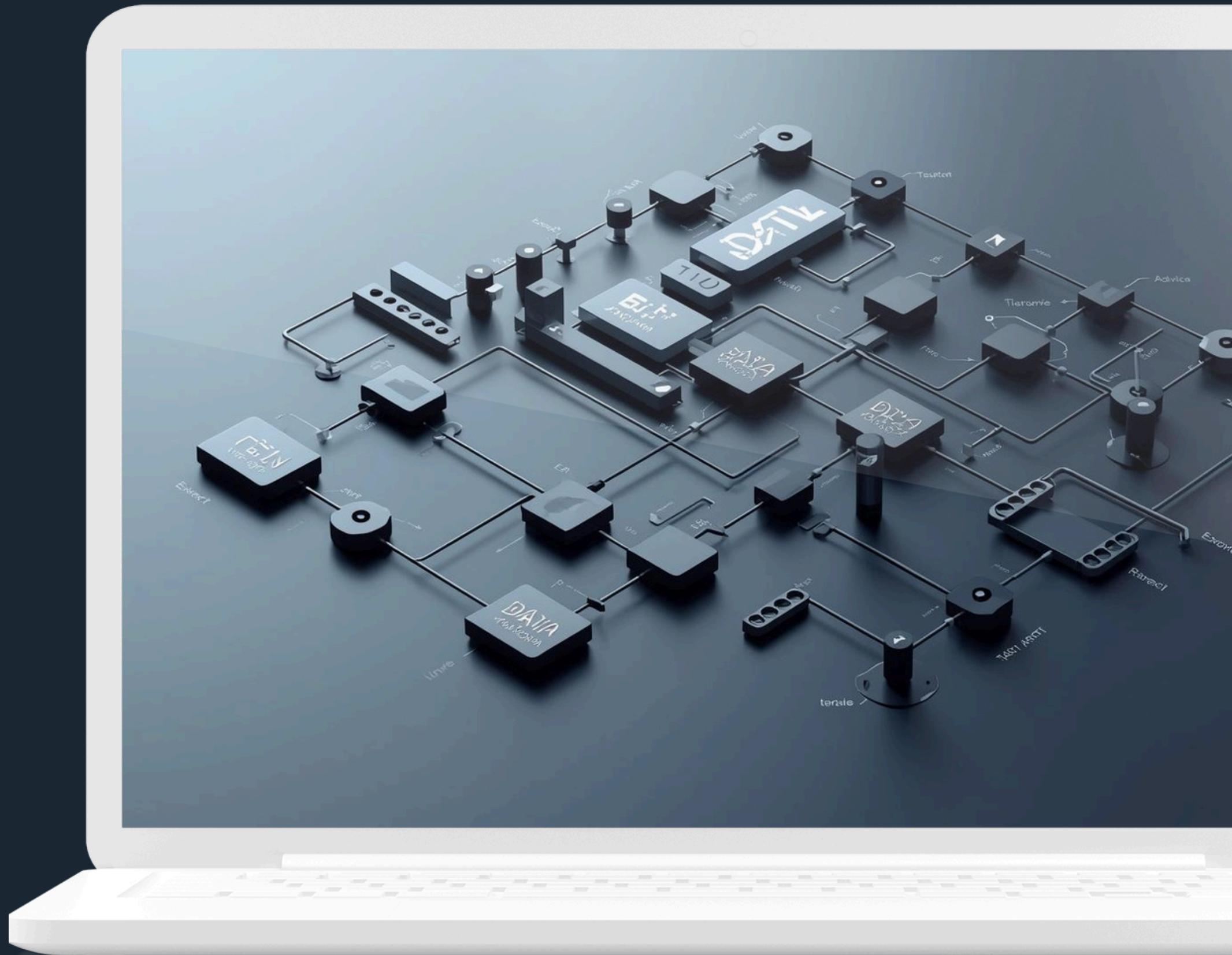
Centralized storage streamlines access and management of vast customer data.



