# ABOUT MYSELF
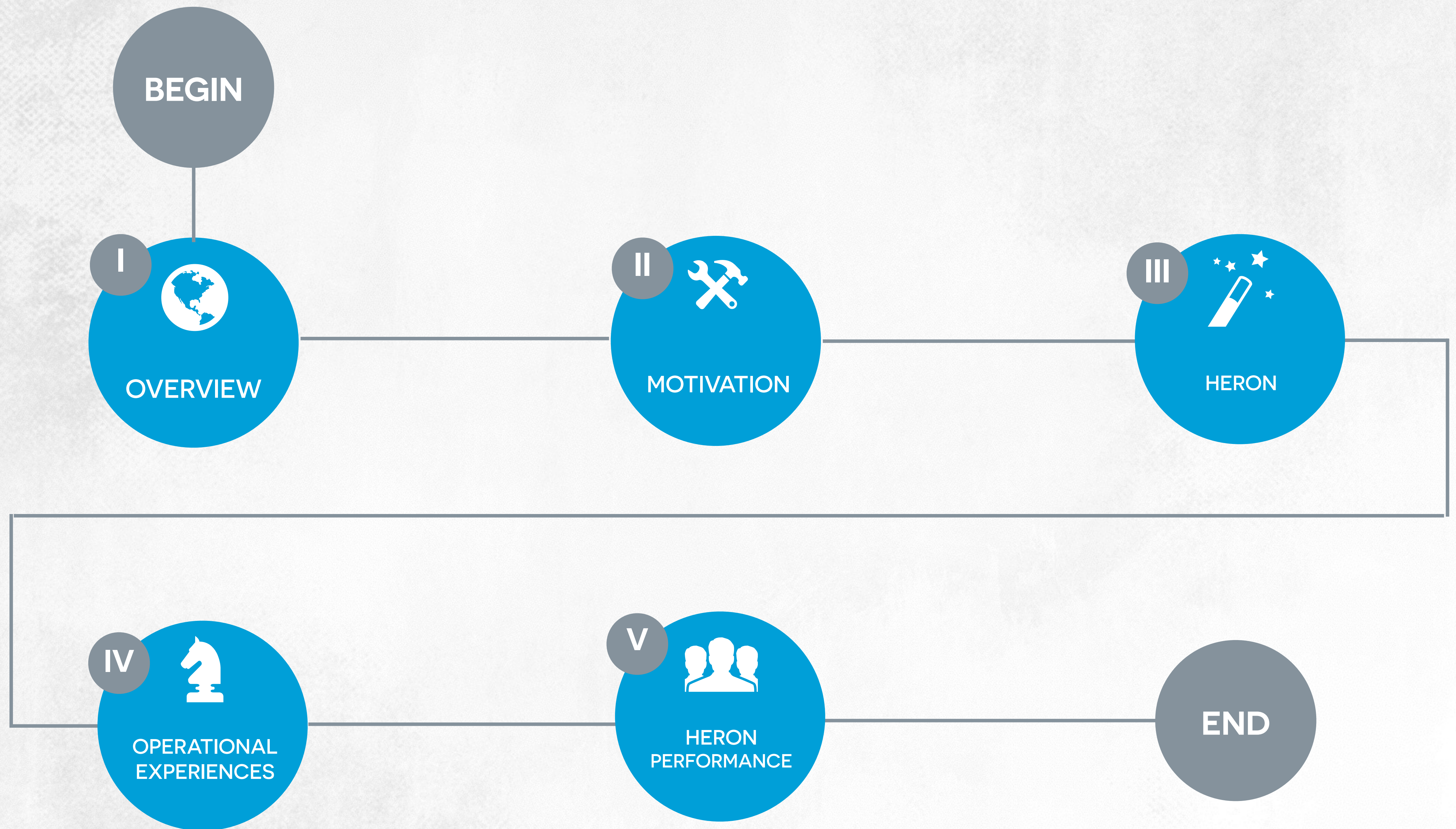## intro

Lead for Twitter Heron project

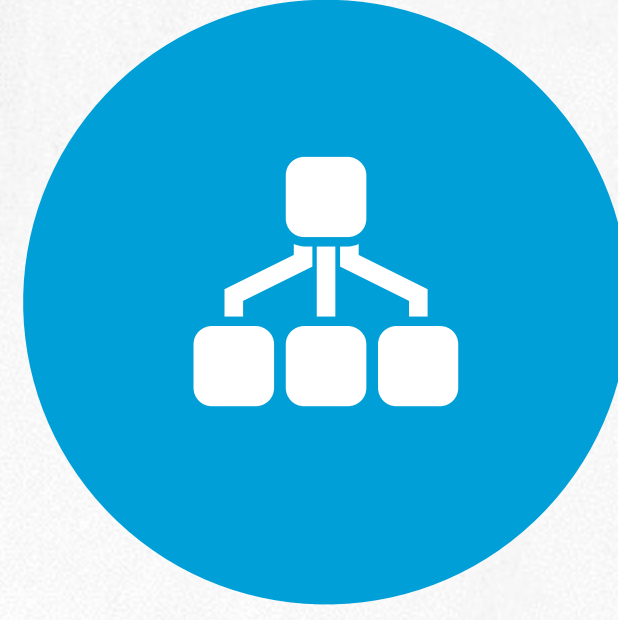Senior Engineer for Real Time Analytics

# OVERVIEW

# TWITTER IS REAL TIME

### REAL TIME TRENDS

Emerging break out trends in Twitter (in the form #hashtags)

### REAL TIME CONVERSATIONS

Real time sports conversations related with a topic (recent goal or touchdown)

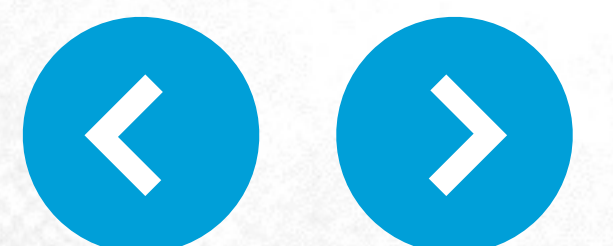### REAL TIME RECOMMENDATIONS

Real time product recommendations based on your behavior & profile

### REAL TIME SEARCH

Real time search of tweets

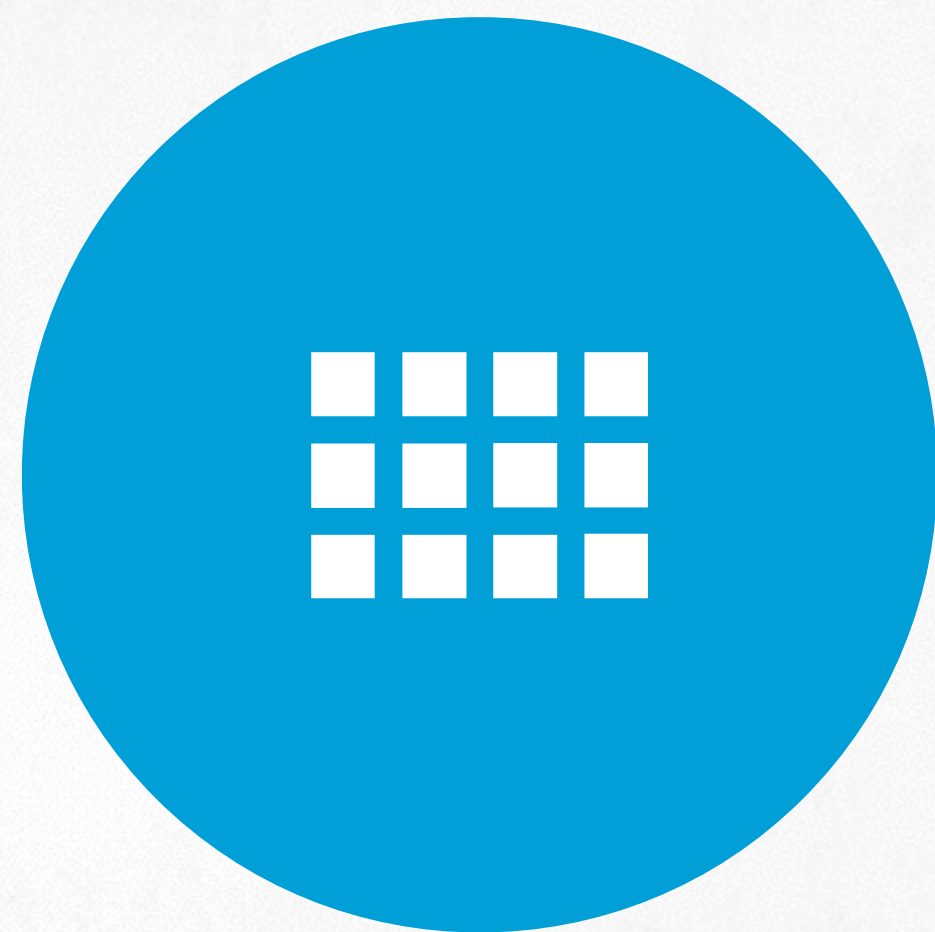## ANALYZING BILLIONS OF EVENTS IN REAL TIME IS A CHALLENGE!
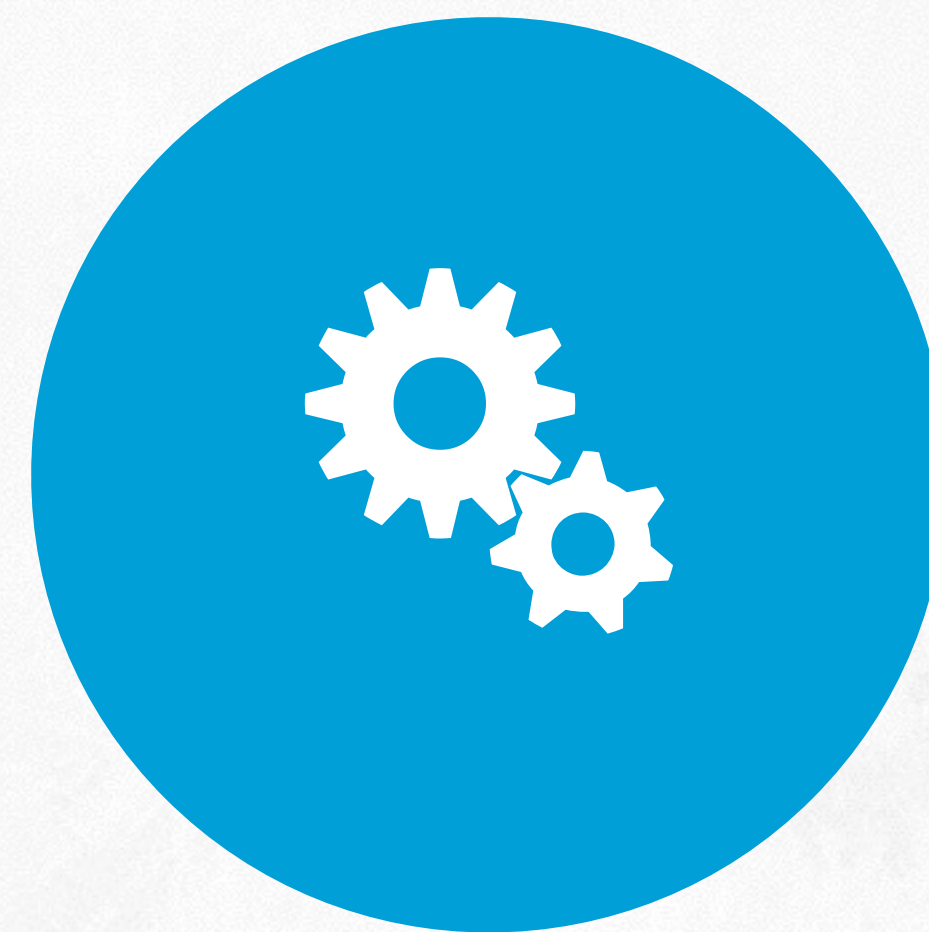
# TWITTER STORM

**Streaming platform** for analyzing **realtime data** as they arrive, so you can react to data **as it happens**.
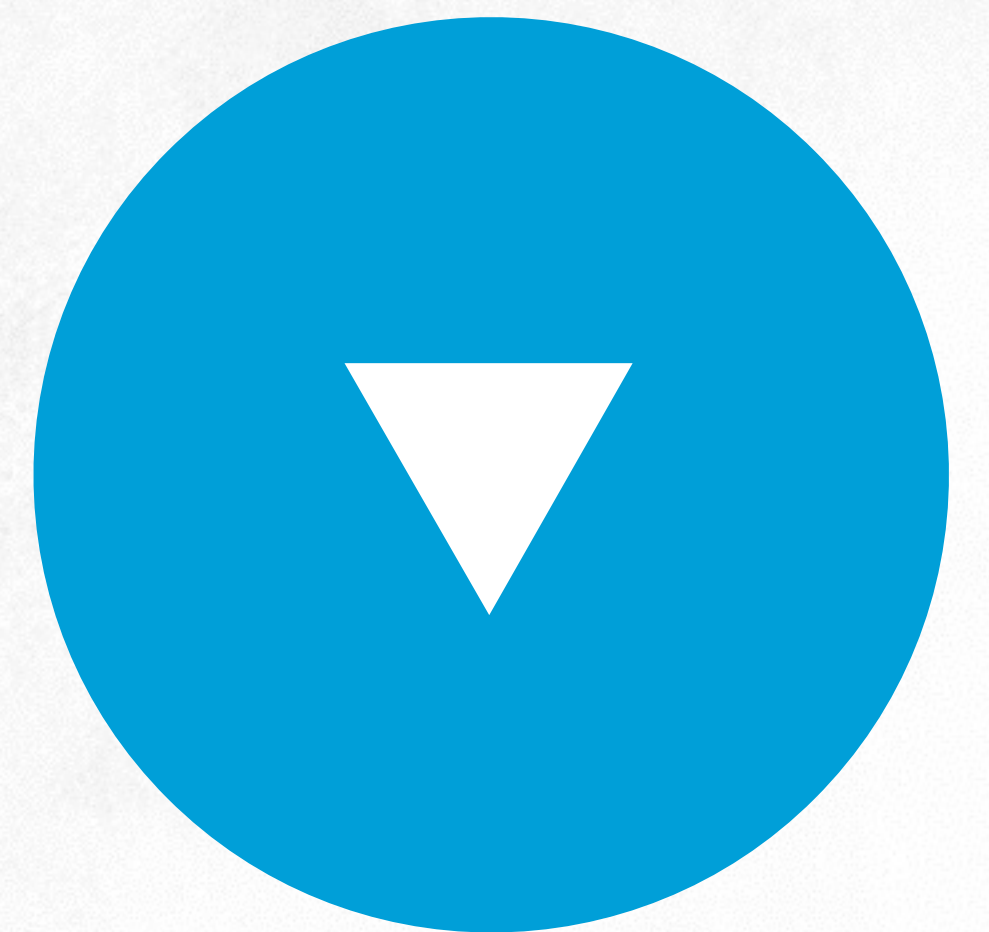
**GUARANTEED MESSAGE PROCESSING**

**HORIZONTAL SCALABILITY**

**ROBUST FAULT TOLERANCE**

**CONCISE CODE– FOCUS ON LOGIC**

# STORM TERMINOLOGY

### TOPOLOGY

Directed acyclic graph

Vertices=computation, and edges=streams of data tuples

### SPOUTS

Sources of data tuples for the topology
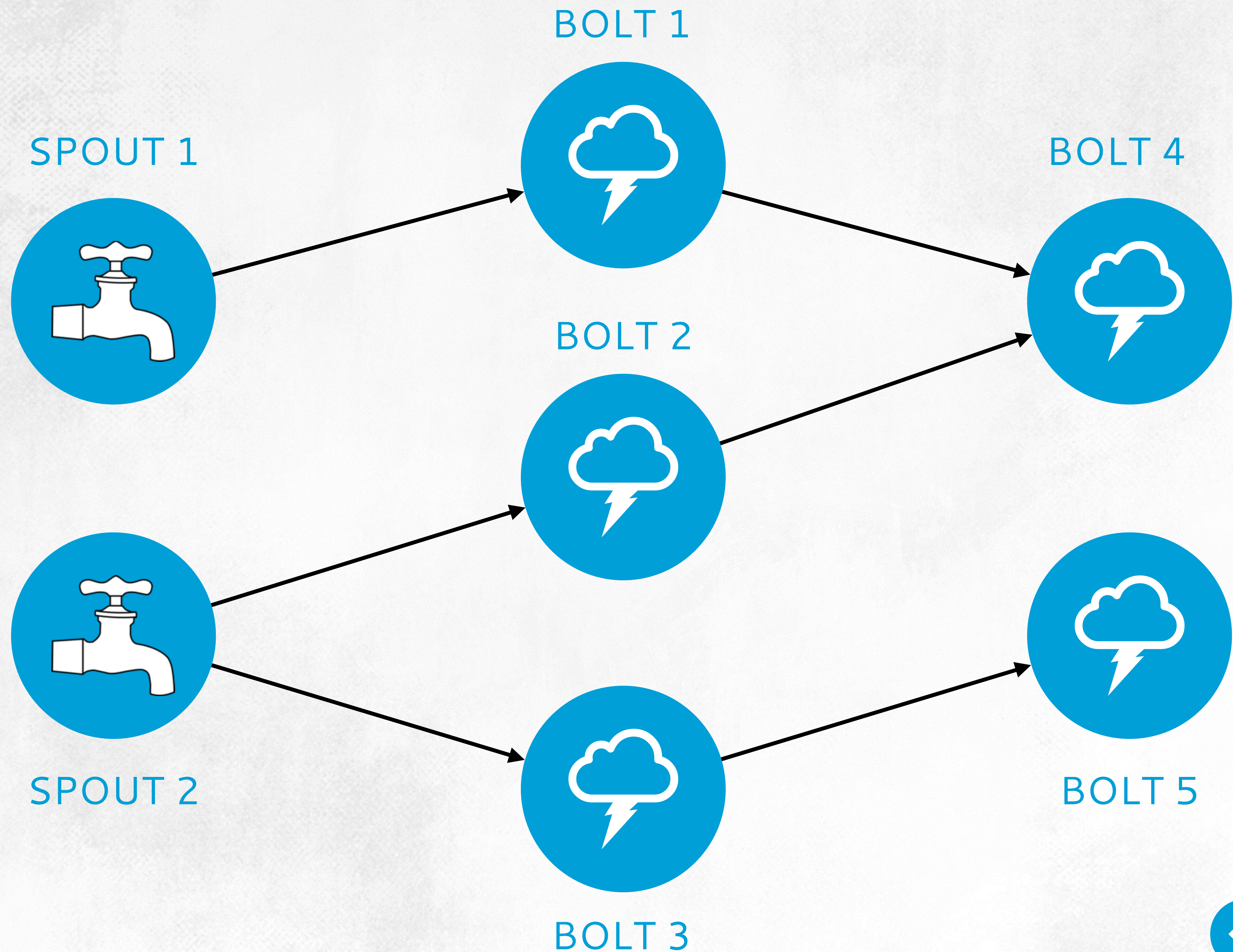
Examples – Kafka/Kestrel/MySQL/Postgres

