# Micro-Batch, Streaming and Serveless
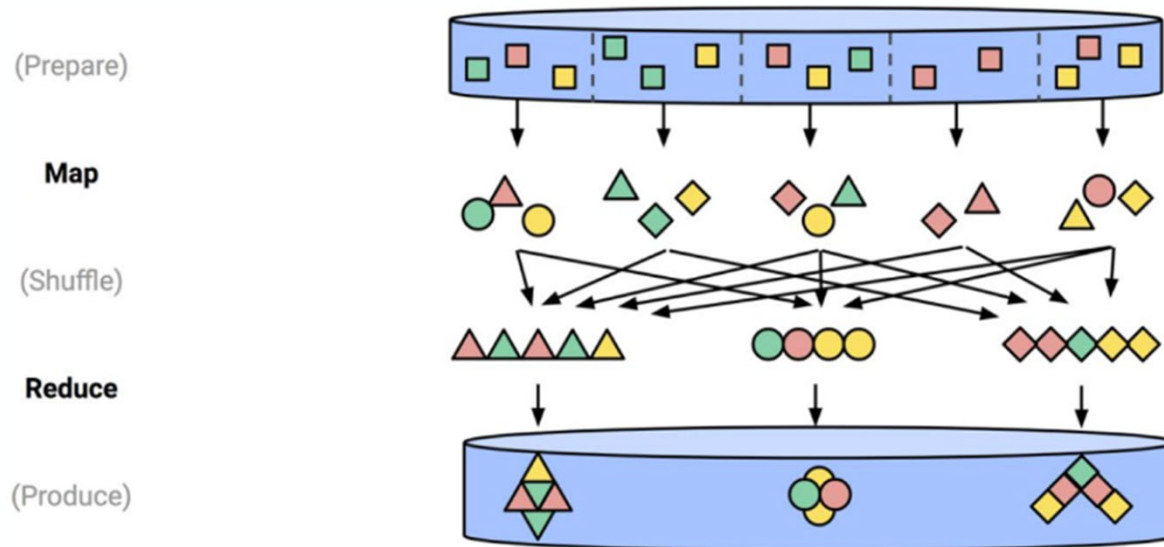
@Luminous Moonlight

Presented by KONY128

2020-11-21

# Map Reduce Review



(Prepare)

Map

(Shuffle)

Reduce

(Produce)

*Google MapReduce https://research.google/pubs/pub62/*

# RDD Review

Narrow Dependencies:

map, filter

union

join with inputs co-partitioned

Wide Dependencies:

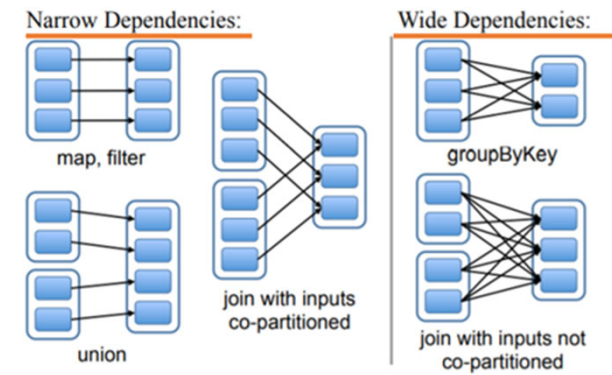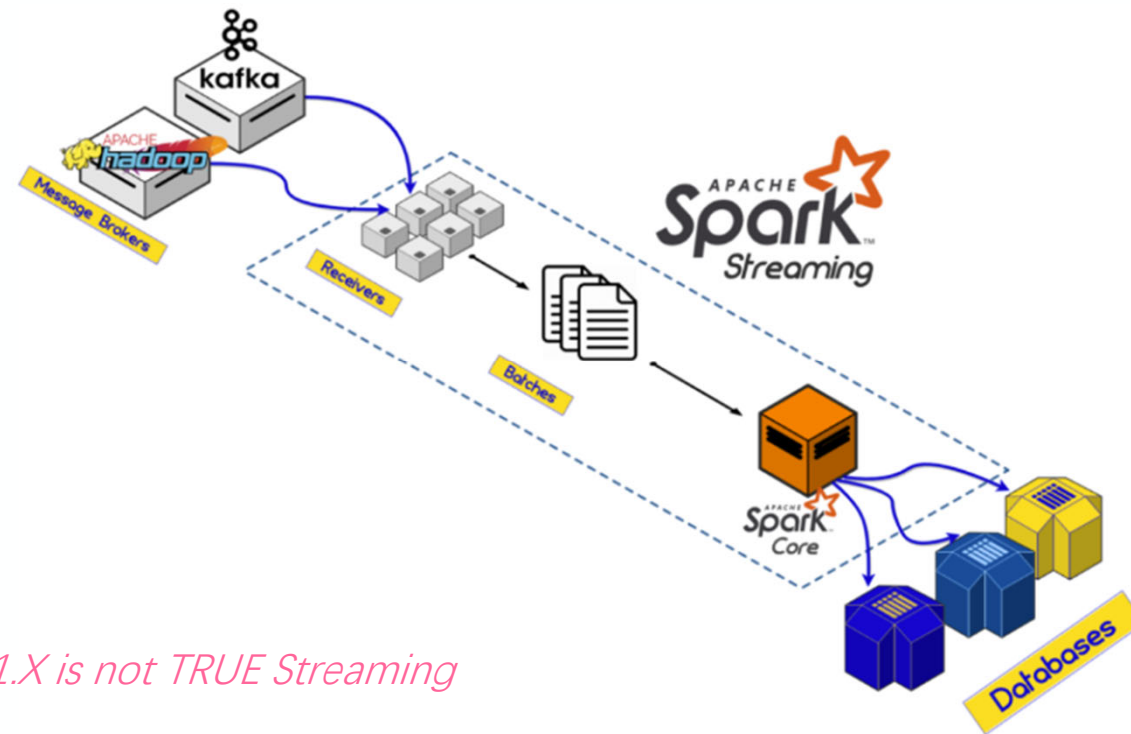groupByKey

