# Real-Time Data Streaming and Processing using Synapse Analytics

## Hari Prasad Bomma

Data Engineer, USA haribomma2007@gmail.com

### Abstract

Extract, Transform, Load (ETL) is a traditional method widely used for data integration, involving extracting data from various sources, transforming it to meet operational needs, and loading it into a target data warehouse. Regular ETL processes typically scheduled at intervals like daily or weekly, offer advantages such as simplifying data processing and reducing resource usage during off peak hours. However, they also present significant drawbacks, including latency and difficulty in scaling with large data volumes, which can lead to processing delays and potential system failures. The paper will explore these challenges and the increasing demand for real time data processing requirements, focusing on high throughput and low latency data streams, and the need for scalable and reliable infrastructure. The paper will present Microsoft's Synapse Analytics as a comprehensive solution, detailing its unified capabilities for data engineering, warehousing, and exploration to meet contemporary data processing needs.

Keywords: Data Streaming, Real Time data processing, Synapse, Apache Spark, Event hub, Blob Storage, IoT

#### Introduction:

Extract, Transform, Load (ETL) is a traditional method used for data integration and processing. The ETL process involves extracting data from various sources, transforming it to fit operational needs, and loading it into a target data warehouse or database. This approach is widely used in many industries for consolidating data, ensuring data quality, and enabling comprehensive analytics. Regular ETL processes are typically scheduled to run at specific intervals, such as daily, weekly, or monthly, depending on the organization's requirements.

While batch loads in regular ETL processes have their advantages, such as simplifying data processing and reducing resource usage during off peak hours, they also come with significant drawbacks. One key challenge is the latency involved, as batch processing often delays the availability of up to date data for analysis, leading to potential delays in decision making. Furthermore, batch loads can become problematic when dealing with large volumes of data, as the time and resources required for processing may increase exponentially. This can lead to longer processing times, potential system failures, and difficulty in scaling with growing data demands. To address these limitations, organizations are increasingly looking towards real time data processing and modern ETL approaches.

In the era of Big Data, the demand for real time data processing and analysis has become increasingly important for businesses to make informed decisions and gain a competitive advantage [5]. The

growth of the Internet of Things has been a significant driver of this trend, as the proliferation of IoT devices has led to an exponential increase in the volume and velocity of data streams that need to be processed and analyzed in near real time [5].

In order to overcome these challenges, modern data processing systems must be designed to handle high throughput and low latency data streams, while also providing scalable and robust infrastructure for data ingestion, processing, and storage. Synapse, Microsoft's cloud based data analytics platform, offers a comprehensive solution for these requirements, with its Synapse Analytics component providing a unified experience for data engineering, data warehousing, and data exploration.

#### Synapse Analytics Capabilities for Real Time ETL:

Synapse Analytics offers a range of features and capabilities that enable efficient and scalable real time ETL workflows in the cloud. Some of the key capabilities include:

Scalable Data Ingestion: Synapse Analytics can handle high volume and high velocity data streams, thanks to its integration with Azure Event Hubs, Azure IoT Hub, and Apache Kafka.

Real Time Data Processing: The Synapse platform utilizes Apache Spark Streaming to process data streams in real time, enabling complex transformations, aggregations, and analytical computations on the incoming data.

Unified Data Analytics: Synapse Analytics provides a unified platform for data engineering, data warehousing, and data exploration, allowing data teams to collaborate and work seamlessly across different stages of the data lifecycle.

Elastic and Scalable Infrastructure: The cloud based nature of Synapse Analytics ensures that the platform can scale up or down based on the changing data processing requirements, providing a highly elastic and scalable infrastructure for real time ETL workflows.

