- Overview

- Autotuning

- GCP PubSub Integration

- Observability

- Other Projects

# Overview
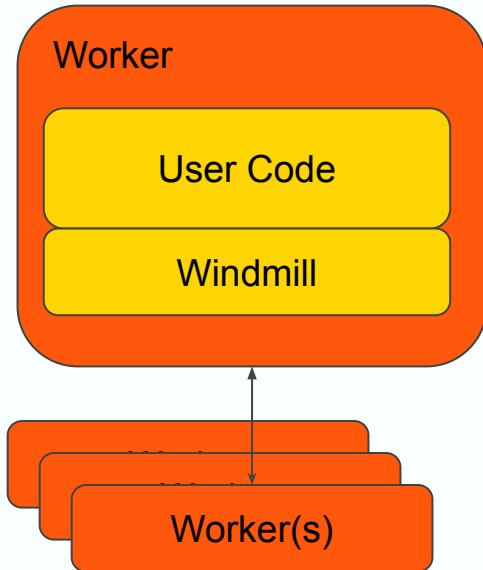
- History of Streaming @ Google

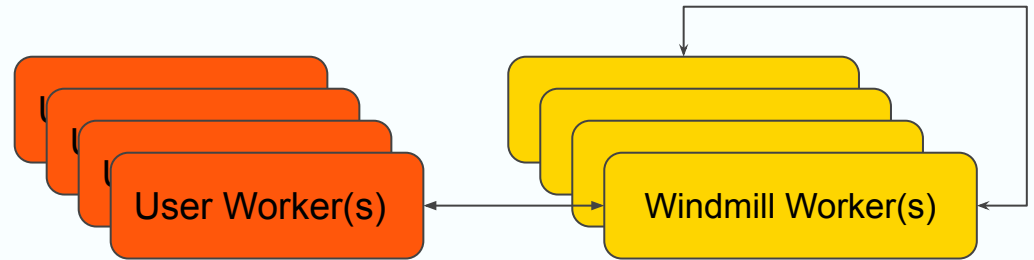- Streaming Appliance vs Streaming Engine

- Streaming Basics

- Everything was batch

- MapReduce

- First streaming systems were designed for Ads

- Streaming MapReduce

- MillWheel

- Streaming Flume

- Windmill (Dataflow)

# History of Streaming @ Google

MapReduce

2004

Simple distributed data processing

Flume

2010

Logical pipelines & optimization

Millwheel

2013

Low-latency streaming

Cloud Dataflow + Apache Beam

2015

Batch + Streaming Serverless Cloud

Streaming Engine vs Streaming Appliance

Streaming Appliance

Worker

User Code

Windmill

Worker(s)

Streaming Engine
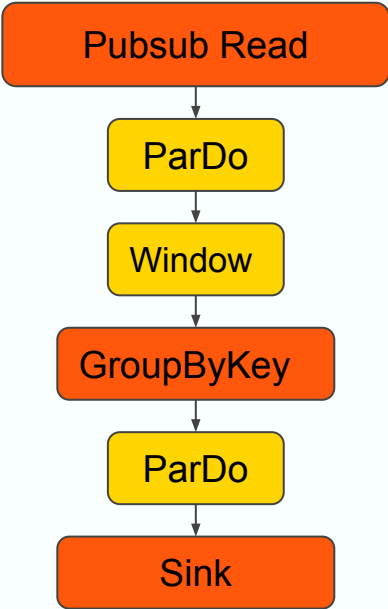
User Worker(s)

Windmill Worker(s)
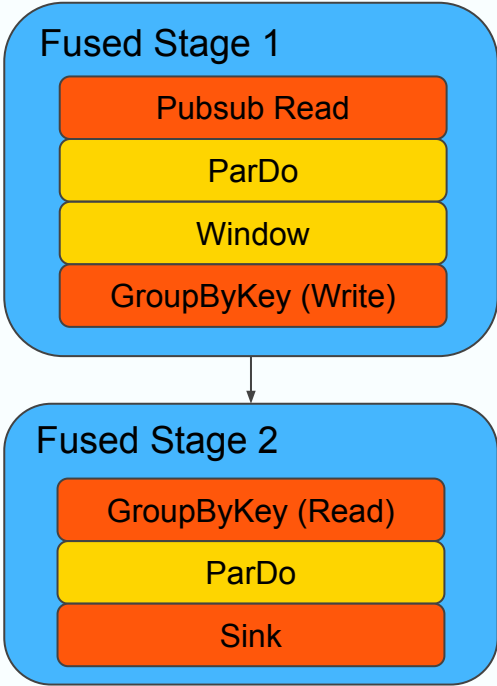
# Streaming Engine vs Streaming Appliance