join with inputs not co-partitioned

Figure 4: Examples of narrow and wide dependencies. Each box is an RDD, with partitions shown as shaded rectangles.

No Scheduler Fault-Tolerance
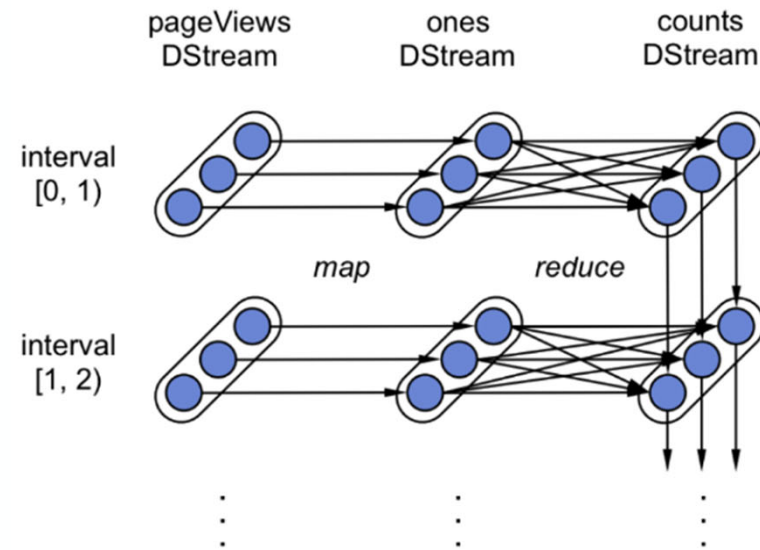
*Resilient Distributed Datasets:*
*https://www.usenix.org/system/files/conference/nsdi12/nsdi12-final138.pdf*

# Spark Streaming



*Spark Streaming 1.X is not TRUE Streaming*

# Discretized Streams

```
pageViews = readStream("http:// ... ", "1s")
ones = pageViews.map(event ⇒ (event.url, 1))
counts = ones.runningReduce((a, b) ⇒ a + b)
```
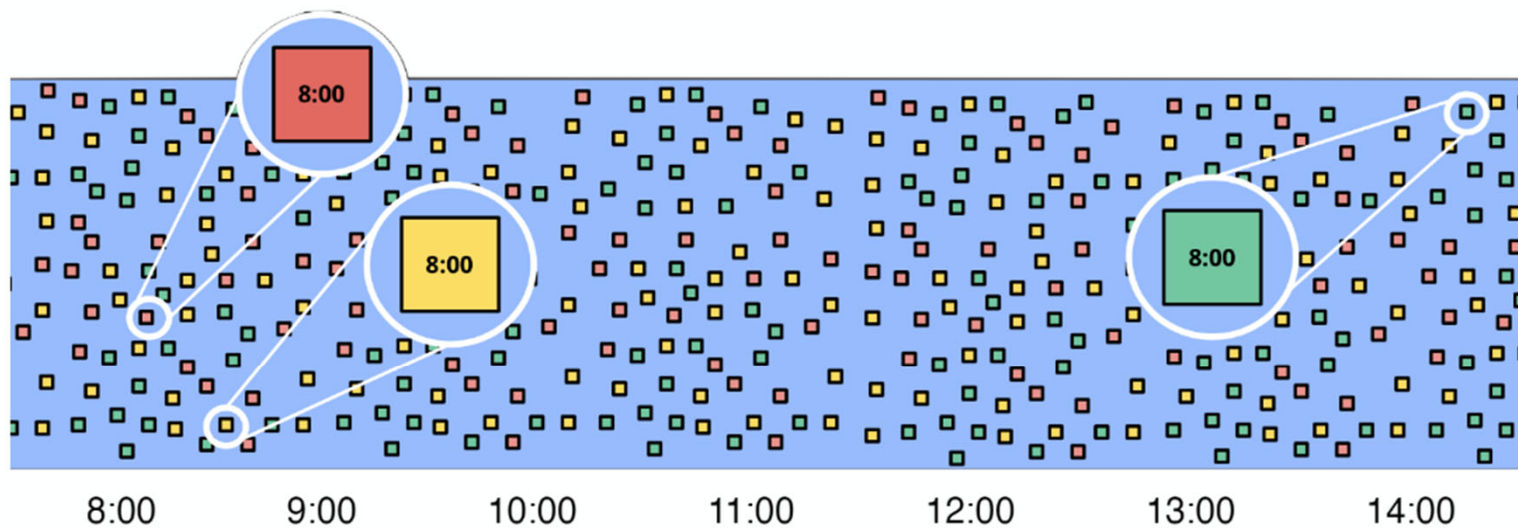
# The Dataflow Model

