# A Project Report on

# ETL Pipeline for Data Processing using AWS S3, PySpark, and Snowflake

**PREPARED BY:**

| ASSOCIATE ID | ASSOCIATE NAME |
|---|---|
| **2387739** | **Aryan Khokle** |

**GUIDED BY:**

| **Mentor** | **Coach** |
|---|---|
| Govindaraj, Kavin Kumar | Yalin Maria |
| Associate | CPO L&D GenC Training, |
| AIA Cloud Data Integration | Human Resource – GenC |

# Contents

# List Of Figures

# Introduction

**ETL Pipeline for Data Processing using AWS S3, PySpark, and Snowflake**

This project aims to develop a strong ETL (Extract, Transform, Load) pipeline for efficient data processing and loading into a cloud data warehouse. The pipeline retrieves raw CSV data from AWS S3, cleans and transforms it with PySpark on AWS Glue, and then loads it into Snowflake using its robust COPY command. The primary objectives are to maintain data quality by eliminating duplicates, managing null values, rearranging columns, casting data types, and truncating excessively long text values. This report details the system architecture, data flow, database design, ETL implementation, and additional aspects.

# Tools and Technologies Used

- **AWS S3:** Used as the data storage layer for both raw and processed data.
- **AWS Glue:** A managed ETL service that supports PySpark based transformations.
- **PySpark:** Used for data cleaning and transformation.
- **Snowflake:** The cloud data warehouse where the final processed data is loaded for further analytics.

# Purpose of Project

The main goal of this project is to create a robust ETL pipeline for data processing that can:

1. **Efficiently Process and Store Employee Data:**

   - The system automates the extraction, transformation, and loading of sales data from various sources stored in AWS S3. This ensures quick and accurate processing, reducing manual effort and minimizing errors. The processed data is stored in a structured format, ready for analysis.

2. **Ensure Data Quality:**

   - Using PySpark on AWS Glue, the system performs data cleaning and transformation tasks such as removing duplicates, handling null values, reordering columns, performing data type casts, and truncating overly long text values. This ensures high data quality and consistency.

3. **Enable Comprehensive Analytics and Reporting:**

   - The system supports advanced analytics by loading processed data into Snowflake. It generates detailed reports on key performance indicators, such as top-performing employees and department-wise employee performance. These insights help managers make informed decisions.

4. **Support Data-Driven Decision-Making:**

   - By providing real-time insights and historical analysis, the system empowers decision-makers to respond quickly to market changes, optimize operations, and develop effective business strategies. Data-driven decisions enhance competitiveness and drive growth.

**About the Software System**

The ETL Pipeline for Data Processing is a comprehensive solution that integrates various AWS services and Snowflake to manage sales data effectively. The system performs the following key functions:

- **Data Extraction:** Extracts raw CSV data from AWS S3 for initial storage.

- **Data Transformation:** Cleans and transforms the raw data using PySpark on AWS Glue, ensuring data quality by removing duplicates, handling null values, reordering columns, performing data type casts, and truncating overly long text values.

- **Data Loading:** Loads the transformed data into Snowflake using its powerful COPY command for efficient storage and analysis.

- **Reporting:** Generates detailed reports on key performance indicators, such as top-selling products and store-wise sales performance, to support comprehensive analytics and informed decision-making.

**Scope of the System**

1. **Data Ingestion:**

   - Automating the process of uploading employee data to AWS S3. The data is initially stored in the "starting" folder as untransformed CSV files.

2. **Data Processing:**

   - Implementing data validation, transformation, and aggregation using AWS Glue with PySpark. The transformed data is stored in the "transformed" folder in S3. The transformation process includes:

3. **Data Storage:**

   - Loading the processed data into Snowflake as a cloud warehouse. The connection setup includes:

4. **Analytics and Reporting:**

   - Providing SQL queries and reports to analyze employee data and support business decisions.

5. **Scalability:**

   - Ensuring the system can handle increasing volumes of employee data as the business grows.