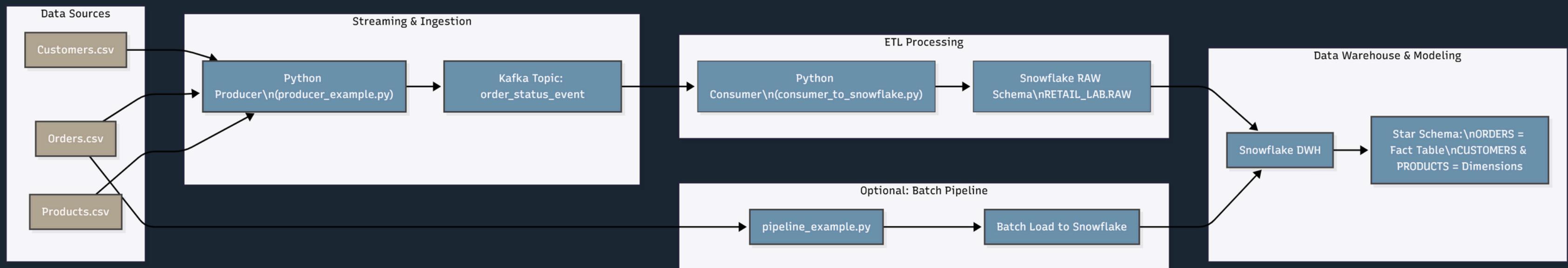
DATA PROCESSING

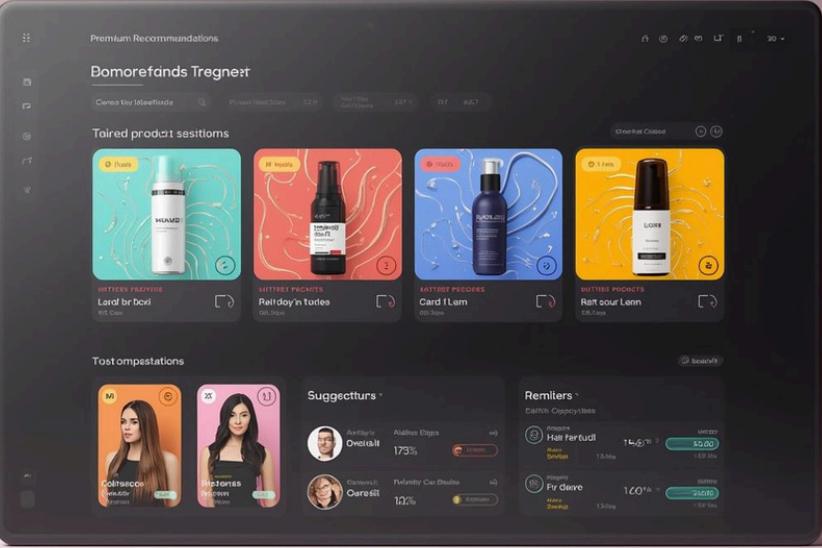
ETL and ELT Explained

This section explores the **data transformation pipeline**, detailing the processes of ETL and ELT that enable Hairvana to derive actionable insights from vast amounts of data.



