

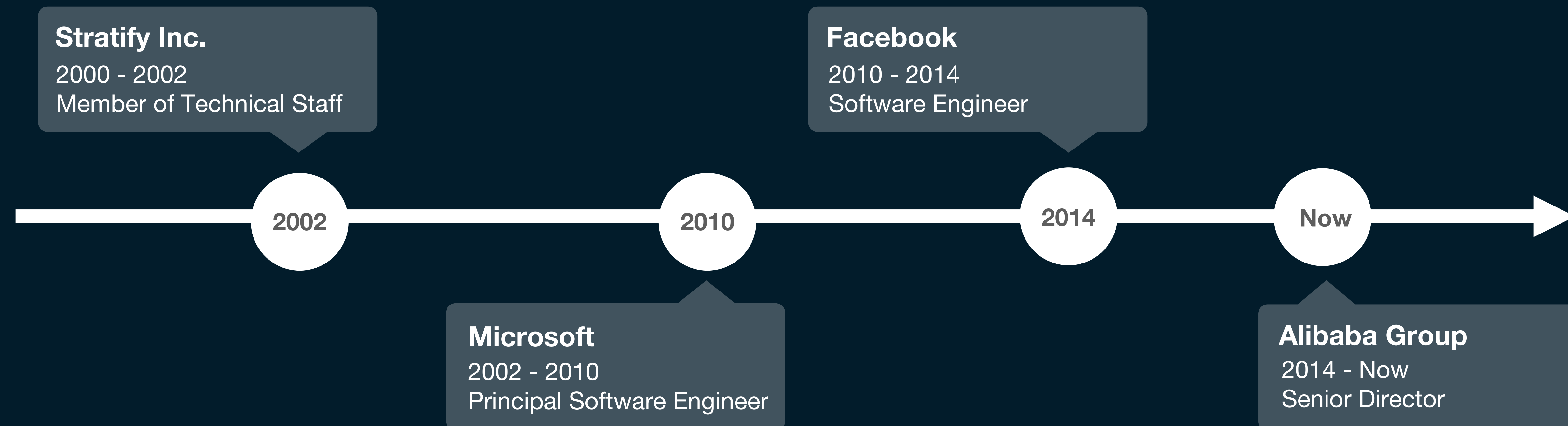
Apache Flink[®] - Redefining Computation

Xiaowei Jiang

Senior Director

Alibaba





FLINK
FORWARD

About Alibaba



EB Total

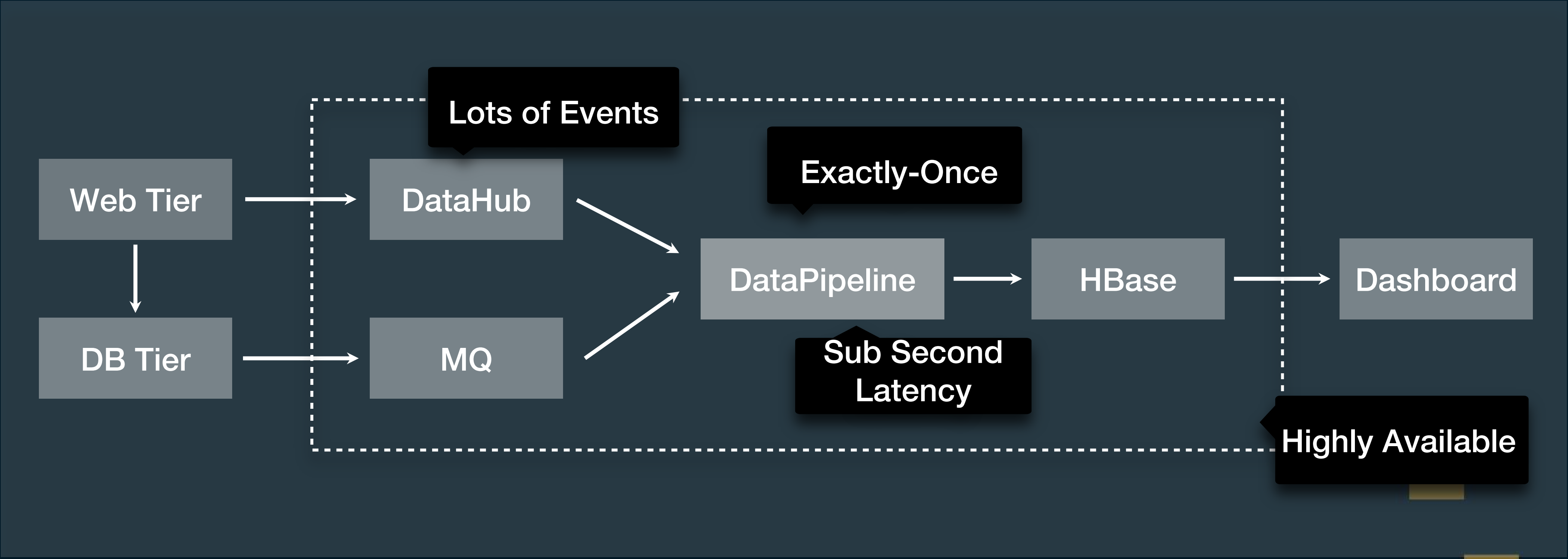
PB Everyday

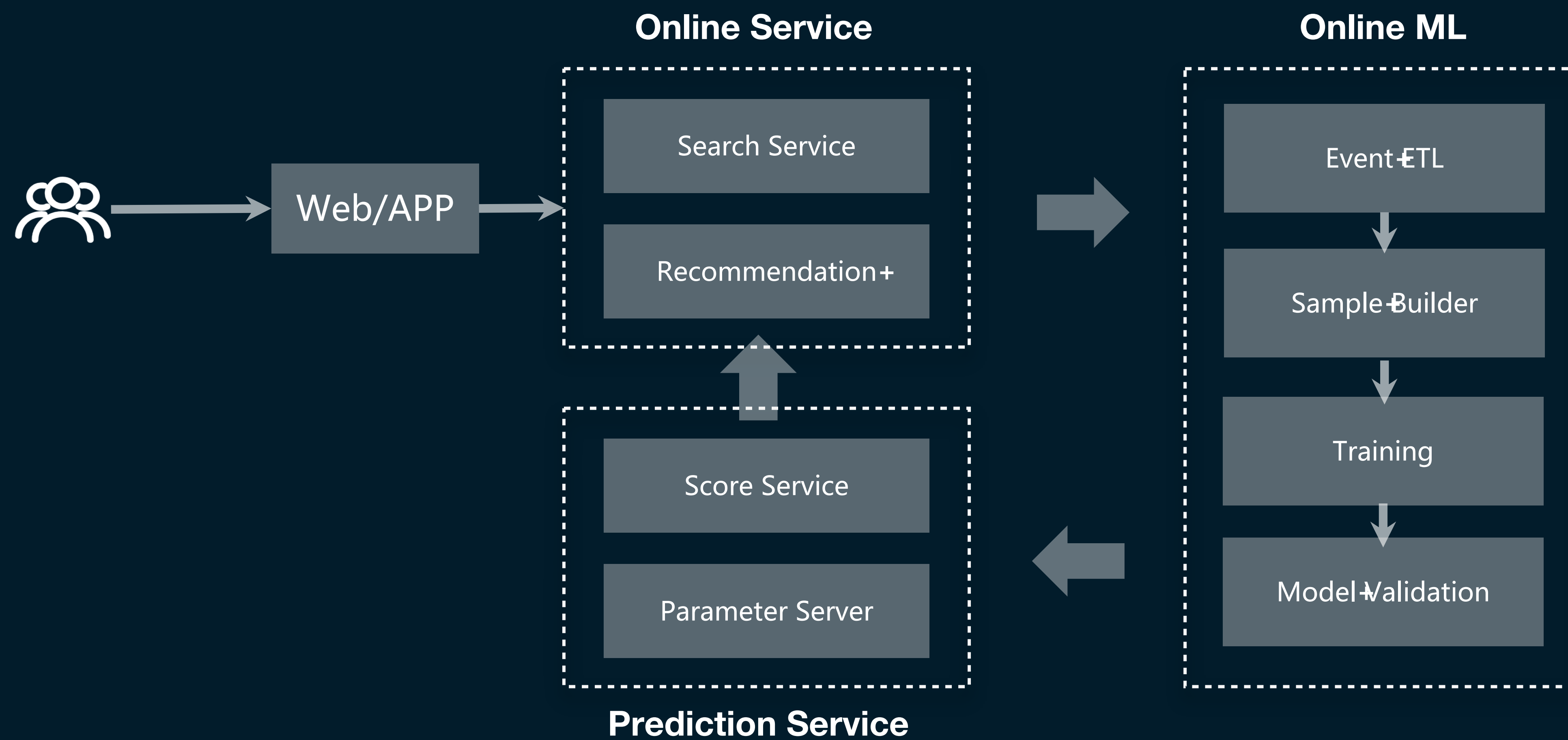
1T Event/Day

1.7B Events/sec



1. Stream Processing





Real-Time Personalization

Large Scale: 100 Millions of Events, 100 Billions of Features

Low Latency: Second Latency from End to End

Complex Logic: Real-Time Training, Feature/Model Update

What is Flink



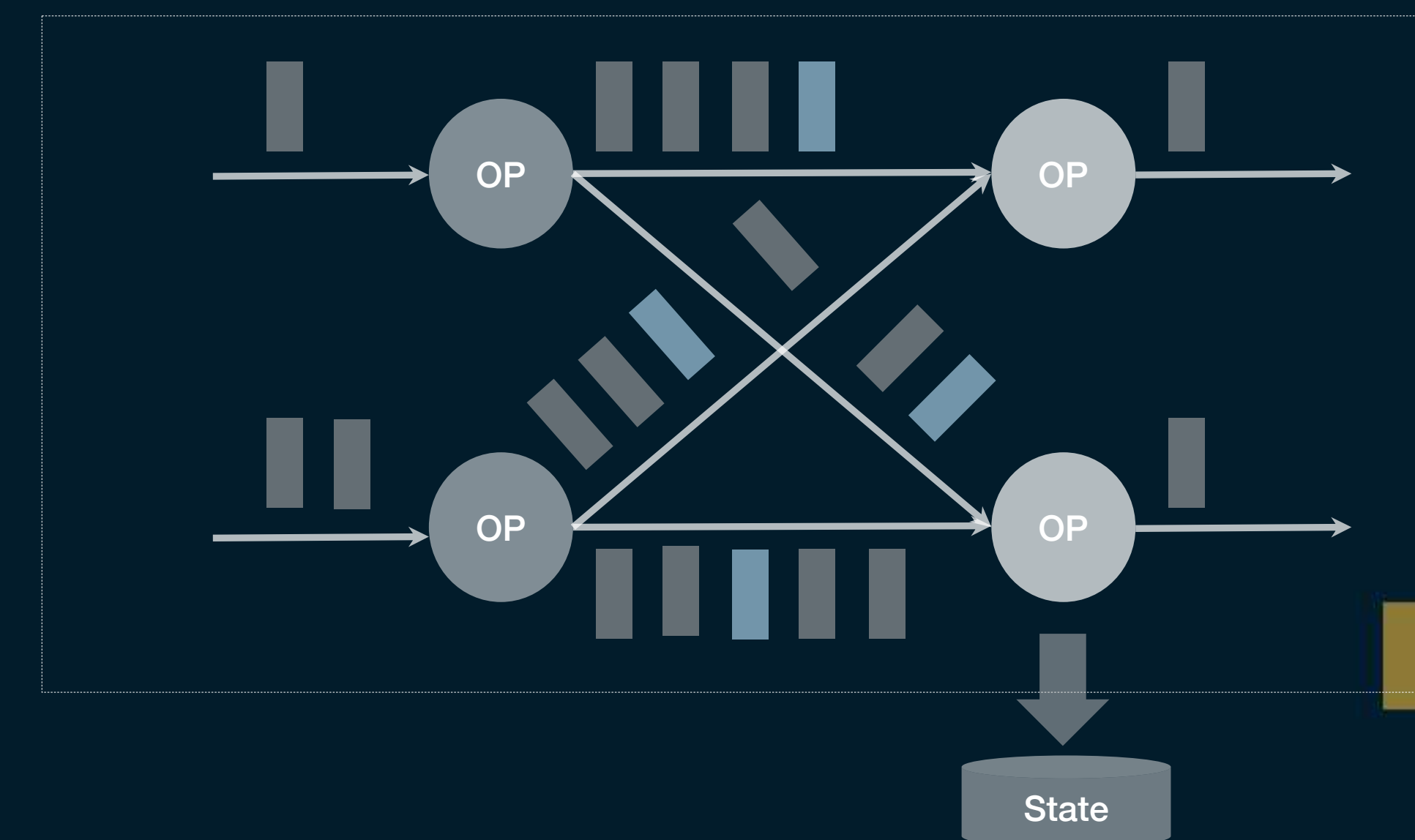
An open source stream processing framework that unifies real-time event-driven applications and real-time analytics.