6. **Security:**

   - Implementing appropriate access controls and permissions to protect sensitive employee data.
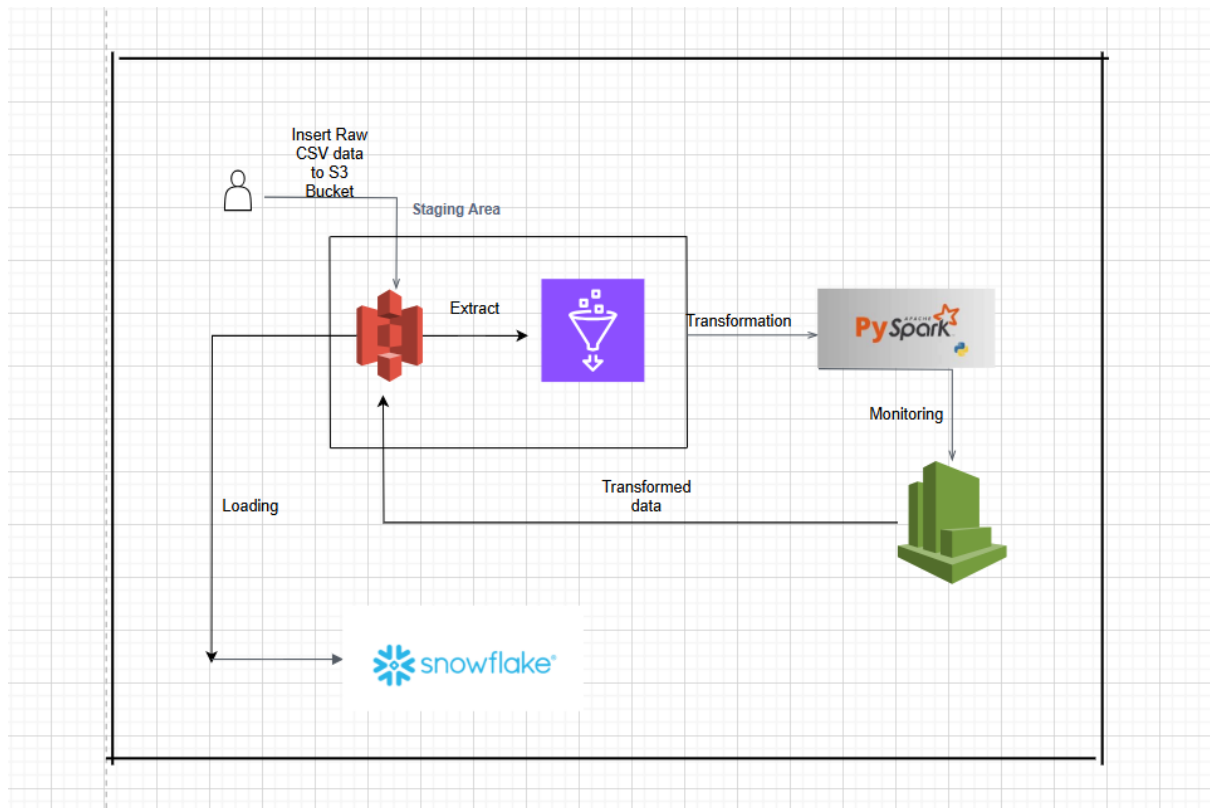
# Flow Chart

Fig 1. Flow Chart

Chapter 1: Data Extraction

## 1.1 Definition

Data extraction is the initial step in the data processing pipeline. In this phase, an employee CSV file is read from a local directory and uploaded to AWS S3 for raw data storage. This process ensures that the raw employee data is securely stored in a centralized location, making it readily available for subsequent processing and analysis.

.

## 1.2 Types of Data Extraction

Manual Extraction: This method involves manually uploading files to the storage system. While it is straightforward, it is labour-intensive and prone to human error, making it less suitable for handling large volumes of data.

Automated Extraction: This method uses scripts or tools to automate the data upload process. It is efficient, reduces the risk of errors, and can handle large volumes of data without manual intervention, making it ideal for modern data processing needs.

## 1.3 Categories of Data Extraction

Batch Extraction: This approach involves extracting data in batches at scheduled intervals. It is suitable for periodic data uploads, such as daily or weekly sales data, and helps in managing data loads effectively.

