



Detecting Change-Points in Apache Beam

By Devon Peticolas - Oden Technologies

<https://github.com/x/slides/tree/master/beam-summit-2022>



BEAM
SUMMIT

Austin, 2022



Devon Peticolas

Principal Engineer @
Oden Technologies
“Beam Guy”



In This Talk



- What is Change-Point Detection?
- Why and how does Oden use Change-Point Detection to deliver features?
- Methods of doing Change-Point Detection in Beam.
- Methods of doing Change-Point Detection with Smoothing.
- Impacts of event sparsity, lateness, and order.

A little about Oden



BEAM
SUMMIT

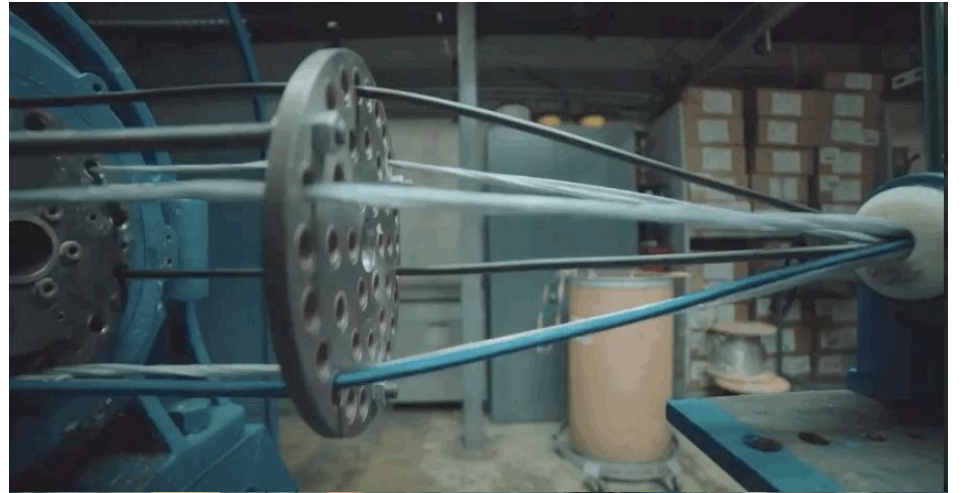
Austin, 2022

Oden's Customers

Medium to large manufacturers in plastics extrusion, injection molding, and pipes, chemical, paper and pulp.

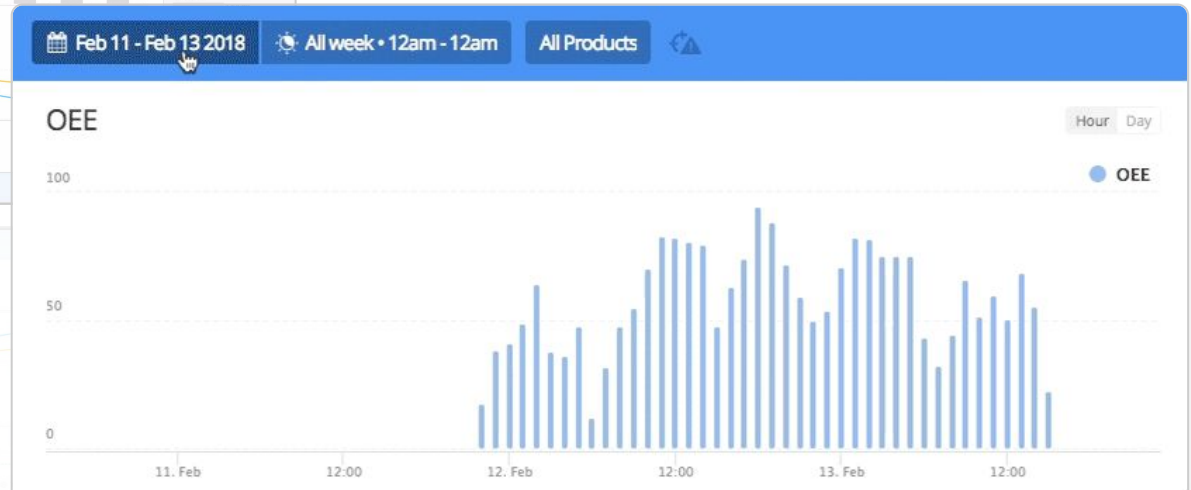
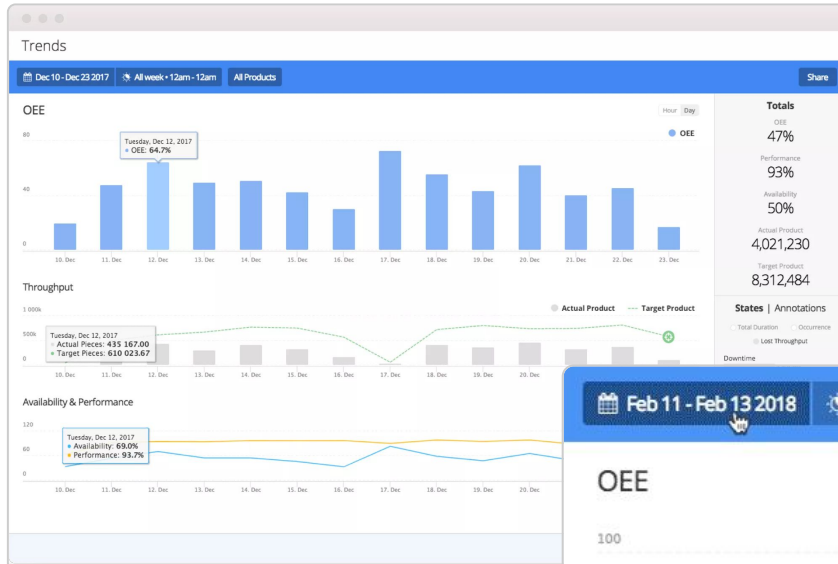
Process and Quality Engineers looking to centralize, analyze, and act on their data.

Plant managers who are looking to optimize logistics, output, and cost.



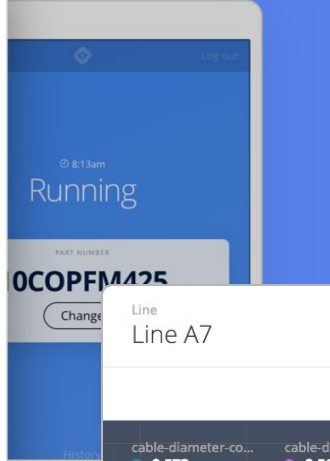
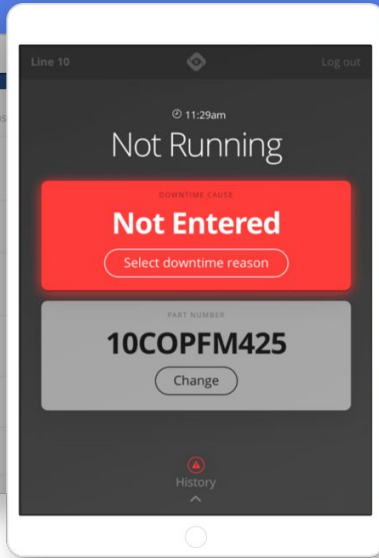
Interactive Time-series Analysis

- Compare performance across different equipment.
- Visualize hourly uptime and key custom metrics.
- Calculations for analyzing and optimizing factory performance.



Real Time Manufacturing Data

- Streaming second-by-second metrics
- Interactive app that prompts on production state changes and collects user input.



Background:

How Oden Uses Beam



BEAM
SUMMIT

Austin, 2022

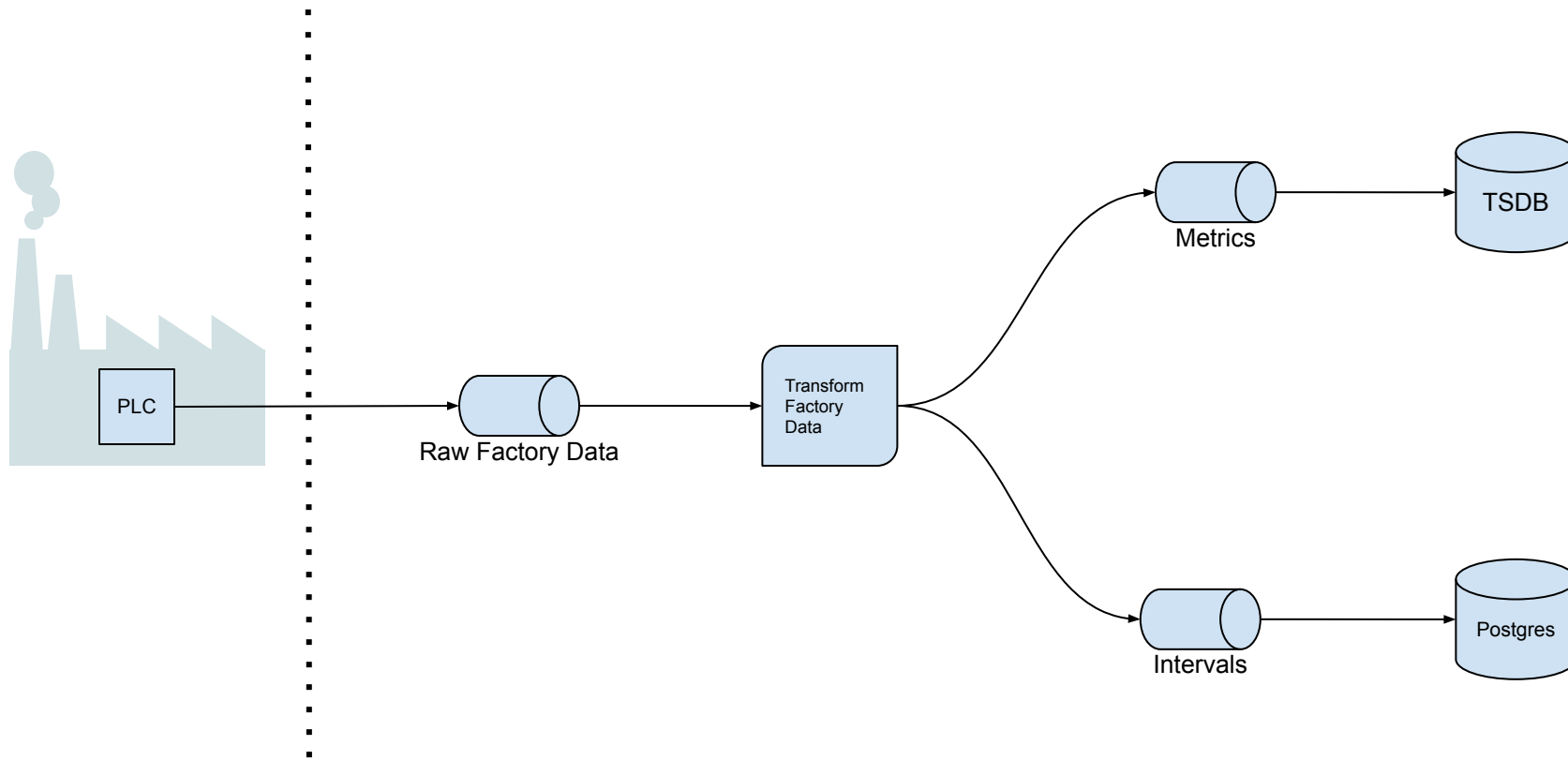
How Oden Uses Beam



- Ingesting “raw” manufacturing data and mapping it into Oden “events”
- Combining events using streaming joins
- Making customer-configured transformations to events
- Transforming metric events into contextual interval events

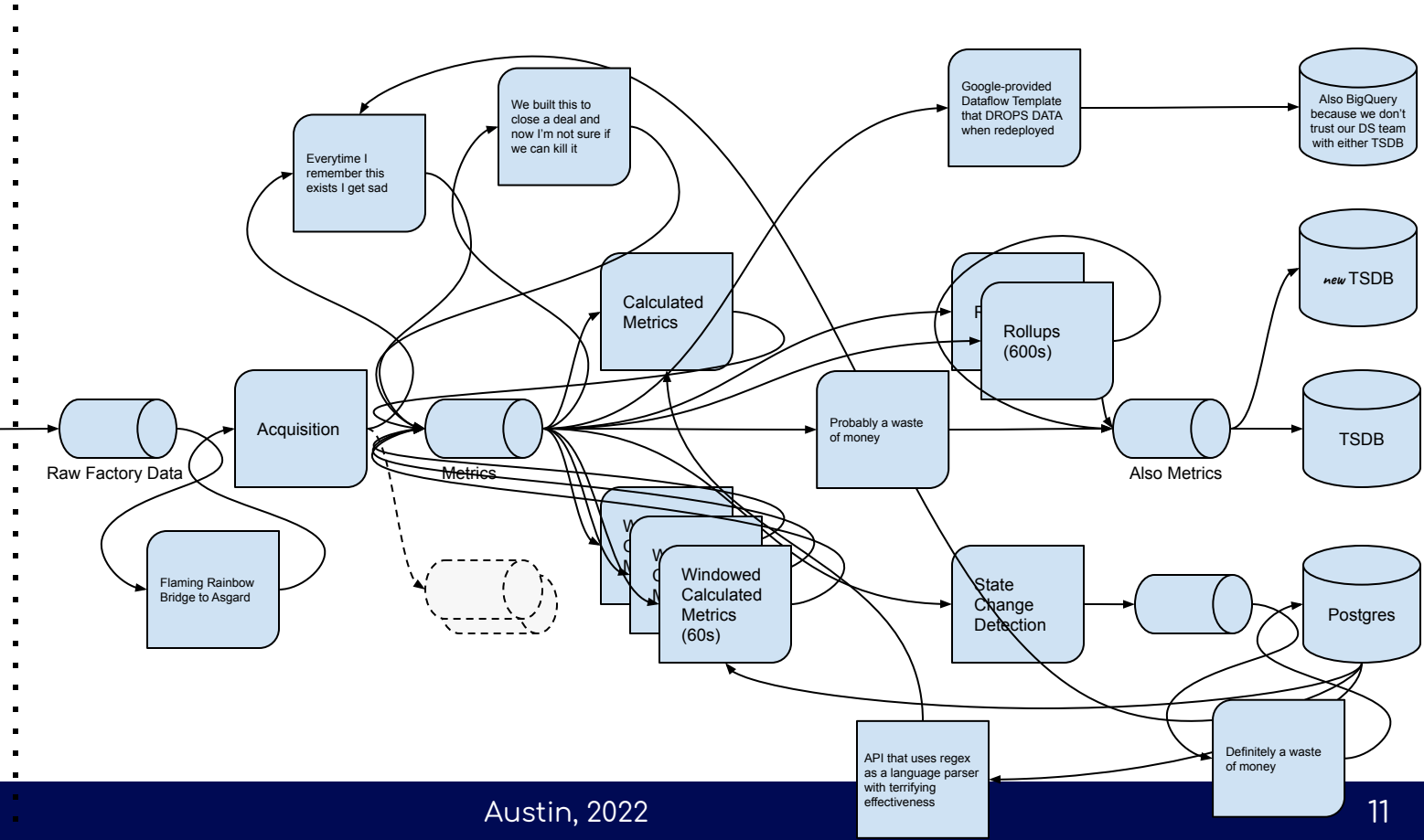
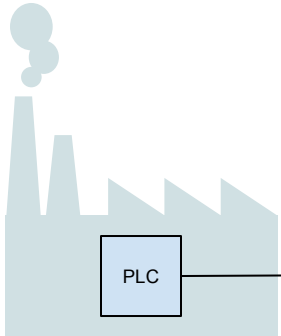
- * *Lots of Side-Input Joining*
- * *Lots of Complex Windowing*
- * *Lots of Performance Concerns*

