# Apache Flink: State, checkpointing and fault tolerance

Kommineni Anvitha

11708336

Lovely Professional University

Abhiram Kumar

11709022

Lovely Professional University

Apurva

11709264

Lovely Professional University

Deeksha Tripathi

11710413

Lovely Professional University

*Abstract*— **Apache Flink is an open-source system for processing streaming and batch data processing platform. It is a distributed to the data processing platform for use in big data applications, primarily involving analysis of data stored in Hadoop clusters. Flink can be easily deployed on a single Java virtual machine (JVM) in standalone mode or YARN-based Hadoop clusters, or on cloud systems Flink is built on the philosophy that many classes of data processing the applications and including the real-time analytics, continuous data pipelines, historic data processing that which means to the batch processing, and iterative algorithms that machine learning, graph analysis can be to expressed and executed as pipelined the fault tolerant data flows. The present Flink's architecture and expand it on how to a diverse the set of use cases can be done to unified under a single of the execution model.**

## I. INTRODUCTION

Data-stream process as example and identified by advanced event process systems and static batch information processing as exemplified by MPP databases and Hadoop were historically thought. They were programmed terribly programming models and API's, and were executed by terribly systems dedicated streaming systems like Apache Storm, IBM Information sphere Streams ,Microsoft Stream Insight, and Stream base versus relative databases or execution the engines for Hadoop, as well as Apache Spark and Apache Drill.The Batch processing information analysis the created up for the share of the employment cases, data to the sizes, and market, whereas streaming information analysis principally to the served specialised. These

continuous streams of information return for instance from net logs, the sevral application logs, sensors, or as changes to application state in databases transaction log records. Instead of treating the streams as streams to the, setups ignore the continual and timely nature of information production.instead the, information records square measure often artificially batched into static information sets hourly, daily, or monthly then processed in a very time agnostic the fashion information assortment tools,and advancement managers, and schedulers the creation and process of batches, in really never-ending processing pipeline. Patterns like the "lambda architecture" mix batch and stream process systems to implement multiple ways of computation: a streaming quick path for timely approximate results, and of these from high latency imposed by batches.



## II. FAULT TOLERANCE

Fault tolerance Flink offers reliable and exempted for the execution with strict exactly to the once data processing for

the consistency for guarantees and deals with the failures through the near check pointing and partial re execution.The final assumption the system to makes to effectively offer these guarantees is that the sources area unit persistent and re playable. Samples of such sources area unit files and in follow, non persistent supplies that may be incorporated by keeping a write ahead log to for among the state of the source operators .The check pointing mechanism of Apache Flink builds on the notion of the distributed consistent snapshots to attain exactly once processing guarantees. The maximum infinite nature of a knowledge stream makes re computation upon recovery to the improper impractical, as months of computation can to be replayed for a long running job. To sure the recovery time, Flink takes a snap of the state of operators, as well as the present the position of the input for streams at regular intervals. The core challenge lies in taking a standardized for the snap of all parallel operators while not halting the execution of the topology. In essence, the snap of all operators ought to discuss with identical logical time within the computation.

The mechanism employed in the Flink is named Asynchronous Barrier Snapshotting .Barriers area unit management records injected into the input streams that correspond to a logical and approach time and logically separate the stream to the in part whose effects are enclosed within the current snap and therefore the half which will be snapshotted to later. Operator receives barriers from upstream associated first performs an alignment part, ensuring then, the operator writes for its state contents of a window, or custom information structures to sturdy storage the storage backend may be associate external system like HDFS.

Once the state has been protected from the faults, the operator forwards the barrier down stream register a snap of their formal state and a worldwide snap are complete to the orienting part to use the all their effects to the operator states.

This guarantees that the information that must be written in to reliable the multiple storage is unbroken to the theoretical minimum and the present state of the operators. Recovery is

from the failures and reverts all operator states to their several states for taken from the last booming snapshot and restarts the input streams ranging from the most recent barrier that the utmost quantity of re computation required upon recovery is proscribed to the quantity of input records between to consecutive barriers, moreover the partial recovery of a unsuccessful subtask is feasible by in addition replaying unprocessed.
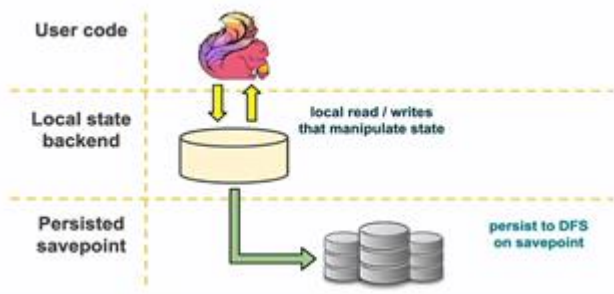
## III. STATE

When the state operators that may contain any form of state, this state must be the part of the snapshots as well. Operator state comes in different forms:

- User defined state: This is the user defined state that is created and modified directly by the transformation functions like map() or filter() operators. User defined state can either be a simple variable in the function's java object, or the associated key or value state of a function.