## A Practical Approach to Balancing Correctness, Latency, and Cost in Massive- scale, Unbounded, Out-of-order Data Processing
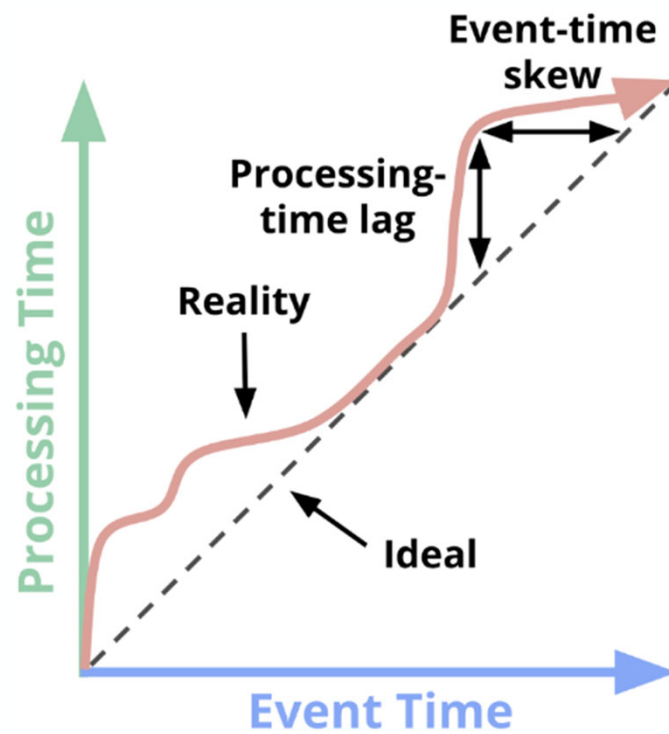
# Data Distribution Of Streaming



Data can be infinitely big with unknown delays.

# Time Distribution Of Events

# Focus

What are you computing?
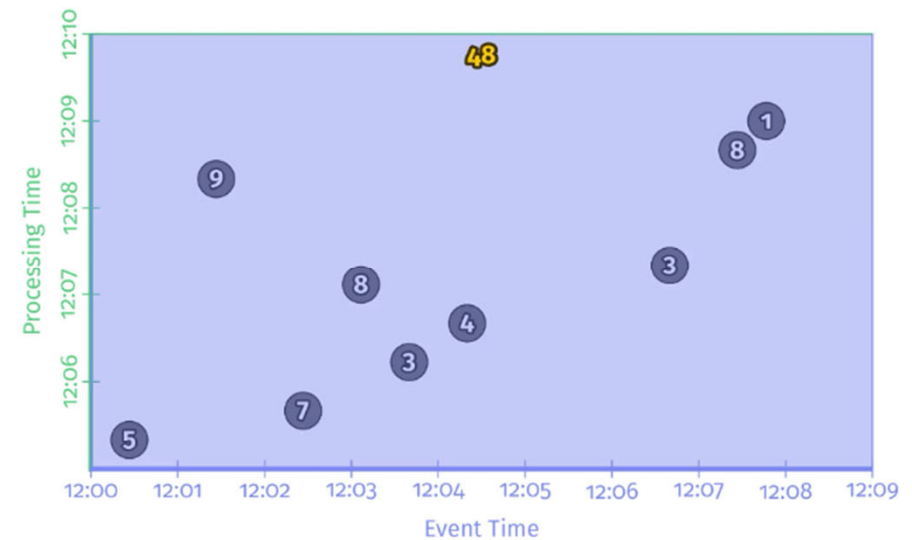
Where in event time?

