



**KubeCon**



**CloudNativeCon**

**North America 2023**





KubeCon



CloudNativeCon

North America 2023

# Building a Scalable and Reliable Change Data Capture for TiKV

*Charles Zheng <charlesz@netflix.com>*

- **Why TiCDC**
- **Challenges and design goals**
- **How TiCDC works internally**
  - Overall Workflow
  - TiCDC Cluster Internal
  - Pipeline: Puller, Sorter, Mounter, Sink
- **Performance improvement**
- **Lessons learned**



KubeCon



CloudNativeCon

North America 2023

# *Why TiCDC*

# Two Common Scenarios



KubeCon



CloudNativeCon

North America 2023

- 01 Incremental data synchronization services for heterogeneous system.
- 02 Cross-region disaster recovery services based on primary and secondary replication

# What TiCDC can do



KubeCon



CloudNativeCon

North America 2023

- **Low-latency incremental data replication for various downstreams**

*TiDB -> TiCDC -> MySQL: Escape Link*

*TiDB -> TiCDC -> Kafka: With Canal-JSON, Avro, Open-Protocol.*

*TiDB -> TiCDC -> S3: With CSV data format*

- **Support database and table filtering**
- **Support most operation through Open API**
- **Support bi-directional replication between DB clusters**



KubeCon



CloudNativeCon

North America 2023

# ***Design Goals and Challenges***

- **High Availability**
  - *Partial nodes crash would not interrupt the data sync*
- **Ensure High Throughput and Low Latency**
  - *Sync large volumes of change events concurrently*
  - *With relatively low latency*
- **Ensure Consistency and Ordering**
  - *Snapshot isolation & Eventual consistency*



# Challenges



KubeCon



CloudNativeCon

North America 2023

- **Capture the change data instantaneously**
- **Aware of and catch up with the schema evolution**
- **Tradeoff between ordering and high throughput**
- **Tradeoff between consistency and low latency**
- **Fetch data spread across multiple nodes**
- **Minimize operational complexity**



