

# Apache Gearpump

## next-gen streaming engine

Karol Brejna, Intel ([karolbrejna@apache.org](mailto:karolbrejna@apache.org))

Huafeng Wang, Intel ([huafengw@apache.org](mailto:huafengw@apache.org))

# Agenda

- What is Gearpump?
- Why Apache Gearpump?
- Apache Gearpump features/internals
- What's next for Apache Gearpump



# What is Gearpump ?

- A super simple pump that consists of only two gears but very powerful at streaming water
- An Akka<sup>[2]</sup> based real-time streaming engine
- An Apache Incubator<sup>[1]</sup> project since Mar.8th, 2016



# Why Gearpump ?

# Stream processing is hard

- Fault tolerance
- Infinite Out-of-order data
- Low latency assurance (e.g. real-time recommendation)
- Correctness requirement (e.g. charge advertisers for ads)
- Cheap to update applications (e.g. tune machine learning parameters)

# Gearpump makes stream processing easier

- fault tolerant stream processing at latency of milliseconds
- handling out-of-order data
- event-time based window aggregation
- Akka-stream DSL and Apache Beam API support
- runtime DAG modification
- responsive UI with abundant metrics information

# Gearpump on TAP

- Gearpump on Trusted Analytics Platform (TAP)
- Stream processing - performance experiments and results

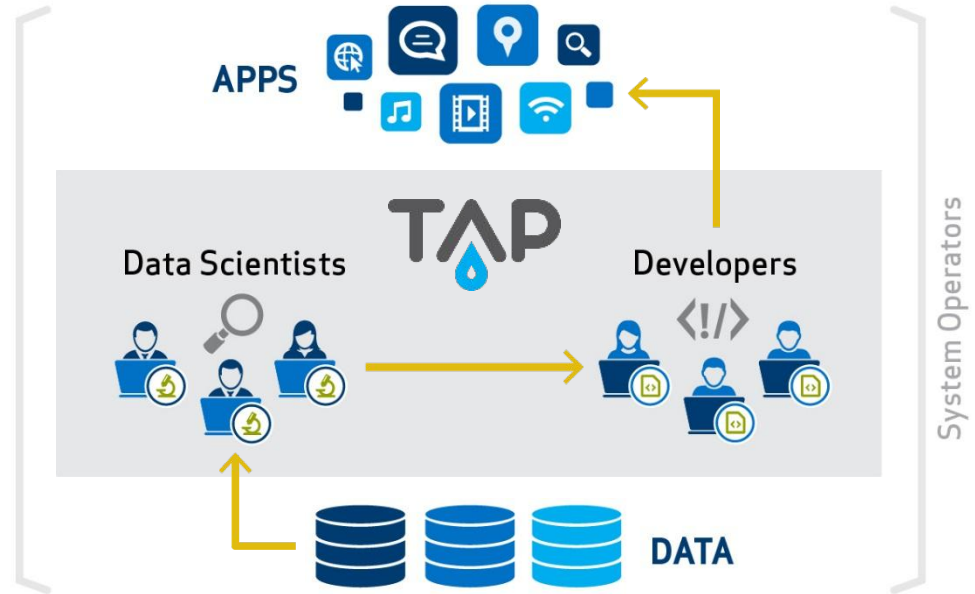
# Gearpump on TAP

- Gearpump on Trusted Analytics Platform (TAP)
- Stream processing - performance experiments and results



# Trusted Analytics Platform (TAP)

- Open Source project
- Collaborative, cloud-ready platform to build applications powered by Big Data Analytics
- Includes everything needed by data scientists, application developers and system operators
- Optimized for performance and security



# Analytics Solutions – Big Data Scale Out

## Applications

Analytics-powered vertical and horizontal solutions



## Analytics

Open source platform for collaborative data science and analytics application development



## Data

Open source, Hadoop-centric platform for distributed and scalable storage and processing



## Infrastructure

Software-defined storage, network and cloud infrastructure optimized for Intel Architecture



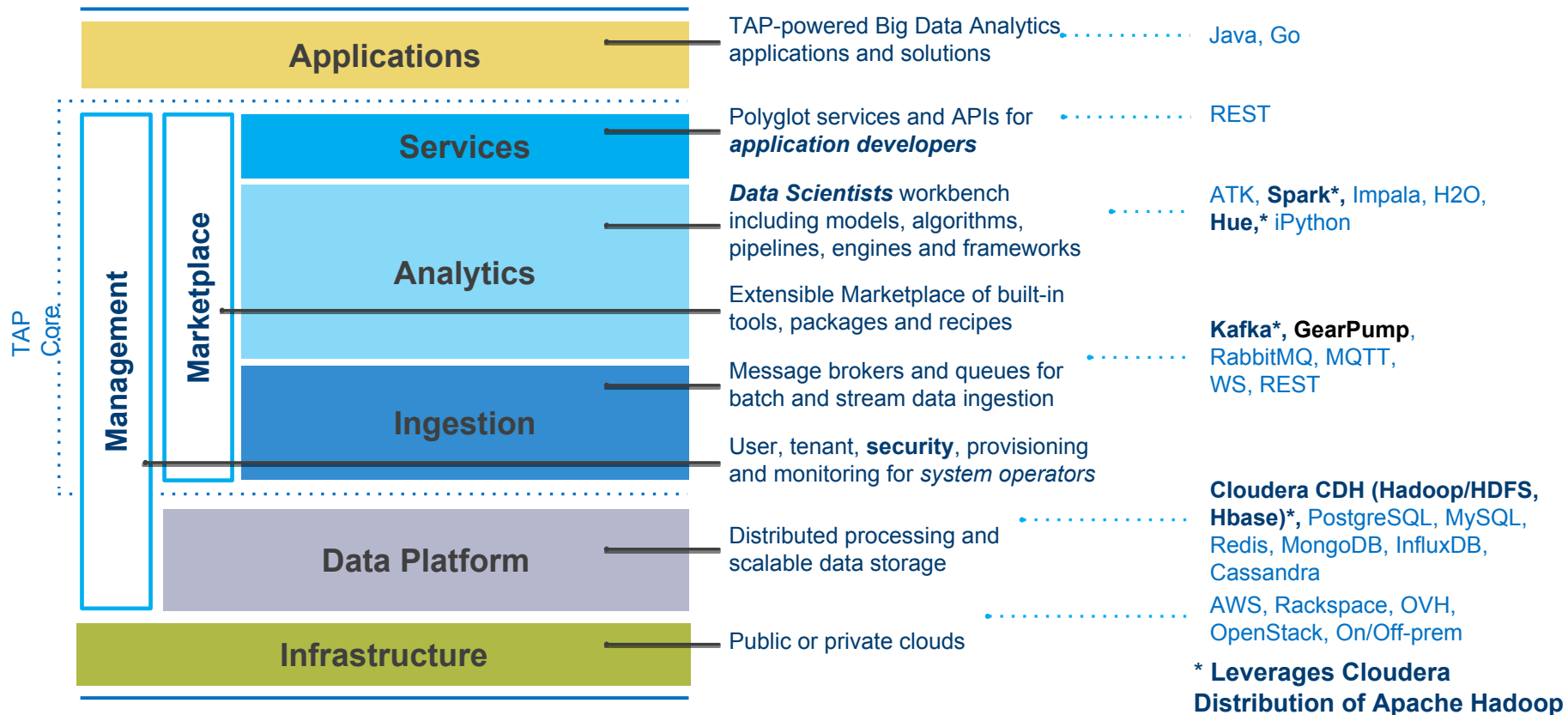
### Machine Learning

Multi-layered, fully-optimized algorithms

### Performance and Security

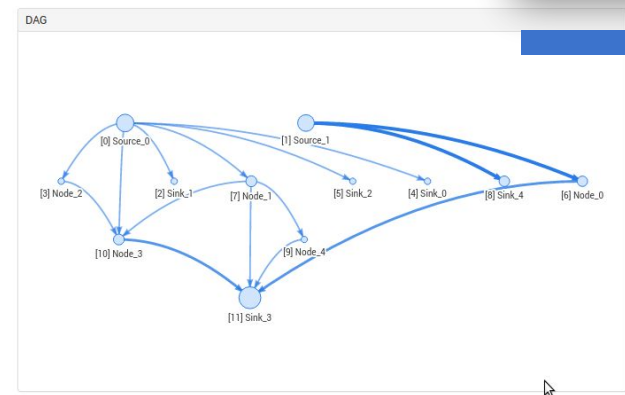
Silicon and software enhancements to protect and accelerate data and analytics

# The Anatomy of Trusted Analytics Platform (TAP)



dag

Overview Metrics **DAG**



**Message Receive Throughput**  
118,445 msg/s  
Total messages received: 9,741,570

**Message Send Throughput**  
99,411 msg/s  
Total messages sent: 8,311,780

**Average Message Processing**  
0.015 ms

**Average Message Receive Lat**  
9.844 ms





## Wearables

Customer behavior analysis using wearable devices



## Healthcare

Predict individual health problems to improve care



## Retail

Asset management with RFID data



## Industrial

Predict equipment failures and optimization based on sensor data



## Genomics

Execute privacy-preserving analytics on diverse distributed data sets



## IoT Developer Platform

Enable development of data-generating IoT apps



## Security

Detect threats – IT, grid, machines, sites



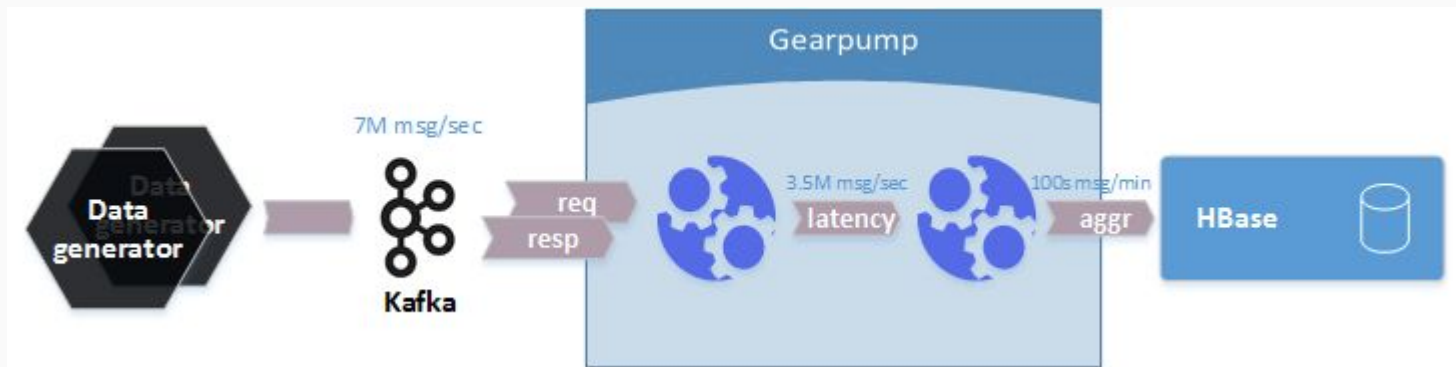
## Forecasting

Financial, inventory, supply chain, etc.