When in processing time?

How do refinements relate?

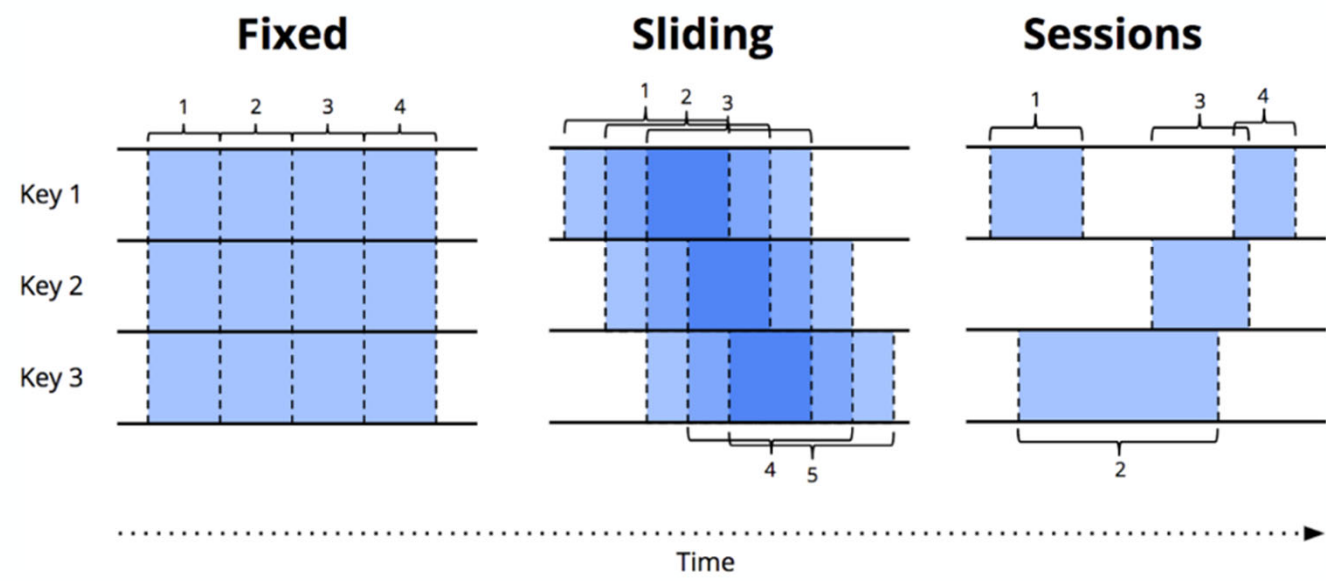# What Are You Computing

```
PCollection<KV<String, Integer>> input = IO.read( ... );
PCollection<KV<String, Integer>> output = input
    .apply(Sum.integersPerKey());
```
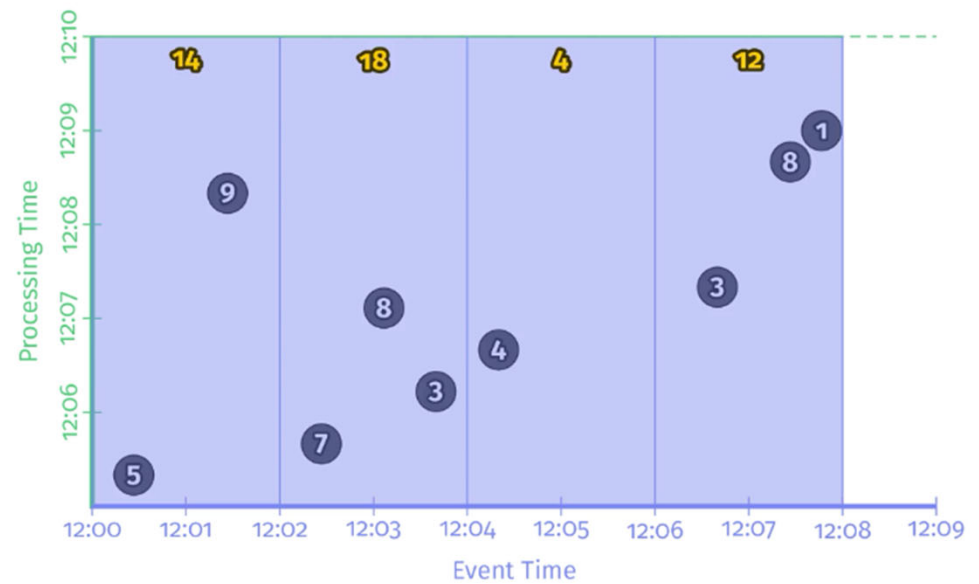
# Where in Event Time