KubeCon



CloudNativeCon

North America 2023

# *How CDC Works*

# Overall Architecture

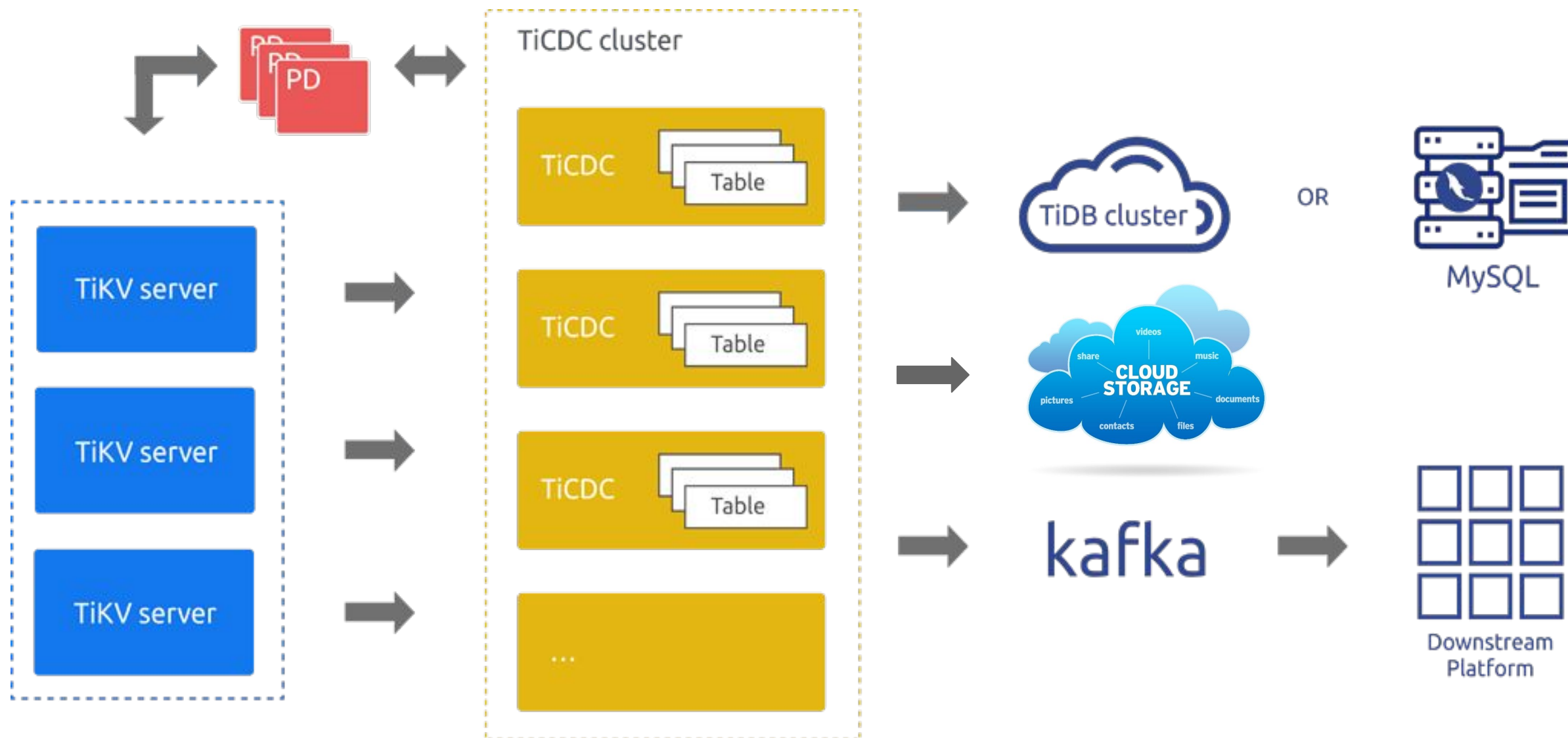


KubeCon



CloudNativeCon

North America 2023



# TiCDC Cluster Internal

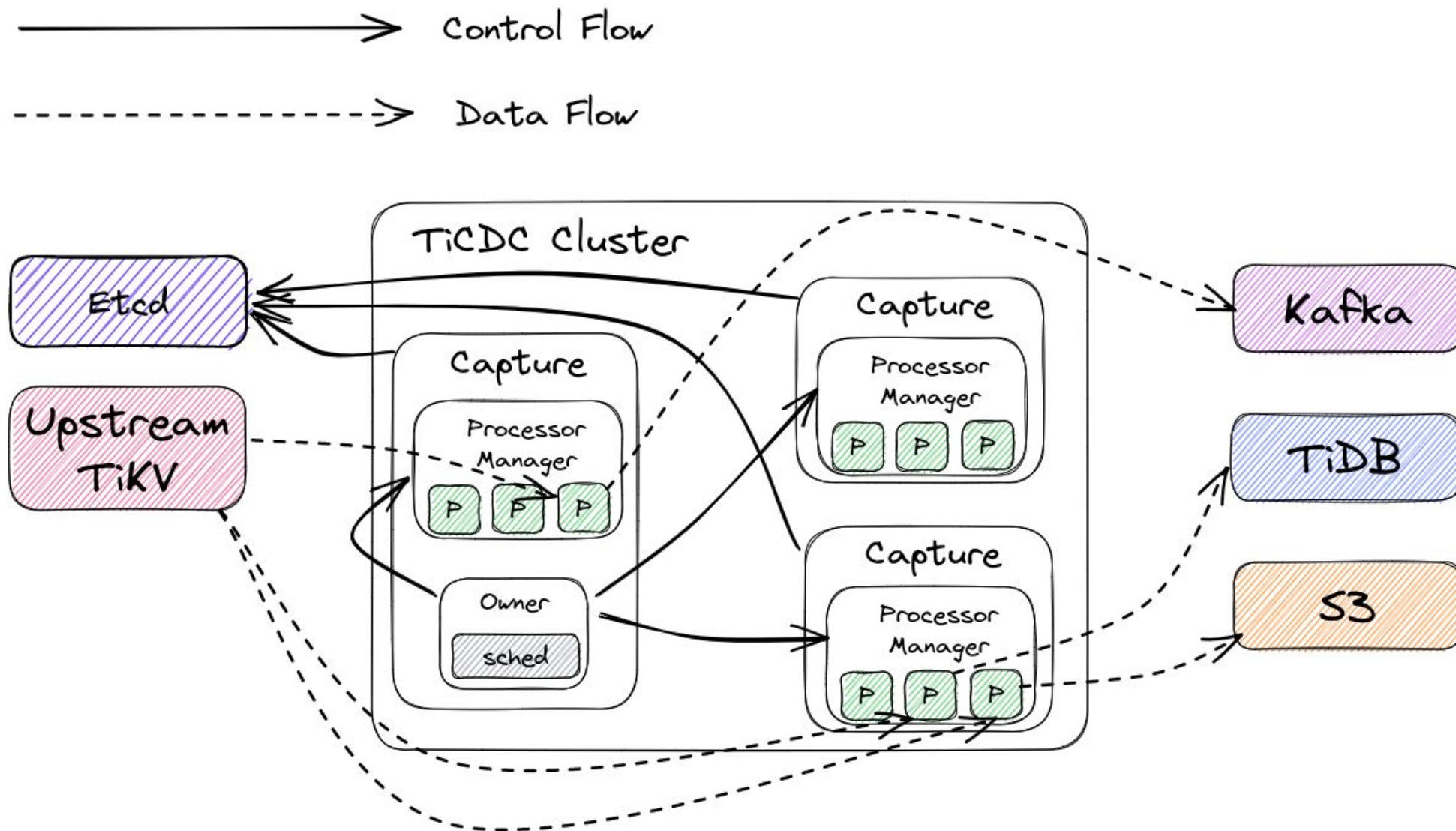


KubeCon



CloudNativeCon

North America 2023



