# Cruise Control: Effortless Management of Kafka Clusters



**Adem Efe Gencer**

Senior Software Engineer
LinkedIn

# Kafka: A Distributed Stream Processing Platform
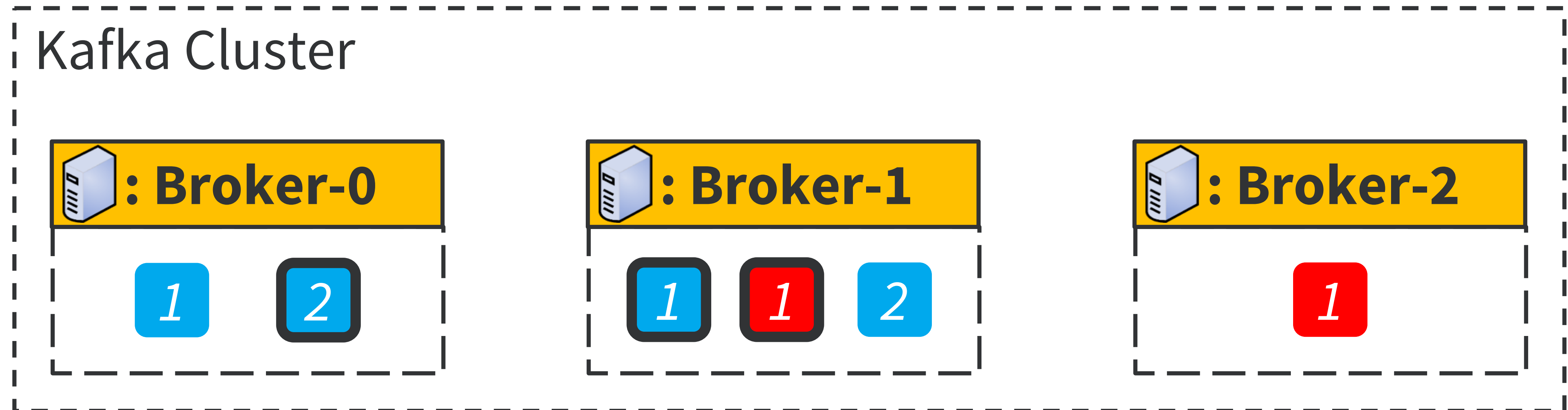
 : High throughput & low latency

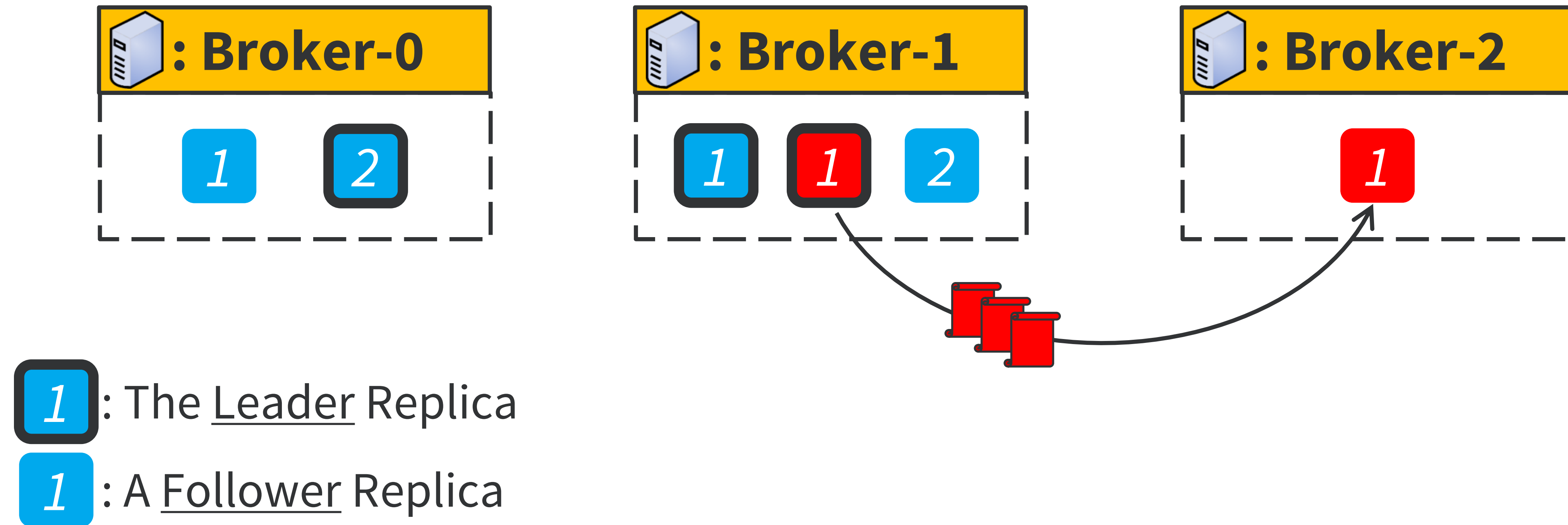 : Message persistence on partitioned data