Benefits of Streaming Engine:

- More efficient use of User Workers

- No need for Persistent Disks

- More responsive Horizontal Autoscaling
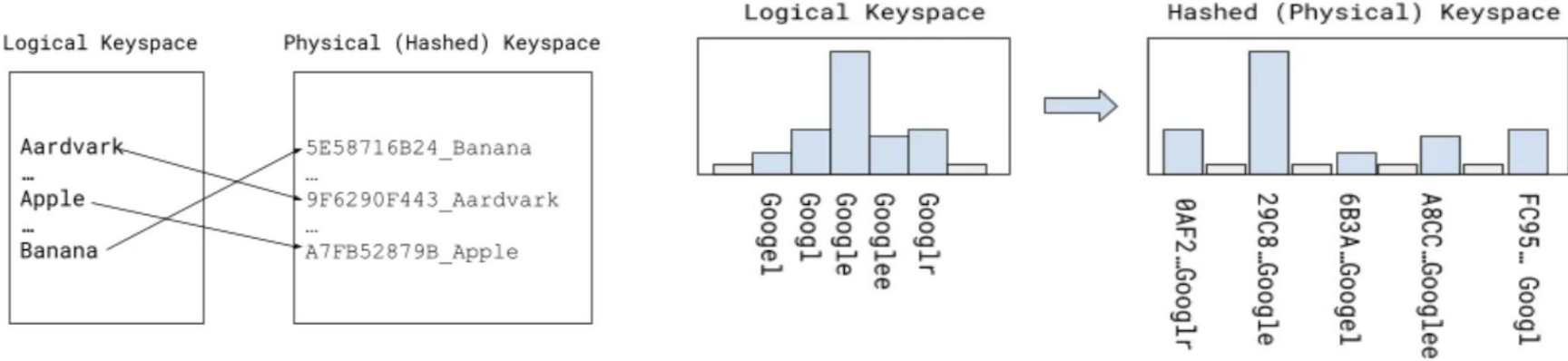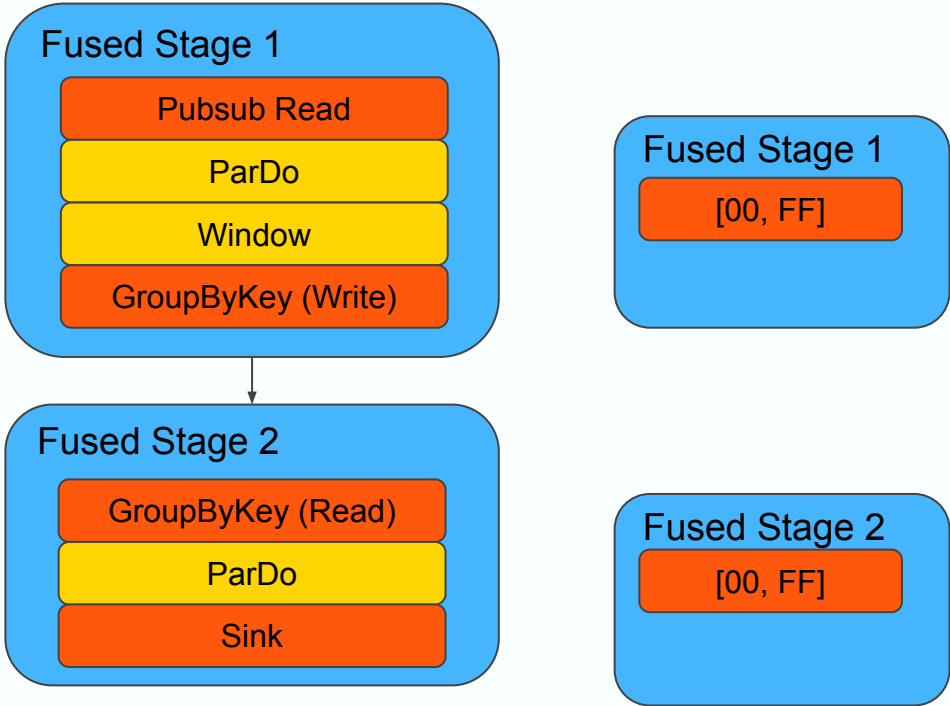
- Improved supportability and visibility