#### Synapse Analytics Architecture for Real Time data:

Synapse Analytics, the Microsoft Synapse platform, is designed to handle the challenges of real time data processing and streaming. The platform's streaming capabilities are powered by Apache Spark Streaming, which provides a robust and scalable framework for processing unbounded streams of data. Real time data processing is a key feature of Synapse Analytics which provides seamless integration with a variety of data sources, including IoT devices, web applications, and other cloud based services. This integration allows for the efficient ingestion of data streams, which can then be processed and transformed in near real time.

The architecture of Synapse Analytics for real time data processing and ETL workflows can be broadly divided into three main components: data ingestion, data processing, and data storage.

**Data Ingestion:** Synapse Analytics employs various mechanisms like Azure Event Hubs, Azure IoT Hub, and Apache Kafka to handle high volume and high velocity data streams.

1. Azure Event Hubs: A fully managed, real time data ingestion service that can ingest millions of events per second. *Example:* A financial services company uses Azure Event Hubs to capture and analyze stock price movements in real time, streaming data from various stock exchanges. This allows them to process and analyze the data within Synapse Analytics to provide insights for trading strategies.

- 2. Azure IoT Hub: A cloud service that manages bi directional communication between IoT applications and devices. Service based Example: A smart home service provider uses Azure IoT Hub to manage data from various smart devices installed in customers' homes. Devices like thermostats, security cameras, and smart lighting systems send real time data to IoT Hub, where it's ingested into Synapse Analytics. The service provider then analyzes this data to offer personalized recommendations, such as optimizing energy usage, enhancing security measures, and providing maintenance alerts to improve overall customer experience.
- 3. Apache Kafka: An open source stream processing platform for handling real time data feeds. Example: An e-commerce platform uses Apache Kafka to manage customer activity streams, capturing events like searches, clicks, and purchases. Synapse Analytics processes these streams in real time to personalize user experiences and recommend products based on consumer behavior.

Data Processing: The Synapse Analytics platform utilizes Apache Spark Streaming for the real time processing of data streams. Spark Streaming provides a scalable and fault tolerant framework for processing unbounded streams of data, with the ability to perform complex transformations, aggregations, and analytical computations on the incoming data. Example: A banking service company uses Apache Spark Streaming to process real time data from customer enrollment activities. As customers sign up through various channels (online, branch offices, mobile apps), the data is streamed into Spark Streaming. This data undergoes immediate transformations, validating customer information, detecting duplicates, and enriching the profiles with additional insights. Aggregated reports on new enrollments are generated in real time, enabling the bank to understand enrollment patterns, segment customers, and tailor marketing strategies promptly. Daily Transactions: Similarly, transactions from ATMs, POS systems, online banking, and mobile apps are streamed to Spark Streaming. It performs real time transformations and aggregations like calculating spending patterns, detecting potential fraud, and updating account balances. This allows for the creation of real time financial dashboards and alerts for any significant transactions, giving customers immediate insight and ensuring financial security.

Data Storage: Synapse Analytics integrates with various data storage solutions, such as Azure Blob Storage, Azure Data Lake Storage, and Azure SQL Data Warehouse, to provide a reliable and scalable data storage infrastructure for the processed data.

- 1. Azure Blob Storage: A highly scalable and cost effective object storage solution for unstructured data. It is ideal for storing raw and processed files, such as logs, images, and backup data, and can easily work with Synapse for further processing.
- 2. Azure Data Lake Storage: A scalable and secure data lake for big data analytics workloads. It allows processing and storage of vast amounts of structured and unstructured data, enabling advanced analytics and machine learning applications within Synapse Analytics.
- 3. Azure SQL Data Warehouse: A cloud based, fully managed data warehouse solution that provides MPP (Massively Parallel Processing) architecture. This facilitates the storage and analysis of large datasets, supporting high performance queries, and complex data transformations within Synapse.



Figure 1: Azure Synapse Streaming Analytics

#### Methodology:

To conduct this research, we will leverage the information provided in the given sources. [1][2][3]

This source discusses the need for a sophisticated and customizable platform that can handle large scale data streams and extract valuable insights through various analytical techniques. The paper emphasizes the importance of a multi layered architecture that includes data ingestion, processing, and storage components to address the challenges of real time data processing.