### BOLTS

Process incoming tuples and emit outgoing tuples

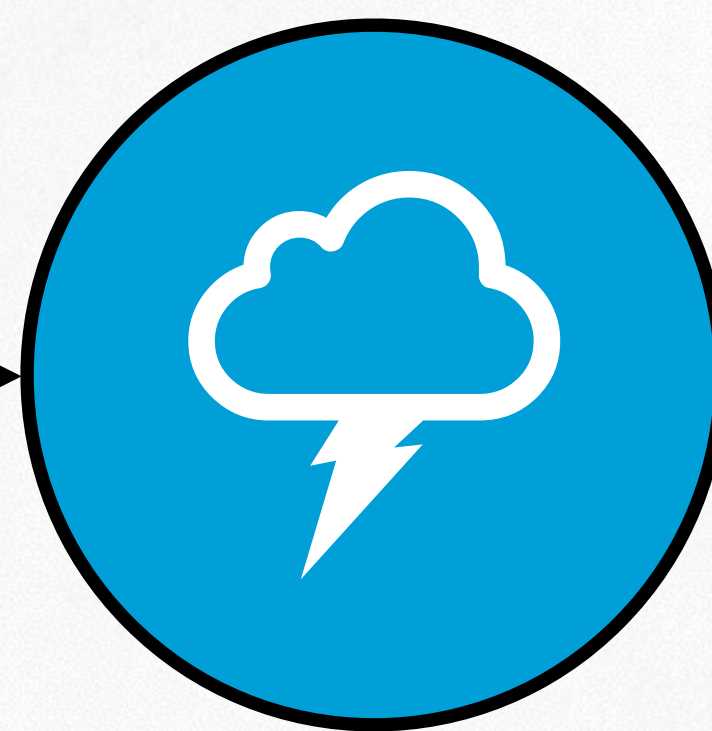Examples – filtering/aggregation/join/arbitrary function

# STORM TOPOLOGY

SPOUT 1

BOLT 1

BOLT 2

BOLT 4

SPOUT 2

BOLT 3

BOLT 5

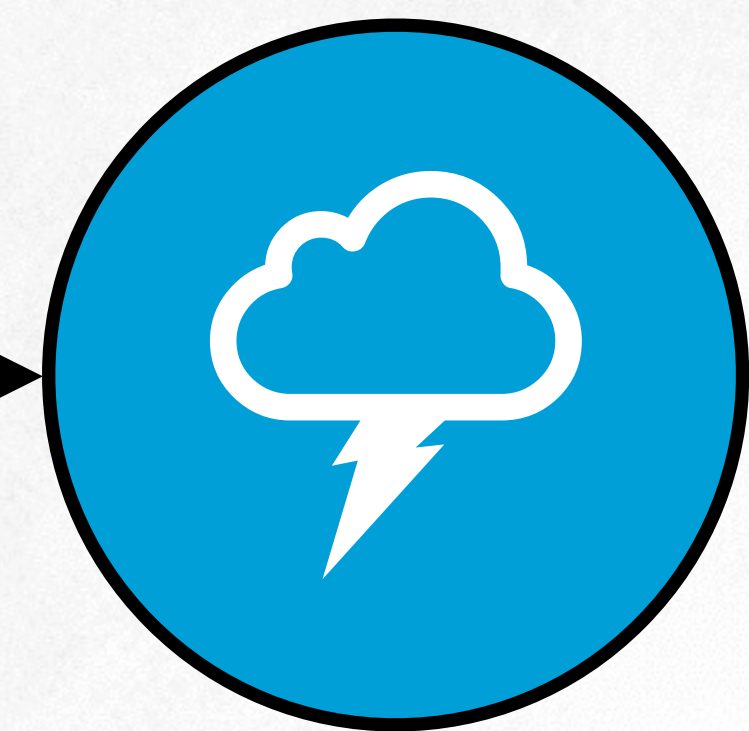# WORD COUNT TOPOLOGY

Live stream of Tweets



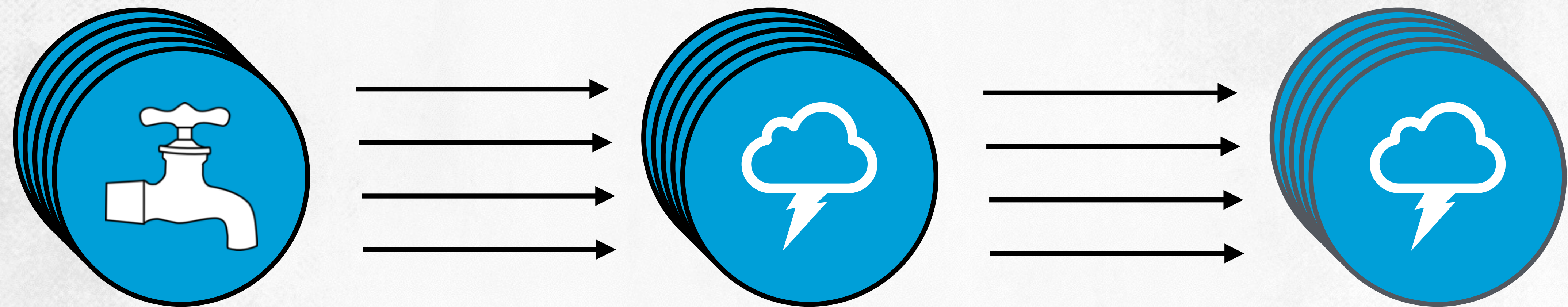TWEET SPOUT       PARSE TWEET BOLT       WORD COUNT BOLT

## LOGICAL PLAN

# WORD COUNT TOPOLOGY



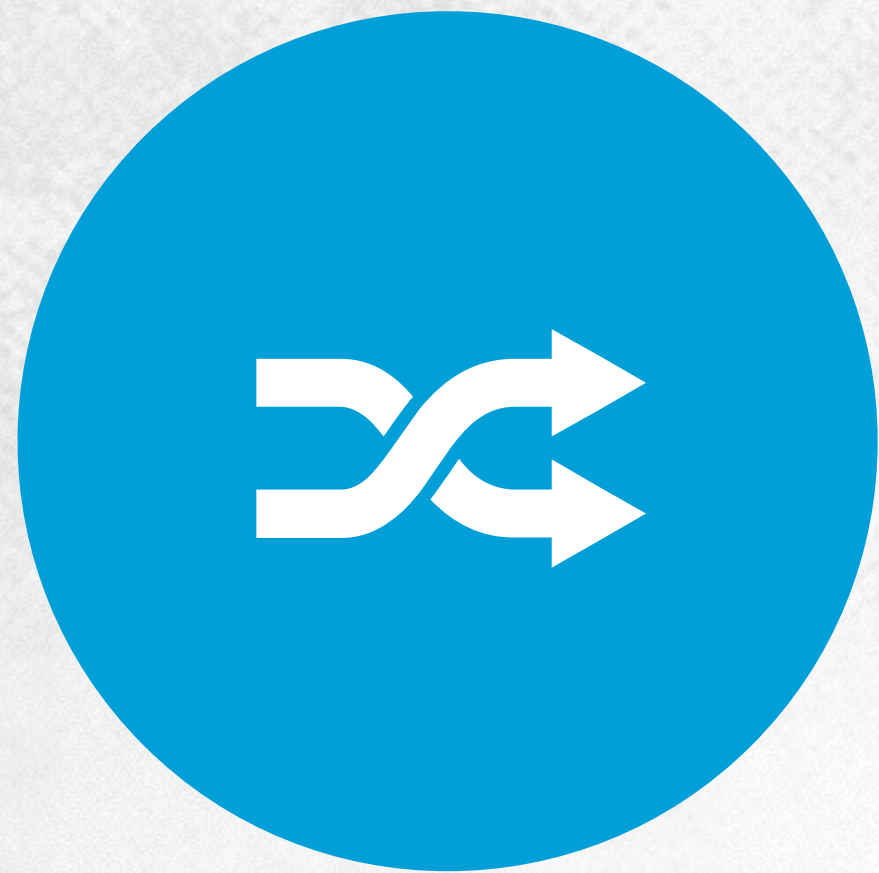**TWEET SPOUT TASKS**

**PARSE TWEET BOLT TASKS**

**WORD COUNT BOLT TASKS**

When a parse tweet bolt task emits a tuple which word count bolt task should it send to?
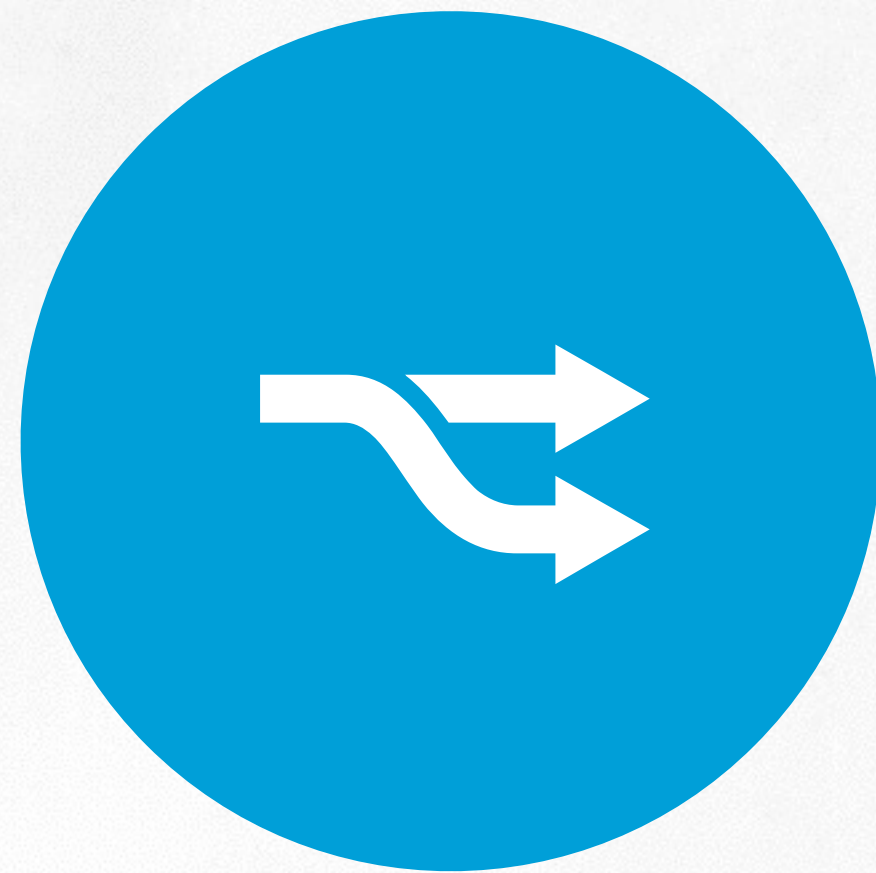
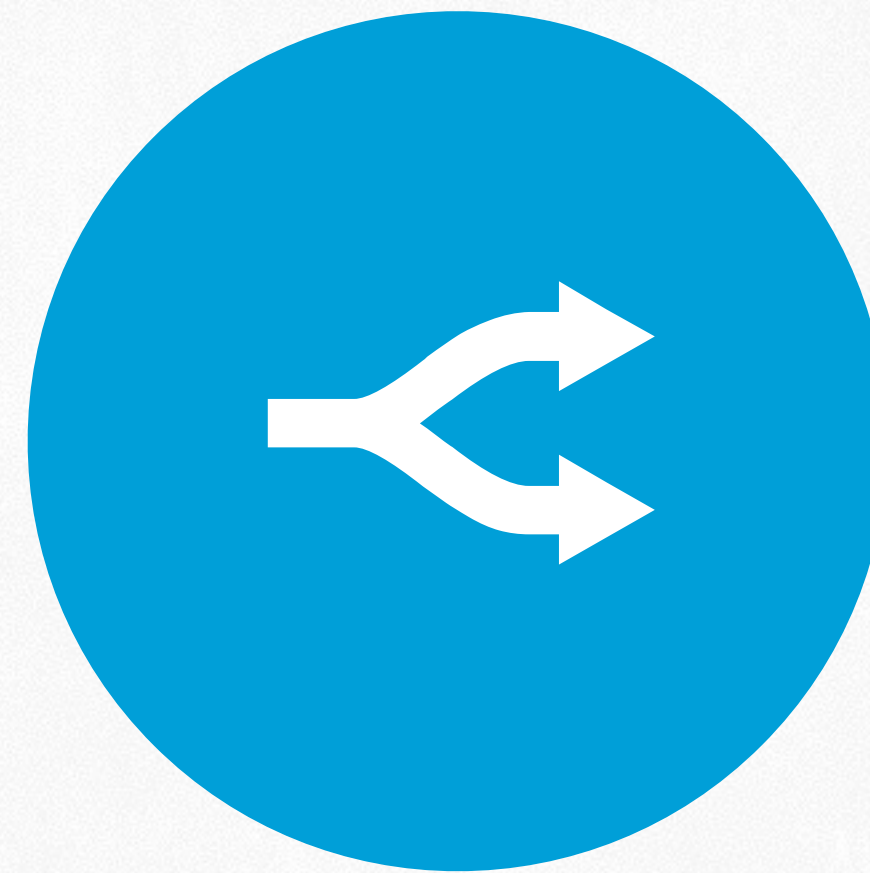# STREAM GROUPINGS

## SHUFFLE GROUPING
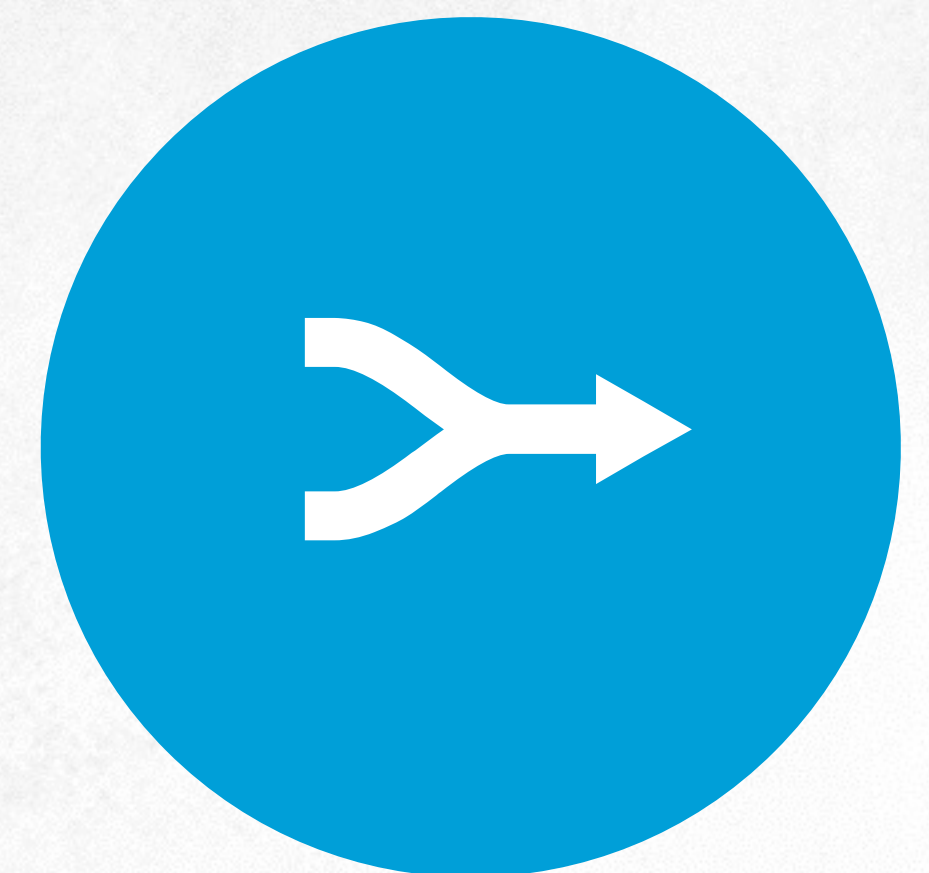
Random distribution of tuples

## FIELDS GROUPING

Group tuples by a field or multiple fields
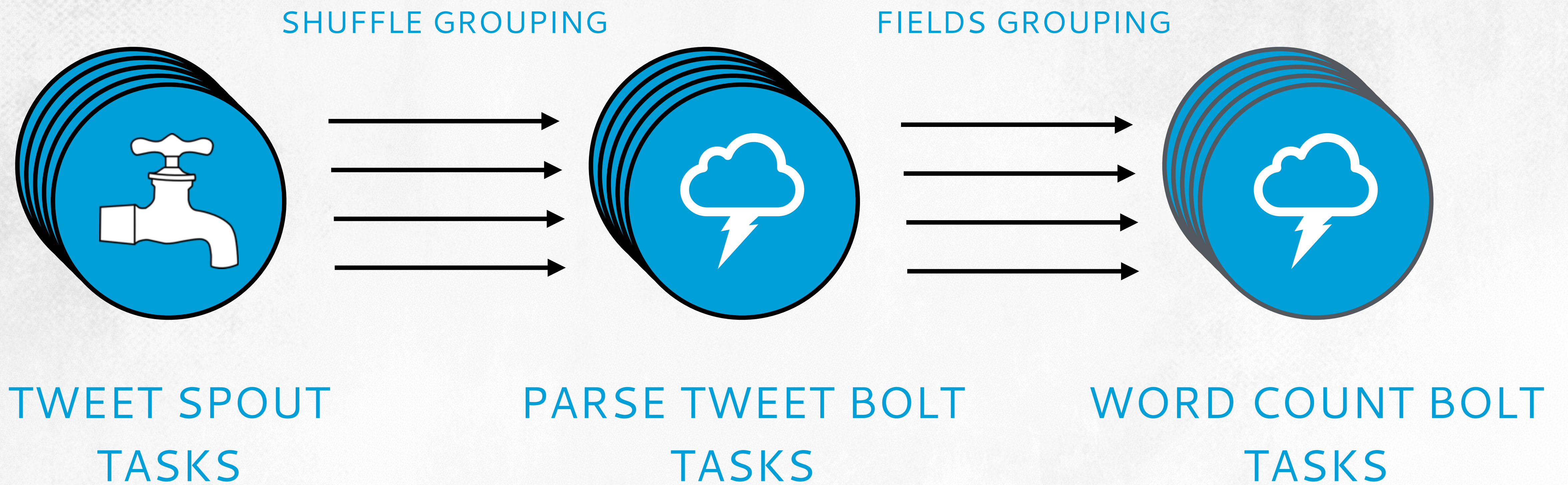
## ALL GROUPING