Flink Program

```
> _
```

```
SELECT
  word, count(*)
FROM
  stream
GROUP BY
  word;
```

Flink DAG





Flink Stream Process



Exactly-Once



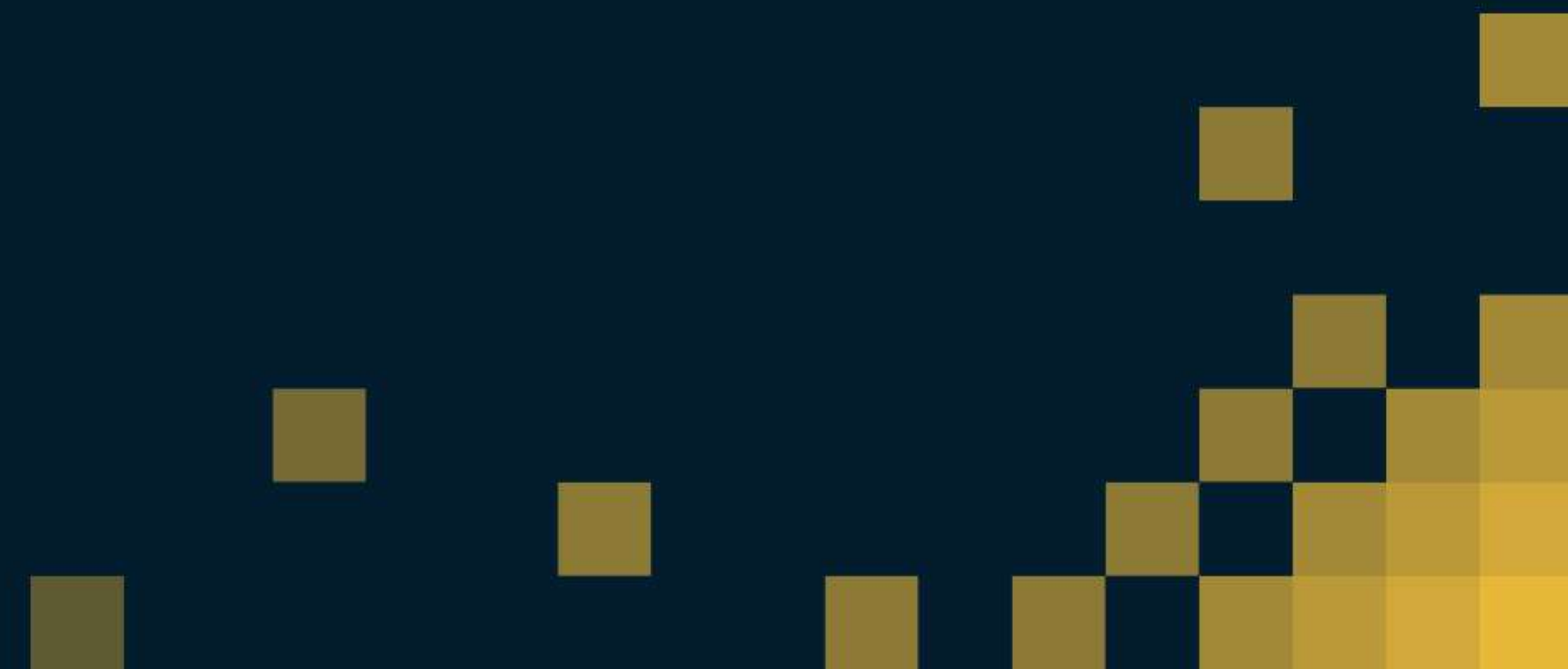
High Throughput

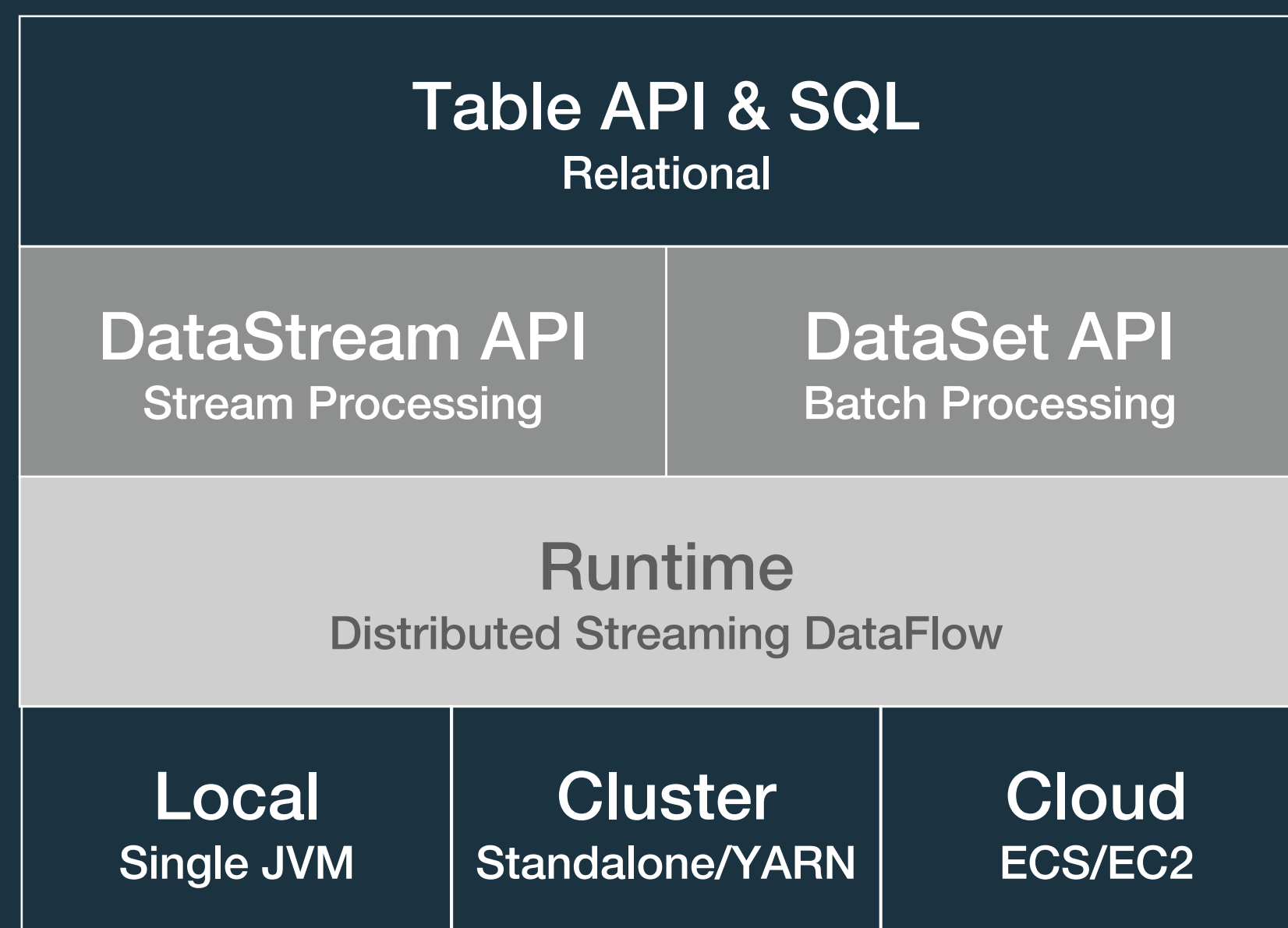


Low Latency



Fault Tolerant





Flink Architecture

Flink Runtime Improvements

Distributed Architecture

Rework Cluster Management [FLIP-6/FLINK-4319]

Fault Tolerance

JobManager Failover [FLINK-4911]

Region-based Task Failover [FLIP1/FLINK-4256]

Performance

Incremental Checkpoint [FLINK-5053]

Async I/O [FLIP12/FLINK-4391]

Credit-based Flow Control [FLINK-7282]



