

+ 25

# Using Distributed Trace for End-to-End Latency Metrics

KUSHA MAHARSHI



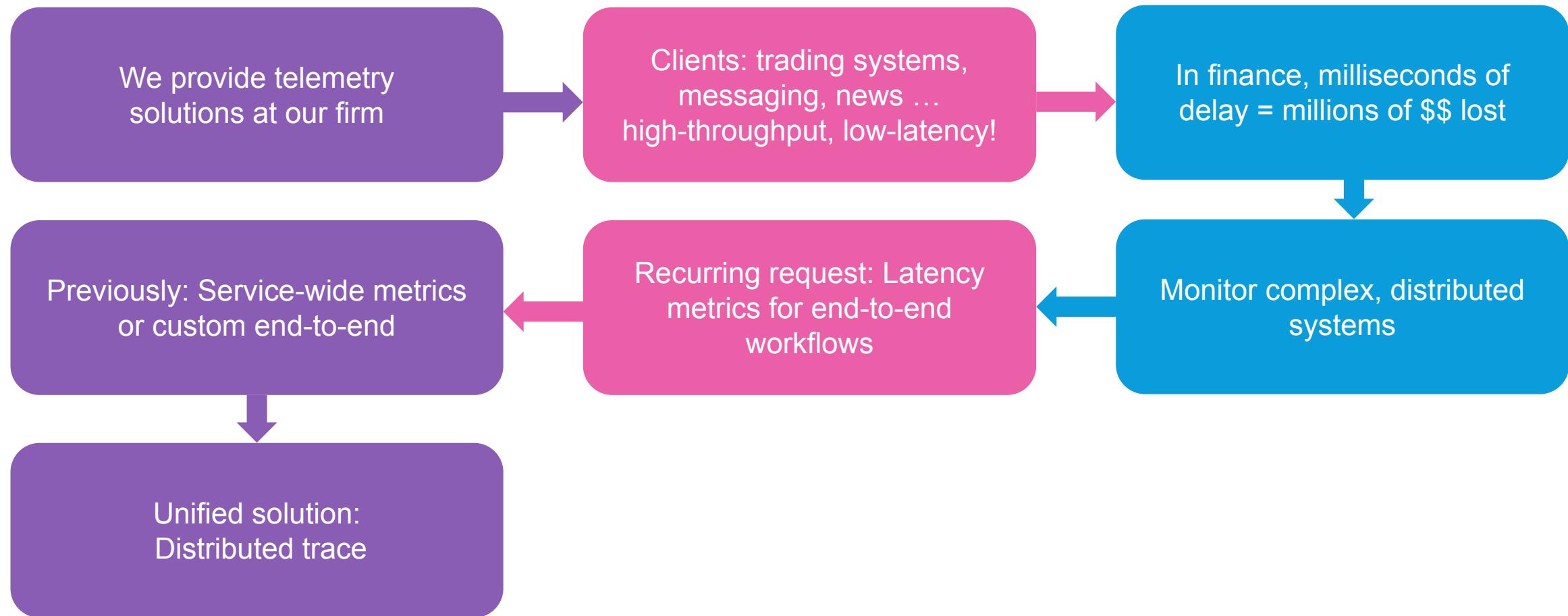
**Cppcon**  
The C++ Conference

20  
25

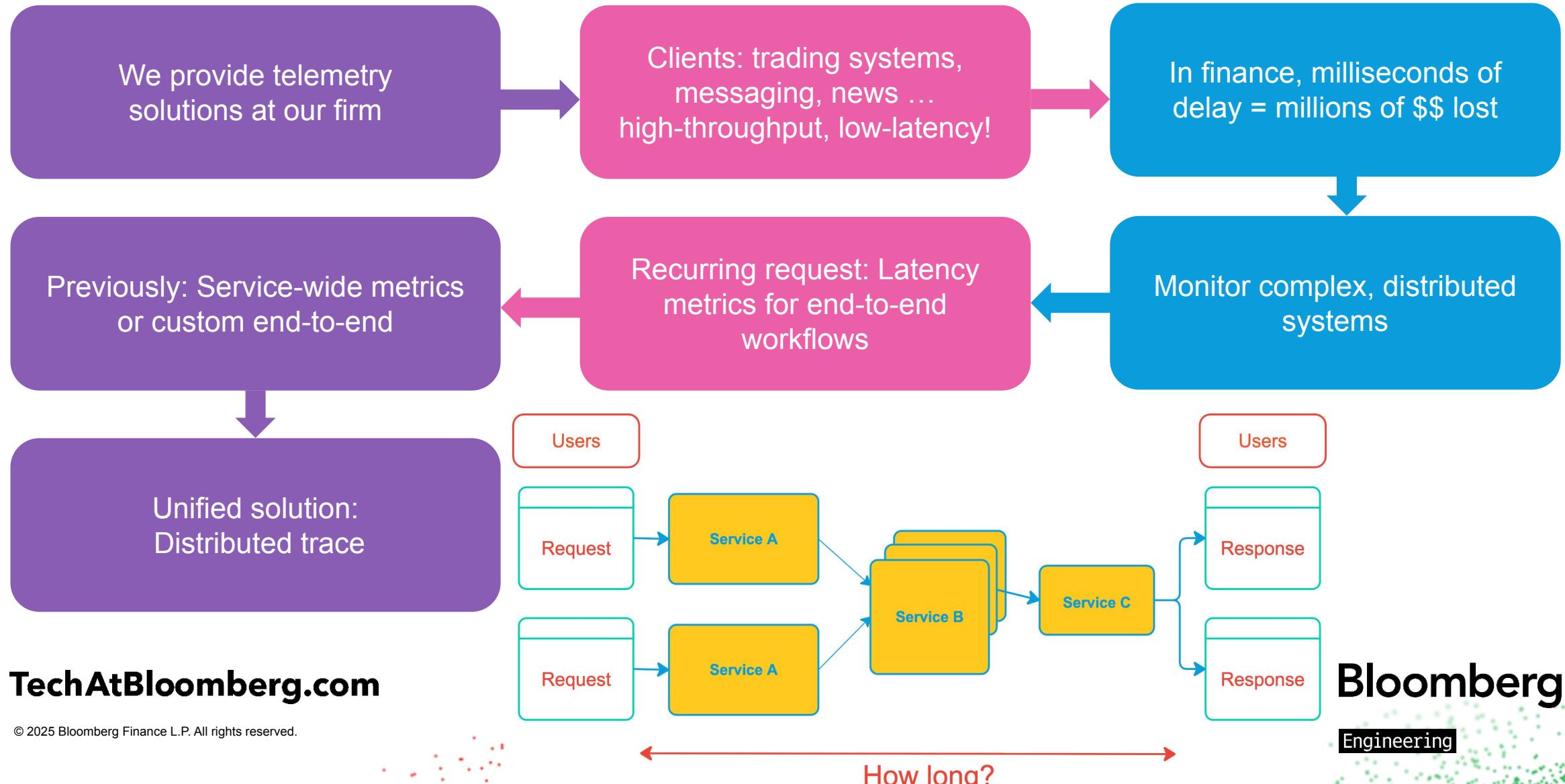


September 13 - 19

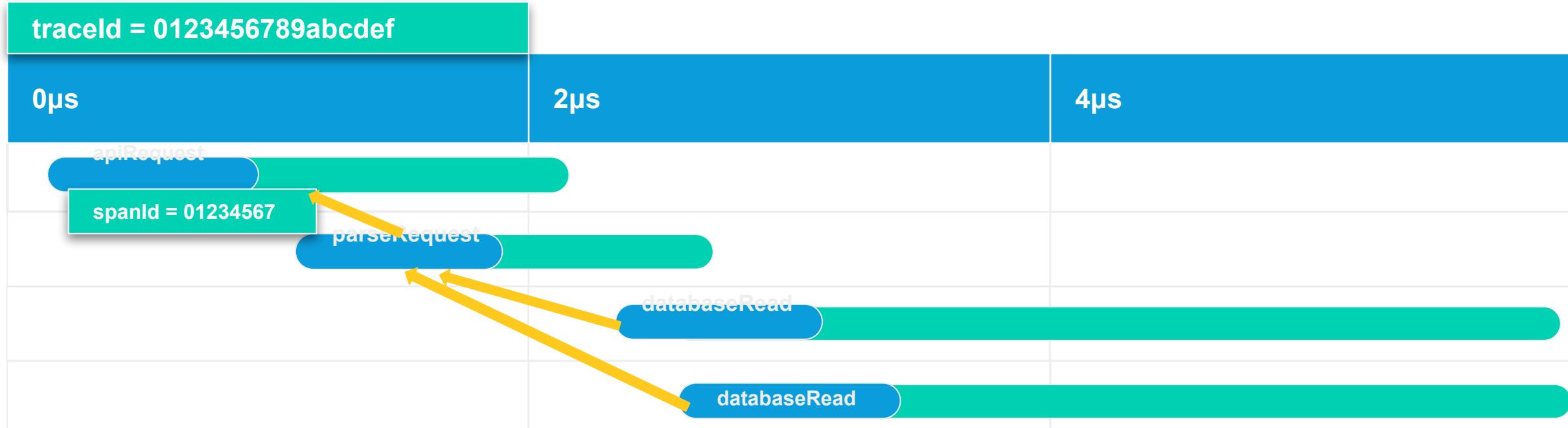
# Storytime



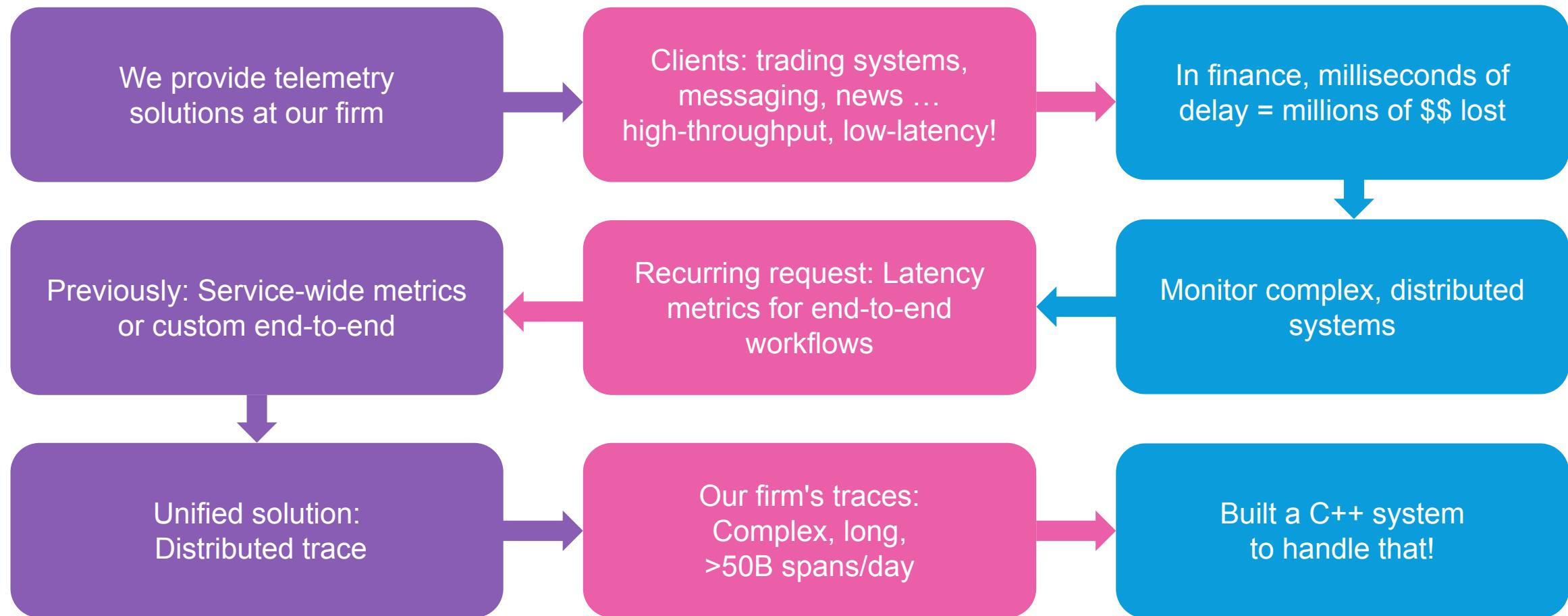
# Storytime



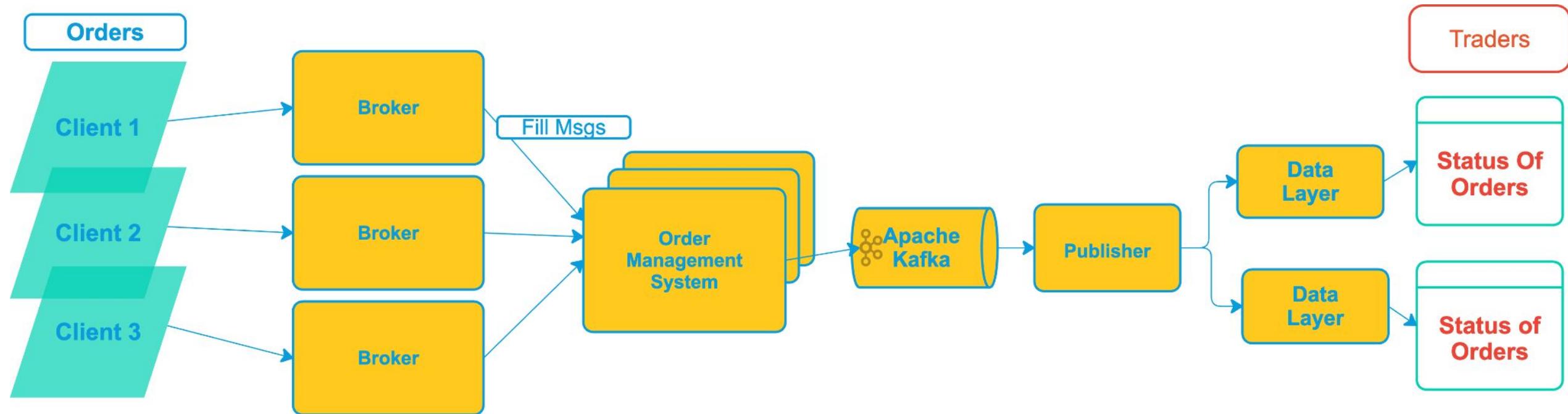
