# Aditi Verma

**Sr Software Engineer**

@aditiverma89 <averma@branch.io>

# Ramesh Shanmugam

**Sr Data Engineer**

@rameshs01 <rshanmugam@branch.io>

branch

# Agenda

- Background

- Moving data with Flink @ Branch

- Scale & Performance

- Flink on Kubernetes

- Auto Scaling & Failure Recovery

branch

# branch

**12B** requests per day (+70% y/y)

**3B** user sessions per day

**50 TB** of data per day

**200K** events per second

**60+** Flink pipelines
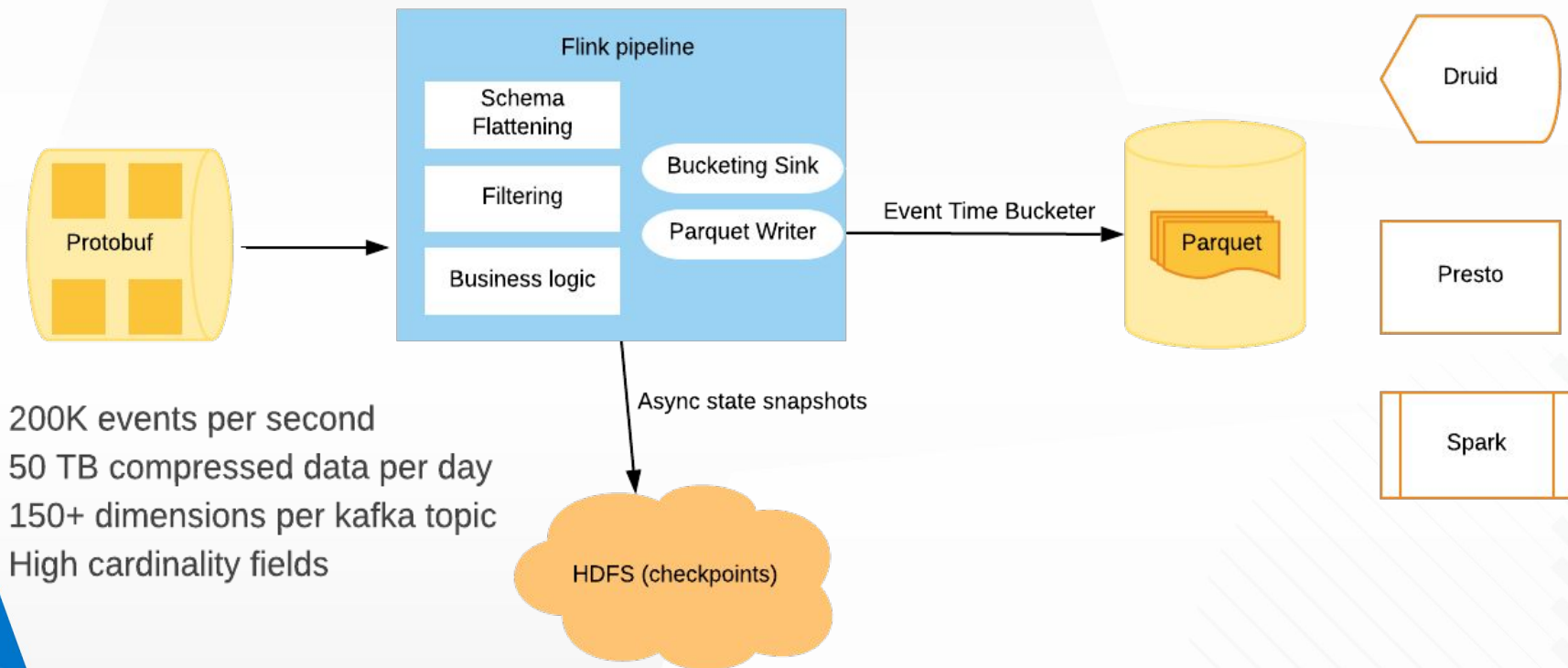
**5+** Kubernetes cluster

# Moving data with Flink @ Branch

*"Life is 10% what happens to you and 90% how you react to it."*
**— Charles R. Swindoll**