# Inside a Capture Server

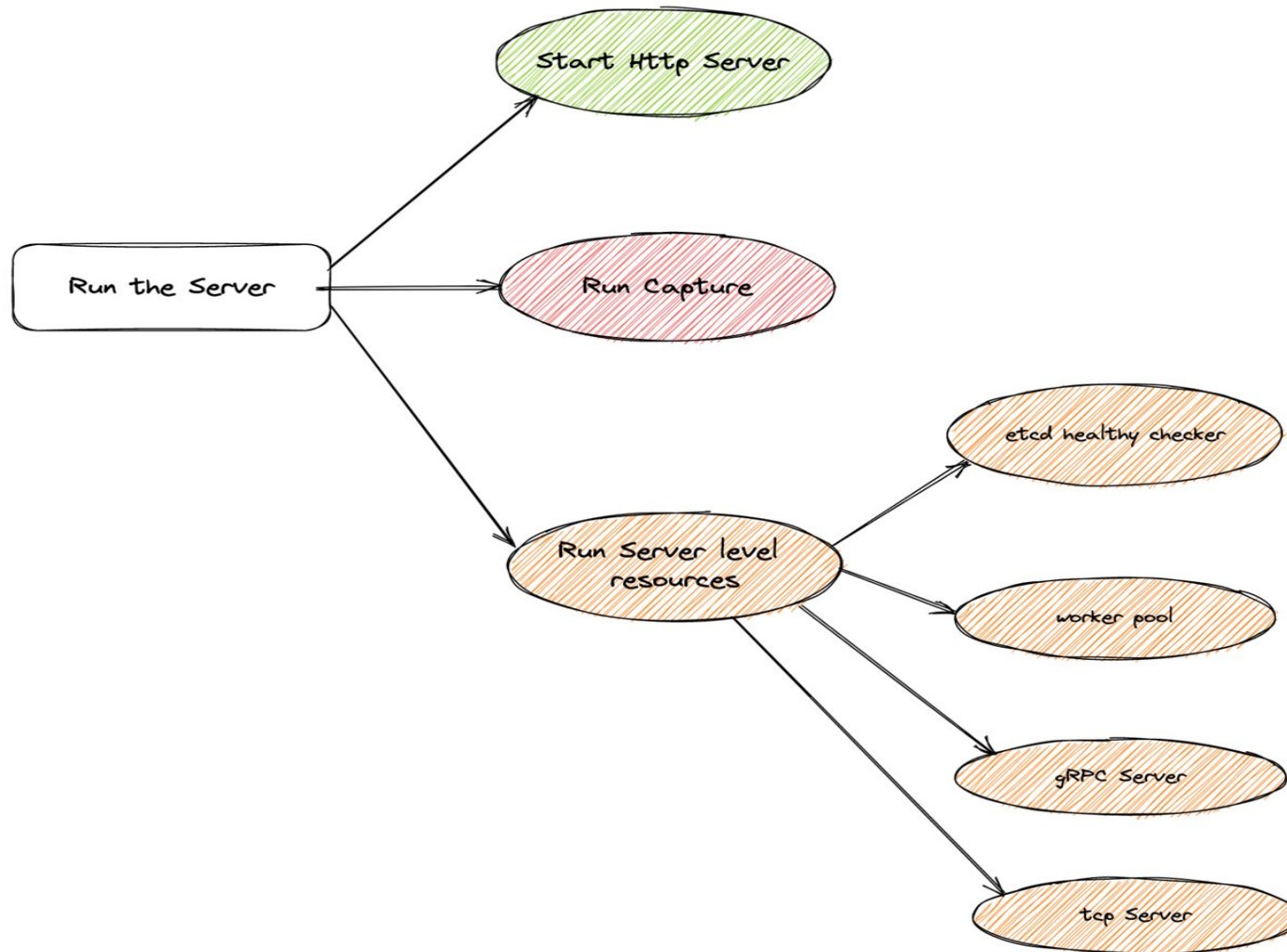


KubeCon



CloudNativeCon

North America 2023



# Campaign

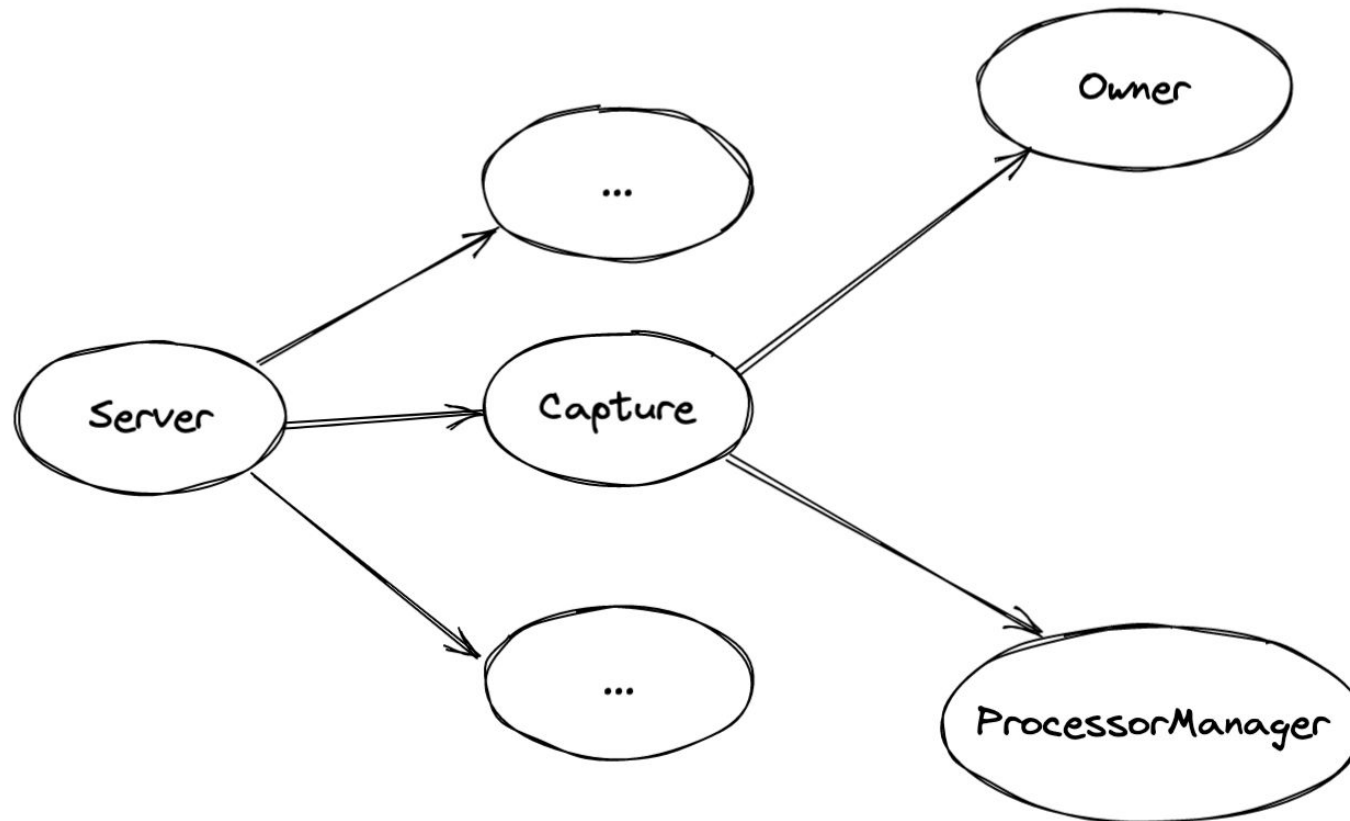


KubeCon



CloudNativeCon

North America 2023



# ChangeFeed Scheduling

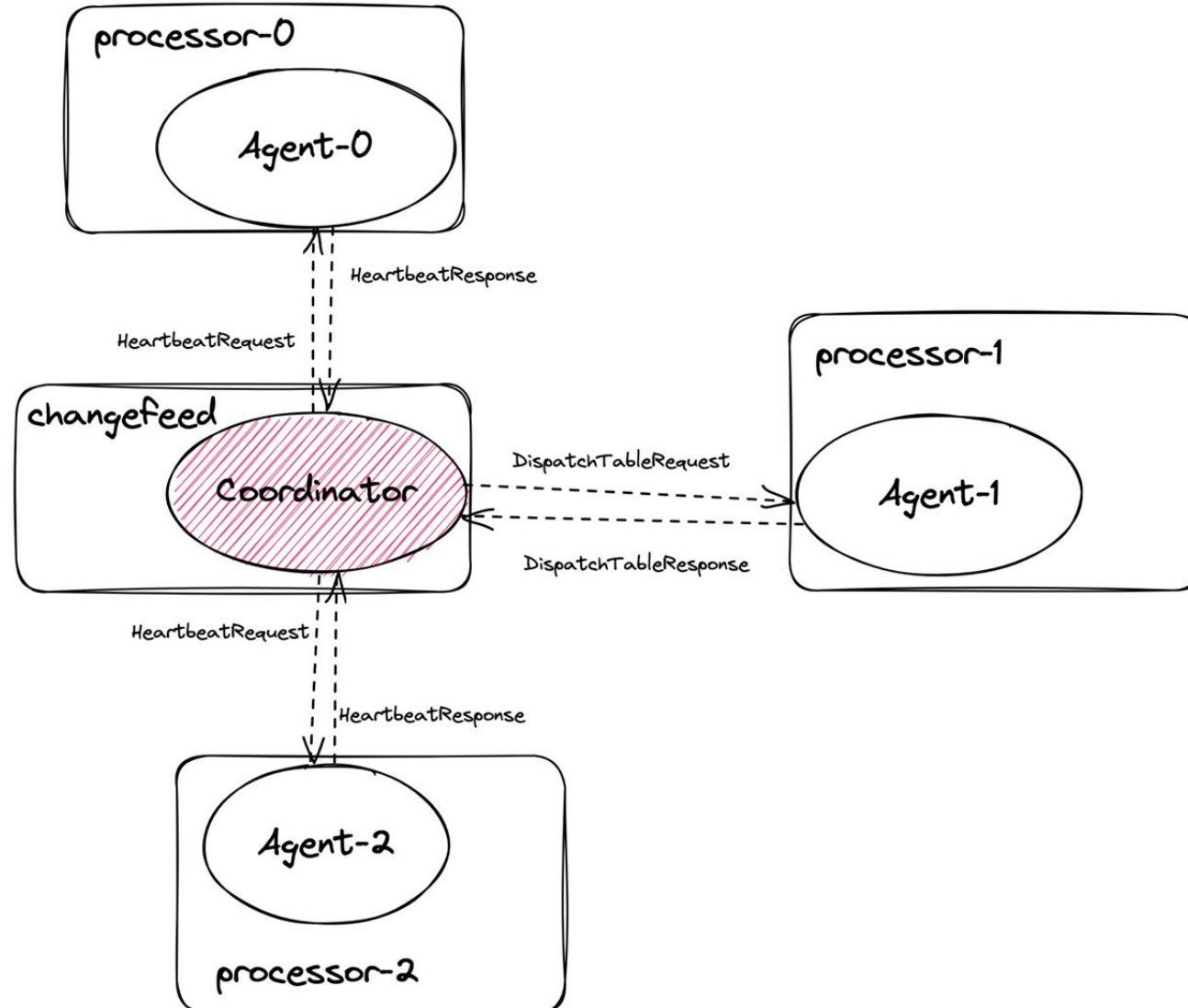


KubeCon