Flink SQL Improvements

Semantics

Functionality

Agg/w Retraction

Window

UDX Support

DDL Support

Connector Support





Status of Production



**Subsecond
Latency**



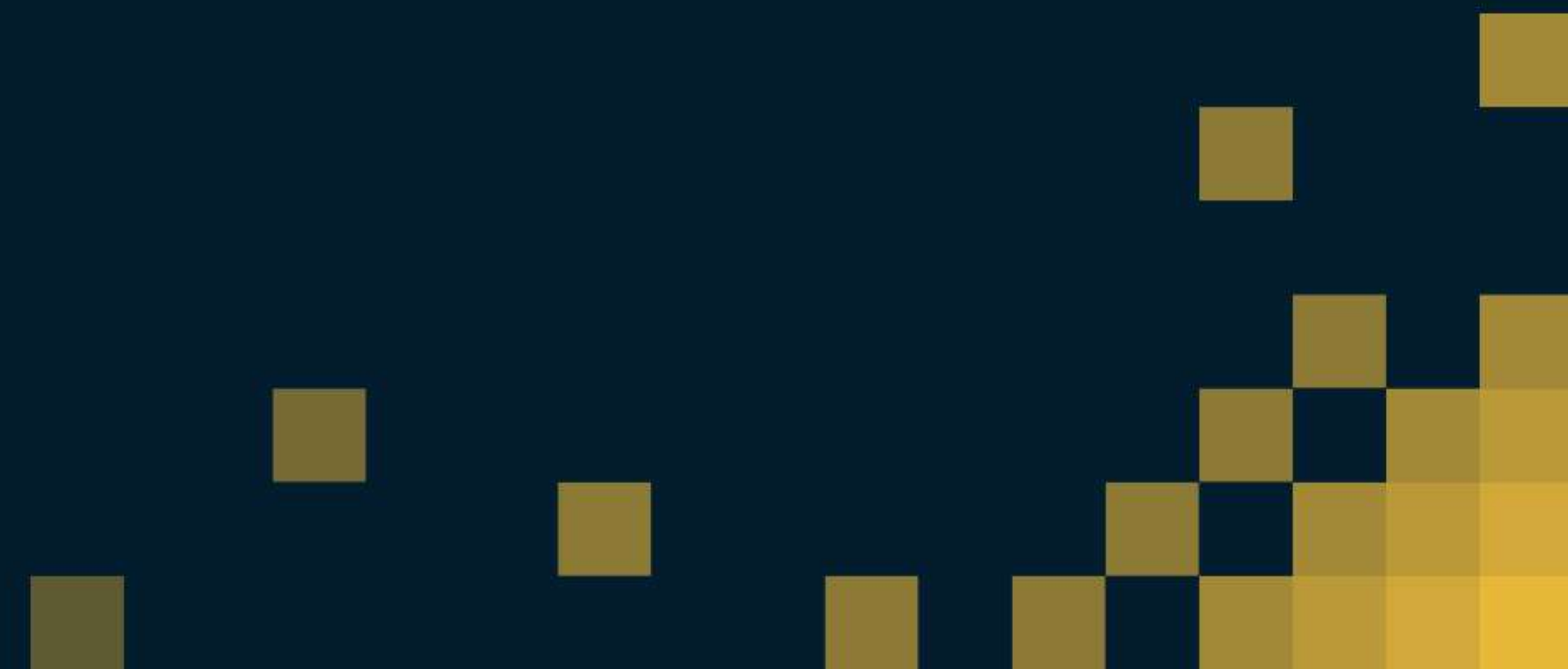
**Billions
Events/sec**



**Ten Thousand
Nodes**

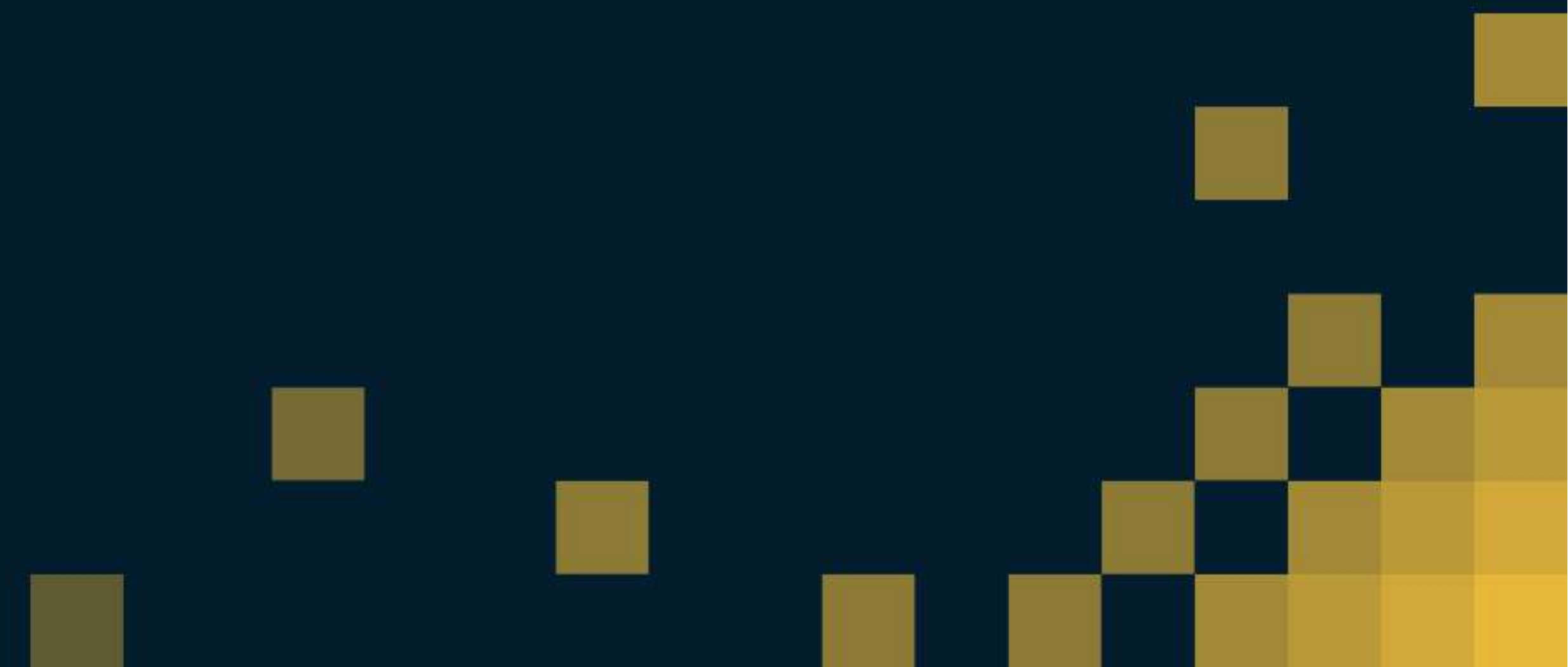


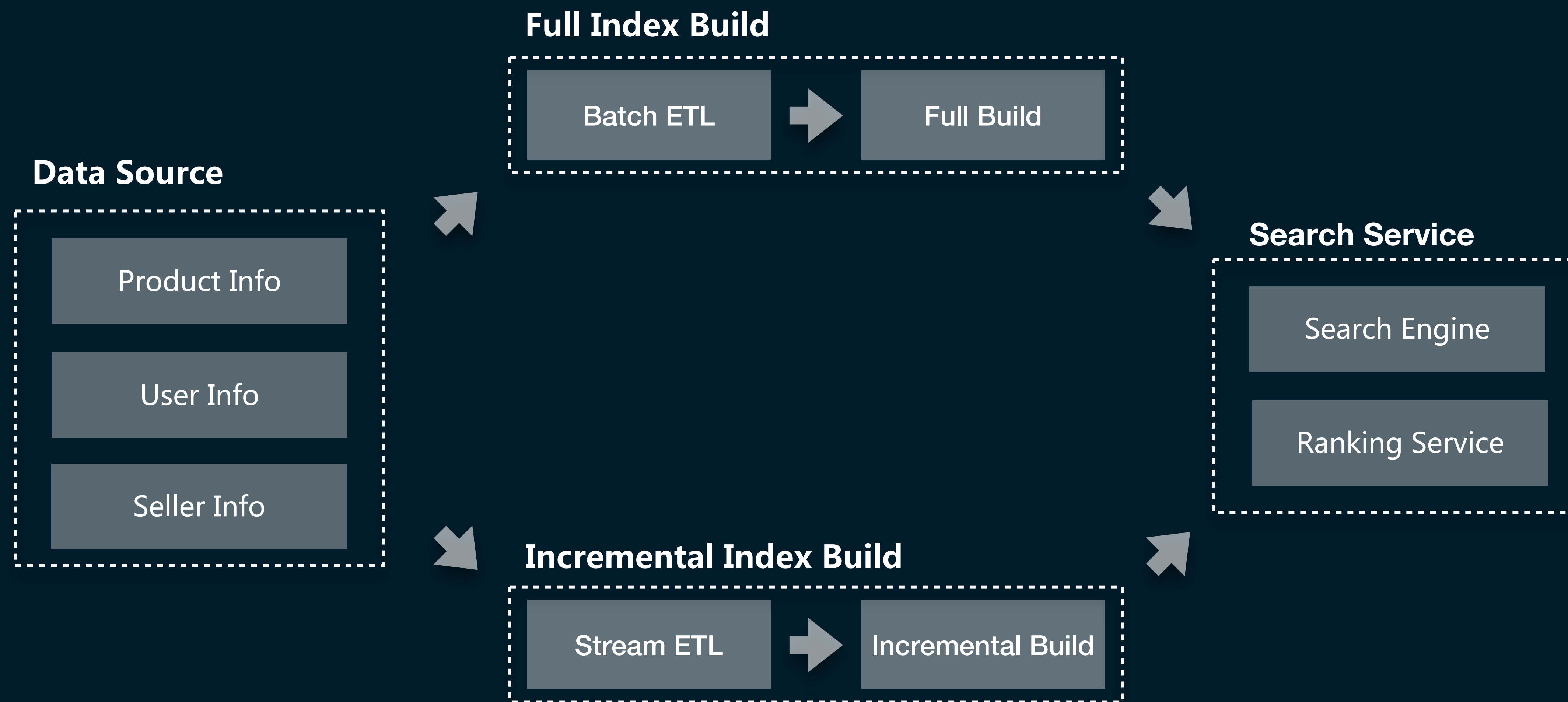
**Ten Thousand
Jobs**





2. Unified Engine





Index Pipelines for Search

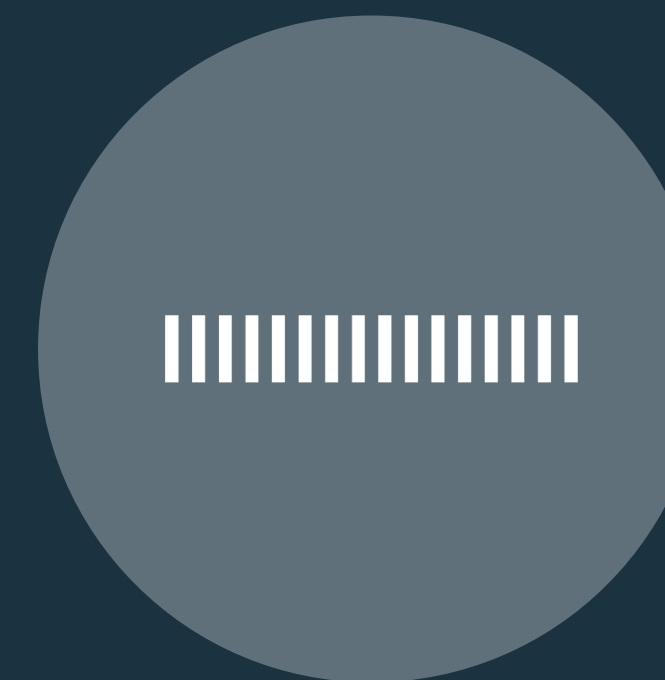
Development Efficiency: Full/Incremental Index Build