- Every message has a key assigned to it
- Keys can be user defined or system defined
- Keys are hashed
- Elements are processed in the context of a key
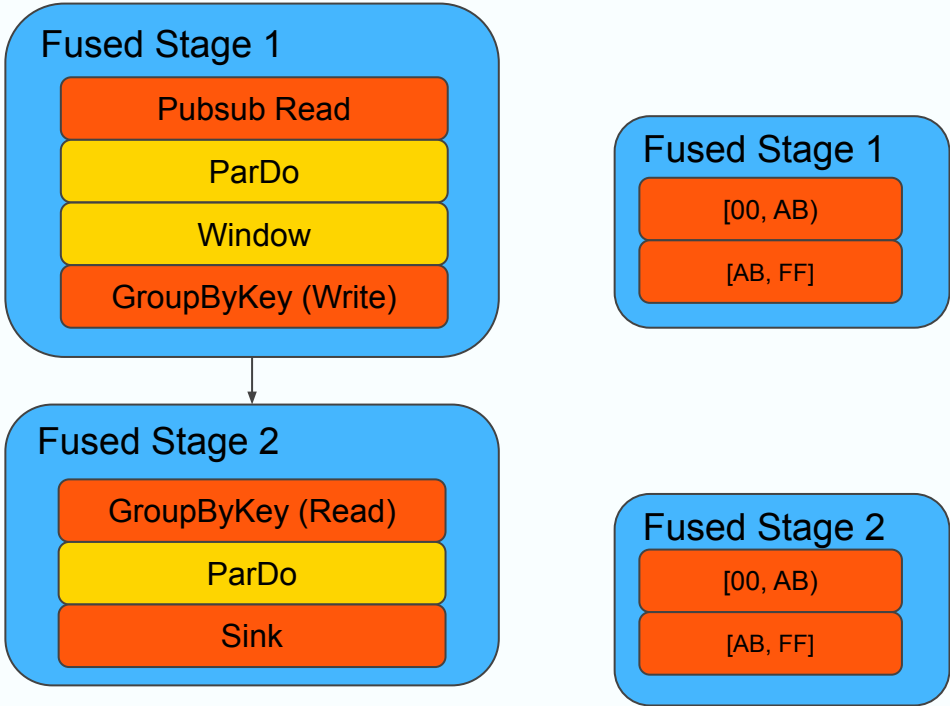- Keys are the basic unit of parallelism

Streaming Basics

- Keys belong to key-ranges
- Key ranges are assigned to workers
- Key ranges can be split and sent to different workers
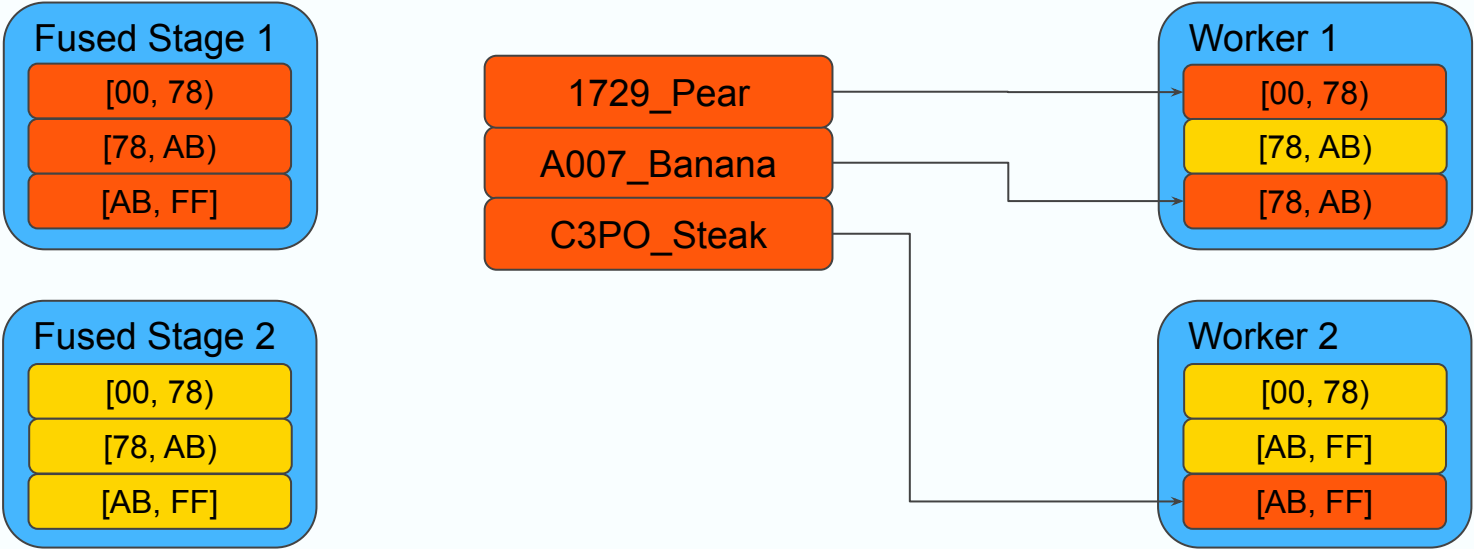
Fused Stage 1
- Pubsub Read
- ParDo
- Window
- GroupByKey (Write)