This source provides a broad overview of the common methods and techniques used for data analytics and processing platforms in Cyber Physical Systems, including the handling of structured, semi structured, and unstructured data streams with high velocity [2].

This source outlines the eight key requirements that a system software should meet to excel at a variety of real time stream processing applications. The requirements cover aspects such as high throughput, low latency, elastic scaling, and fault tolerance, which align with the capabilities of Synapse Analytics.

#### Findings

Based on the review of the provided sources, we can summarize the key findings:

- 1. The exponential growth in data volume and velocity has led to an increasing demand for real time data processing and analytics platforms that can handle high throughput and low latency data streams [5][4].
- 2. Modern data processing systems must deliver insights with minimal latency and high throughput, which has given rise to the streaming data analytics paradigm.
- 3. Effective real time data processing platforms should meet a set of key requirements, including high throughput, low latency, elastic scaling, and fault tolerance.
- 4. Synapse Analytics, with its integration of Apache Spark Streaming, Azure Event Hubs, and Azure SQL Data Warehouse, is well equipped to address the challenges of real time data streaming and processing in the cloud [1][4][5]

5. The unified and scalable architecture of Synapse Analytics enables efficient and seamless ETL workflows, allowing data teams to collaborate and work across different stages of the data lifecycle.

### **Conclusion:**

Synapse Analytics is a scalable cloud based platform that offers a comprehensive set of capabilities for real time data processing and ETL workflows. The platform's integration with various data ingestion mechanisms, such as Azure Event Hubs and Apache Kafka, allows for the efficient handling of high volume and high velocity data streams. The use of Apache Spark Streaming in Synapse Analytics enables real time data processing, with the ability to perform complex transformations, aggregations, and analytical computations on the incoming data. Its multi layered architecture, which encompasses data ingestion, processing, and storage components, allows it to effectively handle the challenges posed by high volume and high velocity data streams. Synapse Analytics' elastic and scalable infrastructure ensures that the platform can adapt to changing data processing requirements, providing a highly efficient and reliable solution for real time ETL. Overall, the features and capabilities of Synapse Analytics make it a compelling choice for organizations looking to implement scalable and robust real time data processing solutions in the cloud.

#### **References:**

- Bohlouli, M., Schulz, F., Angelis, L., Pahor, D., Brandić, I., Atlan, D., & Tate, A. R. (2020). "Towards an Integrated Platform for Big Data Analysis". In arXiv (Cornell University). Cornell University. https://doi.org/10.48550/arxiv.2004.13021
- Chitu, C., & Song, H. (2019)." Data analytics and processing platforms in CPS". In Elsevier eBooks (p. 1). Elsevier BV. https://doi.org/10.1016/b978-0-12-816637-6.00001-4
- Isah, H., & Zulkernine, F. (2018). "A Scalable and Robust Framework for Data Stream Ingestion". In 2021 IEEE International Conference on Big Data (Big Data) (p. 2900). https://doi.org/10.1109/bigdata.2018.8622360
- 4. Stonebraker, M., Çetintemel, U., & Zdonik, S. (2005). "*The 8 requirements of real-time stream processing*". In ACM SIGMOD Record (Vol. 34, Issue 4, p. 42). Association for Computing Machinery. https://doi.org/10.1145/1107499.1107504
- Tantalaki, N., Souravlas, S., & Roumeliotis, M. (2019)." A review on big data real-time stream processing and its scheduling techniques". International Journal of Parallel Emergent and Distributed Systems, 35(5), 571.Taylor & Francis. https://doi.org/10.1080/17445760.2019.1585848
- 6. Xhafa, F., Naranjo, V., Caballé, S., & Barolli, L. (2015). "A Software Chain Approach to Big Data Stream Processing and Analytics" (p. 179). https://doi.org/10.1109/cisis.2015.24