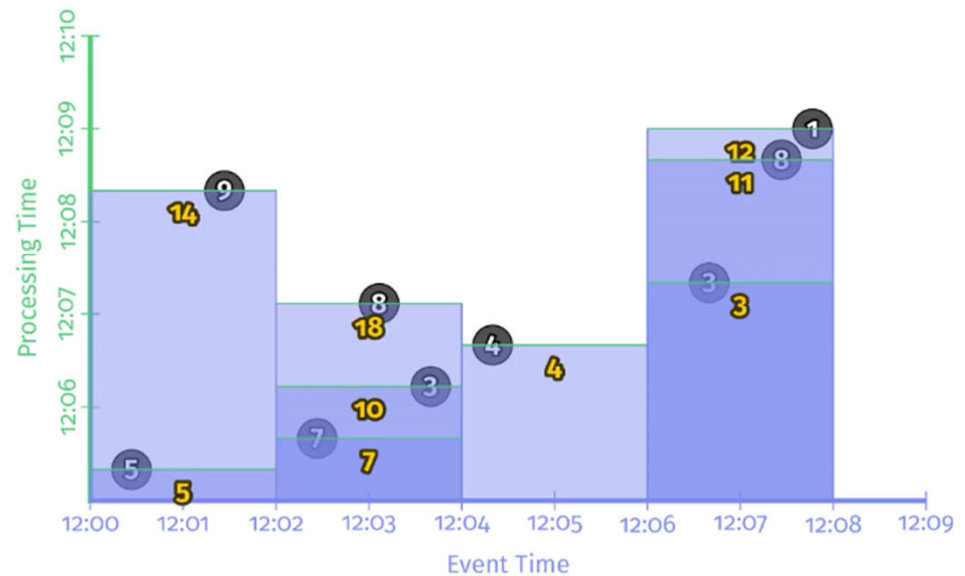Windowing

# Where in Event Time

```
PCollection<KV<Team, Integer>> scores = input
  .apply(Window.into(FixedWindows.of(TWO_MINUTES)))
  .apply(Sum.integersPerKey());
```
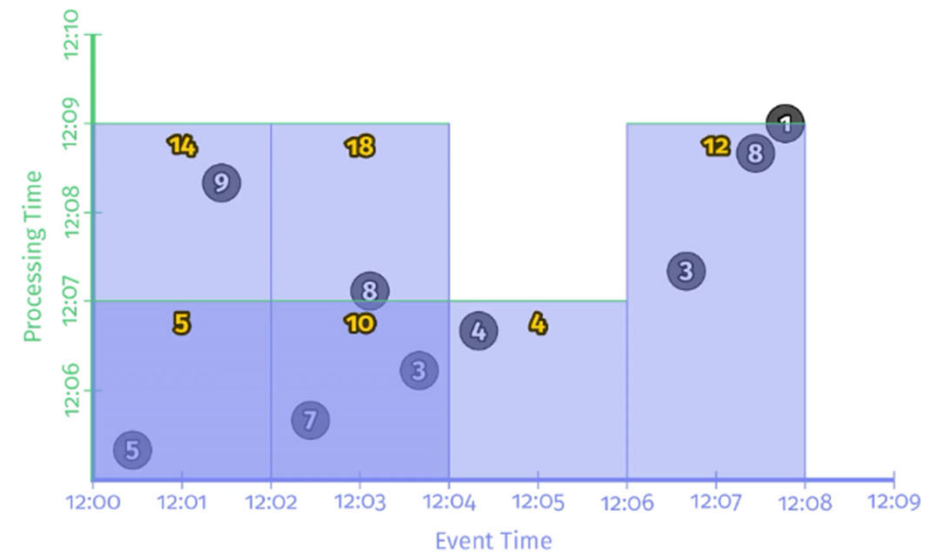
# When in Processing Time