# Gearpump on TAP

- Gearpump on Trusted Analytics Platform (TAP)
- Stream processing - performance experiments and results

# The problem



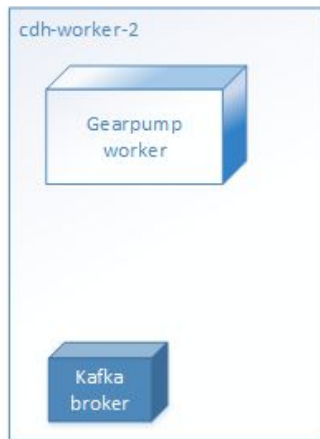
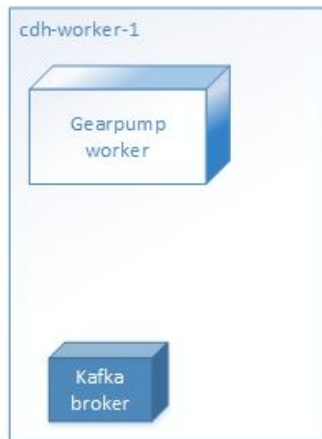
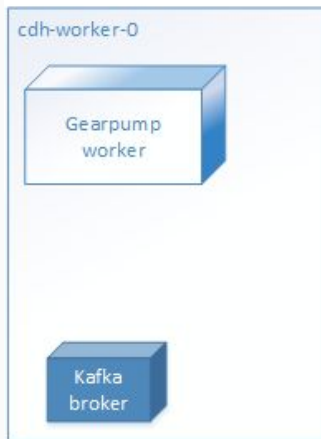
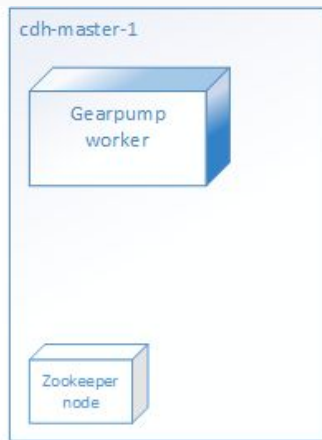
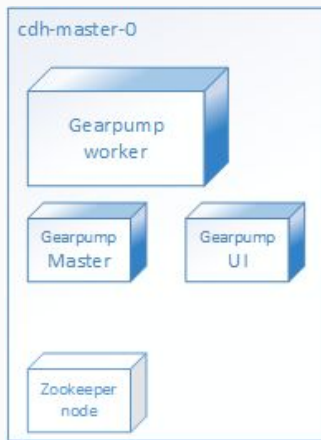
- correlate messages using a key in one second sliding window and produce latency stream messages
- consume latency messages and compute average latency per firm in one minute buckets
- send the aggregate message to HBase

## The expectations

- Handle load of 0.5M msg per second all the time
- Handle load of 7M msg per second for peaks of 1 hour
- Message size 250-500 bytes
- Be able to scale for even more

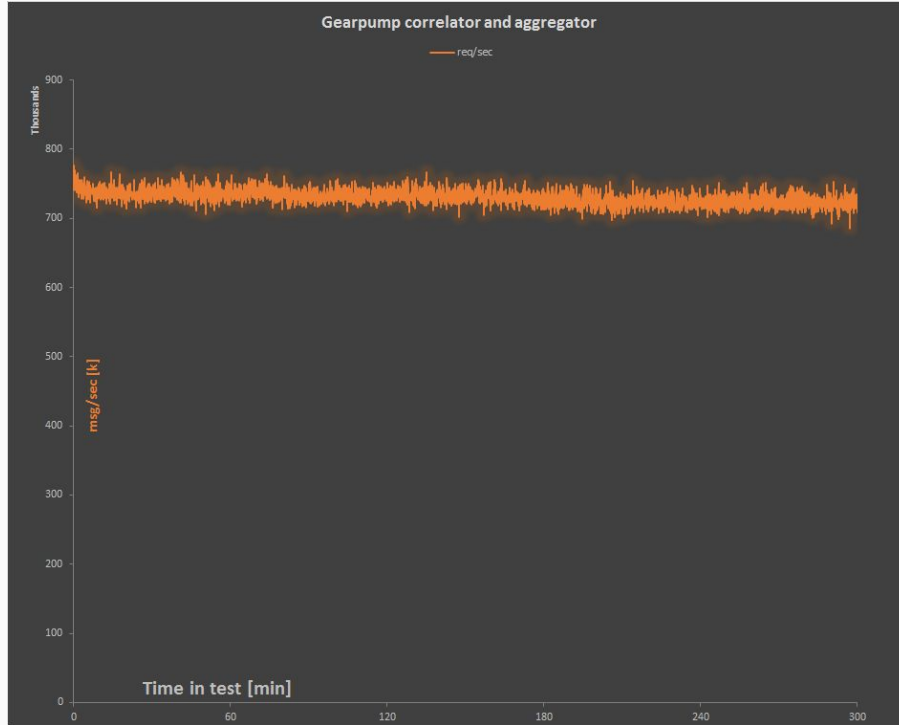


# The hardware



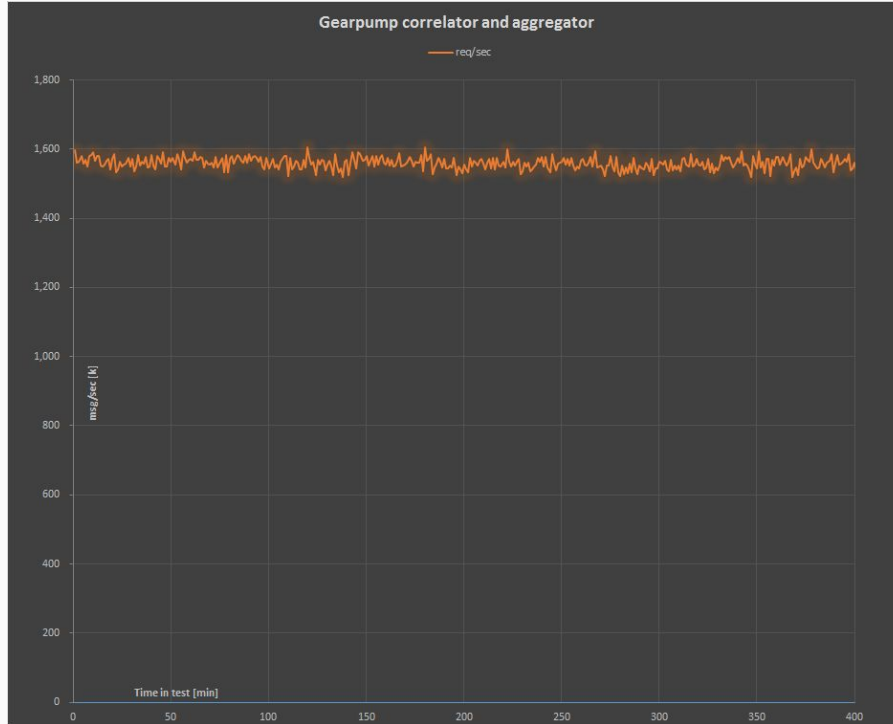
- CPU: Intel(R) Xeon(R) CPU E5-2695 v3 @ 2.30GHz
- Memory: 256 Gbytes DDR4
- Storage: 8 SATA SSDs

# The results (1) - let's start small: ~700k msg/sec



Initial attempt

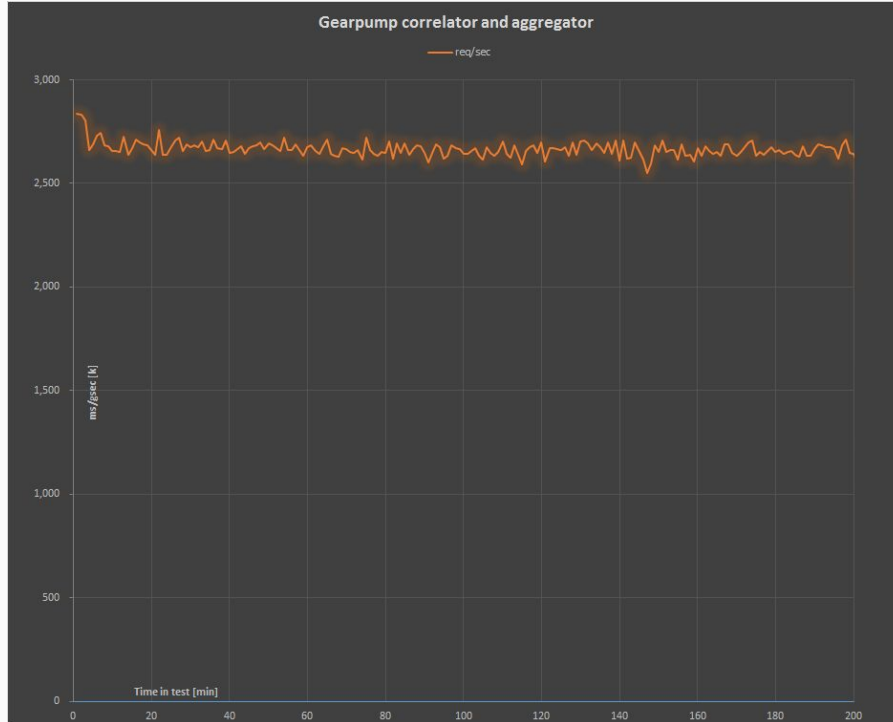
## The results (2) - 8 executors: ~1.6M msg/sec



### Findings:

- We need to improve Kafka Source - queue size, fetch frequency
- Improve Kafka partitions design for concurrency
- Network throughput may be a bottleneck ( $1.6\text{M msg/sec} * 0.5\text{ k} * 8\text{ bit}$ ) - compression