Streaming Factory Data - In Summary

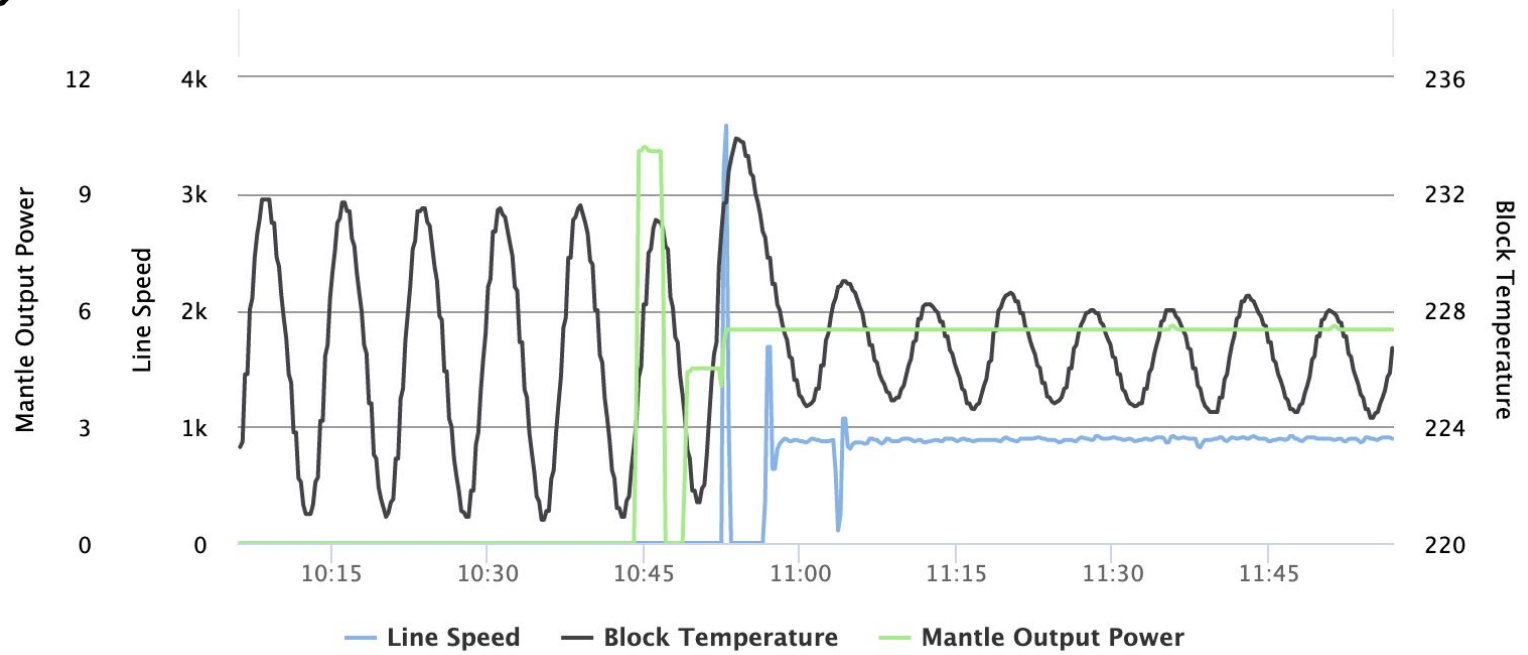
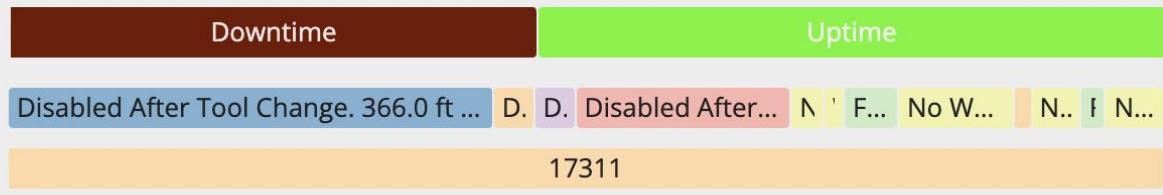