Repeated Update Triggers

```
PCollection<KV<Team, Integer>> scores = input
  .apply(Window.into(FixedWindows.of(TWO_MINUTES))
              .triggering(Repeatedly(AfterCount(1))));
  .apply(Sum.integersPerKey());
```
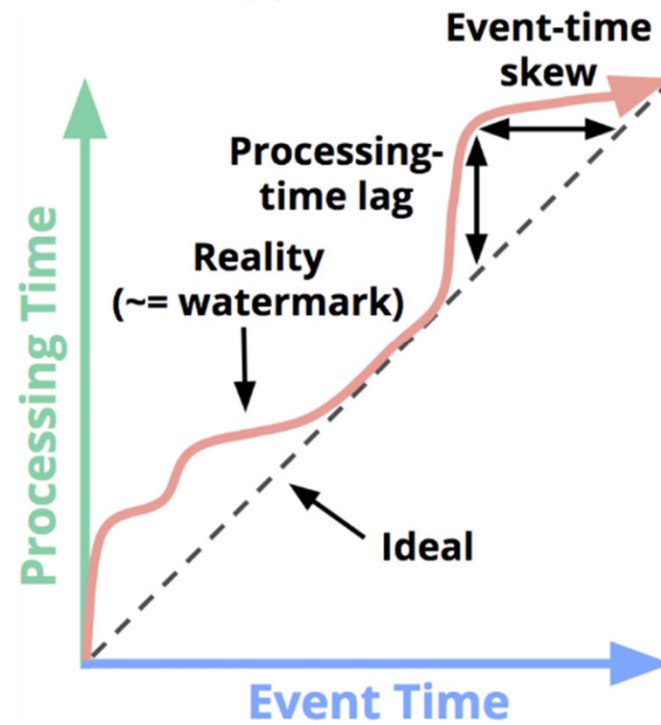
# When in Processing Time

Repeated Update Triggers

```
PCollection<KV<Team, Integer>> scores = input
  .apply(Window.into(FixedWindows.of(TWO_MINUTES))
              .triggering(Repeatedly(AlignedDelay(TWO_MINUTES))))
  .apply(Sum.integersPerKey());
```
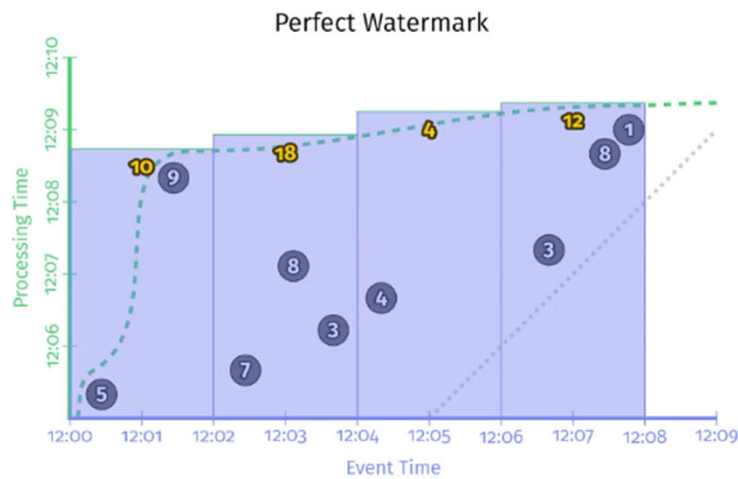
# When in Processing Time

