

## Written Report

### Understanding of the Business Context and Data Sources

TelcoCorp is a leading telecommunications company that generates vast amounts of data from various sources, including network traffic logs, customer transactions, product catalogs, and marketing campaigns. The company aims to build a centralized data platform to consolidate and process data from multiple sources, enabling advanced analytics and business intelligence capabilities.

The data sources include:

- Network Traffic Logs: Raw web server logs stored in a data lake, containing information about user visits, page views, and click events.
- Transactional Data: Customer orders and payment data stored in a relational database.
- Product Catalog: Product information, including descriptions, prices, and inventory levels, stored in CSV files.
- Marketing Campaign Data: Real-time data stream of marketing campaign interactions (e.g., email opens, click-throughs) from a messaging queue.

### Rationale behind the Chosen Solution and Design Decisions

The chosen solution is to design and implement a scalable data platform using AWS Cloud, Apache Spark, Apache Airflow, and Power BI. The design decisions were driven by the need for scalability, performance, and flexibility.

- AWS Cloud provides a scalable and secure infrastructure for data processing and storage.
- Apache Spark is used for data transformation and integration due to its ability to handle large-scale data processing and its compatibility with AWS Glue.
- Apache Airflow is used for workflow orchestration and automation to ensure efficient and reliable data processing.
- Power BI is used for data visualization and business intelligence to provide insights and enable data-driven decision-making.

The design decisions were influenced by the following factors:

- Scalability: The solution needs to handle large volumes of data and scale horizontally to accommodate increasing data volumes.

- **Performance:** The solution needs to process data in near real-time to support timely business decisions.
- **Flexibility:** The solution needs to be flexible to accommodate changing business requirements and new data sources.

## **Scalability and Performance Considerations**

The solution is designed to scale horizontally to accommodate increasing data volumes. The use of AWS Cloud and Apache Spark enables the solution to handle large-scale data processing and scale up or down as needed.

Performance considerations include:

- **Data ingestion:** The solution uses AWS Kinesis Data Firehose and AWS Lambda to ingest data in near real-time.
- **Data processing:** The solution uses Apache Spark to process data in parallel, reducing processing time and improving performance.
- **Data storage:** The solution uses Amazon S3 and Amazon Redshift to store data, providing scalable and performant data storage.

## **Potential Challenges and Limitations**

Potential challenges and limitations include:

- **Data quality issues:** Poor data quality can affect the accuracy of insights and decision-making.
- **Data security:** The solution needs to ensure data security and compliance with regulatory requirements.
- **Complexity:** The solution involves multiple technologies and components, which can increase complexity and require specialized skills.

## **Future Enhancements or Improvements**

Future enhancements or improvements include:

- **Real-time analytics:** Implementing real-time analytics capabilities to support timely business decisions.
- **Machine learning:** Integrating machine learning algorithms to enable predictive analytics and automate decision-making.
- **Data governance:** Implementing data governance policies and procedures to ensure data quality and security.

- Cloud optimization: Optimizing cloud resources and costs to ensure cost-effectiveness and efficiency.

In conclusion, the proposed solution addresses the business requirements of TelcoCorp by providing a scalable and performant data platform for advanced analytics and business intelligence. The solution is designed to accommodate changing business requirements and new data sources, ensuring flexibility and adaptability.