

Sparrow

Distributed Low-Latency Spark Scheduling

Kay Ousterhout, Patrick Wendell, Matei Zaharia, Ion Stoica



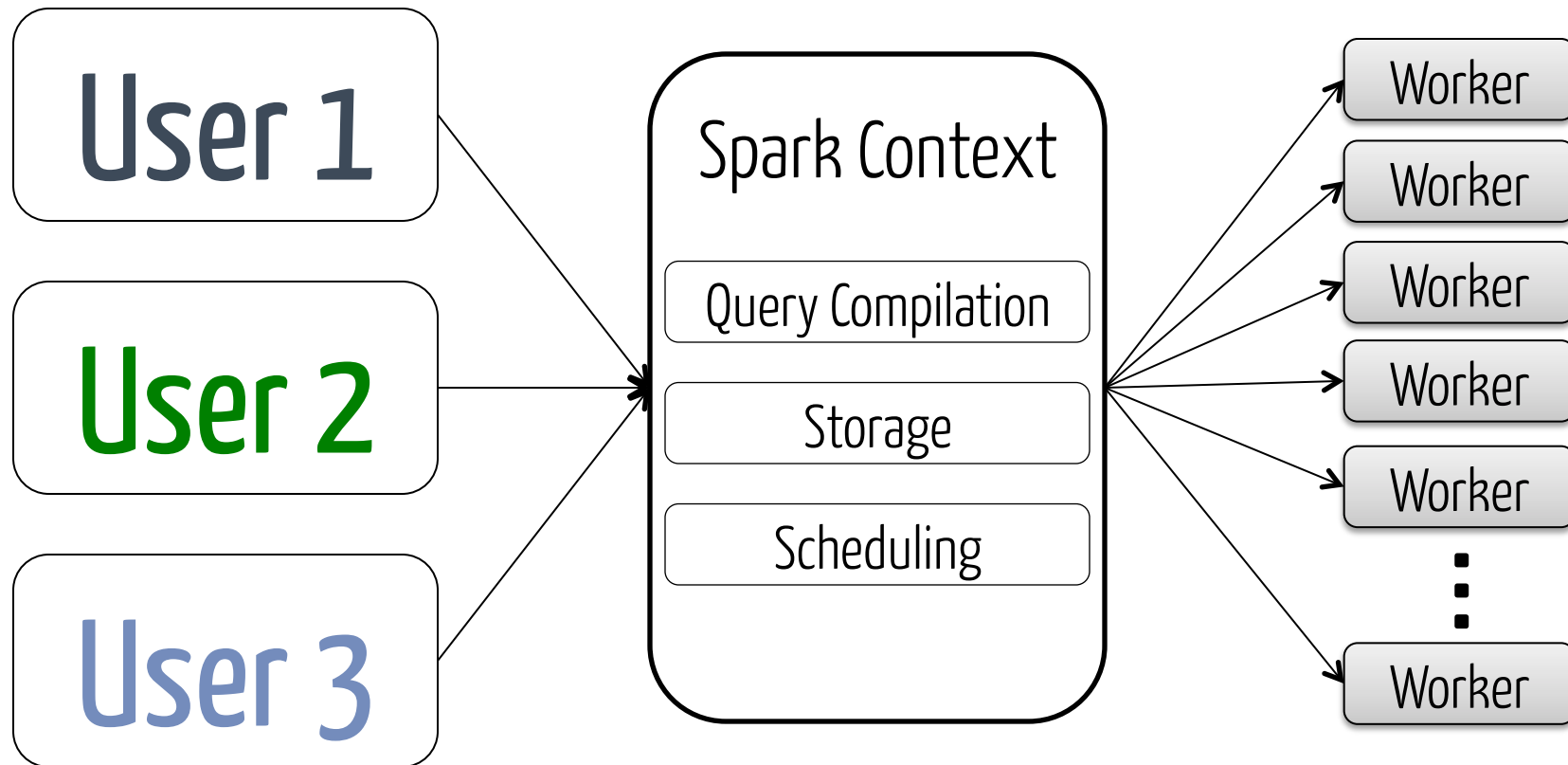
Outline

The Spark scheduling bottleneck

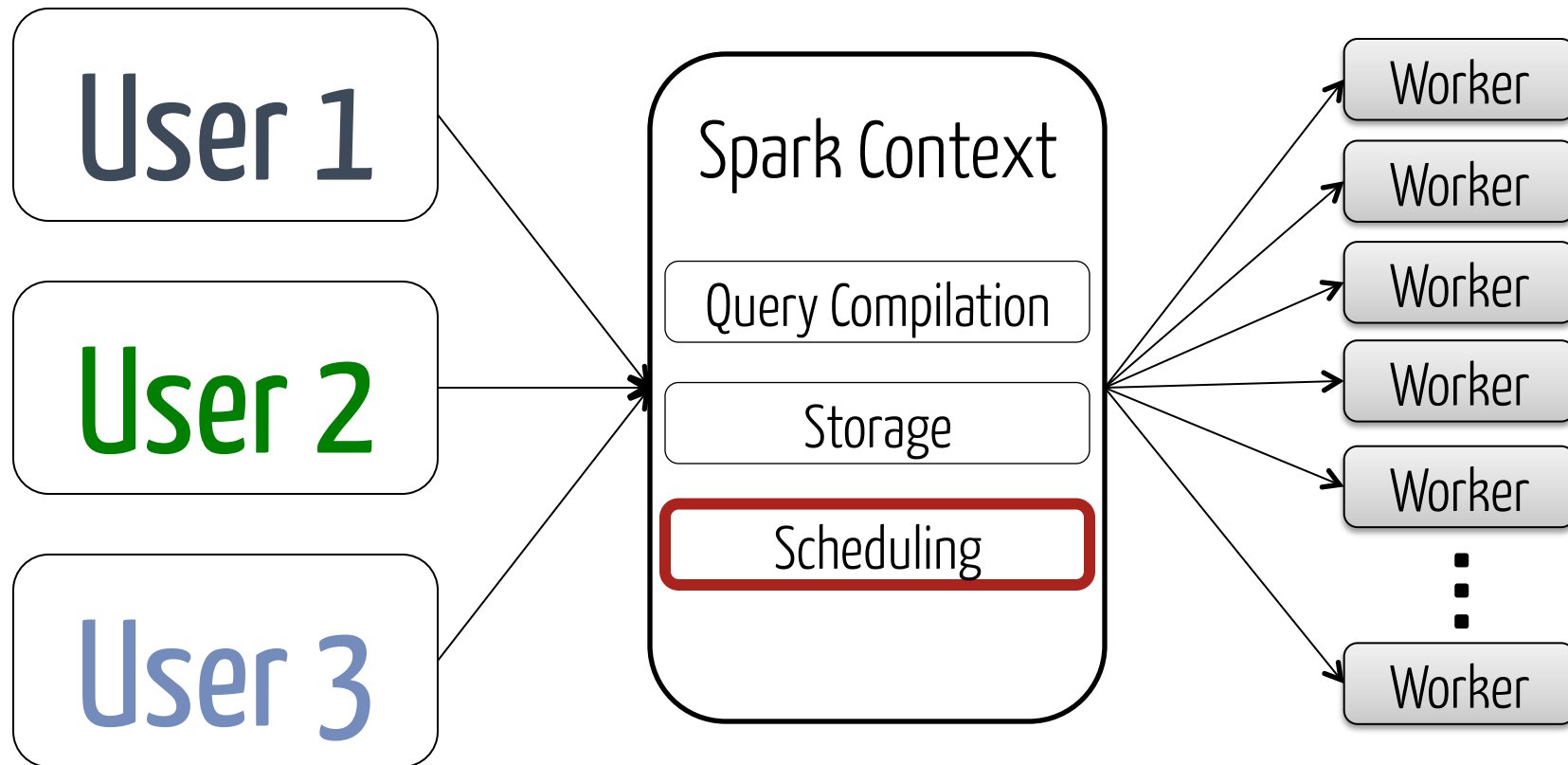
Sparrow's fully distributed, fault-tolerant technique