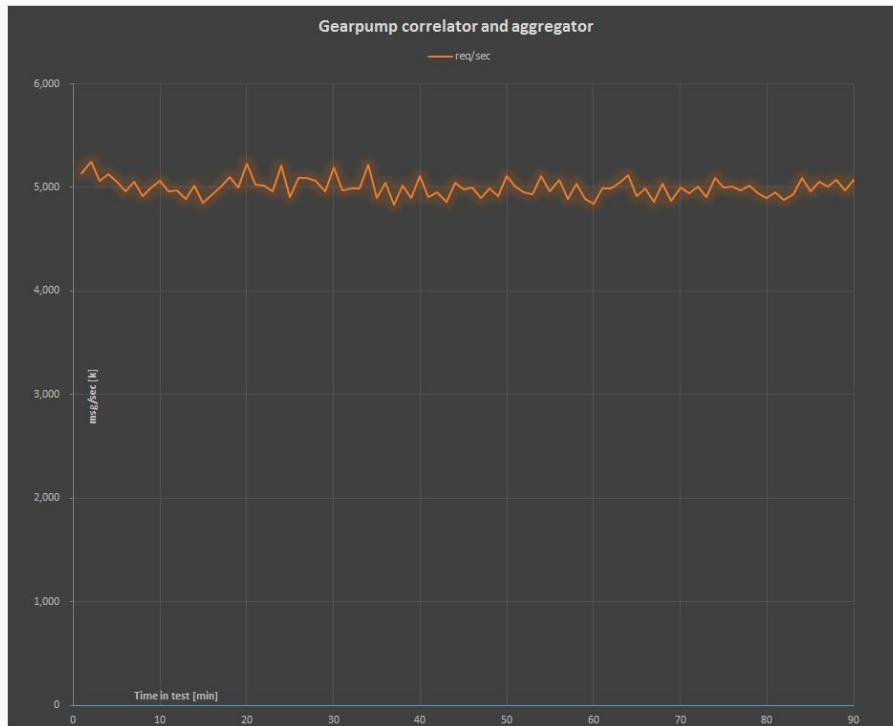
## The results (3) - 16 executors: ~2.7M msg/sec



### Findings:

- JVM defaults designed for moderate workloads - we need to pump them up
- Message marshalling starts to play significant role in performance - look for better alternatives

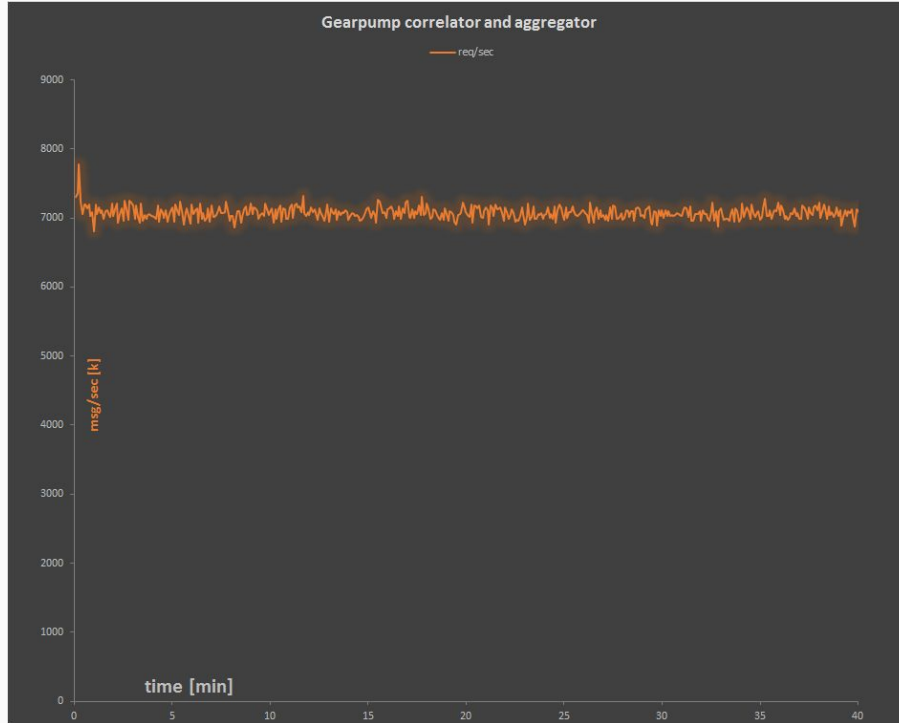
# The results (4) - 32 executors: ~5M msg/sec



## Findings:

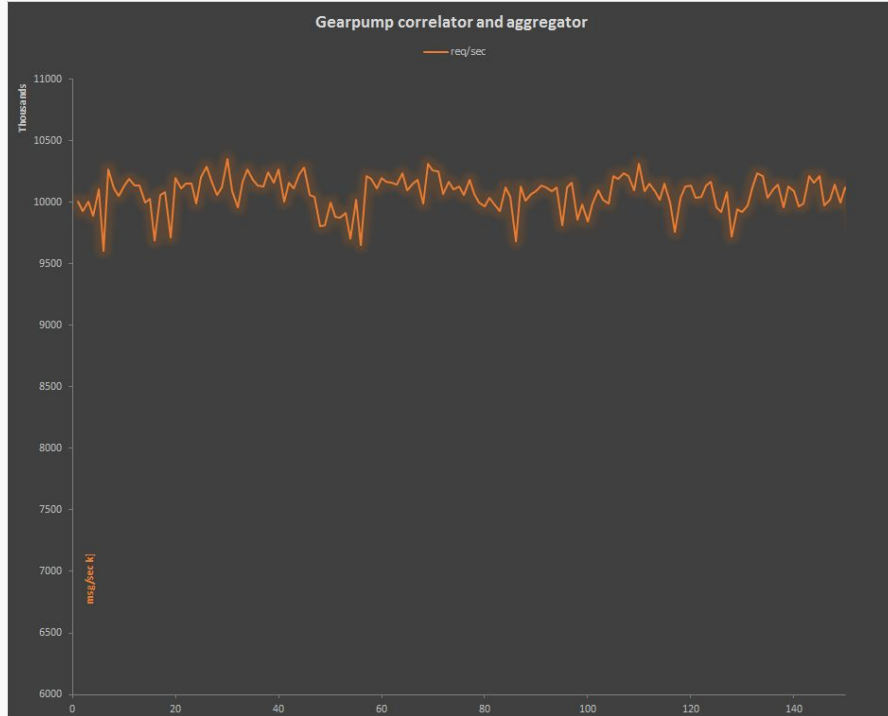
- Backpressure introduced by JVM  
⇔ JVM communication - use task fusing

# The results (5) - 48 executors: 7.4M msg/sec



Mission accomplished!!!

# The results (6) - 64 executors



We can go even further..

## The results - summary

- Great performance numbers on decent hardware
- **Predictable** scalability

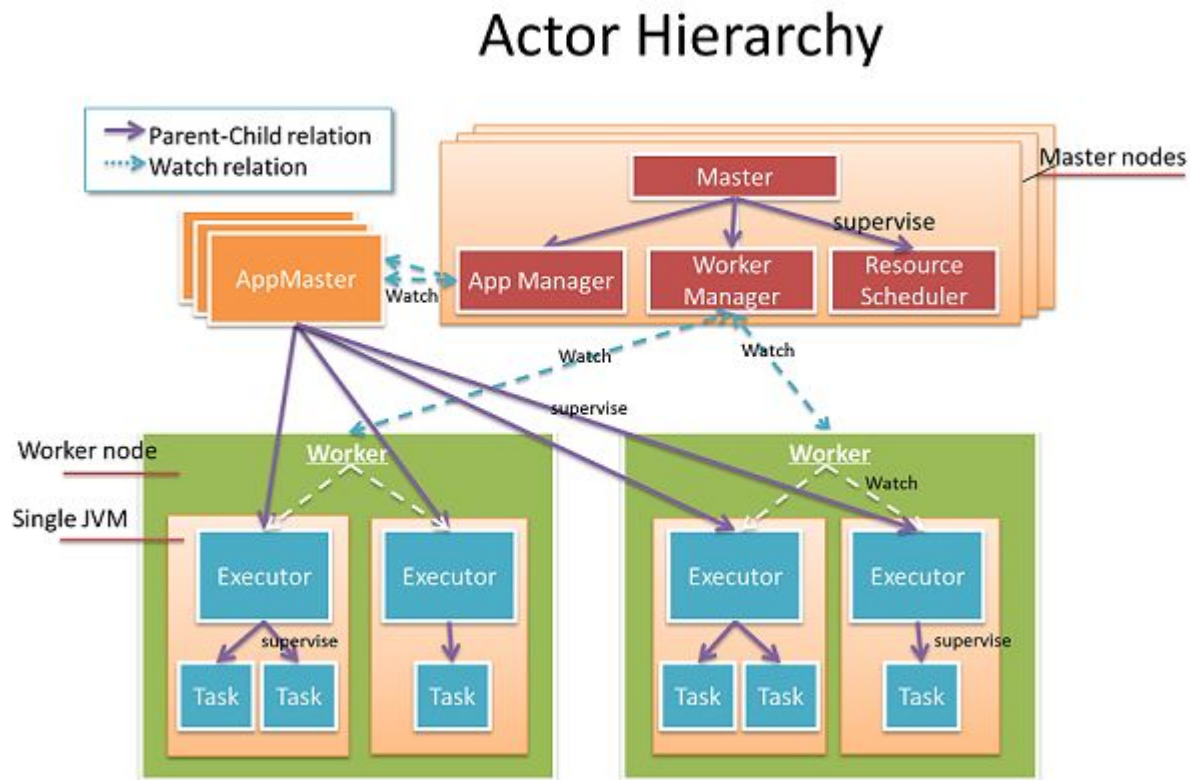
Executors number	Req/sec
8	1.6 M
16	2.7 M
32	5 M
48	7,4 M
64	10 M



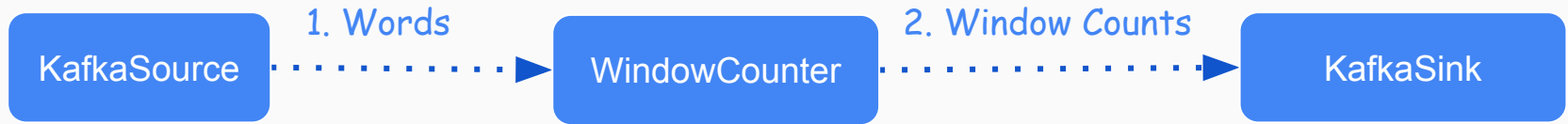
# Gearpump features

# Gearpump Architecture

- Actor concurrency
- Message passing communication
- error handling and isolation with supervision hierarchy
- Master HA with Akka Cluster