Fused Stage 2
- GroupByKey (Read)
- ParDo
- Sink

Fused Stage 1
- [00, FF]

Fused Stage 2
- [00, FF]

*NOTE: all range boundaries are hexadecimal values.*

- Keys belong to key-ranges
- Key ranges are assigned to workers
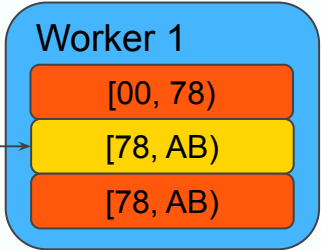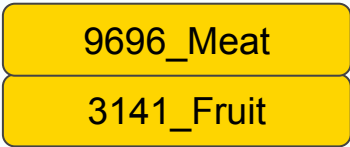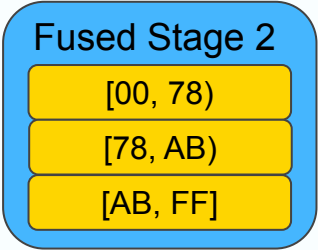- Key ranges can be split and sent to different workers

**Fused Stage 1**
- Pubsub Read
- ParDo
- Window
- GroupByKey (Write)

**Fused Stage 2**
- GroupByKey (Read)
- ParDo
- Sink

**Fused Stage 1**
- [00, AB)
- [AB, FF]

**Fused Stage 2**
- [00, AB)
- [AB, FF]

# Autotuning

# Autotuning: Asymmetric Autoscaling



**Past**: Scaling backend workers linearly with user workers.

**Present**: Scaling each worker pool independently.



Baseline



Asymmetric

# Autotuning: Key-Based Throttling

**Past**: Unconditionally throttling user worker upscale if < 20% CPU utilization.
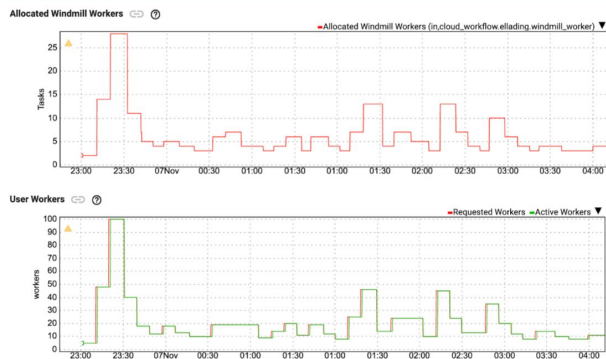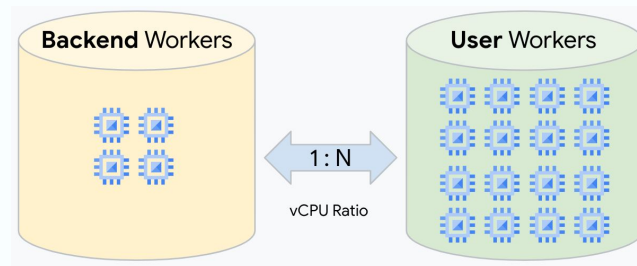


**Present**: Throttle user worker upscale on key parallelism limits (number of keys).

# Autotuning: Downscale Dampening

**Past**: Only consider the current state (backlog, throughput, etc.)

**Present**: Track scaling frequencies, downscale slower when yo-yoing detected (frequent up/down scaling in short time frame).



Before        After

# Autotuning: Scaling Actuation Latencies

**Past**: When autoscale events happen, new workers need to load the pipeline state from persistence. This can take time and lead to backlog and latency.

**Present**: Transfer info directly from workers, reducing latency

Worker

[0, FF]

# Autotuning: Scaling Actuation Latencies

**Past**: When autoscale events happen, new workers need to load the pipeline state from persistence. This can take time and lead to backlog and latency.

**Present**: Transfer info directly from workers, reducing latency

Worker
[0, AB)

New Worker
[AB, FF]

# Autotuning: Scaling Actuation Latencies

**Past**: When autoscale events happen, new workers need to load the pipeline state from persistence. This can take time and lead to backlog and latency.