Streaming Factory Data - In Summary

Reality



Data is Grouped by Production Line

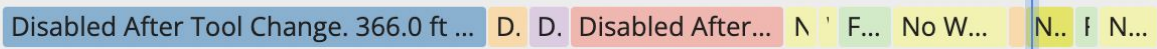


Process Context

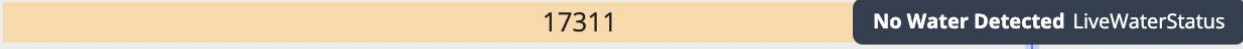
Availability



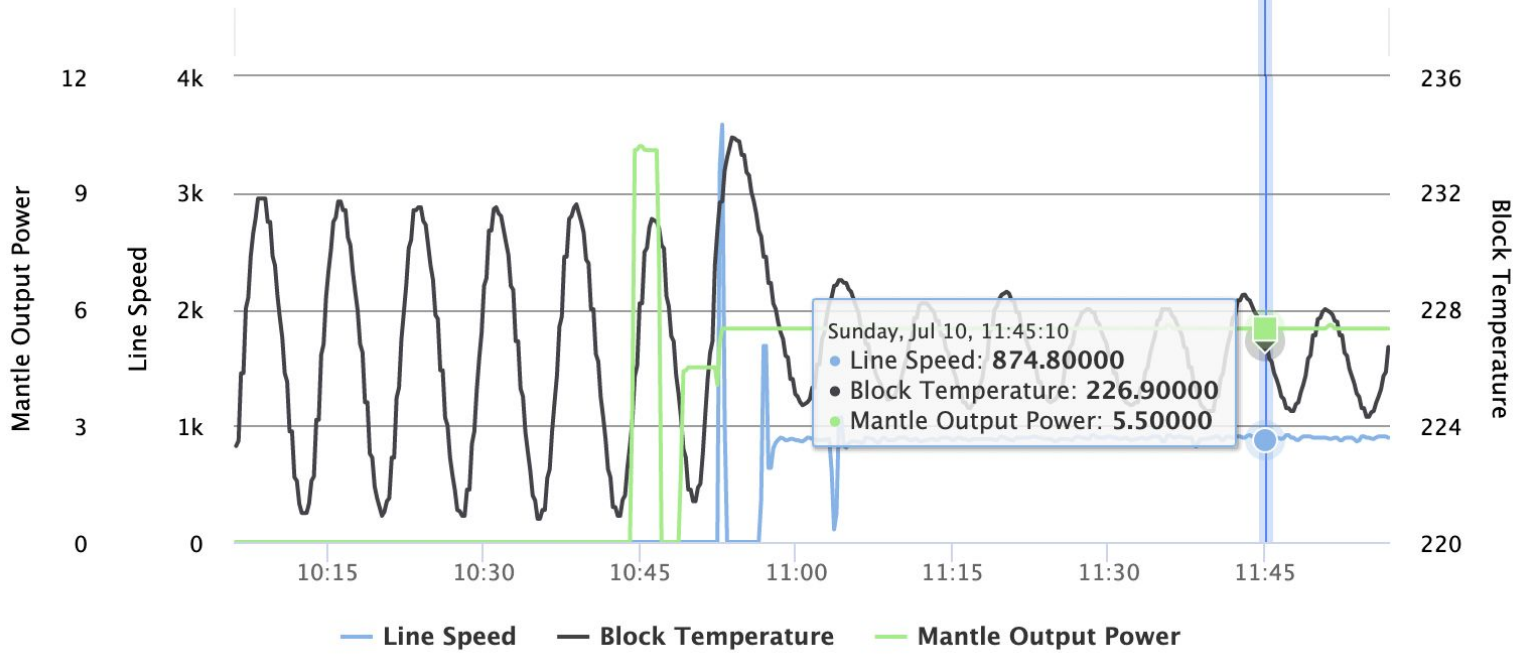
Water Status



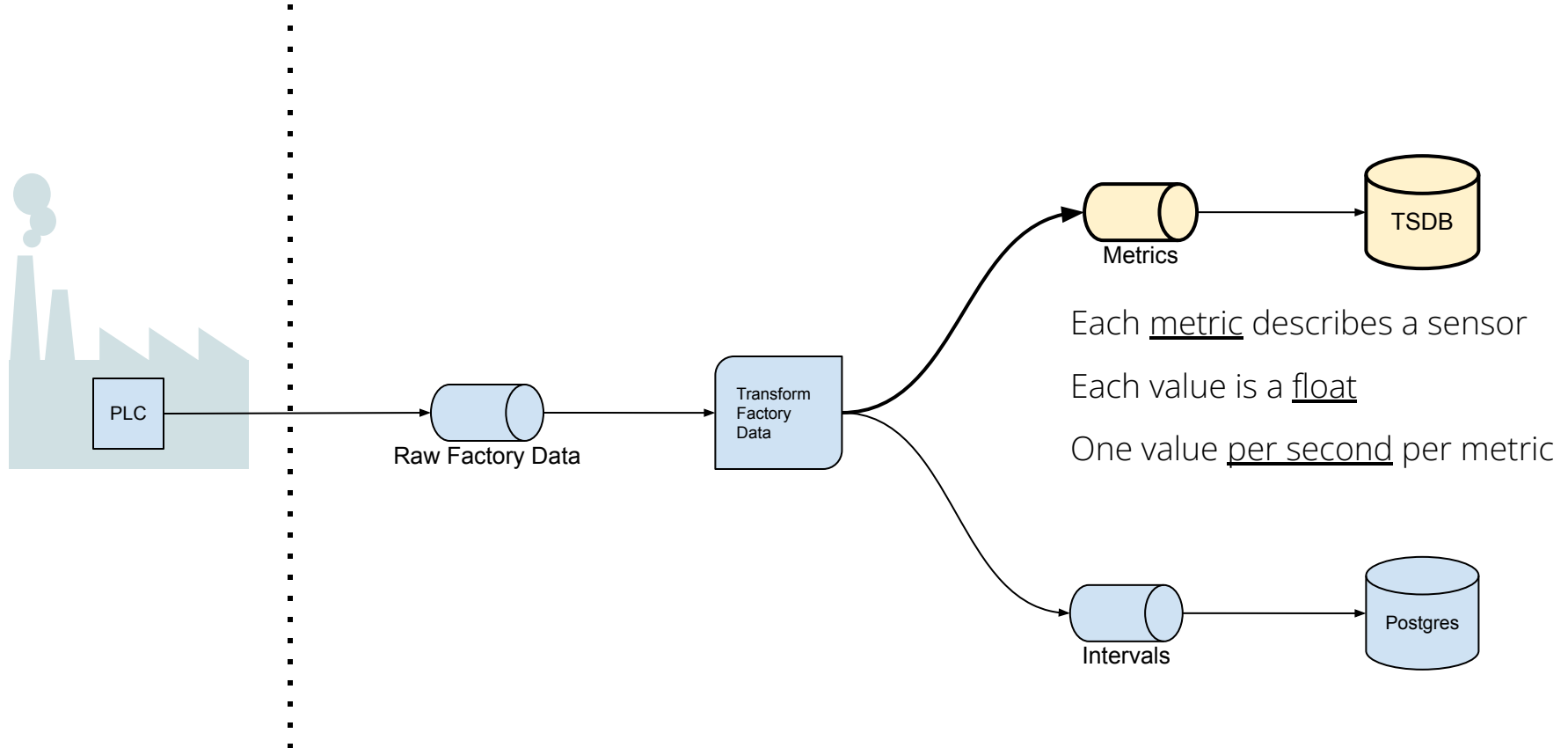
Operator Number



Metrics



Metrics

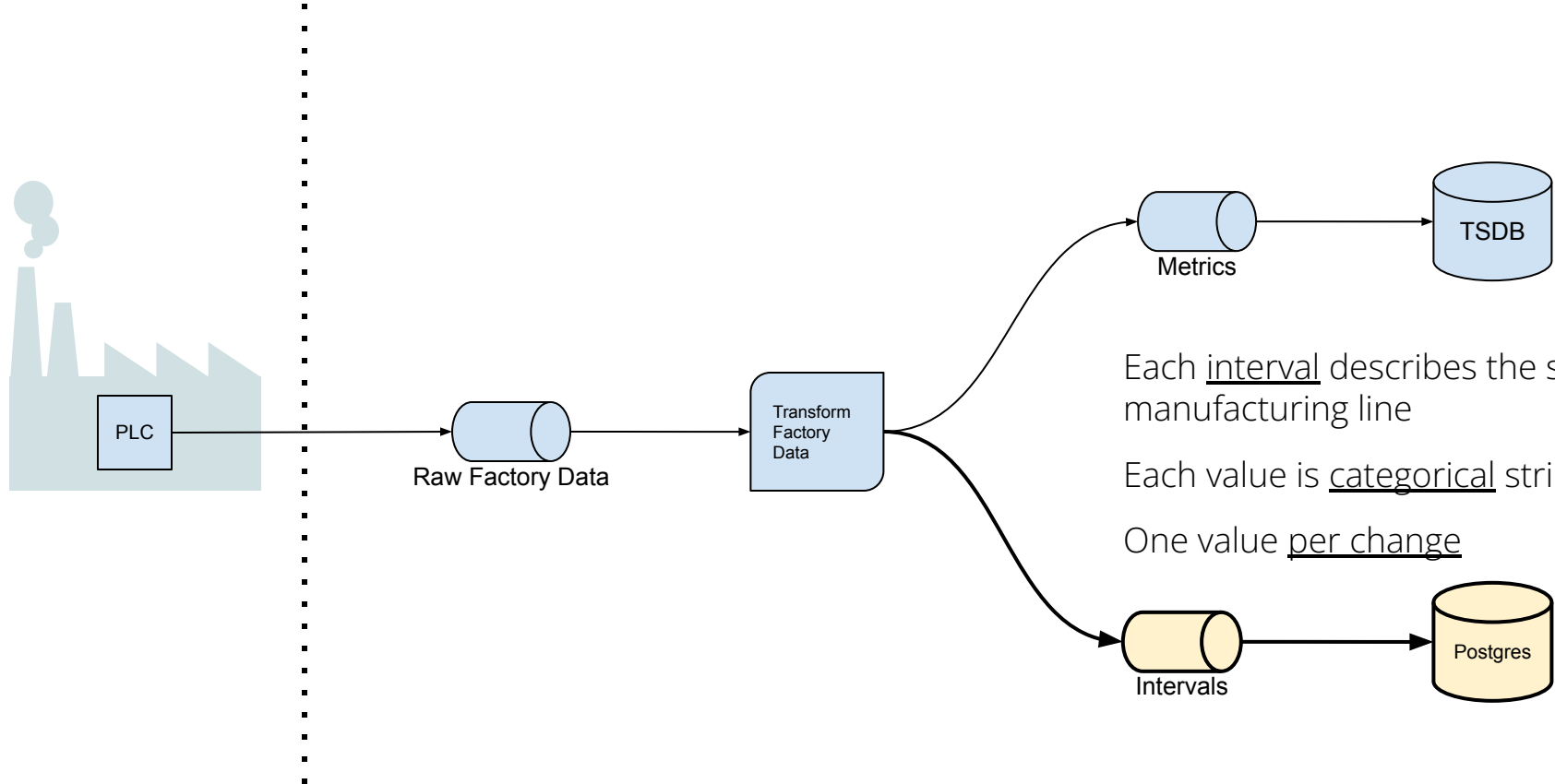


Each metric describes a sensor

Each value is a float

One value per second per metric

Intervals

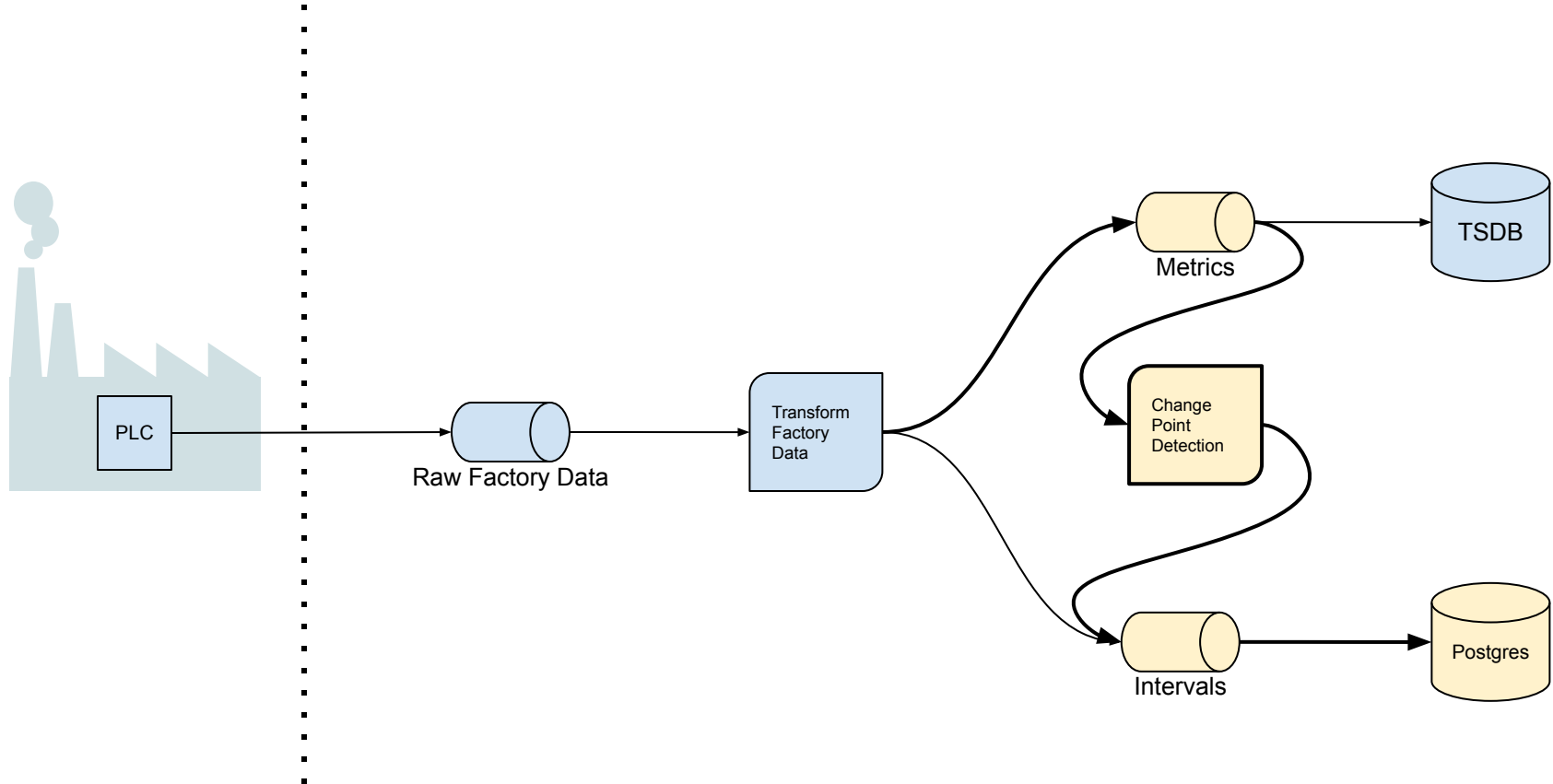


Each interval describes the state of a manufacturing line

Each value is categorical string

One value per change

Creating Intervals from Metrics



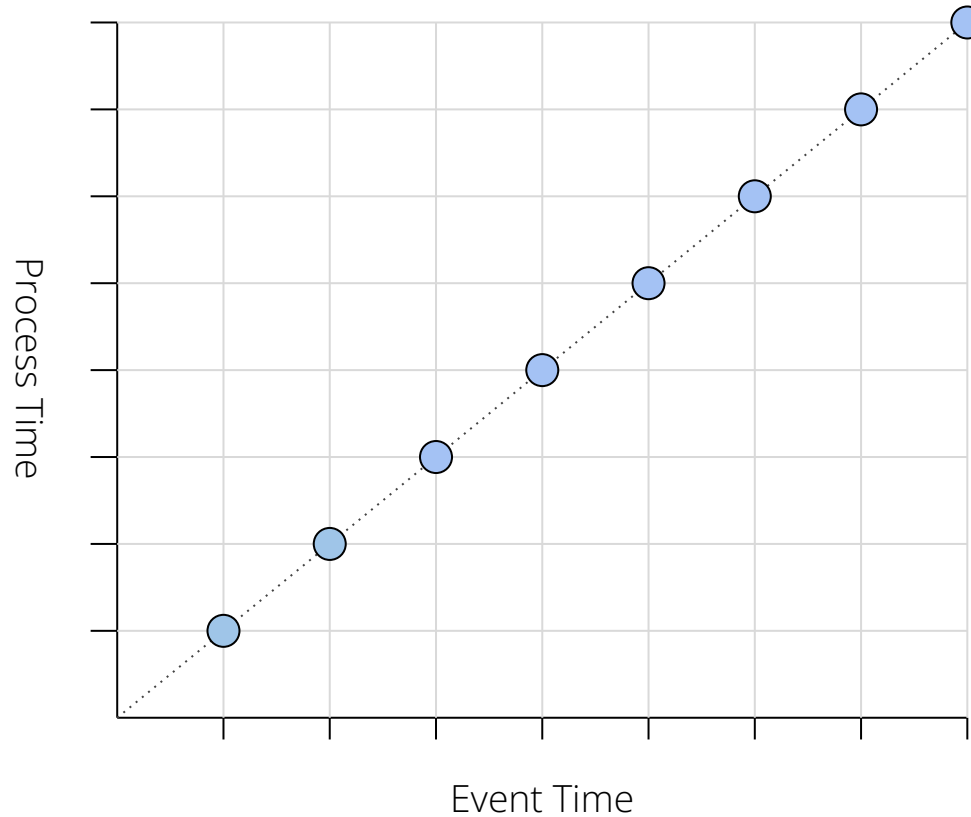
Use Case: Creating Intervals from Metrics



