

Exactly Once vs At Least Once (Dataflow)

Ihaffa Murtopo
Data Engineer
Google



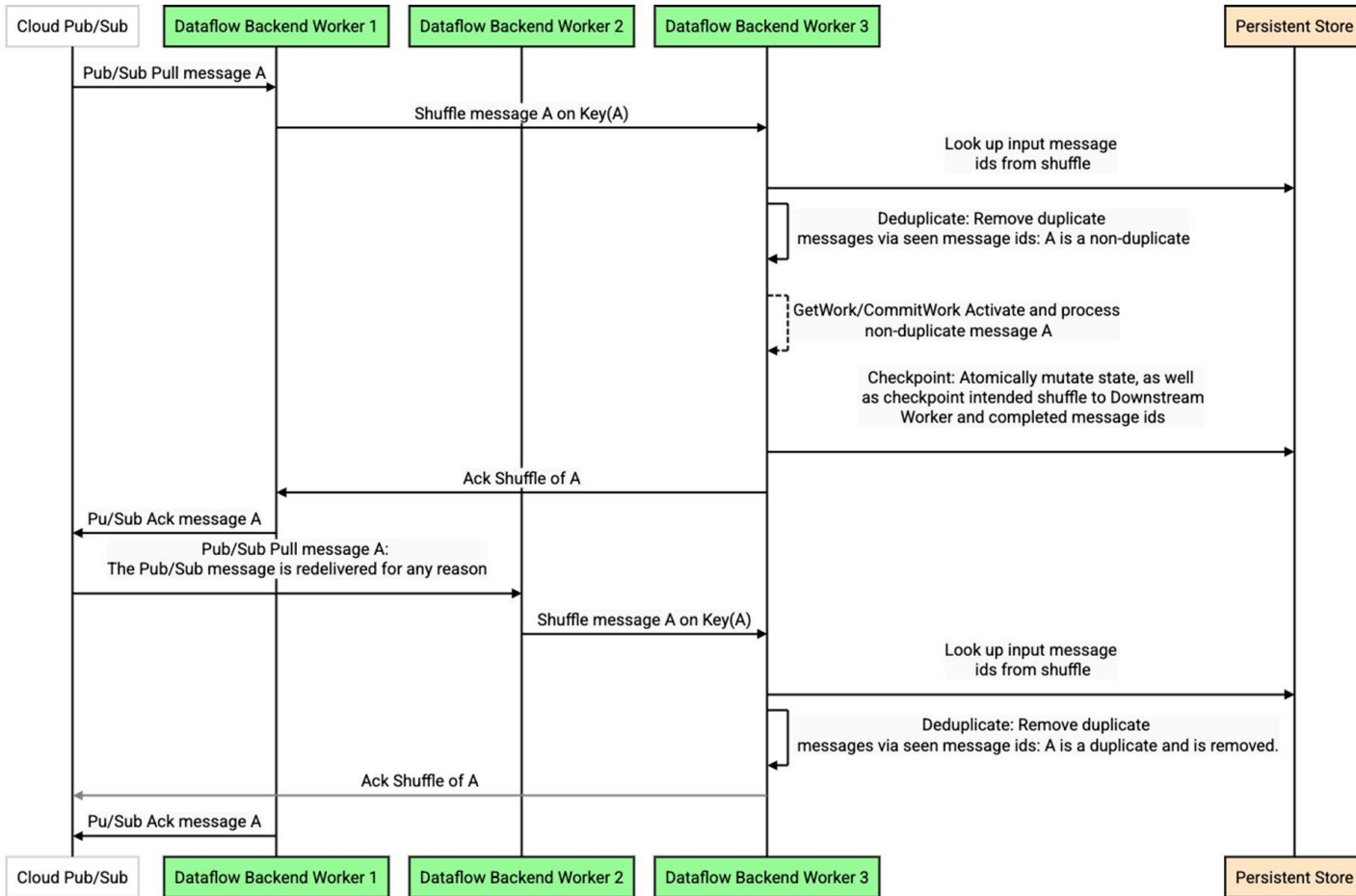
BEAM
SUMMIT

September 4-5, 2024
Sunnyvale, CA. USA

What Is?