Watermarks — Completeness Triggers
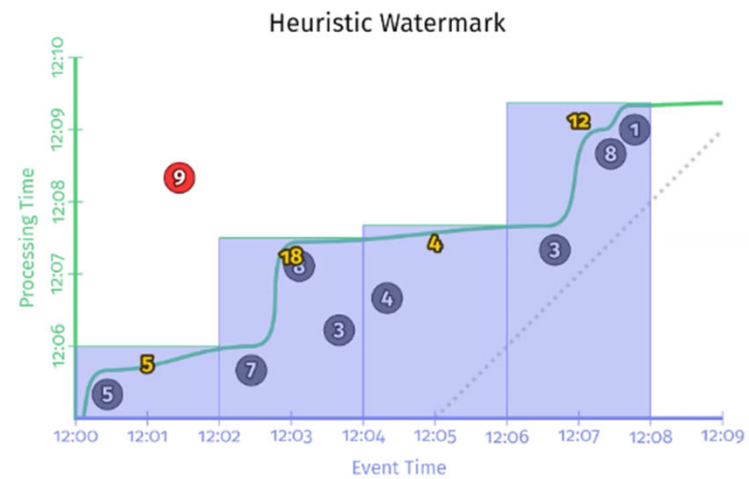
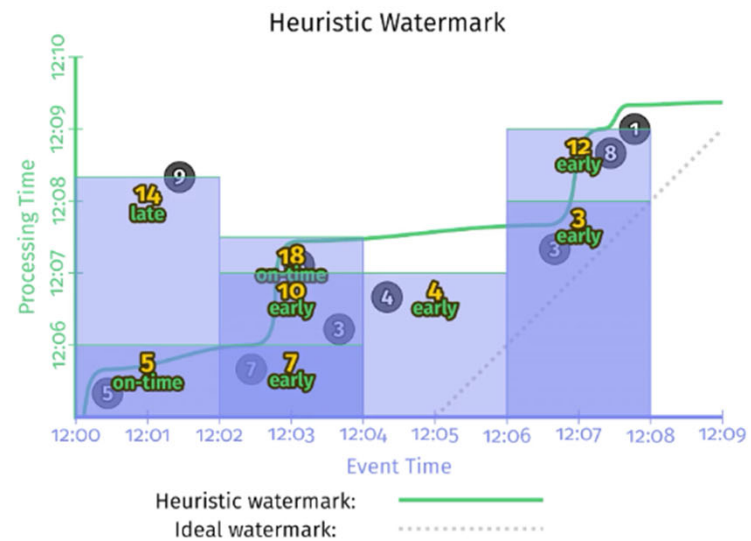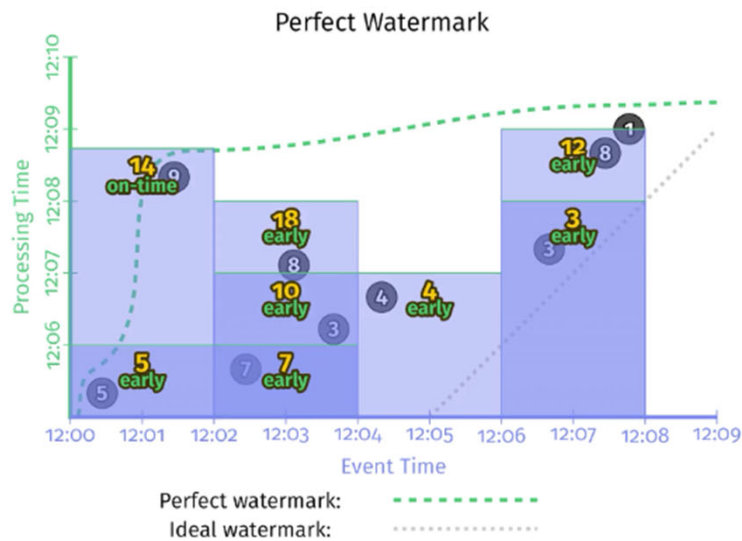# When in Processing Time



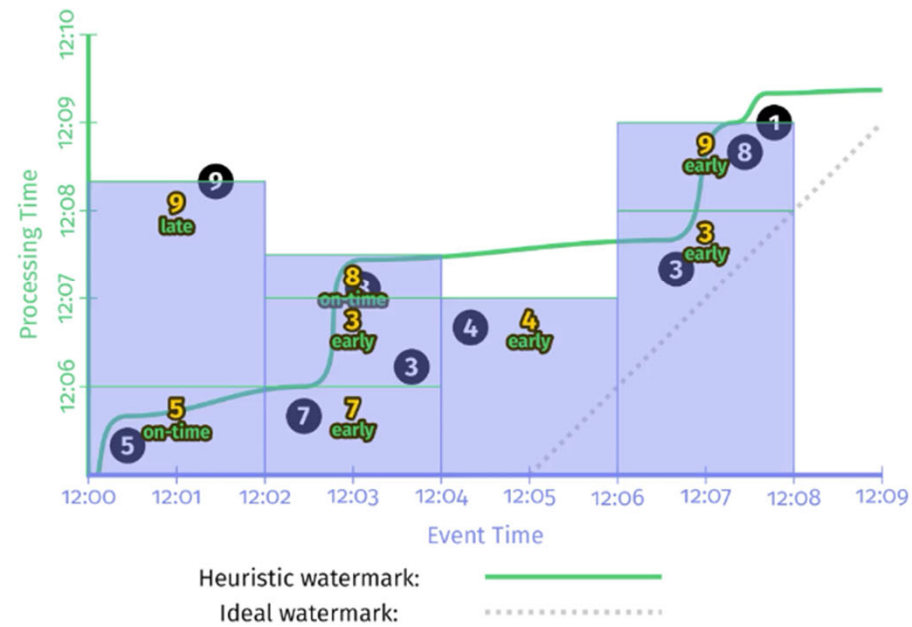Watermarks — Completeness Triggers

# When in Processing Time



Early/On-time/Late Trigger

# How do refinements relate

Discarding

# How do refinements relate



Accumulating & Retracting

# The Dataflow Model Summary

What are you computing?
- Pipeline Code

Where in event time?
- Windowing

When in processing time?
- Triggers & Watermark

How do refinements relate?
- Discarding, Accumulating and Accumulating&Retracting

