

Comparative Evaluation of Spark and Flink Stream Processing

Supervisor: PD Dr. Michael Mock

Contents

| | | |
|----------|---|----------|
| 1 | Introduction | 3 |
| 2 | Technical Background | 4 |
| 2.1 | Apache Spark | 4 |
| 2.2 | Apache Flink | 4 |
| 2.3 | Apache Kafka | 5 |
| 3 | Experiment Setup and Implementation | 5 |
| 3.1 | Data Stream Setup | 5 |
| 3.2 | Evaluation Streaming Workloads | 5 |
| 3.2.1 | Statistics Computation per Trajectory | 6 |
| 3.2.2 | Air Sector Change Detection | 7 |
| 4 | Performance Results | 8 |
| 5 | Conclusion | 8 |

Abstract

Recent years have witnessed a big development improvements of Big Data frameworks such as Apache Spark and Apache Flink. For example, both frameworks support real-time streaming processing. While which platform is the best still an open question. In this paper, we provide a performance comparison of streaming data processing in Spark and Flink, by measuring different performance metrics, namely latency and throughput. The experiment is done using multiple streaming workloads over real-world datasets.

1 Introduction

Big Data is now one of the most appealing topics in Information Technology, it has gained increasing attention of professionals in industry and academics, since the amount of data we generate is increasing exponentially [1] that generated from social media platforms, mobile phones, IoT and financial transactions, etc. Moreover, recent years showed a rapid development of Big Data frameworks such as Spark and Flink that provide more and efficient data processing capabilities.

The Big Data is a method to process not just only just large datasets (volume), but also that arrives in high rate (velocity), and it comes in all kinds of data format (variety) (i.e., 3Vs of Big Data (Volume, Velocity, Variety))[2].

The processing models of Big Data are Batch and Stream processing, Batch processing is a way to handle large, finite volumes of data (e.g., processing of historical transaction records), is more concerned with throughput. While the Stream processing is a method to manage fast and continuously incoming data (real-time), process it as it arrives (e.g., credit card fraud detection and network monitoring applications), where the low latency is required. The variety of data (structured and semi-structured) is included in the two processing models.

TODO: add paragraph about Yahoo benchmark and why do we still need evaluation 1- continuous dev of frameworks 2- explore more complex operation like groupby that shuffle the data.

The purpose of this paper is to provide a comparative experimental evaluation of throughput, latency (i.e., performance) of stream processing in Apache Spark and Apache Flink. We use a datasets of airplanes trajectories provided in the context of Datacron project¹, in order to simulate real-world use cases. Using multiple streaming data workloads we try to cover the main programming API differences between the both frameworks. Informally, this work aims to help developer to determine which platform to choose in production.

The remainder of this paper is organized as follows. In Section 2, we present the main characteristics, APIs, and main data abstraction of Spark and Flink. Section 3 describes the experiment workloads and their implementation in each framework. Section 4 provides and compares the performance results of the different workloads. Finally, Section 5 gives the overall conclusion.

¹<http://www.datacron-project.eu/>

2 Technical Background

In this section, we present the architecture, key characteristics, programming APIs of Apache Spark and Apache Flink. Moreover, a brief overview of the Apache Kafka platform is provided.

2.1 Apache Spark

Apache Spark is an open source project that provide general framework for large-scale data processing [3]. It offers programming API in Java, Scala, Python and R. Its stack includes set of built-in modules including Spark SQL, MLlib for machine learning, GraphX for graph processing, and Spark Streaming to process real-time data streams. It can access different data sources including Hadoop Distributed File System (HDFS), Apache Cassandra database, Apache HBase database etc.

The main data abstraction in Spark core is the Resilient Distributed Datasets (RDDs), which is in-memory collection of elements partitioned across cluster of computers that can be processed in parallel [4]. RDDs supports two categories of operations: transformations that drive new RDDs from the operated ones (e.g., map), and actions which return a certain value (e.g., count). On the other hand, the Discretized Stream (DStream) is the data abstraction in Spark Streaming, which is series of RDDs that represent an input stream of data [5]. In other words, the Spark Streaming is batch-based processing model, which process the continuous stream of data by divided into micro-batches to be processed by the Spark engine. Figure 1 illustrates the general process flow of Spark Streaming.

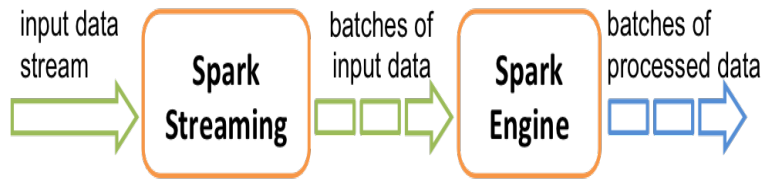


Figure 1: Process flow of Spark Streaming [5].