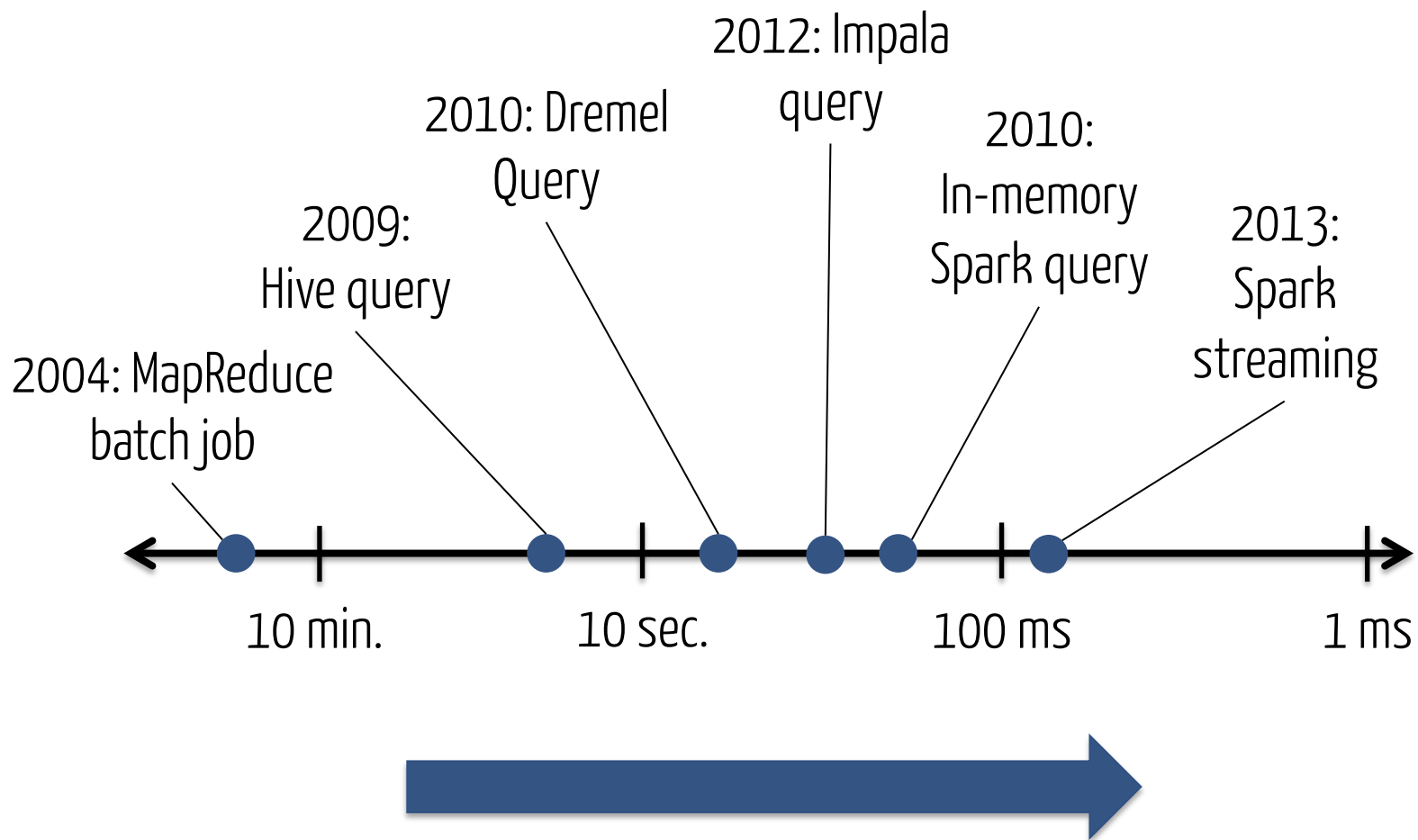
Sparrow's near-optimal performance

Spark Today



Spark Today

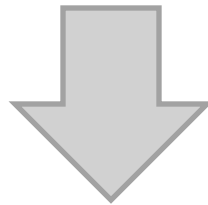




Job Latencies Rapidly Decreasing

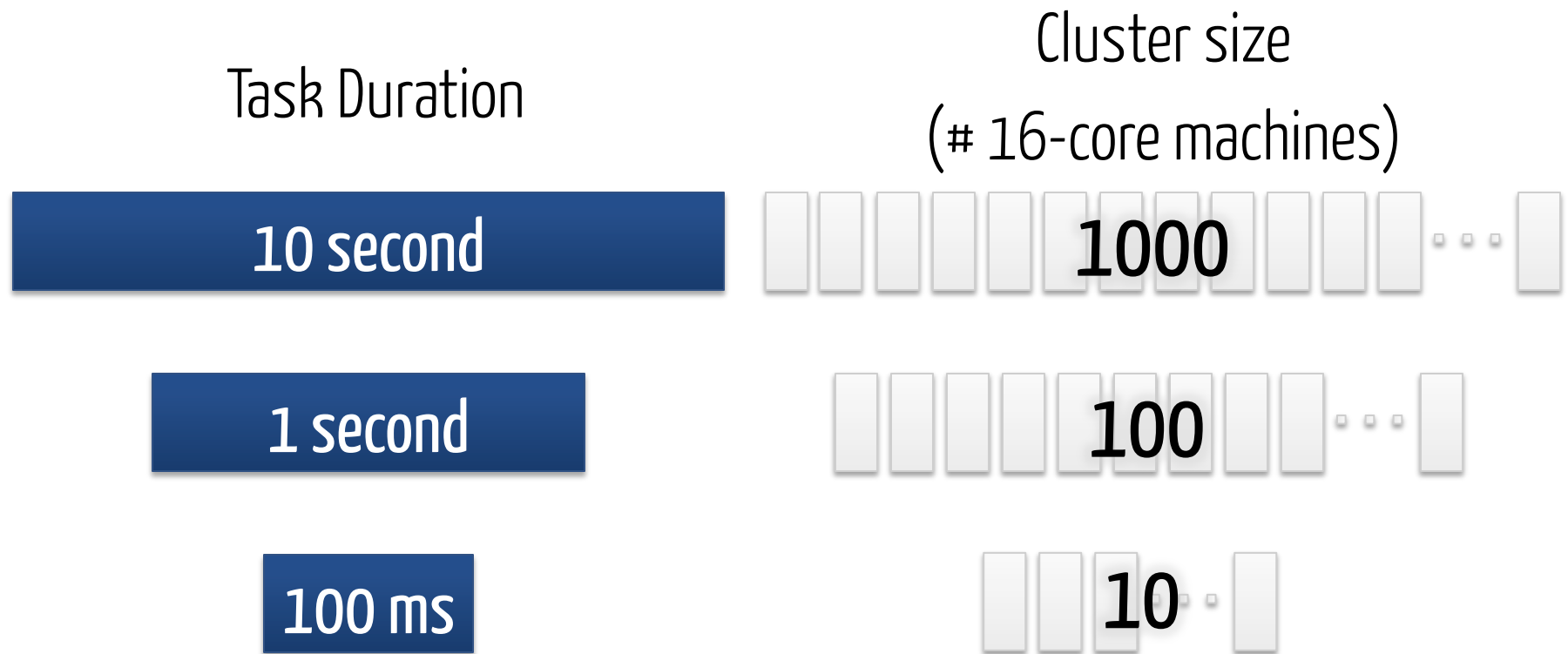
Job latencies rapidly decreasing

Job latencies rapidly decreasing
+
Spark deployments growing in size



Scheduling bottleneck!