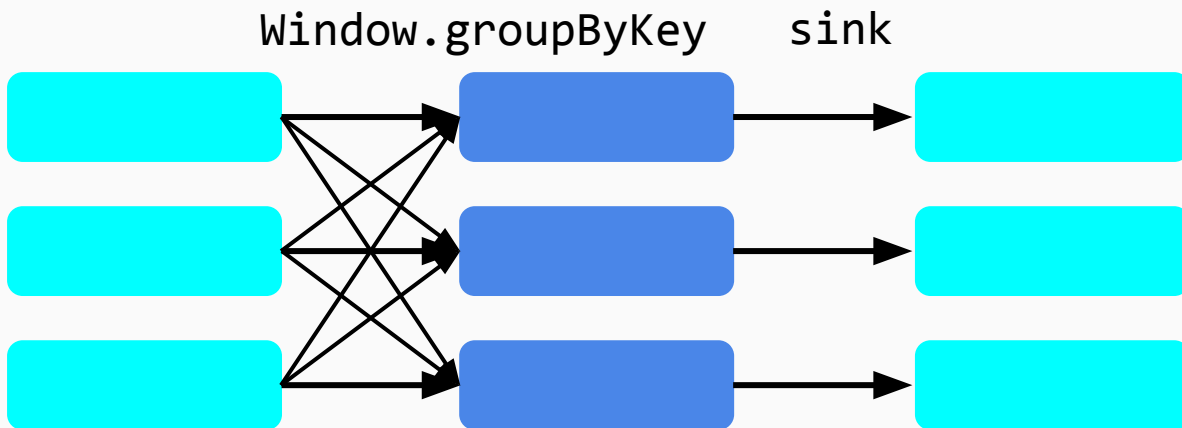
## Use case - Windowed word count



# How Gearpump solves the hard parts

- User interface
- Flow control
- Out-of-order processing
- Exactly once
- Dynamic DAG

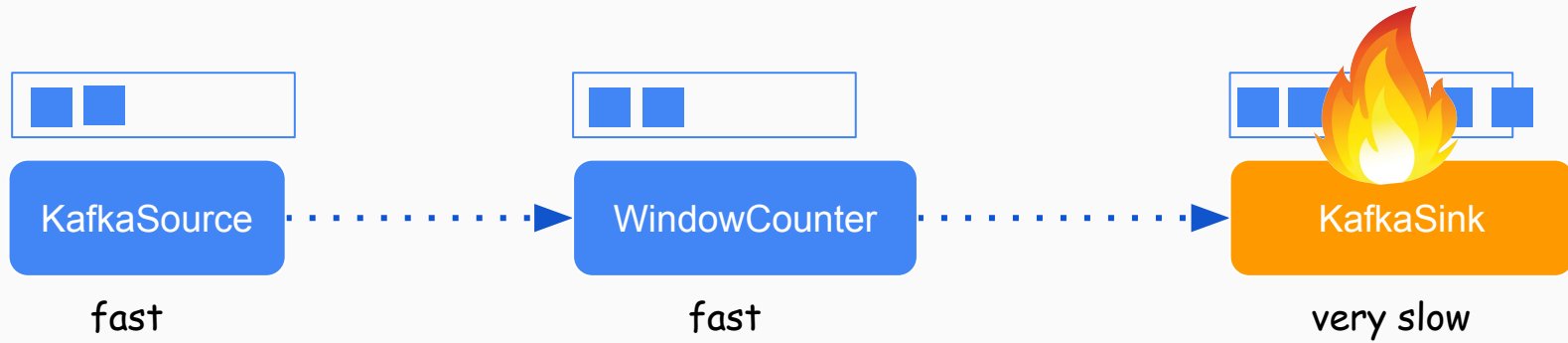
```
val app = StreamApp("dsl", context)
app.source[String](kafkaSource).
  flatMap(line => line.split("[\\s]+")).map((_, 1)).
  window(FixedWindow.apply(Duration.ofMillis(5L))
    .triggering(EventTimeTrigger)).
  // (word, count1), (word, count2) => (word, count1 + count2)
  groupBy(_._1).sum.sink(kafkaSink)
```



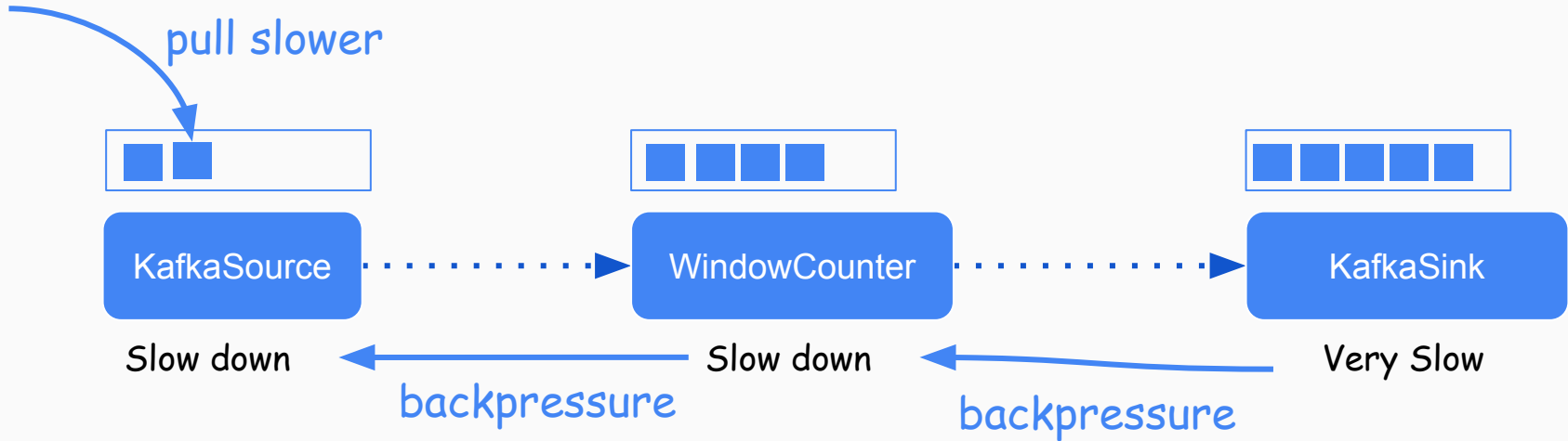
# How Gearpump solves the hard parts

- User interface
- Flow control
- Out-of-order processing
- Exactly once
- Dynamic DAG