2.2 Apache Flink

Apache Flink is an open source project that provide large-scale, distributed stream processing platform, while the batch processing treated as special case of streaming applications [6]. It offers programming API in Java, Scala. Its stack includes the core DataStream and DataSet APIs with additional libraries

such as Complex event processing for Flink (FlinkCEP), Machine Learning for Flink (FlinkML), and Flink Graph API (Gelly).

The main data abstraction of Flink are `DataStream` and `DataSet` that represents read-only collection of elements of the same type. The list of elements is bounded (i.e., finite) in `DataSet`, while it is unbounded (i.e., infinite) in case of `DataStream`.

2.3 Apache Kafka

Apache Kafka is scalable, fault-tolerant distributed publish/subscribe streaming framework [7]. It manages the stream records in different categories (i.e., topics) that are partitioned and distributed over the servers in the Kafka cluster. It allows the data producers to publish stream of records to certain Kafka topic. While the consumer applications can subscribe to one or more topic to read the data stream. It has been widely adopted, for example, Spark and Flink can receive data stream from Kafka.

3 Experiment Setup and Implementation

This section describes the data stream setup, design of evaluation streaming workloads, and implementation details in Spark and Flink.

3.1 Data Stream Setup

Our streaming workloads read input data stream from Kafka. In order to simulate real-world use cases we use datasets of Automatic Dependent Surveillance Broadcast (ADS-B) messages that present the aircraft’s position over time, comprise 22 fields of data such as aircraft ID, date message generated, longitude, latitude, and altitude.

In our experiment a data stream producer component reads the ADS-B datasets, then publishes it to Kafka cluster to be consumed by the workloads in Spark and Flink. Figure 1 illustrates the data producer and the Kafka cluster setup. The data stream is partitioned into four partitions over two servers, while the data producer publishes the stream records randomly to Kafka partitions.

3.2 Evaluation Streaming Workloads

In our experimental evaluation, we developed two real-time stream processing workloads that read the ADS-B messages stream from Kafka in order to perform basic trajectories analysis methods. The workloads are designed in such way that cover general key aspects of real-time streaming processing tasks, and evaluate the corresponding solutions in Spark and Flink. The following are the aspects of streaming data processing covered by our workloads’s design:

- Handling parallel input streams (e.g., Kafka Stream).

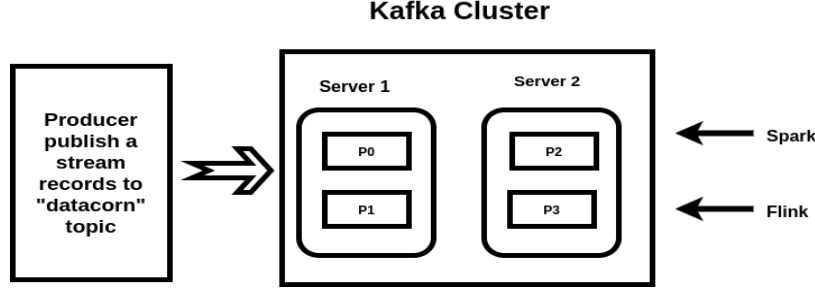


Figure 2: Data Stream Producer & Kafka Cluster Setup.

- How to aggregate the state of input stream.
- Manage the order of stream records.
- How to provide and update global data model in stream processing task.
- Evaluate the performance by measuring the latency and throughput.

We present the description of workloads and the relation with these aspects, and th implementation details of Spark and Flink solutions in Section 3.2.1 & 3.2.1. Afterward, In Section 4 we discuss the performance evaluation and analysis of workloads.

3.2.1 Statistics Computation per Trajectory

In first stream processing workload, we construct stream of trajectories by considering the aircraft’s positions (i.e., ADS-B messages) that belong to the same aircraft as trajectory. Moreover, we continuously compute and aggregate statistics for each arriving position in a trajectory . As example of computed statistics indexes speed mean, mean of location coordinates, min and max altitude, etc. In this workload, we cover the parallel receiving of input data stream, state-full aggregation on input data stream, and preserving the right order of the stream’s records. The implementation details of this workload in Spark and Flink are presented in the following:

- **Implementation in Spark**

The Spark Streaming implementation of trajectory statistics computation read the Kafka stream using `KafkaUtils.createDirectStream` that creates `DStream` instance, in order to scale up and have parallel Kafka stream receivers, a multiple instances of `DStream` are created and combined together using `union` operation. The output stream of the `union` is processed using `mapToPair`, to map each position message to tuple of aircraft ID and position message. Then, the irrelevant tuples are filtered based on the position message type using a `filter` transformation. Afterward, the tuple with same aircraft ID are groped together to construct

trajectories (tuples of Id and list of positions), by applying the `groupByKey` operation. Since streaming model of Spark is micro-batch based and stateless, `updateStateByKey` operation is must be used to manage the state between the batches. In context of this workload, the statistics computation is preformed within the custom `updateStateByKey` function, which uses the last position of trajectory from previous batch to aggregate and compute statistics values to each trajectory’s position in the new arriving batch. The list of new positions must be sorted to preserve the correct order, since the tuples after `groupByKey` operation are shuffled across the cluster’s nodes.