Real-Time Extraction: This approach continuously extracts data as it becomes available. It is ideal for systems that require immediate data availability, such as real-time analytics and monitoring systems.

## 1.4 Advantages

Efficiency: Automating the data upload process saves significant time and effort, allowing the system to handle large volumes of data seamlessly.

Scalability: Automated extraction can scale to accommodate increasing data volumes without requiring additional manual effort.

Reliability: By reducing the risk of human error, automated extraction ensures that data is consistently and accurately uploaded to the storage system.

## 1.5 Disadvantages

Initial Setup: Setting up automated data extraction requires initial configuration of scripts and tools, which can be time-consuming and require technical expertise.

Maintenance: Regular maintenance is needed to ensure the smooth operation of the extraction process, including updates and troubleshooting.

Dependency: The extraction process relies on the availability of the source data and network connectivity, which can be a limitation in some scenarios.

## 1.6 Examples

Employee Records Management: A company uploads daily updates of employee records, including new hires, terminations, and changes in employment status, to AWS S3. This ensures that all employee information is centrally stored and accessible for further processing and compliance.

Employee Activity Logs: A company extracts activity logs from its internal systems, such as login times, application usage, and project contributions, and stores them in a centralized location for monitoring and analysis. This helps in identifying productivity trends and troubleshooting issues promptly.

## 1.7 Steps Involved

Create S3 Bucket: Set up an S3 bucket to store raw sales data, ensuring that the bucket is properly configured for data storage and access.

Generate Policies: Define access policies for the S3 bucket to control who can read and write data, ensuring data security and compliance.

Fig 1. Raw Employee Data CSV File



Fig 2. Creating S3 bucket

Fig 3. Creating IAM role for Glue



Fig 4. Created IAM role for Snowflake



Fig 5.Creating Policy for Snowflake named policyforgluee

# Chapter 2: Data Transformation

## 2.1 Definition

Data transformation is a critical step in the data processing pipeline, where raw sales data is validated and aggregated in AWS Glue using PySpark. This process ensures data quality by removing or correcting invalid data and generating processed datasets that are ready for analysis.

## 2.2 Types of Data Transformation

- **Data Cleaning:** This involves removing or correcting invalid data to ensure accuracy. It includes tasks such as handling missing values, correcting data types, and removing duplicates.
- **Data Aggregation:** This involves summarizing data to provide insights, such as calculating total sales per store or product. Aggregation helps in reducing data complexity and making it more manageable for analysis.
- **Data Enrichment:** This involves adding additional information to the data, such as product categories or customer demographics, to enhance its value and provide deeper insights.

## 2.3 Categories of Data Transformation

- **Synchronous Transformation:** This approach processes data in real-time as it is received, ensuring that the data is immediately available for analysis. It is suitable for applications that require up-to-date information.
- **Asynchronous Transformation:** This approach processes data in batches at scheduled intervals, allowing for more efficient use of resources. It is suitable for applications that can tolerate some delay in data availability.

## 2.4 Advantages

- **Data Quality:** By validating and cleaning the data, transformation ensures that the data is accurate and consistent, which is essential for reliable analysis.
- **Insights:** Aggregating data provides valuable insights by summarizing complex data into meaningful metrics, such as total sales or average order value.
- **Automation:** Automating data transformation tasks reduces manual effort and ensures that the data is processed consistently and efficiently.

## 2.5 Disadvantages

- **Complexity:** Setting up and maintaining data transformation processes requires technical expertise and can be complex, especially for large datasets.

- **Resource Intensive:** Data transformation can consume significant computational resources, especially for large datasets or complex transformations.
- **Latency:** Depending on the approach used, data transformation may introduce delays in data availability, which can be a limitation for real-time applications.

## 2.6 Examples

- **Employee Data Aggregation:** Summarizing employee data by department and role to calculate total headcount and identify top-performing employees.

- **Data Validation:** Checking for missing values in employee records and correcting data types to ensure that the data is accurate and consistent.

## 2.7 Steps Involved

1. **Create ETL Job:**
   - Configure the job name and select the appropriate Glue version (e.g., Glue 3.0) to match the specific requirements.
   - Adjust timeout settings to avoid errors, ensuring that the job has enough time to process the data.