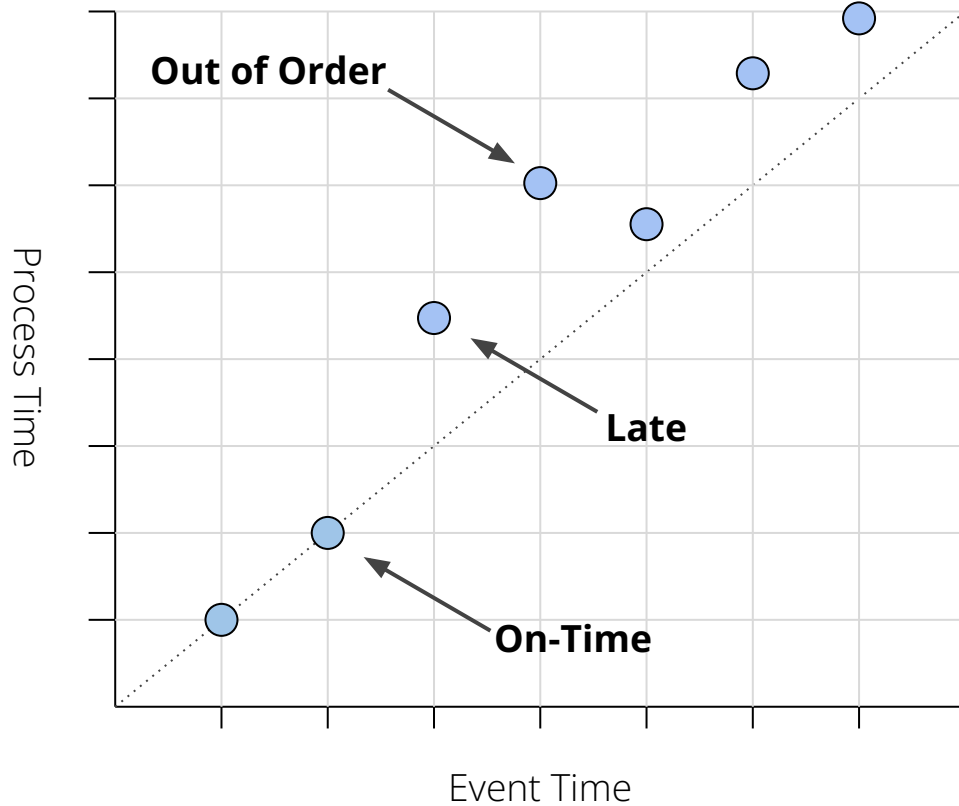
BEAM
SUMMIT

Austin, 2022

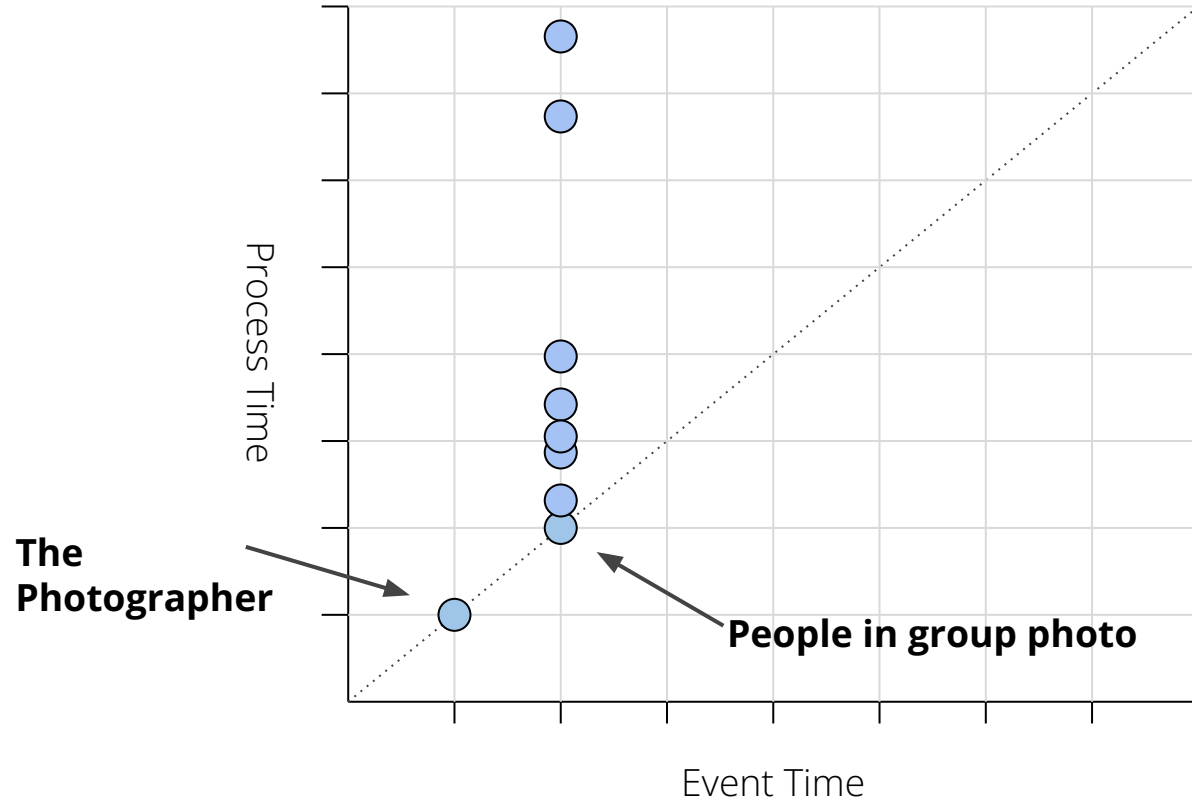
Visualizing Events



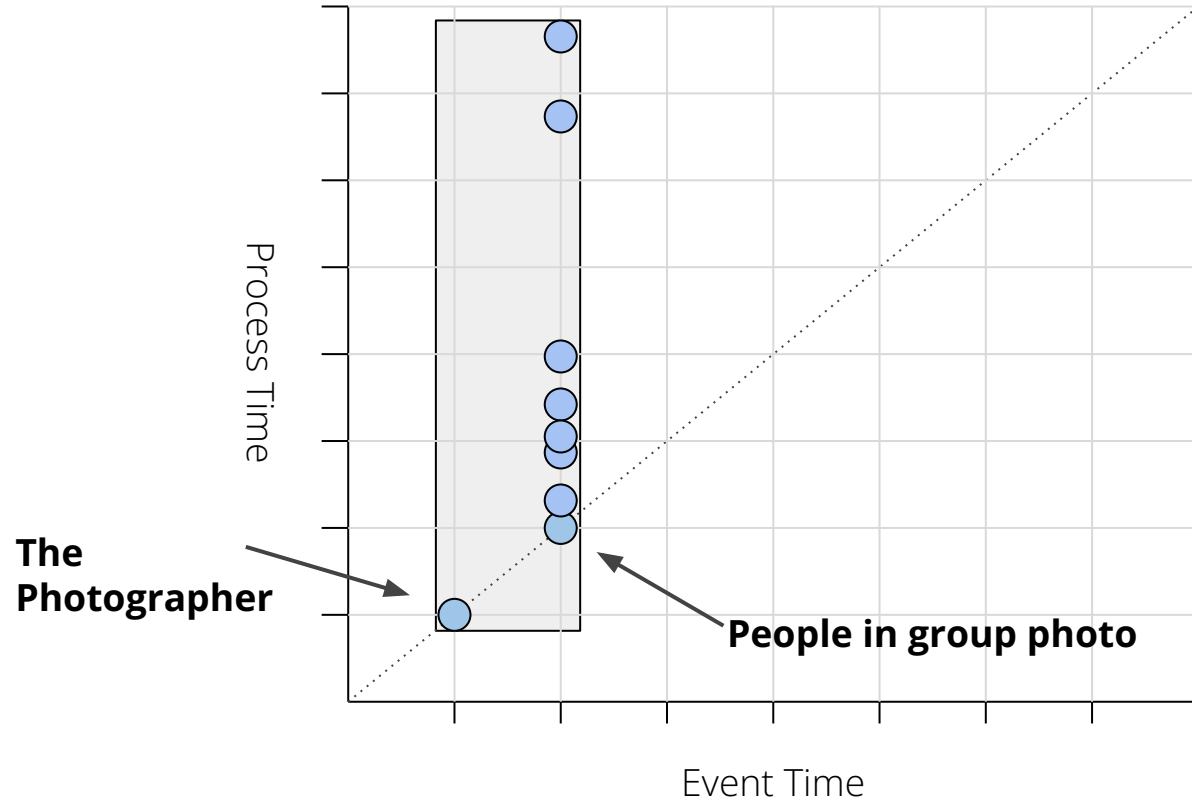
Visualizing Events



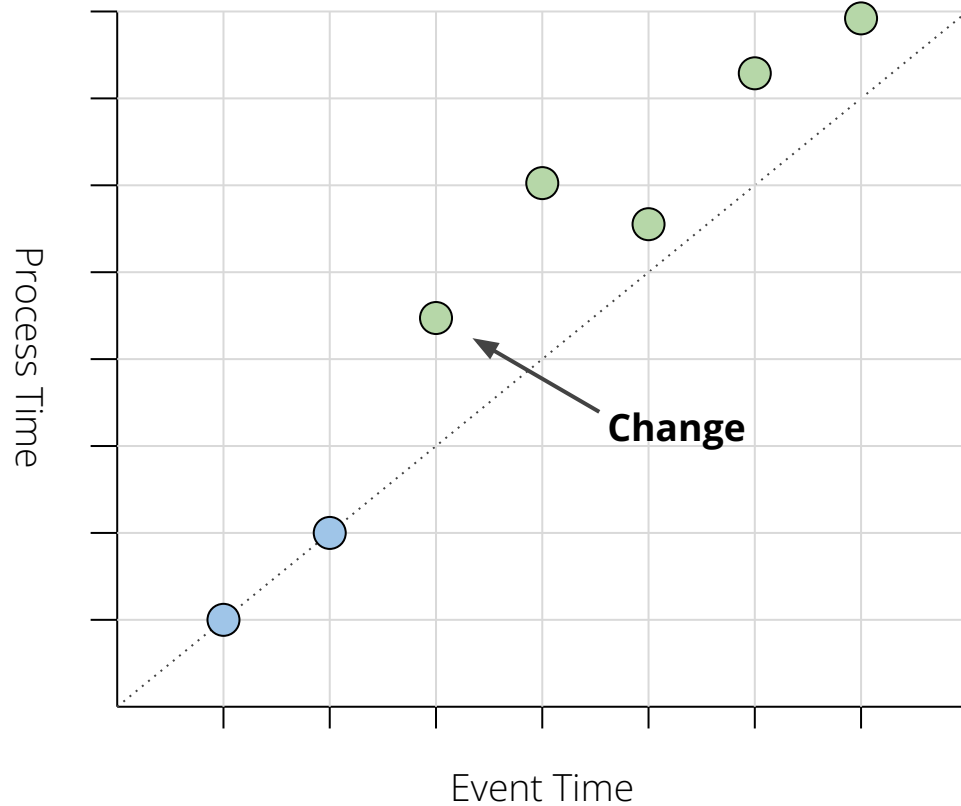
Visualizing Events

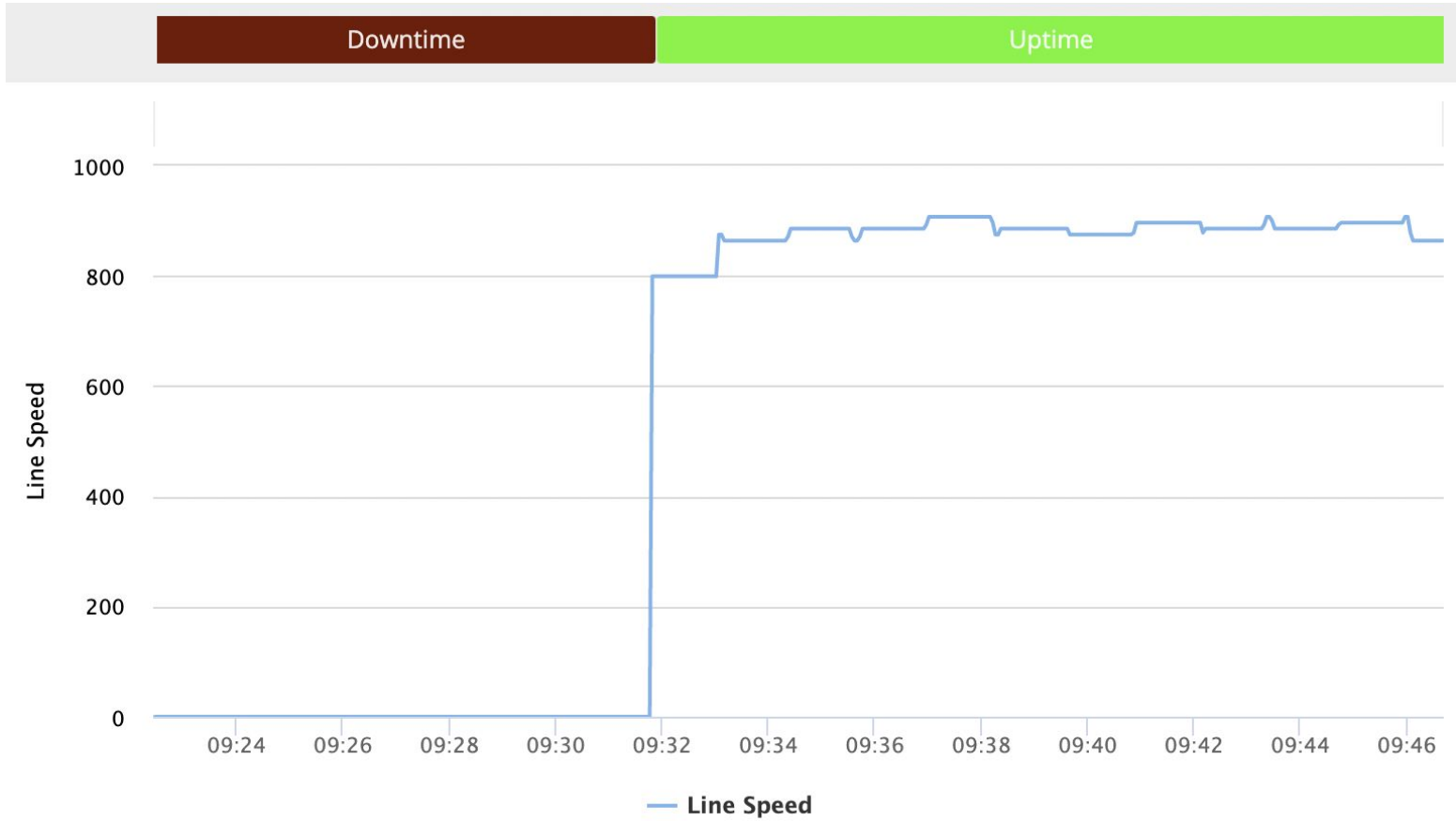


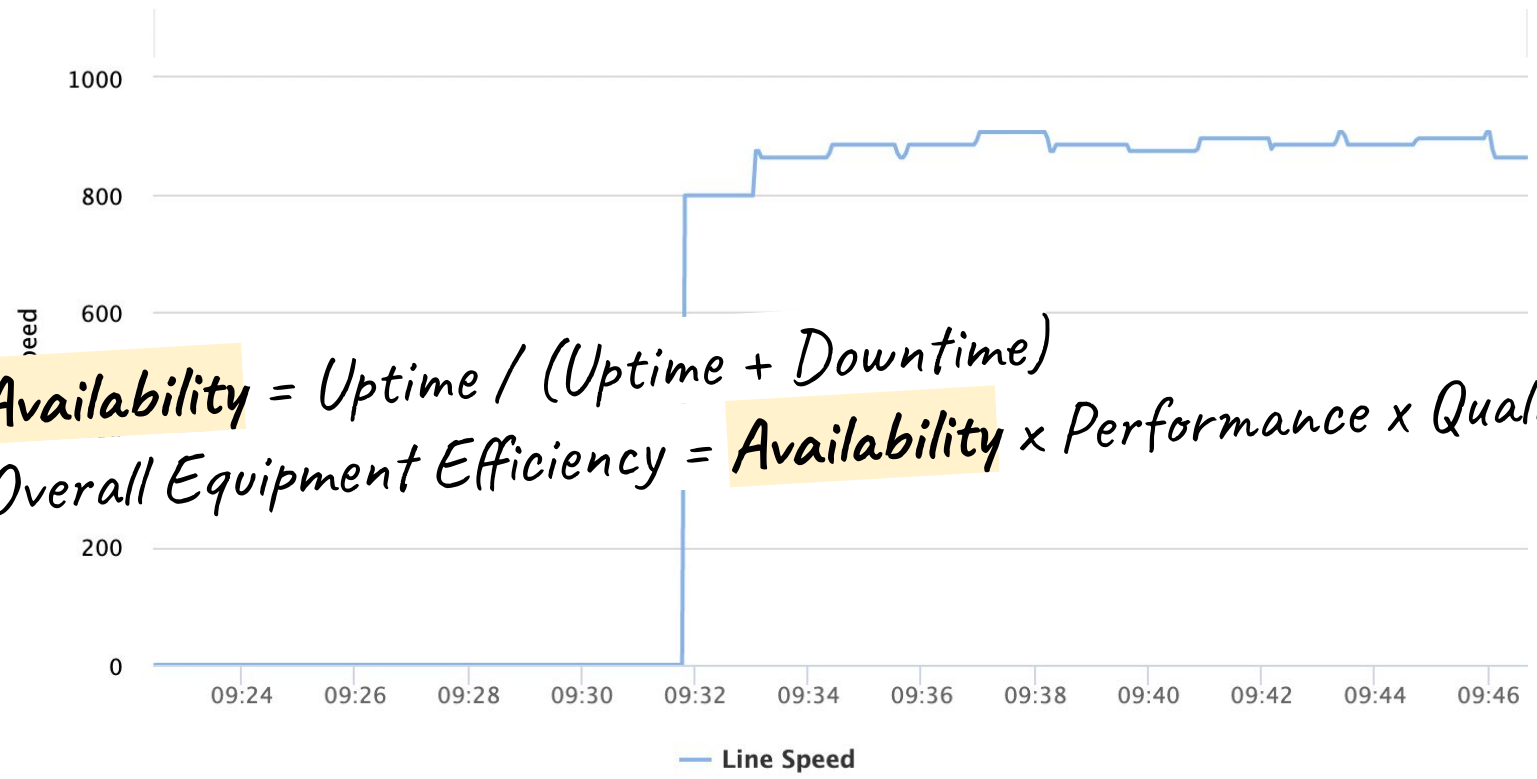
Visualizing Events



Visualizing Events w/ Change





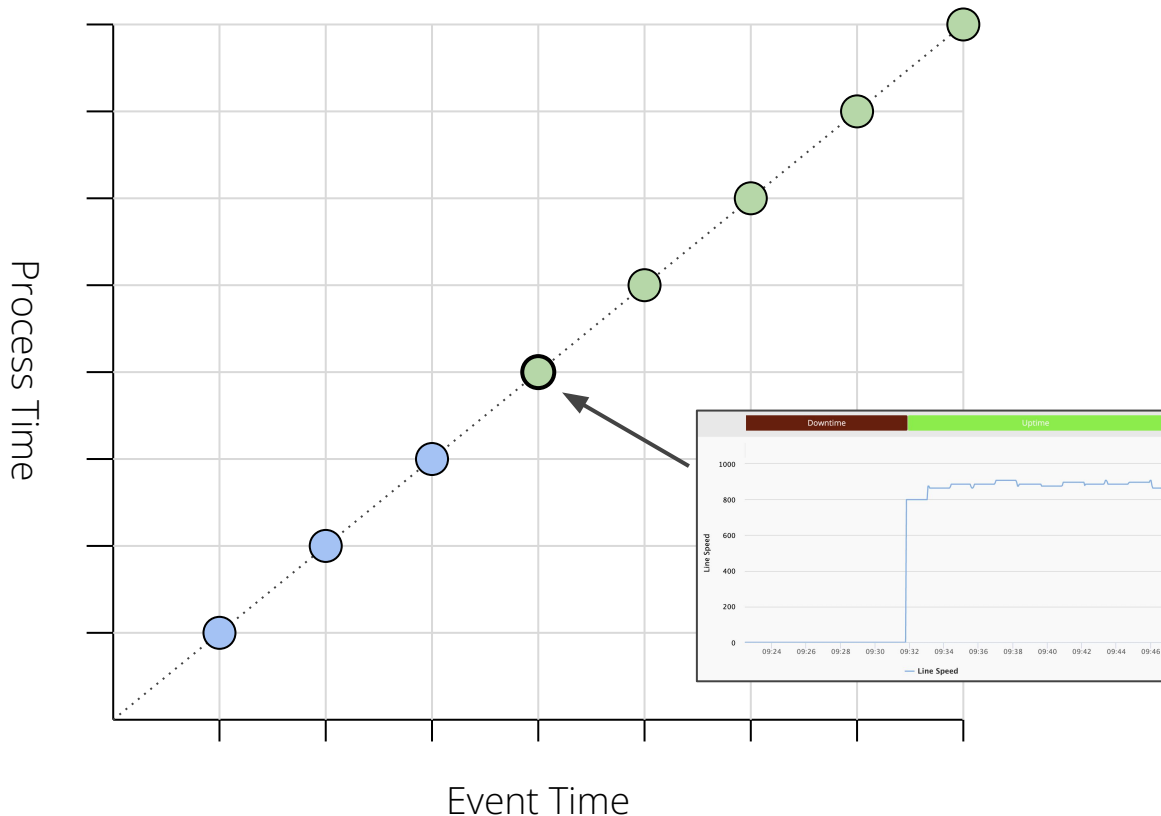


$Availability = Uptime / (Uptime + Downtime)$
 $Overall\ Equipment\ Efficiency = Availability \times Performance \times Quality$

Metrics into Categorical Values

```
// Load configuration every 5-minutes.
PCollectionView<Config> configView = p
    .apply(
        GenerateSequence
            .from(0)
            .withRate(
                1, Duration.standardMinutes(5)))
    .apply(MapElements(...)) // API call
    .apply(View.asSingleton());

// Map metric values to categorical
// values using config side-input.
p.apply(ParDo
    .of(new DoFn<Metric, String>() {
        public void processElement(
            Metric m, ProcessContext c
        ) {
            Config config = c.sideInput(configView)
            if (m.value > config.forMetric(m)) {
                c.output("up");
            } else {
                c.output("down");
            }
        }
    })
    .withSideInputs(configView))
```



Factory A



Last 2 hours ▾

5 lines ▾

Print

List ▾

Line	LINE AGGREGATE State History	Util	Perf	Production	CURRENT RUN ON LINE Work Order	Current Speed	Current Production
Compounding 1 >	Uptime 4h 50m	100%	0%	15,000 lb	UHK-71120000 25h 13m Product 6D50D9B2 Batch 3C64E95D	7,500	174,222 lb
Extruder 1 >	Downtime 2m 55s	81%	181%	36,407 ft	Test Run 347d 1h Product 109810-P-A	0 207 +3.56 / -17 ⚙	12,385,545 ft
Extruder 2 >	Uptime 49m 2s	95%	79%	5,066 ft	XYZ-123 656d 9h Product 85EB50E1 Batch 22CD2E30	50 56 +4 / -6 ⚙	1,753,394 ft
Extruder 3 >	Uptime 3m 44s	63%	78%	12,468 ft	UHL-64957000 2h 56m Product D48199B2	183 210 ± 15 ⚙	17,830 ft
Line 01 >	Uptime 53m 33s	97%	0%	0 ft	line_demo_in_PROD 53d 12h Product 2AFBCA5E Batch D824D563	330 320 +10 / -20 ⚙	0 ft

OEE

0%

Utilization

87.3%

Performance

80.6%

Quality

0%

Solution: Using Beam State



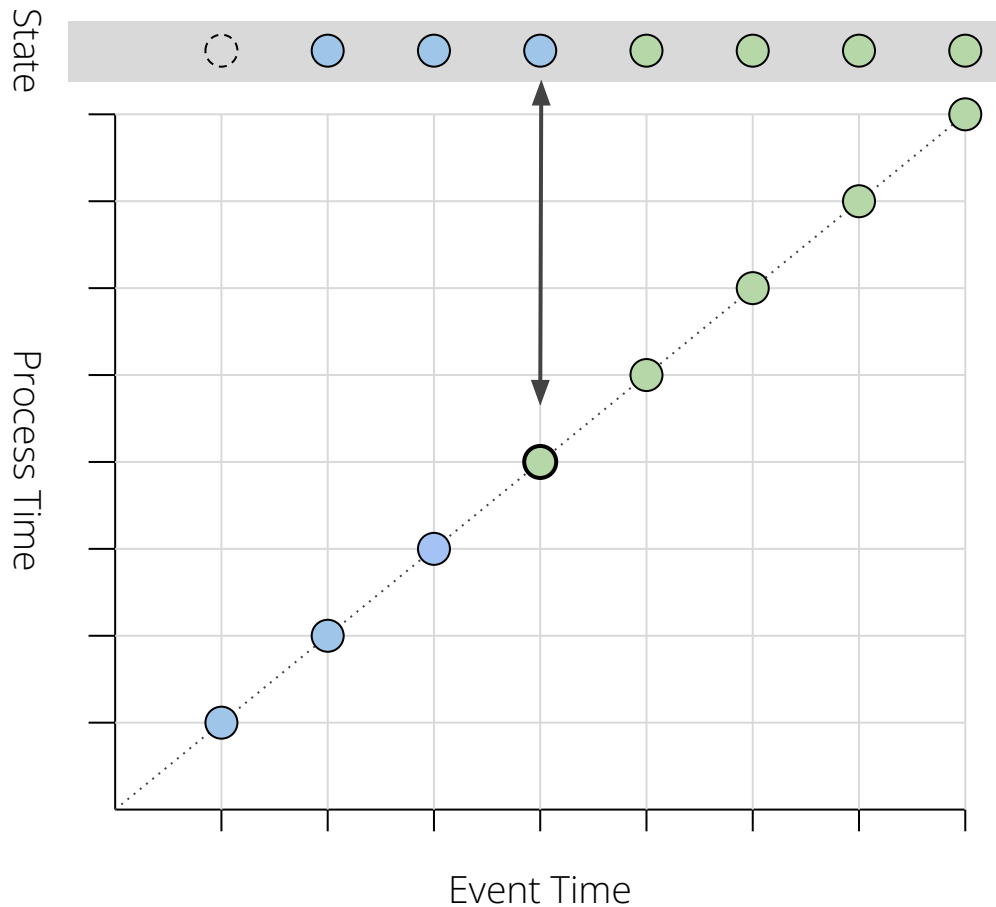
BEAM
SUMMIT

Austin, 2022

Beam State to Detect Change-Points

```
DoFn<KV<String, T>, T> {
  StateSpec<ValueSpec<T>>
  prevSpec =
    StateSpecs.value(...);

  public void processElement(
    ProcessContext c,
    ValueState<T> prev) {
    T curr =
c.element().getValue();
    T last = prev.read();
    if (curr != last) {
      c.output(curr);
    }
  }
}
```



Issues: Using Beam State



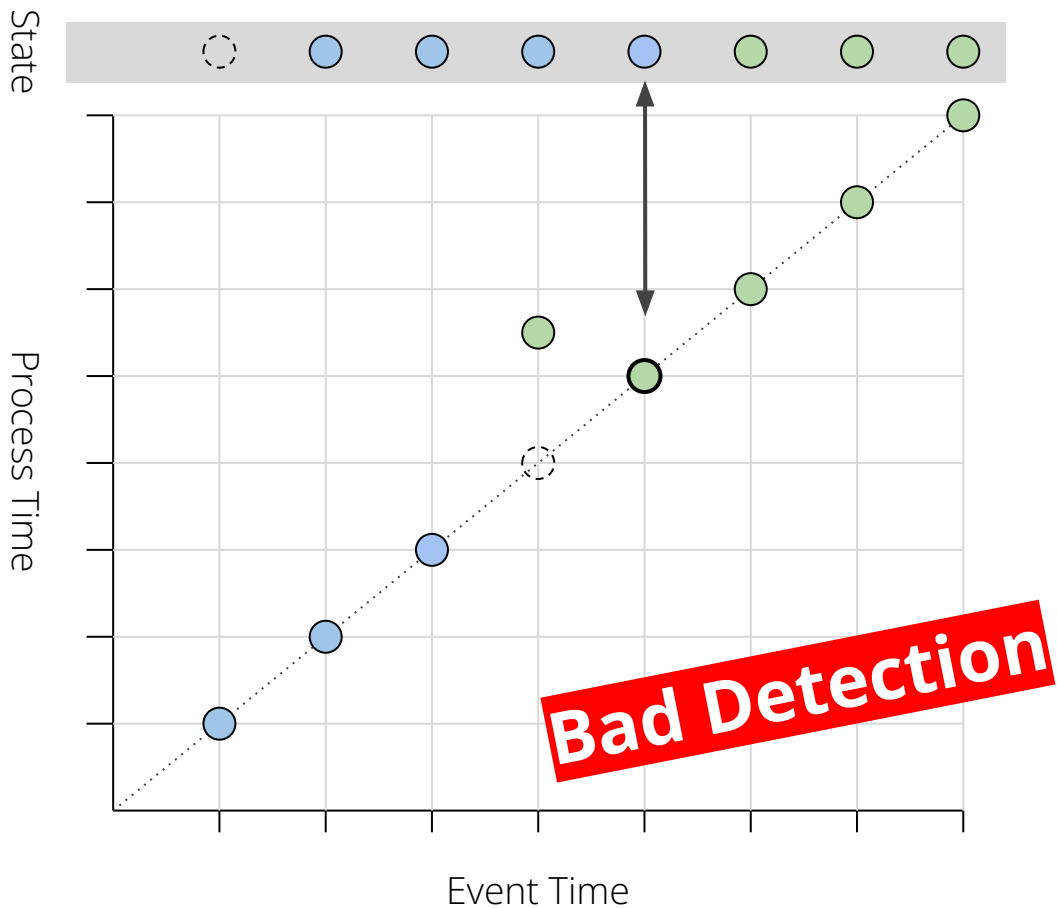
BEAM
SUMMIT

Austin, 2022

Beam State to Detect Change-Points

```
DoFn<KV<String, T>, T> {
  StateSpec<ValueSpec<T>>
  prevSpec =
    StateSpecs.value(...);

  public void processElement(
    ProcessContext c,
    ValueState<T> prev) {
    T curr =
c.element().getValue();
    T last = prev.read();
    if (curr != last) {
      c.output(curr);
    }
  }
}
```