Replicates tuples to all tasks

## GLOBAL GROUPING
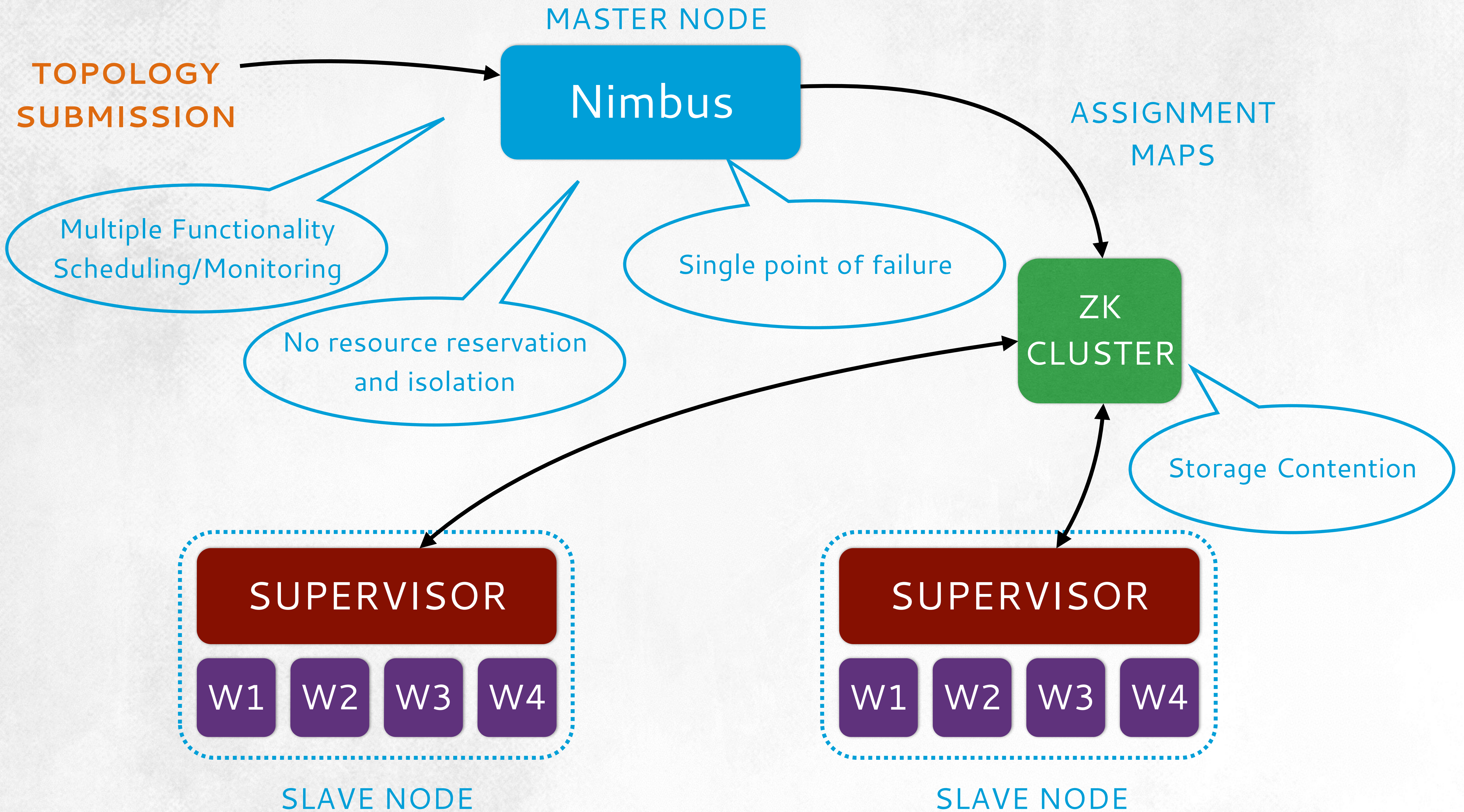
Sends the entire stream to one task

# WORD COUNT TOPOLOGY
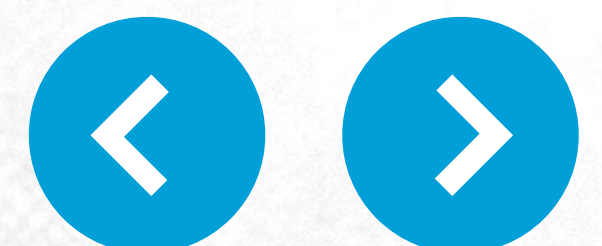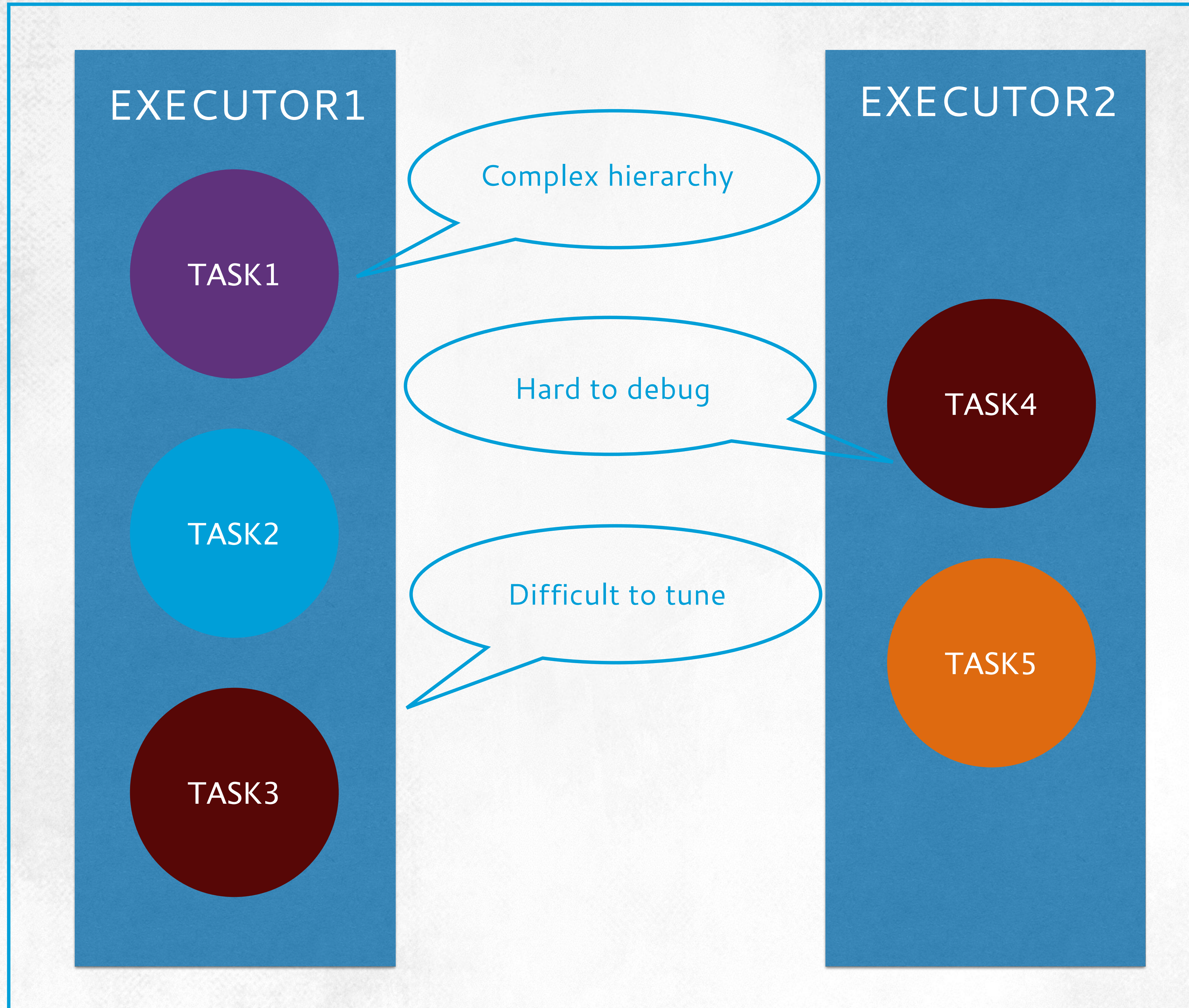
SHUFFLE GROUPING

FIELDS GROUPING

TWEET SPOUT
TASKS

PARSE TWEET BOLT
TASKS

WORD COUNT BOLT
TASKS

MOTIVATION
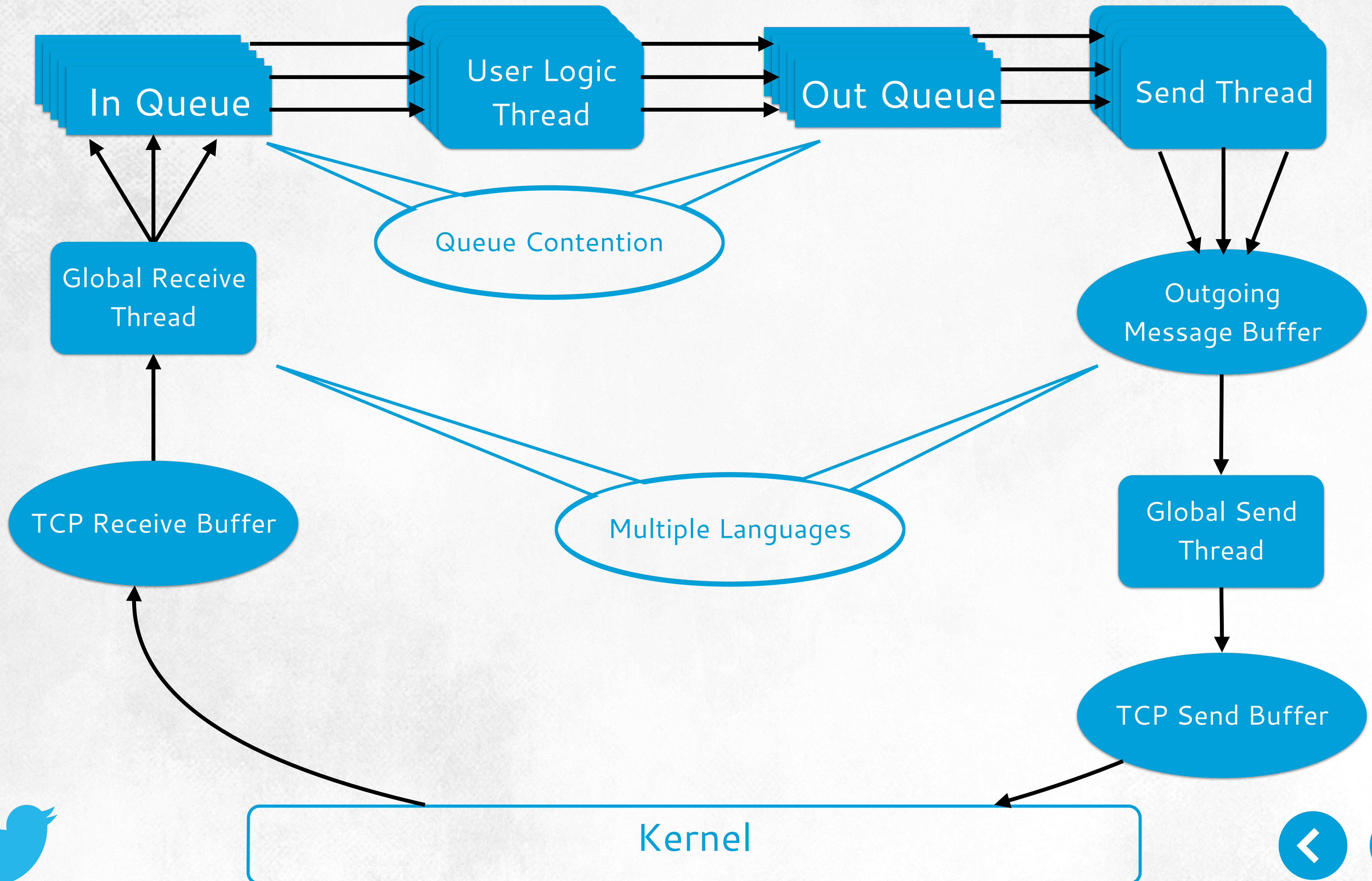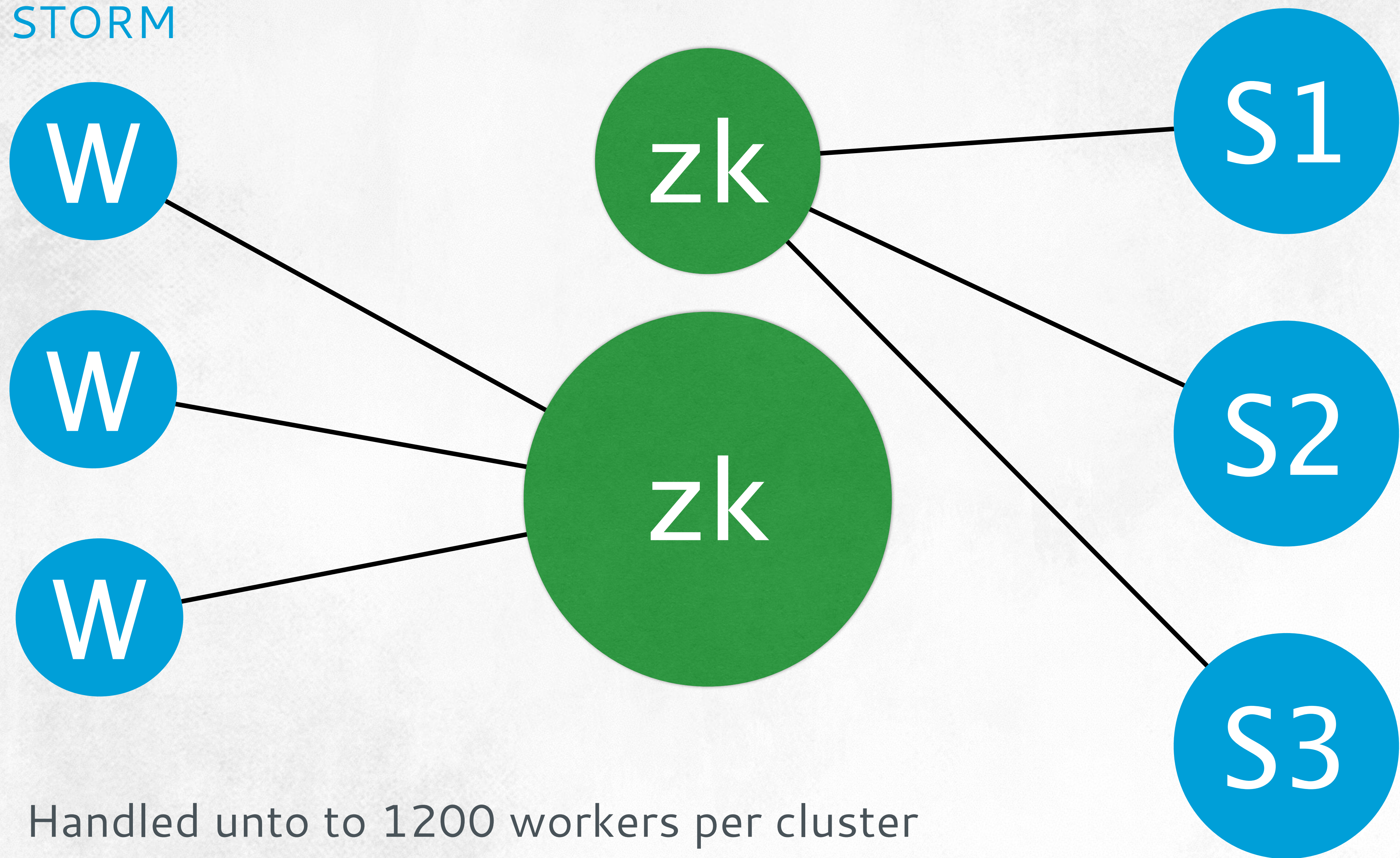
# DATA FLOW IN STORM WORKERS