 : Total ordering within each partition

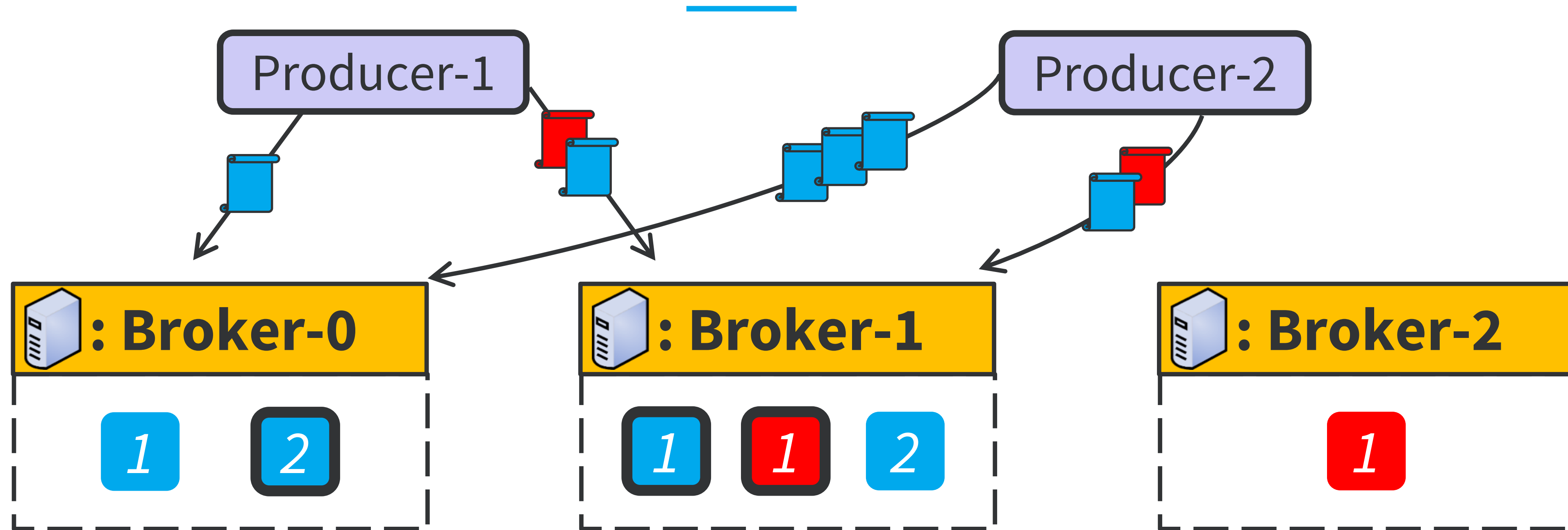# Key Concepts: Brokers, Topics, Partitions, and Replicas