2. **Set Permissions:**
   - Attach policies such as AWSGlueConsoleFullAccess, AmazonS3FullAccess, CloudWatchFullAccess and to grant the necessary permissions for data processing.

3. **Add Trigger:**
   - Choose the S3 bucket as the trigger source and acknowledge the configuration to ensure that the Glue job is triggered when new data is uploaded.

4. **Add Libraries:**
   - Ensure that the Glue job has access to these libraries by specifying them in the job configuration.

5. **PySpark Script:**
   - Edit the script to specify the bucket name, folder, and default date.
   - Generate cleaned employee CSV and ensuring that the data is validated and aggregated correctly.
   - Deploy and test the job to ensure that it processes the data as expected.

```
Script  Info
1  import sys
2  from pyspark.context import SparkContext
3  from awsglue.context import GlueContext
4  from awsglue.utils import getResolvedOptions
5  from awsglue.job import Job
6  from pyspark.sql.functions import col, regexp_replace, udf
7  from pyspark.sql.types import DateType
8  from dateutil import parser
9
10  # Initialize Glue context and Spark session
11  sc = SparkContext.getOrCreate()
12  glueContext = GlueContext(sc)
13  spark = glueContext.spark_session
14  job = Job(glueContext)
15
16  # Load the CSV data from S3
17  input_path = "s3://project-bucket-2387739/project_data.csv"
18  df = spark.read.option("header", "true").csv(input_path)
19
20  # UDF to parse dates
21  def parse_date(date_str):
22      try:
23          return parser.parse(date_str)
24      except:
25          return None
26
```

Fig 6. Spark code for data transformation



Fig 7. Run the job



Fig 8. Checking for transformed CSV

Fig 9. Transformed Data Preview



Fig 10. Setting up roles, warehouses and databases

Fig 11. Creating employee table in Snowflake

# Chapter 3: Data Loading

## 3.1 Definition

Data loading is the process of loading the transformed data into Snowflake for analytics and reporting. This step ensures that the data is available in a structured format, ready for querying and analysis.

## 3.2 Types of Data Loading

- **Full Load:** This method involves loading the entire dataset into the target system. It is suitable for initial data loads or when a complete refresh of the data is required.
- **Incremental Load:** This method involves loading only the new or updated data. It is ideal for ongoing data updates, as it minimizes the amount of data transferred and reduces the load on the system.

## 3.3 Categories of Data Loading

- **Batch Loading:** This approach involves loading data in batches at scheduled intervals. It is suitable for periodic data updates, such as daily or weekly loads, and helps in managing data loads effectively.
- **Real-Time Loading:** This approach continuously loads data as it becomes available. It is ideal for systems that require immediate data availability, such as real-time analytics and monitoring systems.

## 3.4 Advantages

- **Performance:** Optimizing data loading processes ensures faster data processing and reduces the time required to make the data available for analysis.
- **Scalability:** Data loading processes can scale to handle large volumes of data efficiently, making them suitable for growing datasets.
- **Flexibility:** Various data loading methods can be used to meet different requirements, such as full loads for initial data setup and incremental loads for ongoing updates.

## 3.5 Disadvantages

- **Complexity:** Setting up and configuring data loading processes requires careful planning and technical expertise, especially for large datasets.
- **Resource Intensive:** Data loading can consume significant computational resources, particularly for large datasets or complex transformations.
- **Dependency:** The loading process relies on the availability of the target system, which can be a limitation if the system is unavailable or experiencing issues.

## 3.6 Examples

- **Employee Data Loading:** Loading processed employee data into AWS Redshift to enable querying and analysis of employee performance and other HR metrics.

- **Activity Log Data Loading:** Loading employee activity log data into a centralized database for monitoring and analysis, helping to identify productivity trends and troubleshoot issues promptly.