Solution:

Watermark-Triggered Windows

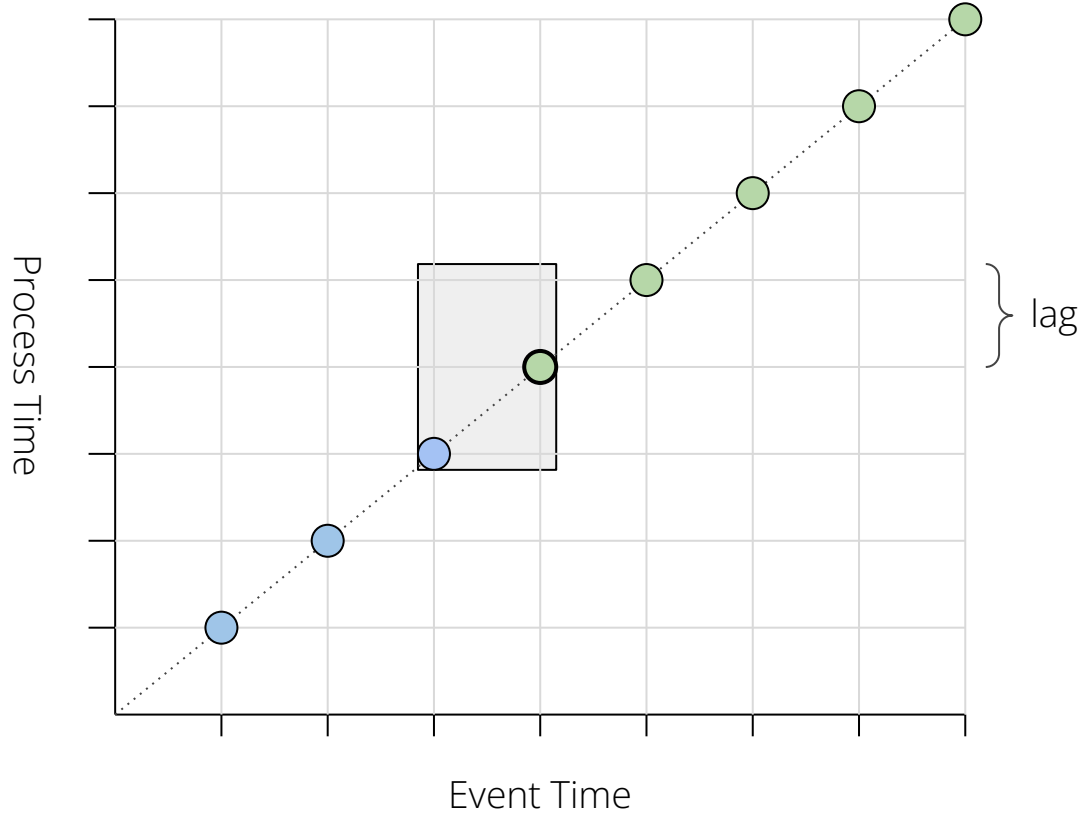


BEAM
SUMMIT

Austin, 2022

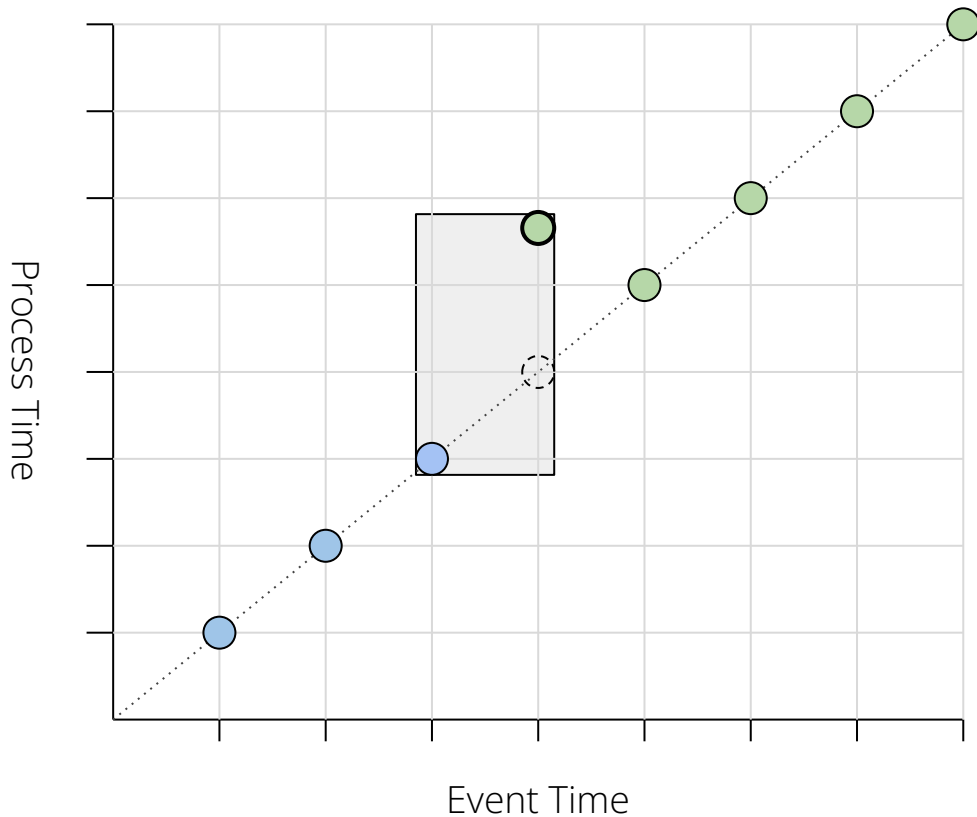
Watermark-Triggered Windows

```
Window
.<T>into(
  SlidingWindows
    .of(TWO_SECONDS)
    .every(ONE_SECOND))
.accumulatingFiredPanels()
.triggering(
  Repeatedly.forever(
    AfterWatermark
      .pastEndOfWindow()))
```



Watermark-Triggered Windows and Out-of-order Data

```
Window
.<T>into(
  SlidingWindows
    .of(TWO_SECONDS)
    .every(ONE_SECOND))
.accumulatingFiredPanels()
.triggering(
  Repeatedly.forever(
    AfterWatermark
      .pastEndOfWindow()))
```



Issues:

Watermark-Triggered Windows and Lag

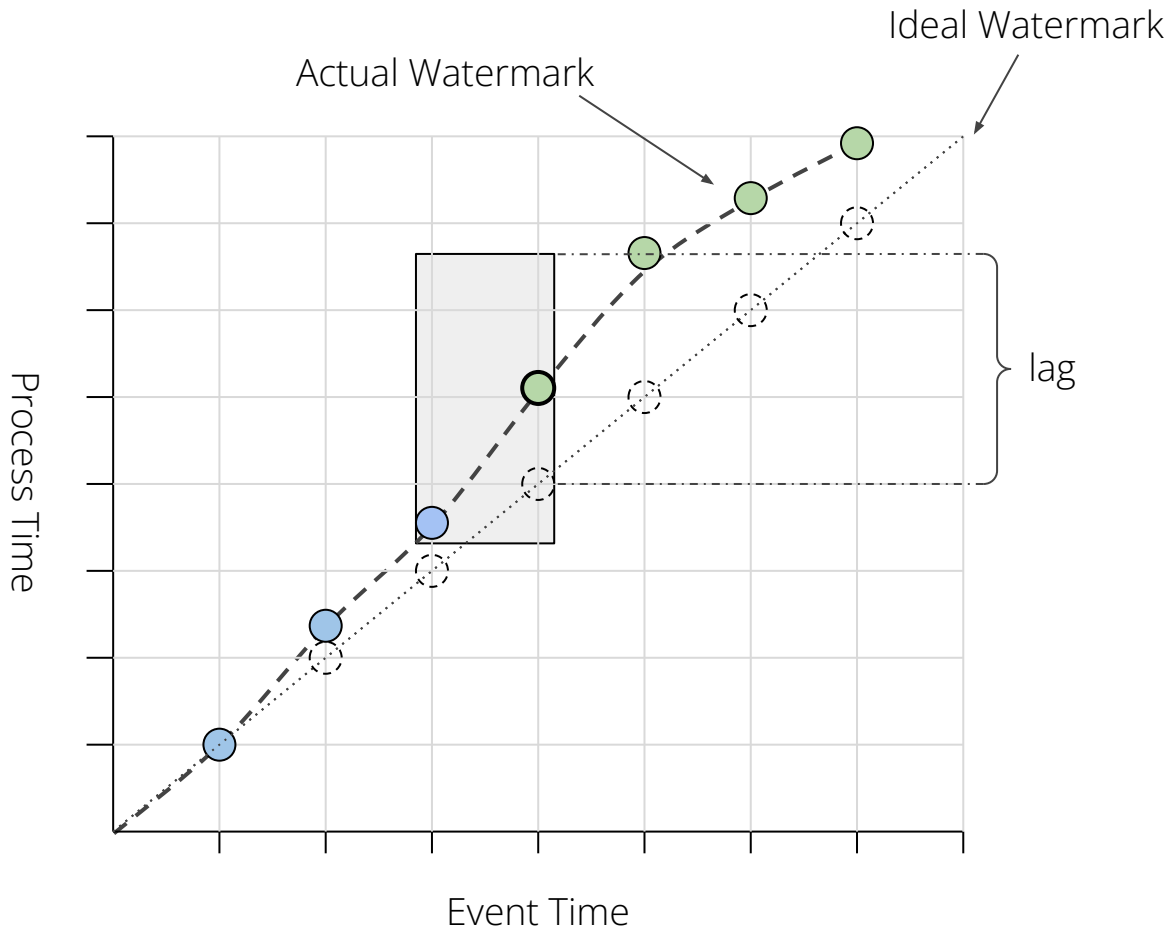


BEAM
SUMMIT

Austin, 2022

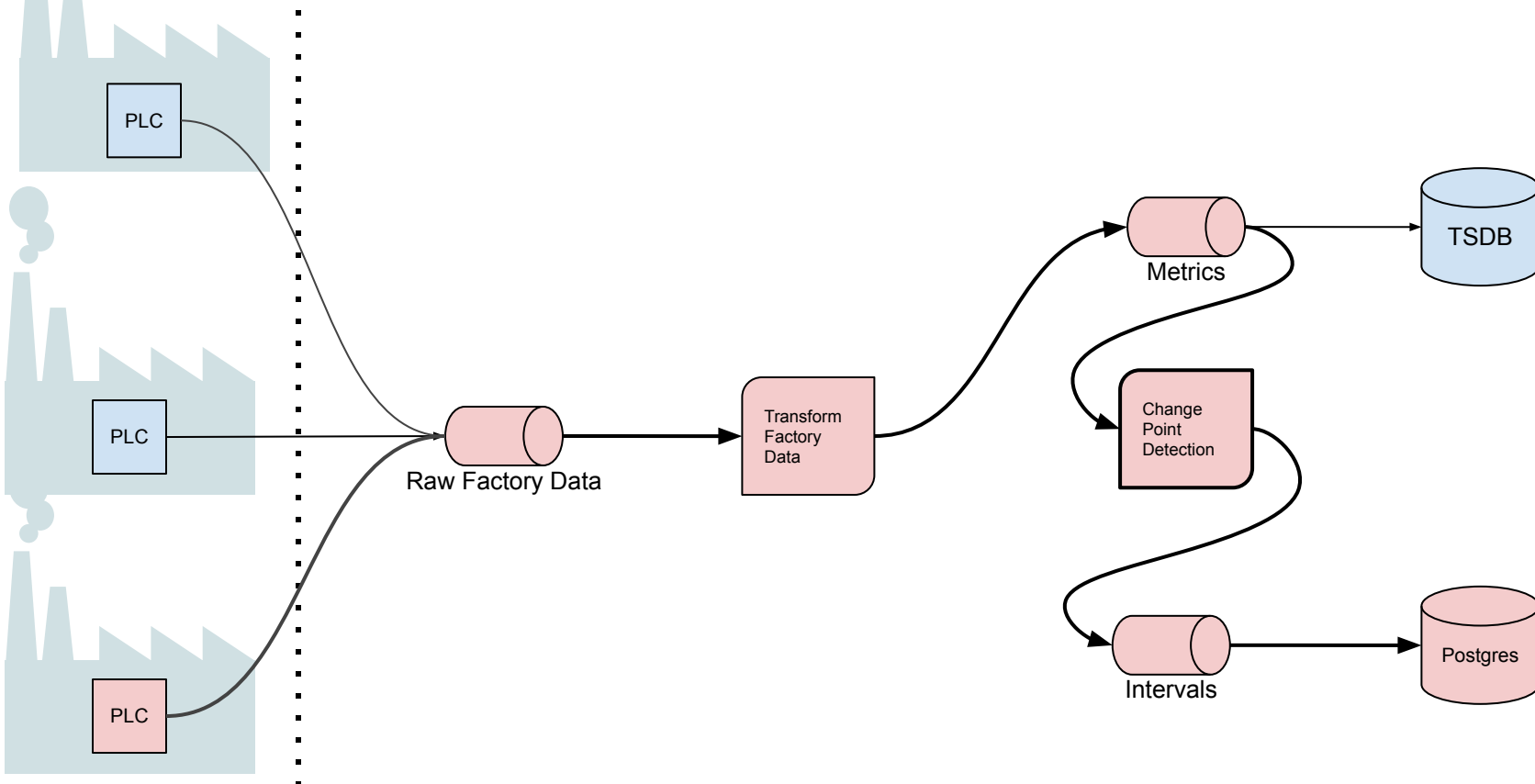
Watermark-Triggered Windows and Lagging Data

```
Window
.<T>into(
  SlidingWindows
    .of(TWO_SECONDS)
    .every(ONE_SECOND))
.accumulatingFiredPanels()
.triggering(
  Repeatedly.forever(
    AfterWatermark
    .pastEndOfWindow()))
```