CloudNativeCon

North America 2023





# Topology

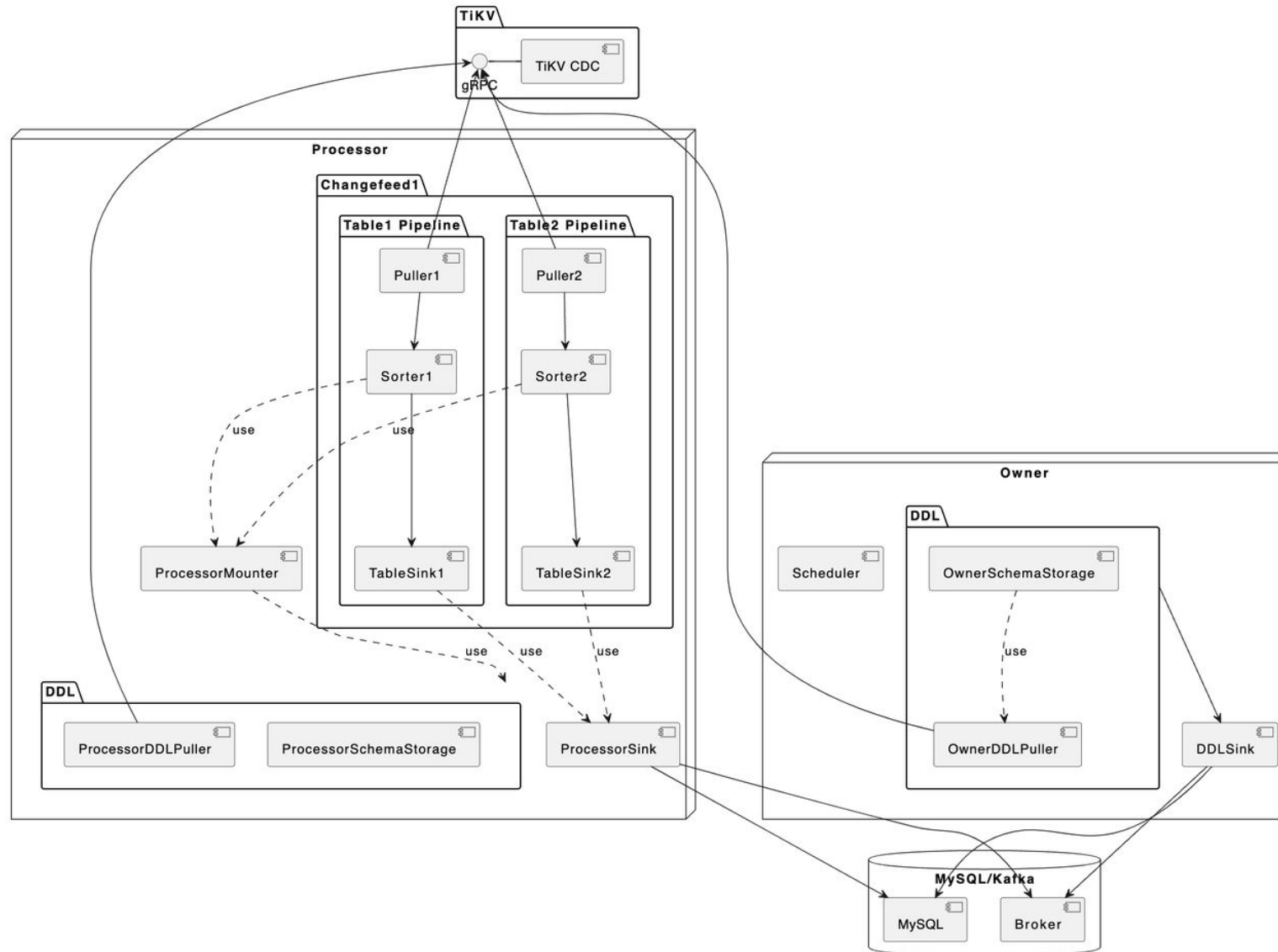


KubeCon



CloudNativeCon

North America 2023





# Pipeline

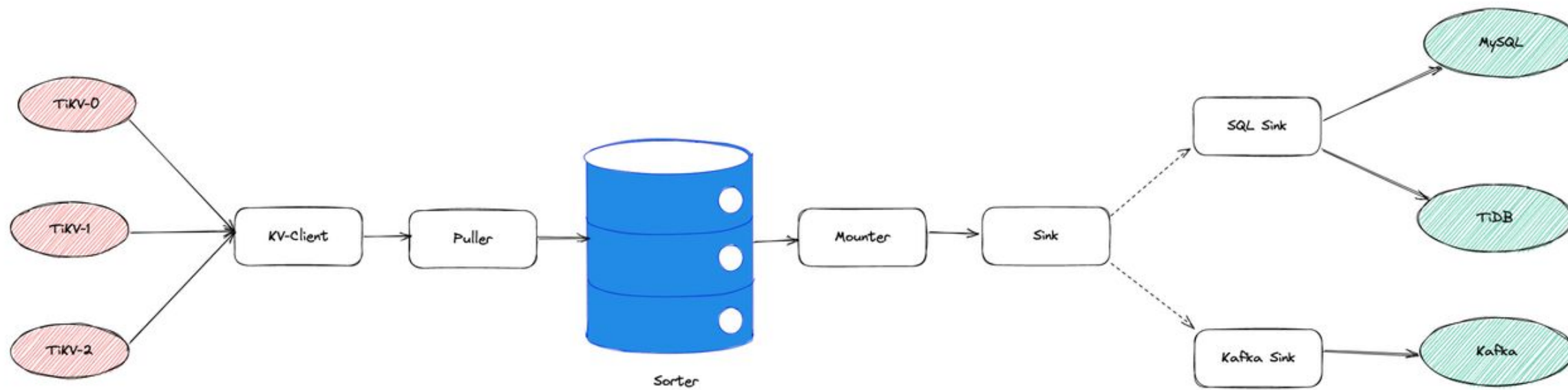


KubeCon



CloudNativeCon

North America 2023



# How Puller Work

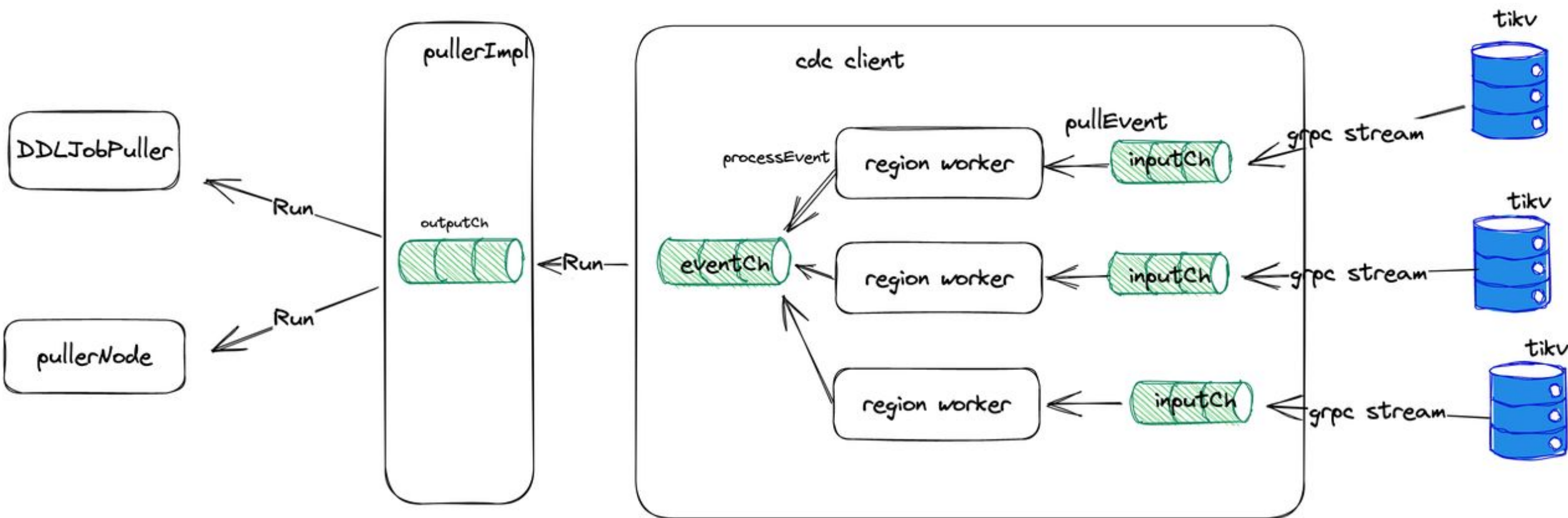