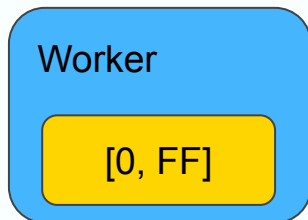**Present**: Transfer info directly from workers, reducing latency

# Autotuning: Scaling Actuation Latencies

**Past**: When autoscale events happen, new workers need to load the pipeline state from persistence. This can take time and lead to backlog and latency.
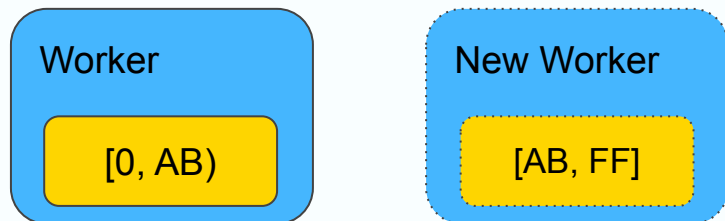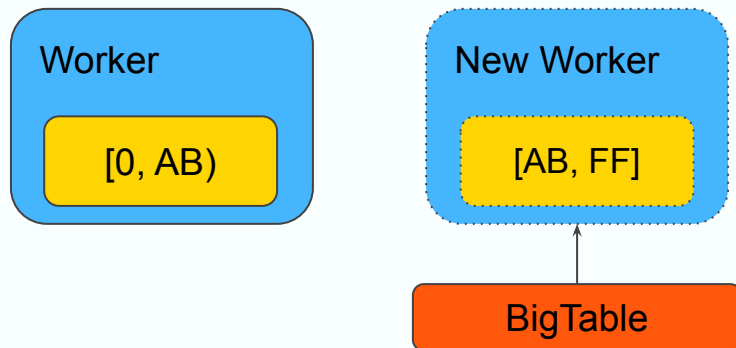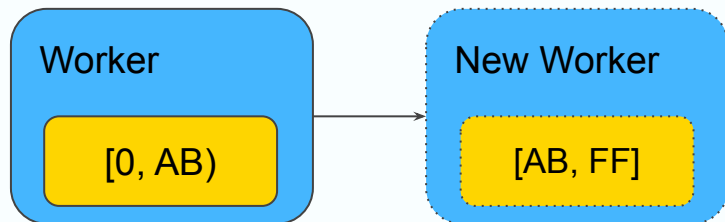
**Present**: Transfer info directly from workers, reducing latency

# Autotuning: Scaling Actuation Latencies

Latency Disabled (top) vs Enabled (bottom)

## User Workers Disabled vs Enabled

# Autotuning: Range Rebalancing

**Past**: If a key range has a disproportionate amount of input rate, its worker would have more load than others, potentially accumulating backlog and wasting resources on other workers.

**Present**: We can split key ranges dynamically and rebalance them across workers based on their throughput

| 0001-key1 | 500mb/s |
| 6002-key2 | 300mb/s |
| BCDF-key3 | 50mb/s |

Worker1

| [00, AB) | 800mb/s |

Worker2

| [AB, FF] | 50mb/s |

**Past**: If a key range has a disproportionate amount of input rate, its worker would have more load than others, potentially accumulating backlog and wasting resources on other workers.

**Present**: We can split key ranges dynamically and rebalance them across workers based on their throughput

| | |
|---|---|
| 0001-key1 | 500mb/s |
| 6002-key2 | 300mb/s |
| BCDF-key3 | 50mb/s |

**Worker1**

| | |
|---|---|
| [00, 55) | 500mb/s |

**Worker2**

| | |
|---|---|
| [55, AB) | 300mb/s |
| [AB, FF] | 50mb/s |

Autotuning: Range Rebalancing

# Autotuning: BigQuery Autosharding

**Past**: Autosharding was only available for Streaming Inserts / File Loads and was load agnostic, which could lead to wasted resources in case of dynamic destinations

**Present**: StorageAPI gets autosharding option, using backlog and throughput as metric.

| Table 1 | 200mb/s | 1000 shards |
|---------|---------|-------------|
| Table 2 | 100mb/s | 1000 shards |
| Table 3 | 1mb/s   | 1000 shards |

# Autotuning: BigQuery Autosharding

**Past**: Autosharding was only available for Streaming Inserts / File Loads and was load agnostic, which could lead to wasted resources in case of dynamic destinations

**Present**: StorageAPI gets autosharding option, using backlog and throughput as metric.

| | | |
|---|---|---|
| Table 1 | 200mb/s | 800 shards |
| Table 2 | 100mb/s | 400 shards |
| Table 3 | 1mb/s | 4 shards |