Kafka Cluster

: Broker-0

| 1 | 2 |

: Broker-1

| 1 | 1 | 2 |

: Broker-2

| 1 |

1 : A Replica of Partition-1 of Blue Topic

# Key Concepts: Leaders and Followers



: Broker-0

: Broker-1
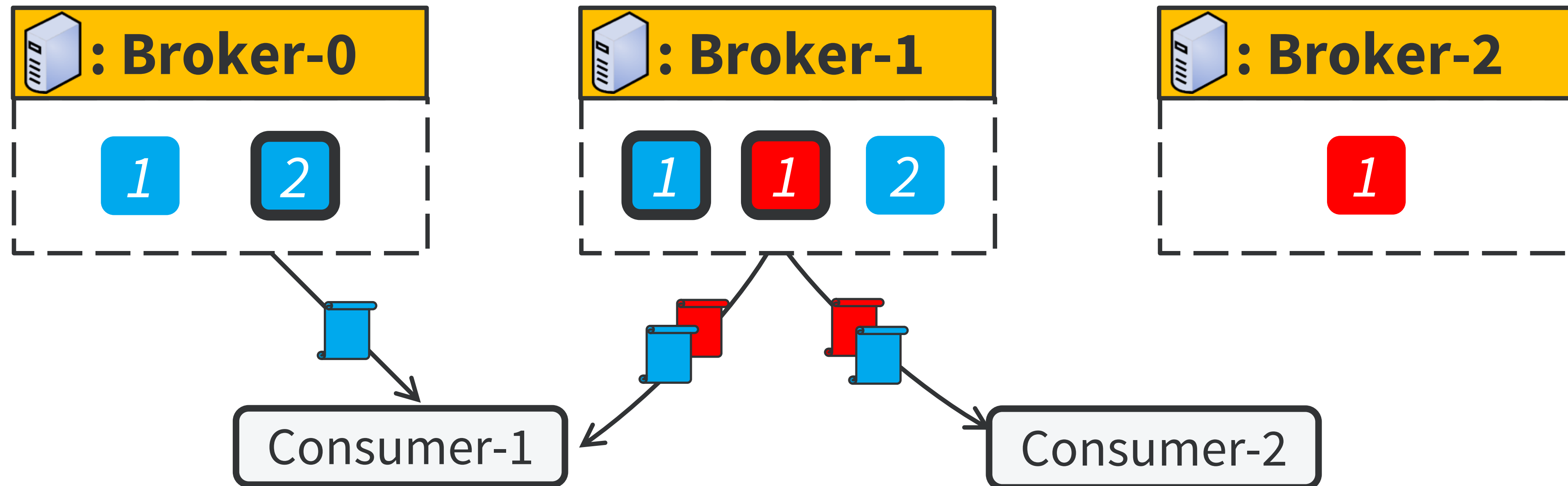
: Broker-2

**1** : The <u>Leader</u> Replica
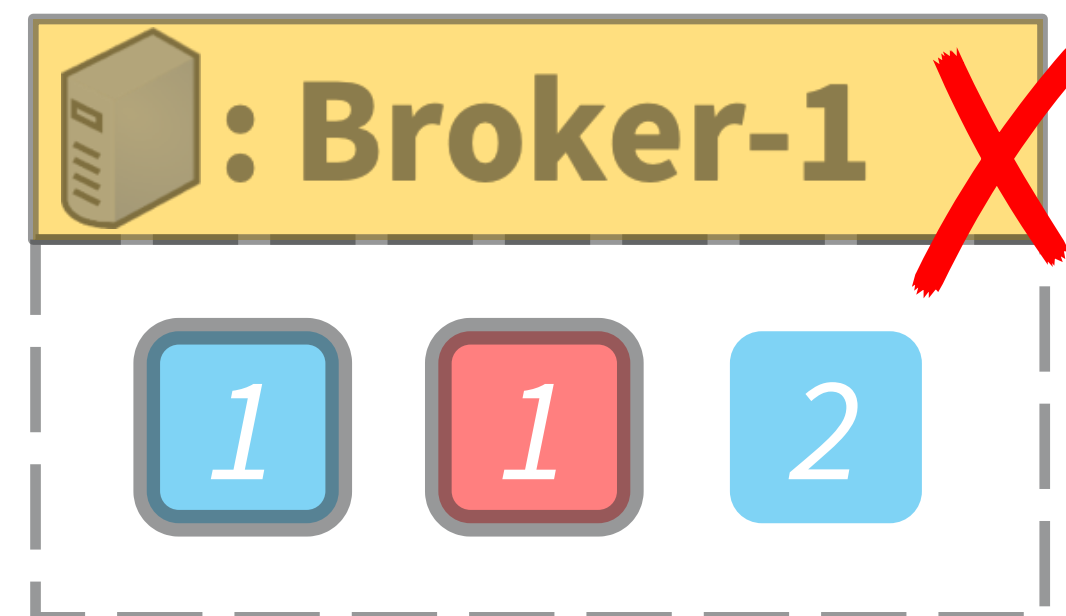
**1** : A <u>Follower</u> Replica
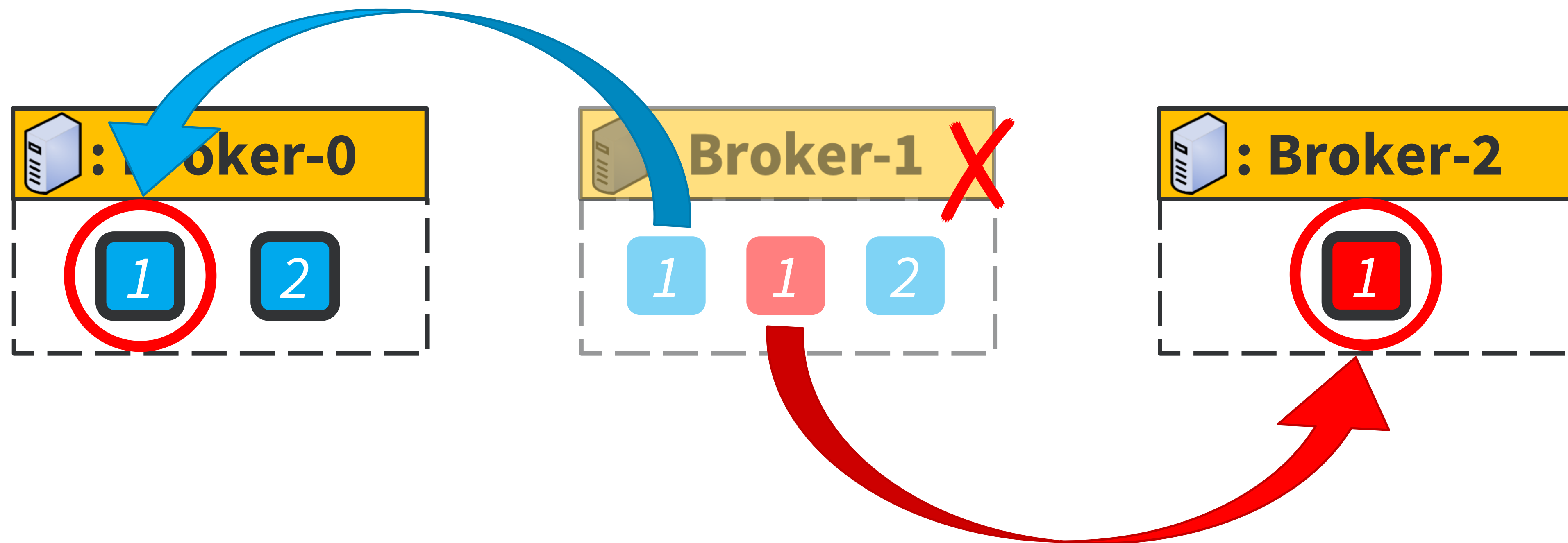
# Key Concepts: Producers