**Receive information**
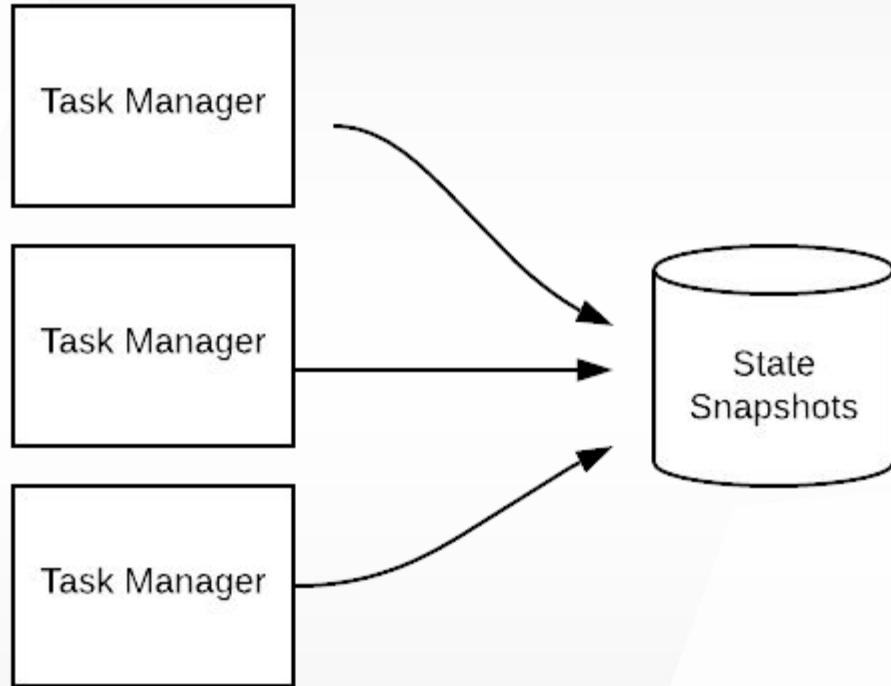**Process it**
**React to it**

**FAST!!**

branch

# Flink @ Branch

**Protobuf**

**Flink pipeline**
- Schema Flattening
- Filtering
- Business logic
- Bucketing Sink
- Parquet Writer

Event Time Bucketer

**Parquet**

Druid

Presto

Spark

Async state snapshots

HDFS (checkpoints)

200K events per second
50 TB compressed data per day
150+ dimensions per kafka topic
High cardinality fields

# State Backend



- Relatively small state backend
- File system backed state

# Parquet

- Higher compression
- Read heavy data set: ingested to Druid and Presto (3M+ queries/day)
- Avro data format
- Memory intensive writes

# Writing parquet with Flink

Two approaches:

1) Close the file with checkpointing

# Writing parquet with Flink

Two approaches:

- a) ~~Close the file with checkpointing~~
- b) Bucketing file sink
    - i) Configured with custom event-time bucketer, parquet writer and batch size
    - ii) Files are rolled out with a timeout of 10 min within a bucket

# Performance and Scale

- 100% traffic increase each year
- Higher parallelism impacts application performance and state size
- Kafka partitions < Flink parallelism requires rebalance on the input stream
- Task manager timeouts