- System state: This state refers to the data buffers that are part of the operator's computation of streaming. A typical example for this state that are the window buffers, inside which the system collects and aggregates records for the windows until the window that has to been it is evaluated and evicted.

Operators snapshot to their state at the point are in the time of streaming when they are received and got the all snapshot barriers from their input streams, before that emitting or existing the stream barriers to their output streams. At that point, for all updates and varities to the state from records that before the barriers will have been made, and no updates that depend on records from to after the barriers have been applied. Because to the state of a snapshot may be potentially large, and it is stored in a configurable state backend. By the default, this is the Job Manager's memory, but for serious setups, a distributed reliable and effectable storage should be configured such as HDFS. After the state has been stored, the operator that

acknowledges the checkpoint, emits the snapshot to barrier into the output streams, and proceeds.



The resulting snapshot now contains:

- For each of the parallel streaming the data source, the offset or position in the stream when the snapshot was started.

- For each operator, a pointer should be the state that was stored as part of the snapshot

Programs written within the information Stream API usually hold state in numerous forms:

1.Windows gather components or aggregates till they're triggered.

2.Transformation functions could use the key/value state interface to store values.

3.Transformation functions could implement the Check pointed Function interface to create their native variables fault tolerant  conjointly state section within the streaming API guide.

When check pointing is activated, such state is persisted upon checkpoints to protect against information loss and recover systematically , and the way and wherever it's persisted.

## IV.  AVAILABLE STATE BACKENDS

Out of the box, Flink bundles situated in the these state backends:

- Memory State Backend

- FsState Backend

- Rocks DBState Backend If nothing else is configured, the system will use the MemoryStateBackend.

1.RocksDBState Backend : If the nothing else is designed in the DBState, the system can use the MemoryStateBackend.It will be use the memory backend state for the streaming.

The MemoryStateBackend holds information internally as objects on the Java heap.Upon checkpoints, this state backend can exposure the state and send it as a part of the stop acknowledgement messages to the  Job Manager (master), that stores it on its heap likewise. The MemoryStateBackend may be designed to use asynchronous snapshots. Out of the box, Flink bundles these state are the backends. Key or the value state and window operators hold hash tables and  that store the values and the triggers, etc. Whereas we have a tendency to powerfully encourage the employment of asynchronous snapshots to avoid block pipelines, please note that this can be presently enabled by default.To disable this feature, users will instantiate . MemoryStateBackend  with the corresponding Boolean flag within the creator set to false this ought to solely used for debug. The size of every individual state is by default restricted to 5 MB. This price may be multiplied within the creator of the MemoryStateBackend. Irrespective of the designed top state size, the state can't be larger than the town frame size see configuration. The aggregate state should match into the Job  Manager  memory. The MemoryStateBackend is inspired for. Jobs that do hold very little state of processing ,they may like jobs that consist solely of record at a time functions like Map, Flat Map,and the Filter.

The FsStateBackend is designed with a classification system address (type, address, path), like

The FsStateBackend holds in-flight information within the TaskManager's memory.Upon checkpointing, it writes state snapshots into files within the designed classification system and directory. Nominal information is keep within the Job
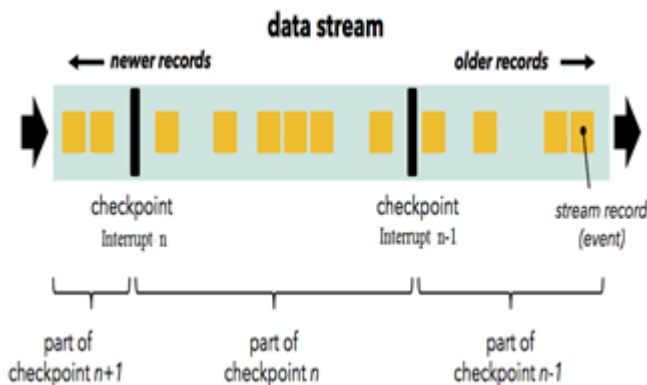
Manager's memory in high-availability mode, within the information checkpoint. The Fs State Backend uses asynchronous snapshots by default to avoid block the process pipeline whereas writing state checkpoints.