# What Is Distributed Tracing?



# Storytime



# Equity Order Execution: An End-to-End Workflow



# Why End-to-End Latency Monitoring Matters

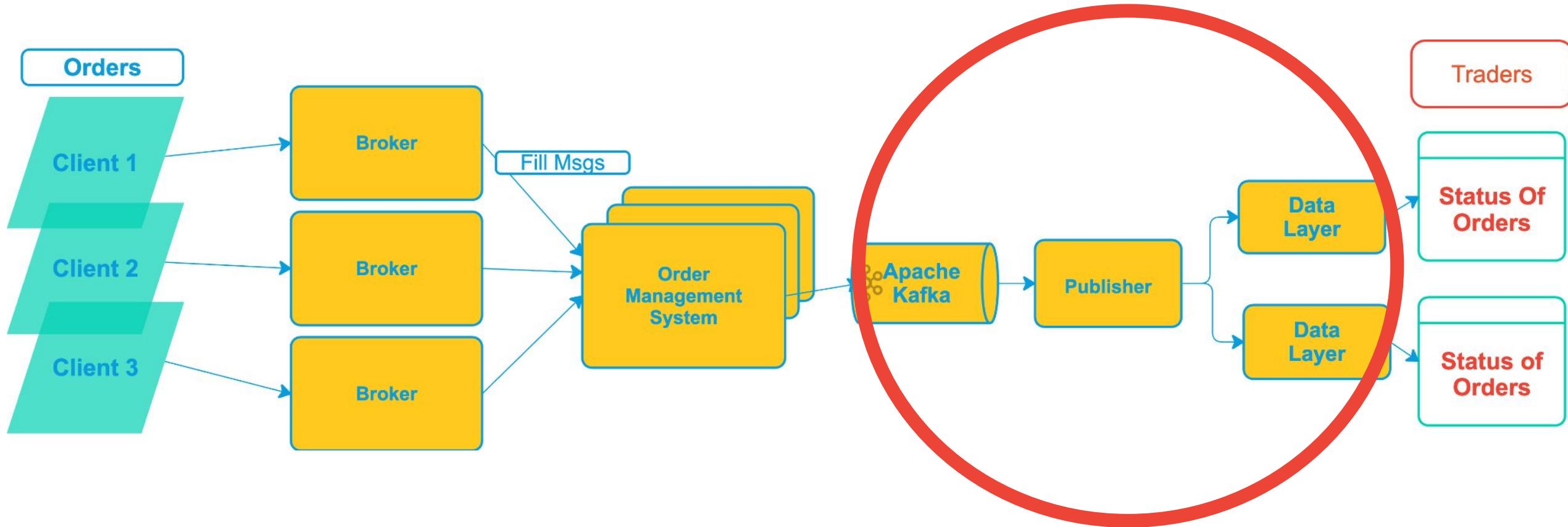
- Measure user experience
  - In complex, distributed architectures
- Alert on spikes
  - Improve time to resolution
- Service Level Objectives (SLOs)
  - Measurable targets for reliability and performance
  - "99.9% of requests should finish within 1s"
  - For engineers, managers, executives
- Design & development