Performance

Scalability          Stability

# Analyzing memory usage

❖ Network Buffers

❖ Memory Segments

❖ User code

❖ Memory and GC stats

❖ JVM parameters

# Containerizing Flink - Mesos

- Longer start-up time on Mesos

- Moved to containerizing Flink application on Kubernetes

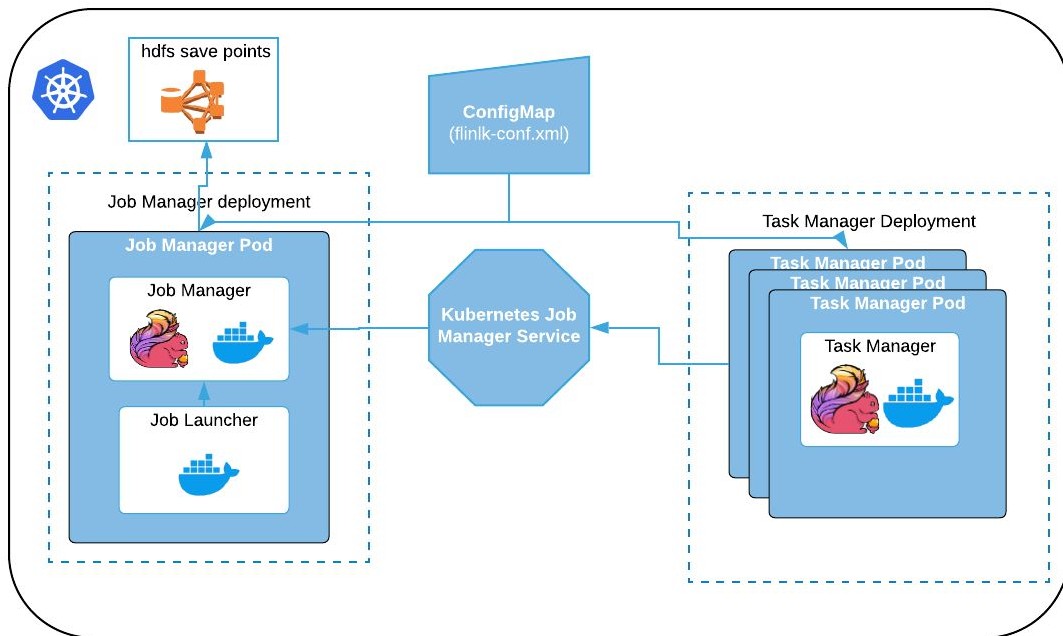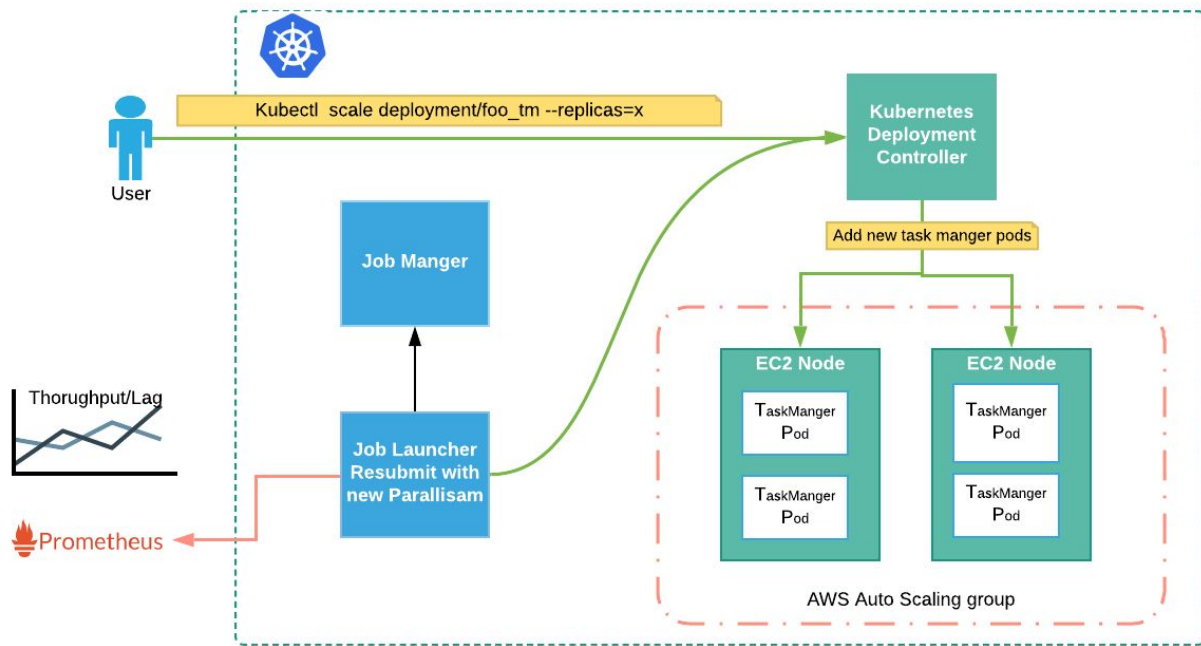- Kubernetes is resource oriented, declarative