The FsStateBackend is encouraged for:

- Jobs with large state and having long windows and large key or the value states.

- All high availability and reliability setups.

## V. CHECKPOINT

Every function and operator in Flink can be **stateful**. Stateful and multi functions store data across the processing of individual elements/events, making the state a critical building block that for any type of more elaborate operation .In order to make the state fault tolerant then the, Flink needs to **checkpoint** the state. Checkpoints allow Flink to recover state to original state and there original positions in the streams to give the application the same semantics as a failure free execution helps to start from restart the streaming and processing. The documentation on streaming that the fault tolerance describes in detail the technique behind Flink's streaming fault tolerance mechanism.



## VI. ENABLING AND CONFIGURATION CHECKPOINTING

The default of check pointing is disabled. To enable the check pointing and call enable Check pointing on the Stream Execution Environment, where $n$ is the checkpoint interval in milliseconds. Other parameters for check pointing include:

- *Exactly once vs. at least once*: There can optionally pass a mode to the enable Check pointing method to choose between the two guarantee levels of the checkpoints. Exactly once is preferable and the useful for the most applications. At least once there may be relevant for certain super low latency consistently few milliseconds applications and configuration of checkpoints

- *Checkpoint timeout*: The check point time after which a checkpoint in progress is aborted and existed, if it did not complete by then the checkpoint is the timeout the data streaming and processing the checkpoint duration have extended and the checkpoint interval.

- *Minimum time between checkpoints*: To make sure that the streaming application makes a certain amount of progress that are between checkpoints, one can define the how much time needs to pass between checkpoints. If this value is set the next checkpoint will be started no sooner than five seconds after the previous checkpoint completed, the regardless of the checkpoint duration have extended and the checkpoint interval. Note that this implies that the checkpoint interval will never be smaller than this parameter of the minimum time between the checkpoints.

  It is often easier to configure the applications on by defining to the time between the checkpoints than the checkpoint interval and space , then the because the time between checkpoints is not susceptible to the fact that checkpoints that may sometimes that they take longer than on average .Note that this value also implies that the number of concurrent checkpoints is rather than it's *1*.

- *Number of concurrent checkpoints*: By default of, the concurrent system check point will not trigger another checkpoint while one is still in the

concurrent progress. This ensures that the topology and structure does not spend too much time on checkpoints and not make progress that with processing the streams. It is possible to allow for multiple overlapping to the checkpoints, which is interesting for pipelines that have a certain processing delay for example because the functions call external services that need some time to respond but that still they want to do very frequent checkpoints 100s of milliseconds to reprocess very little upon failures.

This option cannot be used to see the ,when a minimum time between checkpoints is defined.

- *Externalized checkpoints*: You can configure that the periodic checkpoints to be persisted externally. Externalized checkpoints will write their meta data out to persistent the storage and might are *not* automatically cleaned up when the job fails. This way, you will have a checkpoint around to resume that is from if your job fails. There are more details in the deployment and manage the notes on externalized checkpoints that are registered.

- *Fail or continue task on checkpoint errors*: This determines if a task will be failed if an error occurs in the execution of the task's checkpoint procedure. This is the default are used for behaviour of the fault checkpoint error. Alternatively, when this is disabled and dislocated, the task that will simply decline the checkpoint to the checkpoint that are fail to coordinator and continue running the streaming.

- *Prefer checkpoint for recovery*: This determines the prefer checkpoint for recovery a job checkpoint for recovery from that will fall back to latest checkpoint for even when there are the more of the recent save the points that are available to potentially is reduce to the recovery time.

## VII. CONCLUSION

Flink simplifies the parallel analysis of large amounts of data by using traditional database techniques such as automatic optimization and declarative query languages. Expressive, intuitive APIs allow for batch as well as data stream processing. Flink is scalable and versatile due to ithigh compatibility. Operators are processed in a pipeline,parallel, and free of limitations caused by micro-batching techniques.

## REFERENCES

[1] https://data-flair.training/blogs/flink-tutorial/

[2] https://flink.apache.org/flink-architecture.html

[3]https://www.google.com/search?q=fault+tolerance+in+apache+flink+images&oq=fau&aqs=chrome.0.69i59l2j69i57.3182j0j7&sourceid=chrome&ie=UTF-8

[4] https://ci.apache.org/projects/flink/flink-docs-stable/dev/stream/state/checkpointing.html