TechAtBloomberg.com

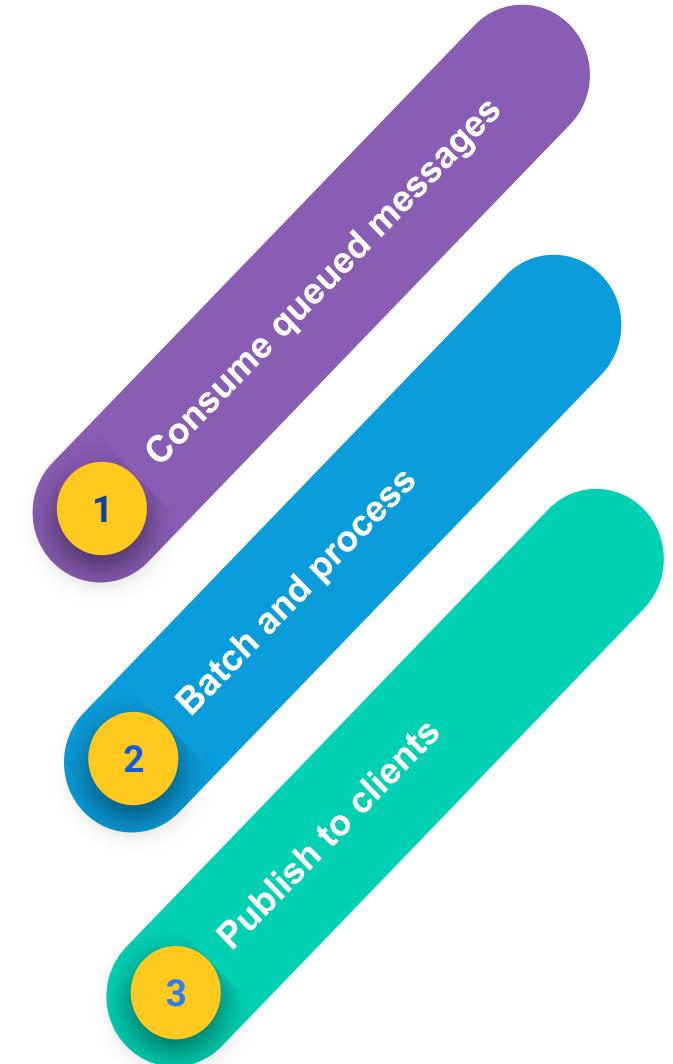
© 2025 Bloomberg Finance L.P. All rights reserved.

Bloomberg  
Engineering

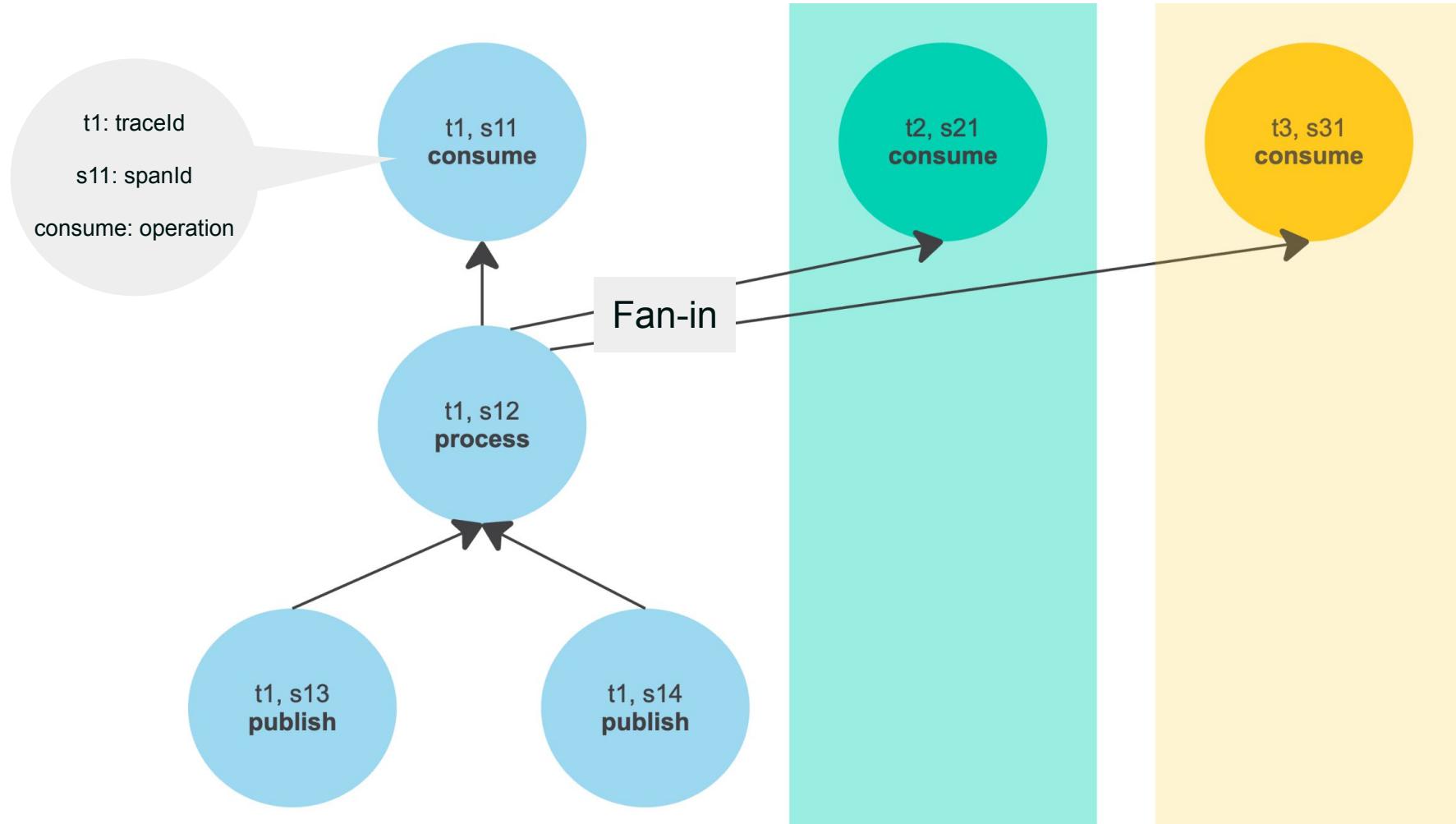
# Equity Order Execution: Zooming In ...



# Example Distributed Trace



# Trace = Directed Acyclic Graph (DAG)

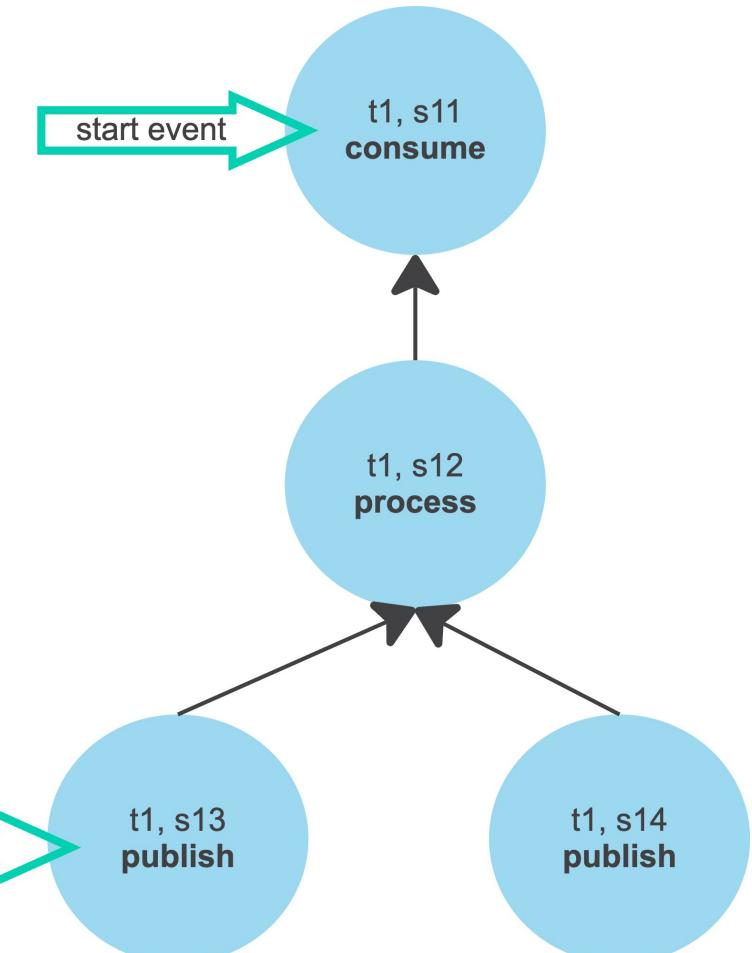
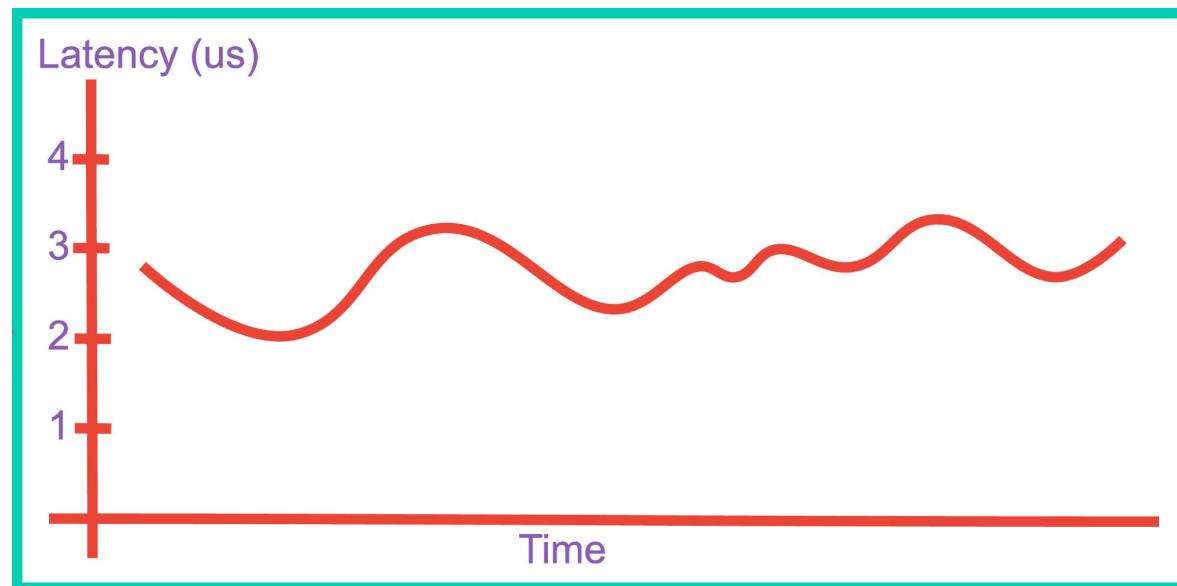