# OVERLOADED ZOOKEEPER

## Scaled up

STORM



Handled unto to 1200 workers per cluster

# OVERLOADED ZOOKEEPER

## Analyzing zookeeper traffic

**67%**

### KAFKA SPOUT

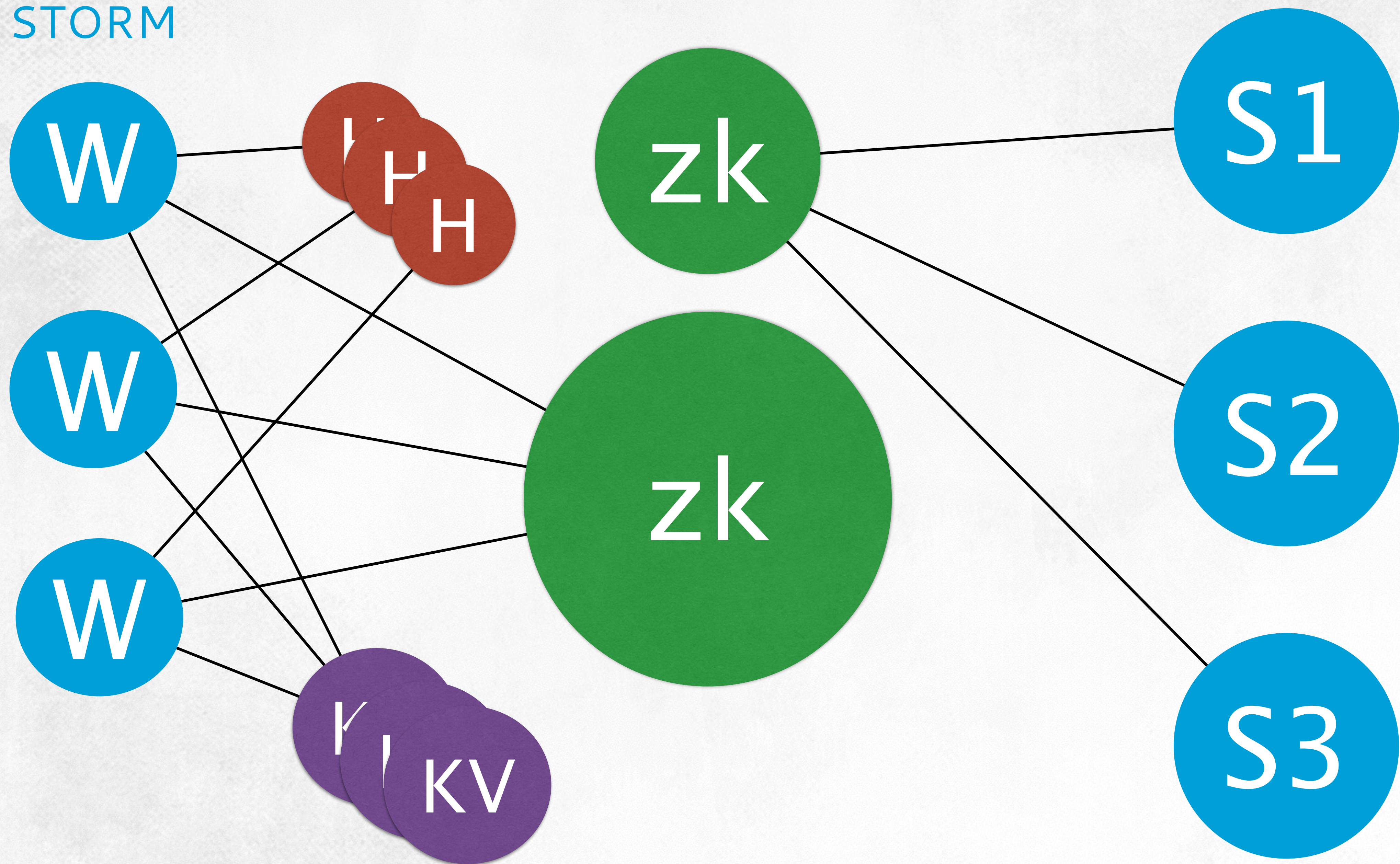Offset/partition is written every 2 secs

**33%**

### STORM RUNTIME

Workers write heart beats every 3 secs

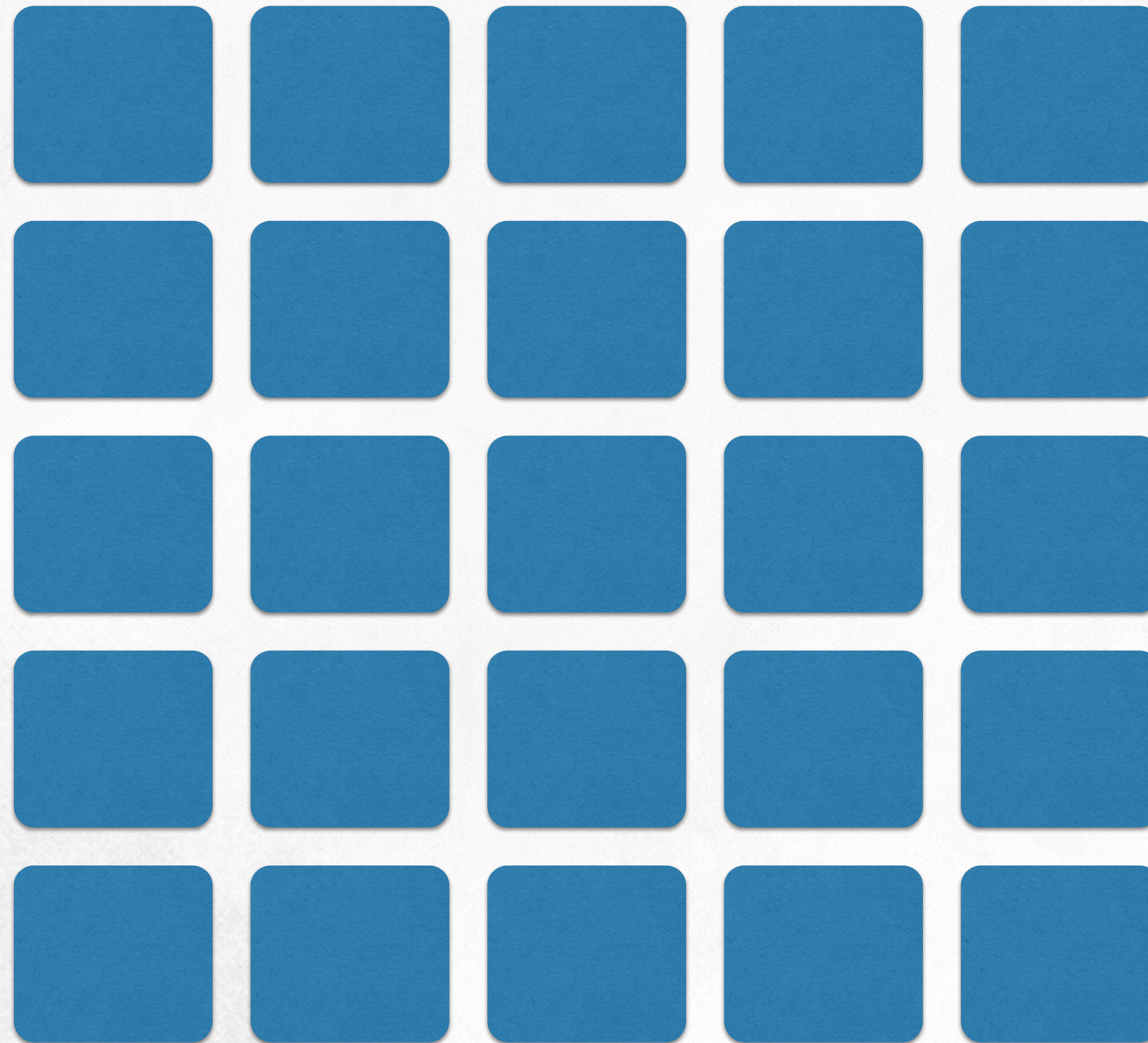# OVERLOADED ZOOKEEPER

## Heart beat daemons

STORM

W  W  W  H H H  zk  zk  KV KV KV  S1  S2  S3

5000 workers per cluster

# STORM – DEPLOYMENT

shared pool

storm cluster

# STORM – DEPLOYMENT
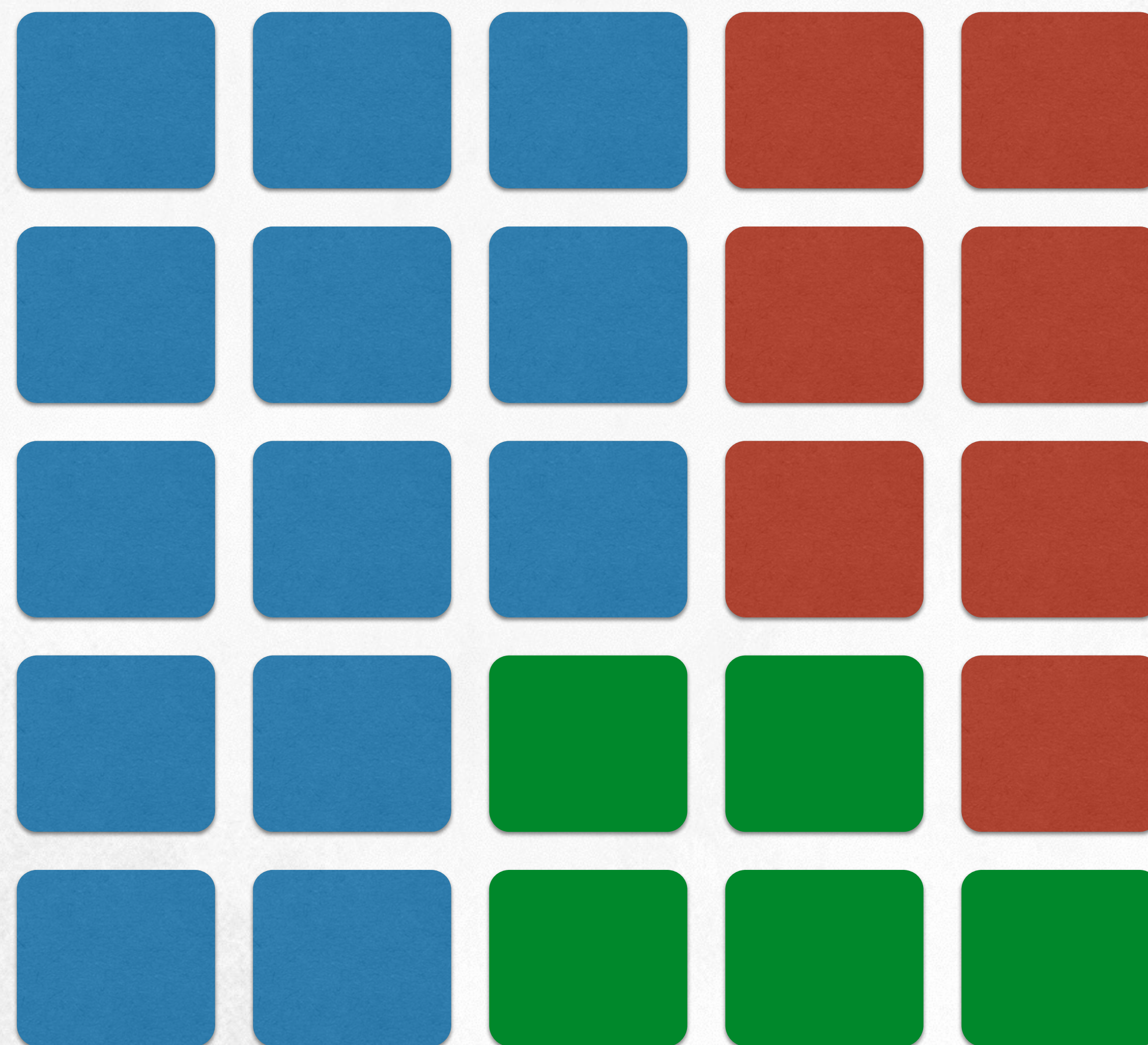
shared pool

isolated pools

joe's topology

storm
cluster

# STORM – DEPLOYMENT

shared pool

isolated pools

joe's topology

jane's topology

storm
cluster
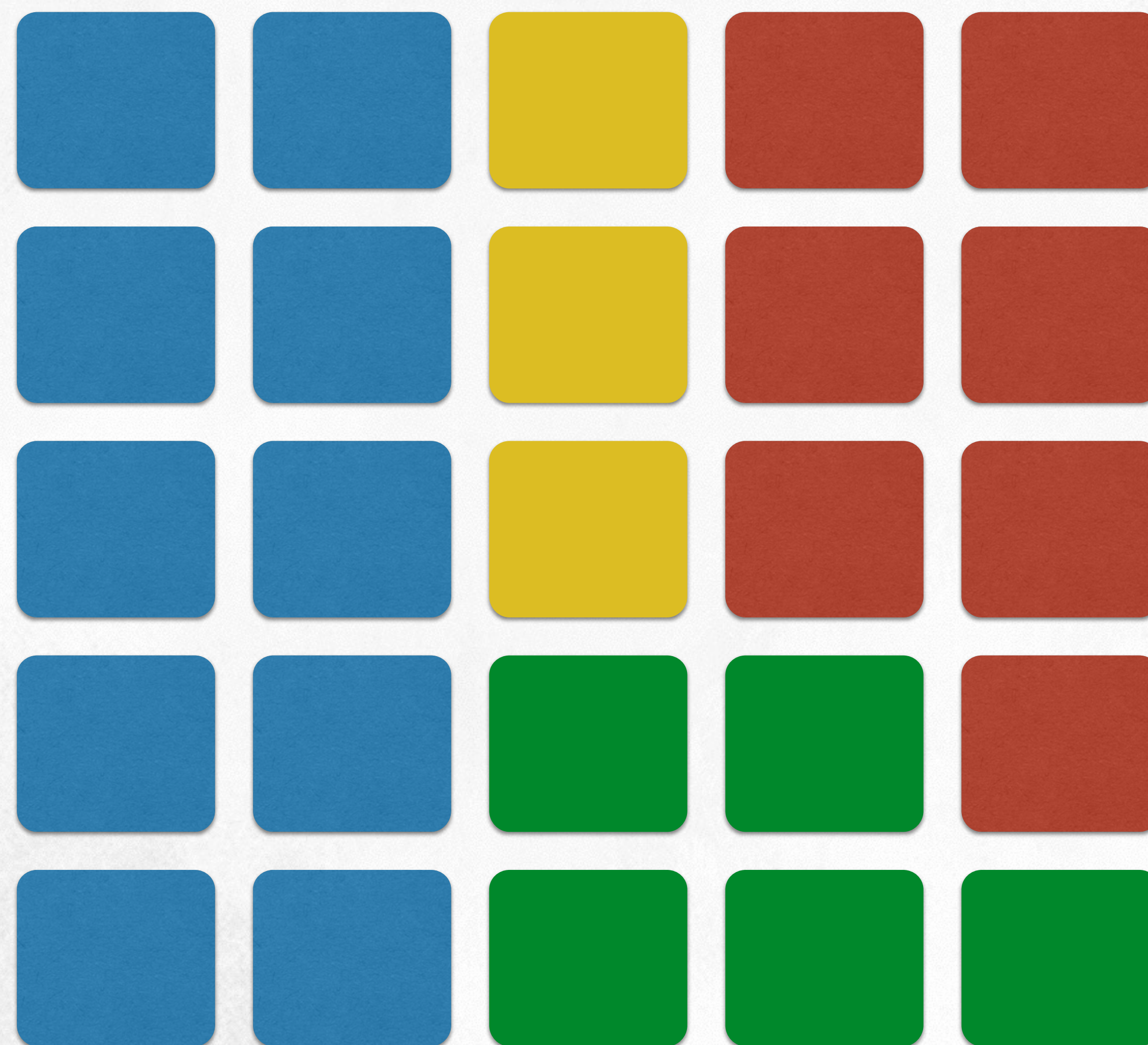
# STORM – DEPLOYMENT
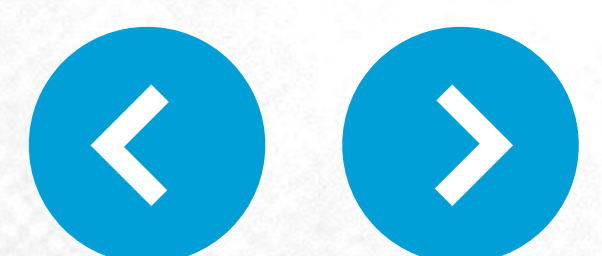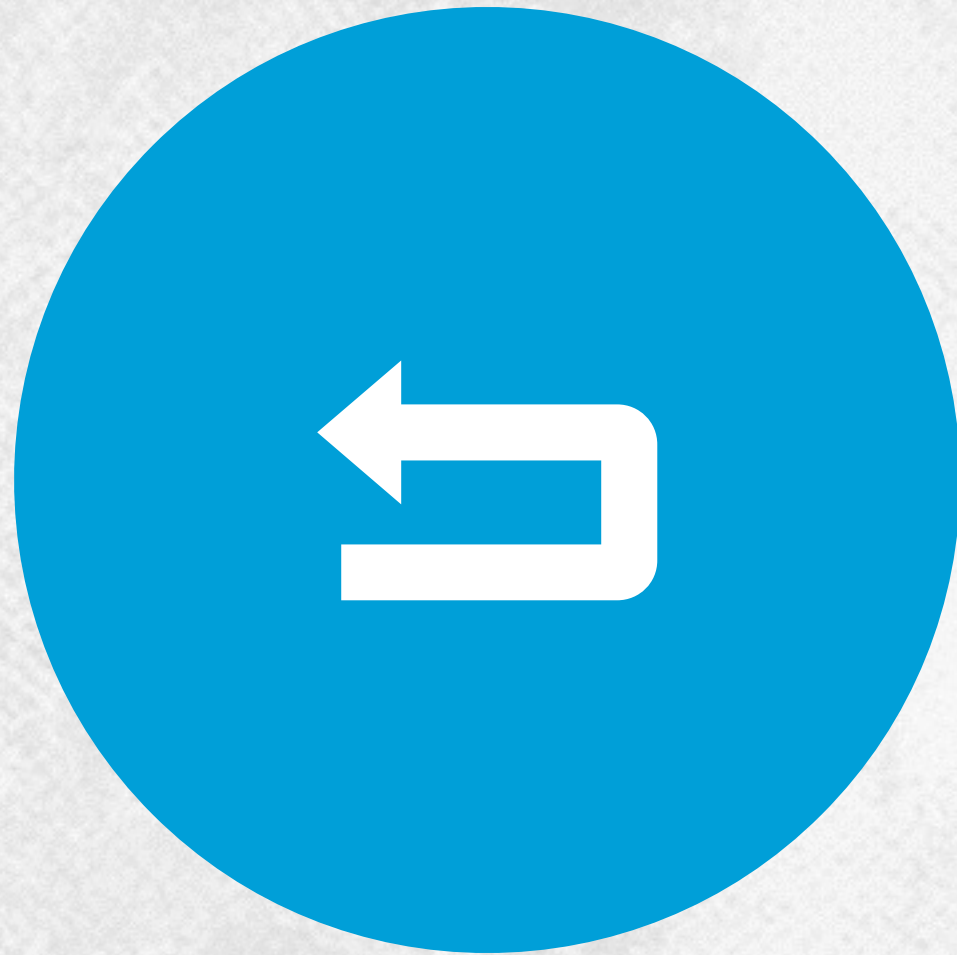
shared pool

