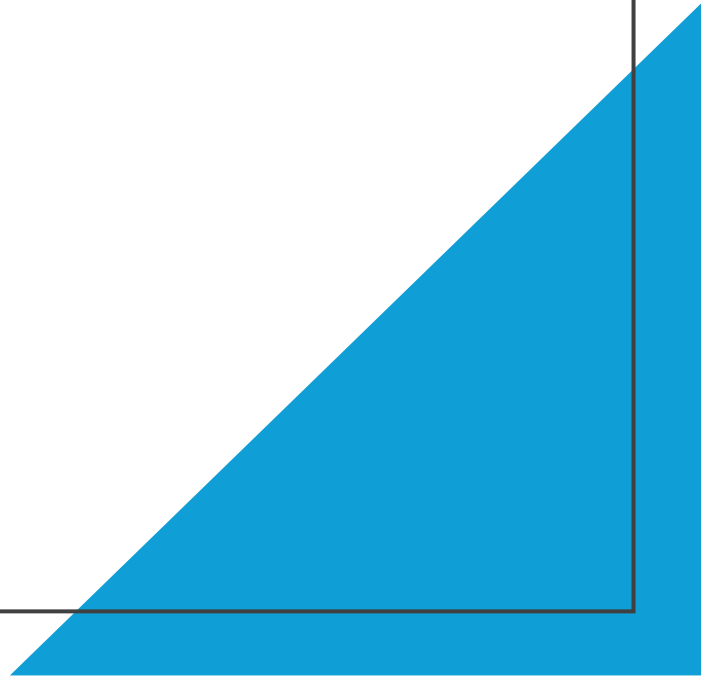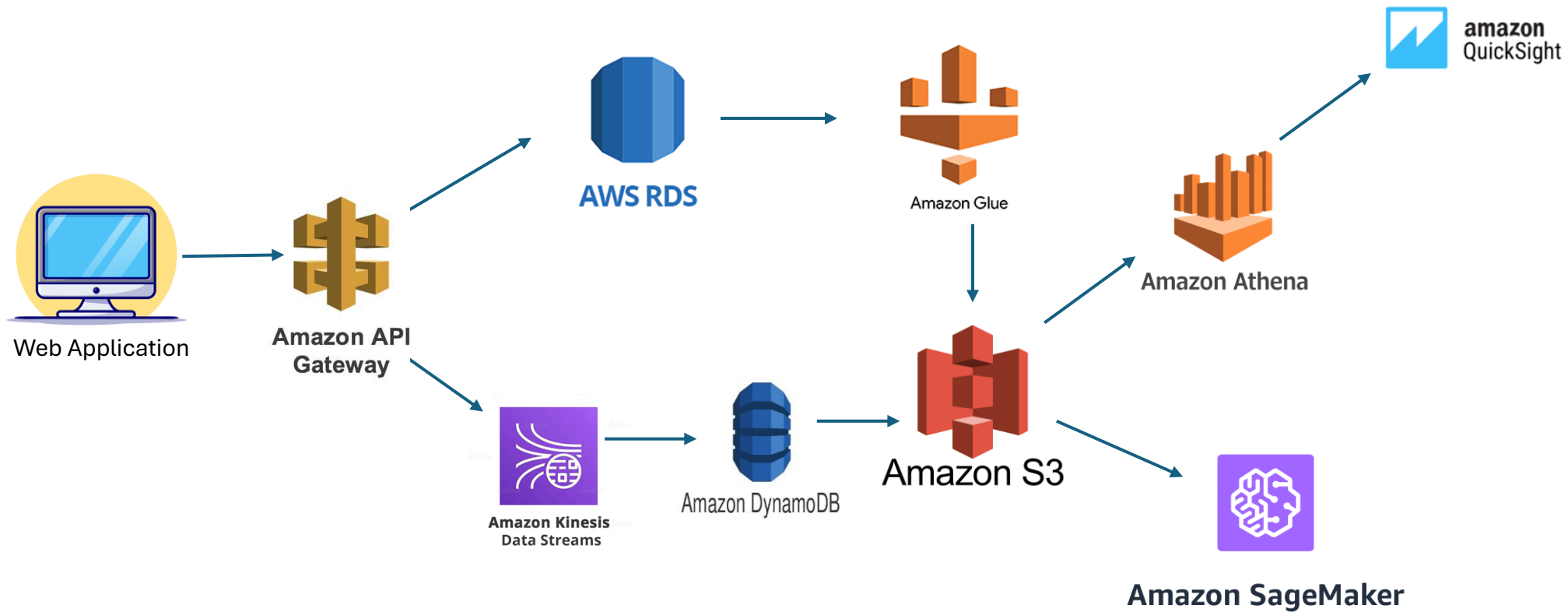