Our Pipeline Architecture





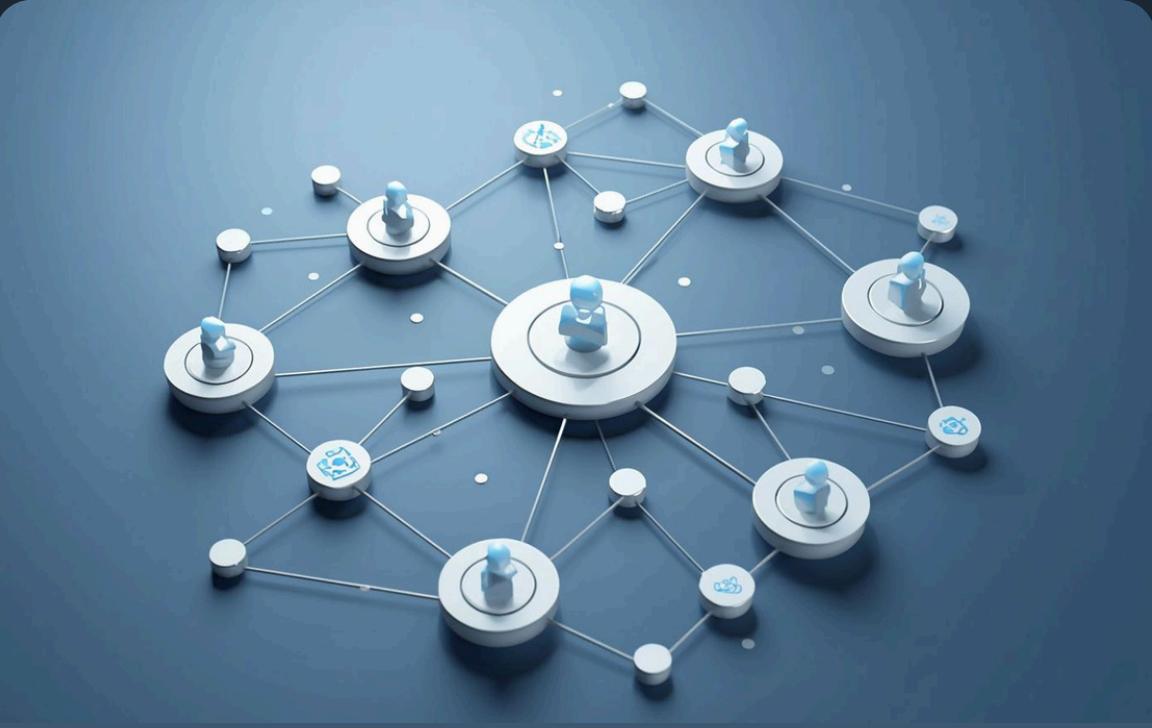
Tailored Recommendations for Individual Customers

Discover Your Perfect Match



Streamlined Inventory for Enhanced Efficiency

Optimize Your Stock Levels



Data-Driven Experience for Every Customer

Elevate Your Service Today

Real-time Insights for Premium Haircare

Driving Marketing Efficiency

Leverage real-time data to tailor marketing strategies and optimize campaign performance for maximum impact.

Enhancing Product Development

Utilize data-driven insights to refine product offerings and develop new formulations that meet customer needs.

Spotting Future Trends

Implement predictive analytics to anticipate market shifts and adapt strategies for emerging customer preferences.

Data-Driven Growth

Foster a culture of data-driven decision making to continuously improve operations and boost overall business performance.

Data Structures

Data Structures

After processing, raw data is organized into clean, structured tables within our Snowflake data warehouse. This makes it easy to query and analyze. Select a data model below to see its structure.

[ORDERS](#)[CUSTOMERS](#)[PRODUCTS](#)

ORDERS

This is the main fact table containing all transactional data. It is optimized for analytical queries on sales performance and trends.

Column Name	Data Type	Description
ORDER_ID	INTEGER	Unique identifier for each order.
CUSTOMER_ID	INTEGER	Foreign key to the Customers table.
PRODUCT_ID	INTEGER	Foreign key to the Products table.
QUANTITY	INTEGER	Number of units sold.
ORDER_TOTAL	DECIMAL(10, 2)	Total value of the order.
ORDER_STATUS	VARCHAR	Current status of the order (e.g., Shipped, Delivered).
SOLD_AT	TIMESTAMP	Date and time of the transaction.
REGION	VARCHAR	Market region where the order occurred.

CUSTOMERS

This dimension table holds all information about our customers, allowing for segmentation and behavioral analysis.

Column Name	Data Type	Description
CUSTOMER_ID	INTEGER	Unique identifier for each customer.
NAME	VARCHAR	Full name of the customer.
EMAIL	VARCHAR	Customer's contact email.
JOIN_DATE	DATE	Date the customer first registered.
ADDRESS	VARCHAR	Street address.
CITY	VARCHAR	City of residence.
COUNTRY	VARCHAR	Country of residence.
POSTAL_CODE	VARCHAR	Postal code.
REGION	VARCHAR	Customer's market region (e.g. Europe, Asia, Americas).