storm cluster

isolated pools

joe's topology
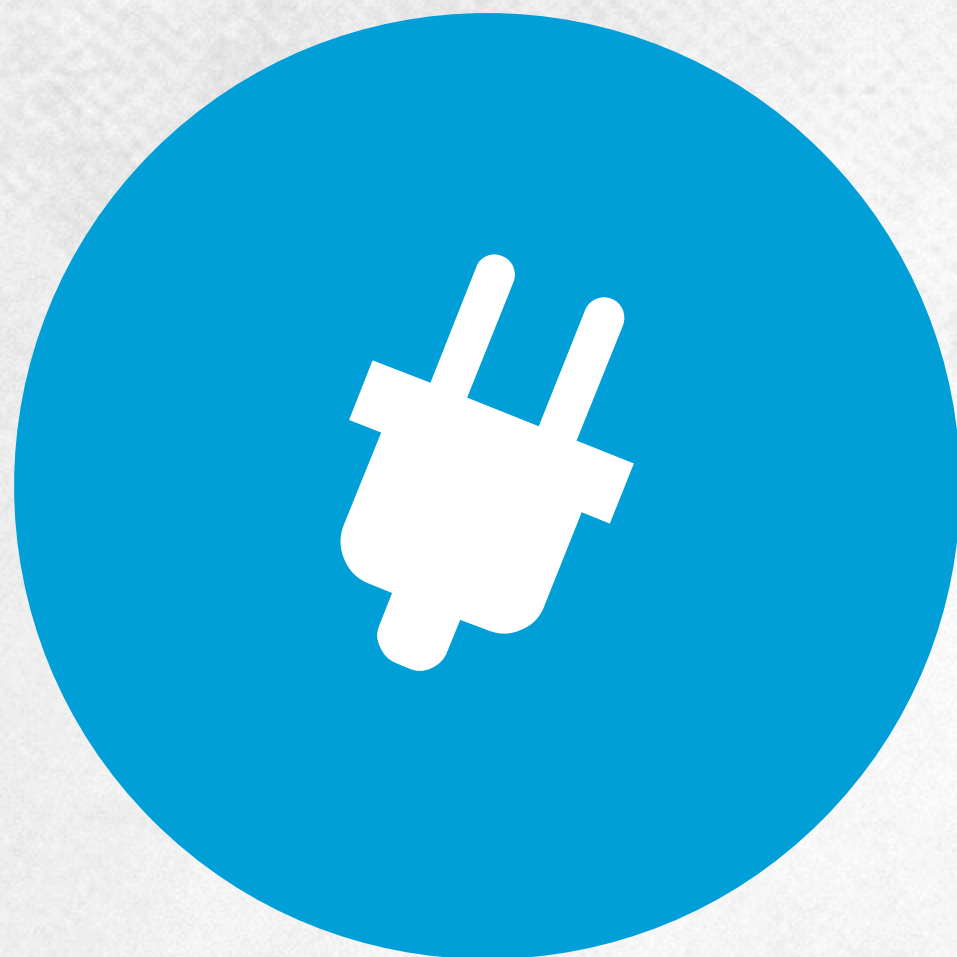
jane's topology

dave's topology

# STORM ISSUES

### LACK OF BACK PRESSURE

Drops tuples unpredictably

### EFFICIENCY

Serialization program consumes 75 cores at 30% CPU

Topology consumes 600 cores at 20–30% CPU

### NO BATCHING

Tuple oriented system – implicit batching by 0MQ

# EVOLUTION OR REVOLUTION?

## fix storm or develop a new system?

**FUNDAMENTAL ISSUES– REQUIRE EXTENSIVE REWRITING**

Several queues for moving data

Inflexible and requires longer development cycle

**USE EXISTING OPEN SOURCE SOLUTIONS**

Issues working at scale/lacks required performance

Incompatible API and long migration process

# HERON

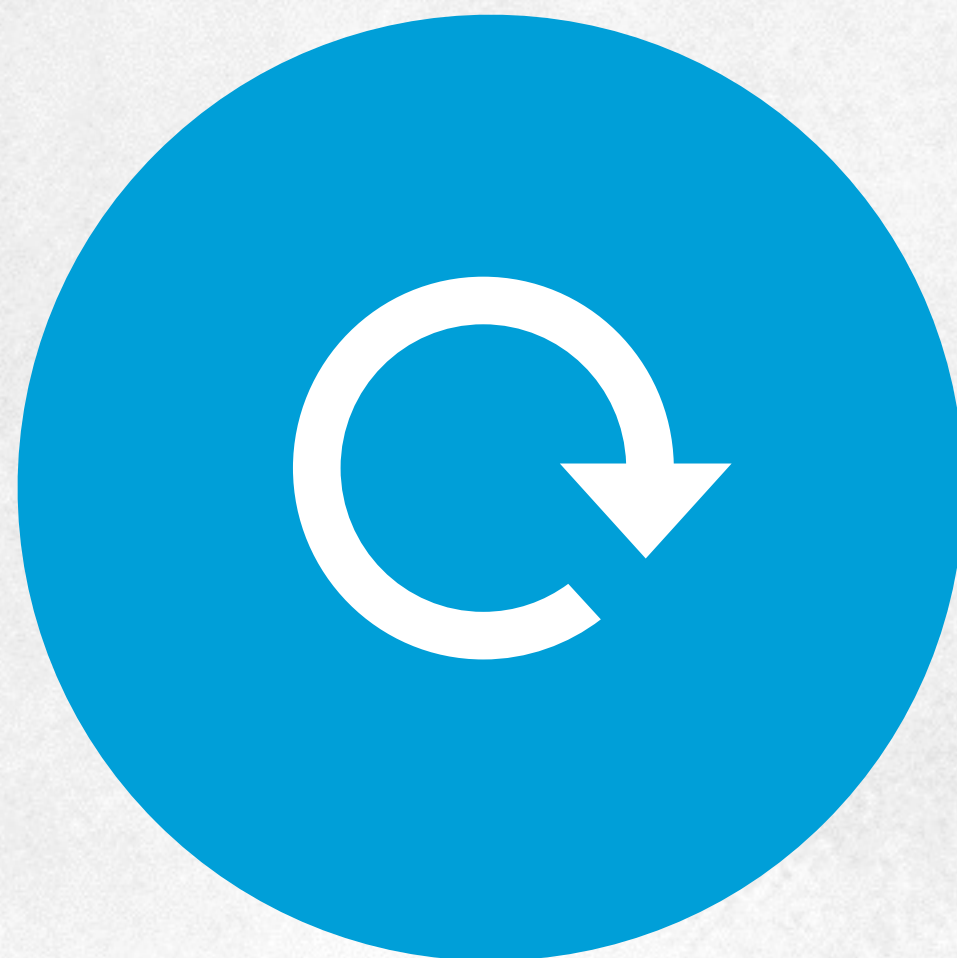# HERON DESIGN GOALS

## FULLY API COMPATIBLE WITH STORM

Directed acyclic graph

Topologies, spouts and bolts

## TASK ISOLATION
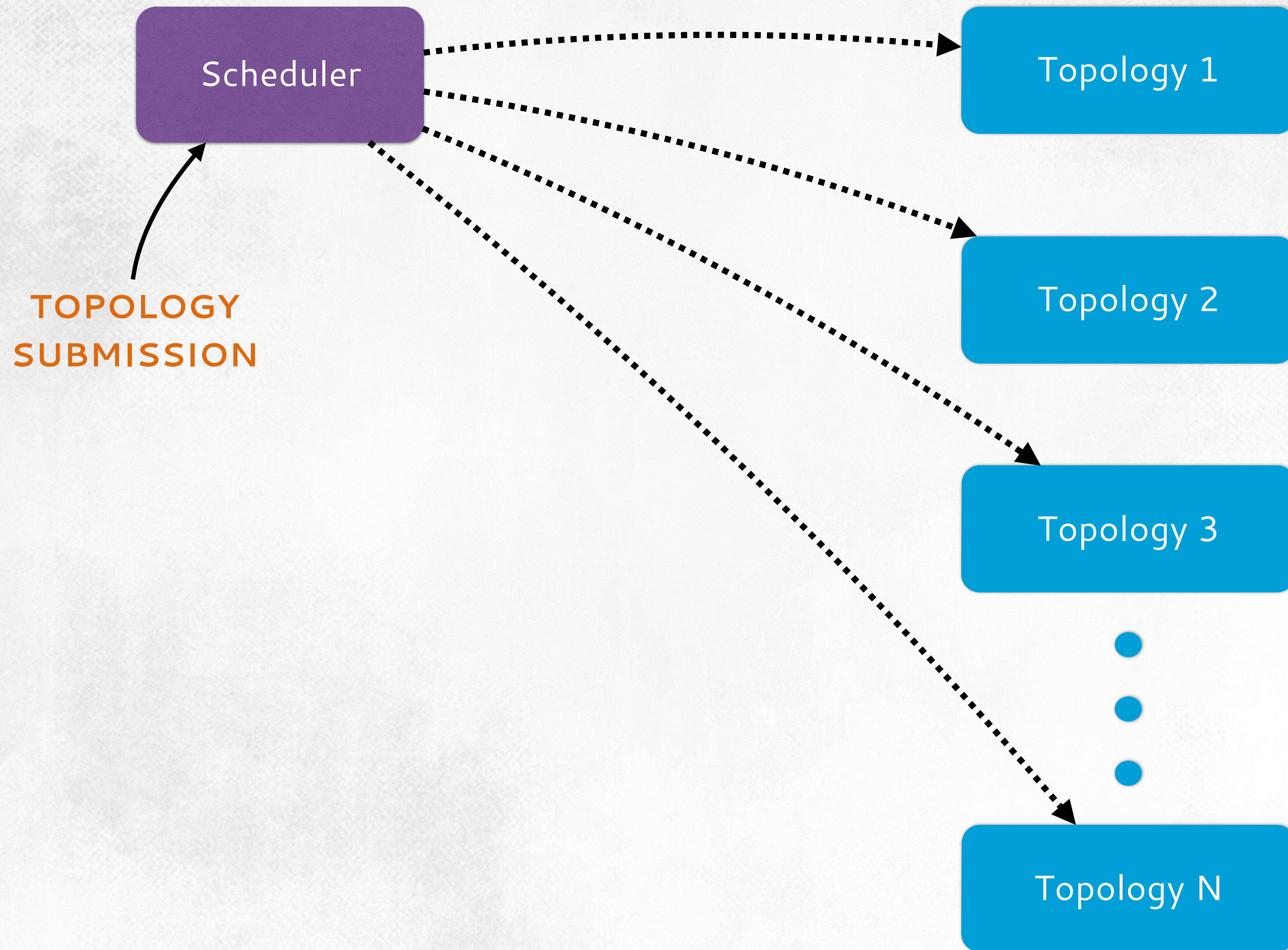
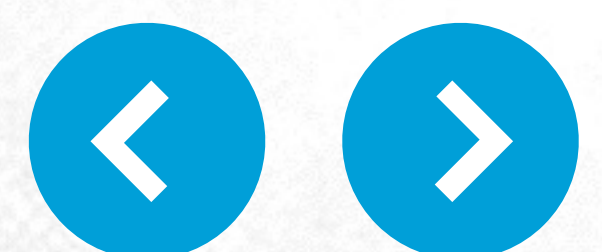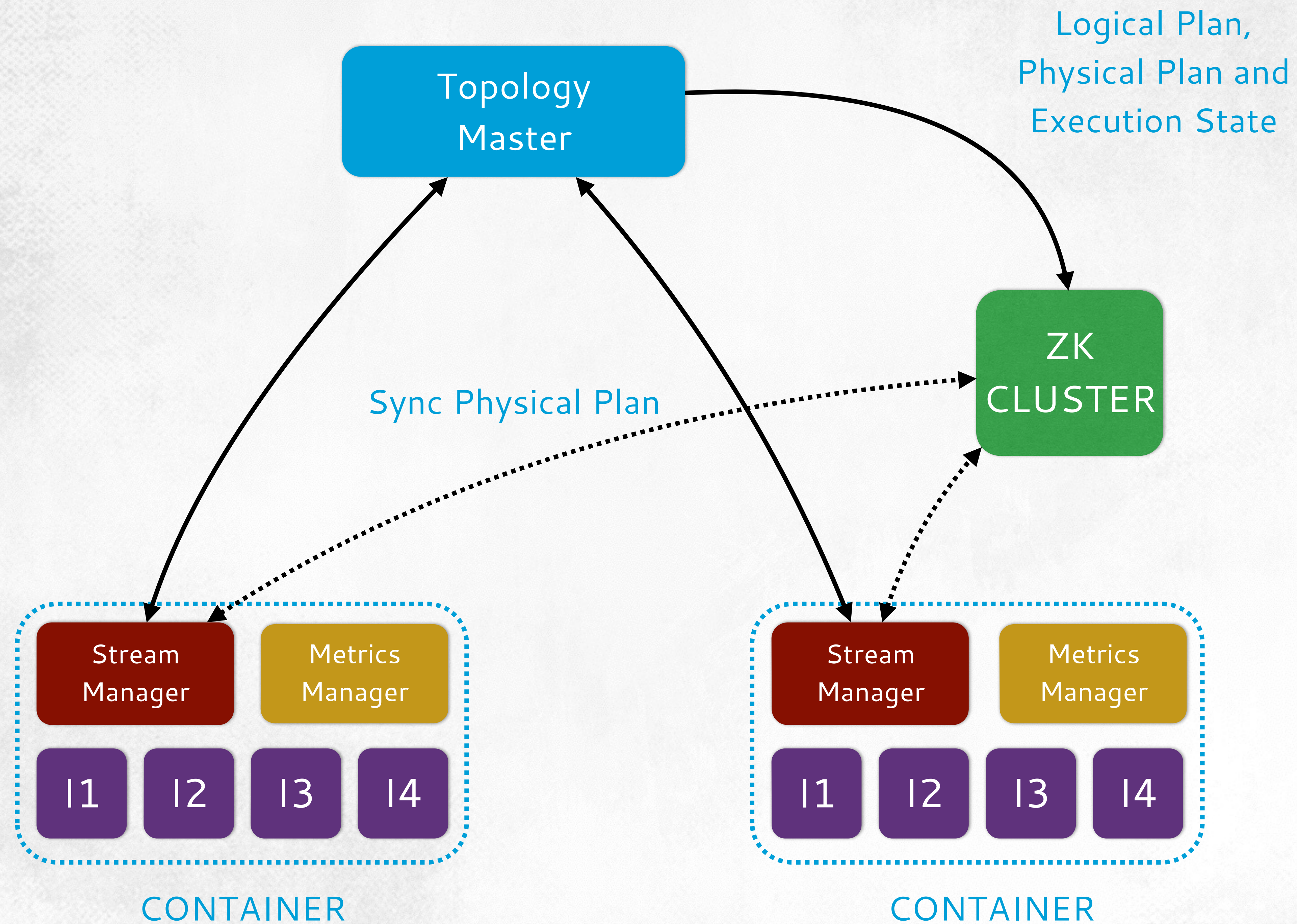Ease of debug ability/resource isolation/profiling

## USE OF MAIN STREAM LANGUAGES

C++/JAVA/Python

# TOPOLOGY ARCHITECTURE

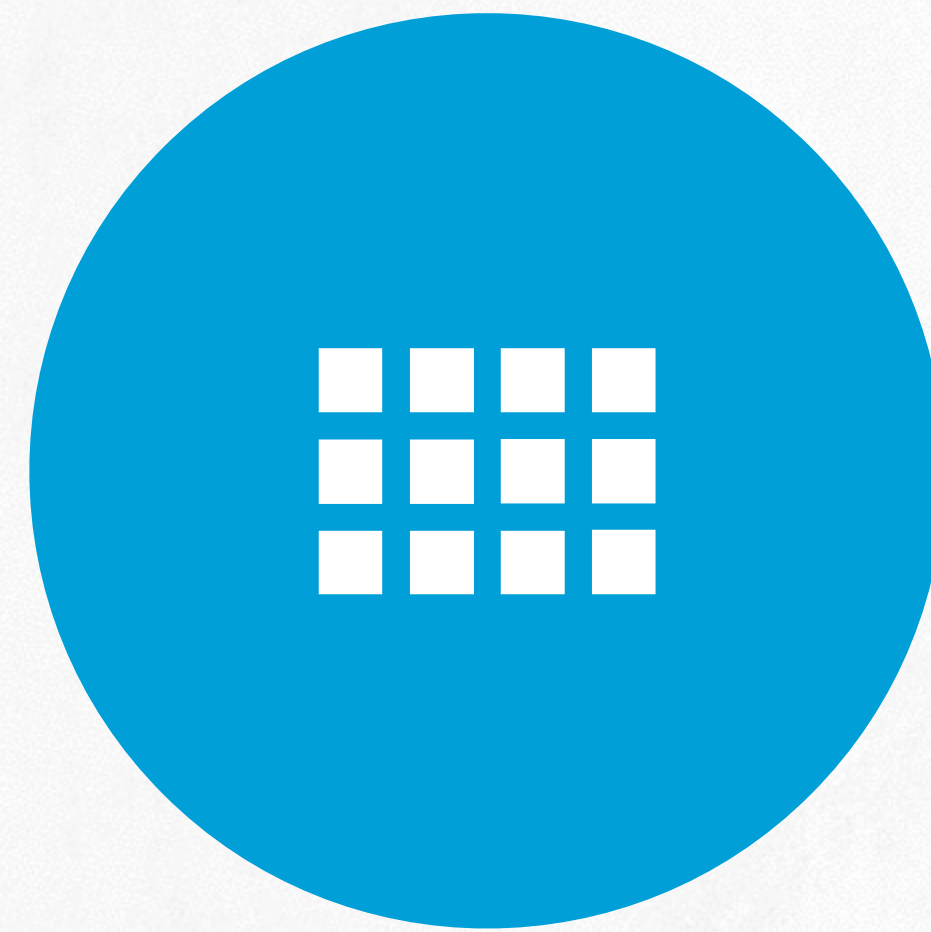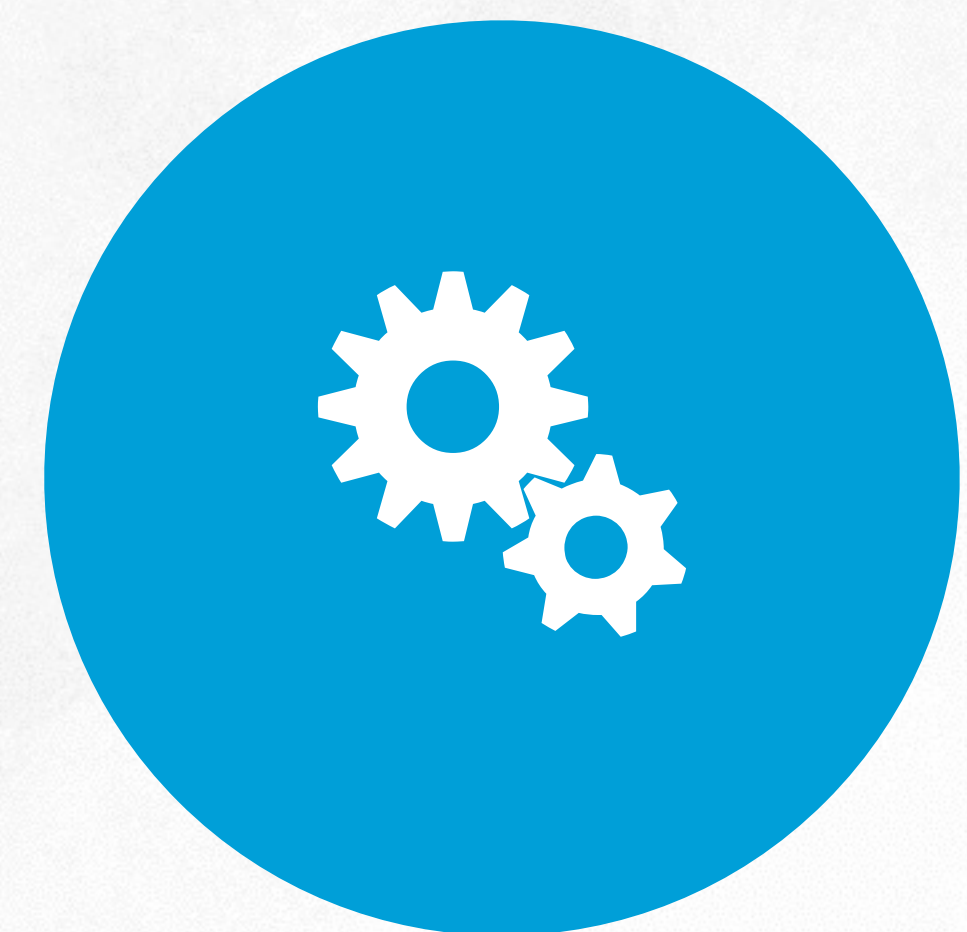# TOPOLOGY MASTER

## Solely responsible for the entire topology

**ASSIGNS ROLE**

**MONITORING**

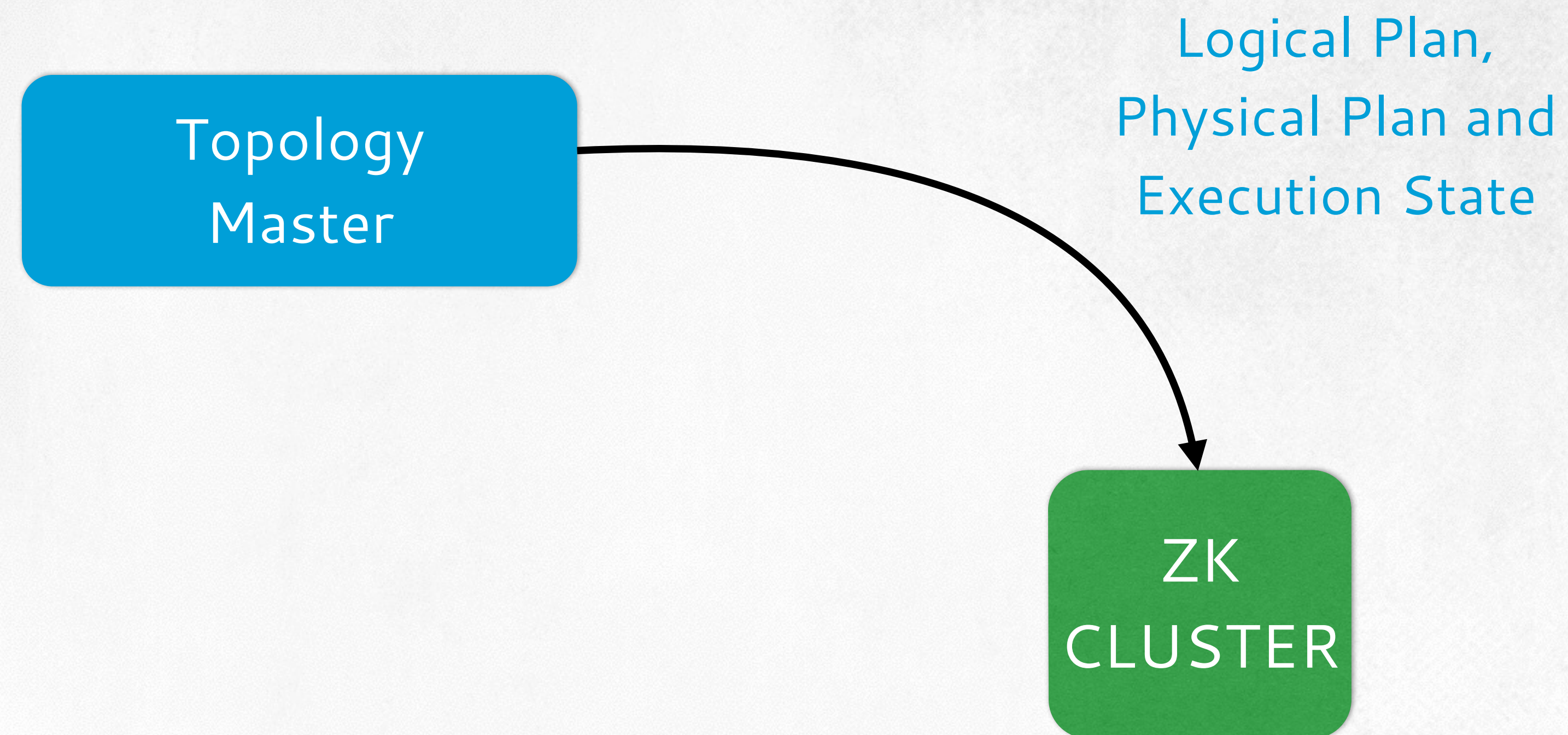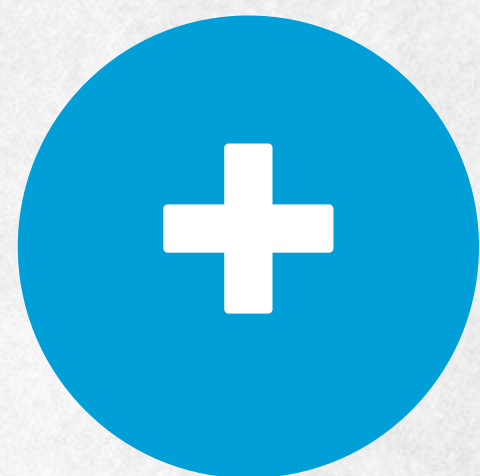**METRICS**