# Rule = Defines End-to-End Workflow

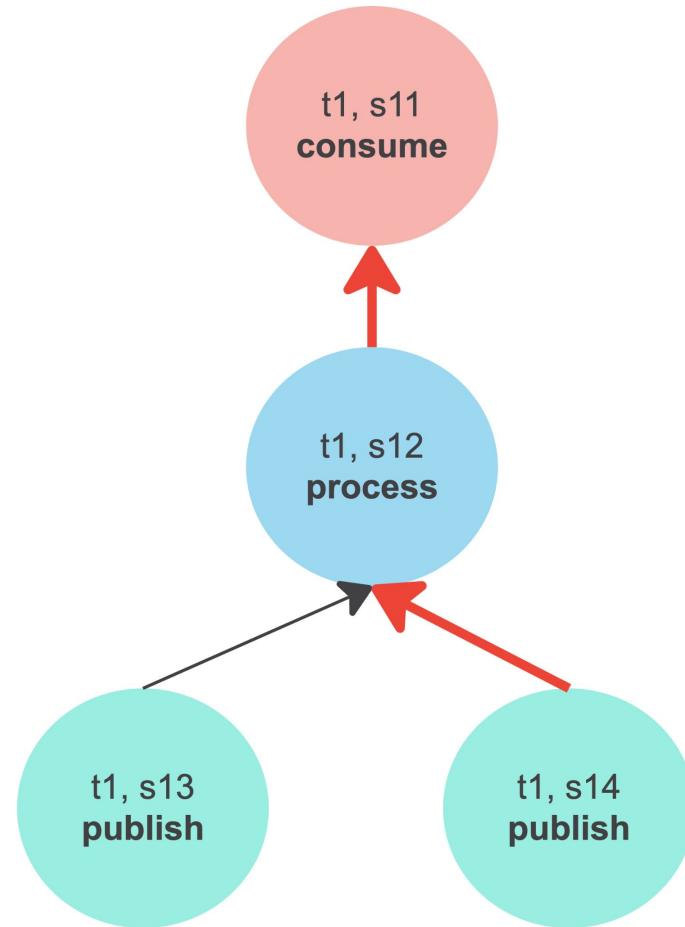
*Get me the latency of requests going from consume to publish*

Rule = consumeToPublish

**start** = consume, **end** = publish

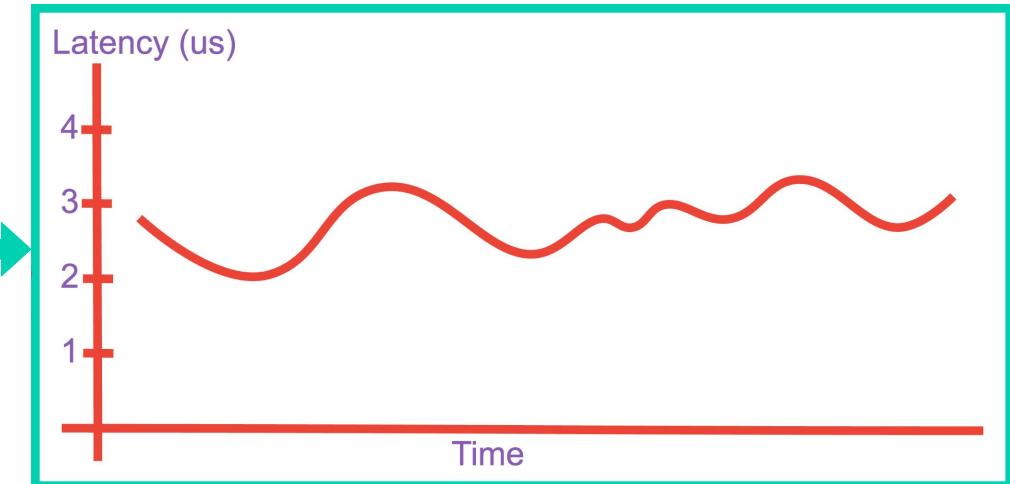


# End-to-End Latency From One Trace



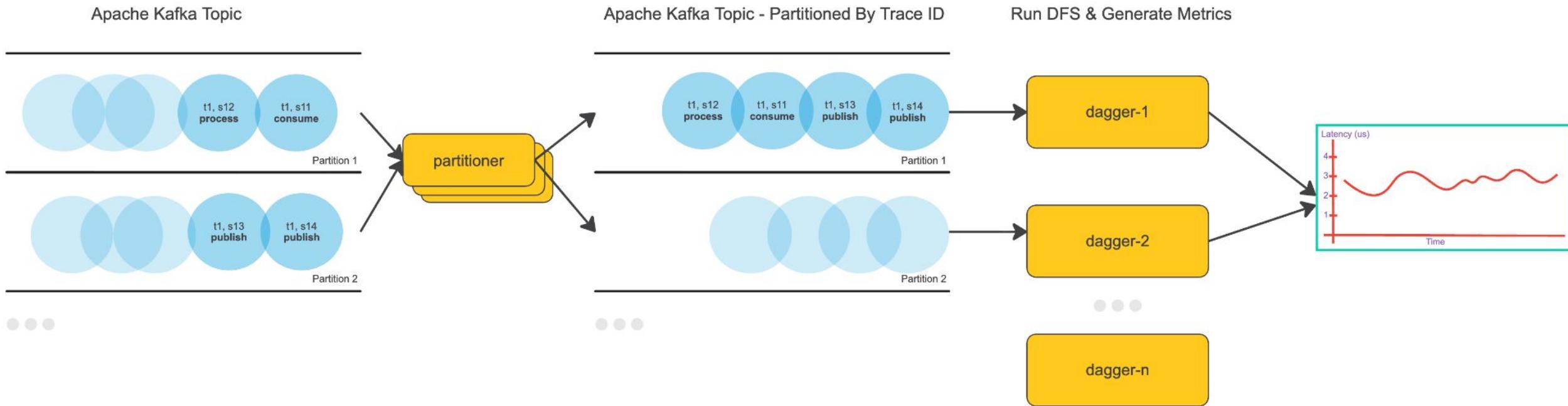
Difference between  
start of start event  
&  
end of end event