Feature	Exactly-Once (Default)	At-Least-Once (New)
Definition	Each record is processed exactly one time, even in the presence of failures. No duplicates or data loss.	Each record is processed at minimum one time, even if there are failures. Duplicates possible.
Advantages	<ul style="list-style-type: none">• Ensures data accuracy	<ul style="list-style-type: none">• Lower Cost• Higher Throughput



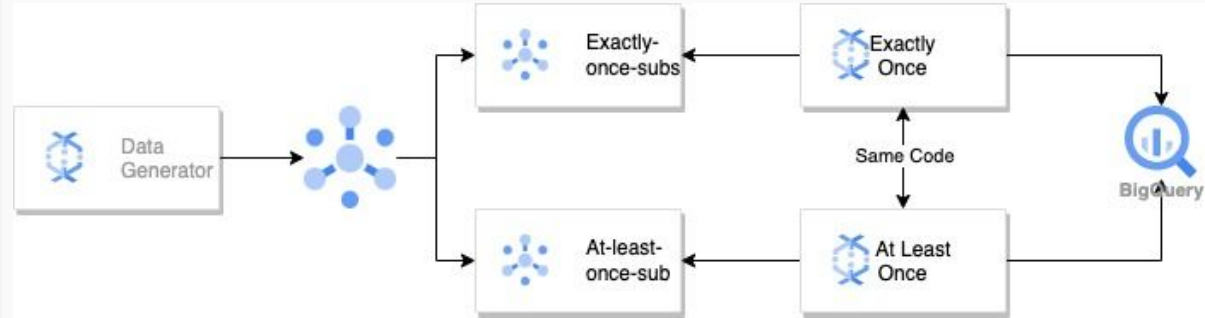


Summary:

- Deterministically assign key to Pubsub Message
- Require to constant interaction with Persistent Store
- If a particular key or worker is slow, impacts unrelated key

Setup Benchmark

- Consider using Dataflow template: Streaming Data Generator
- BQ or Table to count number of duplicates



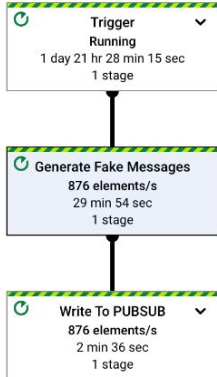
Run At-Least-Once

- Add “streaming_mode_at_least_once”
- Exactly-Once is default configuration

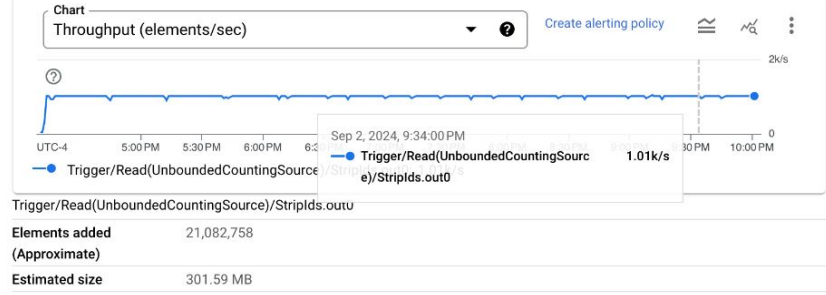
```
mvn compile exec:java -Dexec.mainClass=org.apache.beam.examples.Demo \
-Dexec.args="--runner=DataflowRunner --gcpTempLocation=gs://ihaffa-sandbox
--usePublicIps=false --subnetwork=https://www.googleapis.com/compute/v1/pr
--project=ihaffa-sandbox --region=us-central1 \
--jobName=demo-exactly-once \
--maxNumWorkers=30 \
--streaming --enableStreamingEngine=true \
--subsName=projects/ihaffa-sandbox/subscriptions/exactly-once-subs \
--outputTable=ihaffa-sandbox.beam_demo.exactly_once" \
-Pdataflow-runner
```

```
mvn compile exec:java -Dexec.mainClass=org.apache.beam.examples.Demo \
-Dexec.args="--runner=DataflowRunner --gcpTempLocation=gs://ihaffa-sandbox
--usePublicIps=false --subnetwork=https://www.googleapis.com/compute/v1/pr
--project=ihaffa-sandbox --region=us-central1 \
--jobName=demo-at-least-once \
--maxNumWorkers=30 \
--streaming --enableStreamingEngine=true \
--subsName=projects/ihaffa-sandbox/subscriptions/at-least-once-subs \
--outputTable=ihaffa-sandbox.beam_demo.at_least_once \
--dataflowServiceOptions=streaming_mode_at_least_once" \
-Pdataflow-runner
```