KubeCon



CloudNativeCon

North America 2023



# Why Sorter



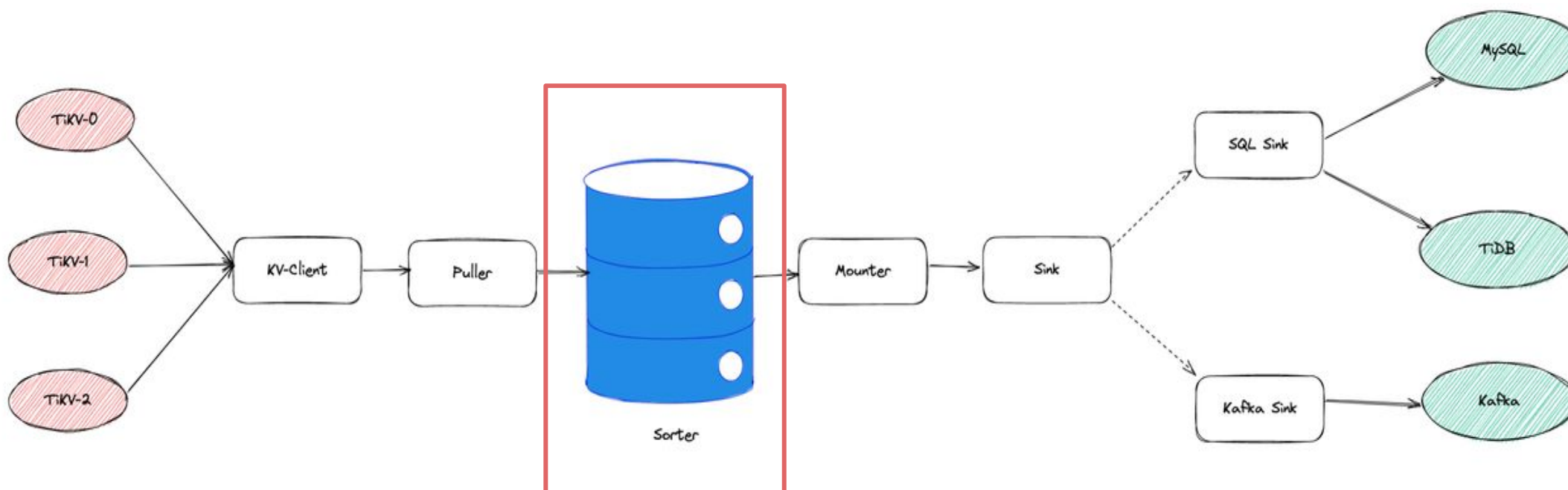
KubeCon



CloudNativeCon

North America 2023

- **Buffer:** Smooth out the peaks and valleys of upstream data flow
- **Sort:** Incoming events may not be in chronological order