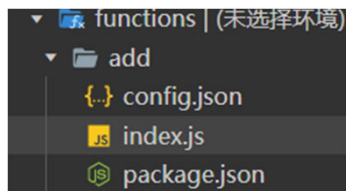# Serveless - FaaS



AWS Lambda

无需预置或管理服务器即可运行代码，您只需为实际使用的资源付费

Amazon S3

以行业领先的可扩展性、数据可用性、安全性和性能存储任意数量的数据

# Exp. Wechat Mini Program Serveless

# Exp. Serveless Function: add

```
functions | (未选择环境)
  add
    {..} config.json
    index.js
    package.json
```
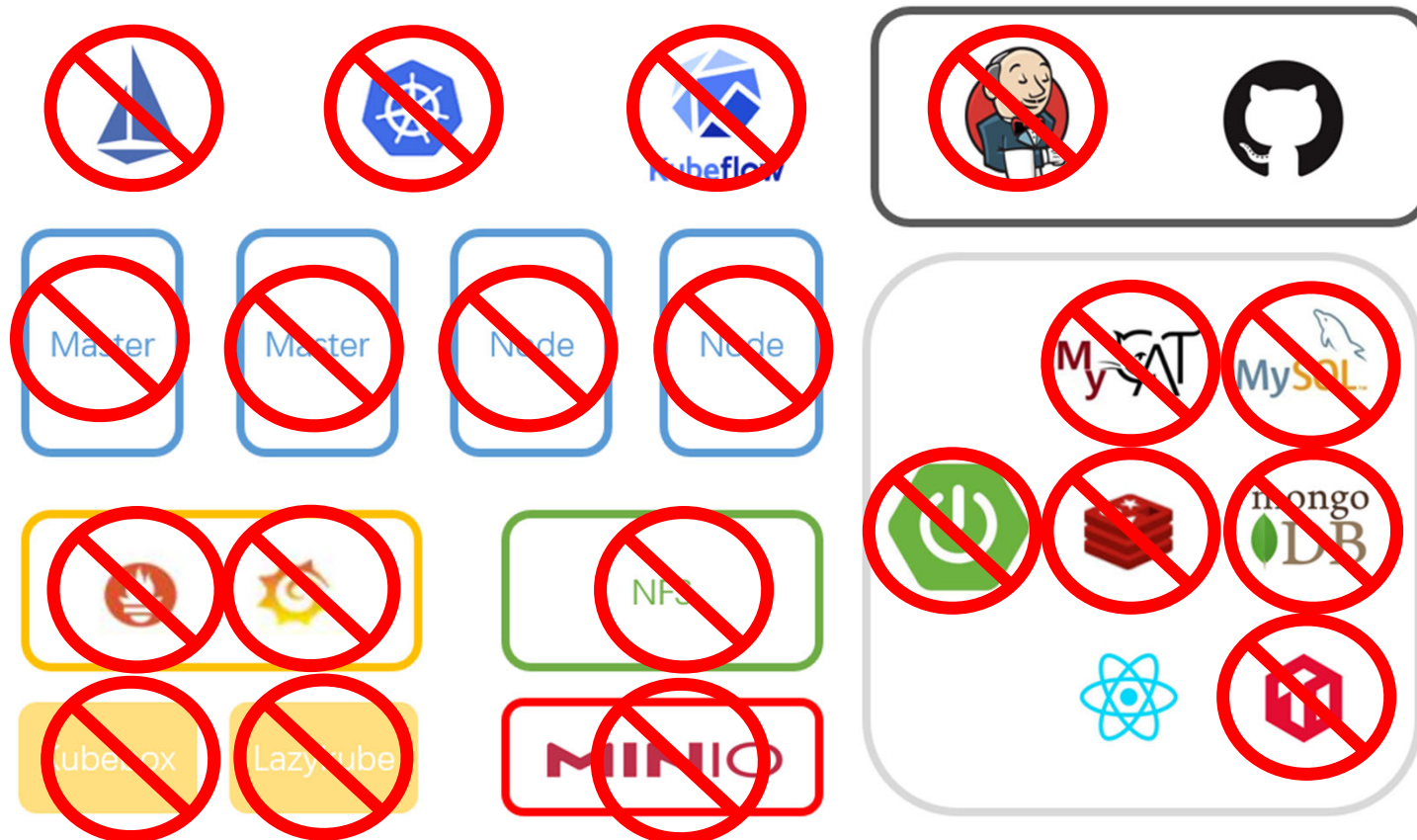
```javascript
// Cloud Function Entrance File
const cloud = require('wx-server-sdk')
...

cloud.init()


// Cloud Function Entrance
exports.main = async (event, context) => {
  const wxContext = cloud.getWXContext()

  return {
    sum: event.a + event.b,
    event,
    openid: wxContext.OPENID,
    appid: wxContext.APPID,
    unionid: wxContext.UNIONID,
  }
}
```

# The Pro of Serveless From A Project Exp.

# The Challenge of Serveless





F1 → F2

F2 → F3

F2 → F4



Latency