

## **Task 2**

### **Part 1 : Unifying various sources of data (e.g. various “views” of a member of The Room) to derive a “master record” (e.g. consolidated member profile).**

In our proposed architecture for The Room, we aim to achieve a cost-effective approach to unify data from various sources and derive a consolidated member profile. Here's how we can do it:

**Data Integration:** We develop a process to bring together data from different sources by extracting and transforming it into a standardized format. This ensures that data from different systems or views associated with a member of The Room can be easily processed and combined.

**Data Mapping and Matching:** We use algorithms to identify and reconcile records across different data sources. By establishing common identifiers or attributes, we can link and merge data points related to the same member. Techniques like fuzzy matching or rule-based matching help us handle variations in the data.

**Data Cleansing and Standardization:** We clean and standardize the data to ensure its quality and consistency. This involves removing duplicates, correcting errors, and formatting the data based on predefined rules or reference datasets.

**Data Fusion and Aggregation:** We merge and aggregate relevant attributes and information from different sources to create a consolidated member profile. By combining the data using predefined rules or algorithms, we prioritize certain sources or attributes based on their quality and reliability.

**Data Governance and Metadata Management:** We establish data governance practices and metadata management to track the origin and lineage of the data. This helps us document the data sources, transformations, and integration rules, ensuring transparency and facilitating future updates.

**Incremental Updates:** We implement mechanisms to handle incremental updates or changes in member data. This allows us to efficiently capture and process updates using techniques like change data capture or event-driven architectures, ensuring that the consolidated member profile remains up-to-date.

**Monitoring and Quality Assurance:** We set up processes to monitor and perform quality assurance checks on the consolidated member profile. This helps us continuously validate and verify the accuracy and completeness of the data, identifying and resolving any discrepancies or anomalies that may arise over time.

By following these steps within the proposed architecture, The Room can effectively unify data from various sources and derive a consolidated member profile. This approach ensures efficient data integration, matching, cleansing, fusion, and ongoing maintenance of the member profile, providing a comprehensive and accurate view of each member's information.

**Part 2: Does a data lake play any part in your architecture? What about a data warehouse? What are the interfaces between the two?**

In the proposed architecture, both a data lake and a data warehouse play important roles in handling data storage and processing:

**Data Lake:** The data lake serves as a central repository for storing raw, unprocessed data from various sources. It acts as a scalable and cost-effective storage solution that can handle diverse data types, including structured, semi-structured, and unstructured data. The data lake is typically built on a distributed file system, such as Hadoop Distributed File System (HDFS), and allows for flexible data ingestion and exploration.

**Data Warehouse:** The data warehouse is designed to support efficient querying, analysis, and reporting of structured and processed data. It is optimized for high-performance analytics and typically employs a schema-on-write approach, where data is transformed, cleaned, and organized into a predefined schema before being loaded into the warehouse. Data warehouses provide a structured and reliable environment for business intelligence and reporting purposes.

The interfaces between the data lake and the data warehouse can be established through data pipelines and data integration processes. Here's how they interact:

**Data Ingestion:** Raw data from various sources is ingested into the data lake. This can include both batch data and real-time streaming data. Data ingestion processes, such as Apache Kafka, Apache Nifi, or AWS Kinesis, capture and store the data in its original format in the data lake.

**Data Processing:** Data processing components, such as Apache Spark, operate on the data lake to perform transformations, aggregations, and data cleansing. This processed data is then loaded into the data warehouse. Apache Spark's capabilities allow for distributed data processing across the data lake, enabling scalable and efficient data transformations.

**Data Transformation and Loading:** As the data is processed within the data lake, it can be transformed into a structured format suitable for loading into the data warehouse. This transformation step involves cleaning, standardizing, and organizing the data according to the predefined schema of the data warehouse.

**Data Loading into Data Warehouse:** Once the data is transformed, it is loaded into the data warehouse. This can be achieved through various methods, such as direct data transfers, batch processing, or ETL (Extract, Transform, Load) workflows. Technologies like Amazon Redshift, Google BigQuery, or Snowflake can be used for efficient loading and storage of structured data in the data warehouse.

In summary, the data lake stores the raw data, and the data warehouse stores the processed data. They work together by bringing in the raw data, processing and organizing it, and then moving it to the data warehouse for analysis and reporting.

### **Part 3 : Integration with data consumption interfaces such as BI and data applications**

In our proposed architecture, we've considered a cost-effective approach to integrate data consumption interfaces like business intelligence (BI) tools and data applications. Here's how we address this:

We create standardized interfaces, called APIs, which act as bridges between the data storage layers (data warehouse or data lake) and the data consumption interfaces. These APIs provide controlled access to the data, allowing users to retrieve and manipulate it without directly interacting with the underlying storage systems.

To simplify data access, we leverage data virtualization techniques. This allows users to seamlessly access and combine data from multiple sources, eliminating the need for duplicating data. Tools like Apache Drill or AWS Glue DataBrew help us achieve this.

We empower data consumers with self-service BI capabilities. By providing user-friendly tools like Tableau or Power BI, users can explore and analyze data on their own, reducing their dependency on IT teams.

Metadata management plays a crucial role. We establish a framework to catalog and document available data assets, relationships, and usage information. This enables users to easily discover and understand the data they need. Tools like Apache Atlas or AWS Glue Data Catalog assist in managing metadata effectively.

We prioritize data security and governance. Implementing measures like access controls, data masking, and encryption ensures that sensitive data is protected and accessed only by authorized users. Additionally, data governance practices help maintain data quality and compliance with regulations.

We leverage scalable infrastructure resources, such as cloud-based solutions like AWS or Azure, to handle varying workloads. Cloud services offer cost-effective scalability and flexibility to accommodate growing data demands.

To optimize performance, we utilize techniques like caching, query optimization, and data indexing. This helps improve query response times and overall system efficiency.

By following these strategies, we achieve seamless integration with data consumption interfaces in a cost-effective manner. Users can access, analyze, and visualize data using BI tools and data applications while ensuring security, governance, scalability, and performance optimization.