Challenge: Consistency



**Low Latency
Fixed Query**

Stream Processing



**Periodic/Continuous
Batch Jobs**

Progressive Processing

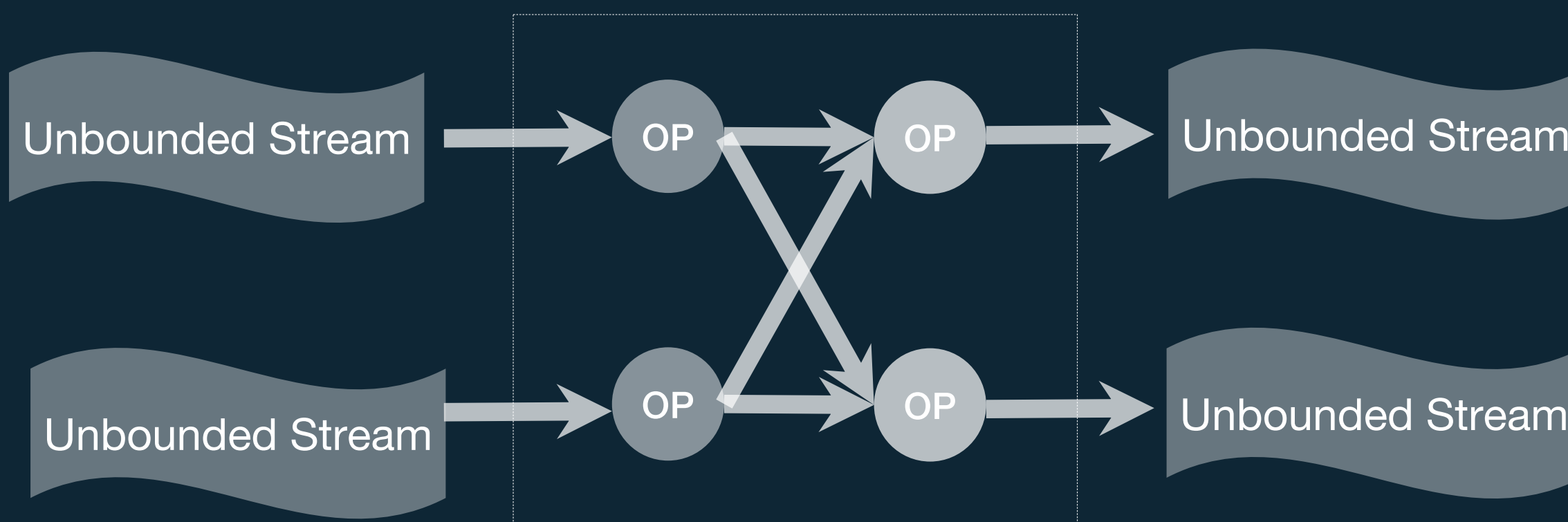


**High Throughput
Flexible Query**

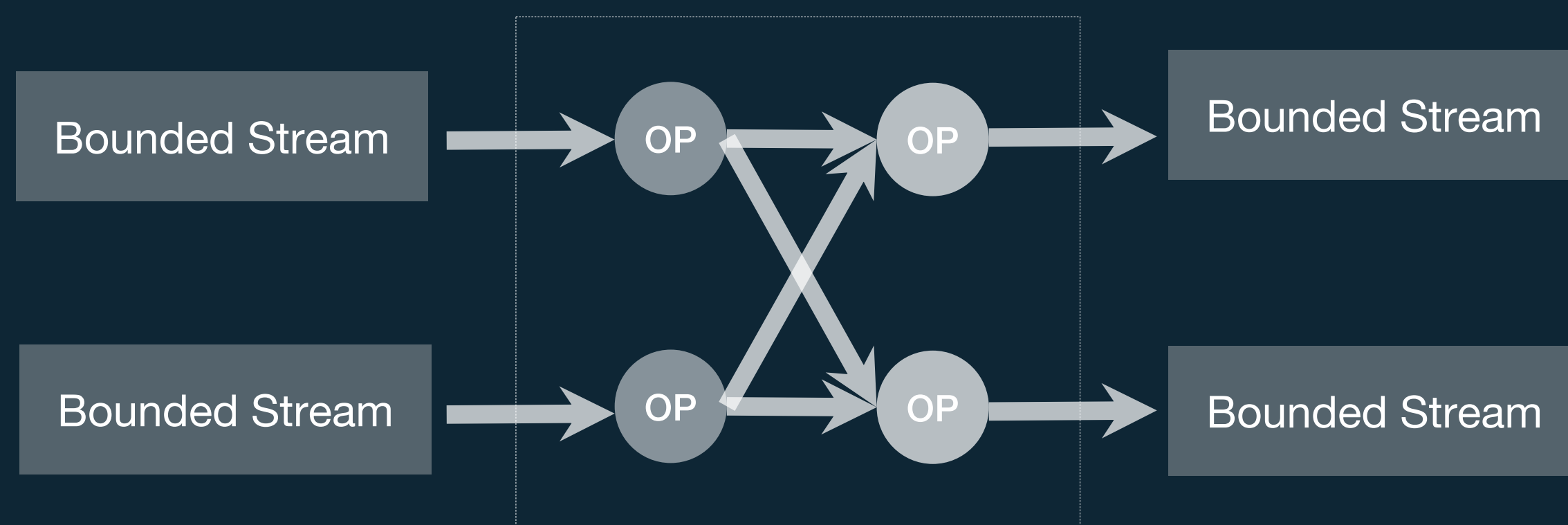
Batch Analytics

Streaming as the **core abstraction**, Batch as a **special case** of streaming

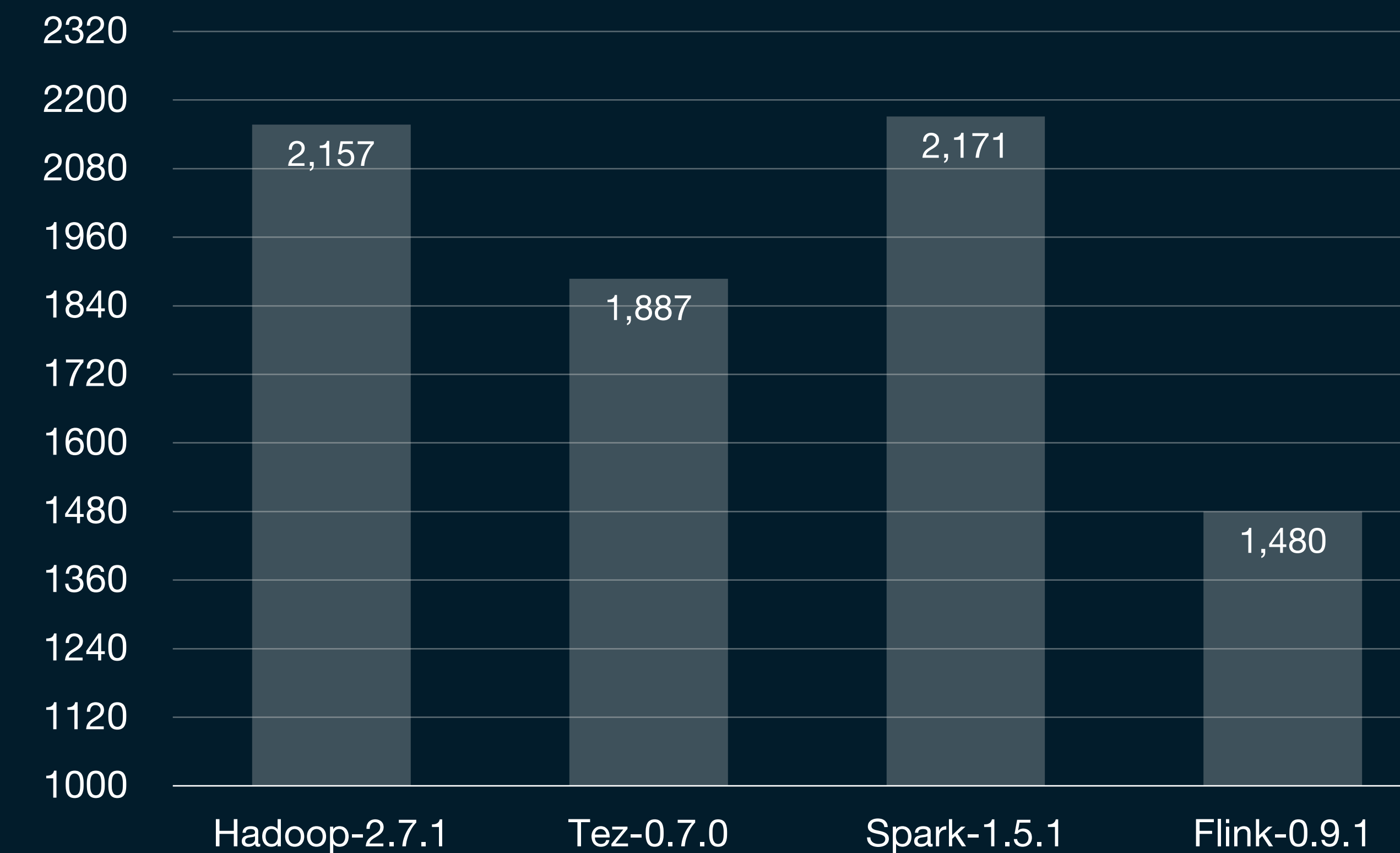
Stream Job



Batch Job



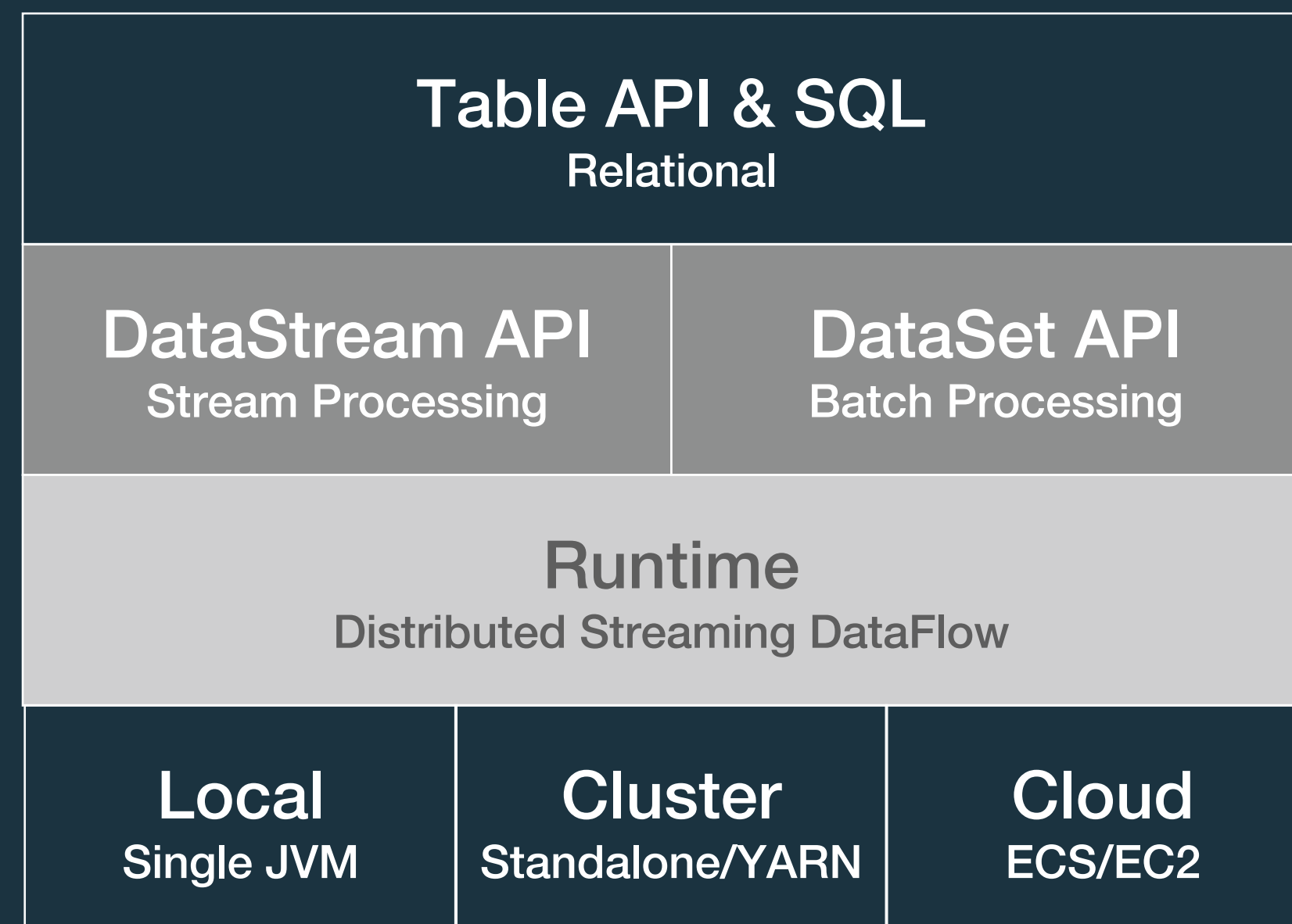
Result of sorting 80GB/node (3.2TB)