Spark scheduler throughput: 1500 tasks / second



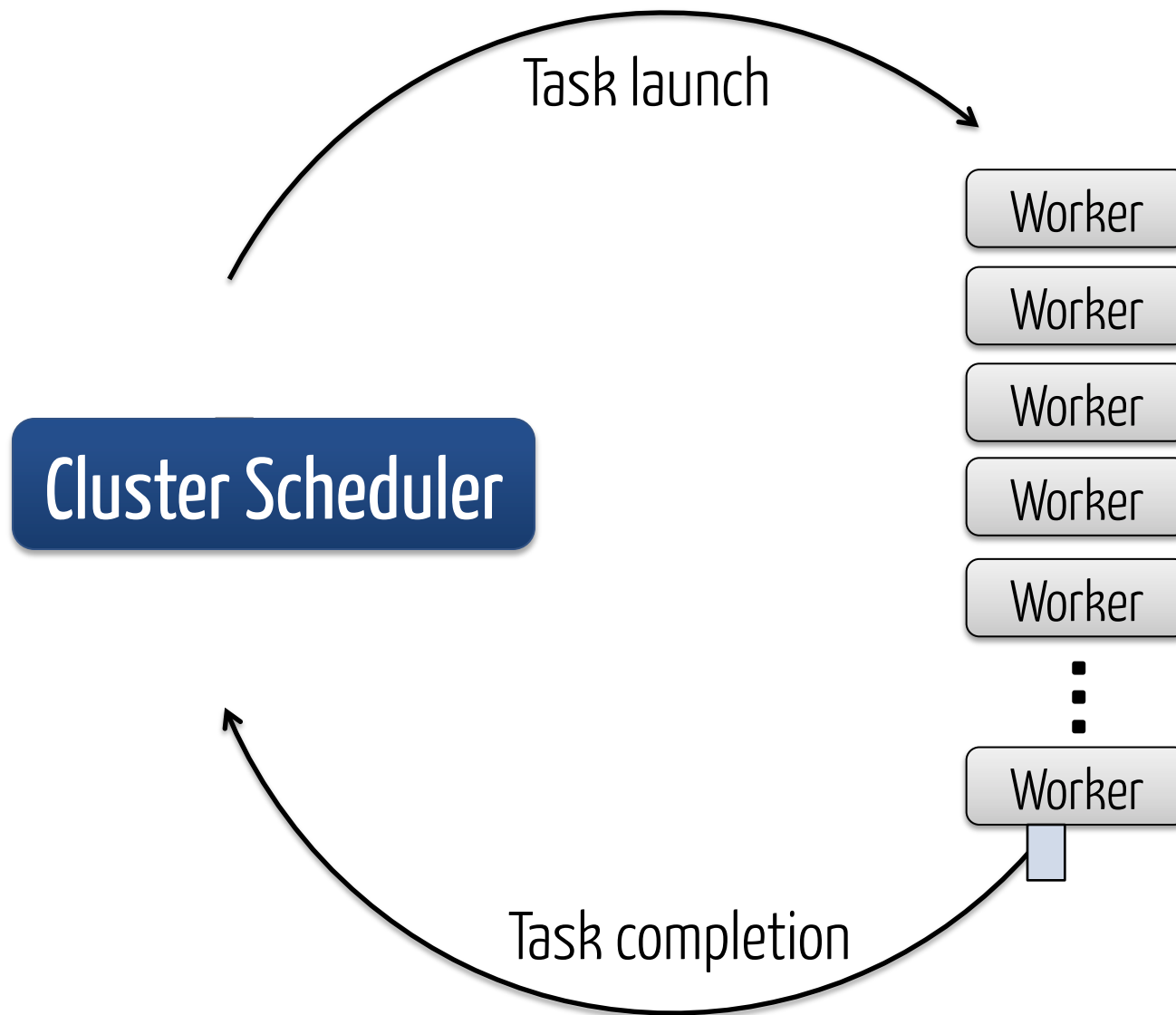
Optimizing the Spark Scheduler

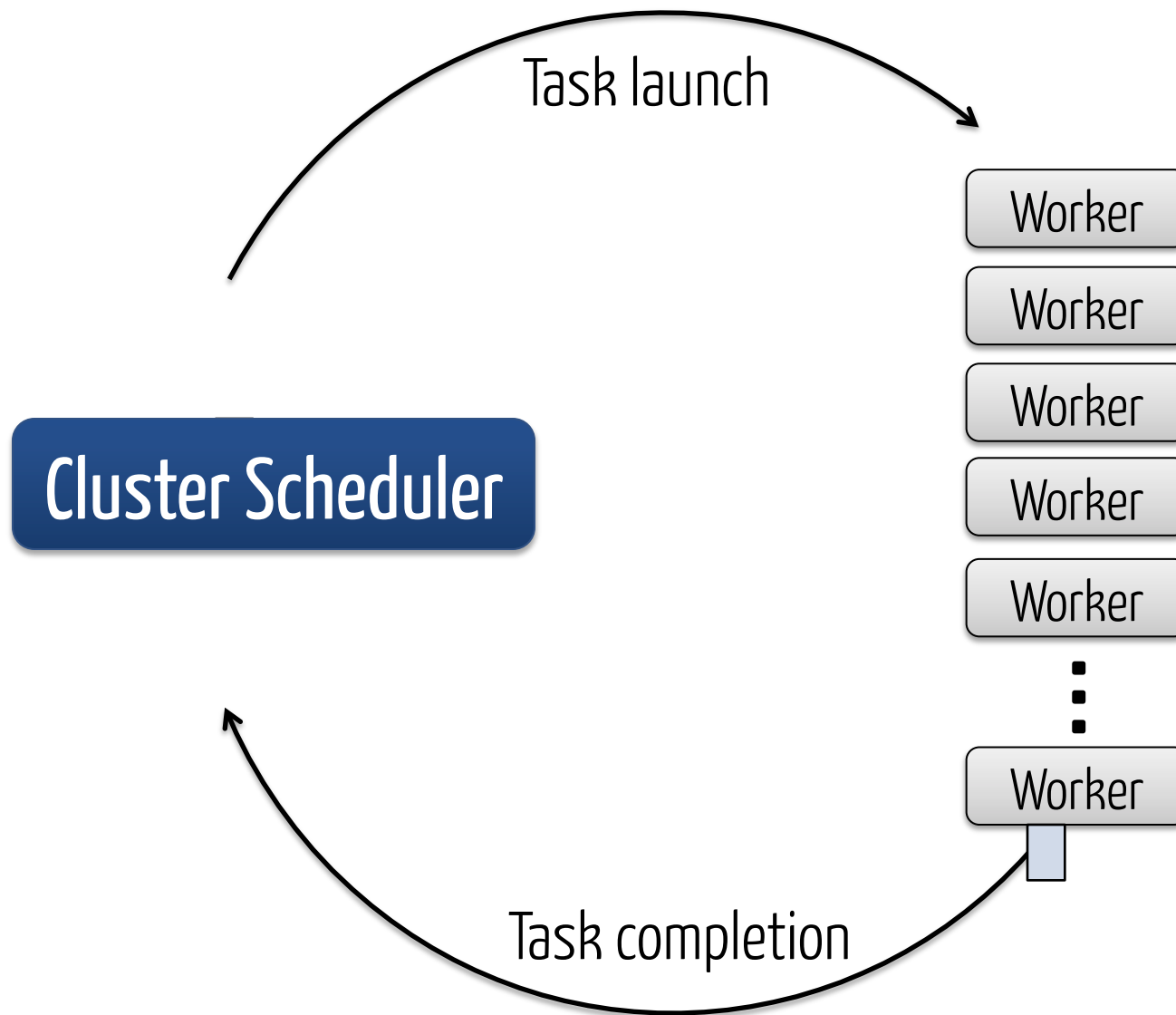
0.8: Monitoring code moved off critical path

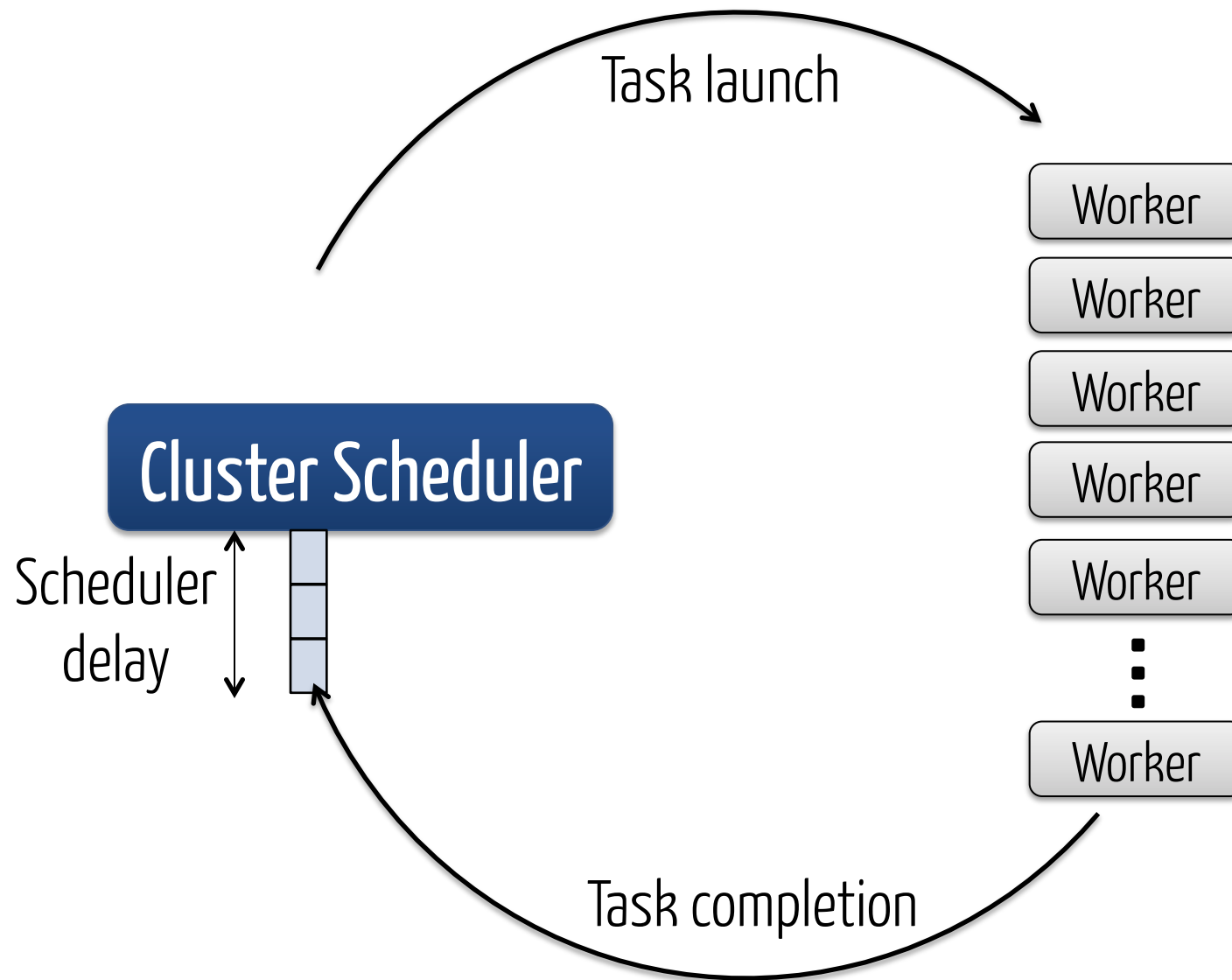
0.8.1: Result deserialization moved off critical path

Future improvements may yield 2-3x higher throughput

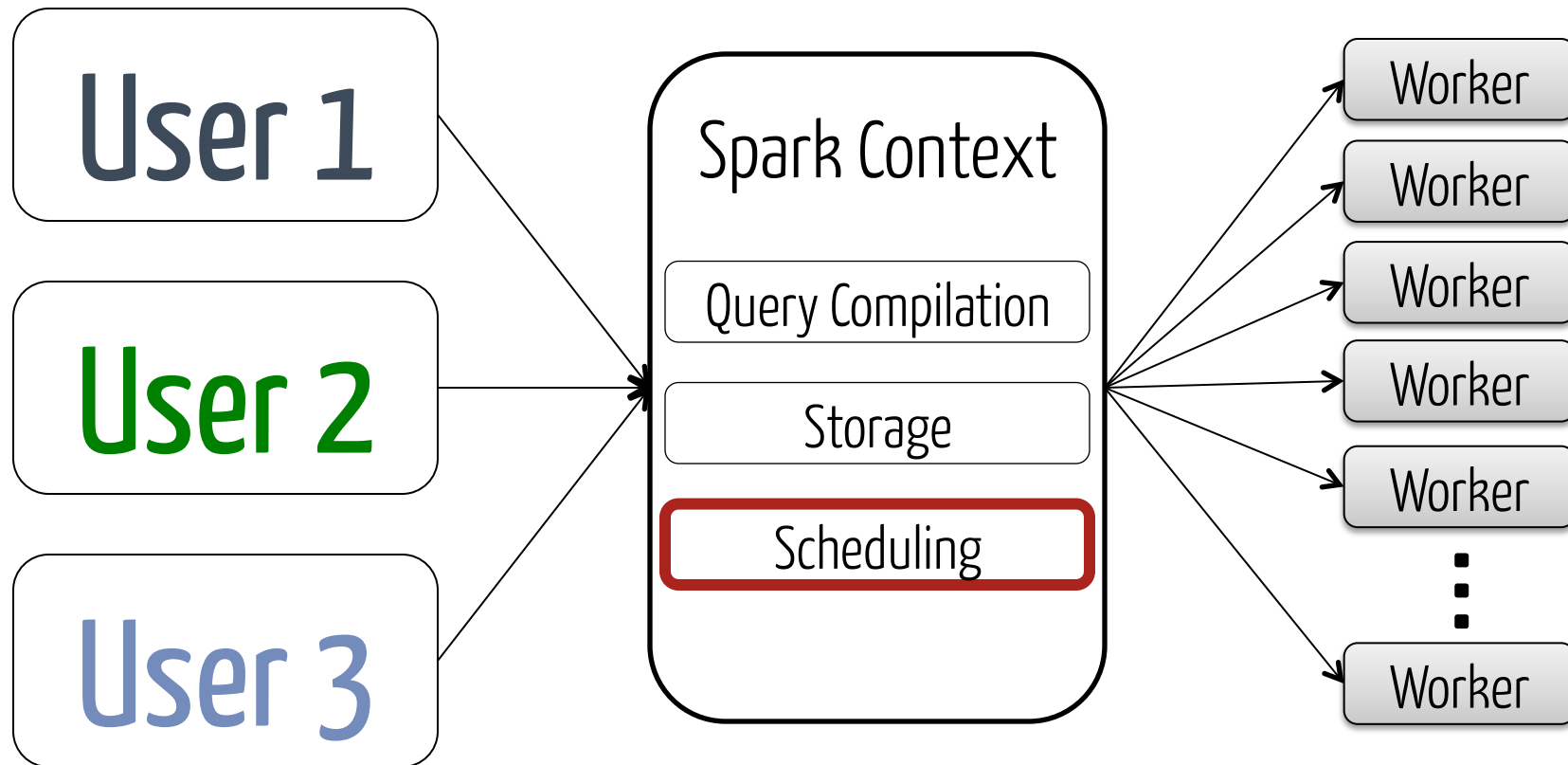
Is the scheduler the
bottleneck in my cluster?