# Kubernetes Terms

# Flink on Kubernetes @ Branch



- Single job per cluster
- Docker image
  - flink image - Task manager + job manager
  - Job launcher - custom launcher + job jar
- Job launcher
  - Application jar
  - Uploads jar
- Config map - flink config.xml
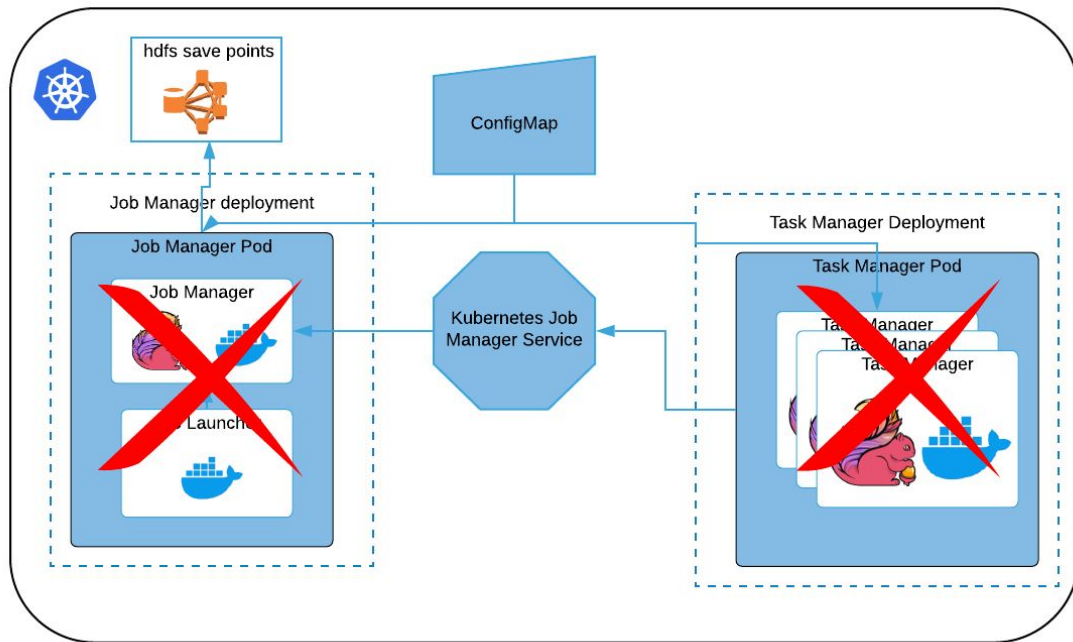  - jobmanager.rpc.address