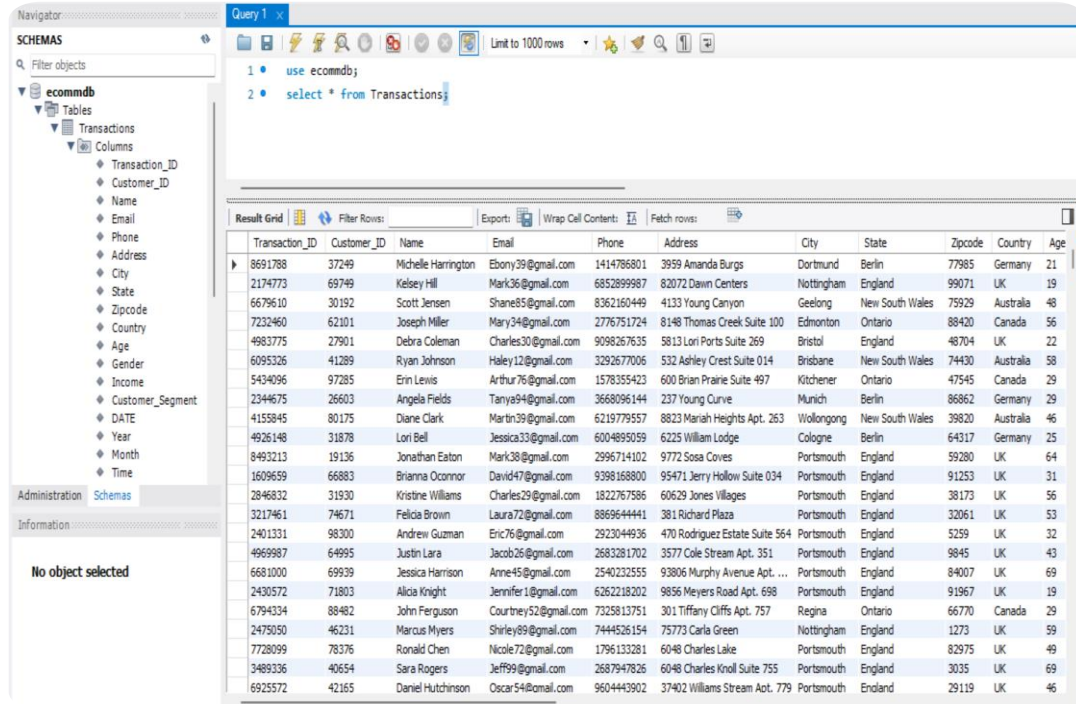# Data Engineering Ecommerce End to End Project