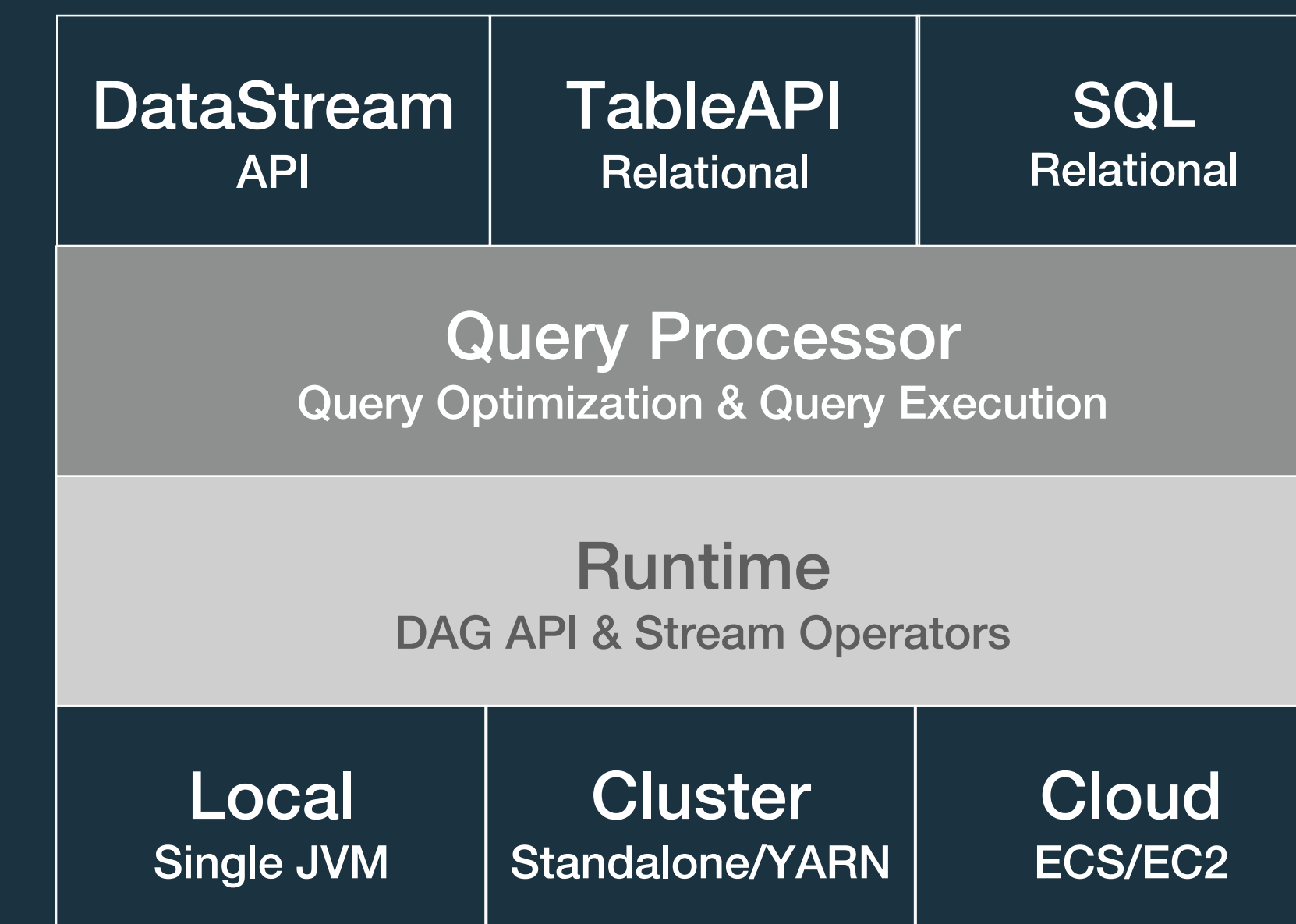
Flink is the fastest due to its pipelined execution