# Without Flow Control - OOM



# With Flow Control - Backpressure

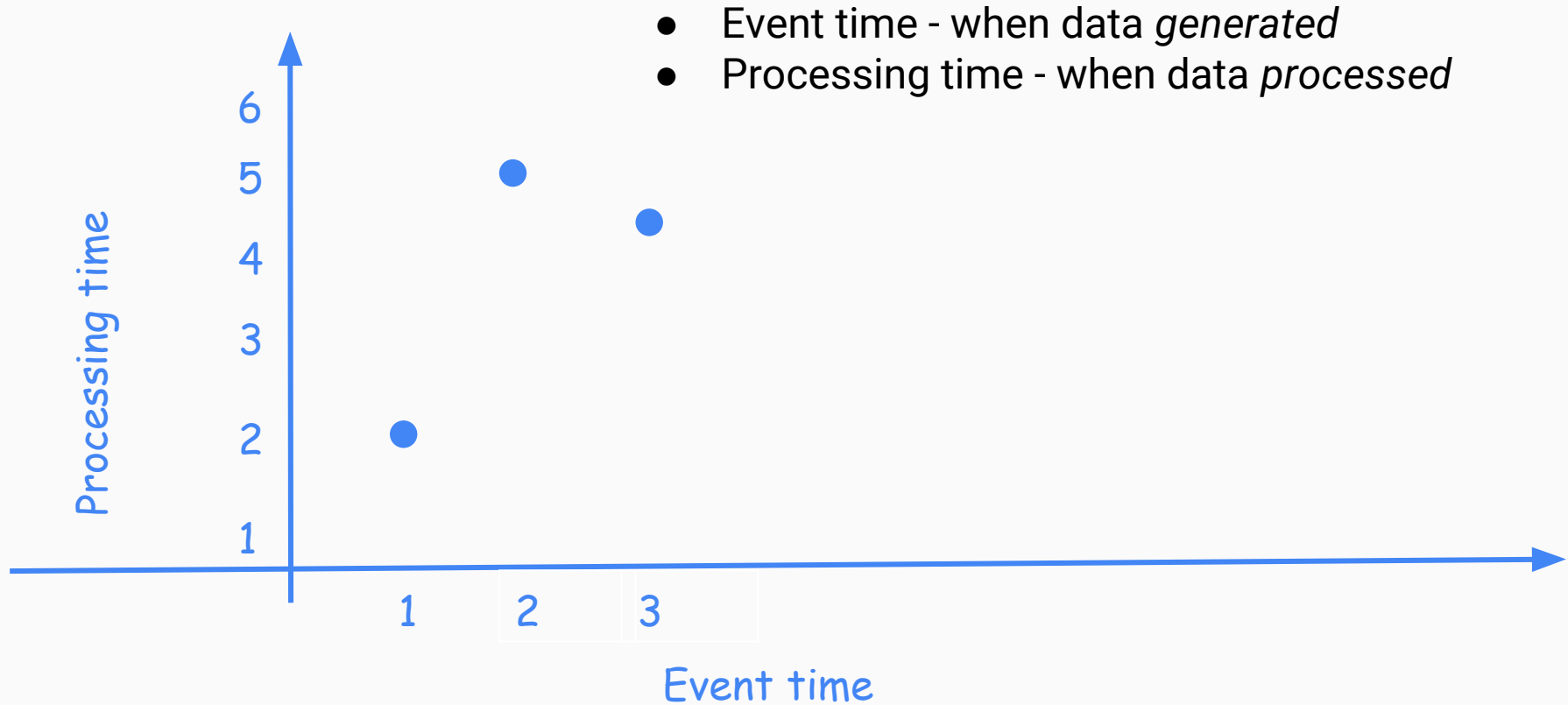


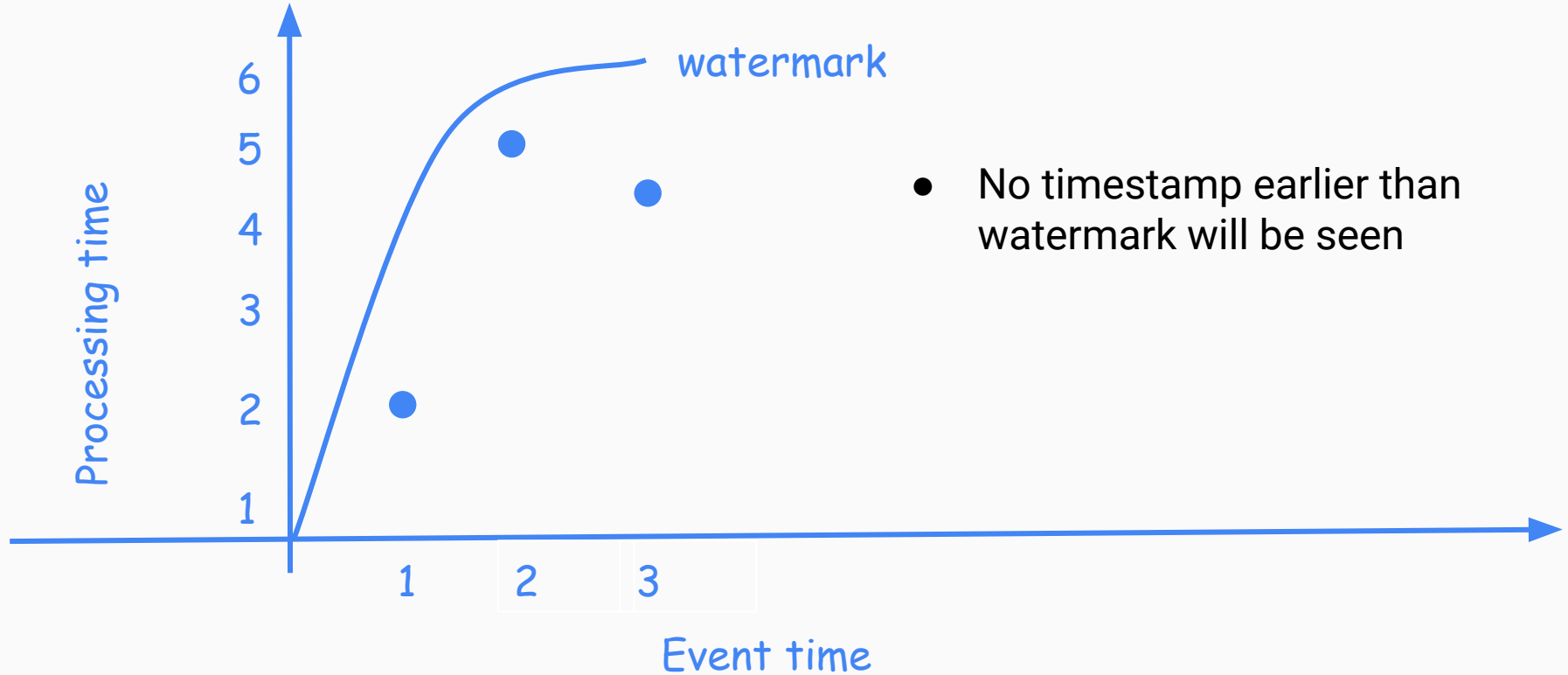


# How Gearpump solves the hard parts

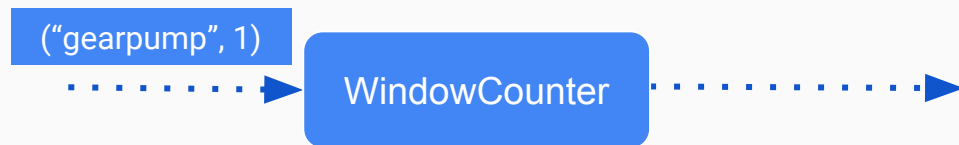
- User Interface
- Flow control
- Out-of-order processing
- Exactly Once
- Dynamic DAG

# Out-of-order data





# When can window counts be emitted ?



- No window can be emitted since message as early as time 1 has not arrived

WindowCounter In-memory Table

window	messages
<code>[0, 2)</code>	<code>("gearpump", 1)</code>
<code>[2, 4)</code>	<code>("gearpump", 3)</code> <code>("gearpump", 2)</code>
<code>[4, 6)</code>	<code>("gearpump", 5)</code> <code>("gearpump", 4)</code>

# Out-of-order processing with watermark

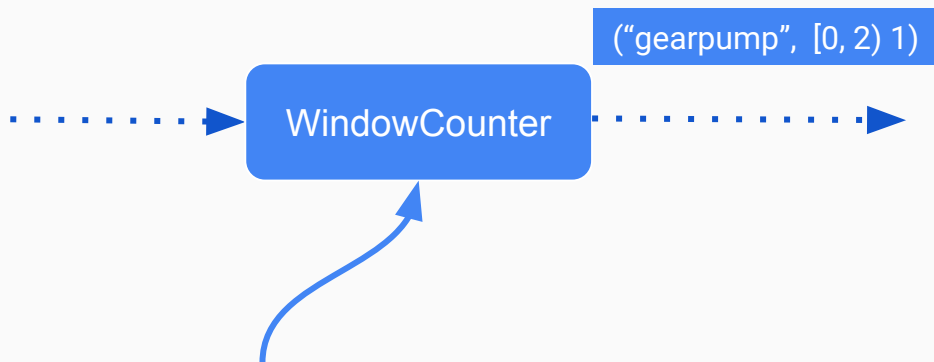


watermark = 0,  
No window can be emitted

WindowCounter In-memory Table

window	messages
[0, 2)	("gearpump", 1)
[2, 4)	("gearpump", 3) ("gearpump", 2)
[4, 6)	("gearpump", 5) ("gearpump", 4)

# Out-of-order processing with watermark

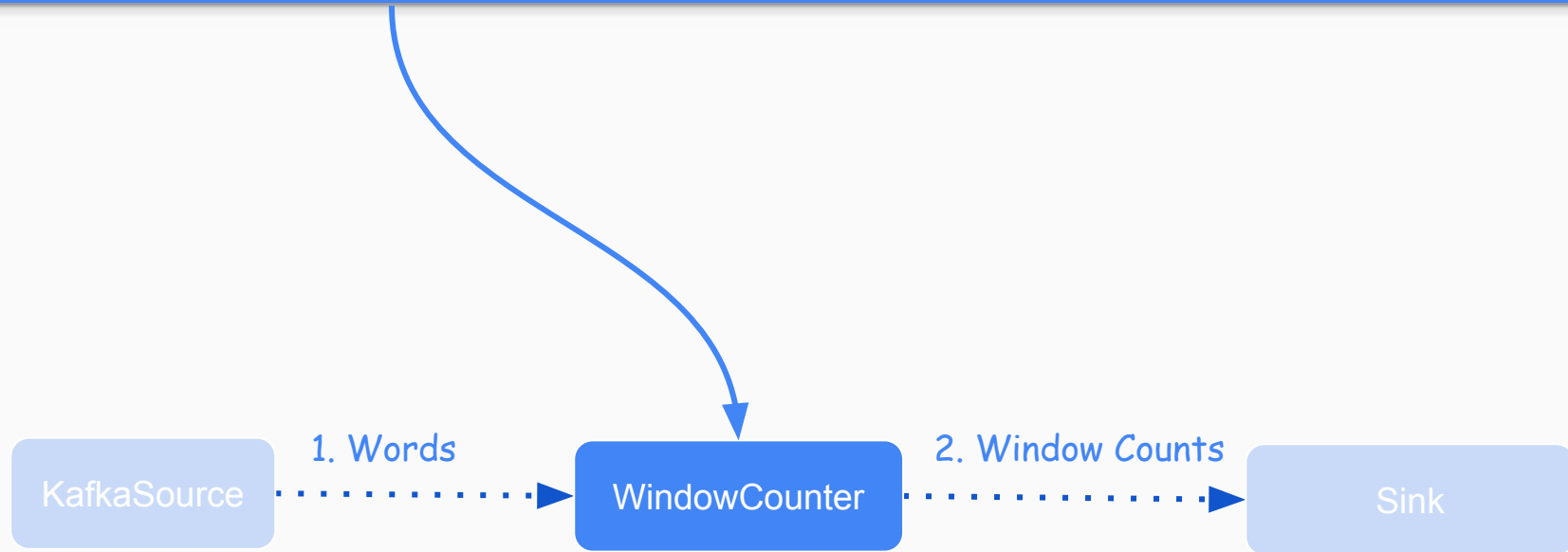


watermark = 2,  
Window [0, 2) can be emitted

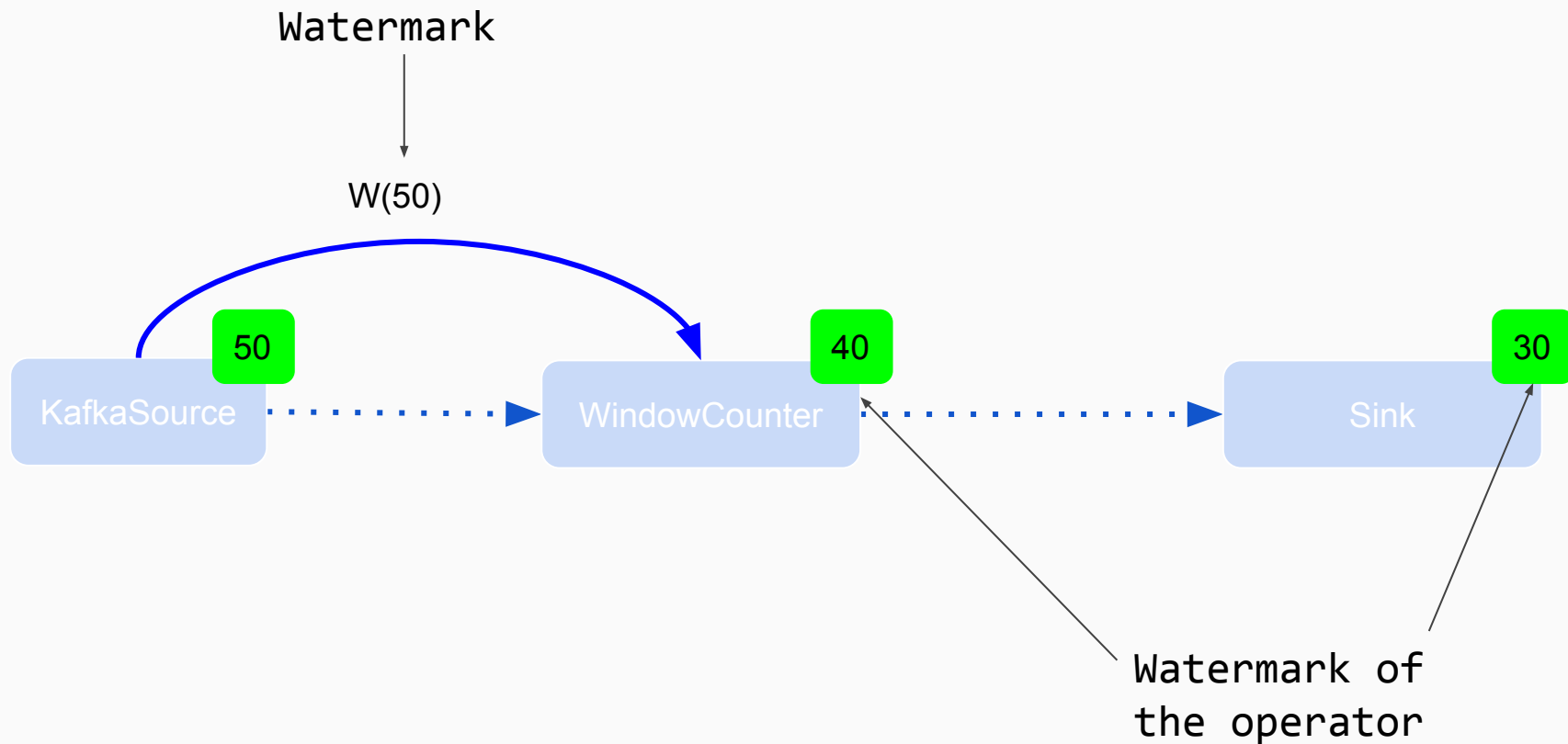
WindowCounter In-memory Table

window	messages
[2, 4)	<code>("gearpump", 3)</code> <code>("gearpump", 2)</code>
[4, 6)	<code>("gearpump", 5)</code> <code>("gearpump", 4)</code>