## 3.7 Steps Involved

1. **Set Up Warehouse and Database:**
   - Create a warehouse to provide the required resources (CPU, memory, and temporary storage) for DML operations:
   - Create a database and schema.
2. **Create IAM Role:**
   - Create an IAM role of AWS Account and attach an S3 inline policy for Snowflake to access the necessary permissions for data loading.
   - Edit the role permissions to ensure that the role has the appropriate access to the required resources.
3. **Create storage integration:**

   - Describe the storage integration after its creation using AWS ARN and S3 bucket path.
   - Use the describe command to get the Amazon Resource Name (ARN).

4. **Create a file format for CSV:**

   - Create a file format for CSV by properly mentioning the header and null options.

5. **Create Stage and Target Table:**

   - Create a stage and target table specifying where the files are stored.
   - Load data into that table using COPY command.

Fig 12. Integrating S3 and Snowflake



Fig 13. Editing IAM arn and external ID

Fig 14. Accessing the external storage



Fig 15. Load data to employee table

Fig 16. Table in snowflake

# Chapter 4: Real Time Data Storage

Real-time data storage is a critical component of modern data processing systems, enabling organizations to capture, store, and analyse data as it is generated. This capability is essential for applications that require immediate insights and rapid response to changing conditions. AWS offers a suite of services that facilitate real-time data storage, ensuring scalability, low latency, and robust data management. In this section, we will explore the concept of real-time data storage, its benefits, challenges, and how AWS services can be leveraged for effective implementation.

## 4.1 Definition

Real-time data storage refers to the continuous and immediate capture and storage of data as it is generated. Unlike batch processing, which involves collecting and processing data at scheduled intervals, real-time data storage ensures that data is available for analysis and decision-making almost instantaneously.

## 4.2 Benefits

1. **Immediate Insights**: Real-time data storage allows organizations to gain immediate insights into their operations. This capability is crucial for applications such as monitoring, fraud detection, and real-time analytics, where timely information is essential for effective decision-making.

2. **Enhanced Responsiveness**: By capturing and storing data in real-time, organizations can respond more quickly to changing conditions. For example, in retail, real-time data storage enables dynamic pricing and inventory management, helping businesses stay competitive and meet customer demands.

3. **Improved Customer Experience**: Real-time data storage supports personalized and timely interactions with customers. For instance, e-commerce platforms can use real-time data to recommend products based on current browsing behaviour, enhancing the customer experience and driving sales.

4. **Operational Efficiency**: Real-time data storage streamlines operations by providing up-to-date information. This capability is valuable in industries such as manufacturing, where real-time data can be used to monitor equipment performance and prevent downtime.

## 4.3 Challenges

1. **Scalability**: Real-time data storage systems must be able to handle large volumes of data generated at high velocity. Ensuring scalability requires robust infrastructure and efficient data processing mechanisms.

2. **Data Quality**: Maintaining data quality in real-time environments can be challenging. Organizations must implement validation and cleansing processes to ensure that the data being stored is accurate and reliable.

3. **Latency**: Minimizing latency is critical for real-time data storage. Any delays in data capture, processing, or storage can impact the timeliness of insights and decision-making.

4. **Cost**: Real-time data storage can be resource-intensive, requiring significant investment in infrastructure and technology. Organizations must balance the benefits of real-time capabilities with the associated costs.

## 4.4 AWS Services for Real-Time Data Storage

1. **AWS Lambda**: AWS Lambda allows you to run code in response to events, enabling real-time data processing without provisioning or managing servers. It can be used to process data streams from Kinesis and trigger actions based on real-time data.

2. **Amazon S3**: Amazon S3 provides scalable object storage for real-time data. It can be used to store raw and processed data, ensuring durability and availability. S3's integration with other AWS services makes it a versatile choice for real-time data storage.

3. **AWS Glue:** AWS Glue is a serverless data integration service that simplifies the process of discovering, preparing, moving, and integrating data from multiple sources for analytics, machine learning, and application development.

4. **Snowflake**: Amazon Snowflake is a data warehousing service that supports real-time data loading and querying for large amount of data. It can be used to analyse real-time data streams and generate insights quickly.

## 4.5 Use Cases

1. **Retail**: Real-time data storage enables dynamic pricing, inventory management, and personalized marketing in retail. By capturing sales and customer data in real-time, retailers can optimize their operations and enhance the customer experience.