Path 1: 3 $\mu$ s  
Path 2: 3.5 $\mu$ s

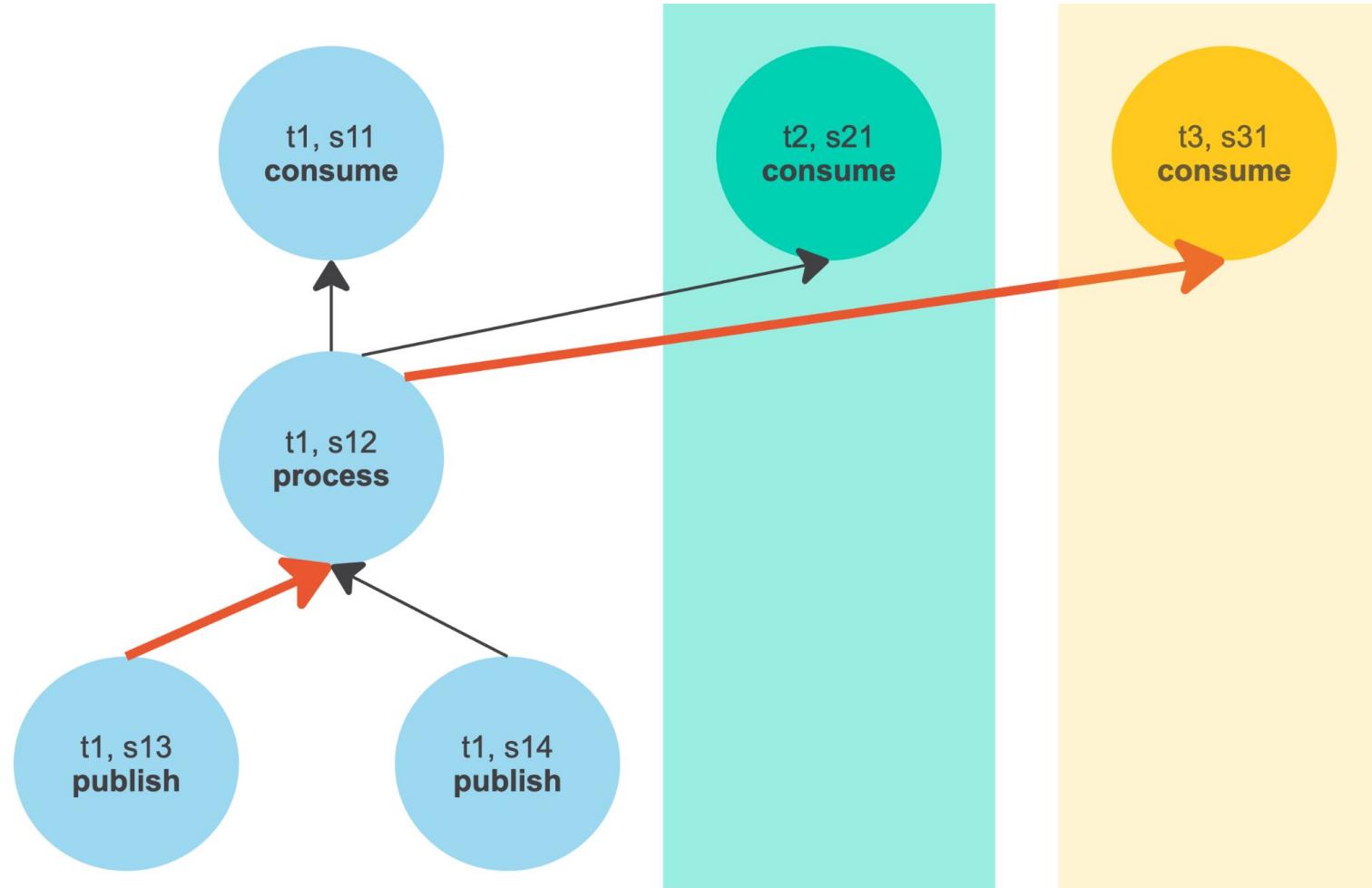


Depth-First Search (DFS) through DAG

# Handling Traces At Scale



# Enter ... (drumroll please) ... Fan-ins!



**Problem**  
Scaling with fan-ins

All spans from the  
same *weakly connected DAG*  
should be handled  
together.

Let's call this a  
trace "**bundle**."

# Two Approaches, One Thing In Common - C++

Memory Control

Concurrency At Scale

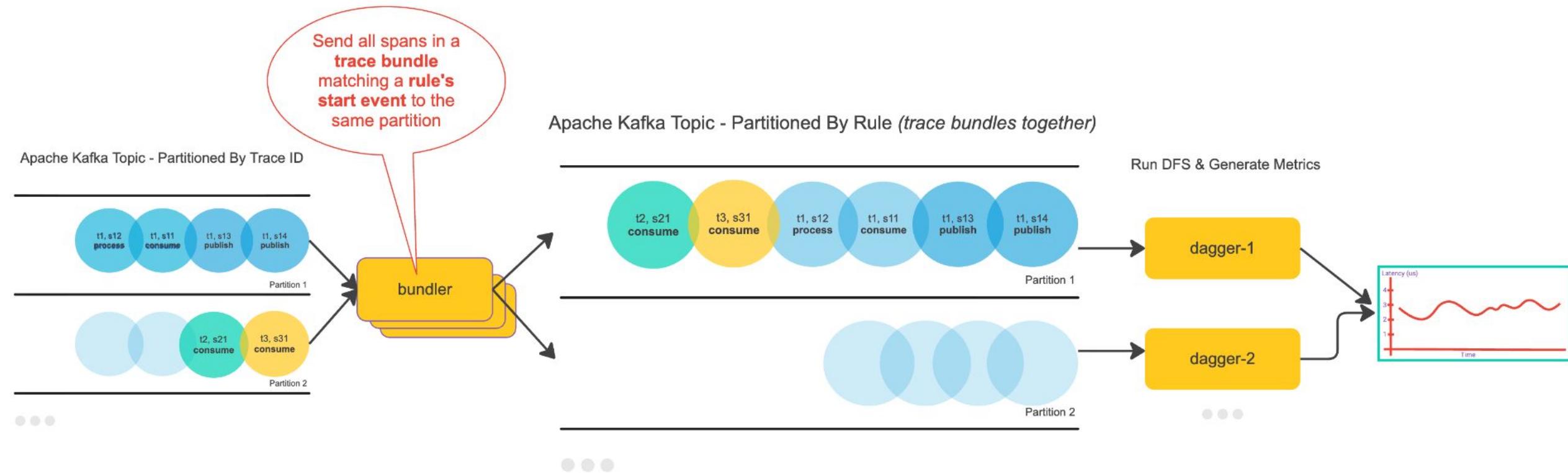
Performant Standard Data Structures

Mature Tooling & Ecosystem

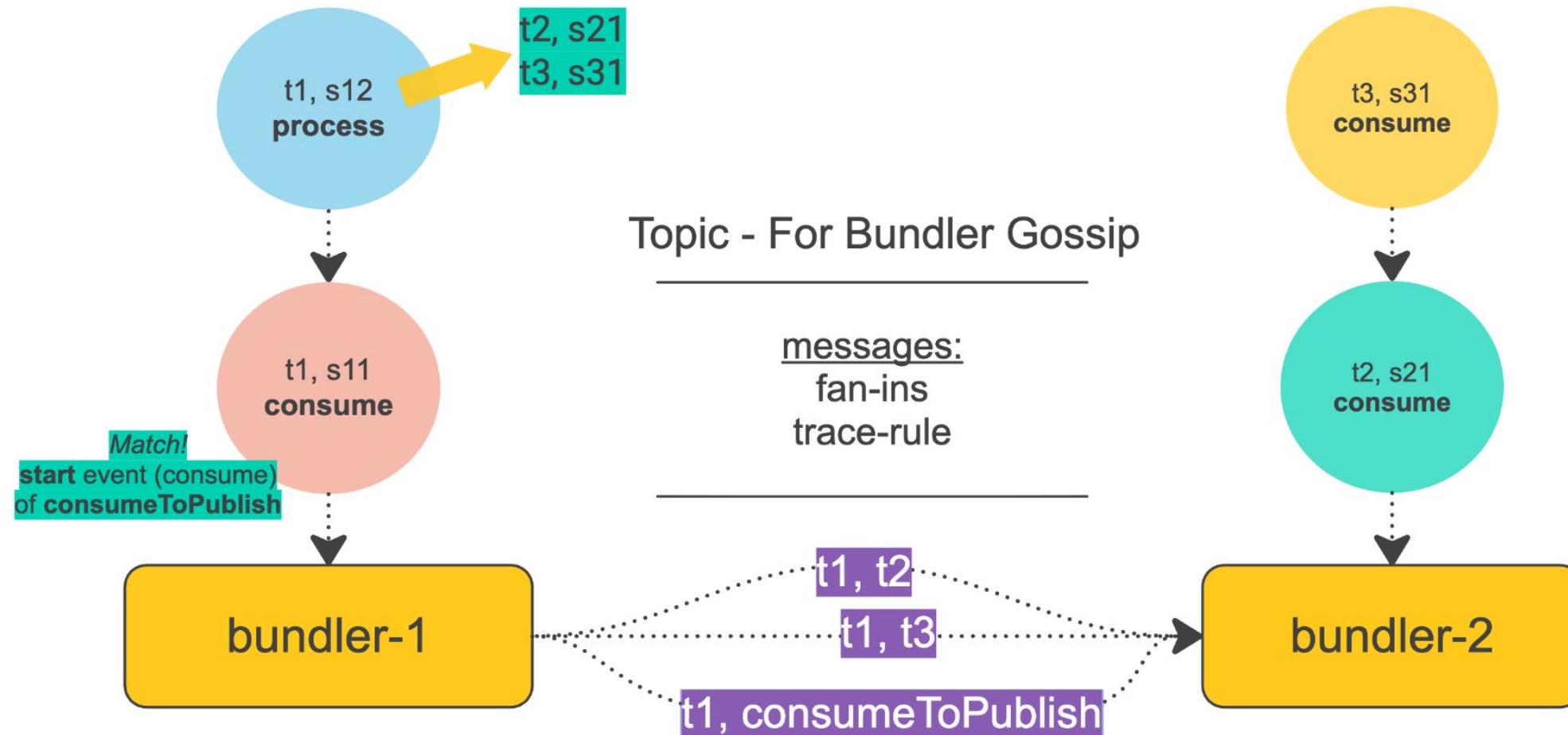
# Agenda: Digging Into the Two Approaches