# Auto Scaling



- When & How much scale
  - Auto - Joblauncher
- Scale
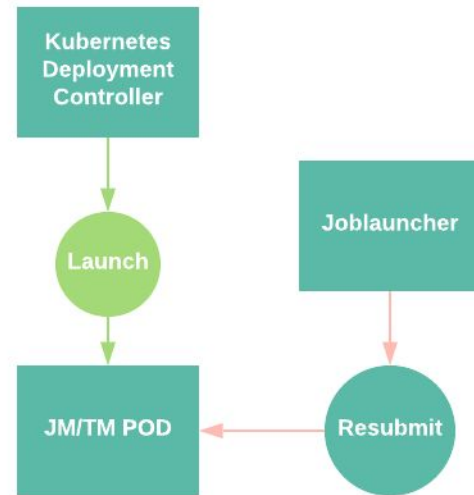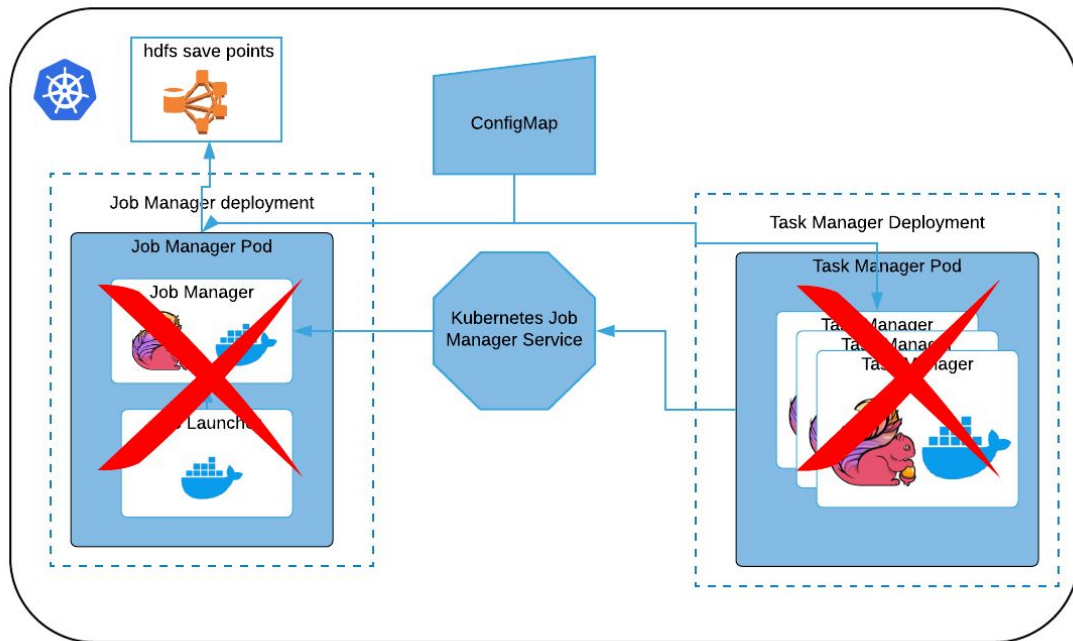  - Replica Set
- Flink job with new parallelism
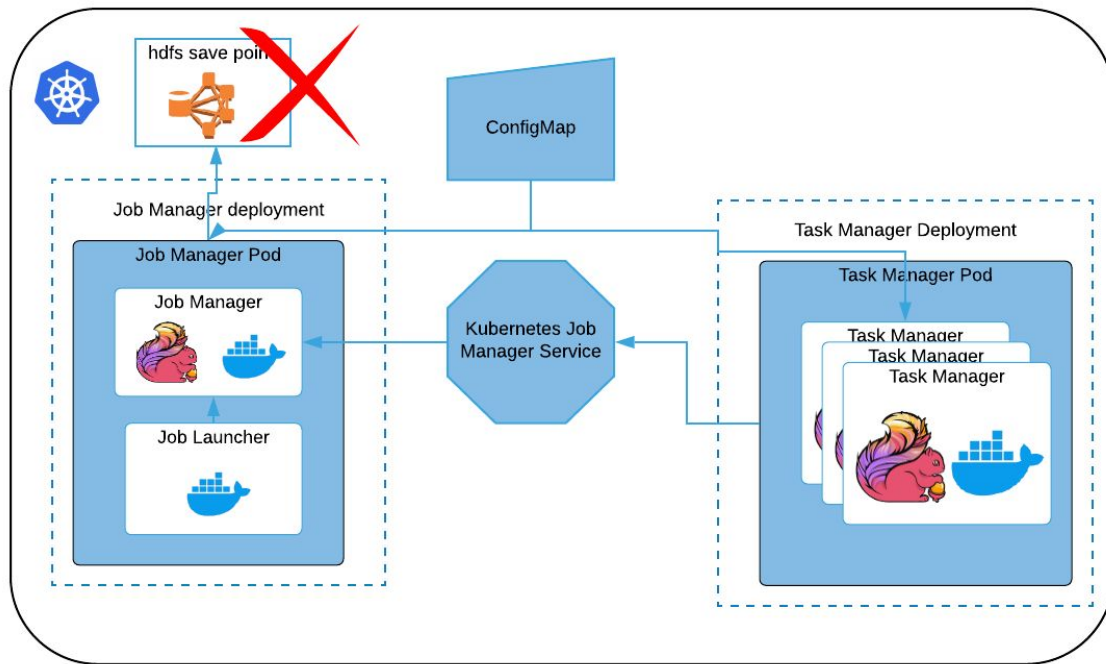
# Failure Recovery