# How Sorter Work

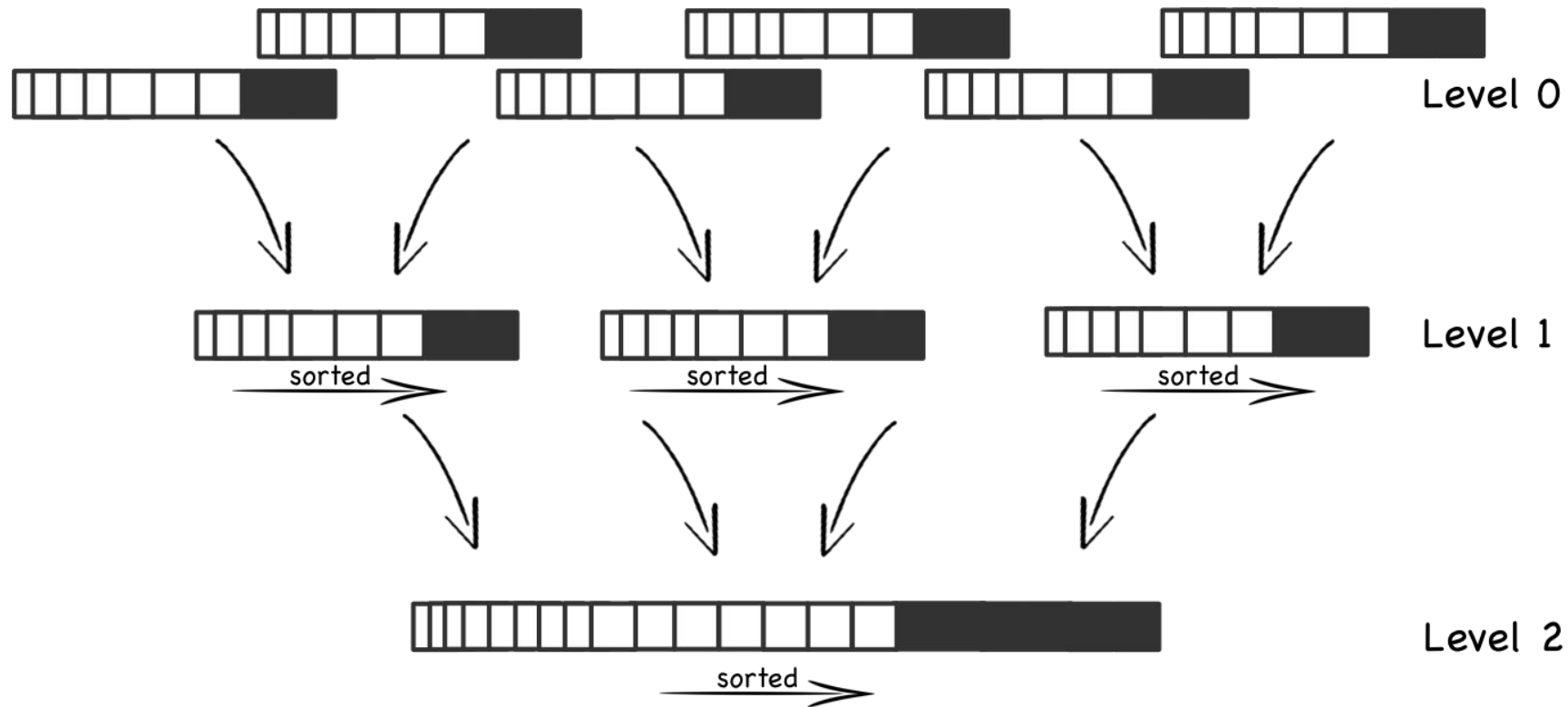


KubeCon



CloudNativeCon

North America 2023



Compaction continues creating fewer, larger and larger files

# How Mounter Works

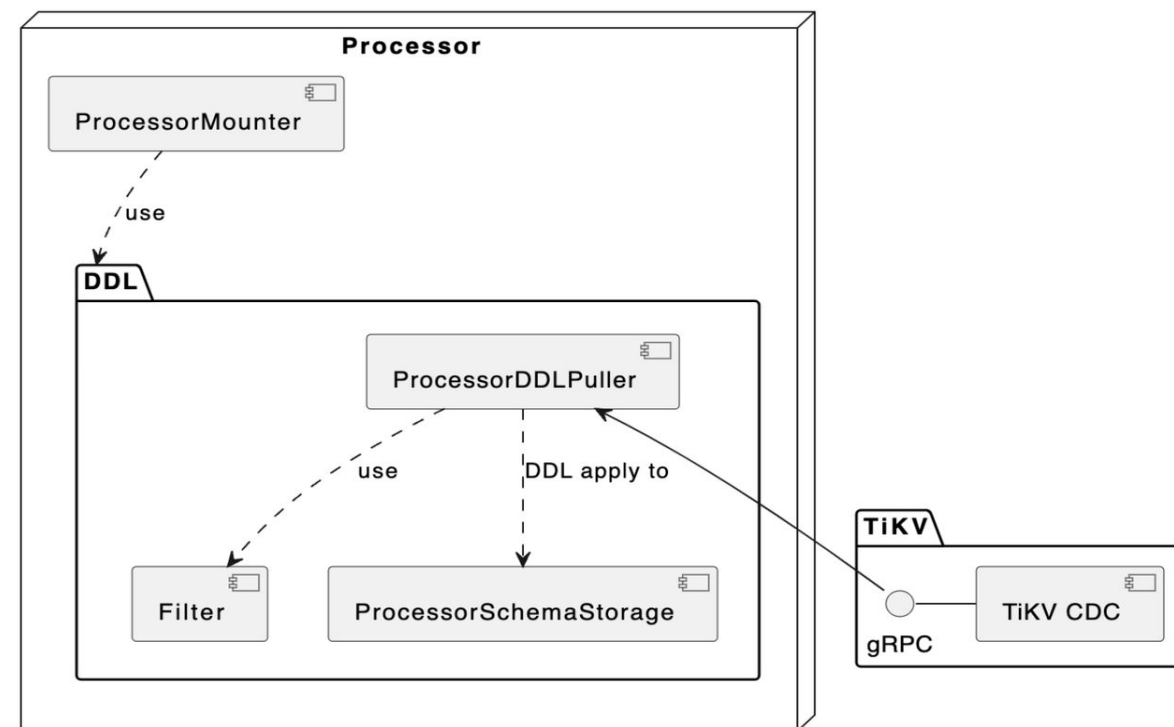
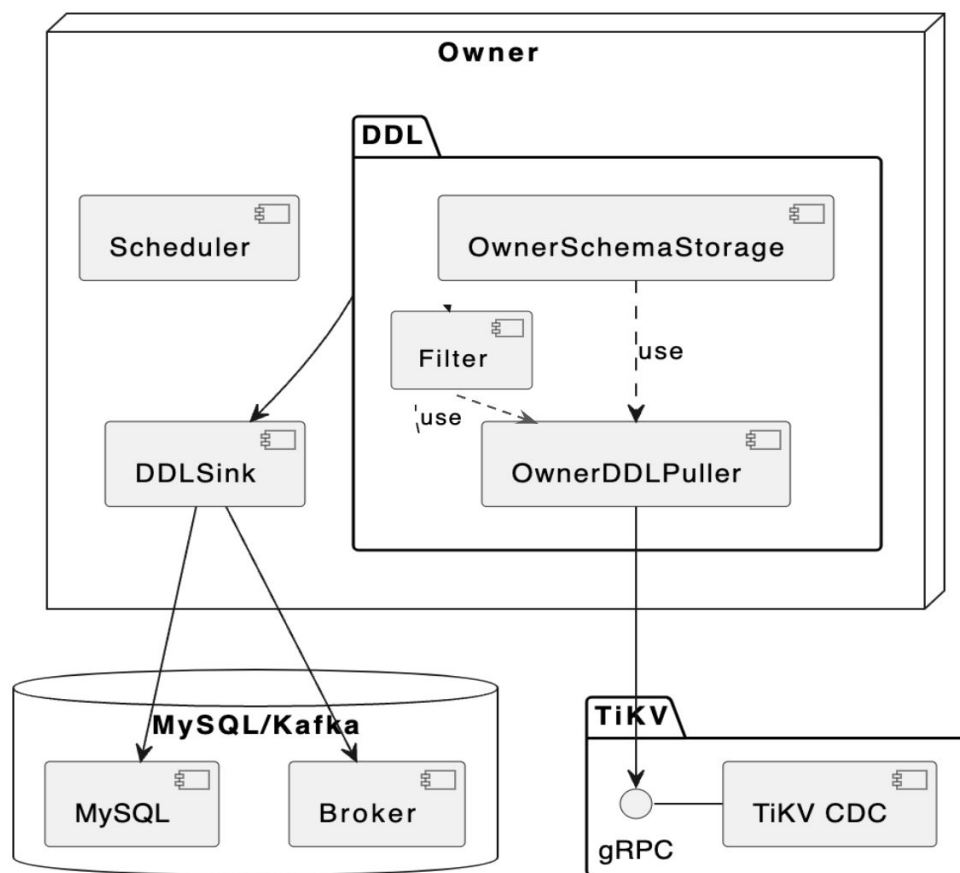