2. **Finance**: In the finance industry, real-time data storage supports fraud detection, algorithmic trading, and risk management. Financial institutions can analyse transactions as they occur, identifying suspicious activity and making informed decisions.

3. **Healthcare**: Real-time data storage is critical for patient monitoring, telemedicine, and medical research. Healthcare providers can capture and analyse patient data in real-time, improving diagnosis, treatment, and patient outcomes.

4. **Manufacturing**: Real-time data storage supports predictive maintenance, quality control, and supply chain optimization in manufacturing. By monitoring equipment performance and production processes in real-time, manufacturers can prevent downtime and improve efficiency.

# Chapter 5: Analytical Data Storage

Analytical data storage is crucial for organizations that need to process and analyze large volumes of data to derive actionable insights. AWS offers a range of services that facilitate efficient and scalable analytical data storage, ensuring that data is readily available for querying and analysis. In this section, we will explore the concept of analytical data storage, its benefits, challenges, and how AWS services can be leveraged for effective implementation.

## 5.1 Definition

Analytical data storage refers to the storage of processed data in a structured format that can be queried using analytical tools. This type of storage is optimized for read-heavy operations and supports complex queries, making it ideal for business intelligence, reporting, and data analysis.

## 5.2 Benefits

1. **Enhanced Query Performance**: Analytical data storage systems are designed to handle complex queries efficiently. By organizing data in a way that optimizes query performance, these systems enable faster data retrieval and analysis.

2. **Scalability**: Analytical data storage solutions can scale to accommodate growing data volumes. This scalability ensures that organizations can continue to analyse data effectively as their data needs expand.

3. **Data Integration**: Analytical data storage systems often support the integration of data from multiple sources. This capability allows organizations to consolidate data from various systems, providing a comprehensive view of their operations.

4. **Advanced Analytics**: By storing data in a structured format, analytical data storage systems enable advanced analytics, such as machine learning and predictive modelling. These capabilities help organizations uncover deeper insights and make data-driven decisions.

## 5.3 Challenges

1. **Data Volume**: Managing large volumes of data can be challenging. Analytical data storage systems must be able to handle high data throughput and ensure data integrity.

2. **Data Quality**: Ensuring data quality is critical for accurate analysis. Organizations must implement data validation and cleansing processes to maintain the quality of their analytical data.

3. **Cost**: Analytical data storage solutions can be resource-intensive, leading to higher costs. Organizations must balance the benefits of advanced analytics with the associated expenses.

4.  **Latency**: Minimizing latency is important for timely data analysis. Analytical data storage systems must be optimized to reduce delays in data processing and querying.

## 5.4 AWS Services for Analytical Data Storage

1.  **Amazon Redshift**: Amazon Redshift is a fully managed data warehousing service that enables fast querying and analysis of large datasets. It supports complex queries and integrates with various data sources, making it ideal for analytical data storage.

2.  **Amazon S3**: Amazon S3 provides scalable object storage for analytical data. It can store raw and processed data, ensuring durability and availability. S3's integration with other AWS services makes it a versatile choice for analytical data storage.

3.  **AWS Glue**: AWS Glue is a fully managed ETL (extract, transform, load) service that prepares data for analysis. It can integrate data from multiple sources and transform it into a format suitable for analytical storage.

4.  **Amazon Athena**: Amazon Athena is an interactive query service that allows you to analyse data directly in Amazon S3 using standard SQL. It is serverless, meaning there is no infrastructure to manage, and you only pay for the queries you run.

5.  **Amazon EMR**: Amazon EMR is a cloud big data platform that provides a managed Hadoop framework. It can process large amounts of data quickly and cost-effectively, making it suitable for analytical data storage and processing.

## 5.5 Use Cases

1.  **Business Intelligence**: Analytical data storage supports business intelligence applications by providing a structured and optimized environment for querying and reporting. Organizations can generate insights from their data to inform strategic decisions.

2.  **Data Warehousing**: Analytical data storage is essential for data warehousing, where large volumes of data are stored and analysed. Data warehouses enable organizations to consolidate data from various sources and perform complex queries.