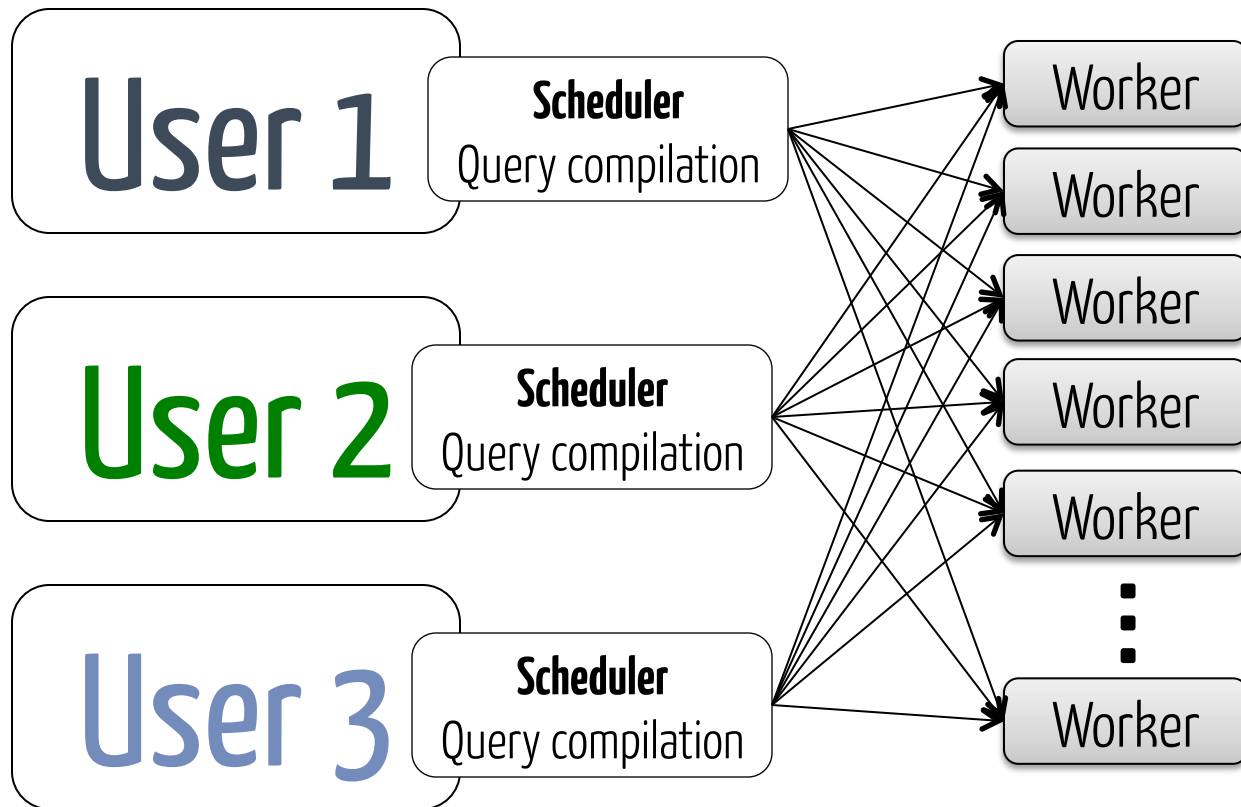




Spark Today

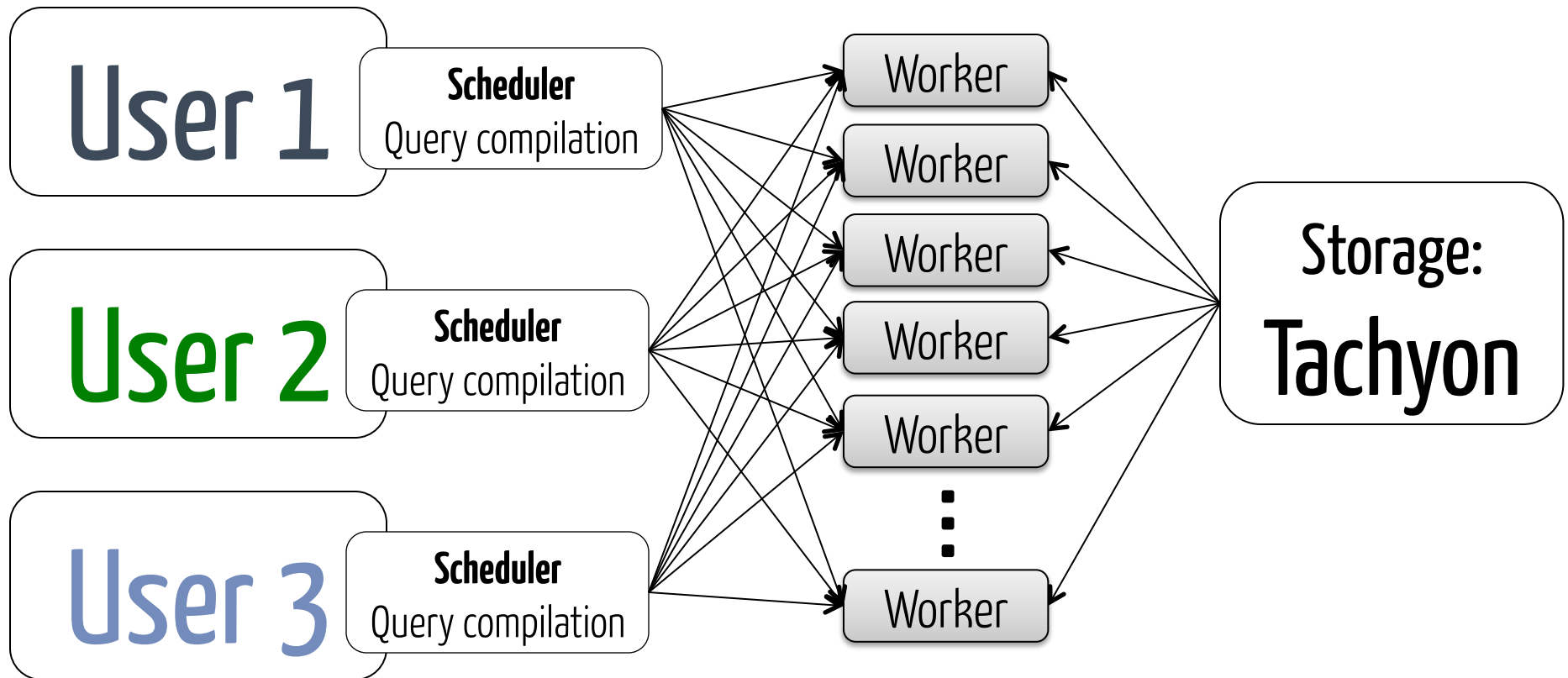


Future Spark

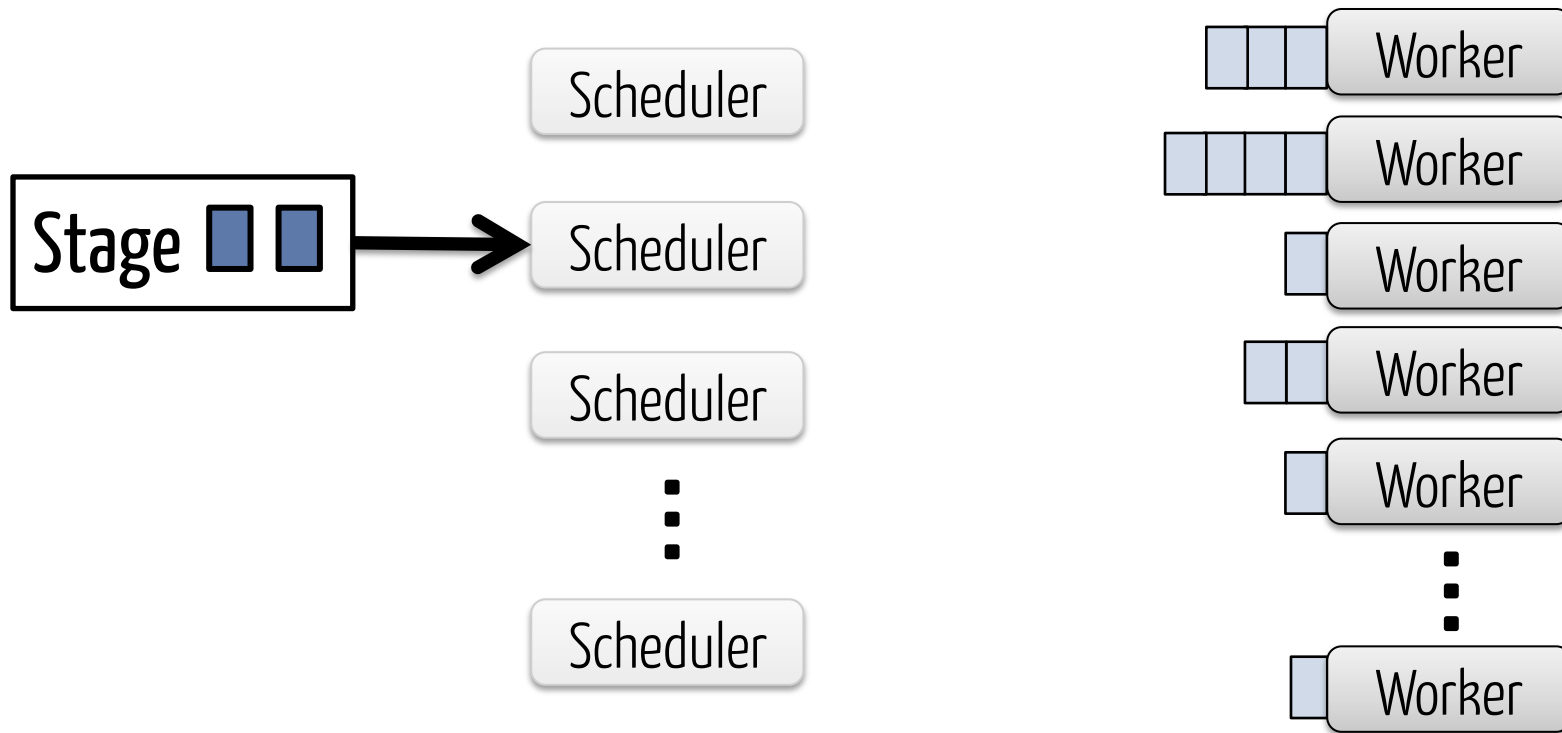


Benefits:
High throughput
Fault tolerance

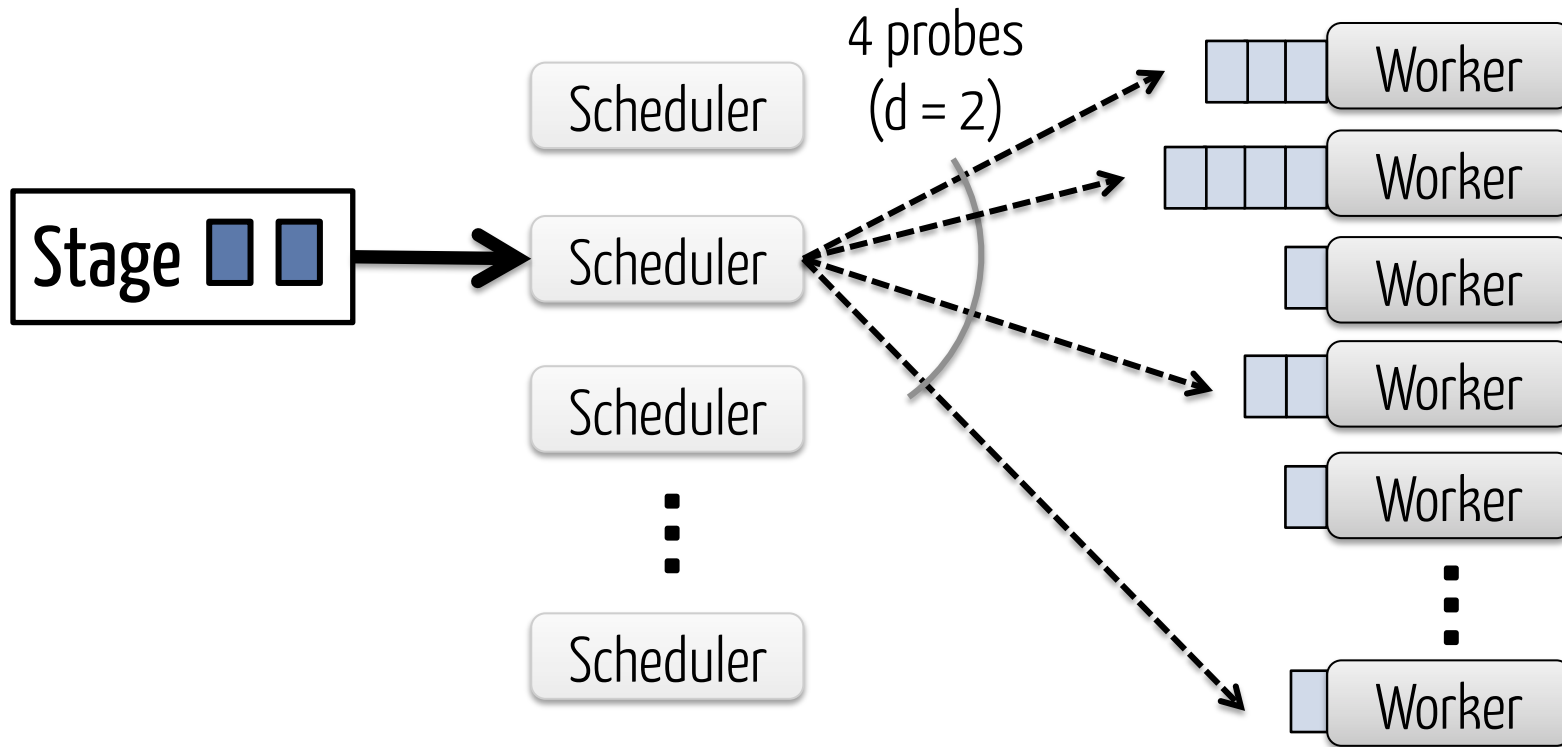