KubeCon



CloudNativeCon

North America 2023



# How Sinkers Works

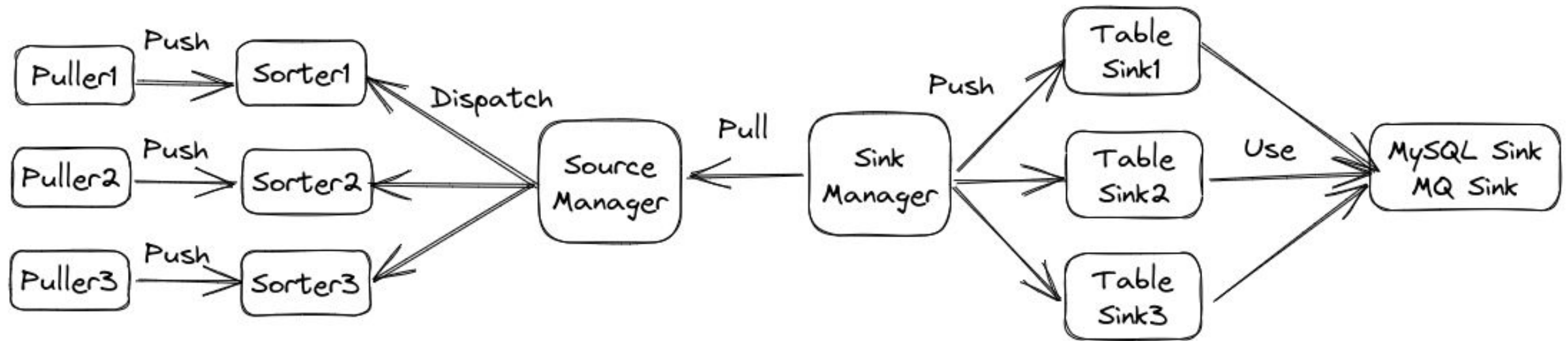


KubeCon



CloudNativeCon

North America 2023





KubeCon



CloudNativeCon

North America 2023

# ***Performance Improvement***

***TiDB -> TiCDC -> TiDB***

1 TiCDC node 16c/64g

Single table throughput	Big table (1200 bytes/row) 80K write QPS (insert only) throughput: 120MB MB/s
Max single table size	Up to 30~40T
Max upstream cluster data size	No limit



## *TiDB -> TiCDC -> Kafka*

1 TiCDC node 16c/64g

Single table throughput	Big table (1200 bytes/row) 35K write QPS (only insert) throughput: 52.8 MB/s
Max single table size	Up to 30~40T
Max upstream cluster data size	No limit

# Throughput by Components



KubeCon  
— North America 2023 —



CloudNativeCon  
— North America 2023 —

	sysbench	ccb	jitv	Average
puller	158k	112k	97k	128k
sorter	522k	214k	131k	250k
mounter	262k	120k	78k	161k
sink	TBD	TBD	22k (kafka sink)	76k

Sink is the bottleneck in most cases

# How to Improve Puller



KubeCon



CloudNativeCon

North America 2023

- **Process the *ResolvedTS* in batch**
- **Use read/write lock instead of the exclusive lock**
- **Optimize how the frontier inspect the region split/merge**
- **Remove unnecessary memory allocation**

