# Assignment

**Scenario:** As a Senior Data Engineer at "DataCo," an e-commerce company, your objective is to design and implement a robust ETL pipeline to handle data ingestion, transformation, and loading from diverse social media and news sources. The data will be delivered as JSON format through APIs from platforms such as Twitter, Facebook, Instagram, TikTok, YouTube, blogs, forums, LinkedIn, etc.

**Requirements:**

1. **Data Sources (APIs):**
   o Identify and integrate APIs from the specified social media and news platforms to simulate data sources.
   o Extract sample JSON data representing various aspects, such as posts, comments, likes, shares, articles, etc.
2. **ETL Processing:**
   o Design an ETL process that dynamically accommodates varying JSON structures from different APIs.
   o Implement the ETL process using a language or framework of your choice, emphasizing modularity and extensibility.
3. **Data Pipeline:**
   o Construct a data pipeline architecture that automates the ETL process for seamless integration of data from diverse sources.
   o Implement mechanisms for handling API rate limits, retries, and failures.
4. **Data Ingestion:**
   o Develop a data ingestion mechanism capable of ingesting real-time data from social media platforms and batch data from news sources.
   o Consider scenarios where data arrives in bursts or steadily throughout the day.
   o The data has to be ingested in AWS open search
5. **Scalability and Performance:**
   o Optimize the ETL process and data pipeline for scalability and performance given the dynamic nature of social media data.
   o Discuss strategies to handle varying data volumes and spikes in activity.
6. **Documentation:**
   o Provide comprehensive documentation detailing the integration of each API, the ETL process, and the overall data pipeline.
   o Include clear instructions for maintaining and updating API connections.
7. **Testing:**
   o Develop test cases to ensure the accuracy and reliability of the ETL process and data pipeline.
   o Include tests for different JSON structures and variations across platforms.
8. **Security Considerations:**
   o Address security considerations in the design and implementation, especially when handling sensitive data from social media platforms.

**Submission:**

- Share the codebase via a version control system (e.g., GitHub).
- Include a README file explaining the architecture, setup instructions, and any additional information.

**Evaluation Criteria:** Candidates will be evaluated based on their ability to:

- Effectively integrate and handle JSON data from diverse APIs.
- Design and implement a flexible and scalable ETL process.
- Build a resilient data pipeline considering API dynamics.
- Efficiently handle real-time and batch data ingestion.
- Optimize for scalability and performance in a dynamic data environment.
- Provide thorough documentation and testing.