Future Spark



Scheduling with Sparrow

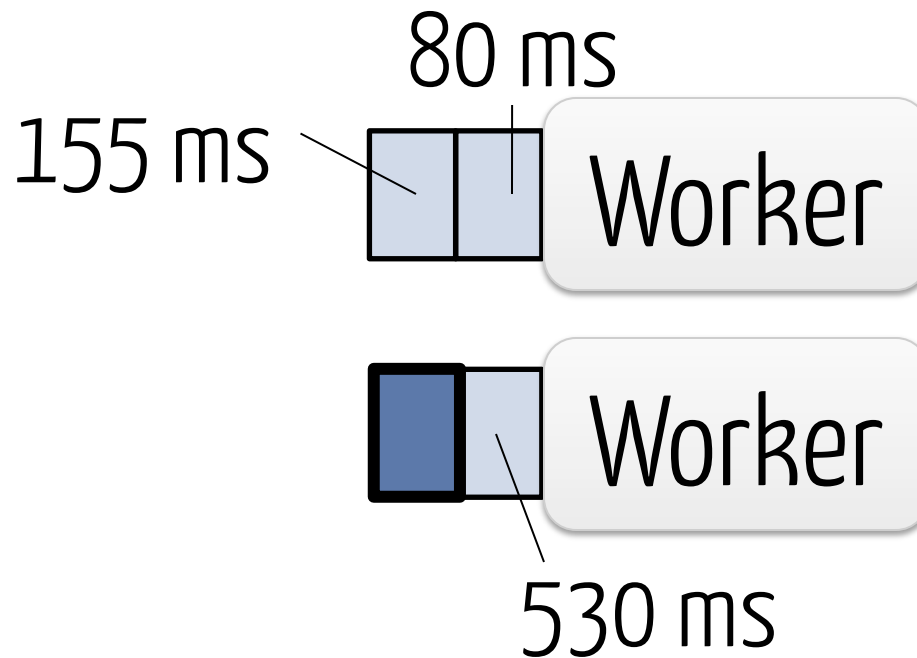


Batch Sampling



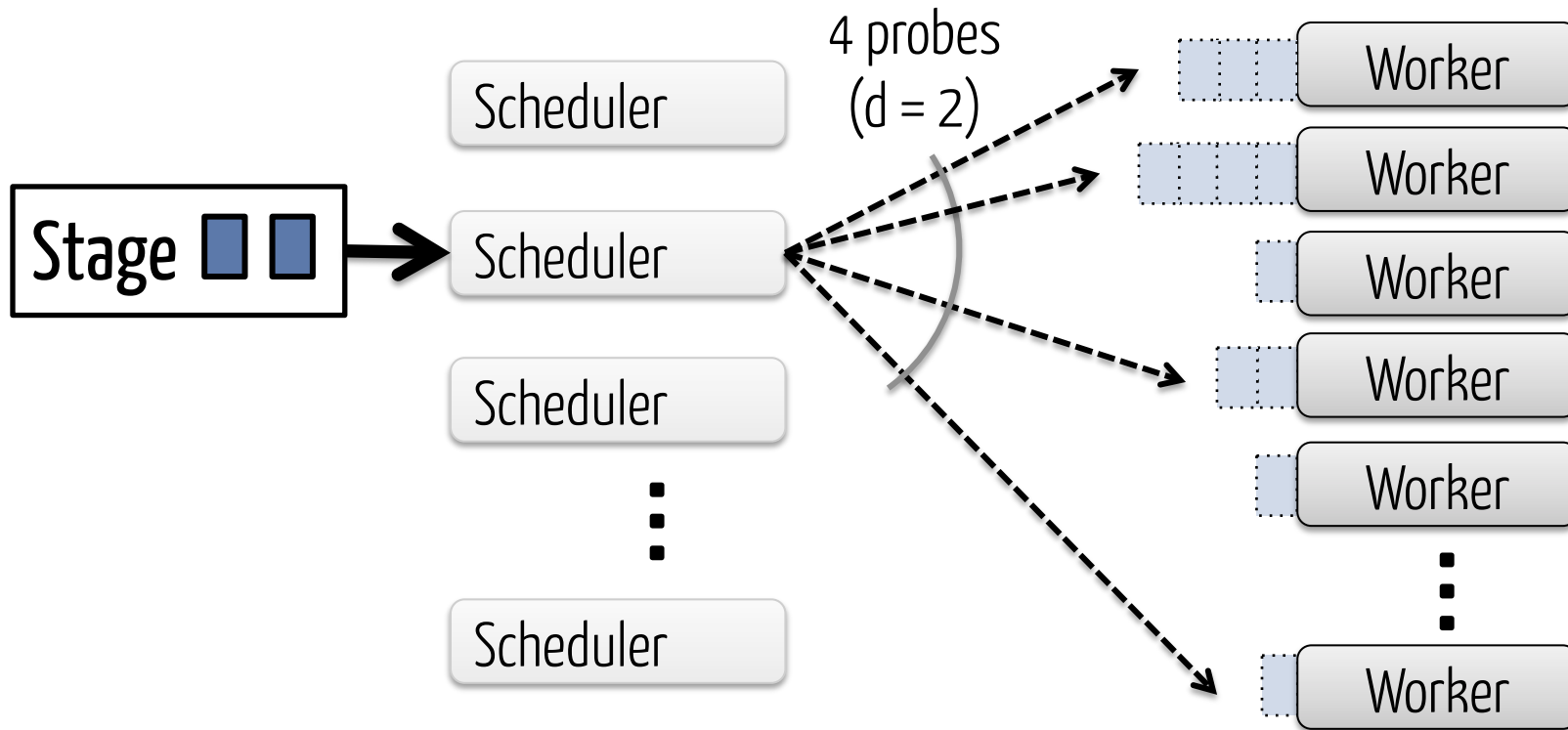
Place m tasks on the least loaded of $d \bullet m$ workers