# Key Concepts: Consumers

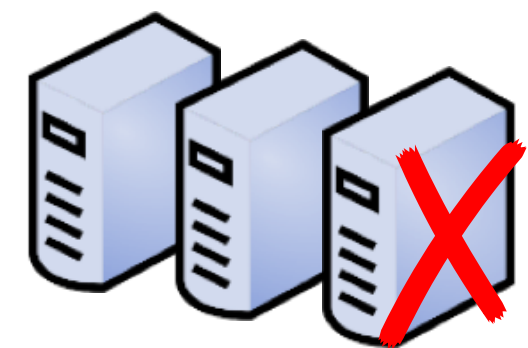# Key Concepts: Failover via Leadership Transfer

# Key Concepts: Failover via Leadership Transfer
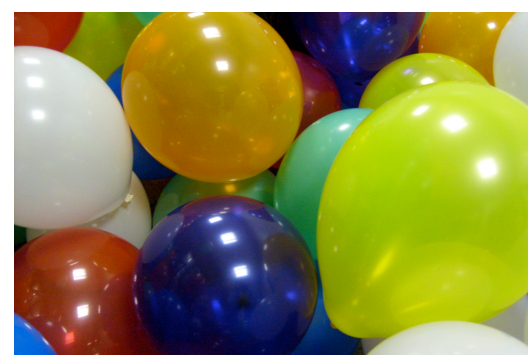
# Kafka Incurs Management Overhead

 : Large deployments – e.g. @**Linked**in: 2.6K+ Brokers, 44K+ Topics, 5M Partitions, 5T Messages / day