Tez and Spark do not overlap 1st and 2nd stages

MapReduce is slow despite overlapping stages



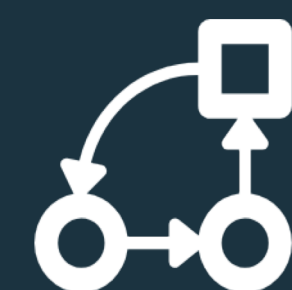
Old Design



New Design



Improvements in New Design



Runtime

New Operator Framework
Customizable Scheduling
Flexible Chaining



Query Execution

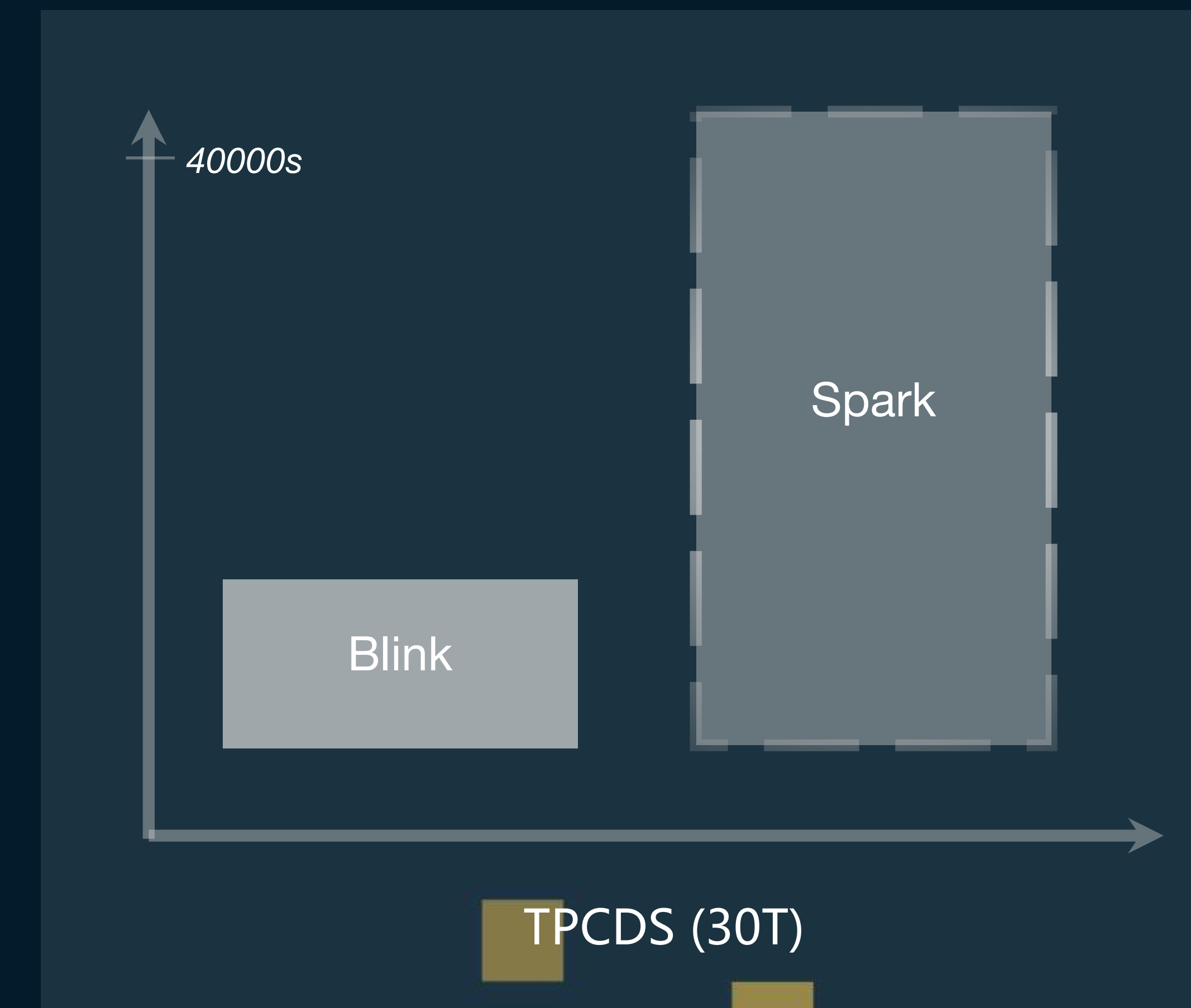
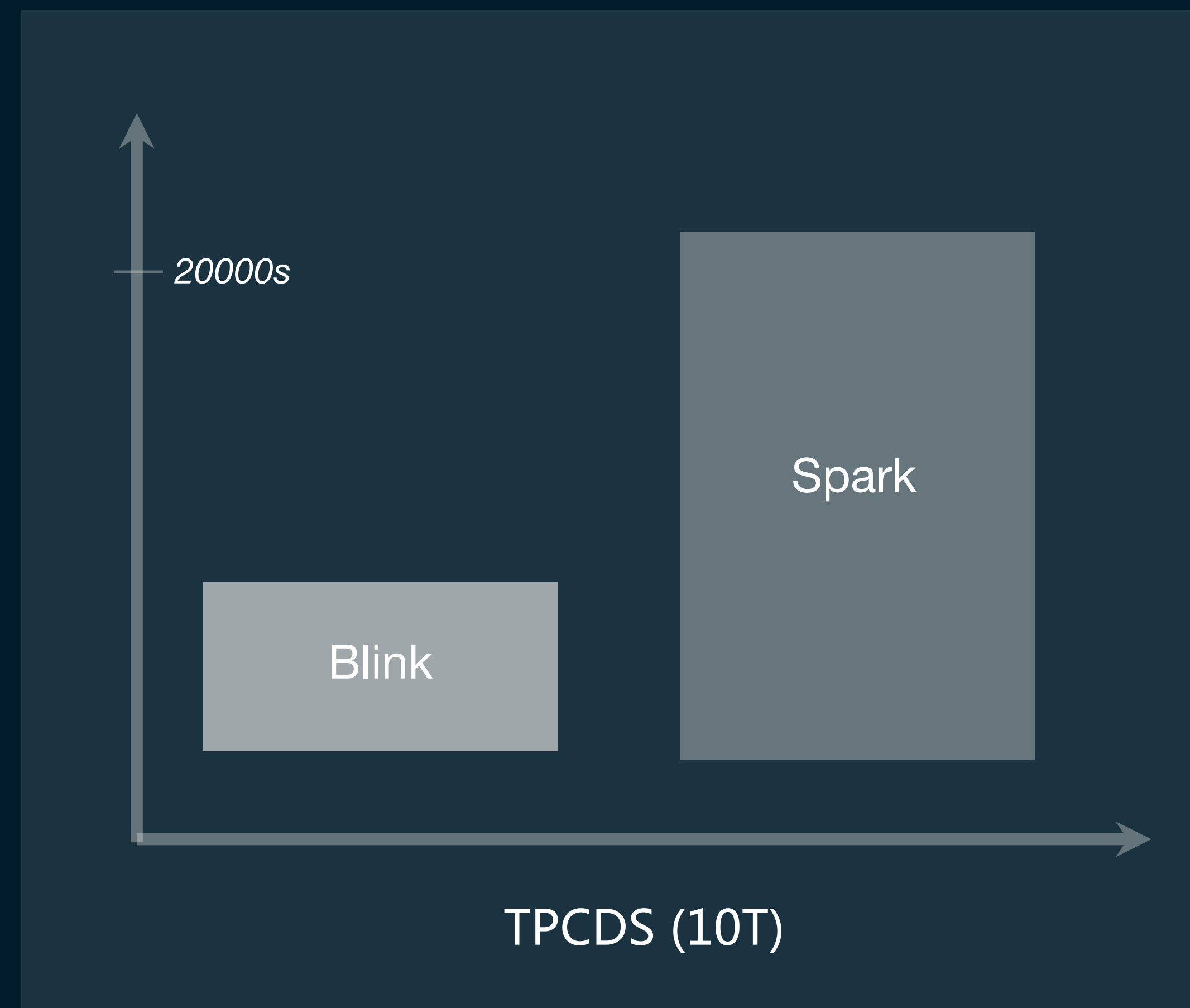
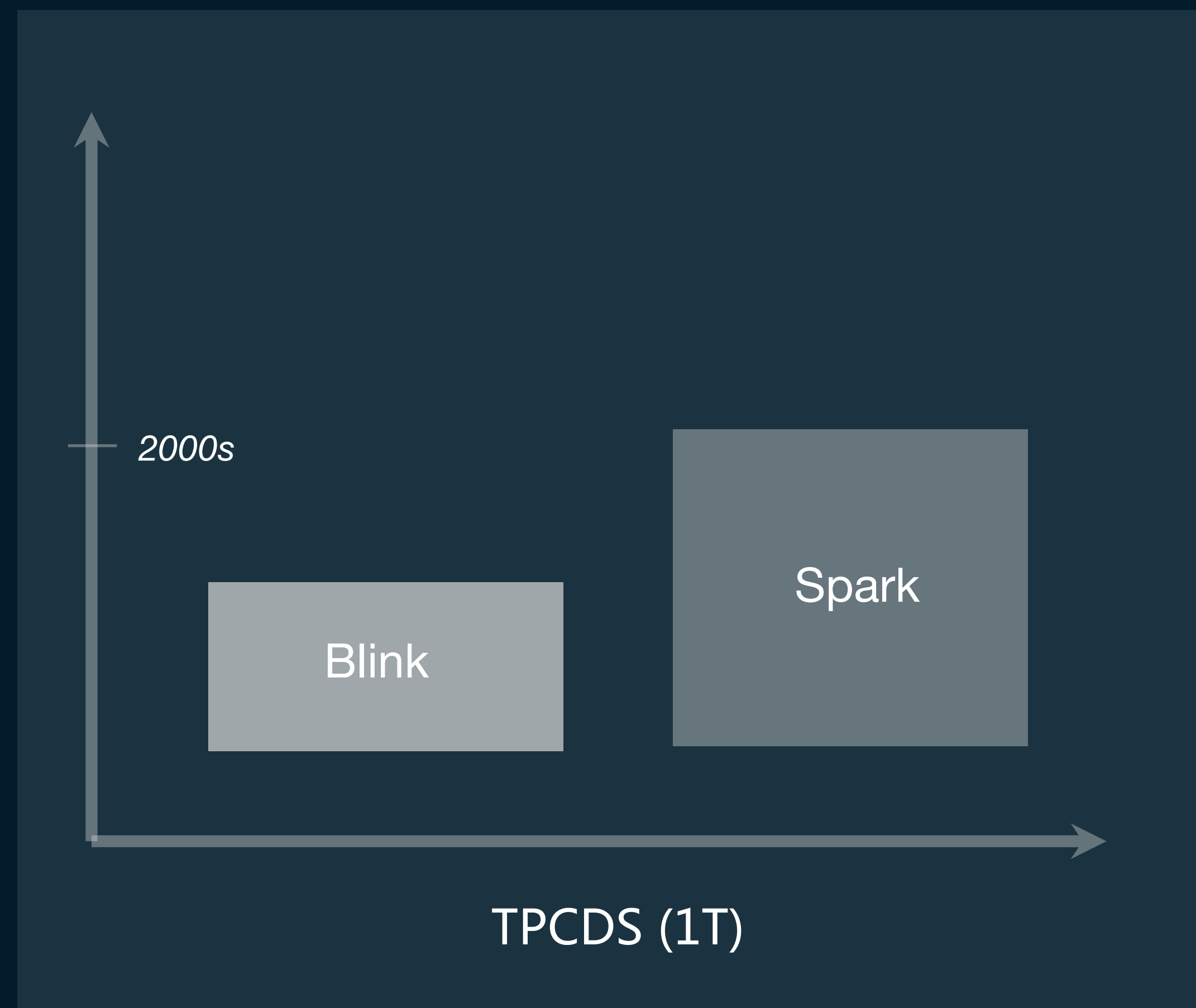
Expression Optimizations
Performant Operators
Resource Optimizations

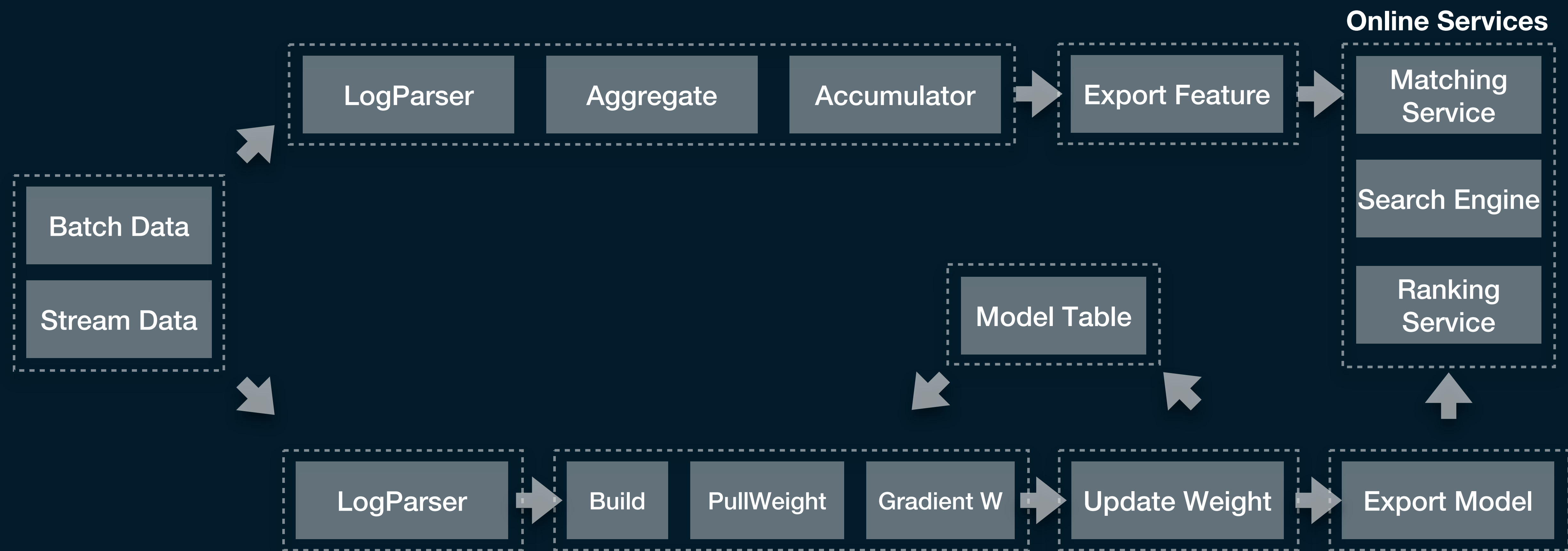


Query Optimizer

Cost Based
Advanced Rules
Rich Stats

TPC-DS Performance *(the Lower, the Better)*





Search' s Algorithm Platform

Unified Pipeline for Batch & Streaming
 Streaming: 100M QPS, 100B features
 Batch: Over 400TB in a single job