# TOPOLOGY **MASTER**

Topology
Master

Logical Plan,
Physical Plan and
Execution State
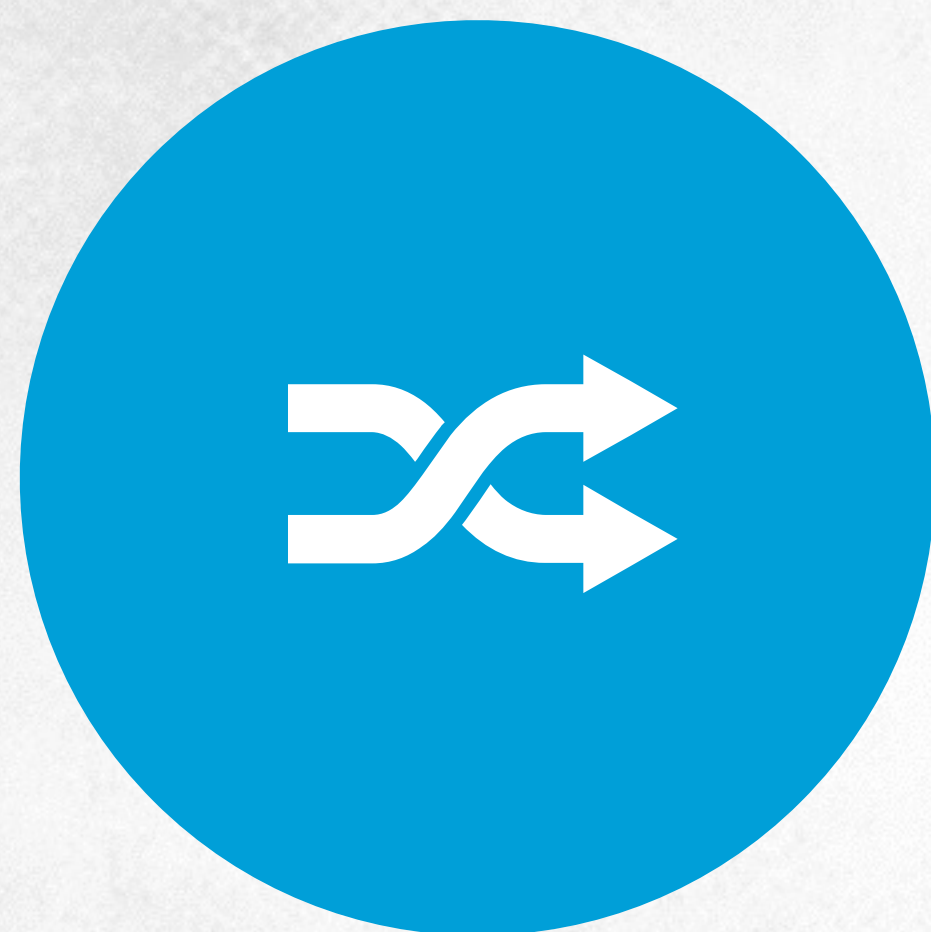
ZK
CLUSTER

**+** PREVENT MULTIPLE TM BECOMING MASTERS

**+** ALLOWS OTHER PROCESS TO DISCOVER TM
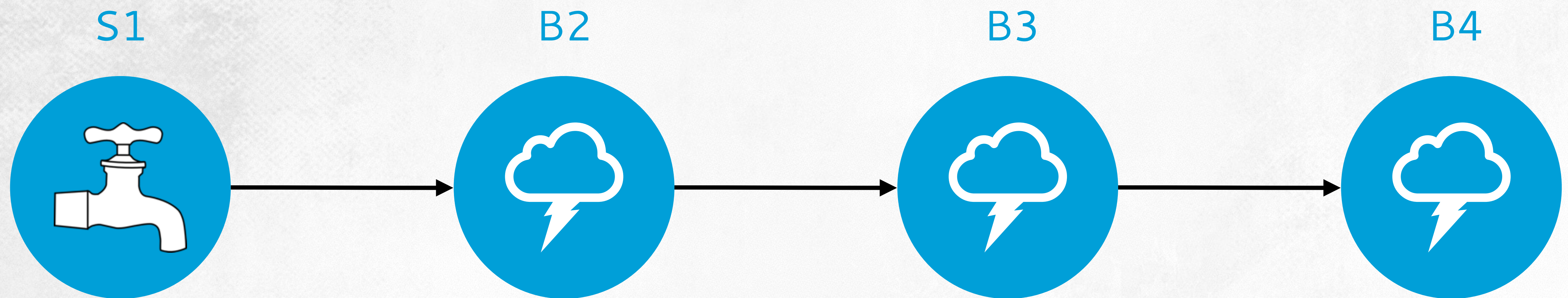
# STREAM MANAGER

## Routing Engine

**ROUTES TUPLES**

**BACK PRESSURE**

**ACK MGMT**

# STREAM MANAGER

S1          B2          B3          B4
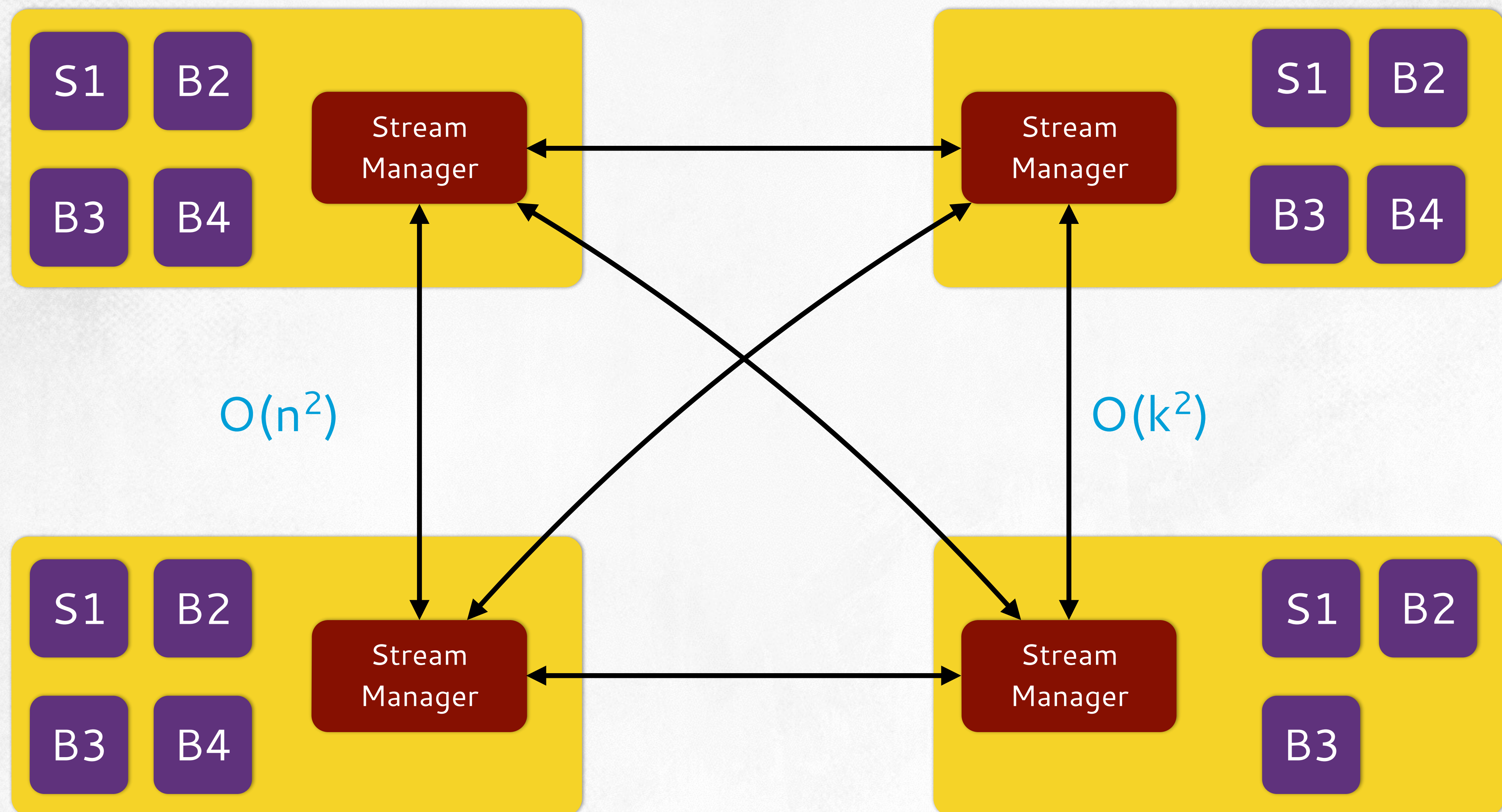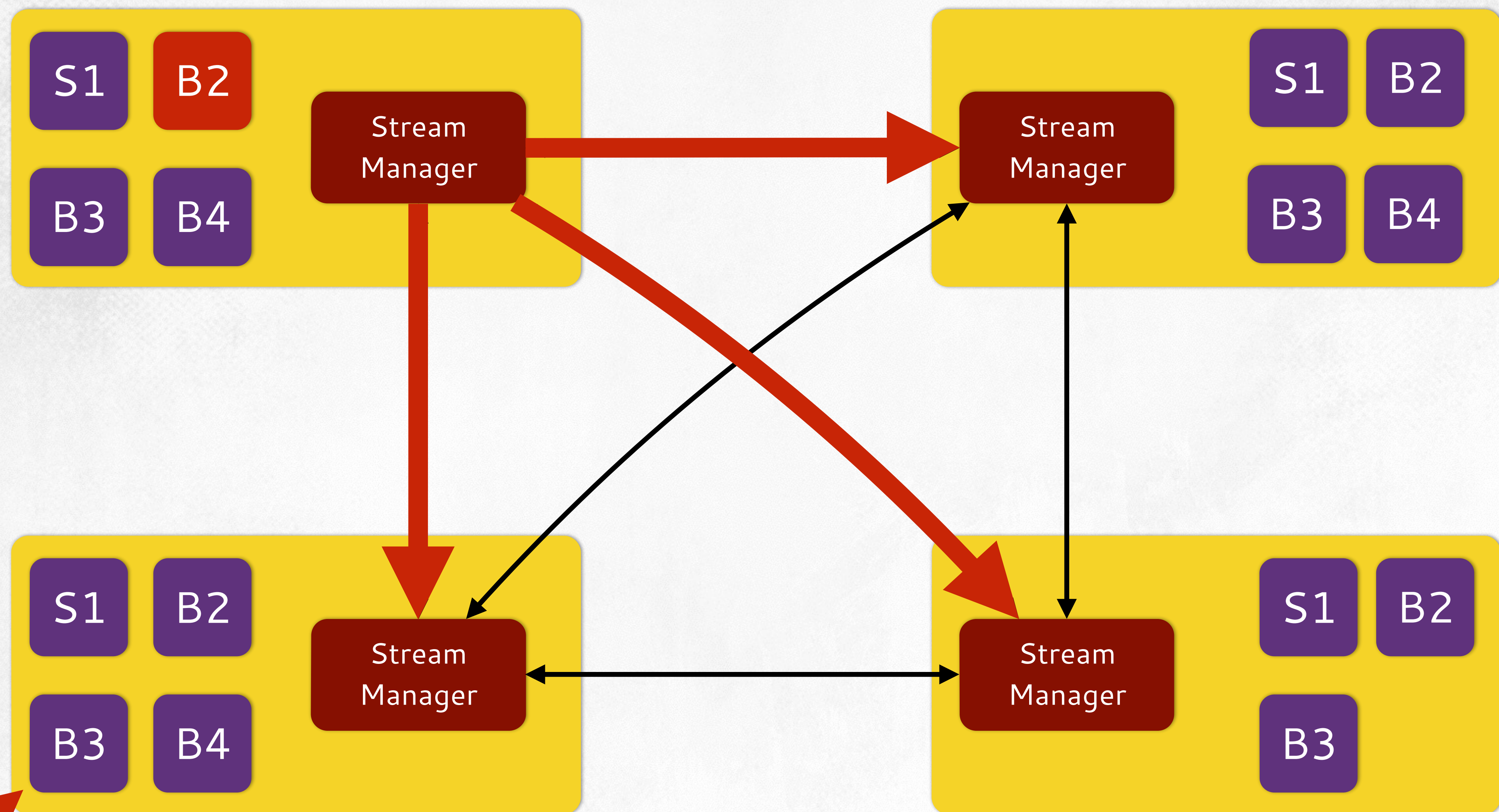
# STREAM MANAGER

# STREAM MANAGER
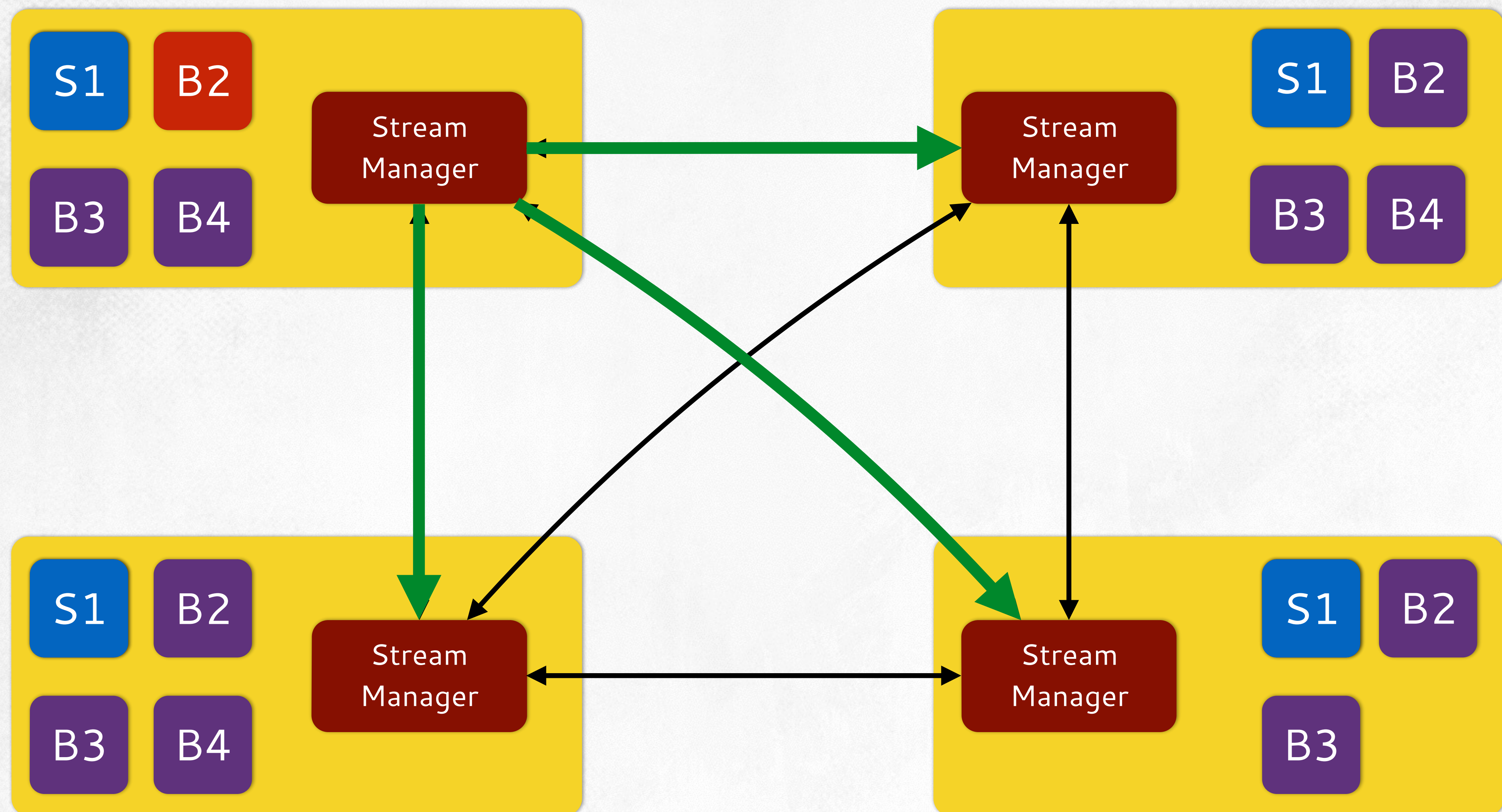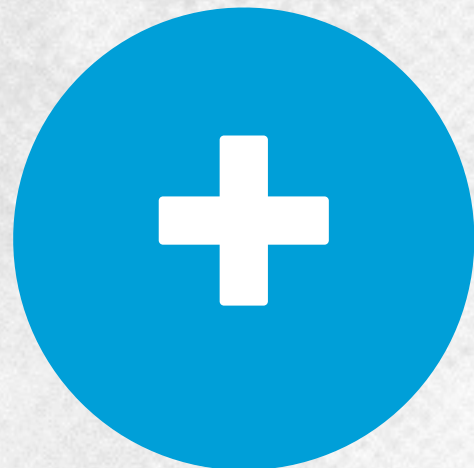## tcp back pressure

SLOWS UPSTREAM AND DOWNSTREAM INSTANCES

# STREAM MANAGER
spout back pressure

# STREAM MANAGER

## back pressure advantages

**+** **PREDICTABILITY**

Tuple failures are more deterministic

**+** **SELF ADJUSTS**

Topology goes as fast as the slowest component
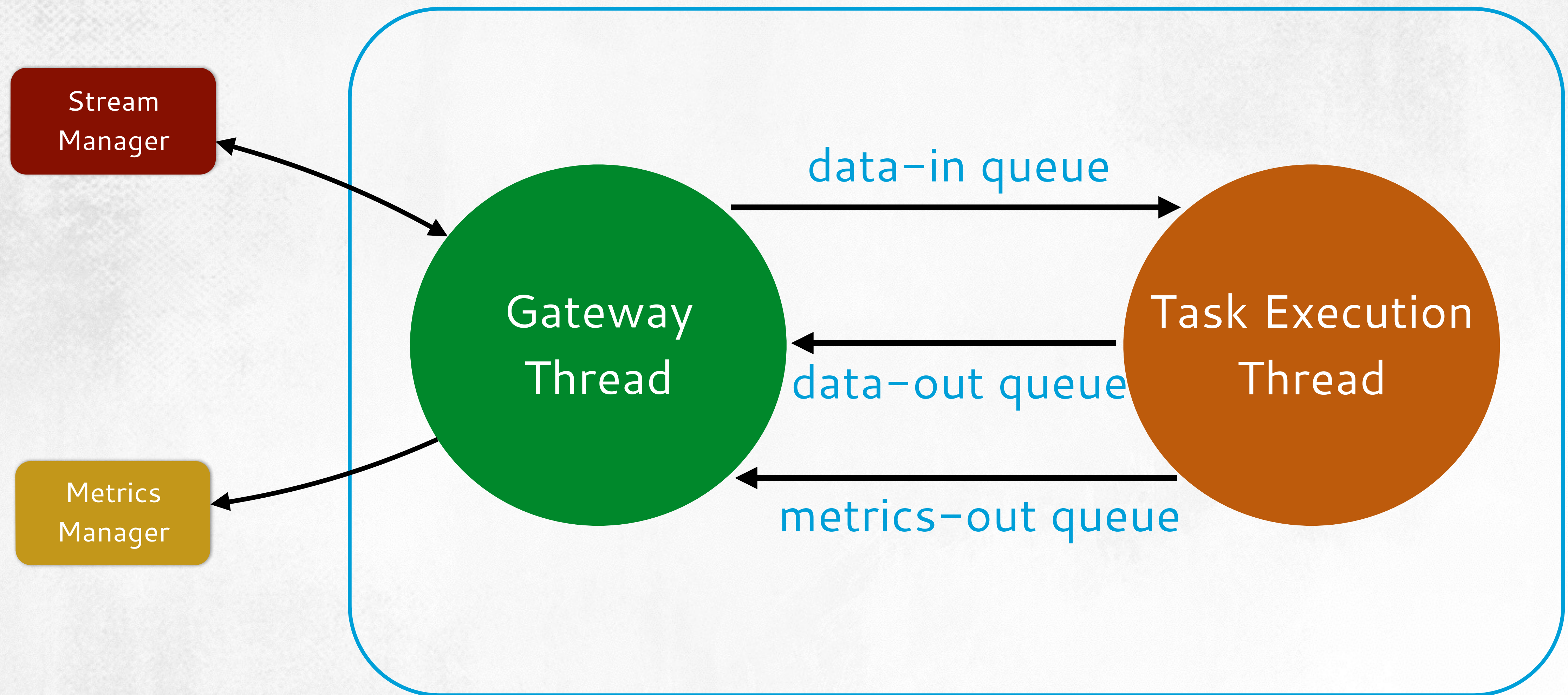
# HERON INSTANCE