branch

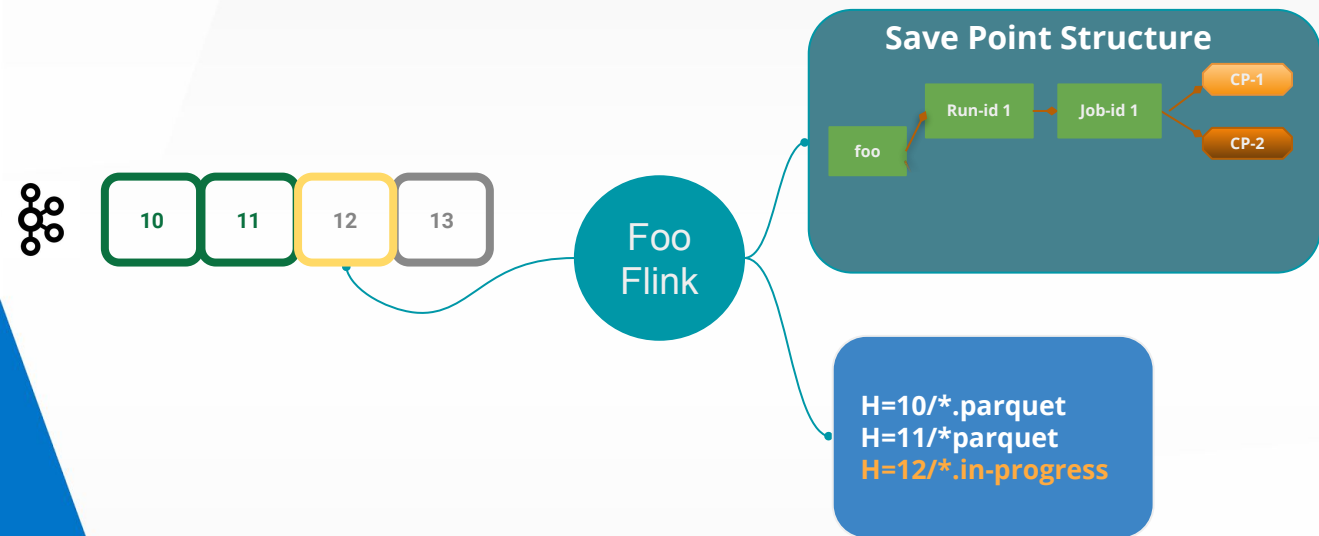# Job / Task Manager Goes Down?

# Job / Task Manager Goes Down?

# Savepoint Failure



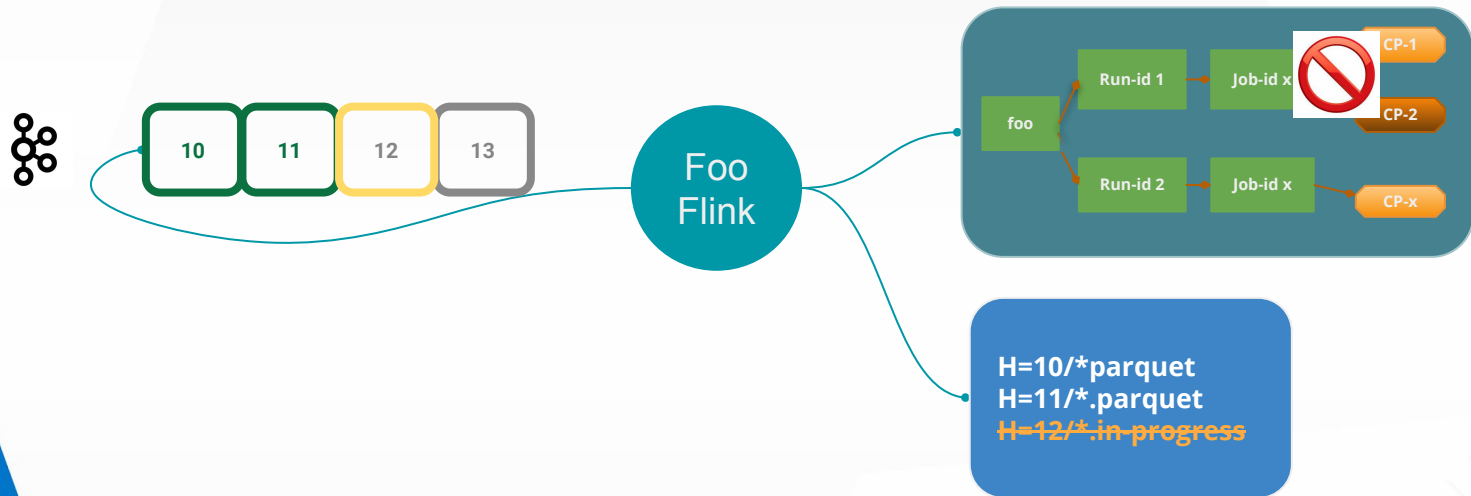- Reasons
  - Truncation
  - Schema mismatch
  - Hdfs outage