The E-commerce Data Pipeline and Analytics project is an end-to-end solution leveraging AWS (Kinesis, Glue, S3, Athena, QuickSight, SageMaker) to ingest, transform, and analyze transactional and clickstream data. It delivers real-time insights, dashboards, and personalized recommendations, enhancing customer engagement and business decisions.

# Process Flow Chart

# AWS RDS



- The web application sends transactional data to secure RESTful endpoints created in AWS API Gateway, ensuring reliable data reception.

- These endpoints are integrated with AWS Lambda functions, facilitating the data processing.

- The processed data is seamlessly and securely loaded into Amazon RDS.

# AWS GLUE

• **Extract**: Data is extracted from the MySQL database containing e-commerce transactions.

•**Transform**:

- Drop duplicates to maintain data quality.
- Drop unnecessary fields to reduce data size.
- Remove null fields to ensure completeness.
- Change schema to standardize data structure.
- Apply custom SQL queries for complex transformations.

•**Load**: The cleaned and transformed data is loaded into Amazon S3, ready for further use.

# S3 Data Lake

- The primary objective of the S3 data lake is to consolidate various data sources into a single, scalable storage solution. This allows for:

- **Enhanced Data Integration**: Combining transactional and clickstream data for a unified view.

- **Improved Data Accessibility**: Making data easily accessible for analytics and machine learning.

- **Scalable and Cost-Effective Storage**: Leveraging the scalability and cost-efficiency of Amazon S3.

# AWS Athena

AWS Glue Data Catalog was used to create and manage metadata for all transactional data stored in S3.

Glue Crawlers were configured to automatically scan the S3 bucket, infer the schema of new or modified data, and update the catalog

Amazon Athena was used to perform interactive SQL queries on the data stored in S3, leveraging the metadata defined in the Glue Data Catalog.

Athena's integration with Glue allowed it to utilize the schema information for efficient query execution. This enabled quick analysis of large datasets without the need to manage any infrastructure.

# AWS Quicksight

- Count by Shipping Method and Customer Segment:
  - Regular customers have the most transactions, especially with Standard shipping.
  - Premium customers follow, with New customers having the least transactions.

- Transactions per Payment Method:
  - Credit Card is the most used payment method, followed by Debit Card and PayPal.
  - Cash transactions are the least common.

# AWS Quicksight

- Monthly Sales by Product Category:
  - Sales peak in December, especially for Electronics, Grocery, and Clothing.
  - Home Decor and Books have lower sales.

- Sales per Product Type:
  - Top-selling items: Smartphones, Televisions, Furniture.
  - Notable sales in consumables like Water and Juice.
  - Lower sales in Shorts and Thriller books.

- Sum of Sales per Country:
  - USA leads in total sales, followed by the UK and Germany.
  - Australia and Canada have the least sales.

- Sales by Brand:
  - Leading brands: Nike, Adidas, Sony, Samsung.
  - Strong performance in beverages (Coca-Cola, Pepsi) and books (Penguin Books, HarperCollins).
  - Consistent but moderate sales for retail brands (Bed Bath & Beyond, Zara).

# Kinesis

- Real-time Ingestion and Scalability: Collects and ingests real-time clickstream data from the web application, automatically scaling to handle high data volumes and ensuring data durability through replication across multiple Availability Zones.

- Processing and Storage: Processes clickstream data from Kinesis, generates user sessions based on a 30-minute inactivity rule, and writes the processed data to Amazon DynamoDB, scaling automatically with the data load and ensuring cost efficiency by charging only for actual processing time.

| | Name | ▲ | Status | ▽ | Capacity mode | ▽ | Provisioned shards | ▽ | Sharing policy | ▽ | Data retention period | ▽ | Encryption | ▽ | Consumers with enhanced fan-ou |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| ☐ | Ecomm_Clickstream | | ⊘ Active | | On-demand | | - | | No | | 1 day | | Disabled | | 0 |

# Sagemaker ML

- Developed Recommendation System: Leveraged Amazon SageMaker to create a recommendation system using clickstream data stored in Amazon DynamoDB.

- Data Preparation: Extracted and transformed user activity data, including session details and transaction IDs, into a format suitable for machine learning.

- Exploratory Data Analysis (EDA): Used Jupyter notebooks on SageMaker to perform EDA and visualizations, gaining insights into user behavior.

- Model Training: Trained a machine learning model on SageMaker to analyze user behavior and generate personalized recommendations.

- Real-time Recommendations: Deployed the model to provide real-time recommendations, continuously updating with new data.

- Enhanced User Engagement: Delivered tailored content to users, enhancing engagement and satisfaction with the web application.