Queue length poor predictor of wait time



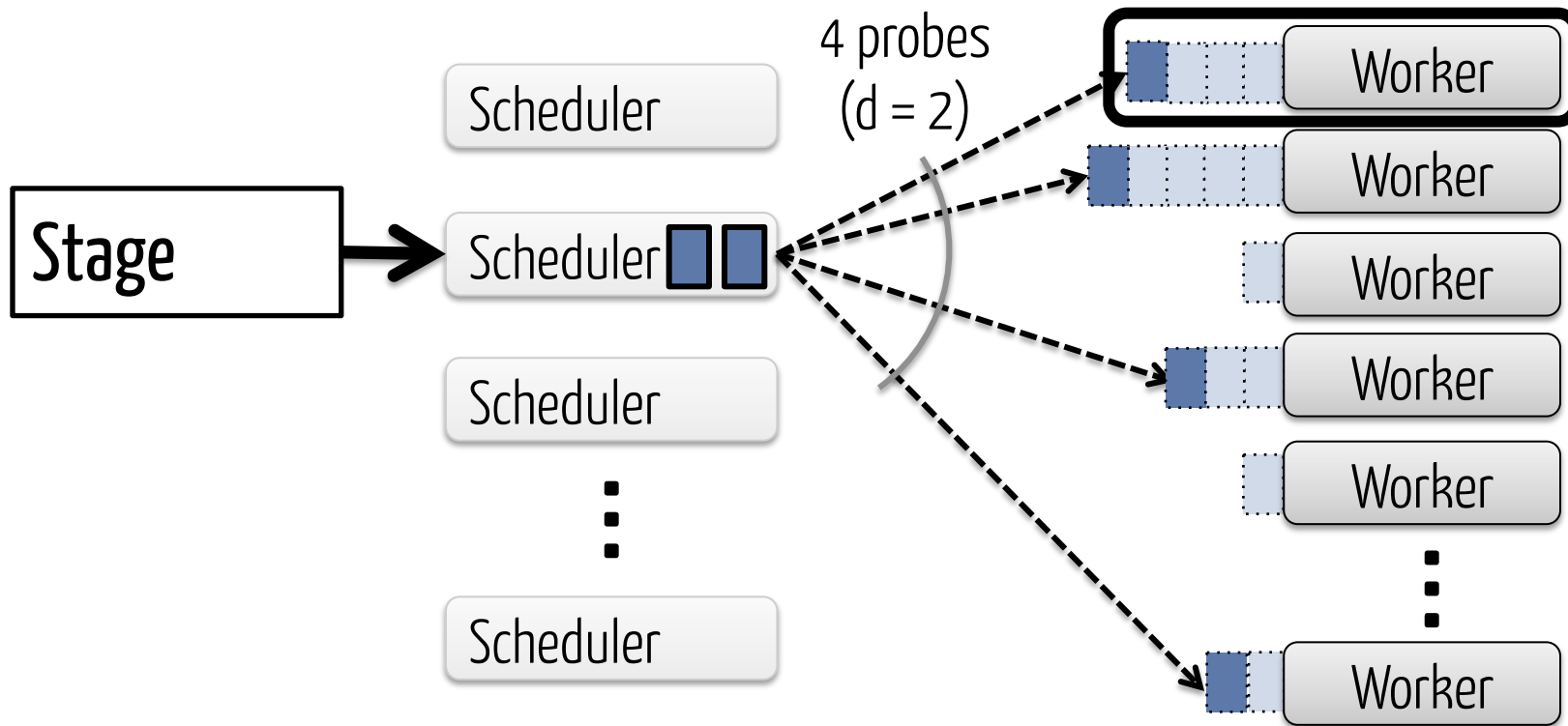
Poor performance on heterogeneous workloads

Late Binding



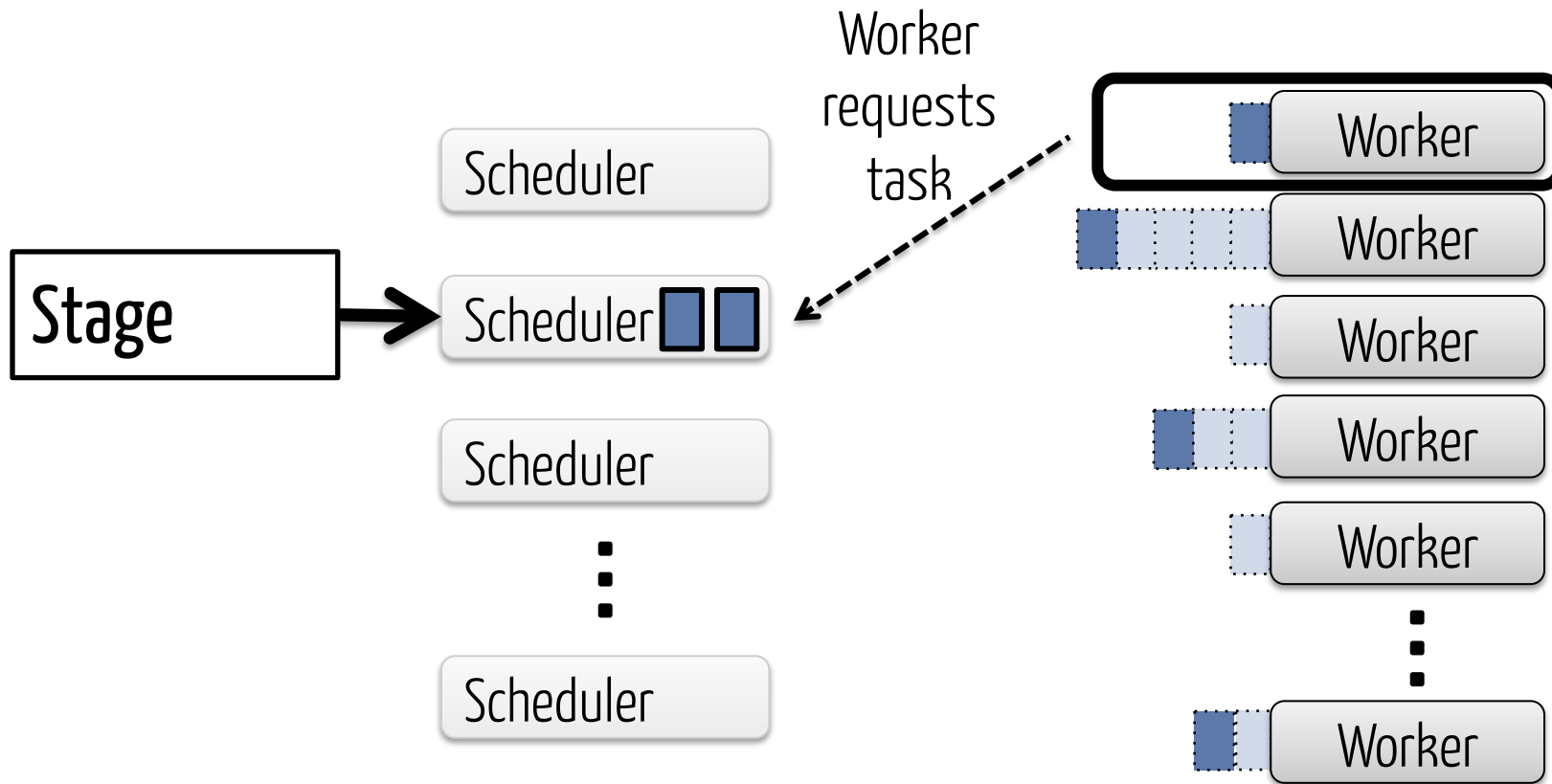
Place m tasks on the least loaded of $d \cdot m$ workers