Flink

+



Hive



Flink

+

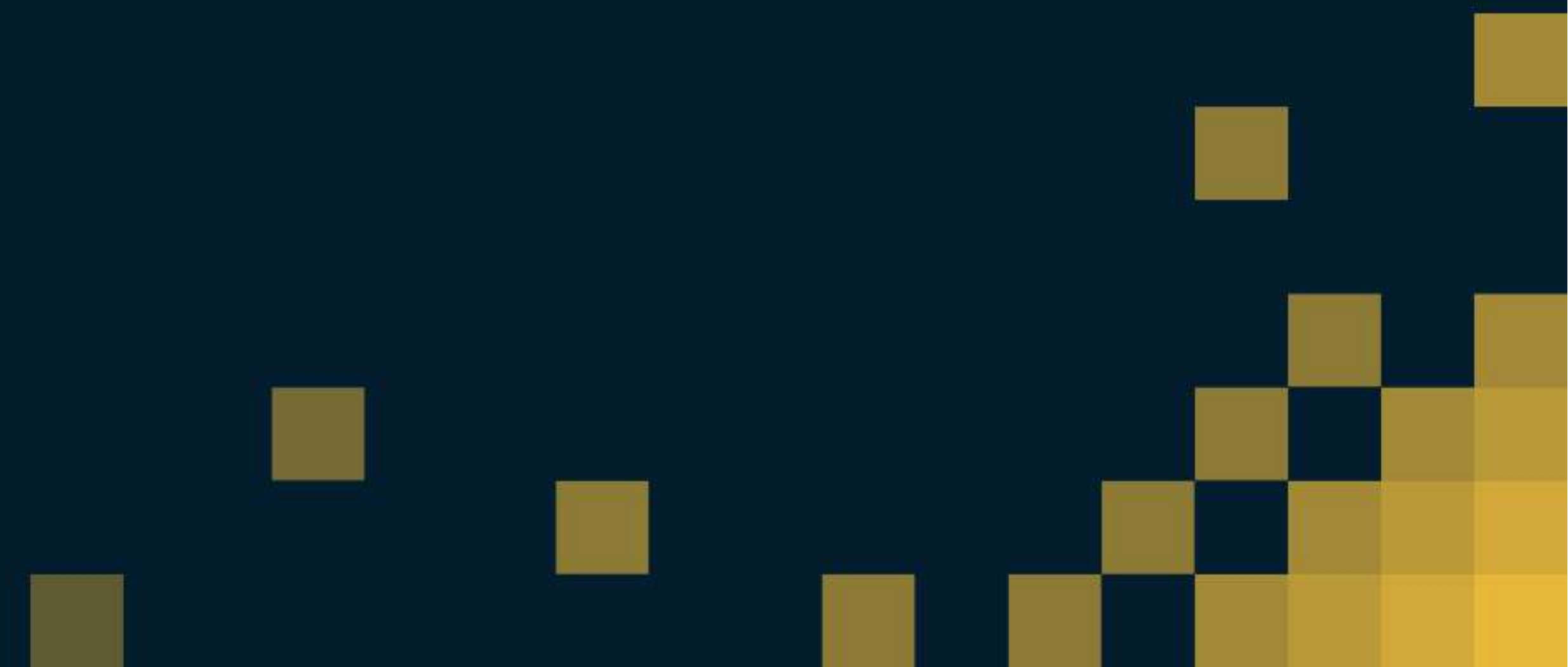


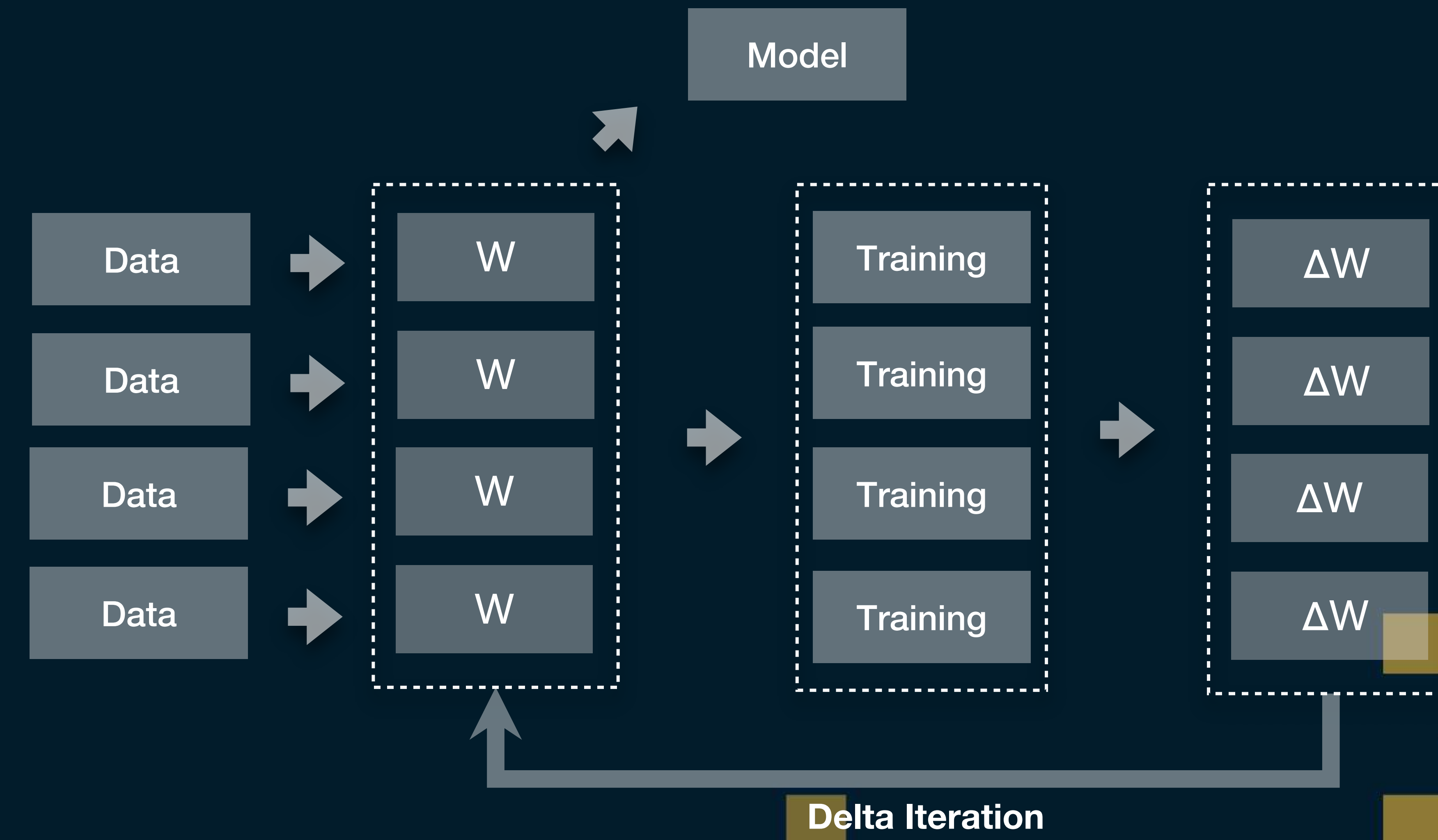
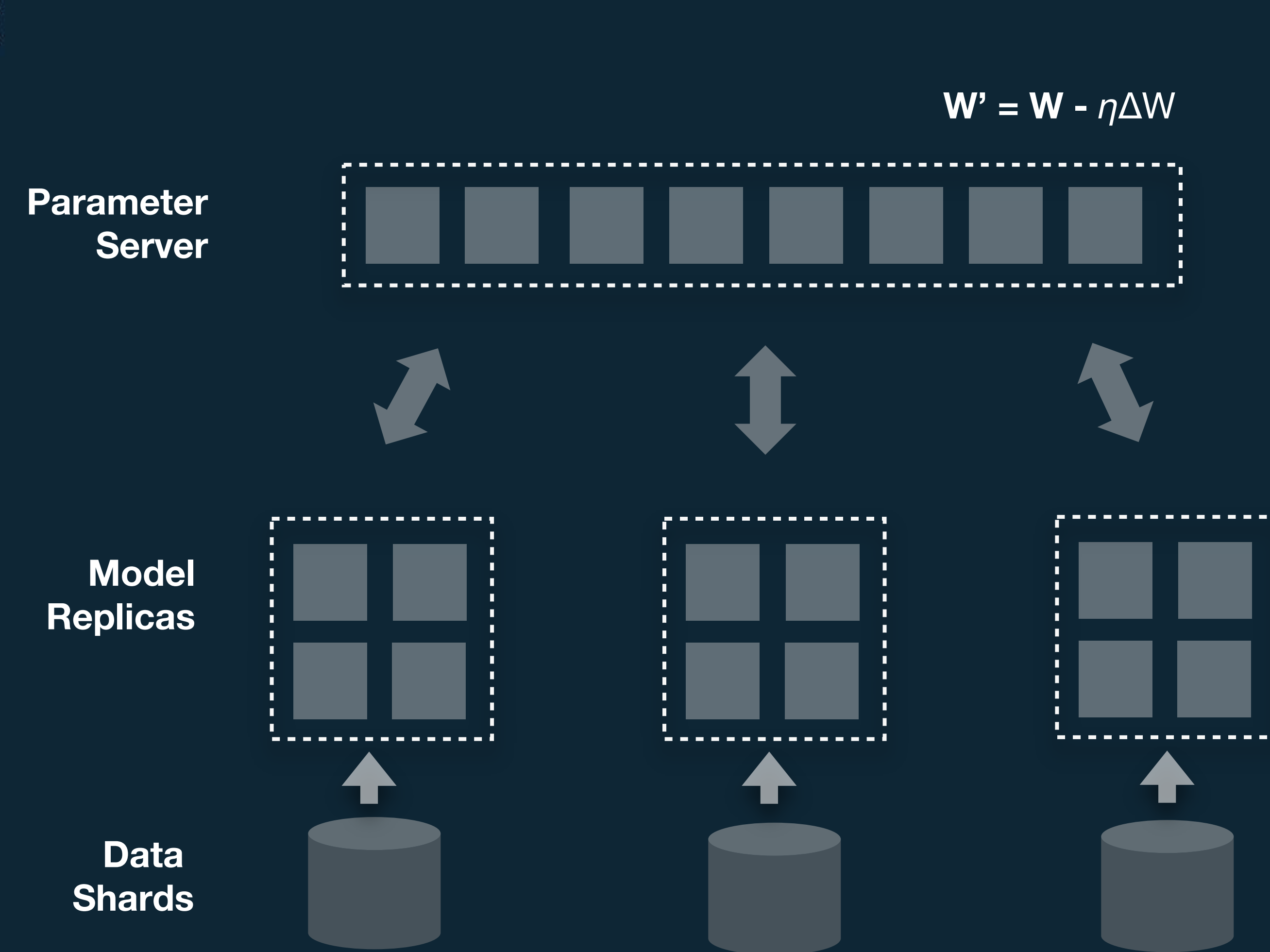
Zeppelin