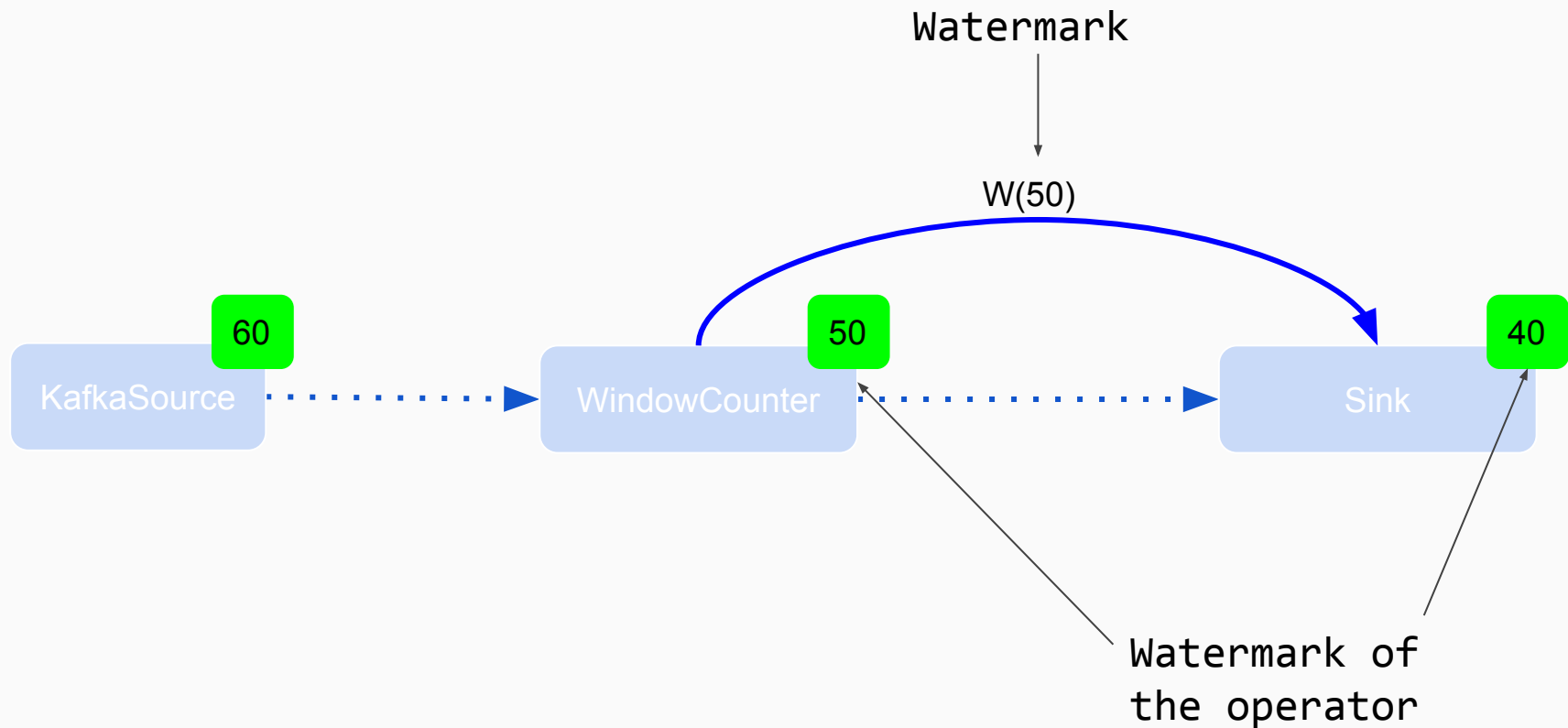
# How to get watermark ?



# From upstream







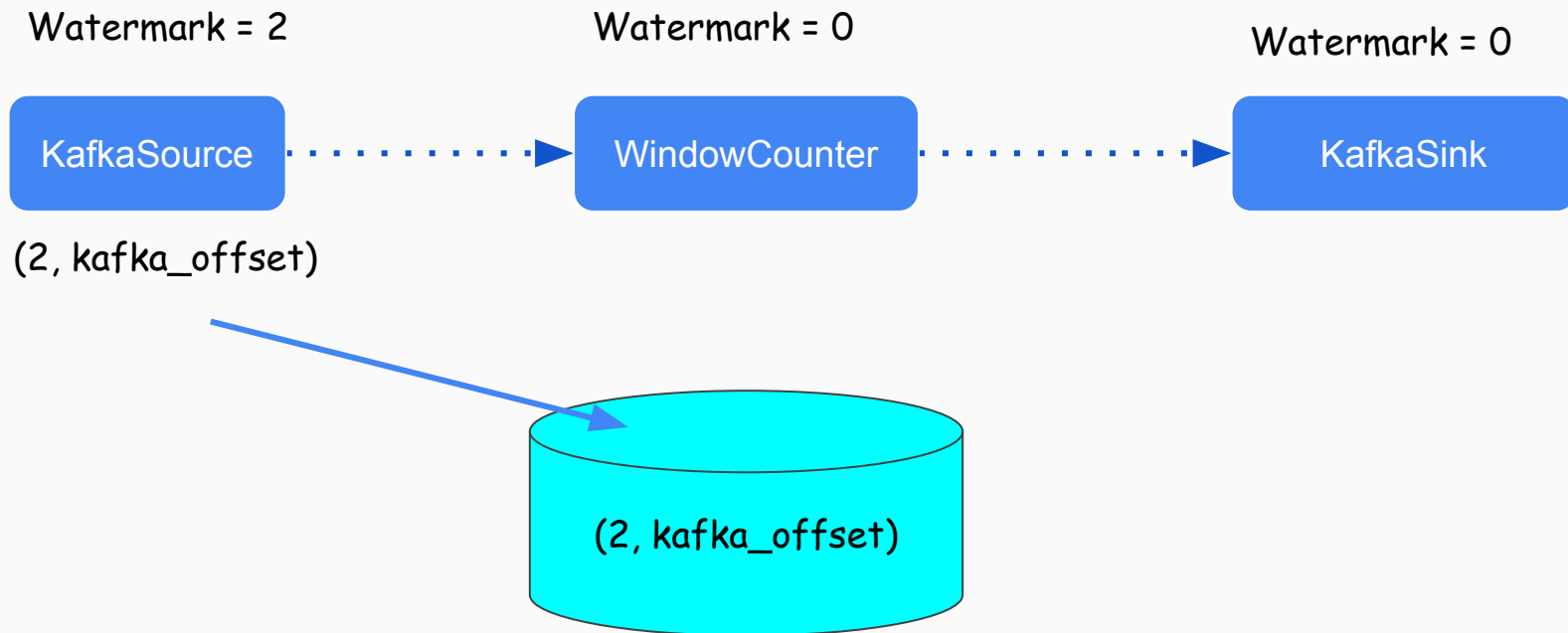
# More on Watermark

- Source watermark defined by user
- Usually heuristic based
- Users decide whether to drop data arriving after watermark

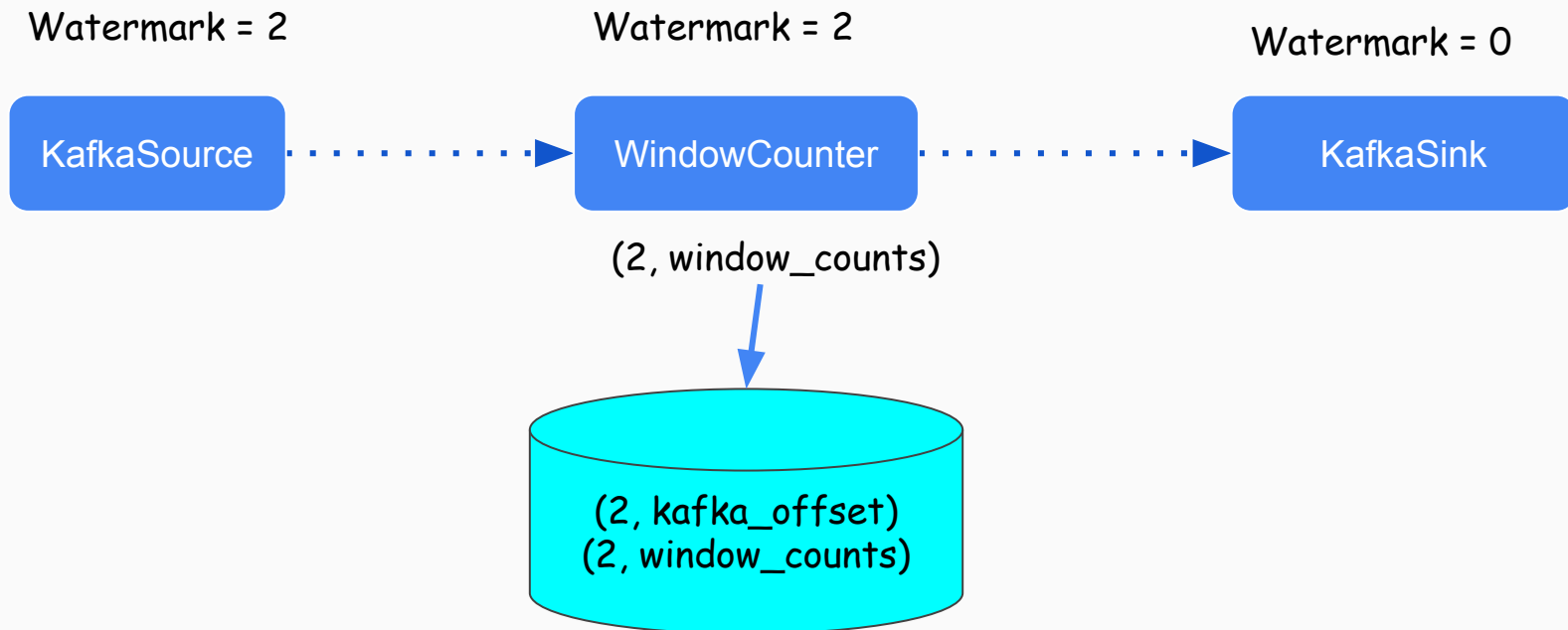
# How Gearpump solves the hard parts

- User Interface
- Flow control
- Out-of-order processing
- **Exactly once**
- Dynamic DAG