# How to Improve Mounter



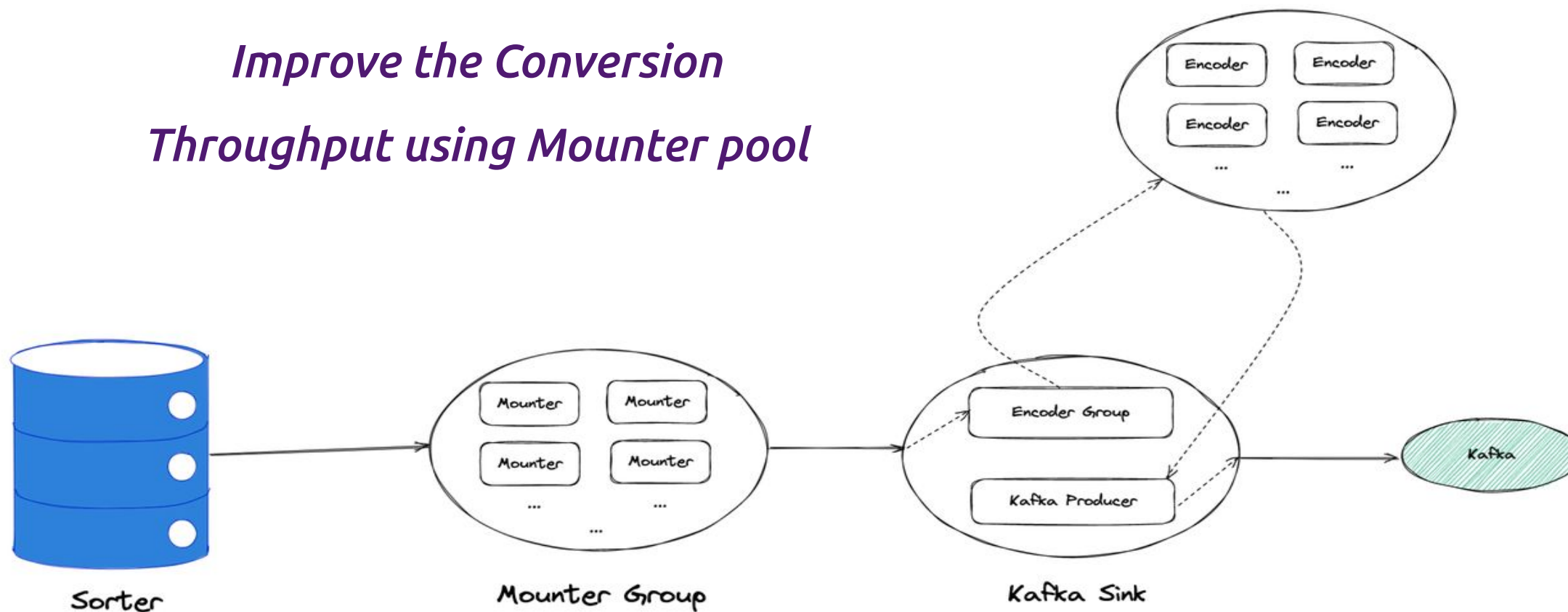
KubeCon



CloudNativeCon

North America 2023

*Improve the Conversion  
Throughput using Mounter pool*



# How to Improve Sink

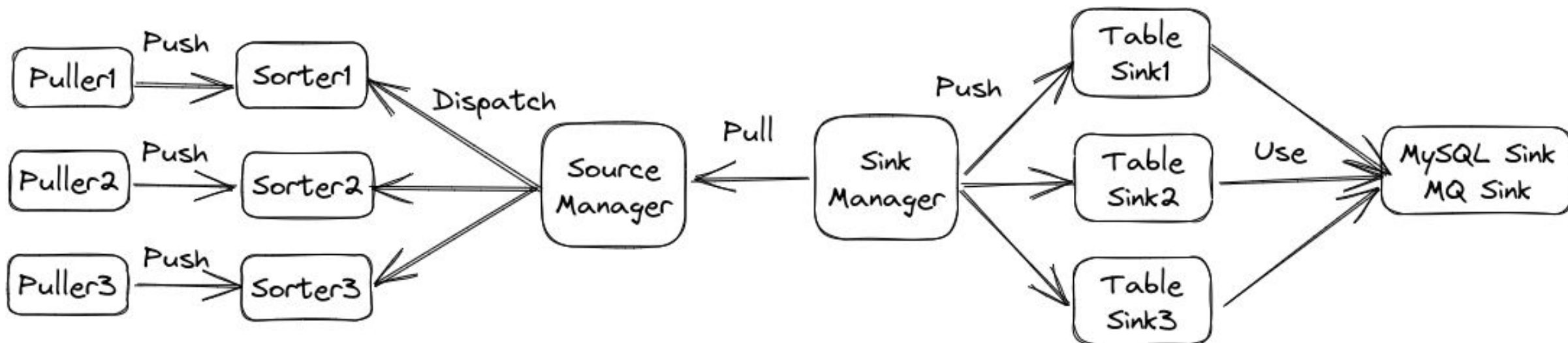


KubeCon



CloudNativeCon

North America 2023



# Performance Improvement (to Kafka)



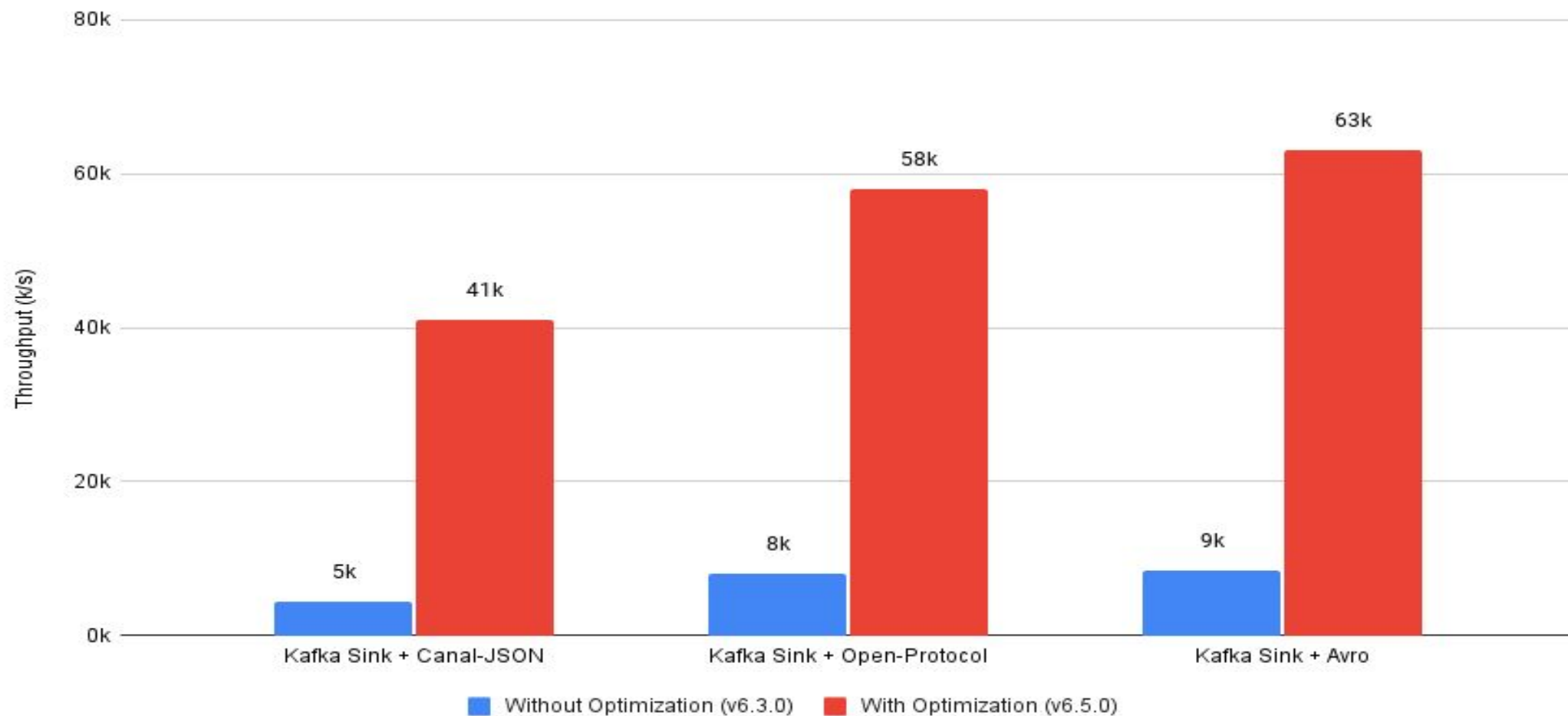
KubeCon



CloudNativeCon

North America 2023

## Comparison Of TiCDC Sink Throughput



# Performance Improvement (to MySQL)



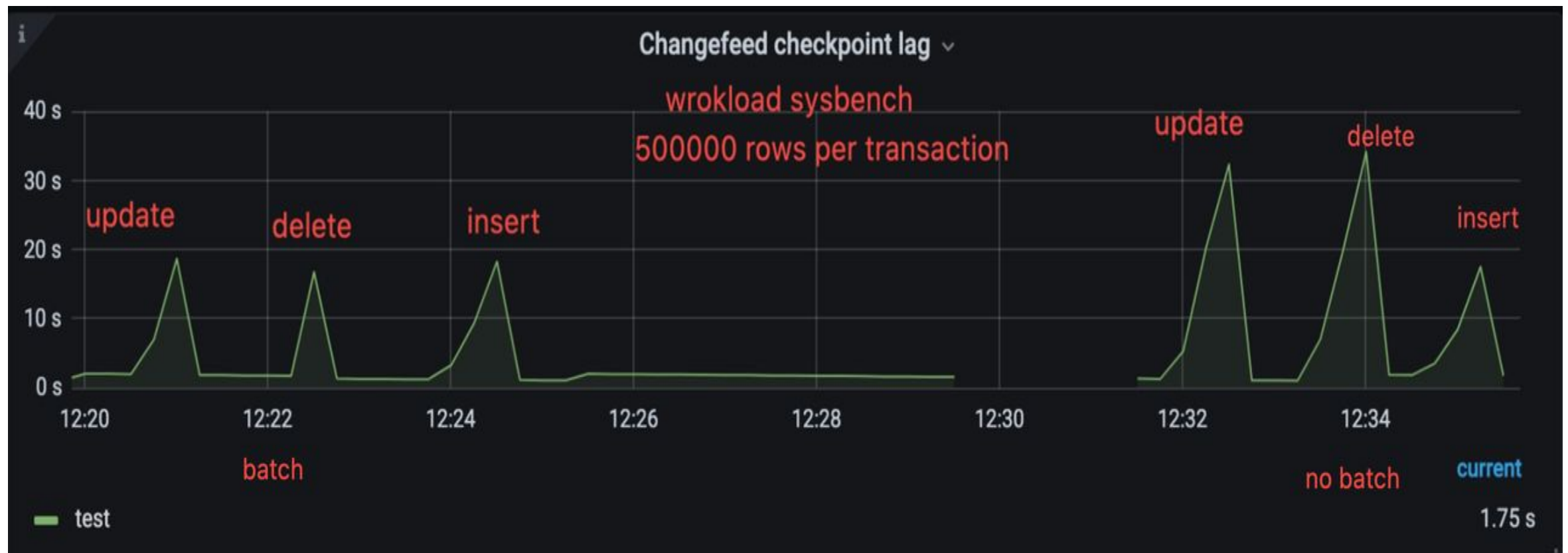
KubeCon



CloudNativeCon

North America 2023

## Sync DML events in batch





KubeCon



CloudNativeCon

North America 2023

# ***Lessons Learned***





- **Design the system architecture align with upstream Database**
- **Clear boundary between subcomponents**
- **Choose the push model and pull model carefully**
- **Implement the old value feature from day 1**
- **Implementing an efficient sorter as well as keeping the scale in mind**



KubeCon



CloudNativeCon

North America 2023

# THANKS !

**TiCDC:** <https://github.com/pingcap/tiflow>

**Any Questions?**

find me at

**Charles Zheng** <charlesz@netflix.com>



PromCon  
North America 2021



**Please scan the QR Code above  
to leave feedback on this session**