# Autotuning: BigQuery Autosharding



Streaming Inserts + Autosharding

Storage API

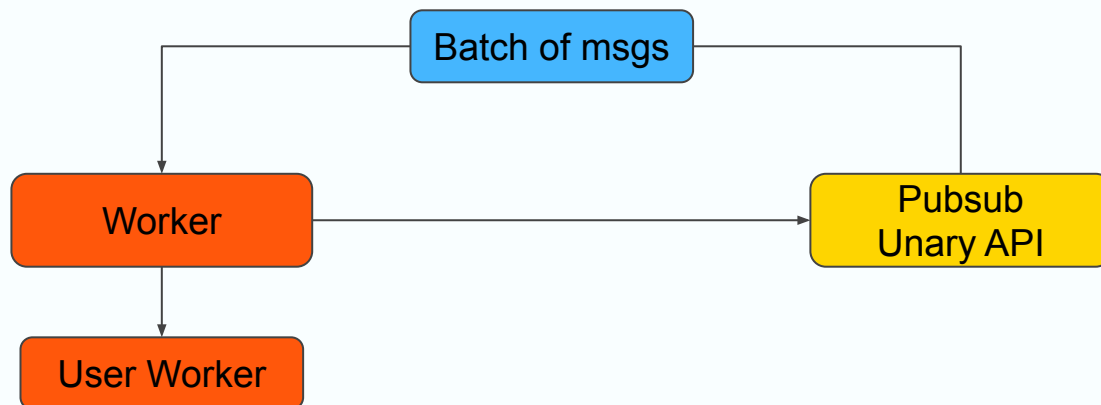Storage API + Autosharding

# GCP PubSub Integration

# PubSub Streaming Pull

**Past**: Pipelines used old Pubsub API Unary Pull

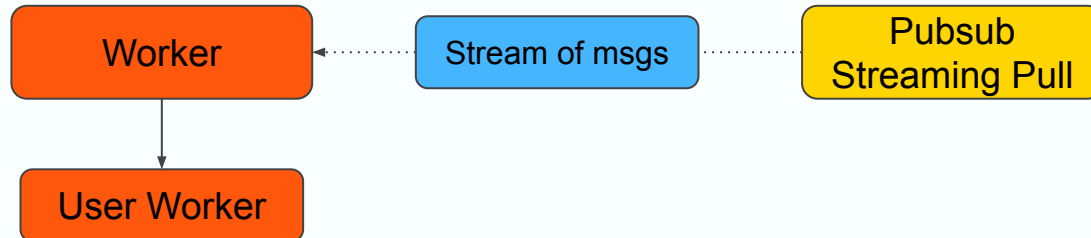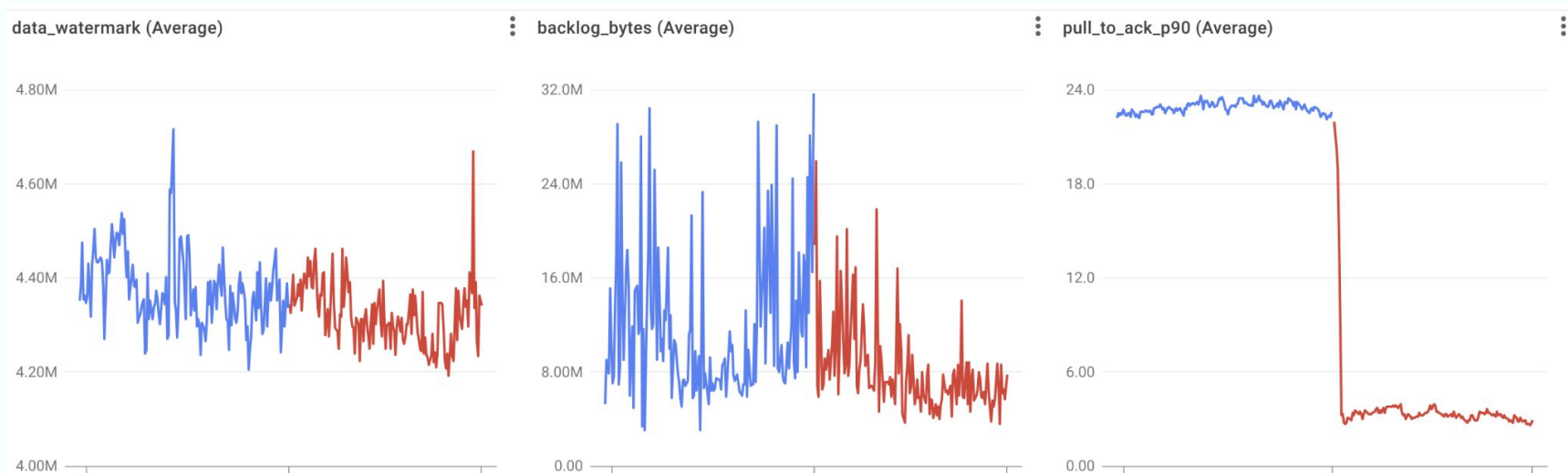**Present**: Pipelines use newer Pubsub API Streaming Pull, improving throughput and latency