- **Implementation in Flink**

The detail of Flink’s implementation of statistics computation workload as follows. First, a `FlinkKafkaConsumer` is used to read the Kafka stream that constructs `DataStream` instance, moreover, it implicitly runs multiple parallel Kafka stream consumers, which makes the Kafka partitions balanced between the instances. Second, the stream records (i.e., ADS-B messages) are parsed and transformed to tuples of ID and position message using a configured `map` transformation. Third, a `filter` function is used to filter irrelevant tuples based on the message type. Fourth, the result tuples after the filtration steps are grouped with similar ID using `keyBy` operator, to constructs the stream of trajectories. Finally, a `reduce` transformation on the trajectories tuples is used to calculate the statistics for new received trajectory’s position, utilizing the aggregated statistics of previous position, while old position is discarded and the new record is used to aggregate the statistics.

//TODO: discuss differences between Flink and Spark.

3.2.2 Air Sector Change Detection

The goal of sector change detection workload is to detect the enter or leaving of aircraft from one air sector to other, by processing the real-time stream of trajectories. Given that the dataset of sectors (i.e., polygons) is available as reference to assign the corresponding sector for certain aircraft’s position. In this workload, we test the aspect of providing the global model of data (sectors dataset) in streaming processing workload. Moreover, the preserving od stream’s records order and consuming of parallel Kafka streams points were covered.

The implementation details of this workload in Spark and Flink are presented in the following:

- **Implementation in Spark**

We implemented the sector change detection workload in Spark by constructing the trajectories stream and assigning the sectors for each position in all trajectories, then a change between the assigned sectors of two consecutive aircraft’s positions is marked as alert case. First, a multiple

instances of `DStream` are created using `KafkaUtils.createDirectStream` to read in parallel from Kafka, after all input data streams are combined using `union` function. Second, the records (i.e., ADS-B messages) of the resulted union stream are parsed and mapped to tuples of (aircraft ID, position message) using defined `mapToPair` transformation. Third, the unrelated tuples are filtered based on the message type through `filter` operation. Fourth, the trajectories stream is constructed by combining all tuples with the same ID using the `groupByKey` function. Then, to manage and aggregate the state of trajectories stream between the mini-batches a configured `updateStateByKey` function is applied, the corresponding sector are assigned, the old positions of certain trajectory from previous batch are discarded, but the last position of old batch is kept to be used alongside with the new positions. The sector data set is provided using the `Broadcast` shared value feature in Spark that distributes the data set to all cluster nodes efficiently. Finally, the detection of sector change is done by extra filtering transformation using `filter`, by checking the sector of the positions list in each tuple.

- **Implementation in Flink**

The Flink implementation of sector change detection workload, uses the `FlinkKafkaConsumer09` that read the ADS-B messages stream from Kafka, which used to construct `DataStream` instance. The `DataStream` records are parsed to tuples of aircraft ID and aircraft's position, by a configured `map` function. The irrelevant tuples are filtered based on the message type using a `filter` operation. Afterward, the tuples with common IDs are grouped together by the `keyBy` operator. Since the Flink streaming processes the stream records as it come, a defined `reduce` operator assign corresponding sector of new tuple, and the old tuple with same ID is used to retrieve the previous sector to attach it to the new tuple, while the old tuple is discarded. Finally, the tuples with different previous and current sector that represent change in sector case, they are filtered using defined `filter` transformation.

4 Performance Results

This section provides and compares the performance results of the different workloads.

5 Conclusion

This section gives the overall conclusion.

References

- [1] Gantz, John, and David Reinsel. "The digital universe in 2020: Big data, bigger digital shadows, and biggest growth in the far east." IDC iView: IDC Analyze the future 2007.2012 (2012): 1-16.
- [2] Laney, D 2001 3D Data Management: Controlling Data Volume, Velocity, and Variety. META Group.
- [3] Apache Spark. Available: <https://spark.apache.org/> [Accessed February 2017].
- [4] Zaharia, Matei, et al. "Resilient distributed datasets: A fault-tolerant abstraction for in-memory cluster computing." Proceedings of the 9th USENIX conference on Networked Systems Design and Implementation. USENIX Association, 2012.
- [5] Spark Streaming. Available: <http://spark.apache.org/docs/latest/streaming-programming-guide.html> [Accessed March 2017].
- [6] Apache Flink. Available: <https://flink.apache.org/> [Accessed February 2017].
- [7] Apache Kafka. Available: <https://kafka.apache.org/> [Accessed February 2017].