- Designs for the approaches
  - Interaction between different C++ microservices
- C++ data structures used
- Pinpoint implementation-level bottlenecks
  - Use of profiling for improving performance
- Summarize pros and cons
- Key takeaways from this journey

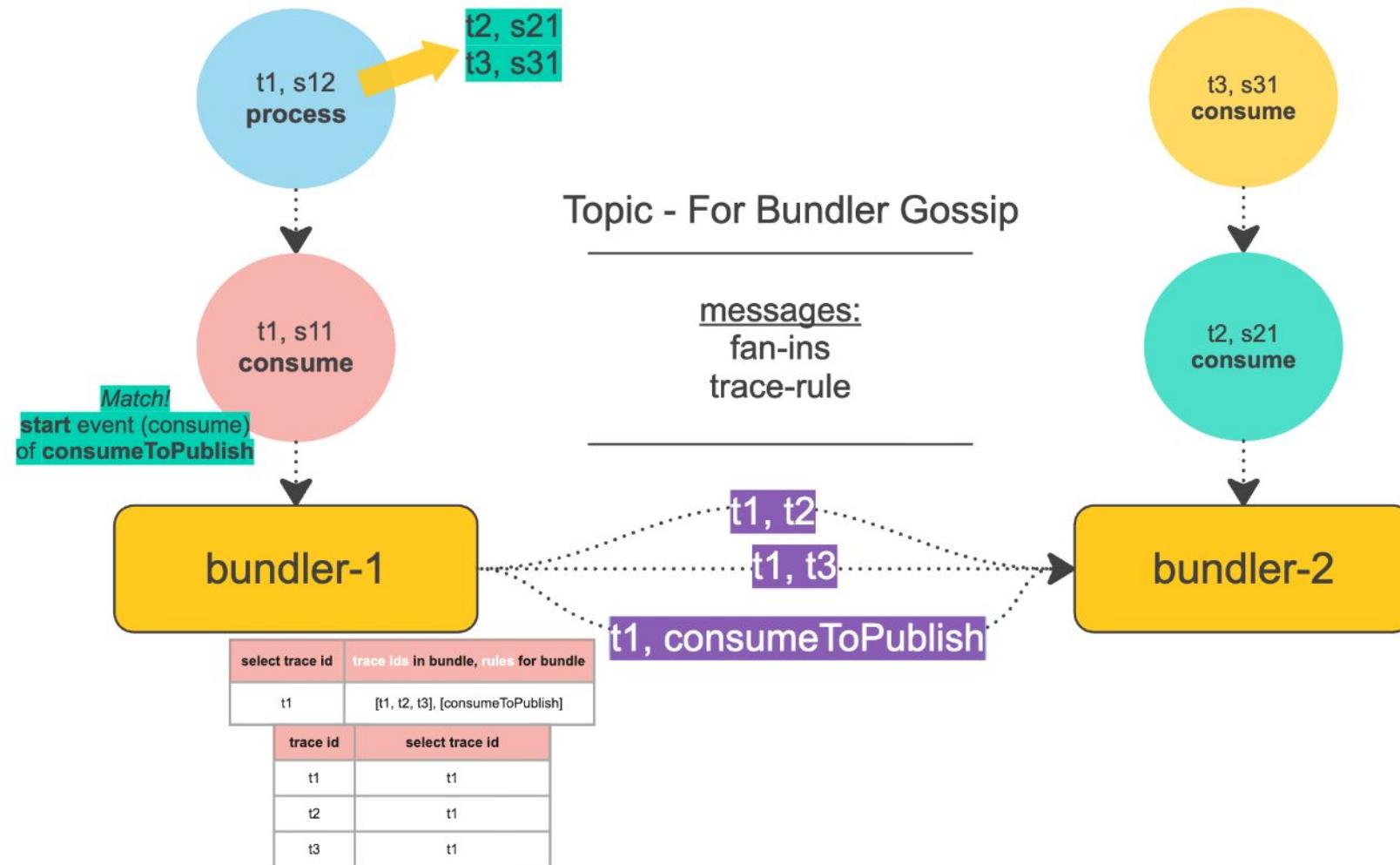
# Approach #1: Aggregating Trace "Bundles" By Rule



# Bundler for Merging Traces Bundles



# Bundler's Data Structures



# Zooming Into Bundler

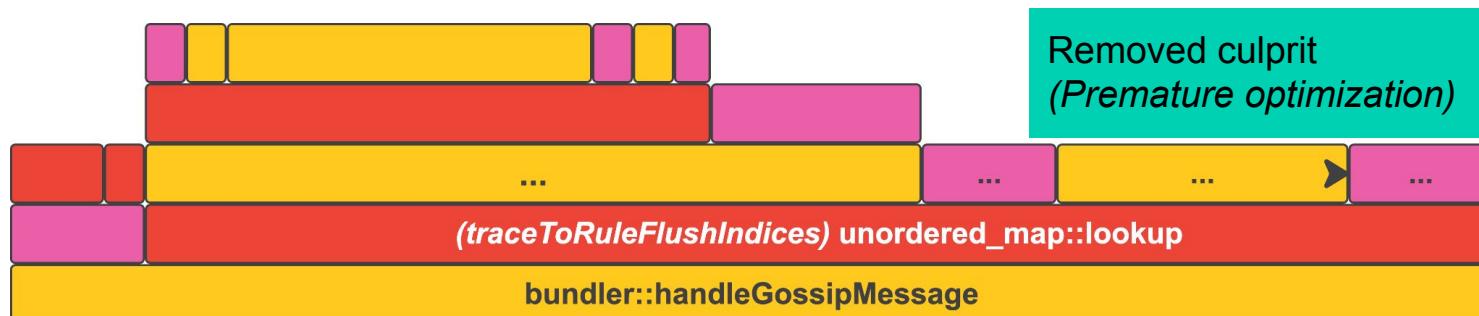
```
unordered_map<Traceld, Traceld> traceToBundleKey;
```

```
struct TracesRules {  
    unordered_set<string> rules;  
    unordered_set<Traceld> tracesInBundle;  
};
```

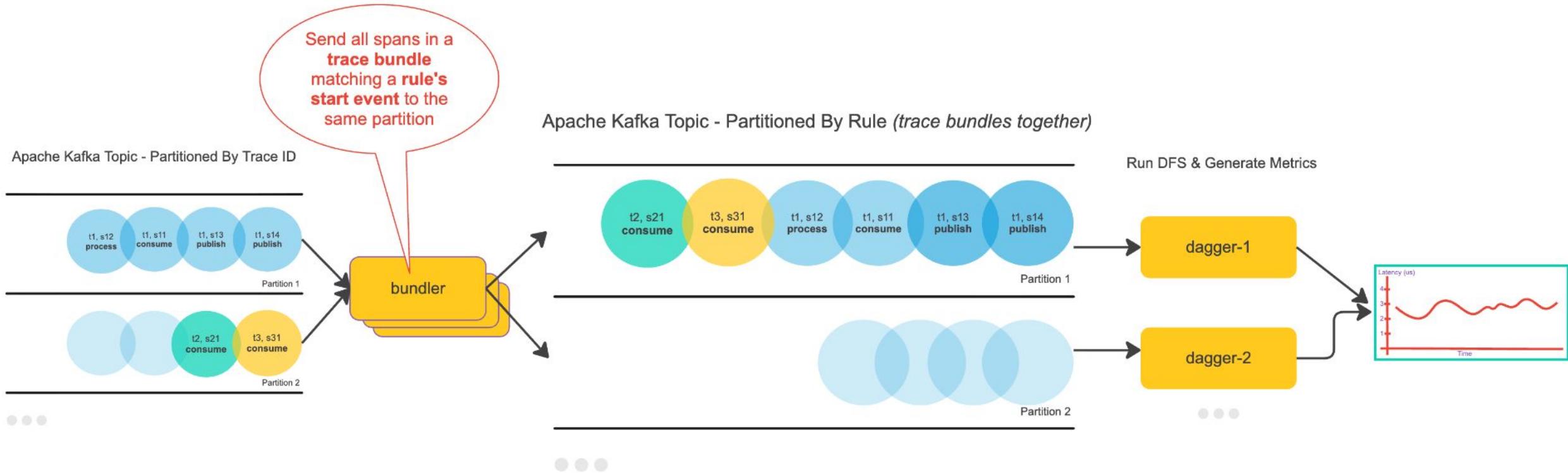
```
unordered_map<Traceld, TracesRules> bundleKeyToTracesRules;
```