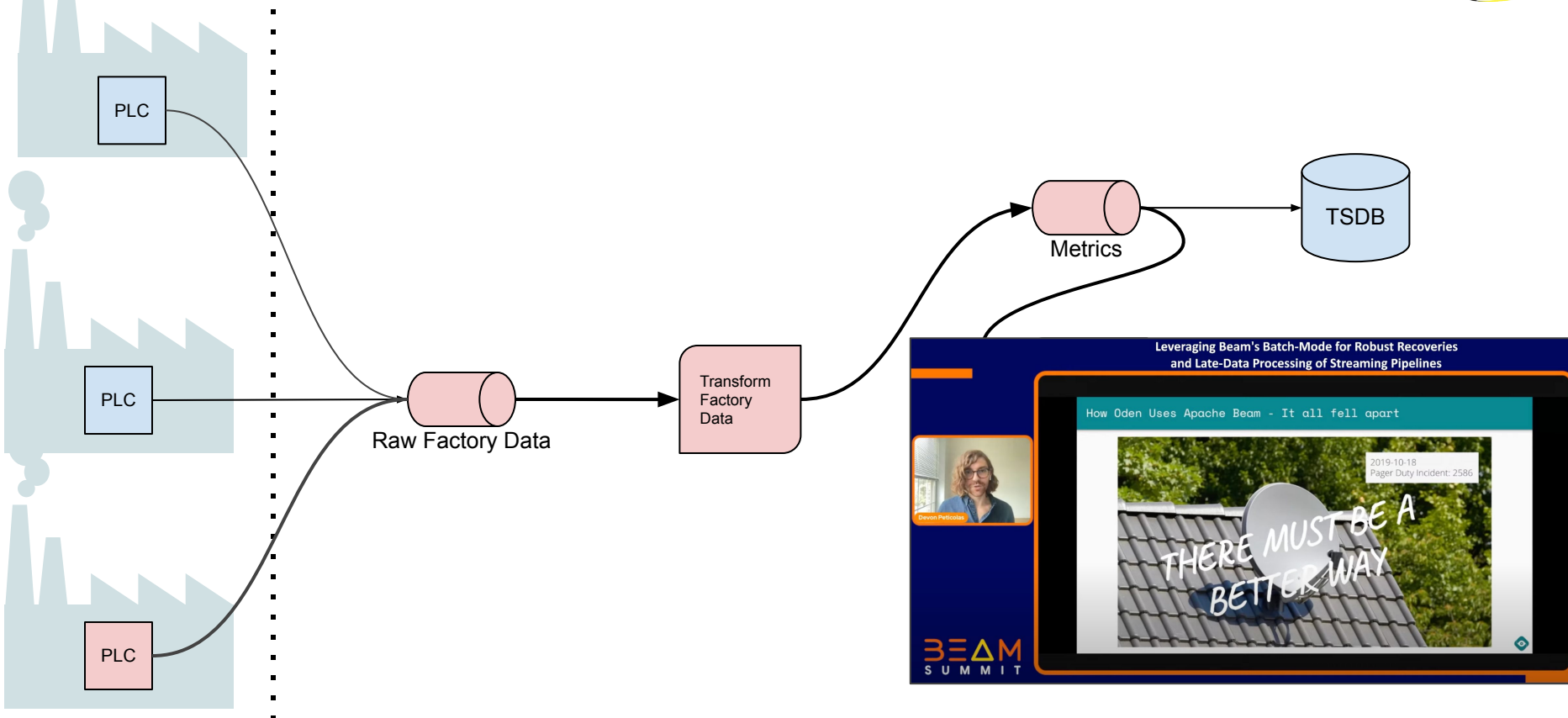


Non-homogeneous Lag



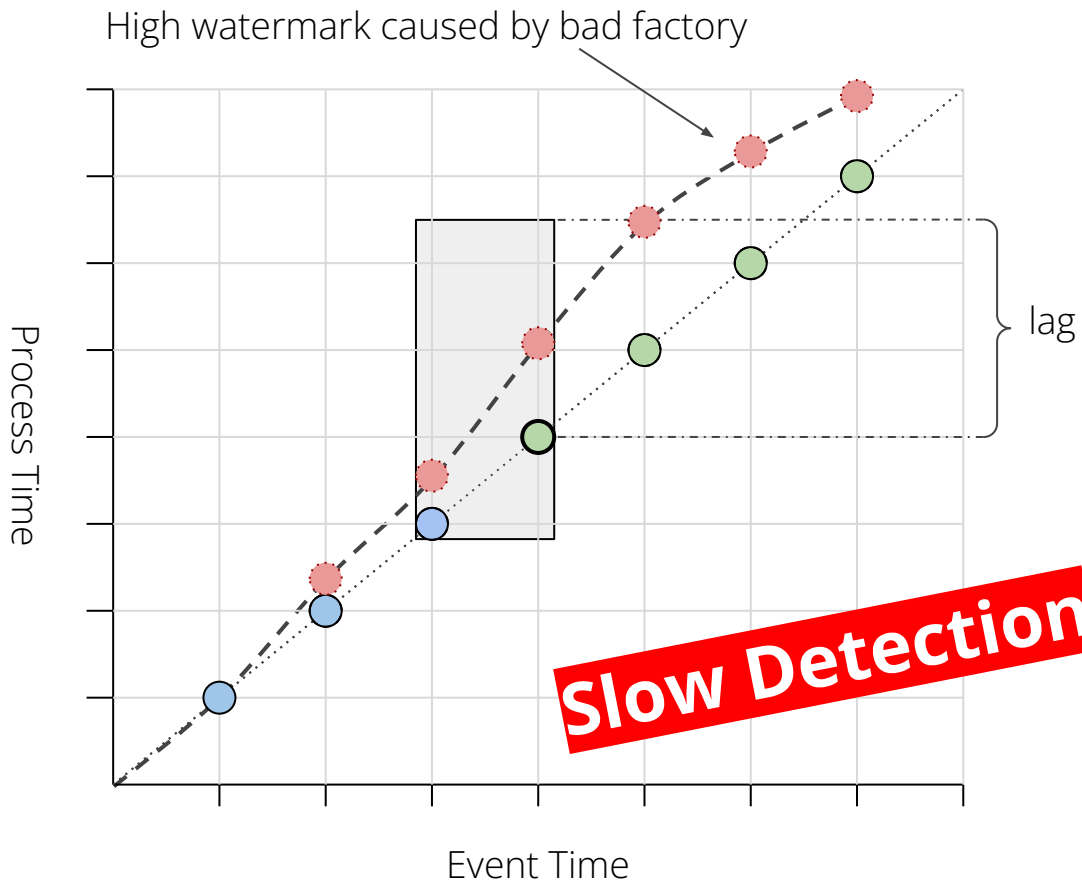


Non-homogeneous Lag



Watermark-Triggered Windows and Lagging Data

```
Window
.<T>into(
  SlidingWindows
    .of(TWO_SECONDS)
    .every(ONE_SECOND))
.accumulatingFiredPanels()
.triggering(
  Repeatedly.forever(
    AfterWatermark
      .pastEndOfWindow()))
```



Solution: Data-Triggered Windows

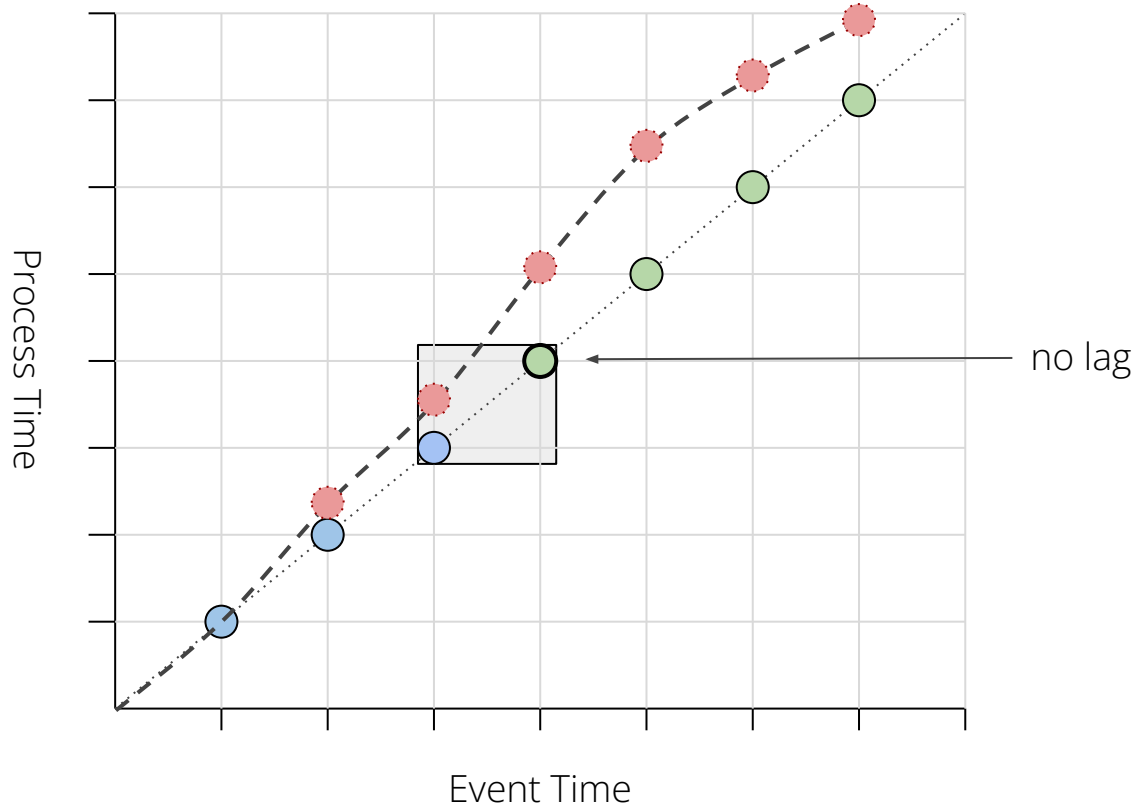


BEAM
SUMMIT

Austin, 2022

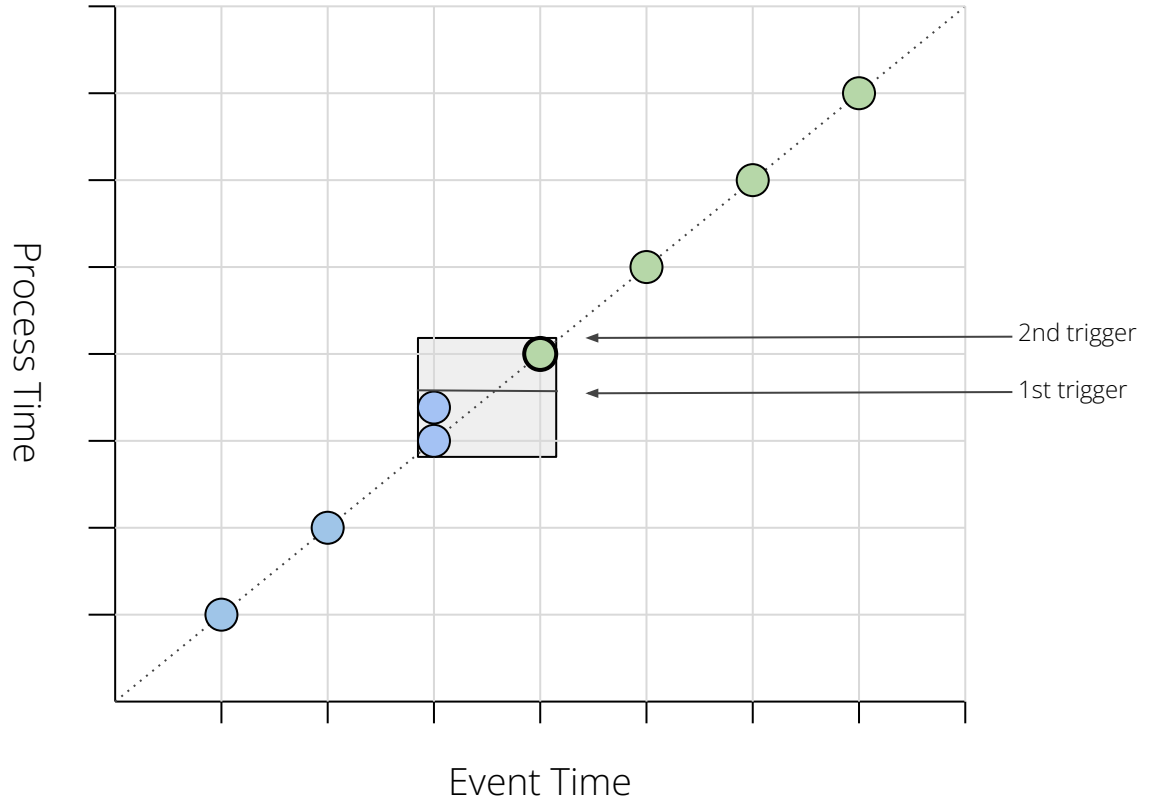
Data-Triggered Windows

```
Window
.<T>into(
  SlidingWindows
    .of(TWO_SECONDS)
    .every(ONE_SECOND))
.accumulatingFiredPanels()
.triggering(
  Repeatedly.forever(
    AfterPane
      .elementCountAtLeast(2)))
```



Data-Triggered Windows

```
Window
.<T>into(
  SlidingWindows
    .of(TWO_SECONDS)
    .every(ONE_SECOND))
.accumulatingFiredPanels()
.triggering(
  Repeatedly.forever(
    AfterPane
    .elementCountAtLeast(2)))
```



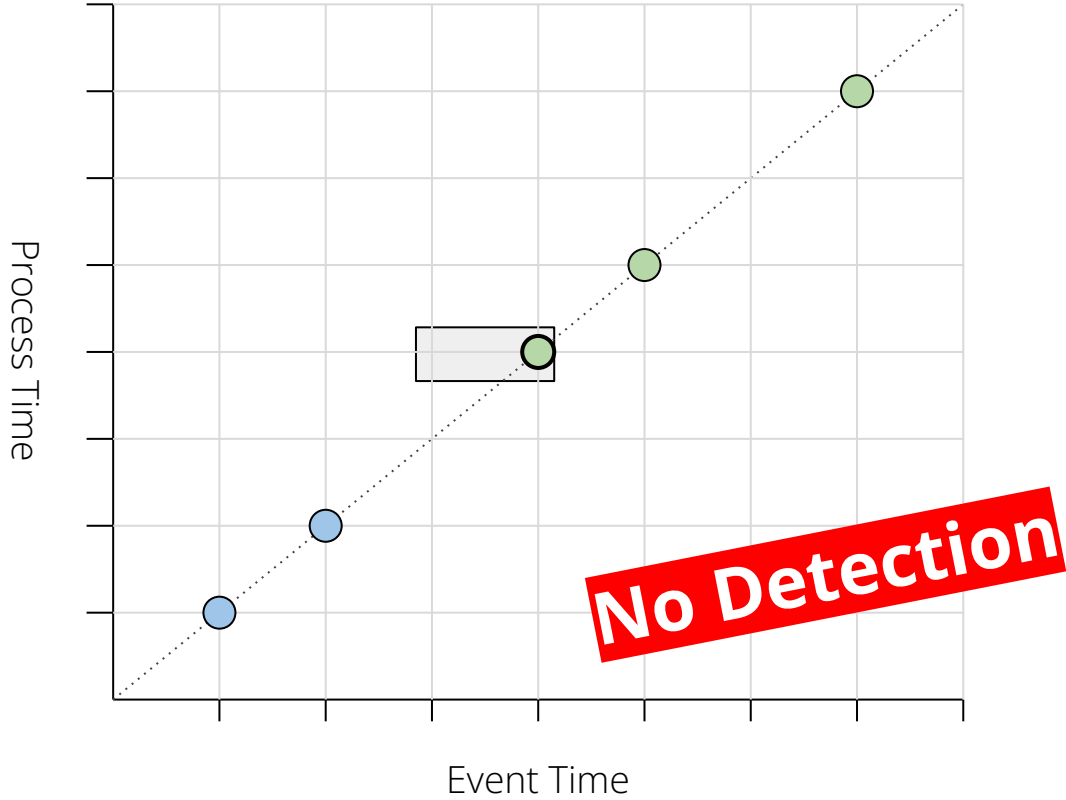
Issues:

Using Windows but Sparse Data

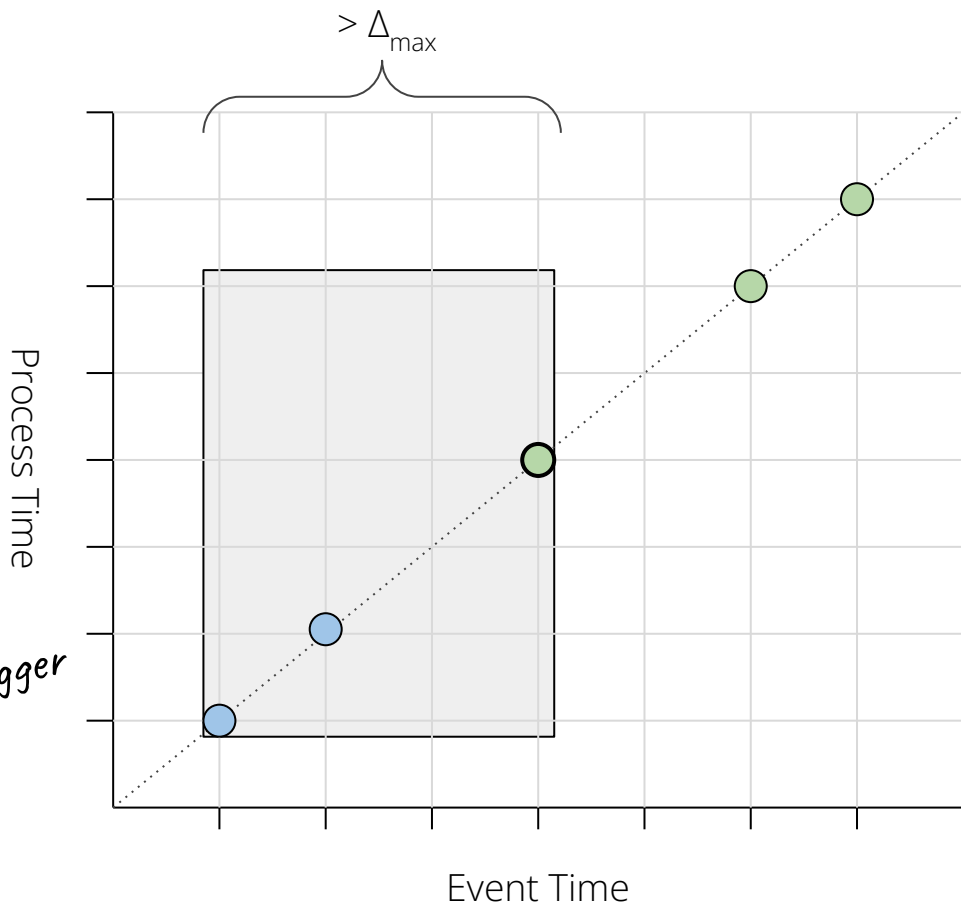


BEAM
SUMMIT

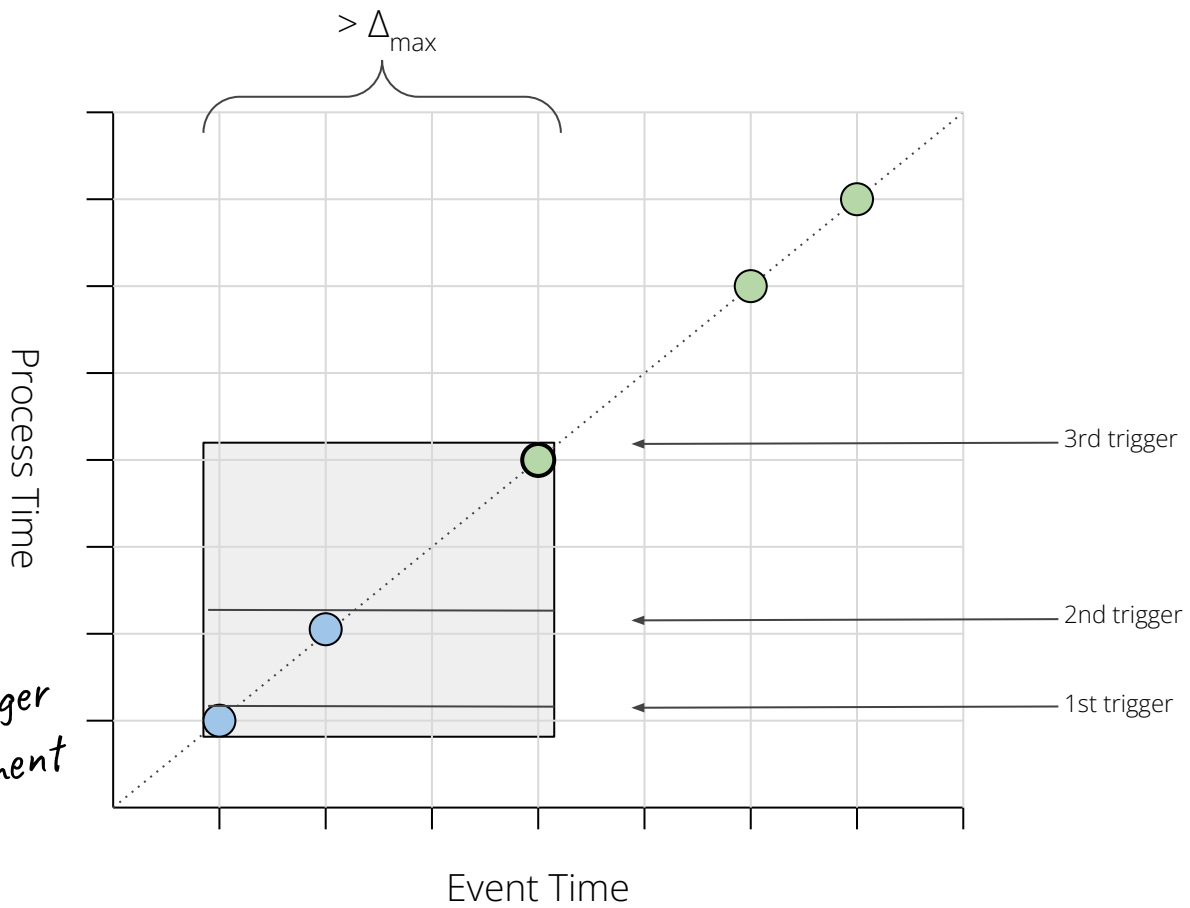
Austin, 2022

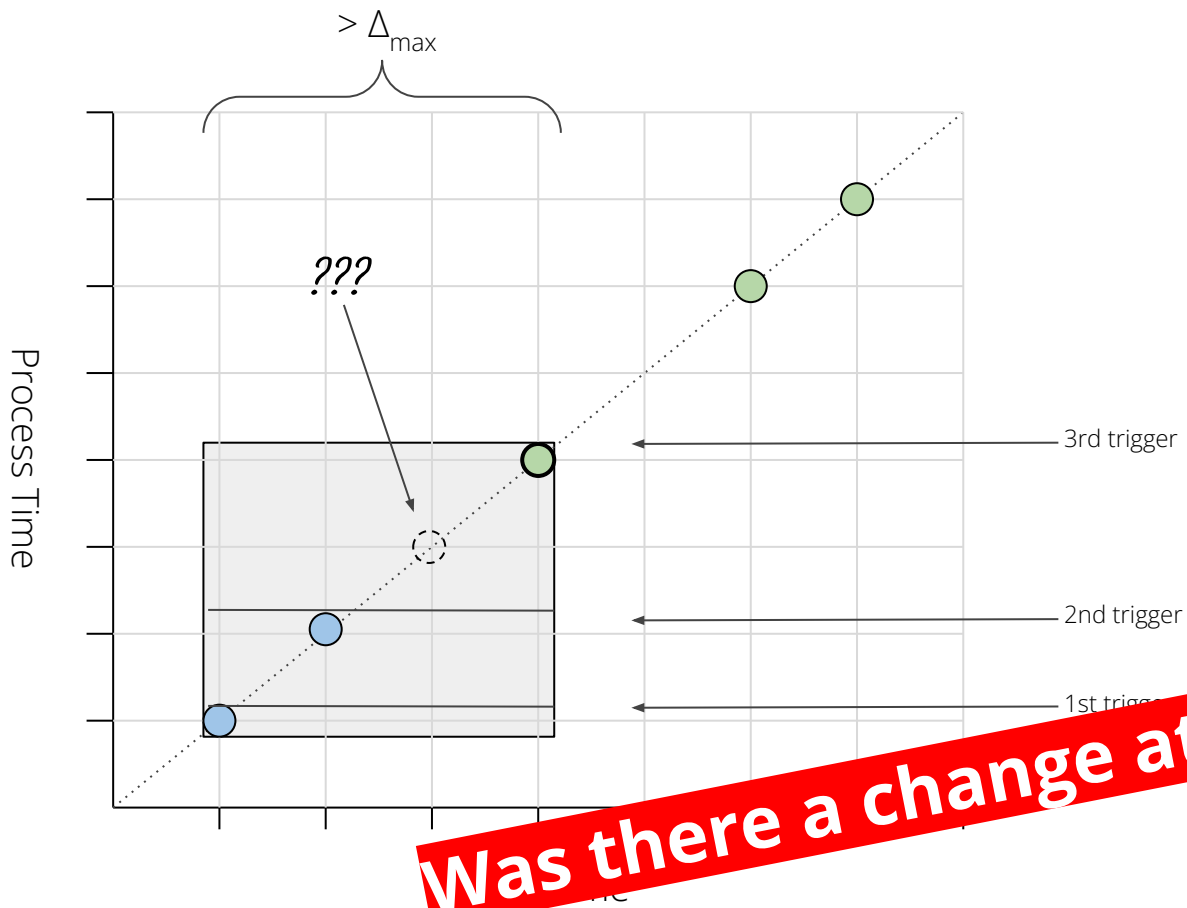


Watermark-driven trigger



*Data-driven trigger
must trigger every element*





Was there a change at t3?