3.  **Machine Learning**: Analytical data storage systems provide the foundation for machine learning applications. By storing data in a structured format, these systems enable the training and deployment of machine learning models.

4.  **Financial Analysis**: In the finance industry, analytical data storage supports activities such as risk assessment, fraud detection, and portfolio management. Financial institutions can analyse large datasets to identify trends and make informed decisions.

# Challenges Faced and Solutions

Implementing a robust Employee Data Analysis System involves overcoming several challenges in data processing. These challenges can impact the efficiency, accuracy, and scalability of the system. Here, we explore some common challenges and their solutions, particularly in the context of using AWS services and Snowflake.

1. **Data Volume**:

   - o **Challenge**: Handling large volumes of employee data generated daily from multiple departments can be overwhelming. The sheer amount of data can create bottlenecks in processing and storage.

   - o **Solution**: Utilize scalable storage solutions like AWS S3 and Snowflake. AWS S3 provides virtually unlimited storage capacity, while Snowflake can handle large datasets with its columnar storage and parallel processing capabilities. Implementing data partitioning and compression techniques can further optimize storage and query performance.

2. **Data Quality**:

   - o **Challenge**: Ensuring the accuracy and consistency of data is critical. Incomplete, incorrect, or inconsistent data can lead to erroneous insights and decisions.

   - o **Solution**: Implement data validation and cleansing processes using AWS Glue. It will ensure high data quality before it is loaded into Snowflake. Regular audits and monitoring can also help maintain data integrity.

3. **Data Integration**:

   - o **Challenge**: Integrating data from various sources and formats can be complex. Employee data may come in different formats, requiring transformation before analysis.

   - o **Solution**: Use AWS Glue for ETL (extract, transform, load) processes. AWS Glue can automatically discover and catalog data. Standardizing data formats and using consistent schemas can also simplify integration.

4. **Latency**:

   - o **Challenge**: Minimizing latency in data processing is essential for timely insights. Delays in data extraction, transformation, or loading can hinder real-time analysis.

   - o **Solution**: Leverage AWS Kinesis for real-time data streaming and processing. AWS Kinesis Data Streams can capture and process data in real-time, reducing latency. Additionally, using AWS Glue for event-driven processing can ensure that data is processed as soon as it is available.

# Conclusion

The Employee Data Analysis System represents a major leap forward in how a company can manage, process, and analyze its employee data. By harnessing the capabilities of AWS services such as S3, Glue, and Snowflake, this system offers a robust and scalable solution tailored to meet diverse data management and analytics requirements.

The process begins with data extraction, where daily employee CSV files are uploaded from a local directory to AWS S3. This ensures that raw employee data is securely stored in a centralized location, ready for subsequent processing. Automating this step allows the system to efficiently handle large volumes of data, minimizing human error and ensuring consistent and accurate data uploads.

The next phase involves data transformation, where AWS Glue with PySpark is used to validate and aggregate the raw employee data. This step is crucial for maintaining data quality by removing or correcting invalid entries and preparing datasets for analysis. Automation of these tasks reduces manual effort and ensures consistent and efficient data processing. AWS Glue's scalability and flexibility make it ideal for handling varying data volumes and complexities.

Once the data is transformed, it is loaded into Snowflake for analytics and reporting. This ensures that the data is structured and optimized for querying and analysis. By streamlining the data loading process, the system can efficiently manage large datasets, making it suitable for growing data needs. Snowflake provides a powerful and scalable data warehousing solution that supports complex queries and analytics.

Additionally, CloudWatch logs are utilized to monitor and log errors during the ETL process, ensuring timely detection and resolution of issues. This enhances the reliability and robustness of the entire data pipeline.

In summary, the Employee Data Analysis System is a comprehensive and scalable solution that addresses various data management and analytics needs for a company. Leveraging AWS services, the system ensures efficient data extraction, transformation, loading, and reporting. This robust system not only enhances operational efficiency but also provides valuable insights that support data-driven decision-making and strategic planning. As the company grows, the Employee Data Analysis System will be instrumental in managing and analyzing employee data, driving business success and growth.

THANK YOU