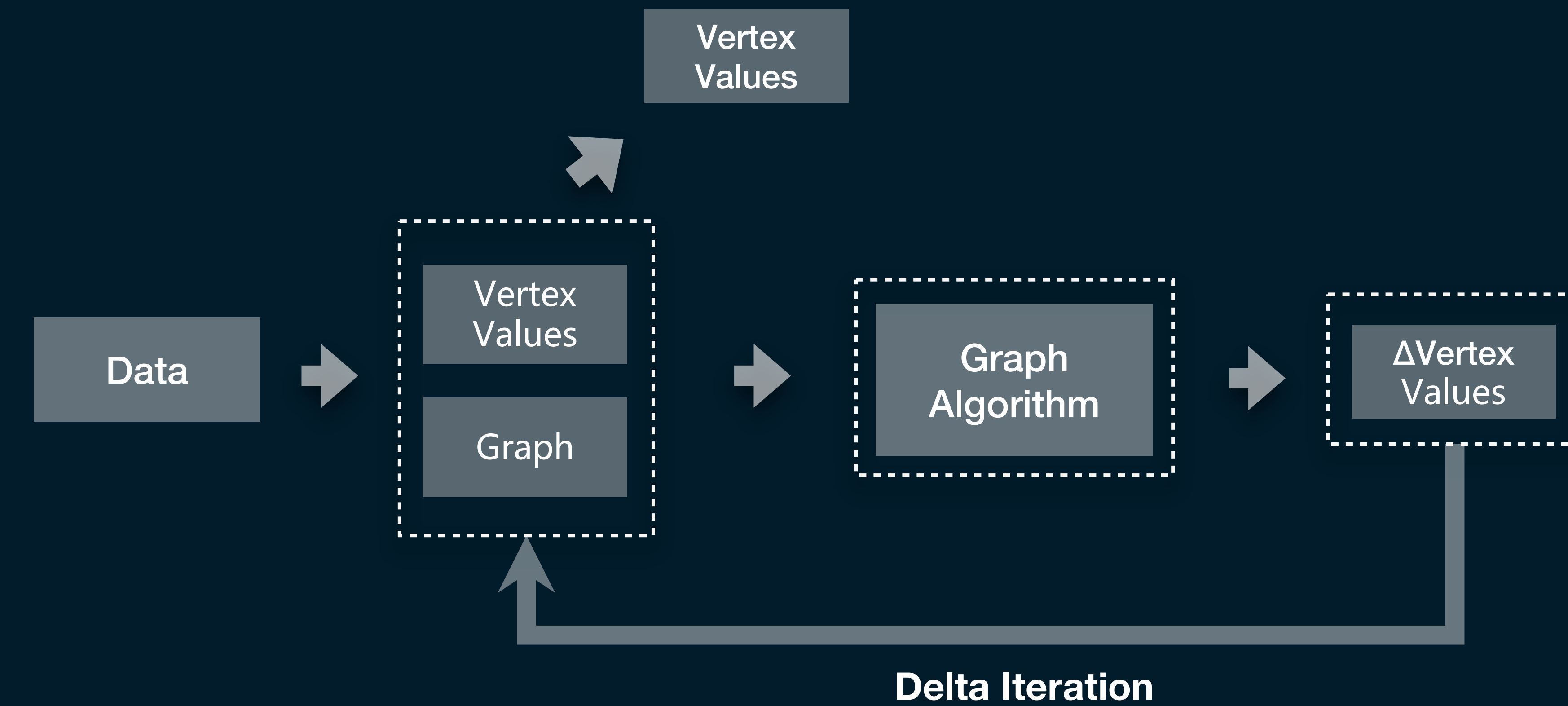
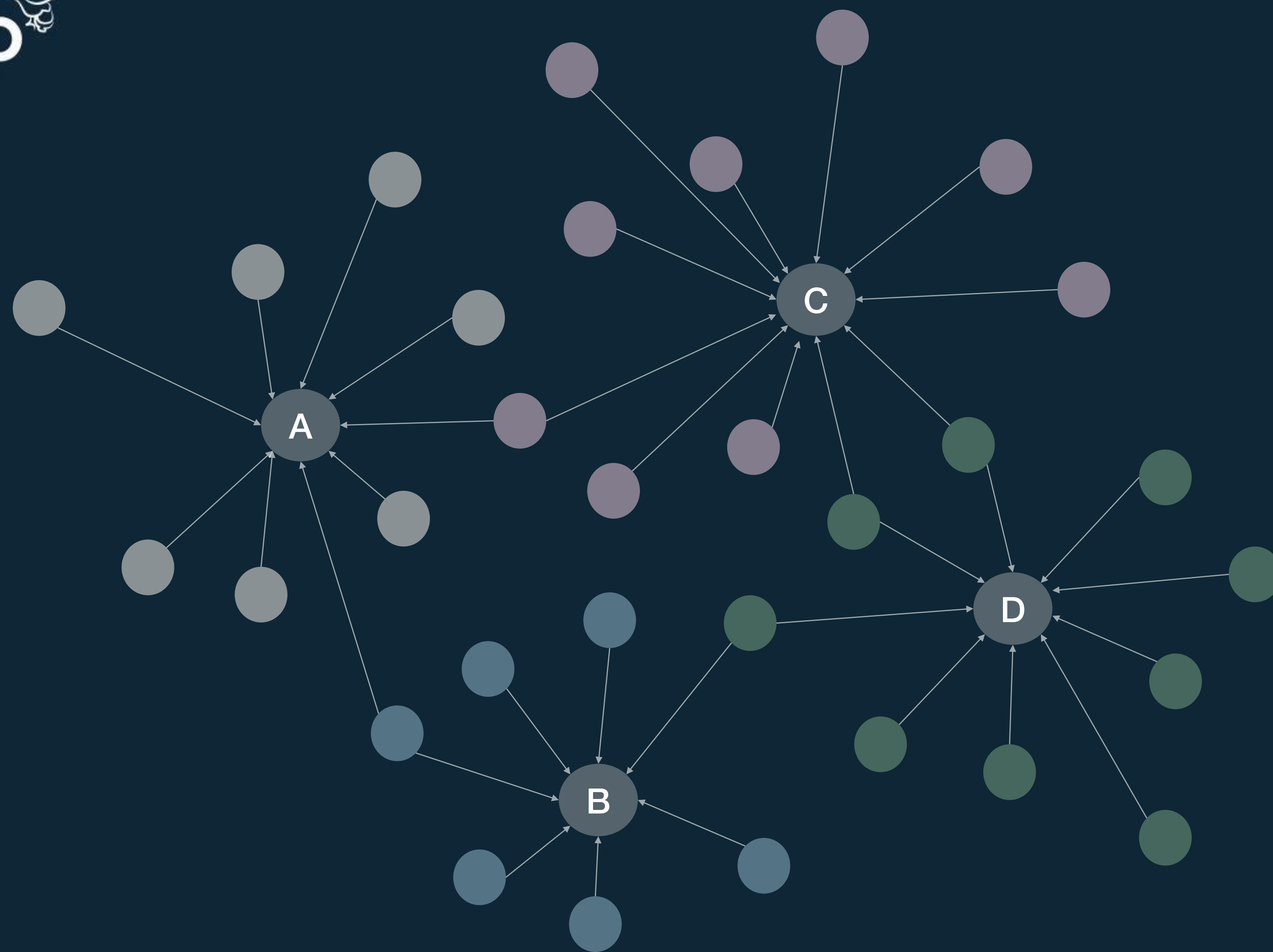


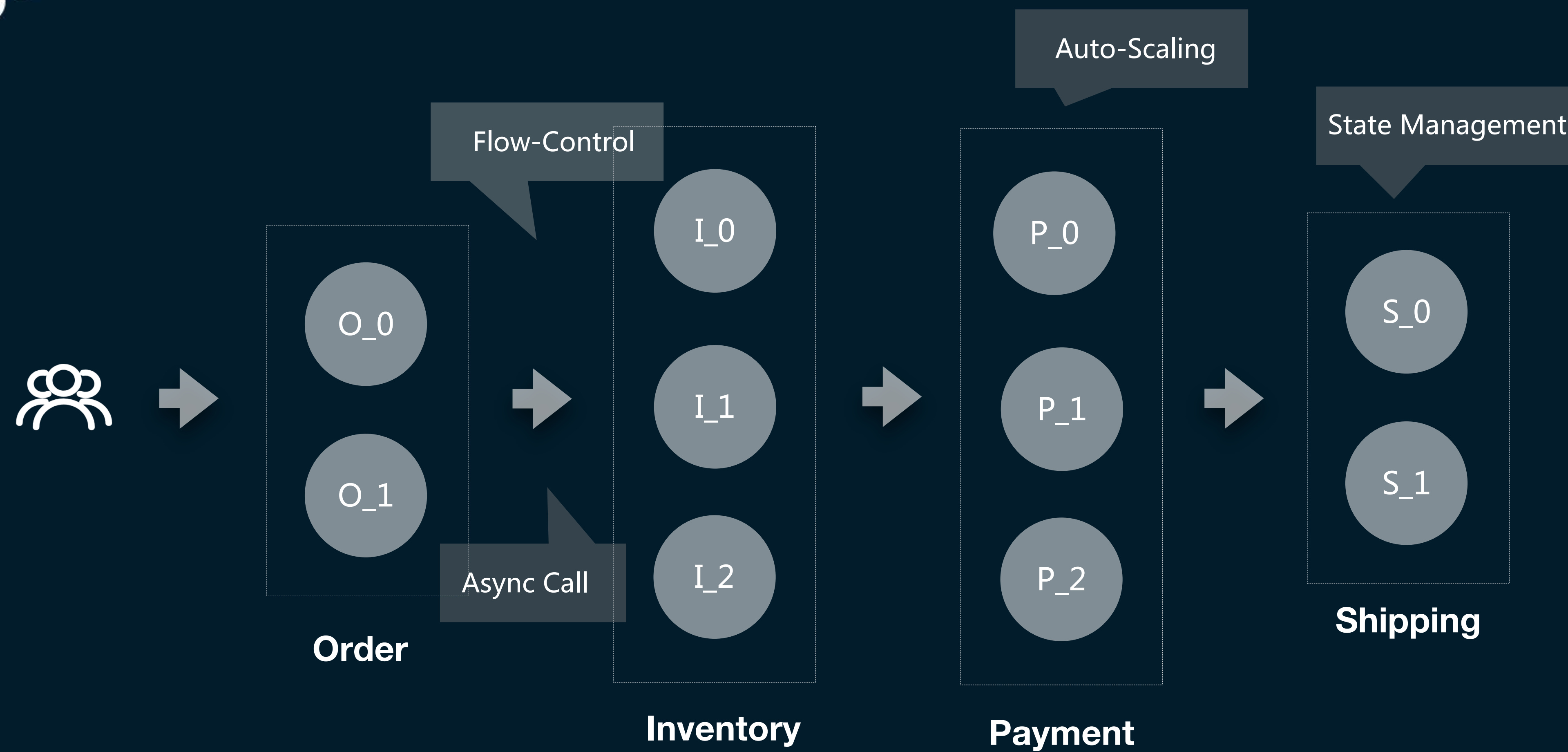


3. Future





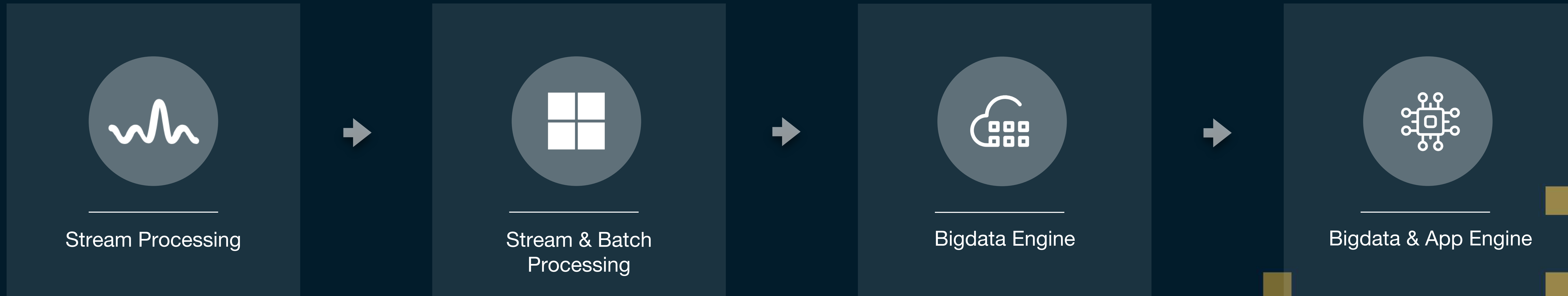




MicroService

Event Driven/Async Processing
Back Pressure & Flow Control
Auto-Scaling
State-Management & Atomicity

Apache Flink – Streaming Technology Redefining Computation



THANKS

