## AWS SERVICES USED IN DATA ENGINEERING (DE) WORKFLOWS

Several AWS services are commonly used in data engineering (DE) workflows. Here are some AWS services that are relevant to data engineering:

1. **S3** (Simple Storage Service): AWS's scalable object storage service is often used as a data lake or data storage solution in data engineering pipelines.

2. **Glue**: AWS Glue is an ETL (Extract, Transform, Load) service that helps with data preparation and transformation for analytics and data processing.

3. **Athena**: A serverless query service that allows you to analyze data directly from S3 using standard SQL, making it useful for ad hoc data exploration and analysis

4. **Redshift**: Amazon Redshift is a fully managed data warehousing service that is optimized for online analytical processing (OLAP) and handling large datasets.

5. **EMR** (Elastic MapReduce): EMR is a managed big data platform that enables the processing of large datasets using popular frameworks like Apache Spark, Hadoop, and Presto.

6. **Data Pipeline**: AWS Data Pipeline is a web service that helps with the orchestration and automation of data-driven workflows, allowing you to move and process data between different AWS services.

7. **Kinesis**: Kinesis offers real-time data streaming capabilities, allowing you to ingest and process large volumes of streaming data for near real-time analytics

8. **Glue DataBrew**: A visual data preparation service that helps clean and normalize data for analytics and machine learning.

9. **Quicksight**: Amazon QuickSight is a business intelligence service that allows you to create interactive visualizations and reports for data analysis and reporting.

10. **Step Functions**: AWS Step Functions is a serverless workflow service that allows you to coordinate and orchestrate different components of a data pipeline or workflow.

These are just a few examples of AWS services commonly used in data engineering workflows. The choice of services may vary depending on specific use cases and requirements.

## PIPELINE (ROUGH SKETCH) OF HOW IS THE DATA PIPELINE FLOW IN AWS.

Certainly! Here's a rough sketch of a data pipeline flow in AWS:

1. Data Ingestion: Data from various sources is ingested into the pipeline. This can include structured, semi-structured, or unstructured data. Common sources include databases, APIs, streaming data, files, and IoT devices.

2. Data Storage: The ingested data is stored in a scalable and durable data storage solution. Amazon S3 (Simple Storage Service) is often used as a central data repository or a data lake for storing the raw data.

3. Data Preparation: AWS Glue can be utilized for data preparation and transformation tasks. Glue allows you to define data schemas, perform data cleansing, normalization, and data format conversions. It also facilitates data enrichment and deduplication.

4. Data Processing: Data processing tasks are performed on the prepared data. AWS offers various services for data processing, depending on the requirements. Amazon EMR (Elastic MapReduce) can be used for distributed processing using tools like Apache Spark, Hadoop, or Presto. AWS Lambda can also be employed for serverless data processing

5. Data Analysis and Querying: For ad hoc data analysis and querying, AWS Athena can be utilized. Athena allows SQL-based querying directly on the data stored in S3, enabling exploration and analysis without the need for infrastructure provisioning.

6. Data Warehousing: If the data needs to be stored in a structured and optimized format for analytics, Amazon Redshift can be used as a fully managed data warehousing solution. Redshift allows high-performance querying of large datasets using SQL

7. Data Visualization and Reporting: AWS QuickSight can be employed for creating interactive visualizations and reports based on the processed and analyzed data. QuickSight integrates with various AWS services and provides dashboards for data exploration and sharing.

8. Data Orchestration: AWS Step Functions can be used for orchestrating and managing the different stages of the data pipeline. Step Functions enables you to define workflows, handle dependencies between tasks, and manage error handling and retries.

9. Monitoring and Logging: Throughout the pipeline, it's important to monitor and log the pipeline activities for troubleshooting, performance optimization, and compliance purposes. AWS CloudWatch can be used for collecting logs and monitoring various metrics.

It's essential to note that the specific services and components used in the data pipeline can vary based on the requirements and architecture of the system. This is a high-level overview, and the actual implementation of the pipeline can involve additional steps, services, and customizations based on the specific use case and data engineering requirements.