# But, It Wasn't Always This Way!

- Handling the messages between *bundler's*
  - Slow: lag in Kafka consumption
- Previously, we had
  - `unordered_map<Traceld, Traceld>` `traceToBundleKey;`
  - `unordered_map<Traceld, unordered_set<Traceld>>` `bundleKeyToTraces;`
  - `unordered_map<Traceld, unordered_map<string, int>>` **`traceToRuleFlushIndices;`**
- Profiled handler
  - Culprit: Lookups/writes in **`traceToRuleFlushIndices`**



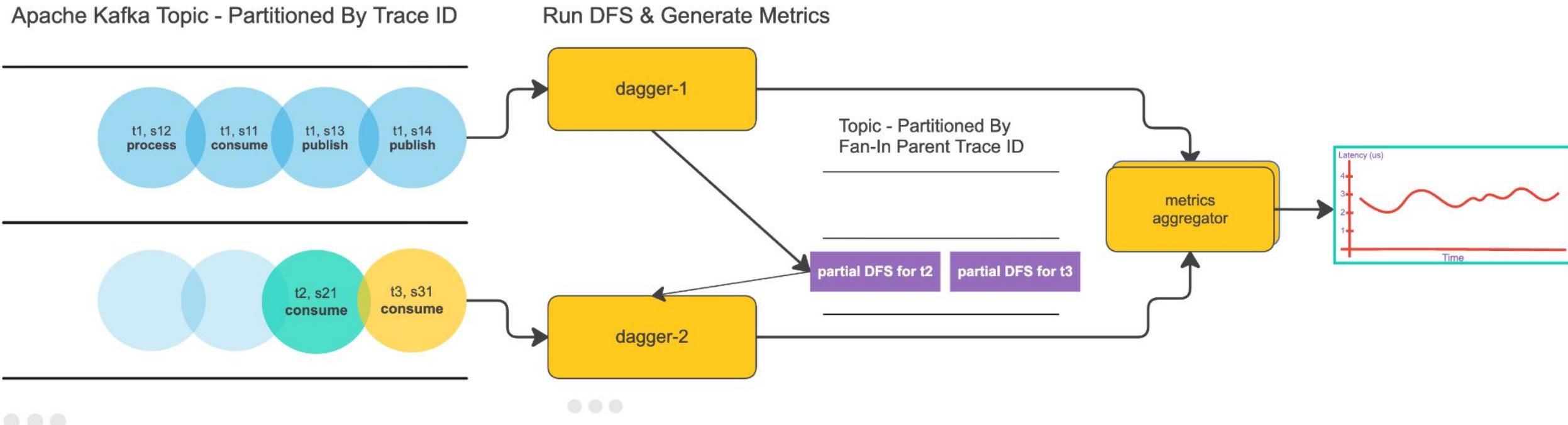
# Review: Aggregating Trace "Bundles" By Rule



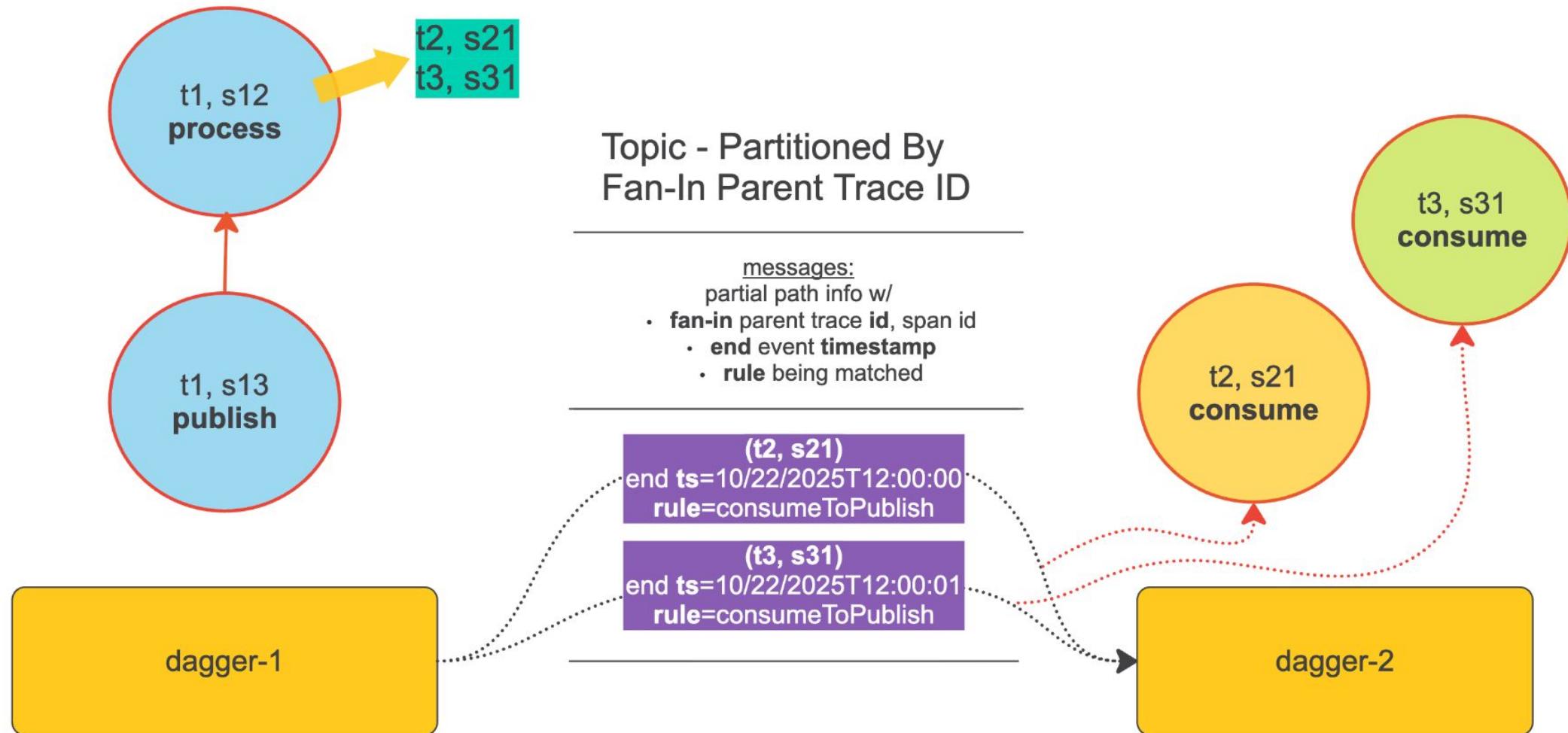
# Summarizing Approach #1

- Puts all spans from a trace bundle in the same partition
- Works well when rules match a small portion of traces
- Constraints:
  - Load imbalance
    - Rule with high volume trace matches bloat *dagger* memory
    - e.g., rules on middleware spans
  - Fan-in memory bloat
    - Fan-in traces not matching any rules exacerbate imbalance