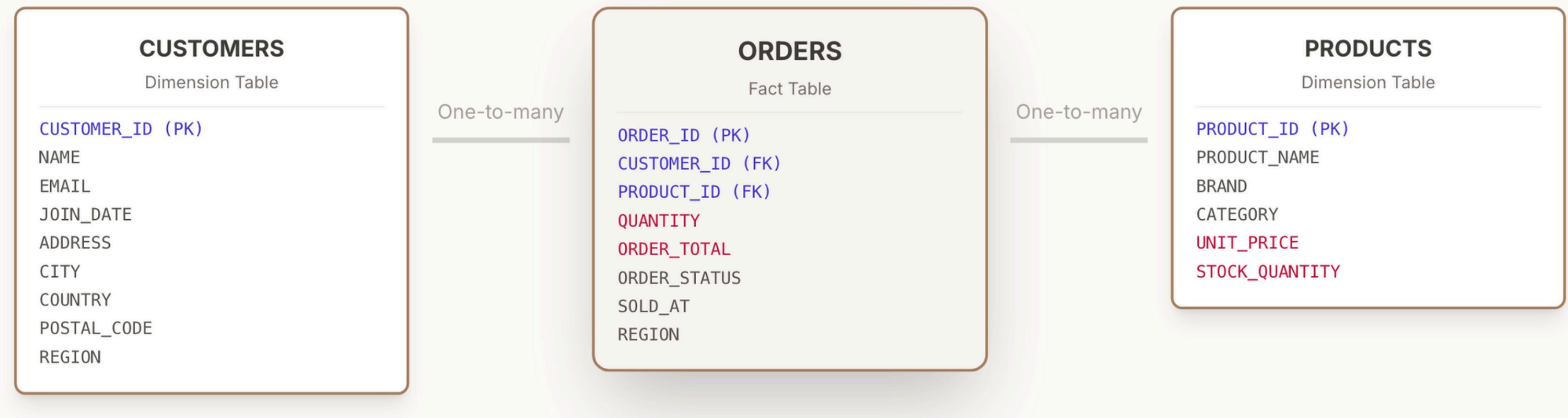
PRODUCTS

This dimension table contains details for every product in our catalog, used for inventory and sales analysis.

Column Name	Data Type	Description
PRODUCT_ID	INTEGER	Unique identifier for each product.
PRODUCT_NAME	VARCHAR	Name of the product.
BRAND	VARCHAR	Brand name of the product.
CATEGORY	VARCHAR	Product category (e.g., Shampoo, Serum).
UNIT_PRICE	DECIMAL(10, 2)	Retail price of one unit.
STOCK_QUANTITY	INTEGER	Current number of units in stock.

Star Model

We model our data using a star schema. This design features a central "fact" table ('ORDERS') containing quantitative data, connected to descriptive "dimension" tables. This structure simplifies queries and improves performance for reporting and analytics.



Git hub

https://github.com/shuhanchang12/ETL_ecommerce

Screenshots of snowflake tables to validate that you ingested the data.