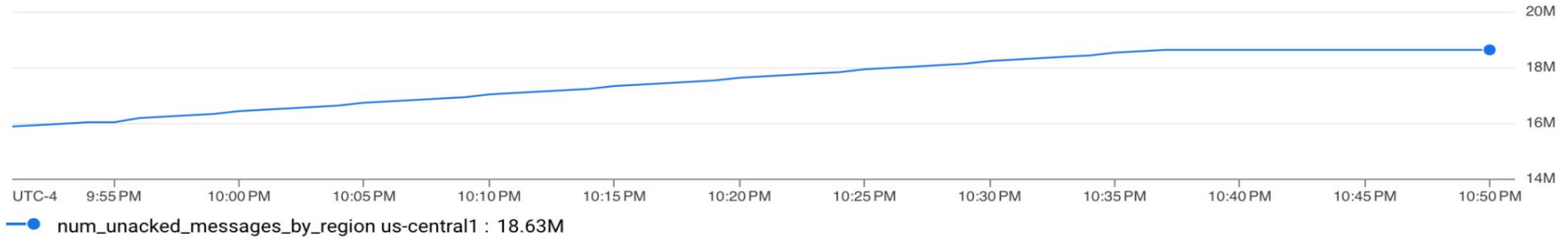




Input collections



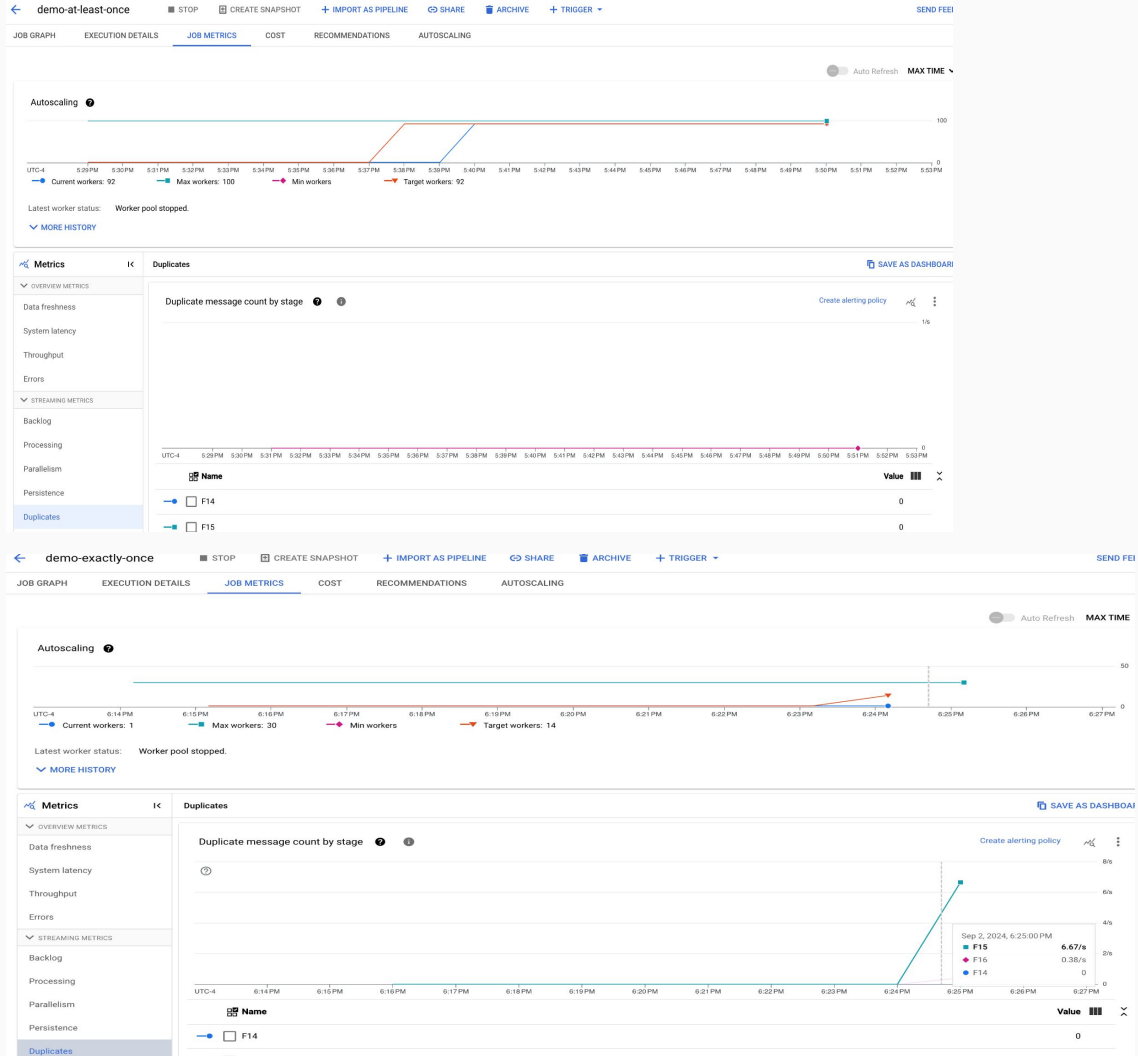
Unacked messages by region



BEAM
SUMMIT

Duplicates Metric

- Duplication will indicate 0 when using At-Least-Once



Queries On BQ

- Count exact number of duplicates

```
SELECT uuid, count(*) as duplicate FROM `ihaffa-sandbox.beam_demo.at_least_once`  
group by uuid  
having duplicate > 1);
```