# Savepoint Structure



**Save Point Structure**

foo → Run-id 1 → Job-id 1 → CP-1, CP-2

H=10/*.parquet
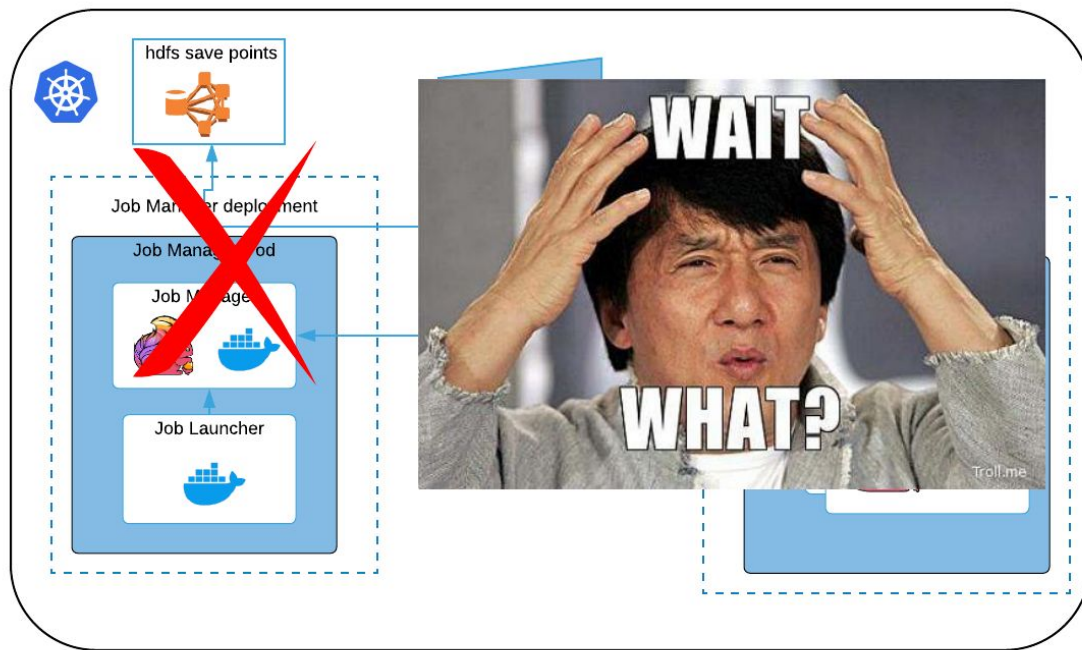H=11/*parquet
H=12/*.in-progress

- job/run-id/flink-job-id/cp-x

- Run id - incremental number

- Job id - flink job name

# Savepoint failure recovery

# Auto Recovery does not work?



- Continuous monitoring and proper alerts
- start job from latest offset
- Have different backfill route

# Next Steps….

- Parquet memory consumption (when too many buckets open)
  - Window + Rocks db => Parquet
  - Two stage process
    - row oriented streaming
    - batch to convert columnar