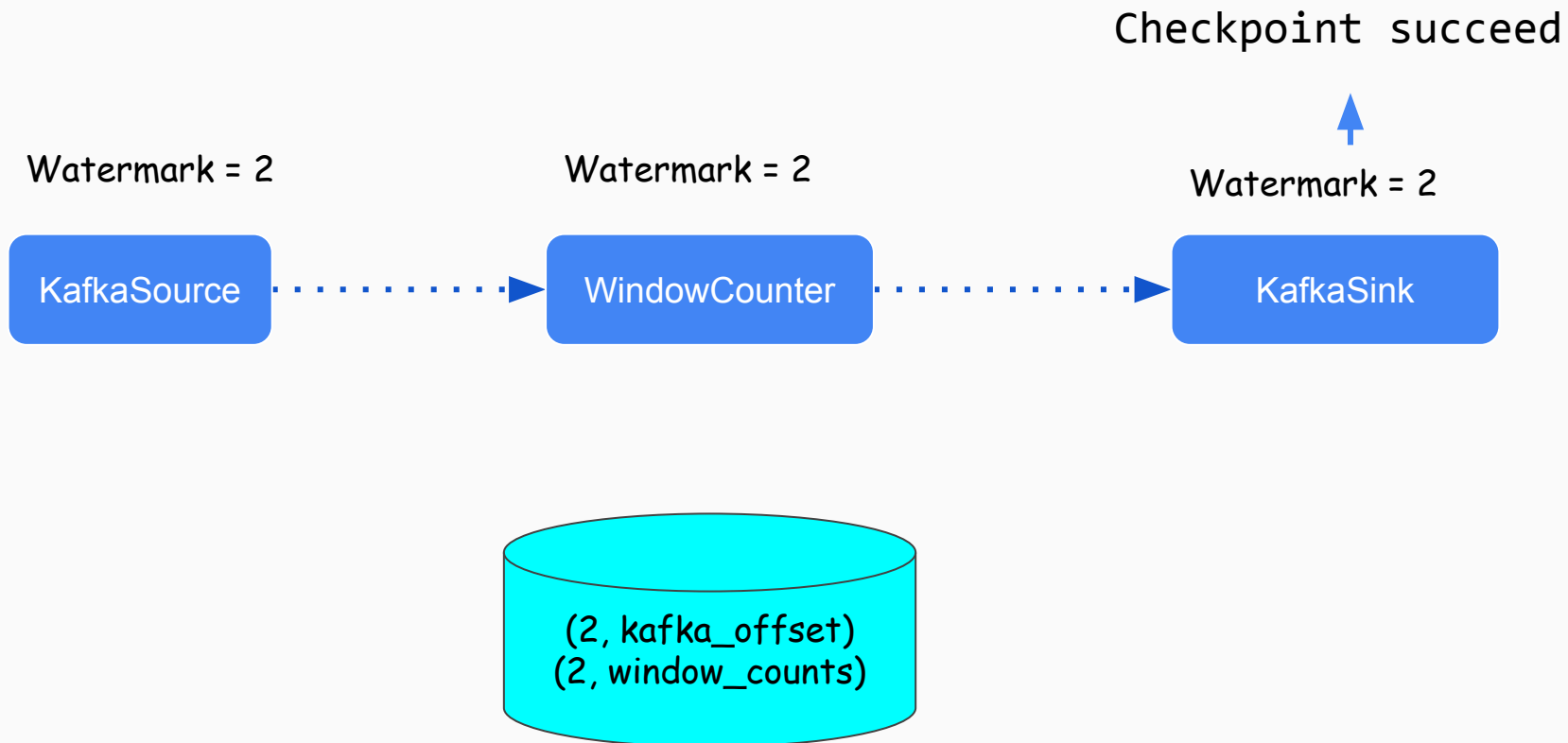
# Exactly Once with asynchronous checkpointing