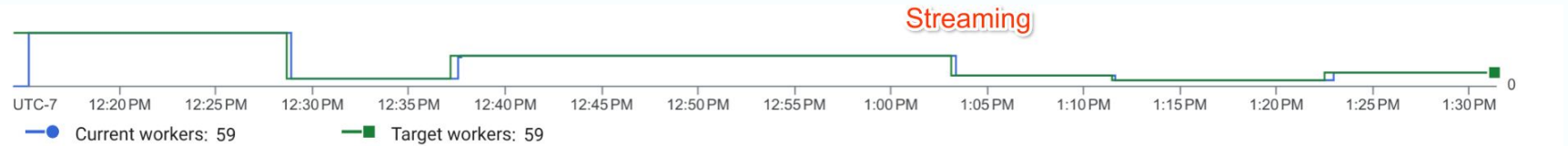**Past**: Pipelines used old Pubsub API Unary Pull

**Present**: Pipelines use newer Pubsub API Streaming Pull, improving throughput and latency
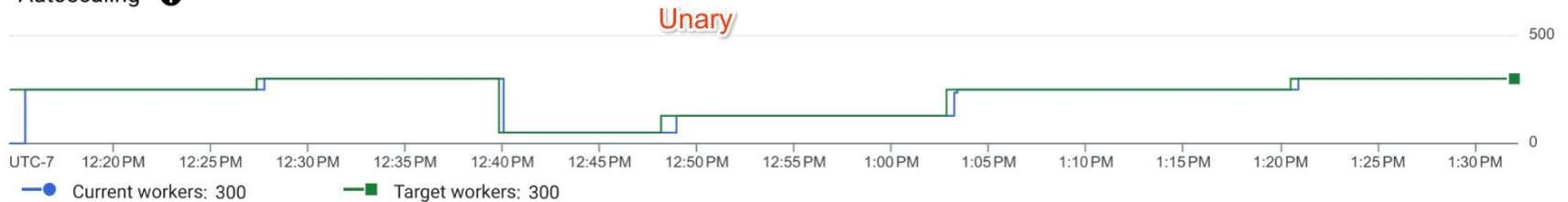
## Latency and Backlog Improvements

## Usage improvements



Streaming

UTC-7   12:20 PM   12:25 PM   12:30 PM   12:35 PM   12:40 PM   12:45 PM   12:50 PM   12:55 PM   1:00 PM   1:05 PM   1:10 PM   1:15 PM   1:20 PM   1:25 PM   1:30 PM

— Current workers: 59    — Target workers: 59

Autoscaling ❓

Unary

UTC-7   12:20 PM   12:25 PM   12:30 PM   12:35 PM   12:40 PM   12:45 PM   12:50 PM   12:55 PM   1:00 PM   1:05 PM   1:10 PM   1:15 PM   1:20 PM   1:25 PM   1:30 PM

— Current workers: 300    — Target workers: 300

Latest worker status:    Autoscaling: Raised the number of workers to 300 so that the pipeline can catch up with its backlog and keep up with its input rate.