Database Explorer

HORIZON CATALOG

Databases

Search

ETL_DB

RETAIL_LAB

INFORMATION_SCHEMA

PUBLIC

RAW

Tables

CUSTOMER

ORDER_STATUS_EVENTS

PRODUCER_STATUS

SNOWFLAKE

SNOWFLAKE_LEARNING_DB

SNOWFLAKE_SAMPLE_DATA

TEST_DB

USER\$SCHANG

RETAIL_LAB / RAW / ORDER_STATUS_EVENTS

Table ACCOUNTADMIN 26 minutes ago 324 16.0KB

Table Details Columns Data Preview Copy History Data Quality PREVIEW Lineage

• COMPUTE_WH 100 of 324 Rows • Updated just now

	EVENT_ID	ORDER_ID	CUSTOMER_ID	NEW_STATUS	STATUS_TS	SOURCE
1	6279b282-cf22-4bd6-a5b5-51f1bc49c77f	11	74	CREATED	2025-10-19T18:55:01.261+0000	order_service
2	e83b94c2-8d78-4068-a783-ad505a96d1b8	11	74	PAID	2025-10-19T18:55:01.784+0000	order_service
3	9d2177cc-4a52-48f0-bb45-7f4c12593cc5	11	74	PACKED	2025-10-19T18:55:02.288+0000	order_service
4	fb95f7c-98ec-4a3a-9157-f9f2cea7a63a	11	74	SHIPPED	2025-10-19T18:55:02.793+0000	order_service
5	13c68753-0b52-419c-ad70-fe08fad7385e	12	99	CREATED	2025-10-19T18:55:03.374+0000	order_service
6	d1094a6b-d4f0-4063-8e21-494873cdb8ee	12	99	PAID	2025-10-19T18:55:03.895+0000	order_service
7	d219f137-1ce6-420e-b05e-eee9a59eabac	12	99	PACKED	2025-10-19T18:55:04.431+0000	order_service
8	c8c365d2-513e-48b4-b0fb-d6846da9ade5	12	99	SHIPPED	2025-10-19T18:55:04.953+0000	order_service
9	fee7082a-0188-4188-829d-4899da65d790	13	49	CREATED	2025-10-19T18:55:05.575+0000	order_service
10	b08d6658-781f-47ae-9e6b-6b73c06b9f7d	13	49	PAID	2025-10-19T18:55:06.085+0000	order_service
11	23b9f49a-bb5d-4a48-9c24-bab0f6c3386f	13	49	PACKED	2025-10-19T18:55:06.661+0000	order_service
12	c7c63a5b-c776-4209-8194-c21935062d94	13	49	SHIPPED	2025-10-19T18:55:07.167+0000	order_service
13	11d5b766-4114-47de-92ba-62dc1a2ddf0d	13	49	DELIVERED	2025-10-19T18:55:07.717+0000	order_service
14	acc66b16-df73-44cb-b565-1386bf1e67a1	13	49	CANCELLED	2025-10-19T18:55:08.361+0000	order_service
15	b8742fd2-0c54-43b0-9da7-1d2291582e99	14	48	PAID	2025-10-19T18:57:30.604+0000	order_service

The data streaming with Redpanda-events created,

Redpanda

order_status_event Refreshing in 4 secs Debug bundle

Size: 89.2 kB | Estimated messages: 355 | Cleanup Policy: Delete | Retention: ~7 days or Infinite

Produce Record | Delete Records

Topics Schema Registry Consumer Groups Security Quotas Connect Transforms Reassign Partitions

Messages Consumers Partitions Configuration ACL Documentation

START OFFSET MAX RESULTS
Newest - 50 50 Add filter

Filter table content ...

TIMESTAMP	KEY	VALUE
2025/10/19 下午8:57:19	Null	{"event_id": "cf5b71e4-aadd-4d9f-af2f-e69100dd0610", "order_id": 12, "customer_id": ..., "status": "PENDING", "timestamp": "2025-10-19T20:57:19Z"} JSON - 188 B
2025/10/19 下午8:57:20	Null	{"event_id": "dd954180-d26d-4e43-861f-608f6ed0ad38", "order_id": 12, "customer_id": ..., "status": "PENDING", "timestamp": "2025-10-19T20:57:20Z"} JSON - 185 B
2025/10/19 下午8:57:21	Null	{"event_id": "1c8b43ac-d7f2-41a4-9493-783c615da49a", "order_id": 12, "customer_id": ..., "status": "PENDING", "timestamp": "2025-10-19T20:57:21Z"} JSON - 187 B
2025/10/19 下午8:57:21	Null	{"event_id": "99cd9bc4-cbca-4341-9175-04422963c3e0", "order_id": 12, "customer_id": ..., "status": "PENDING", "timestamp": "2025-10-19T20:57:21Z"} JSON - 188 B
2025/10/19 下午8:57:22	Null	{"event_id": "e63f8499-a798-4d9f-b550-32aa73386746", "order_id": 12, "customer_id": ..., "status": "PENDING", "timestamp": "2025-10-19T20:57:22Z"} JSON - 190 B
2025/10/19 下午8:57:22	Null	{"event_id": "3abd96cf-aec6-4a3c-b0fd-987cc0575956", "order_id": 13, "customer_id": ..., "status": "PENDING", "timestamp": "2025-10-19T20:57:22Z"} JSON - 188 B
2025/10/19 下午8:57:23	Null	{"event_id": "c3c94096-35c6-4aa9-af95-f6a57113e791", "order_id": 13, "customer_id": ..., "status": "PENDING", "timestamp": "2025-10-19T20:57:23Z"} JSON - 185 B
2025/10/19 下午8:57:23	Null	{"event_id": "60f29be2-9d97-4912-9912-608fb5d049e9", "order_id": 13, "customer_id": ..., "status": "PENDING", "timestamp": "2025-10-19T20:57:23Z"} JSON - 187 B
2025/10/19 下午8:57:24	Null	{"event_id": "ee5e4773-d4f6-4be5-9d09-78ef31d76e43", "order_id": 13, "customer_id": ..., "status": "PENDING", "timestamp": "2025-10-19T20:57:24Z"} JSON - 188 B