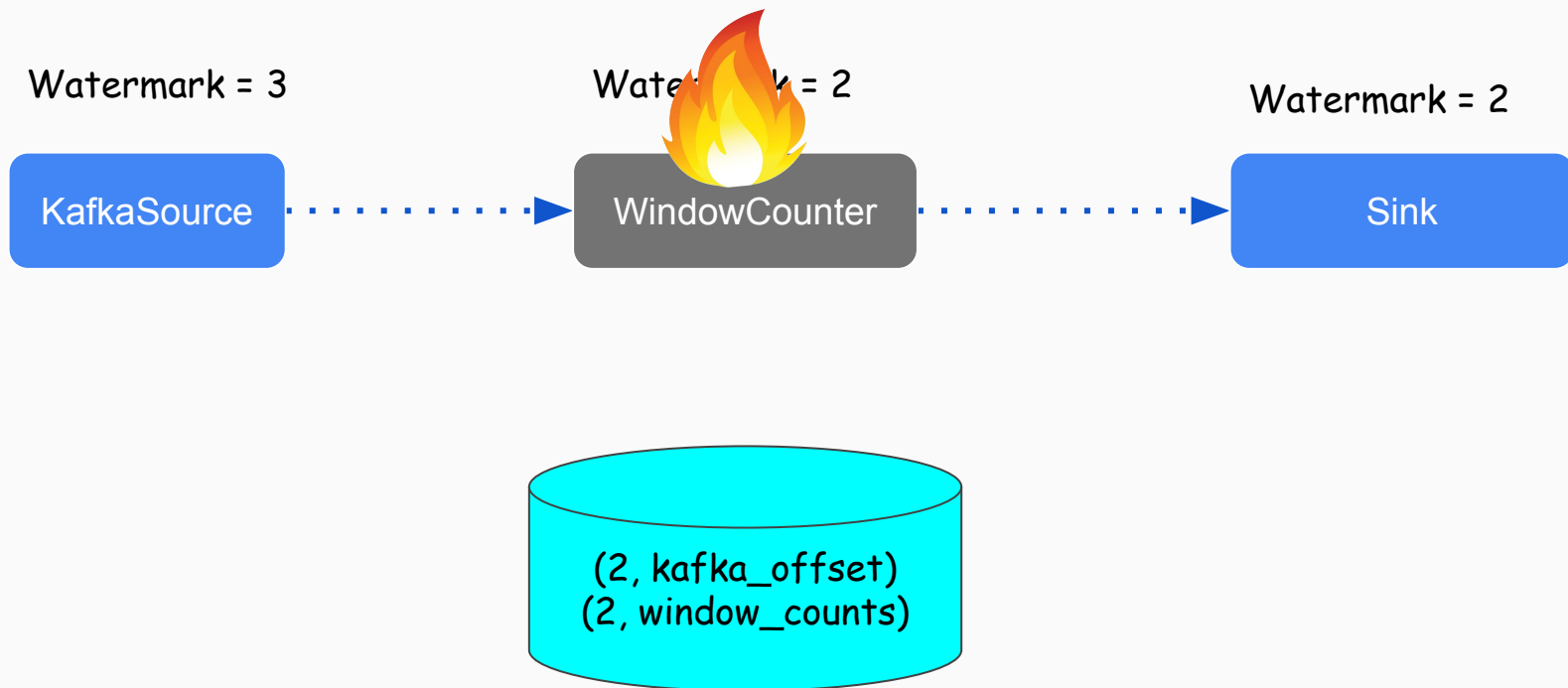


# Exactly Once with asynchronous checkpointing

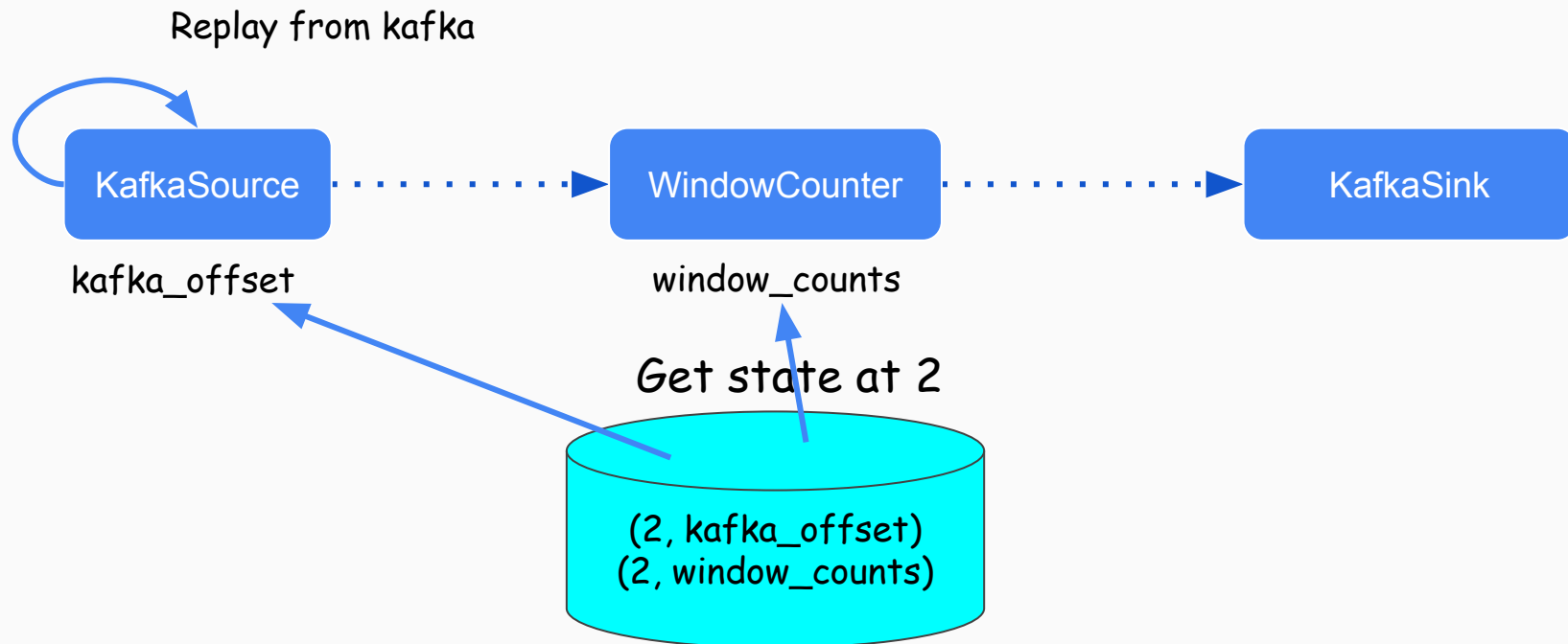


# Exactly Once with asynchronous checkpointing





# Recover to latest checkpoint at 2

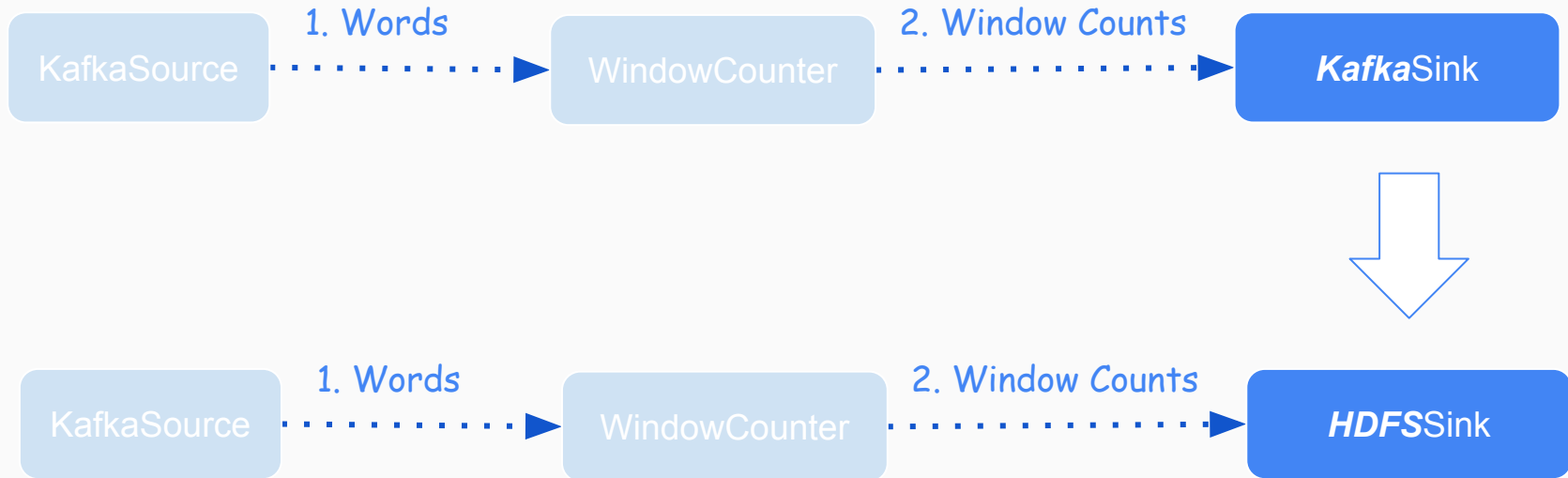




# How Gearpump solves the hard parts

- User Interface
- Flow control
- Out-of-order processing
- Exactly Once
- **Dynamic DAG**

# Update the DAG on-the-fly



**Without Restart**

# Advanced features

# DAG Visualization

AnomalyDetector ■

Uptime ⓘ  
**51 secs**

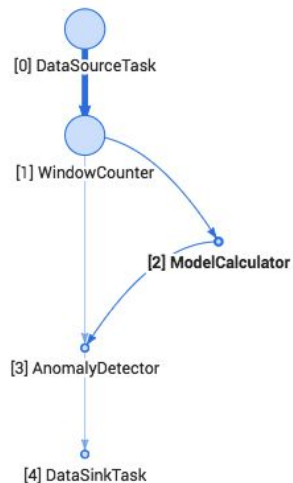
Application Clock ⓘ  
**23:55:21**  
2016/08/03

**5** processors  
**2** executors

Overview Metrics DAG

- Watermark
- Node size reflects throughput
- Edge width represents flow rate

DAG



Source Processors Send Throughput ⓘ  
**870,999** msg/s  
Total: 49,768,641

Sink Processors Receive Throughput ⓘ  
**4.20** msg/s  
Total: 209

End-to-End Latency ⓘ  
**1.67** ms

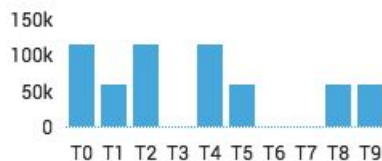
# DAG Visualization

- Data skew analysis

## Tasks ?

Message Receive Throughput

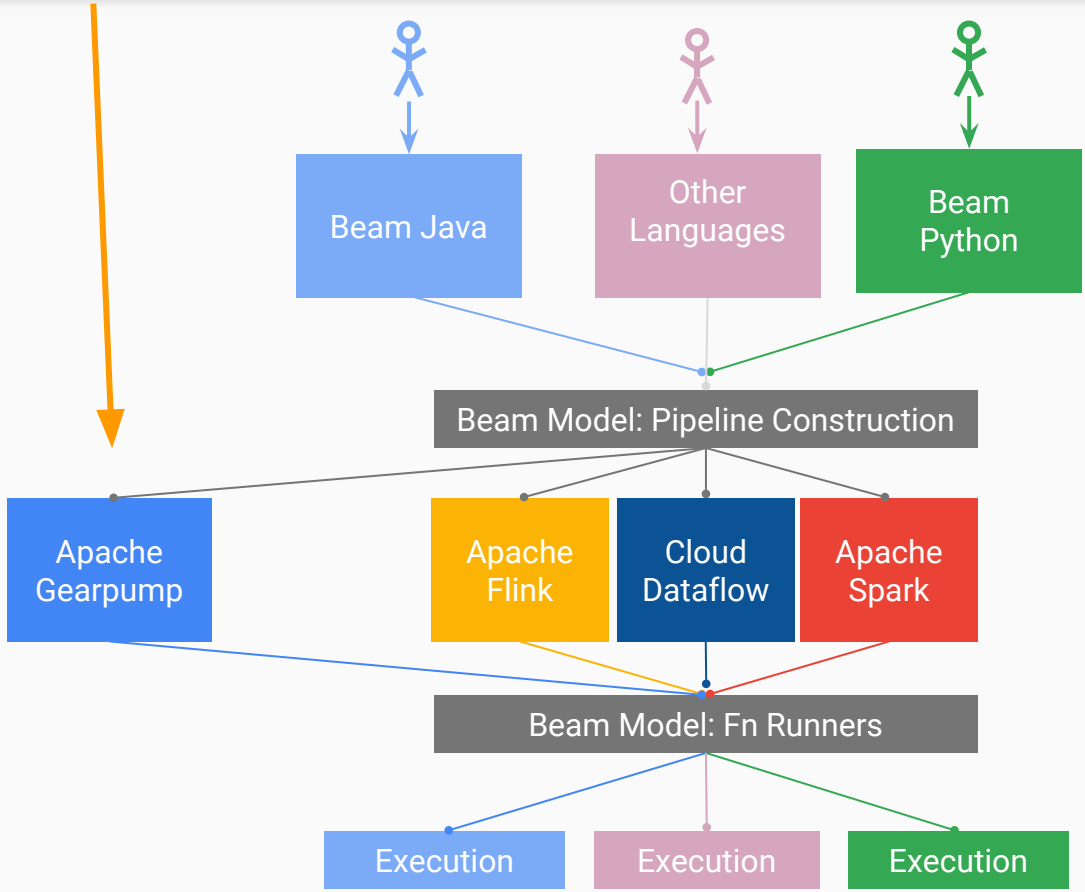
Distribution of Message Receive Throughput



Task <span>^</span>	Total Messages msg	Mean Rate msg/s	MA 1m <span>?</span> msg/s
T0	7,837,232	162,695	112,886
T1	3,915,279	81,277	56,387
T2	7,832,785	162,600	112,825
T3	0	0	0
T4	7,833,889	162,434	112,403
T5	3,917,413	81,223	56,196
T6	0	0	0
T7	0	0	0
T8	3,914,714	81,154	56,106
T9	3,914,623	81,118	56,002

Showing 1-10 of 10 records

# Apache Beam<sup>[6]</sup> Gearpump Runner



# What's next for Gearpump

# Experimental features

- Web UI Authorization / OAuth2 Authentication
- CGroup Resource Isolation
- Binary Storm compatibility
- Akka Streams integration (Gearpump Materializer)



# Summary

- Gearpump is good at streaming infinite out-of-order data and guarantees correctness
- Gearpump helps users to easily program streaming applications, get runtime information and update dynamically

# References

1. [gearpump.apache.org](http://gearpump.apache.org)
2. [akka.io](http://akka.io)
3. <http://www.slideshare.net/SeanZhong/strata-singapore-gearpumpreal-time-dagprocessing-with-akka-at-scale>
4. <https://www.cs.cmu.edu/~pavlo/courses/fall2013/static/papers/p734-akidau.pdf>
5. <https://yahooeng.tumblr.com/post/135321837876/benchmarking-streaming-computation-engines-at>
6. [Apache Beam \[project overview\]](#)
7. [www.trustedanalytics.org](http://www.trustedanalytics.org) - learn more about TAP

# Get involved

Our home: <http://gearpump.apache.org>

Contribute code: <https://github.com/apache/incubator-gearpump>

Report issues: <https://issues.apache.org/jira/browse/GEARPUMP>

The team: Kam Kasravi, Manu Zhang, Huafeng Wang, Weihua Jiang, Sean Zhong, Karol Brejna, Stanley Xu, ..., **YOU?**

TRUSTED ANALYTICS PLATFORM (TAP)

**TAP**

ABOUT BLOG PARTNERS ARCHITECTURE **DEPLOYMENTS** EVENTS RESOURCES



## Accelerate Advanced Analytics on Big Data

**Quick Links**

**Developer Repositories**  
Access projects, code and documentation

[go to GitHub](#)

**Getting Started Guide**  
Learn about installation & administration as well as working with data in TAP and in the toolkit

[go to GitHub](#)

### What is TAP?

**Trusted Analytics Platform (TAP)** is open source software, optimized for performance and security, that accelerates the creation of Cloud-native applications driven by Big Data Analytics.

TAP makes it easier for developers and data scientists at enterprises, Cloud Service Providers and System Integrators, to collaborate by providing a shared, flexible environment for advanced analytics in public and private Clouds.

The project brings together many proven and familiar open-source components, integrating them in one platform. Its capabilities are exposed as easy-to-consume services while preserving the benefit of the many vibrant communities behind these open source projects.

**TAP for Data Scientists**  
Extensible tools, scalable algorithms and powerful engines to train and deploy predictive models

**TAP for Developers**  
Consistent APIs, services and runtimes to quickly integrate these models into applications

**TAP for System Operators**  
Integrated stack can be easily provisioned in a Cloud infrastructure

# Learn More About TAP

[www.trustedanalytics.org](http://www.trustedanalytics.org)

# Engage in Community events

Meetups, workshops, & webinars

<http://trustedanalytics.org/#resources>