Use Case: Creating “Smoothed” Intervals from Metrics

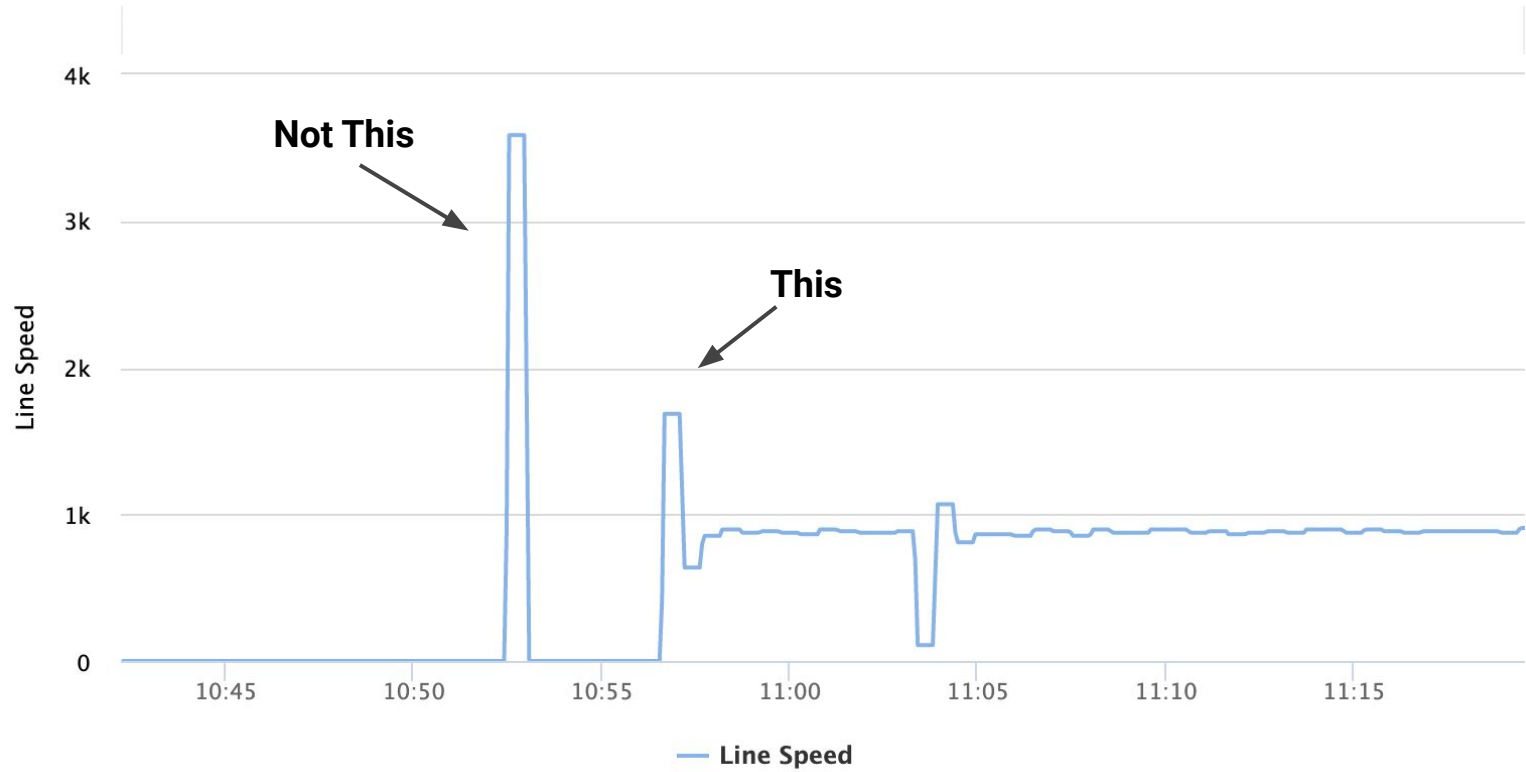


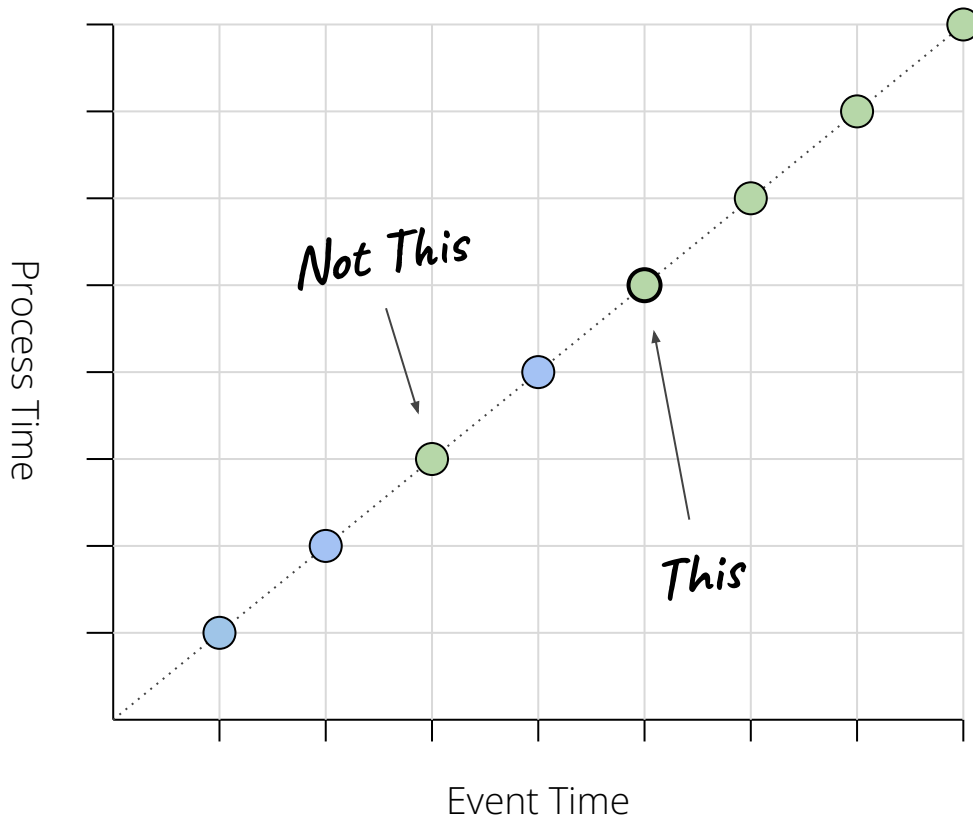
BEAM
SUMMIT

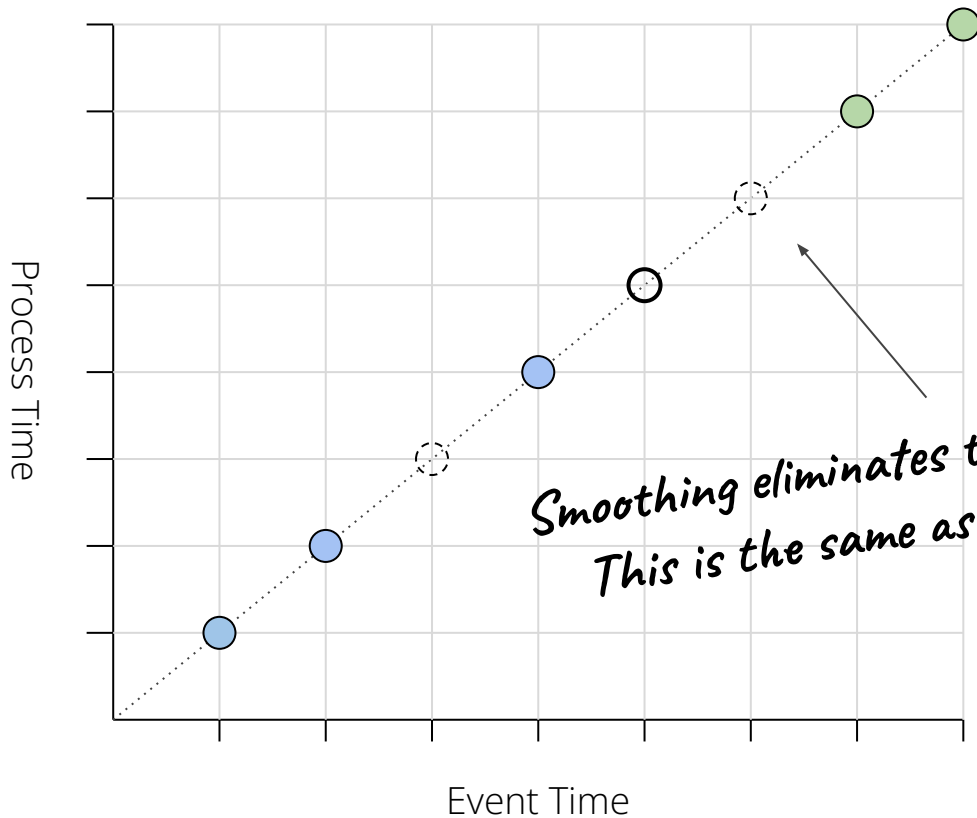
Austin, 2022

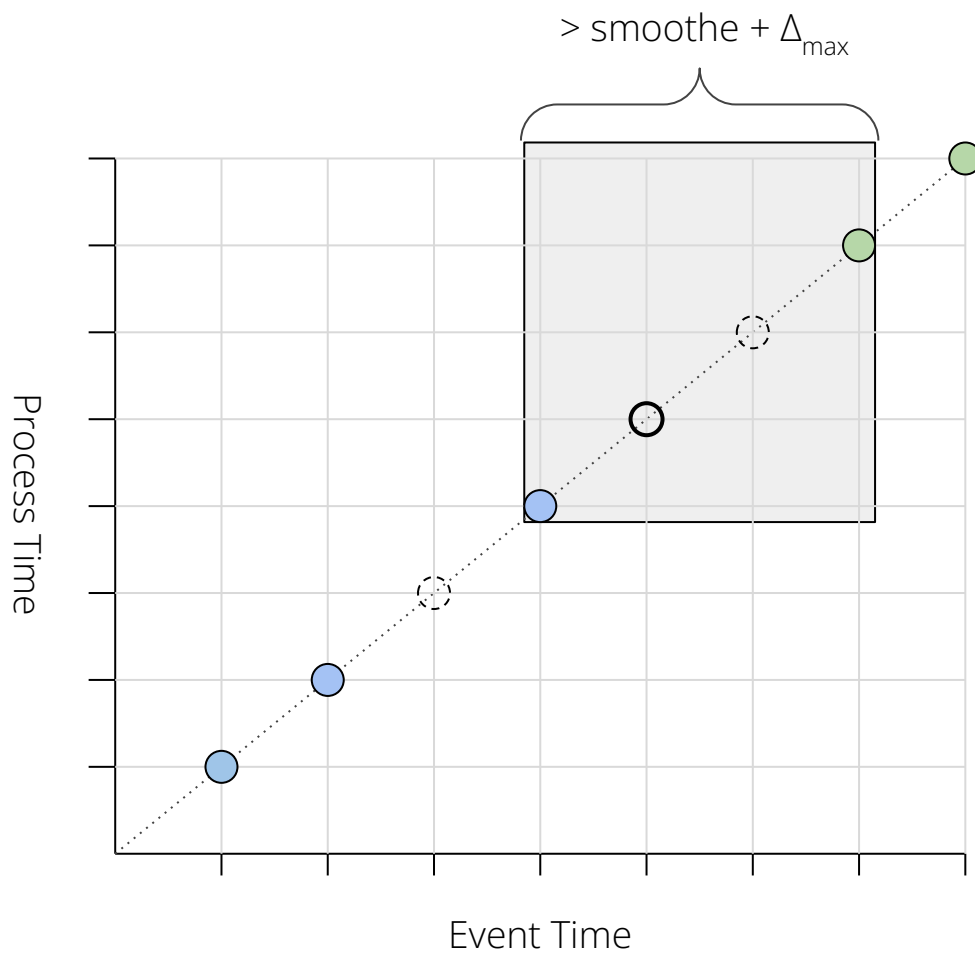
Downtime

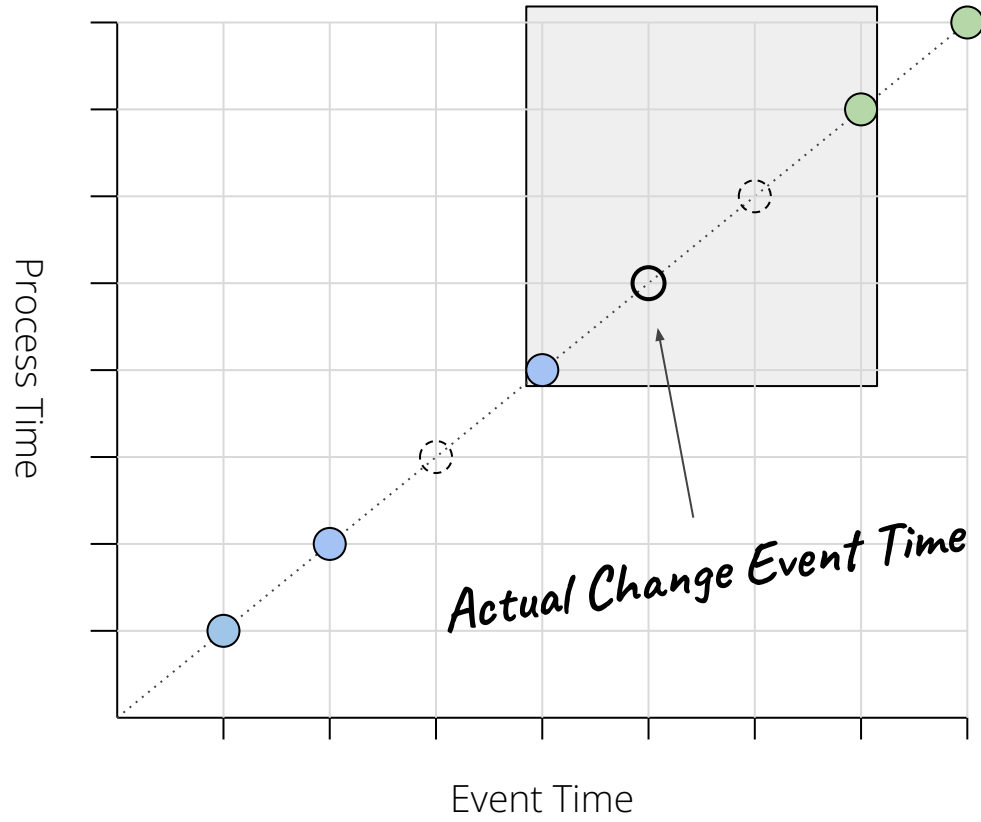
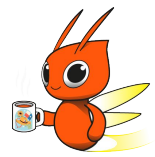
Uptime











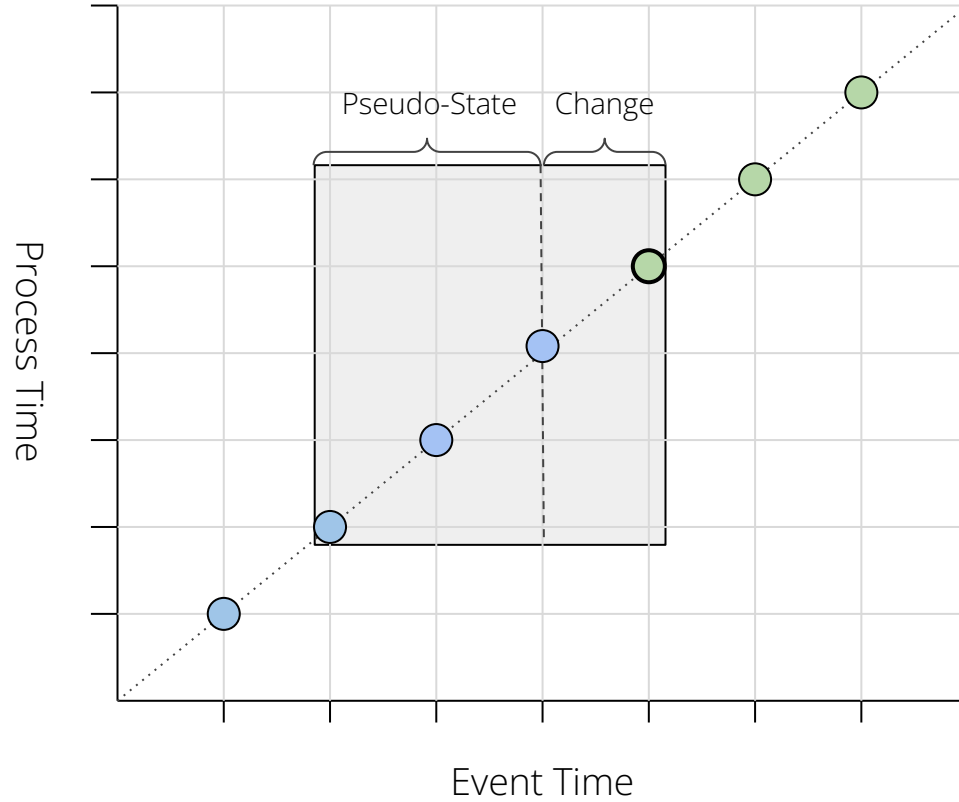
*Created lag is
always
>= smoothe*

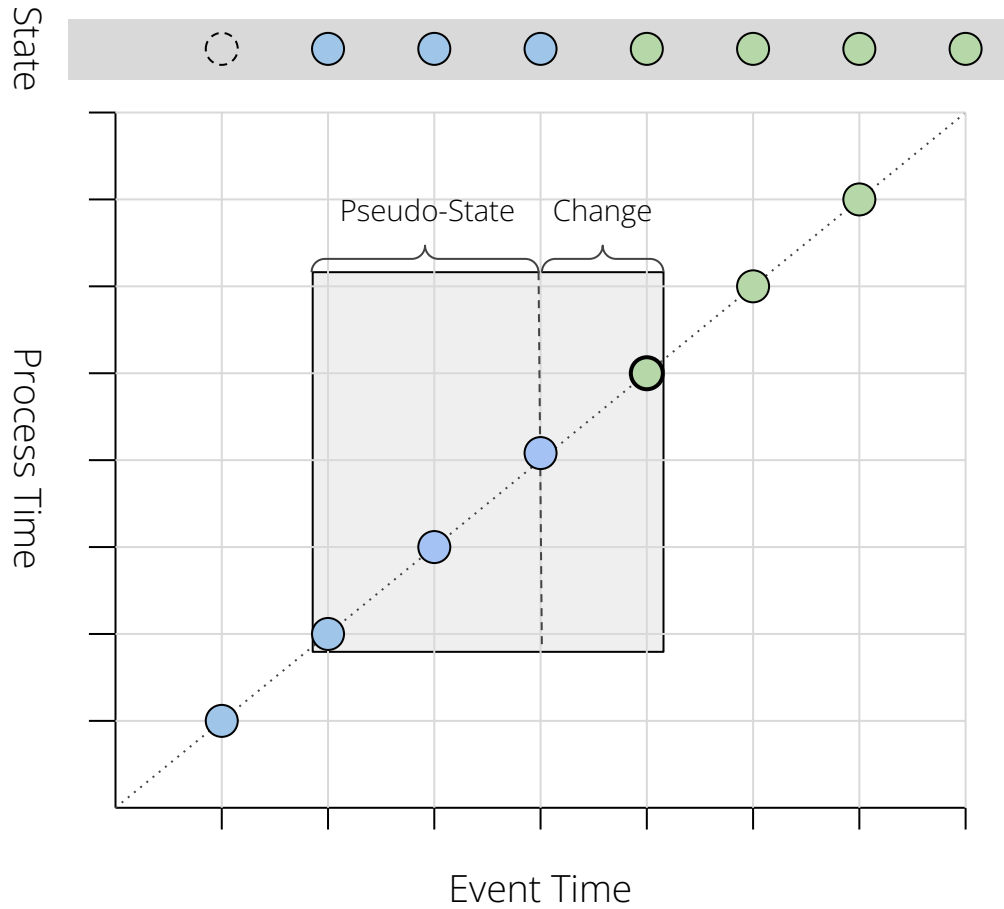
Solution: State + Sliding Windows

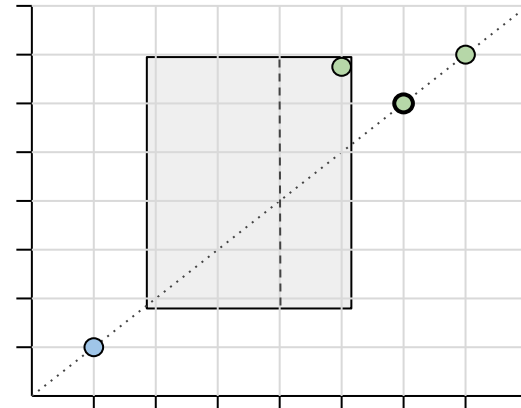
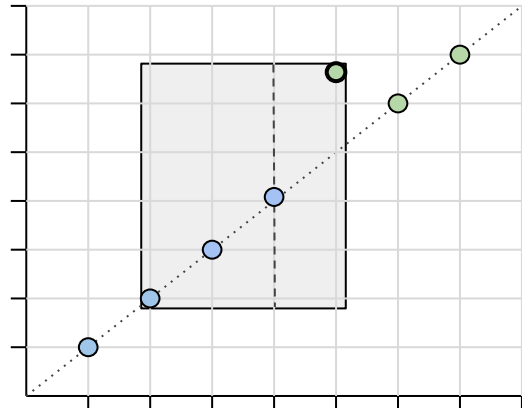
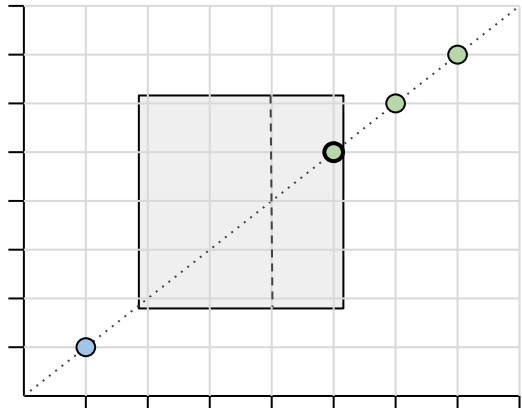


BEAM
SUMMIT

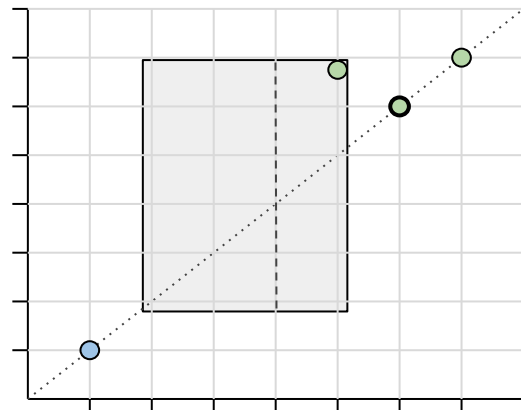
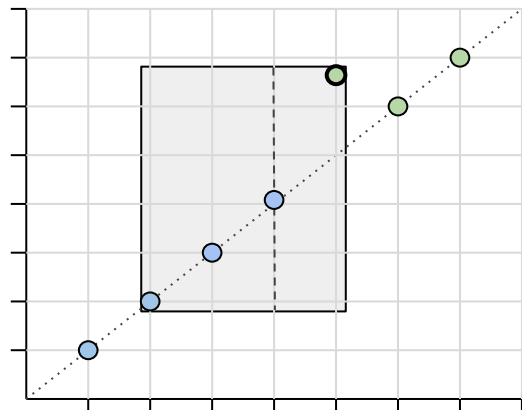
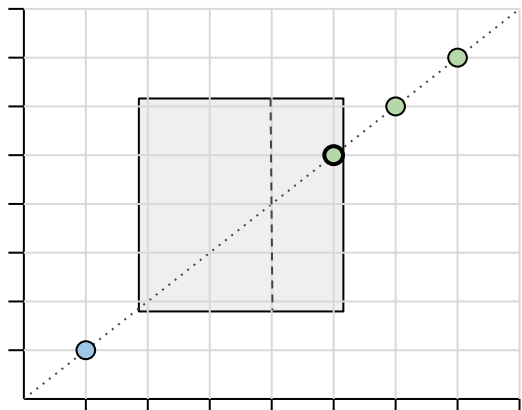
Austin, 2022







I think this is wrong but...



$$\text{smoothe-time} + \text{event-time } \Delta_{\max} + \text{out-of-order } \Delta_{\max} < \text{process-time } \Delta_{\min}$$

Key Takeaways



- Oden Uses Change-Point Detection to transform Metrics into Intervals
- Beam State is fast and good at sparsity, but bad at out-of-order
- Windowing is slower and good at out-of-order, but bad at sparsity
- Combining Beam State and Windowing is good at out-of-order and sparsity
- “Smoothed” Change-Point Detection is just a sparsity problem

Thank You

And a special thanks to Jie Zhang, Jake Skelcy, and Deepak Turaga

Questions?

Email: devon@petiocol.as
Github: github.com/x



BEAM
SUMMIT

Austin, 2022