## Does the real work!

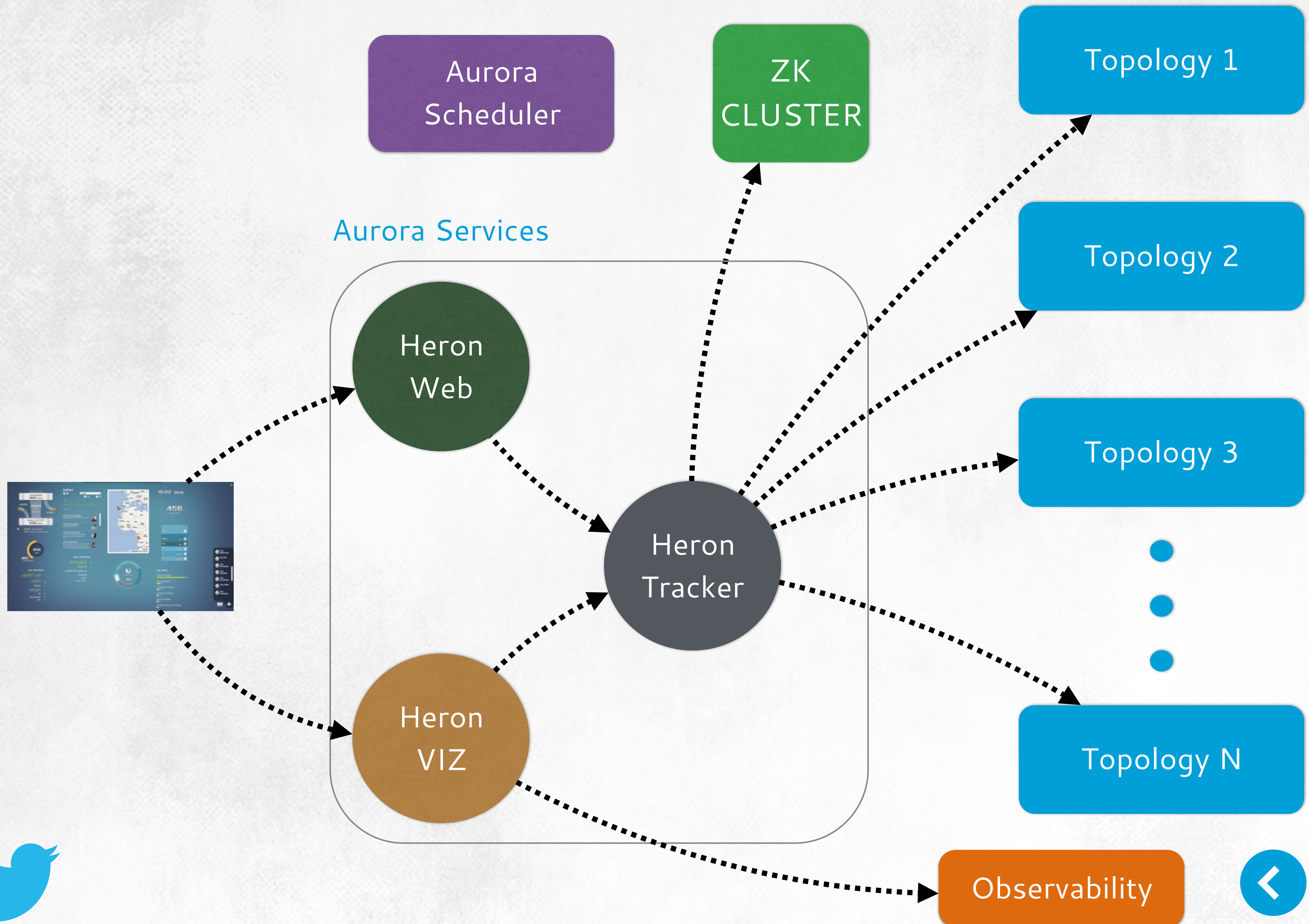**RUNS ONE TASK**

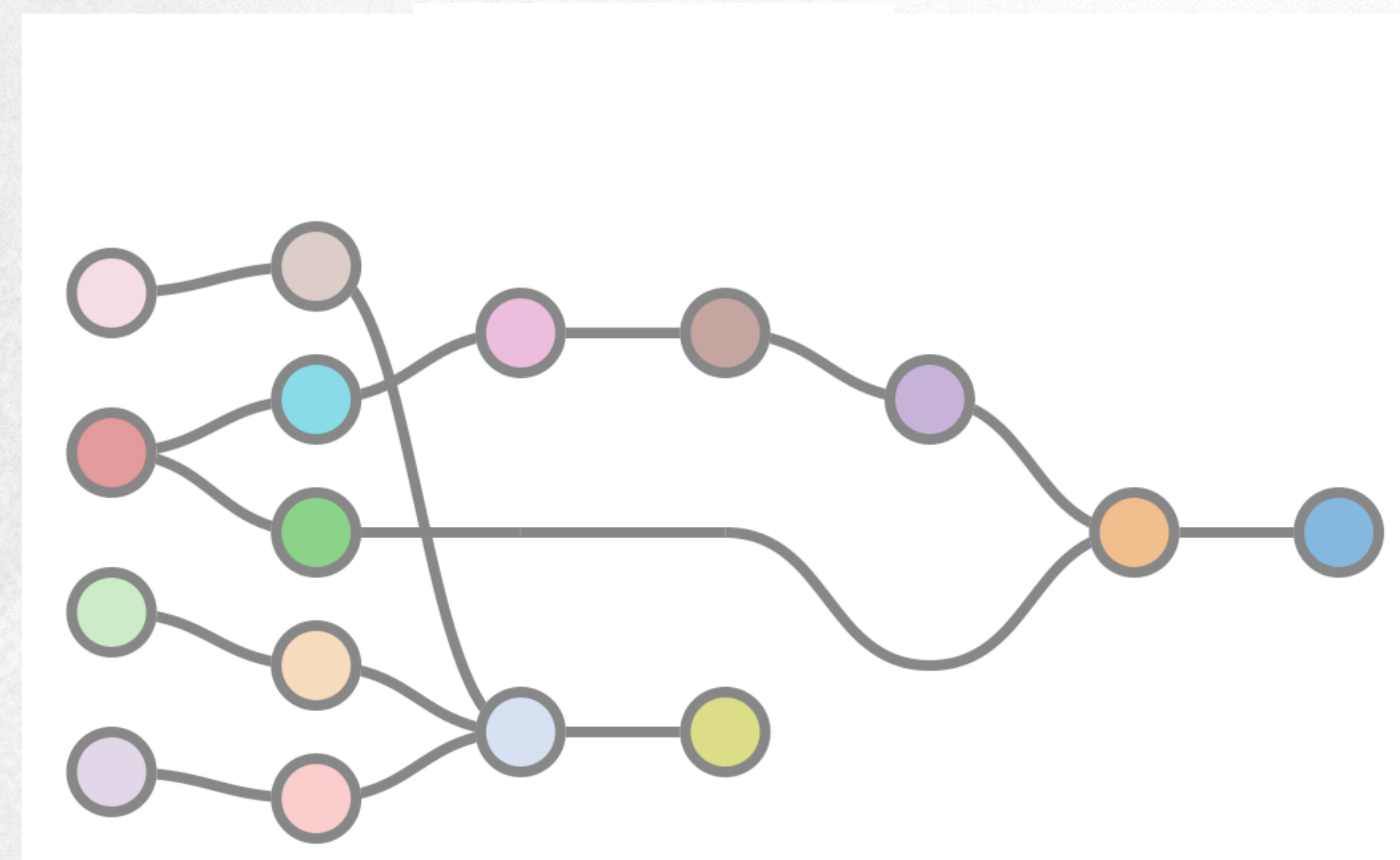**EXPOSES API**

**COLLECTS METRICS**
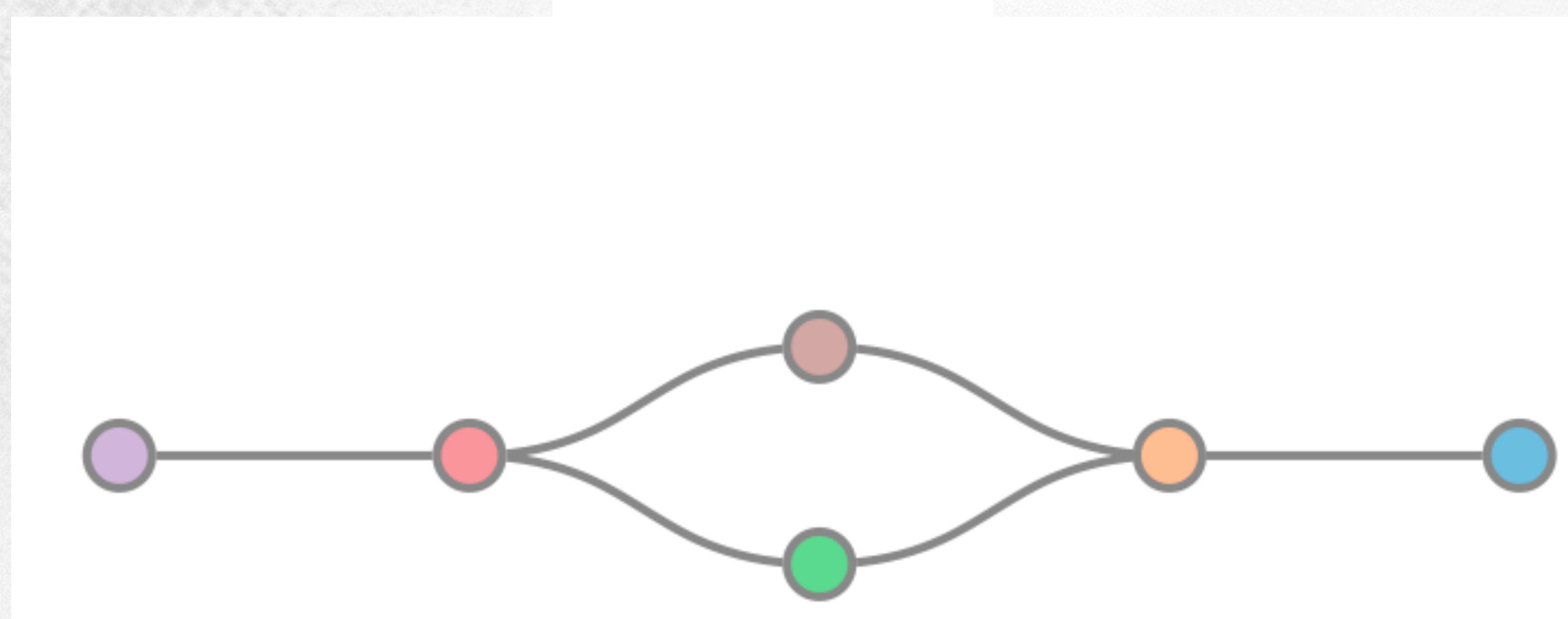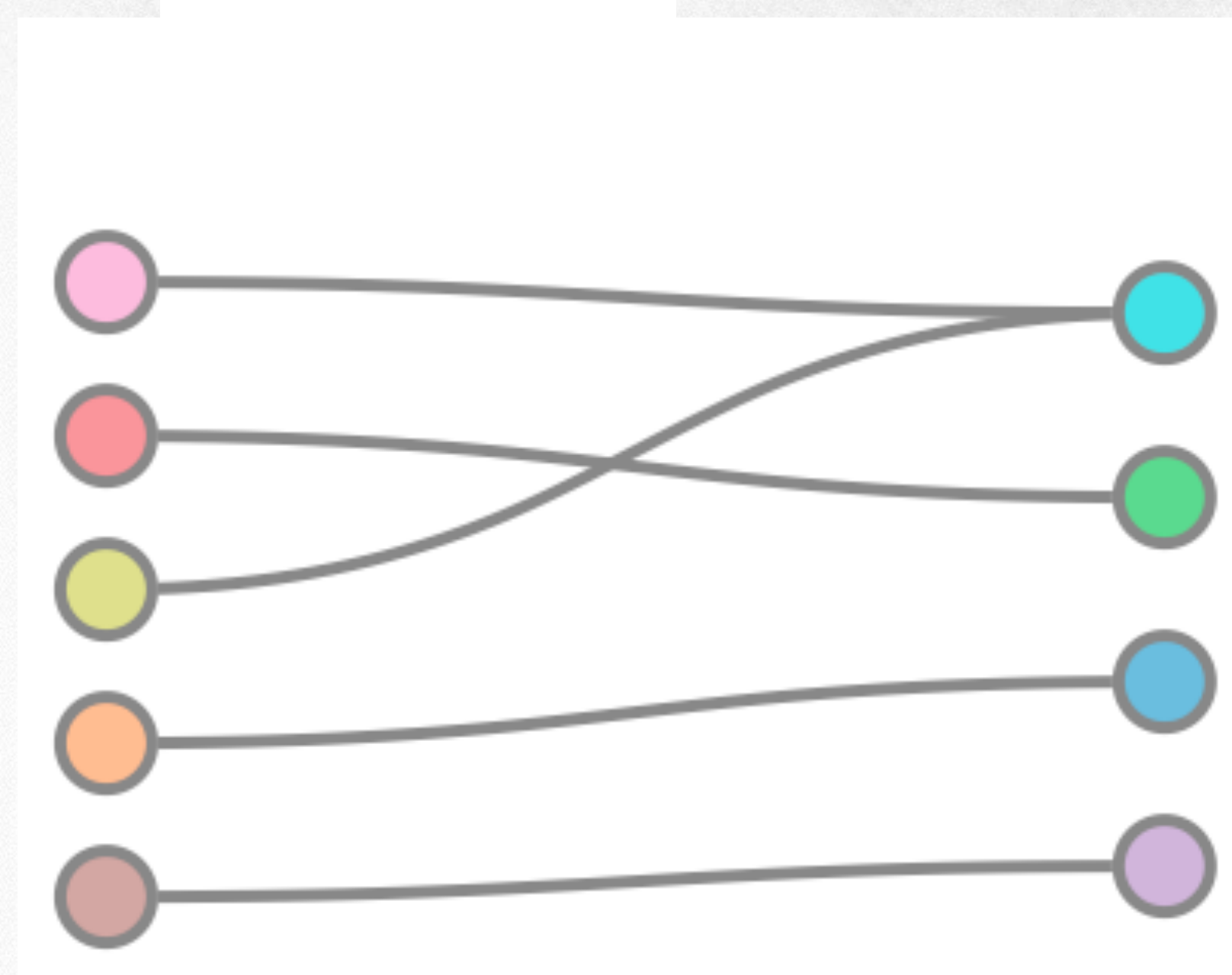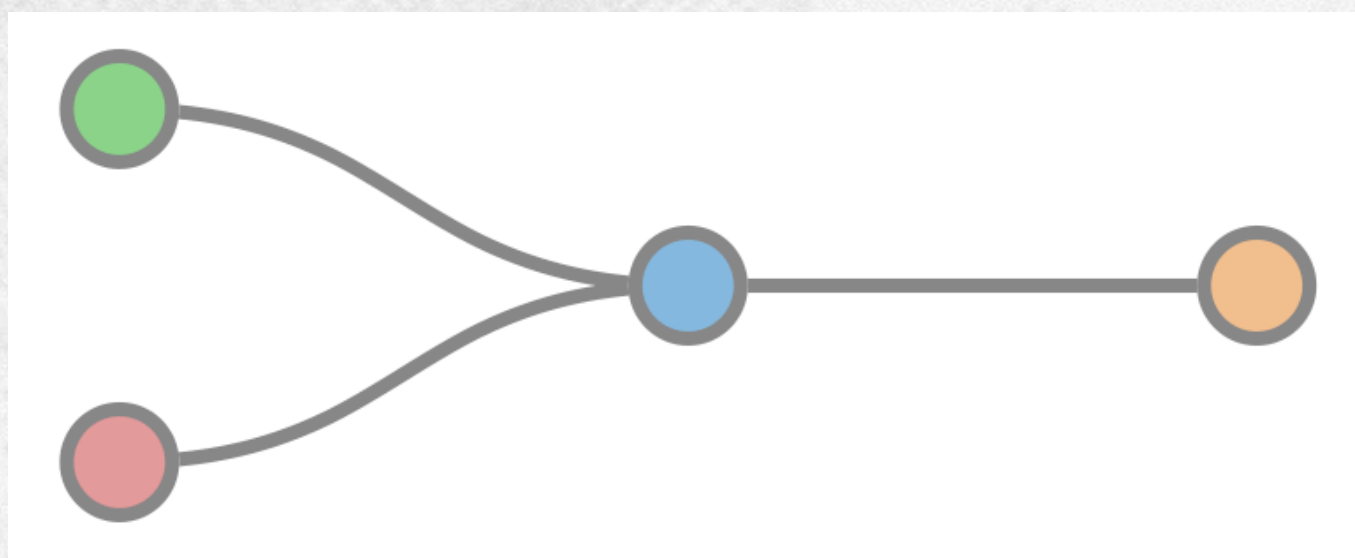
# HERON INSTANCE

