# International Journal of Research Publication and Reviews

# DATA INTEGRATION

*Siddhnath Gharge*

*ASM Institute of Management And Computer Studies, Thane*
*Email: - siddhnath99@gmail.com*

## ABSTRACT

Data integration is the process of combining data from different sources into a single integrated view. Integration begins with the capture process and includes steps such as cleansing, ETL mapping, and transformation. Ultimately, data integration allows analytic tools to create effective and actionable business intelligence.

There is no one-size-fits-all approach to data integration. However, data integration solutions typically include some common elements such as the data source's network, the master server, and clients accessing data from the master server.

In a typical data integration process, the client sends a data request to the master server. The master server then retrieves the required data from internal and external sources. The data is extracted from the source and integrated into a single, consistent dataset. It will be sent back to the client for use.

## 1. DATA INTEGRATION: DEFINITION

In today's world, managing, processing, storing and retrieving data are the most important and necessary skills. This process is all about data storage, and data storage is important. For this part, we will use Snowflake Data Warehouse, a simple and comprehensive data warehouse. The main goal of this project is to get data from various sources, perform operations on that data and pass it to Snowflake's data warehouse using Kafka middleware.

**Speed of data:**

Speed refers to the speed of data processing. Emergency procedures, such as sending and receiving data, require the use of KAFKA middleware during transfer to the organization to increase its value.

**Processing of data:**

This process streams data from various data sources and uses Pandas to convert it to Python-related data. Import the data into Python and use Python to create the Kafka Sink Snowflake Connector. You also need to create a Kafka topic after creating the Kafka Connector. From Kafka topics using the Sink Snowflake Connector to Snowflake databases or data warehouses. Faster than other data sources such as Azure Synapse Analytics, Azure Data Lake Storage, and Oracle Exadata Cloud Service, Snowflake Data Warehouse can also store large amounts of data and easily retrieve data from Snowflake in minimal time.

**Usage:**

The Process gives easy steps to perform operations on large amount of data. It can be passed through middleware easily and also, we can retrieve it from snowflake data warehouse. We can store huge amount of data into snowflake warehouse and can perform operations on it

**Human collaboration:**

We've created a guarantee to take away the impurities in data Associate in Nursing make it clean. Now, the future step is to mix data from different sources to induce a unified structure with additional significant and valuable information. This is often largely used if the info is separated into different sources.

To create it simple, let' assume we've data in CSV format in several places, all talking concerning an equivalent scenario. Say we have some data about an employee during a database. We have a tendency to not expect all the data about the worker to reside within the same table.

It's possible that the worker's personal information is going to be placed on one table, the employee's project history will be at a second table, the employee's time-in and time-out details will be at another table, and so on. So, if we wish to try to do some analysis concerning the employee, we want to induce all the employee data in one common place.

This method of transferring data along in one place is termed data integration. To try data integration, we are able to merge multiple pandas DataFrames victimization the merge function.

## 2.  KAFKA CLUSTER

In a distributed computing system, a cluster is a collection of computers that work together to achieve a common goal. A Kafka cluster is a system consisting of multiple brokers, themes, and both partitions.

The main goal is to balance the workload between replicas and partitions.

Kafka Clusters Architecture especially includes the subsequent five components:

- Topics
- Broker
- ZooKeeper
- Producers
- Consumers

**Kakka topic:**

Kafka subjects are the types used to prepare messages. Each subject matter has a call this is particular throughout the whole Kafka cluster. Messages are despatched to and study from particular subjects.

In different words, manufacturers write facts to subjects, and purchasers study facts from subjects.

Kafka topics are multi-subscriber.

This means that a topic can have 0, 1, or multiple consumers that subscribe to that topic and the data written to it. In

Kafka, topics are split and duplicated across brokers throughout the implementation. The broker references each node in the Kafka cluster. Partitions are important because they parallelize topics and allow for high message throughput.

**Kafka producer:**

Basically, a utility this is the supply of the facts circulate is what we name a producer. In order to generate tokens or messages and in addition submit it to at least one or greater subjects within side the Kafka cluster, we use Apache Kafka Producer. However, to submit a circulate of data to at least one or greater Kafka subjects, this Kafka Producer API lets in to a utility. Moreover, its significant element is Kafka Producer class.

**Kafka consumer:**

Kafka clients are commonly a part of a client institution. When a couple of clients are subscribed to a subject and belong to the equal client institution, every client within side the institution will get hold of messages from a distinctive subset of the walls within the topic.

## 3.  CONCLUSION

In Today's world keeping track of data and processing on it is the toughest thing faced by peoples. Data integration is the process which makes it easier and faster with the help of Apache Kafka and Snowflake Data Warehouse.

## REFERENCES

[1]  https://www.confluent.io/lp/apache-kafka/?utm_medium=sem&utm_source=google&utm_campaign=ch.sem_br.nonbrand_tp.prs_tgt.kafka_mt.xct_rgn.india_lng.eng_dv.all_con.kafka-general&utm_term=apache%20kafka&creative=&device=c&placement=&gclid=EAIaIQobChMI5v3G3PfM-AIVu5NmAh1nlwN1EAAYASAAEgLj5vD_BwE

[2]  https://data-flair.training/blogs/kafka-producer/#:~:text=What%20is%20Kafka%20Producer%3F,we%20use%20Apache%20Kafka%20Producer.

[3] https://kafka.apache.org/26/javadoc/index.html?org/apache/kafka/clients/consumer/KafkaConsumer.html#:~:text=Kafka%20uses%20the%20concept%20of,and%20fault%20tolerance%20for%20processing.

[4] https://www.snowflake.com/guides/what-kafka#:~:text=The%20Snowflake%20Connector%20for%20Kafka,of%20separating%20storage%20and%20compute.