Query results

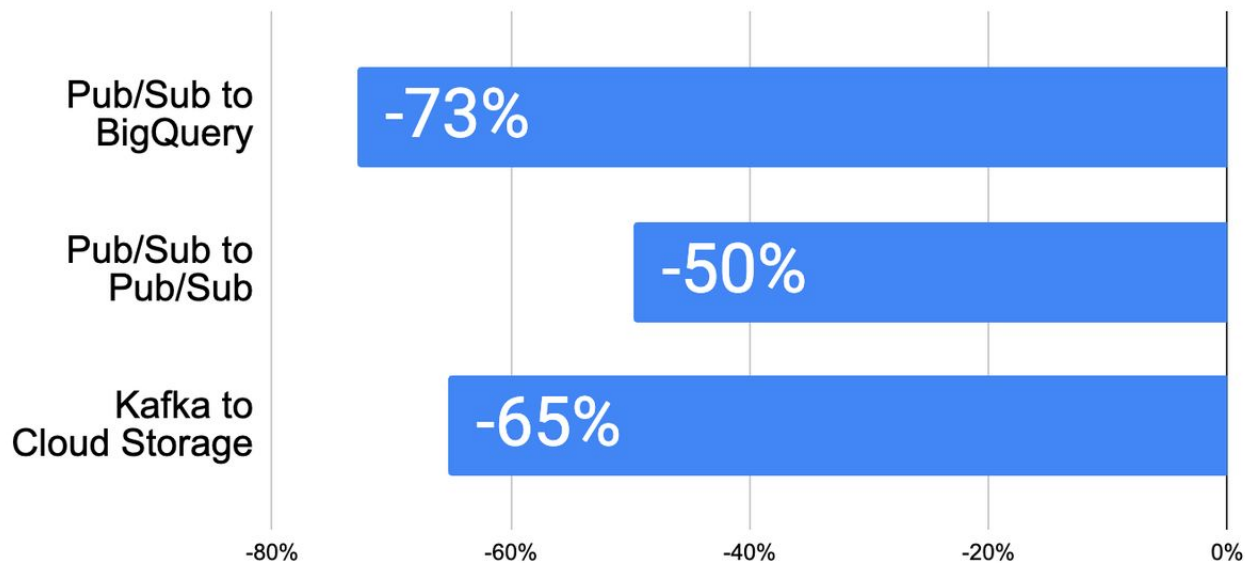
INFORMATION	RESULTS	CHART	JSON	EXECUTION DETAILS	EXPLAIN
	uuid ▾	duplicate ▾			
	57a47b0f-fb55-493e-95c7-d48f...	2			
	a5cbbe2e-c1ff-4272-9631-19b...	2			
	918d6685-0ad9-4713-bdfa-ee0...	2			
	30e2dc58-2084-4917-91e7-7ca...	2			
	76ee7883-4bb4-49a2-8206-13...	2			

```
SELECT COUNT(*) AS `Total Duplicate` from(  
SELECT uuid, count(*) as duplicate FROM `ihaffa-sandbox.beam_demo.at_least_once`  
group by uuid  
having duplicate > 1);
```

Query results

INFORMATION	RESULTS	CHART	JSON	EXECUTION DETAILS	EXPLAIN
	Total Duplicate ▾				
	838332				

Cost Savings (At-Least-Once)



When To Use?

Case	At-Least-Once (New)	Exactly-Once (Default)
Pipeline	<ul style="list-style-type: none">• Map-only Pipeline• Downstream app perform deduplication• Log processing, CDC	<ul style="list-style-type: none">• Aggregation, e.g count, sum, mean• Business-critical: fraud, financing, etc
Other Consideration	<ul style="list-style-type: none">• Cautious of external call• Cautious of high number of duplicates in pipeline	<ul style="list-style-type: none">• Pubsub read is less performant



Thank you!

Questions?

Blog: [Dataflow At Least Once vs Exactly Once](#)

Metric Exporter: [Git Repo, Run Dataflow cost & performance benchmarking \(Streaming & batch\)](#)

Email: lhaffa@google.com



BEAM
SUMMIT