# OPERATIONAL EXPERIENCES

# HERON DEPLOYMENT

# HERON SAMPLE TOPOLOGIES

# HERON @TWITTER

STORM is decommissioned

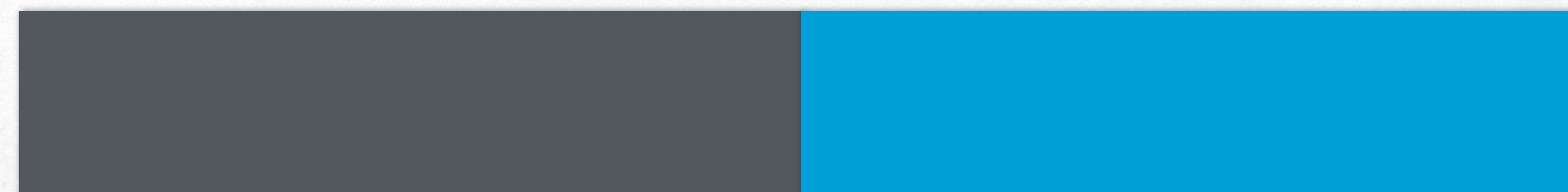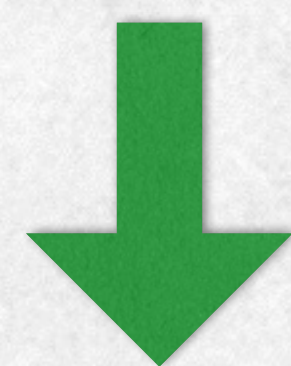| Large amount of data produced every day | Large cluster | Several topologies deployed | Several billion messages every day |

1 stage █████████ 10 stages

↓ 3x reduction in cores and memory
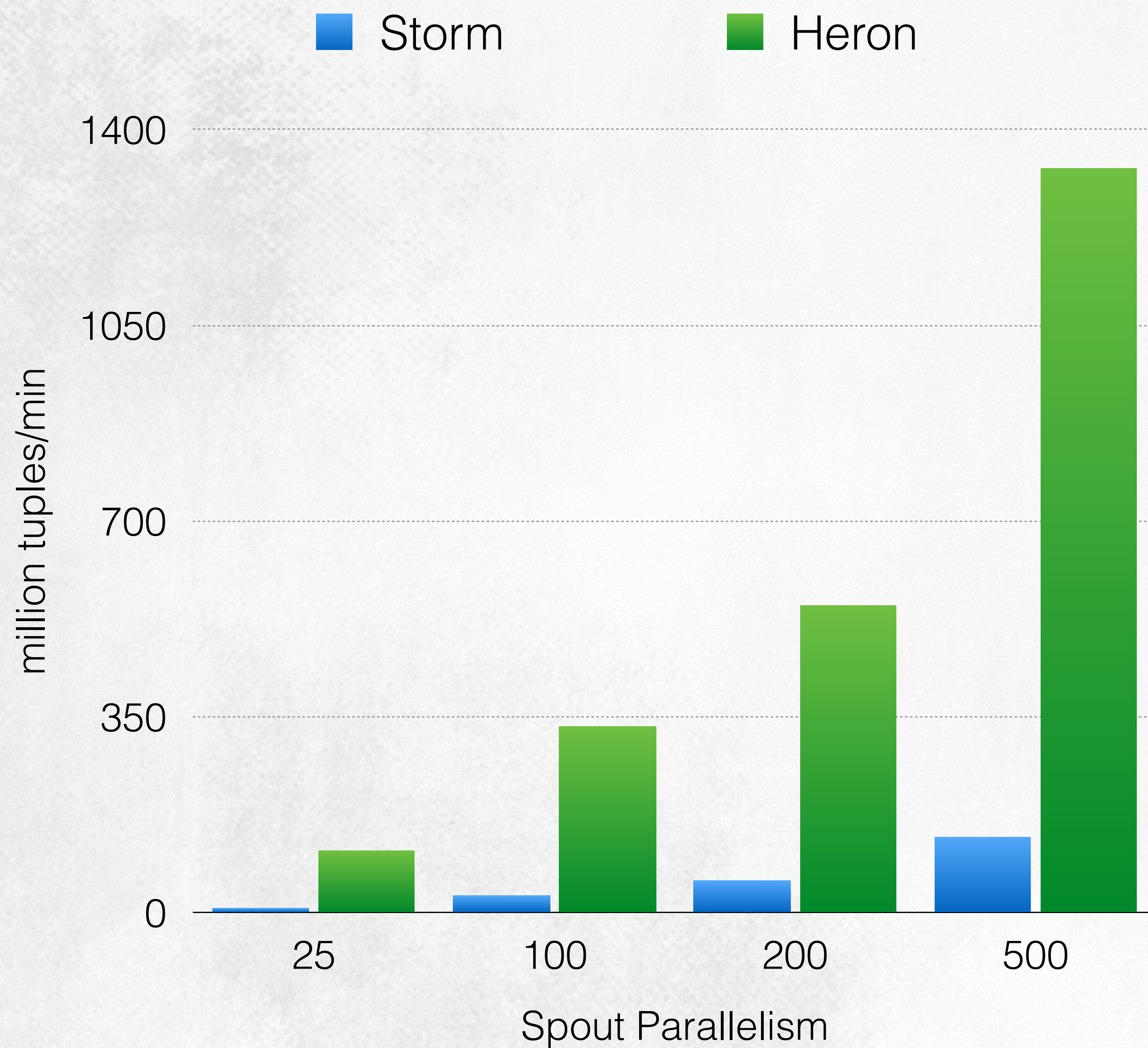
# HERON PERFORMANCE

Settings

| COMPONENTS | EXPT #1 | EXPT #2 | EXPT #3 | EXPT #4 |
|---|---|---|---|---|
| Spout | 25 | 100 | 200 | 300 |
| Bolt | 25 | 100 | 200 | 300 |
| # Heron containers | 25 | 100 | 200 | 300 |
| # Storm workers | 25 | 100 | 200 | 300 |

# HERON PERFORMANCE

## Word count topology – Acknowledgements enabled

### Throughput

Legend: Storm, Heron

million tuples/min vs Spout Parallelism (25, 100, 200, 500)

10-14x

### Latency

Legend: Storm, Heron

latency (ms) vs Spout Parallelism (25, 100, 200, 500)

5-15x

# HERON PERFORMANCE
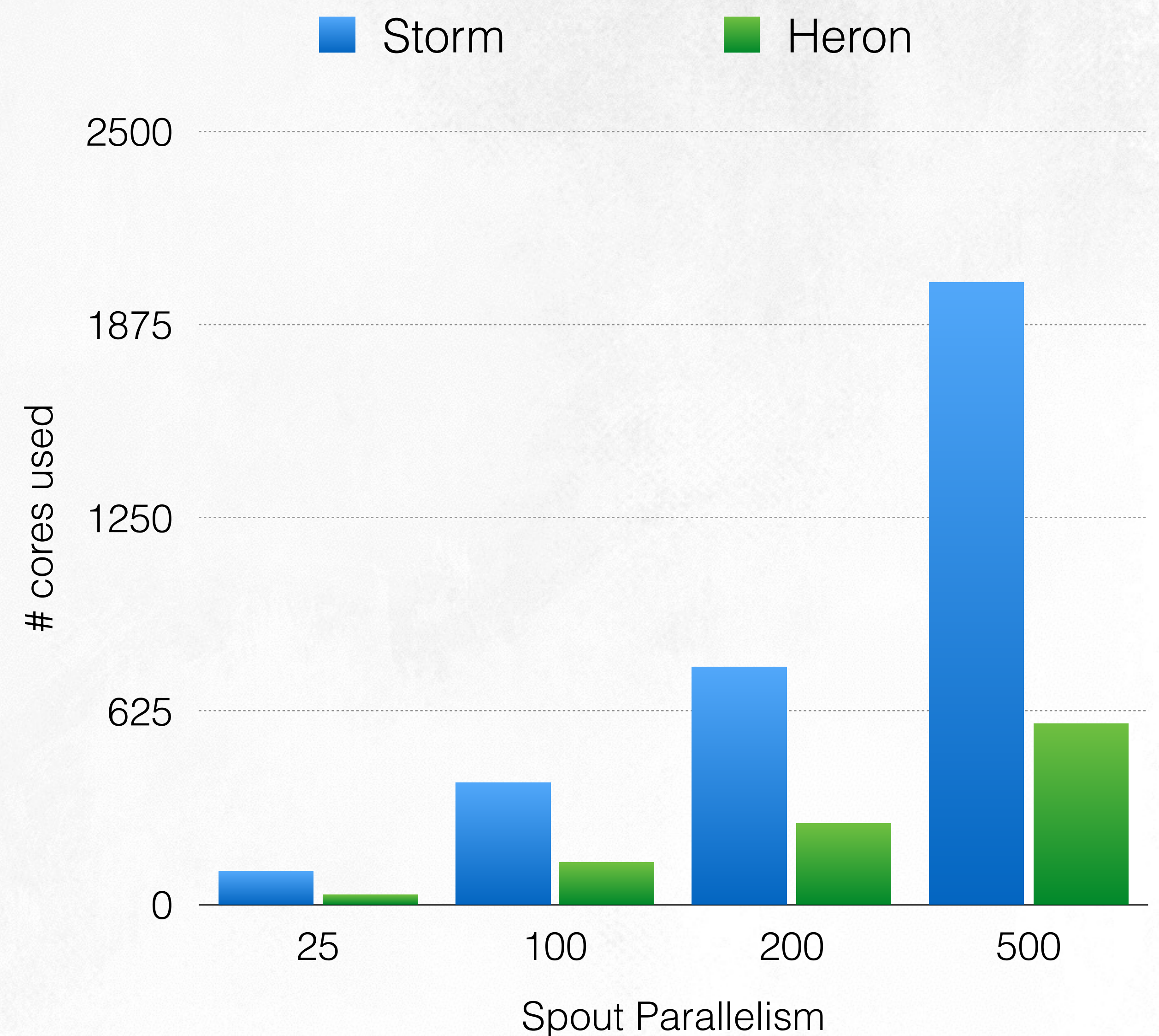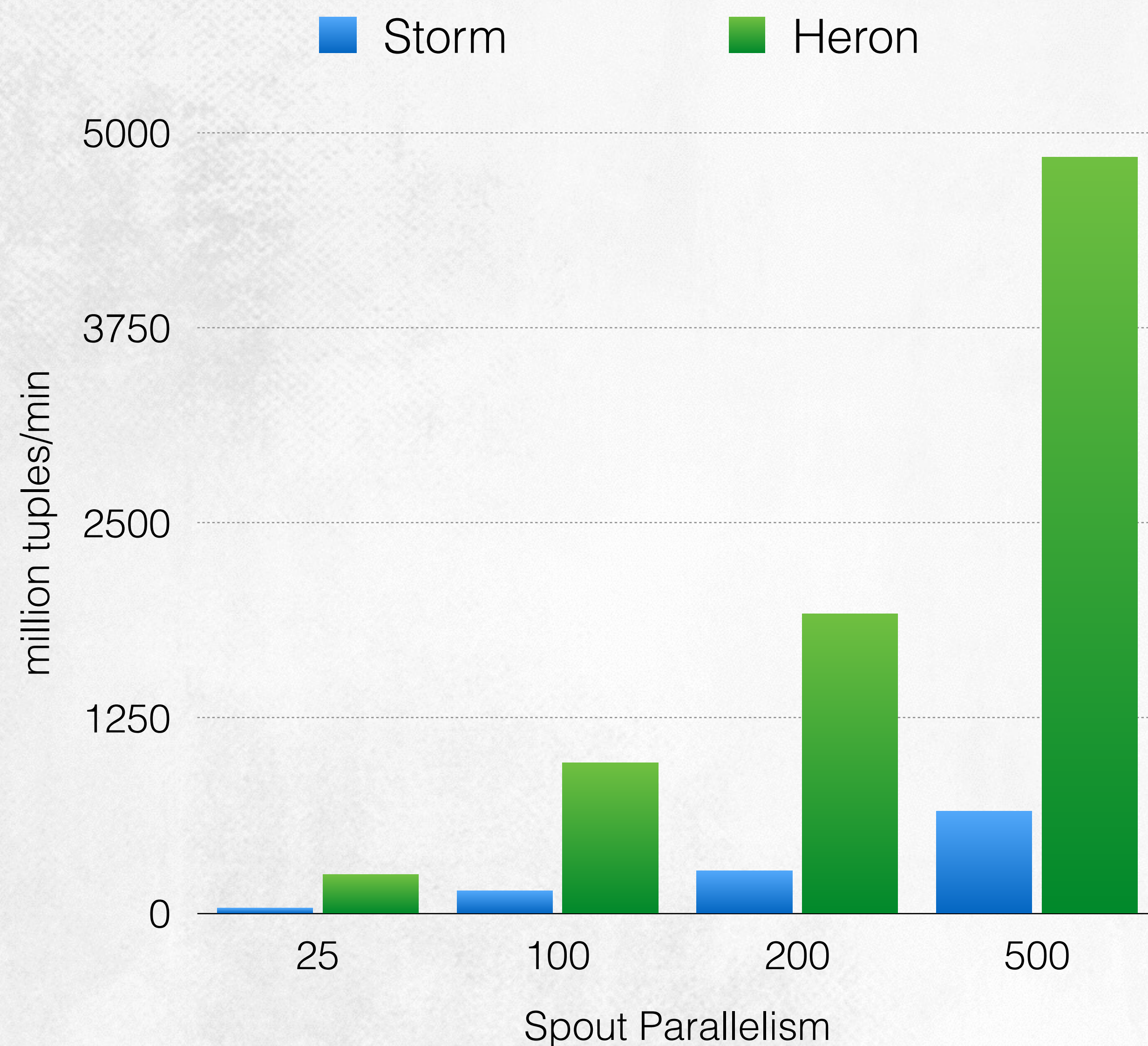## Word count topology – CPU usage

■ Storm   ■ Heron



# cores used

Spout Parallelism

2–3x

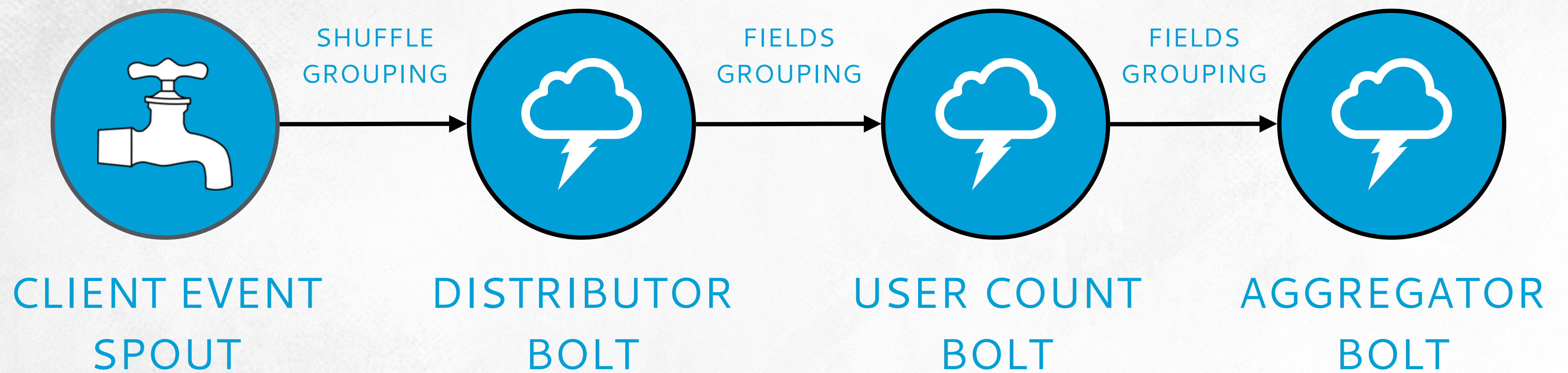# HERON PERFORMANCE

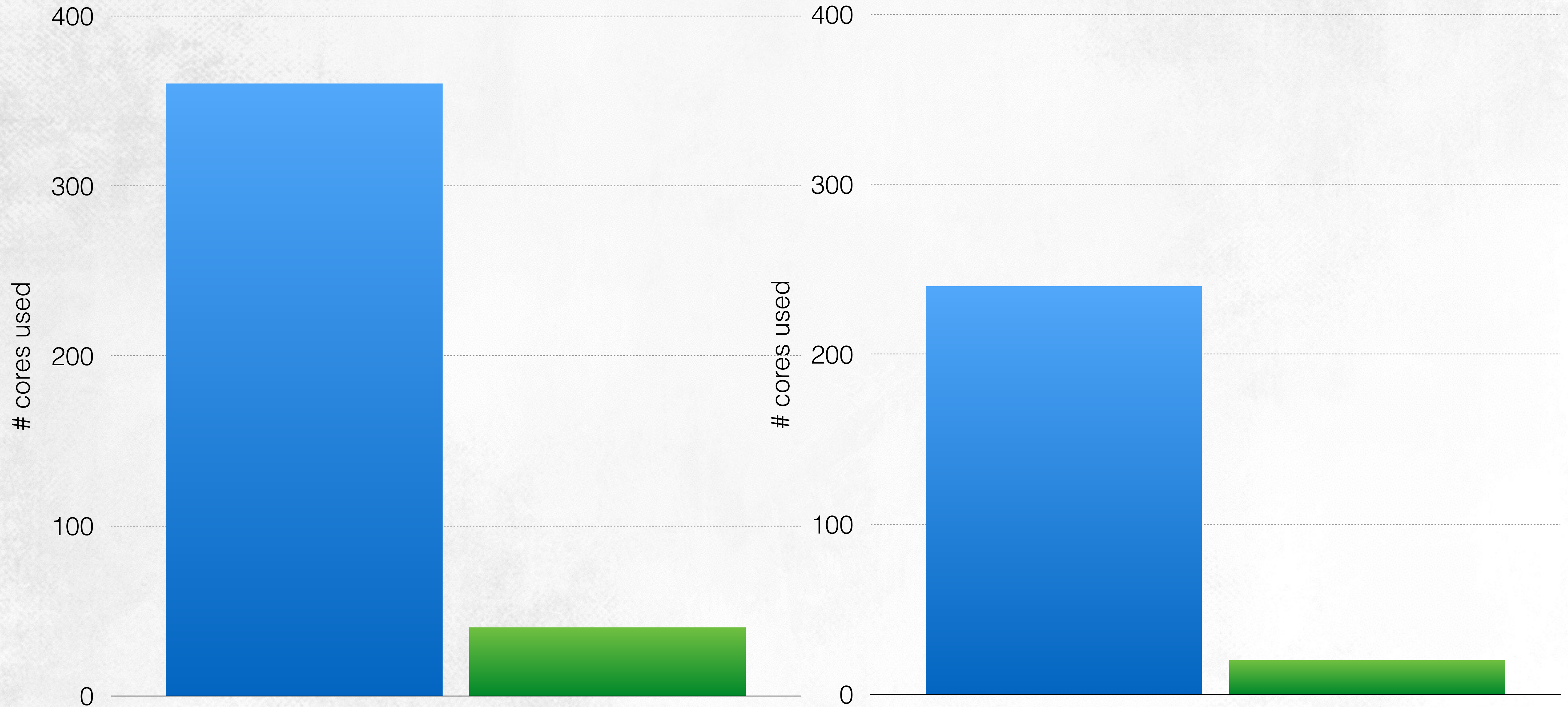Throughput and CPU usage with no acknowledgements – Word count topology

# HERON PERFORMANCE
## CPU usage – RTAC Topology

■ Storm    ■ Heron
Acknowledgements enabled

■ Storm    ■ Heron
No acknowledgements

# CONCLUSION

### SIMPLIFIED ARCHITECTURE

Easy to debug, profile and support

### HIGH PERFORMANCE

7–10x increase in throughput

5–10x improvement in latency

### EFFICIENCY

3–5x decrease in resource usage