9.14 kB 813ms

Collapse sidebar

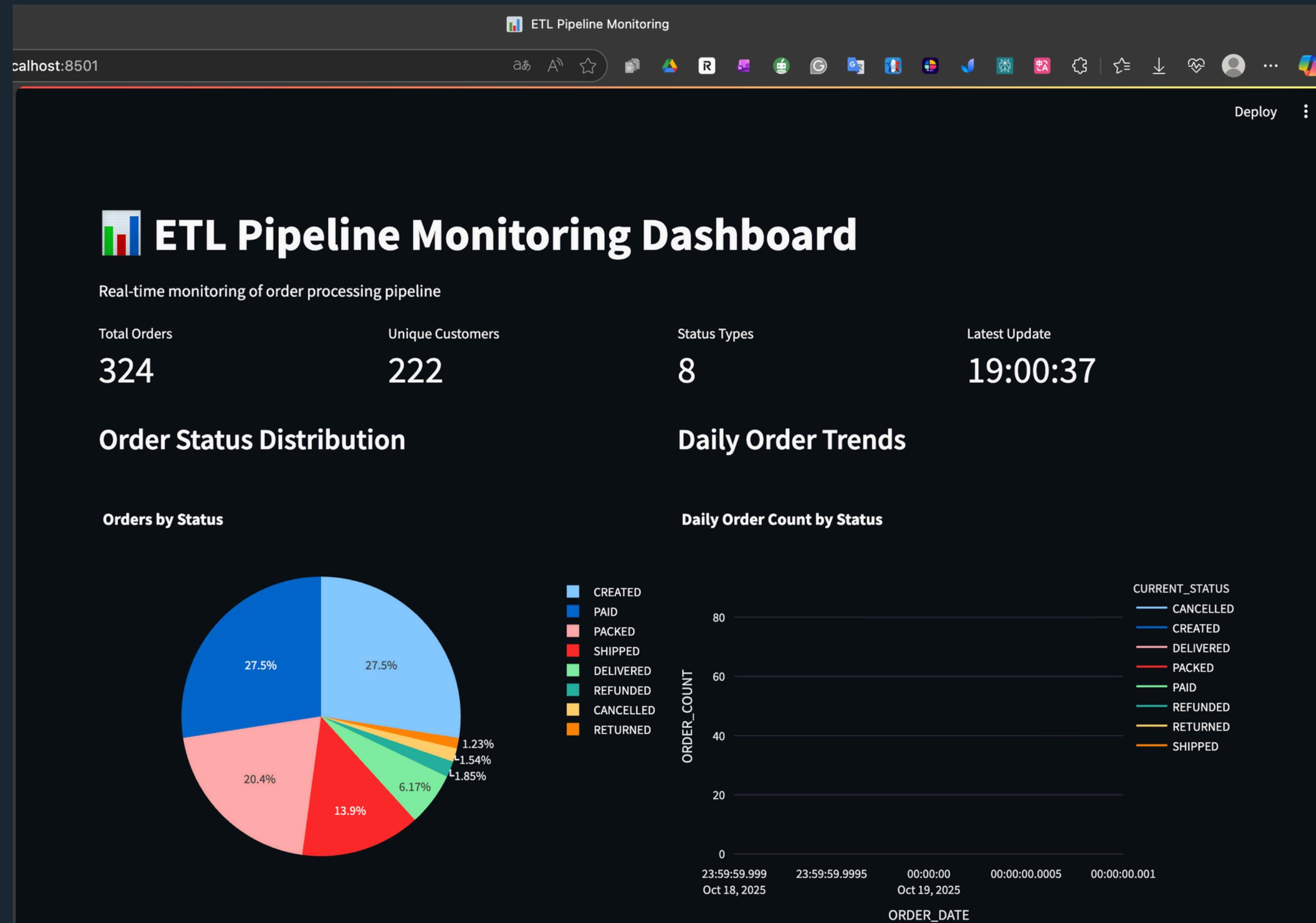
Ingestion in Snowflake

The screenshot shows the Snowflake interface with the following details:

- Left Sidebar:** Shows the navigation menu under "Work with data" with "Projects" selected. Other options include Ingestion, Transformation, AI & ML, Monitoring, Marketplace, Catalog, Data sharing, Governance & security, Compute, Admin, and Manage.
- Top Bar:** Displays the date and time (2025-10-19 8:27pm), user (ACCOUNTADMIN), and workspace (TEACH_WH (X-Small)).
- Central Area:** A code editor window titled "RETAIL_LAB.RAW" containing SQL code for creating tables and inserting data. The code includes:

```
68     CREATED_AT    TIMESTAMP_NTZ
69 );
70 SHOW TABLES IN RETAIL_LAB.RAW;
71 CREATE TABLE IF NOT EXISTS RETAIL_LAB.RAW.PRODUCER_STATUS (
72     EVENT_ID      STRING,
73     PRODUCER_ID   NUMBER,
74     CUSTOMER_ID   NUMBER,
75     NEW_STATUS    STRING,
76     STATUS_TS     TIMESTAMP_NTZ,
77     SOURCE        STRING
78 );
79
80 CREATE TABLE IF NOT EXISTS RETAIL_LAB.RAW.ORDER_STATUS_EVENTS (
81     EVENT_ID      STRING,
82     ORDER_ID      NUMBER,
83     CUSTOMER_ID   NUMBER,
84     NEW_STATUS    STRING,
85     STATUS_TS     TIMESTAMP_NTZ,
86     SOURCE        STRING
87 );
88
89 SELECT COUNT(*) FROM RETAIL_LAB.RAW.ORDER_STATUS_EVENTS;
90 SELECT * FROM RETAIL_LAB.RAW.ORDER_STATUS_EVENTS LIMIT 20;
```
- Results Tab:** Displays the results of the query in a table format. The table has columns: EVENT_ID, ORDER_ID, CUSTOMER_ID, NEW_STATUS, STATUS_TS, and SOURCE. The data shows 20 rows of order status events.
- Query Details:** Shows the duration (1.1s), number of rows (20), and query ID (01bfd211-0001-3bc5-00...).
- Event ID Distribution:** A histogram showing the distribution of Event IDs, with 100% filled.
- Order ID Distribution:** A histogram showing the distribution of Order IDs, with 11 to 16.
- Customer ID Distribution:** A histogram showing the distribution of Customer IDs, with 32 to 99.

monitoring APP–Stramlit



Data Governance & Security

Protecting Hairvana's Data Assets with Enterprise-Grade Solutions

Access Control & Authentication

- Role-Based Access Control (RBAC) in Snowflake
- Multi-Factor Authentication (MFA)
- Least Privilege Principle enforced
- User access auditing
- Automatic credential rotation

Snowflake RBAC

OAuth 2.0

Encryption & Data Protection

- End-to-end encryption (TLS 1.3)
- Data at rest encrypted (AES-256)
- Sensitive data masking (PII)
- Anonymization for dev/test environments
- Secure key management (KMS)

AES-256

TLS 1.3

Data Quality & Lineage

- Automated validation of ingested data
- Complete traceability (source → destination)
- Data lineage visualized in Snowflake
- Quality checks at each pipeline stage
- Alerts for detected anomalies

Snowflake Streams

Tasks

Monitoring & Audit

- Centralized logs of all activities
- Real-time Kafka pipeline monitoring
- Automatic incident alerts
- 24/7 surveillance dashboards
- Quarterly audit reports

Snowflake Query History

Monitoring App

Hairvana complies with international data protection standards

GDPR

SOC 2 Type II

ISO27001

CCPA

**COMPLIANCE &
REGULATIONS**

Data Retention: 7 years for transactional data
Right to be Forgotten: Guaranteed deletion within 30 days
Backup: Daily snapshots with 90-day retention

Merci