# Approach #2: Handing Off "Partial-DFS"



# Gossip About Partial-DFS



# Zooming Into Dagger

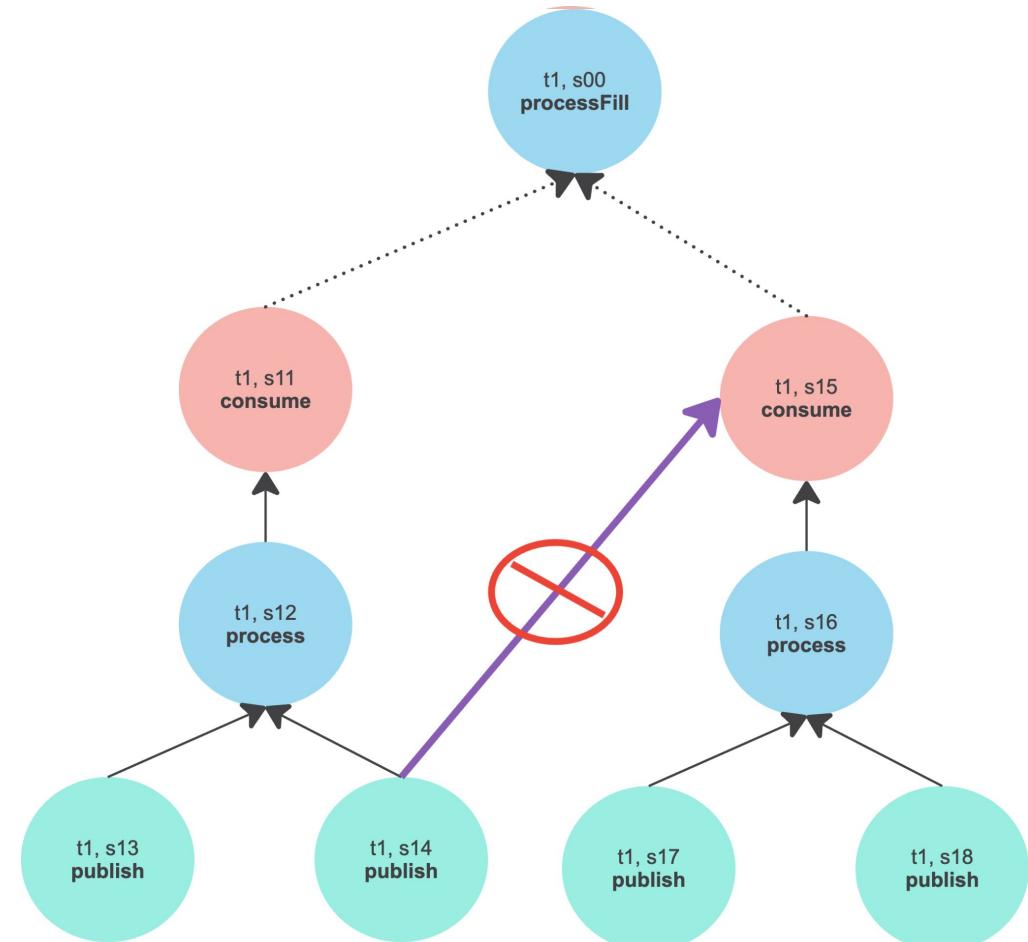
```
struct TraceData {  
    unordered_map<SpanId, Span>           spans;  
    vector<string>                         rules;  
};
```

```
unordered_map<Traceld, shared_ptr<TraceData>>  tracesToProcess;
```

```
unordered_map<Traceld, FanInData>      partialDfsToProcess;
```

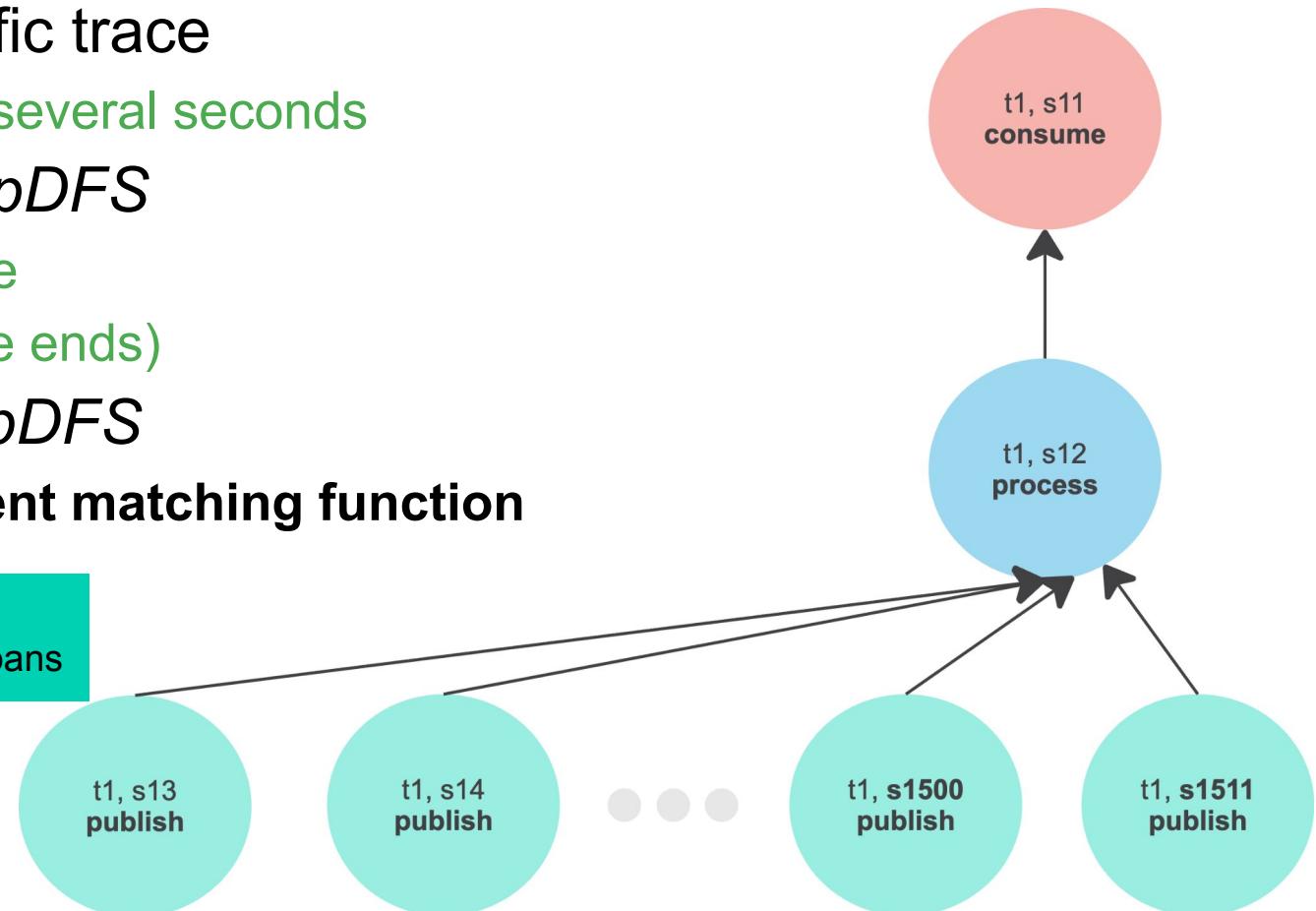
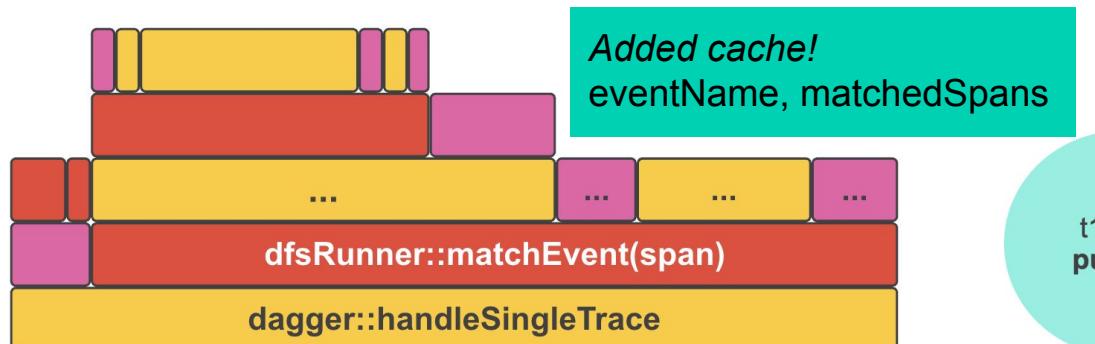
# But, Processing Some Traces Took Long!

- Running the DFS over a specific trace
  - Processing time occasionally in several seconds
- Special user configuration: *skipDFS*
  - Start and end event in the same trace
  - NOT(multiple starts && multiple ends)



# But, Processing Some Traces Took Long!

- Running the DFS over a specific trace
  - Processing time occasionally in several seconds
- Special user configuration: *skipDFS*
  - Start/end event in the same trace
  - NOT(multiple starts AND multiple ends)
- Profiled DFS execution for *skipDFS*
  - Culprit: repetitive calls to the **event matching function**



## **Summarizing Approach #2**

- Lazy approach (as opposed to #1, greedy)
- Send info regarding fan-ins only when relevant to DFS
- Avoids load imbalance
- Constraints:
  - Deep, chained fan-ins
    - Consecutive fan-ins in one path incur gossip delays

# Conclusion

End-to-end latency metrics from complex traces can be generated at scale.

- Pick the best approach for you
- Design for speed from the start
  - Avoid premature optimizations
  - Leverage tooling to iteratively improve code performance
- End-to-end latency monitoring is crucial
  - Trace provides a unified solution
  - Client impact can be directly monitored
- Distributed, streaming DAGs are fun!

Bloomberg

Engineering

# Thank you!

Let's connect:

[kmaharshi@bloomberg.net](mailto:kmaharshi@bloomberg.net)

[www.linkedin.com/in/kusha-maharshi](https://www.linkedin.com/in/kusha-maharshi)

Learn more: [www.TechAtBloomberg.com/cplusplus](https://www.TechAtBloomberg.com/cplusplus)

We are hiring: [bloomberg.com/engineering](https://bloomberg.com/engineering)

TechAtBloomberg.com

# Using Distributed Trace for End-to-End Latency Metrics

Engineering

Bloomberg

CppCon 2025  
September 18, 2025

Kusha Maharshi  
Senior Software Engineer, Telemetry Infrastructure

TechAtBloomberg.com