Late Binding



Place m tasks on the least loaded of $d \cdot m$ workers

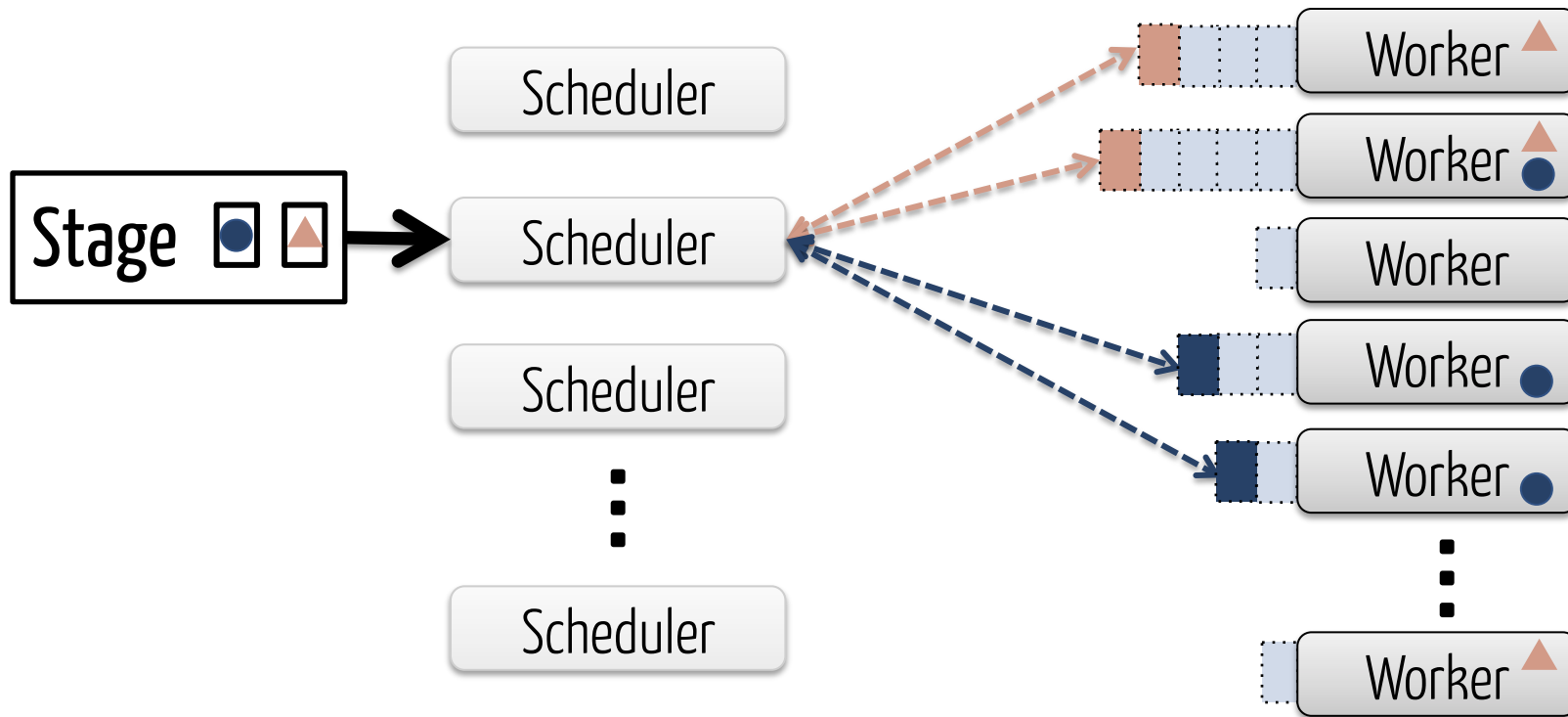
Late Binding



Place m tasks on the least loaded of $d \bullet m$ workers

What about constraints?

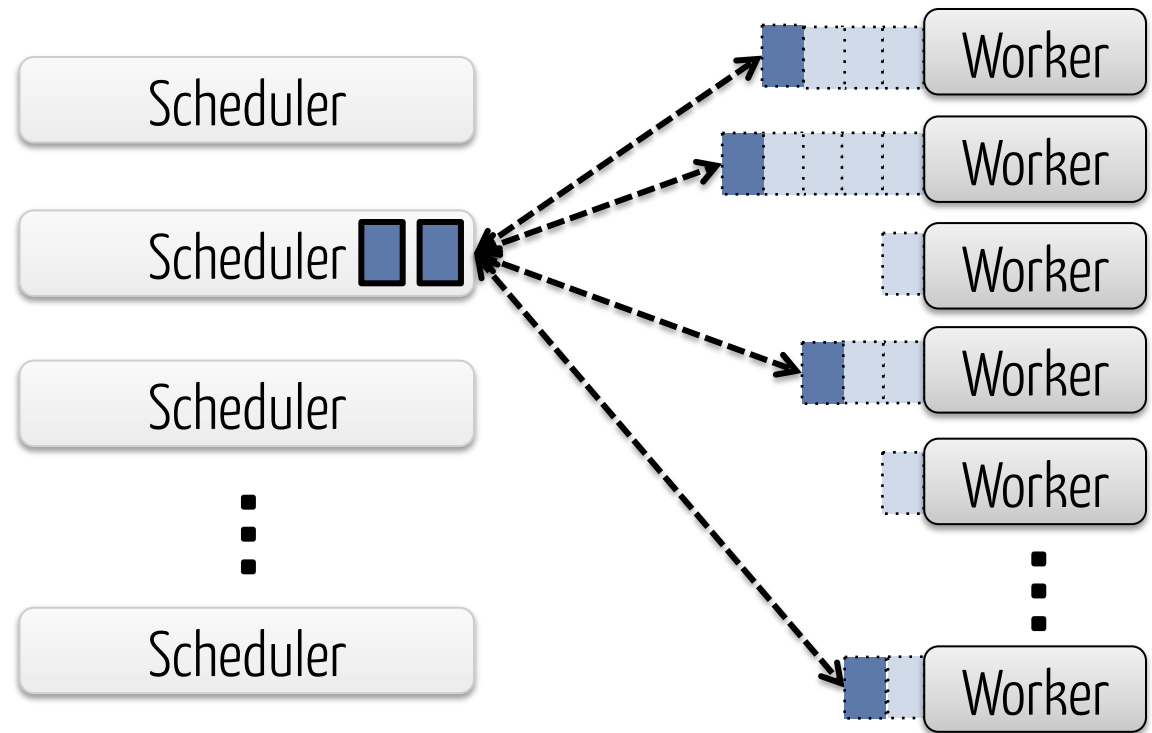
Per-Task Constraints



Probe separately for each task

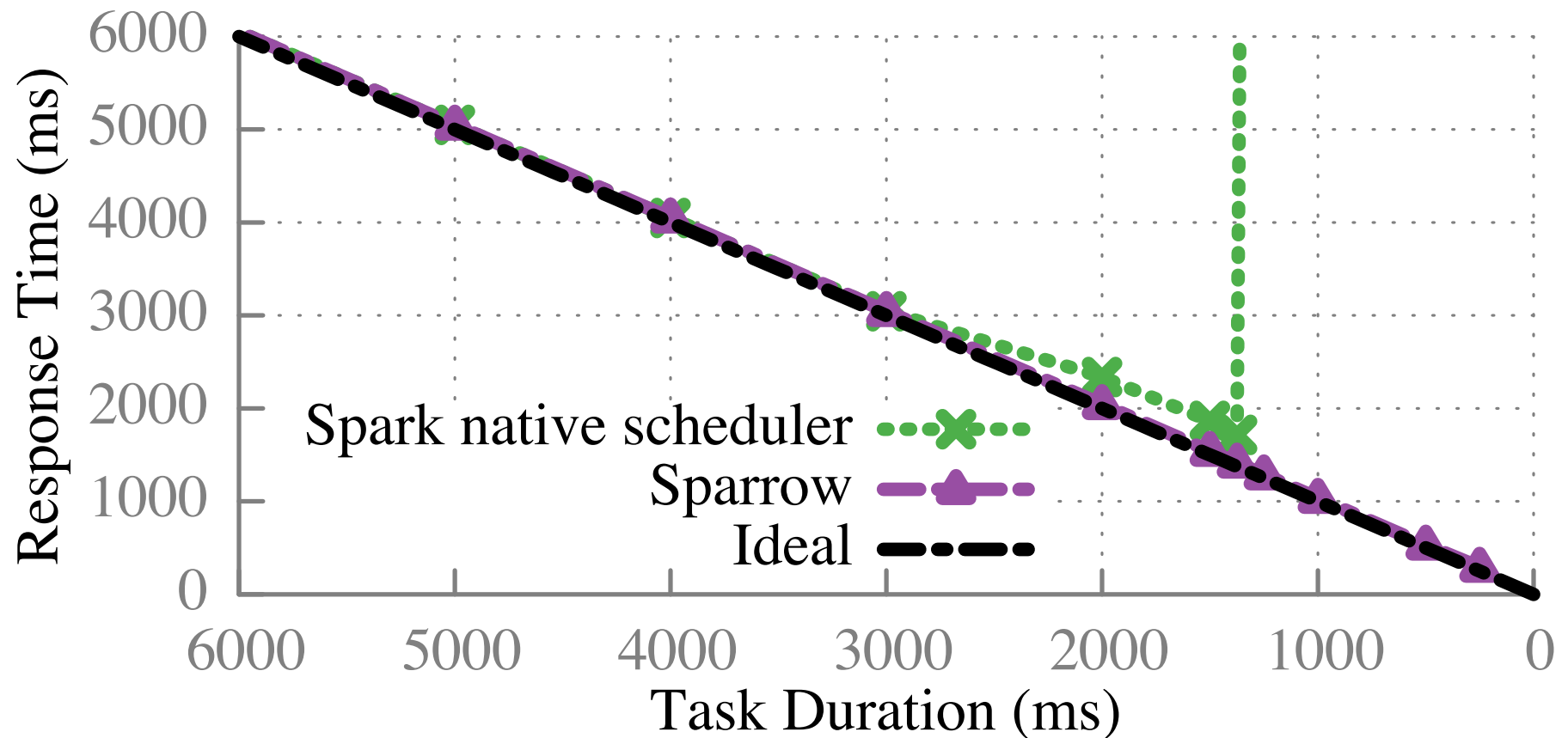
Technique Recap

Batch sampling
+
Late binding
+
Constraints



**How well does Sparrow
perform?**

How does Sparrow compare to Spark's native scheduler?