∨ MORE HISTORY

# Observability

# Observability: New Metrics

- Collecting many new Streaming Engine metrics

- Some integrated into Dataflow UI

- All available in Monitoring UI

- Available dashboard template for easy detailed job performance monitoring

### New Metrics

| Metrics | Path |
|---|---|
| Duplicates Filtered | /job/duplicates_filtered_out_count |
| Processing Parallelism | /job/processing_parallelism_keys |
| Backlog Bytes | /job/backlog_bytes |
| Backlog Seconds | /job/estimated_backlog_processing_time |
| Timers Processed | /job/timers_processed_count |
| Timers Resident | /job/timers_pending_count |
| Status of Streaming Pull connections | /job/pubsub/streaming_pull_connection_status |
| The number of bytes produced by this ptransform | /job/estimated_bytes_produced_count |
| Checkpoint bytes written | /job/streaming_engine/persistent_state/write_bytes_count |
| Checkpoint bytes read | /job/streaming_engine/persistent_state/read_bytes_count |
| Checkpoint Latency | /job/streaming_engine/persistent_state/write_latencies |
| User Processing Latency | /job/bundle_user_processing_latencies |
| Key (Range) Availability | /job/streaming_engine/key_processing_availability |
| The number of bytes consumed by this ptransform | /job/estimated_bytes_consumed_count |
| The number of bytes being processed by ptransform | /job/estimated_bytes_active |
| Pubsub Pull to Ack Latency | /job/pubsub/pulled_message_ages |
| Persistent State Usage | /job/streaming_engine/persistent_state/stored_bytes |
| Late pubsub messages | /job/pubsub/late_messages_count |
| Target workers | /job/target_worker_instances |
| Pubsub Publish Messags/Errors | /job/pubsub/published_messages_count |

Launched

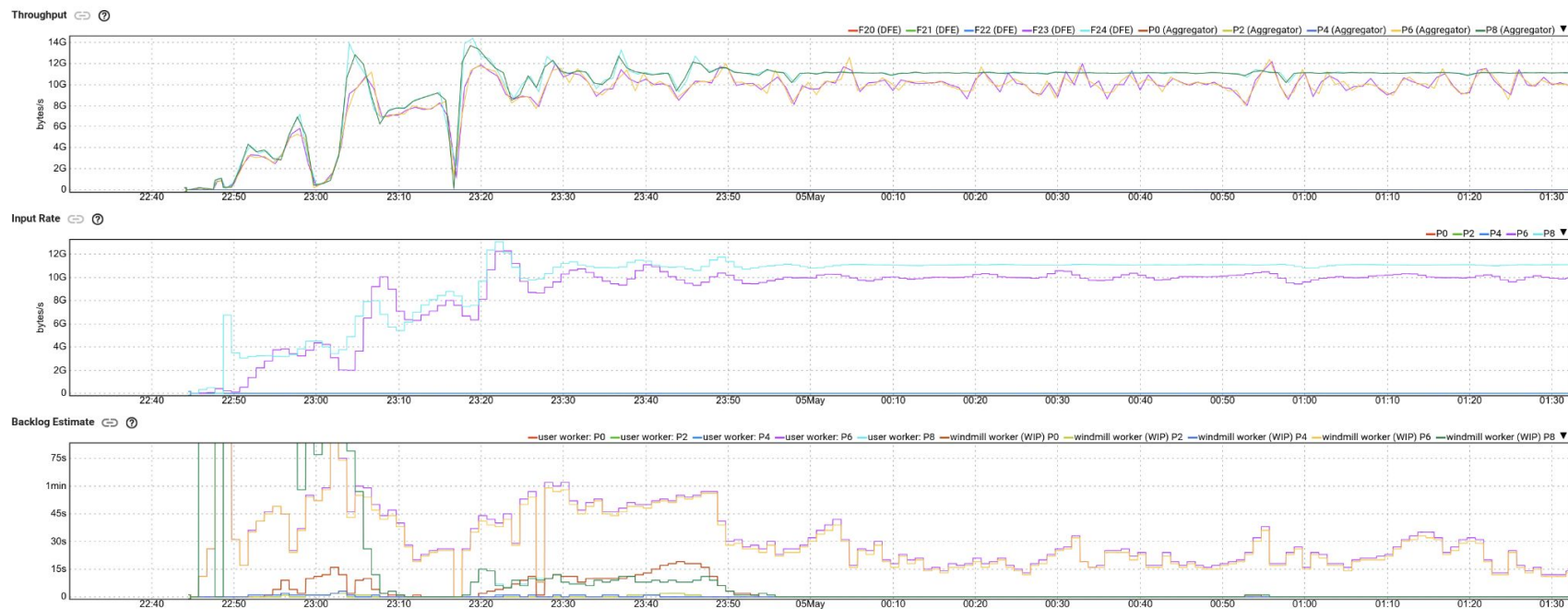# Observability: Dashboard Template

Importing template

Preview of first few graphs

# Other Projects

# Out of the box

We wanted to test the throughput of sources and sinks without any special settings. We got to 10 GB/s for these I/O combos:

- Pubsub to BQ
- Pubsub to Pubsub
- Pubsub to GCS*
- Kafka to GCS*
- Kafka to BQ

## Pubsub to GCS example

Collection of +190 self-contained Dataflow pipelines ready to use, including most common sources, sinks, and use cases.

https://github.com/GoogleCloudPlatform/dataflow-cookbook

Iñigo San Jose, Tom Stepp

# QUESTIONS?

BEAM
SUMMIT
NYC 2023