100 16-core EC2 nodes, 10 tasks/job, 10 schedulers, 80% load

TPC-H Queries: Background

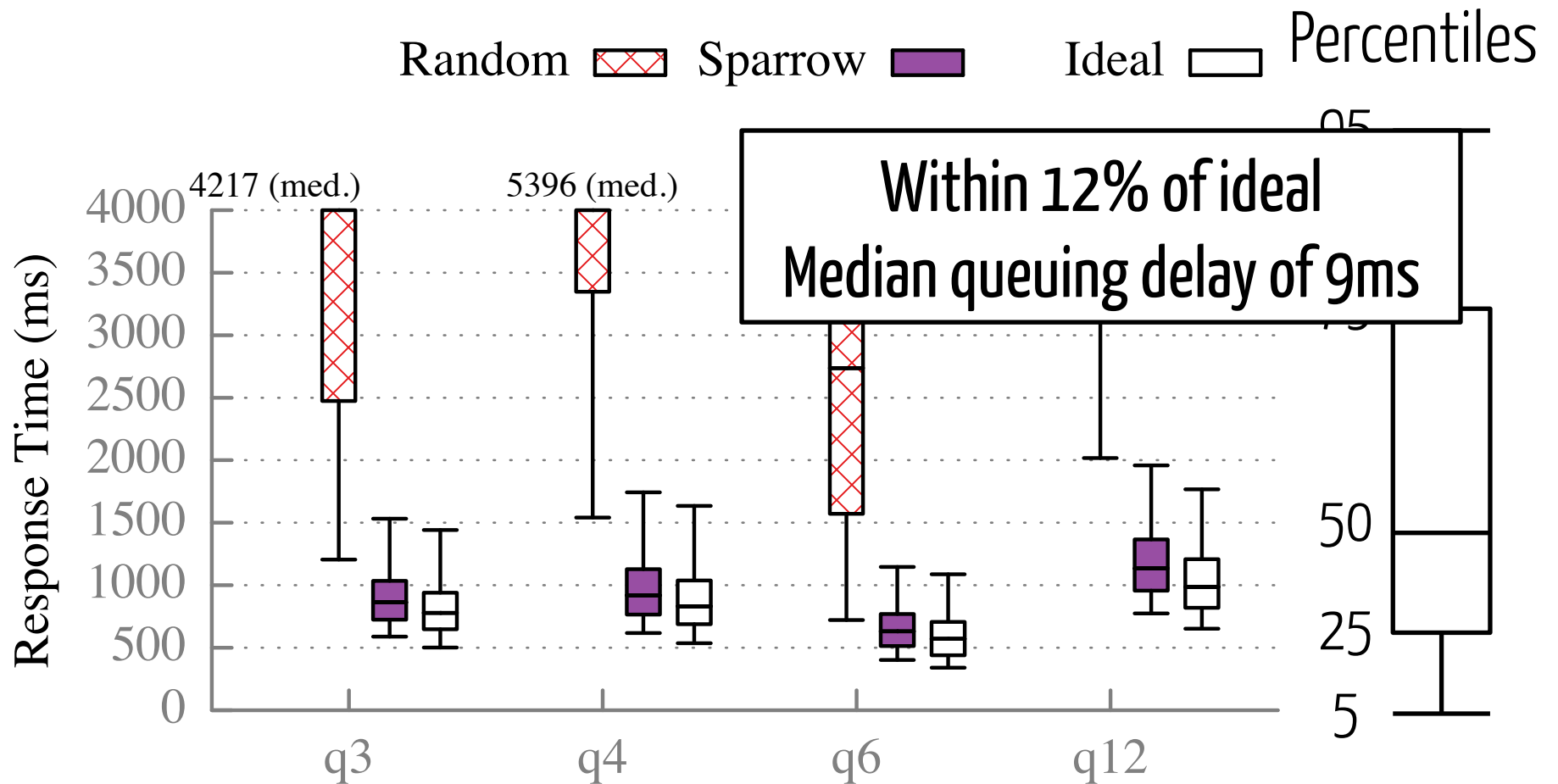
TPC-H: Common benchmark for analytics workloads

Shark: SQL execution engine

Spark

Sparrow

TPC-H Queries



100 16-core EC2 nodes, 10 schedulers, 80% load

Policy Enforcement

Priorities

Serve queues based on strict priorities

High Priority



Low Priority



Worker

Fair Shares

Serve queues using weighted fair queuing

User A (75%)

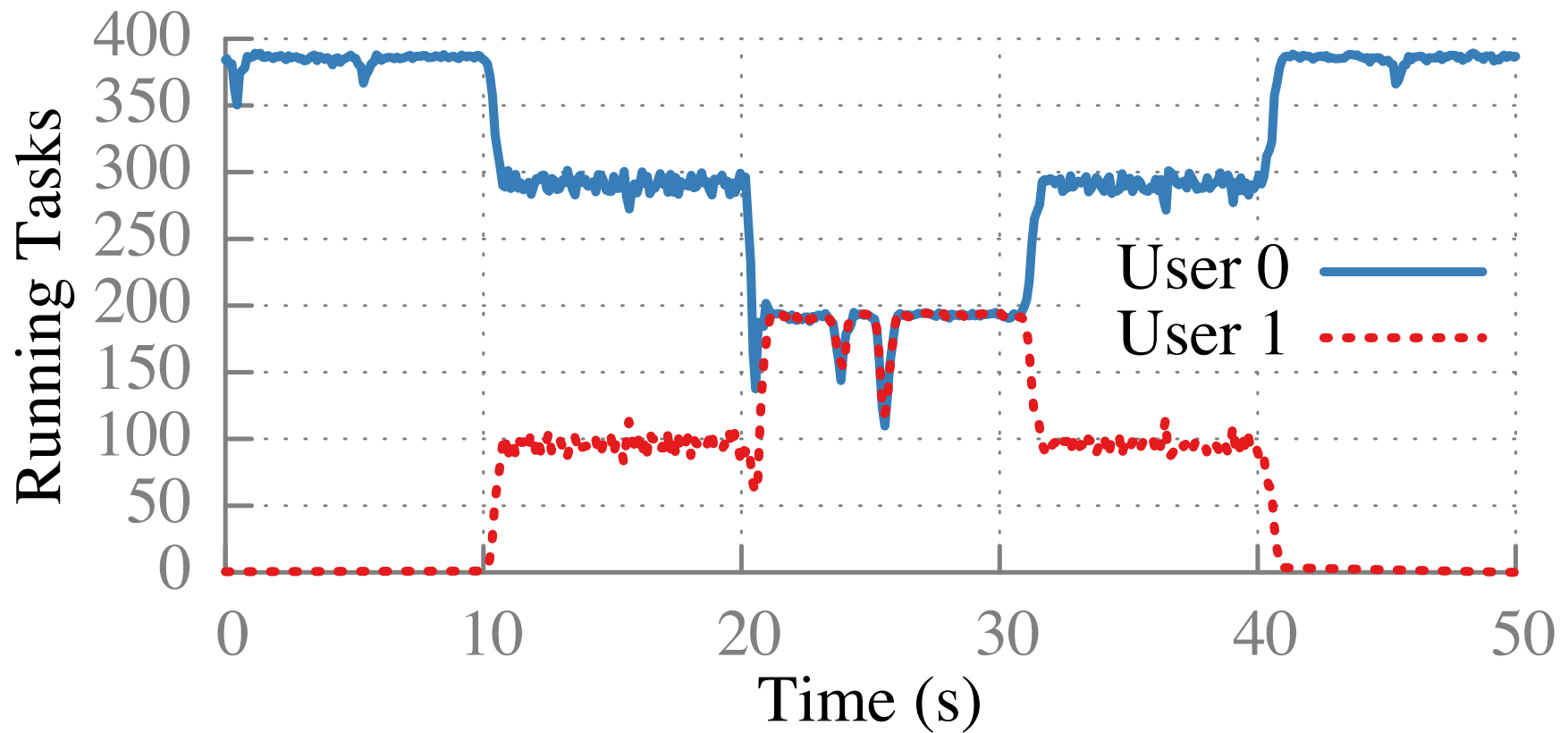


User B (25%)

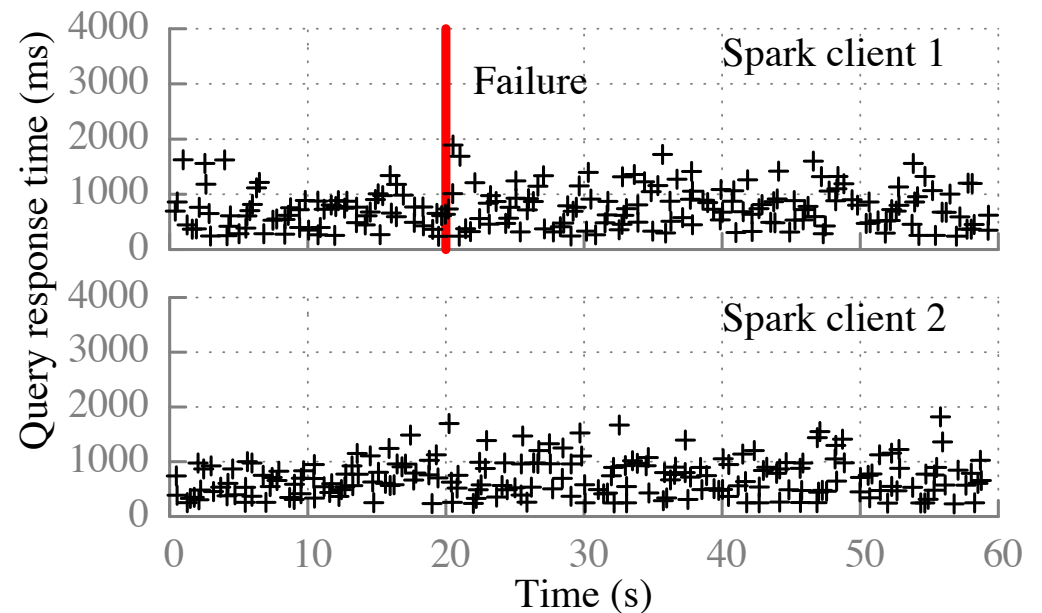
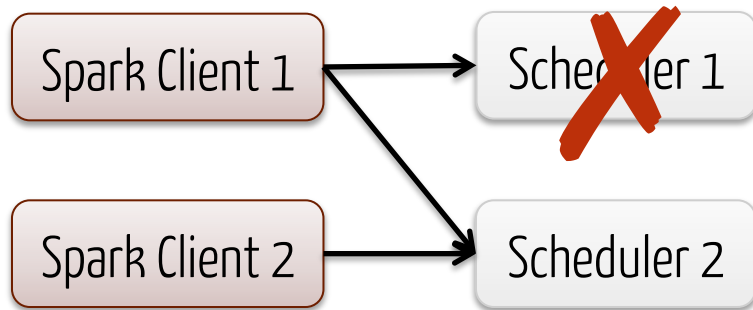


Worker

Weighted Fair Sharing



Fault Tolerance



Timeout: 100ms

Failover: 5ms

Re-launch queries: 15ms

Making Sparrow feature-complete

Interfacing with UI

Delay scheduling

Speculation

(1) Diagnosing a Spark scheduling bottleneck

Metric	Min	25th percentile	Median	75th percent
Duration	1 ms	1 ms	1 ms	1 ms
Time spent fetching task results	0 ms	0 ms	0 ms	0 ms
Scheduler delay	30 ms	81 ms	91 ms	95 ms

(2) Distributed, fault-tolerant scheduling with Sparrow