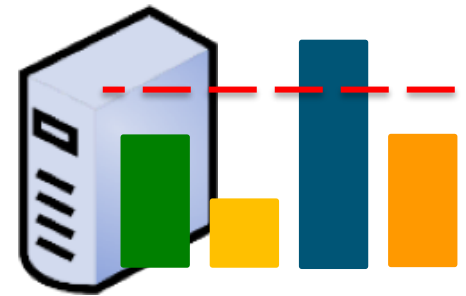 : Frequent hardware failures

 : Load skew among brokers

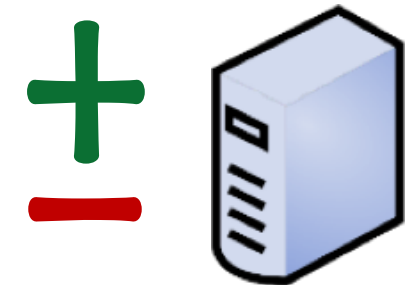 : Kafka cluster expansion and reduction

# Alleviating the Management Overhead

**1** Admin Operations for Cluster Maintenance

**2** Anomaly Detection with Self-Healing

**3** Real-Time Monitoring of Kafka Clusters

# **1** Admin Operations for Cluster Maintenance

 : Dynamically balance the cluster load

 : Add / remove brokers

 : Demote brokers – i.e. remove leadership of all replicas

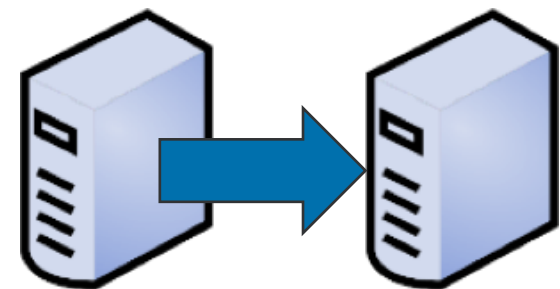 : Trigger preferred leader election

 : Fix offline replicas
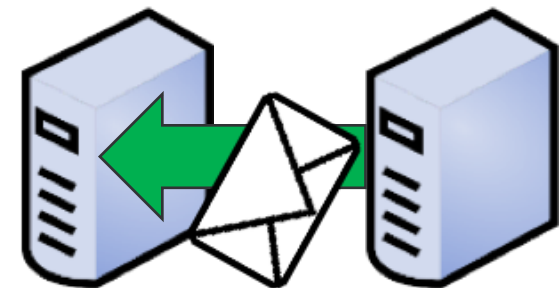
# ① Admin Operations for Cluster Maintenance

: Dynamically balance the cluster load

: Add / remove brokers

: Demote brokers – i.e. remove leadership of all replicas

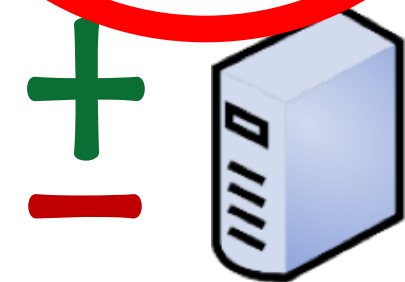: Trigger preferred leader election
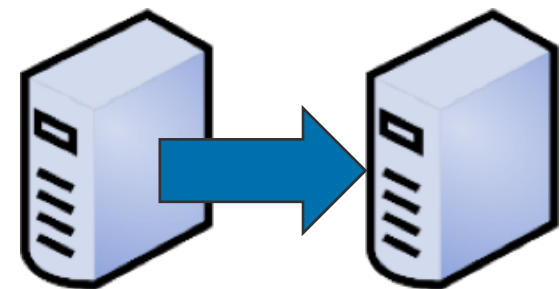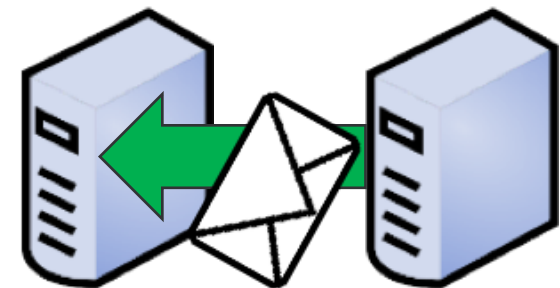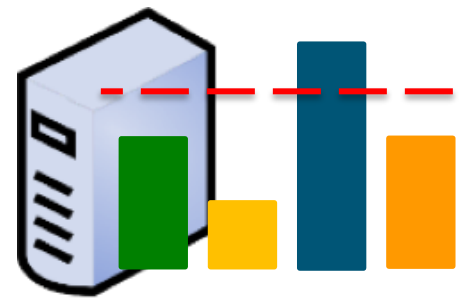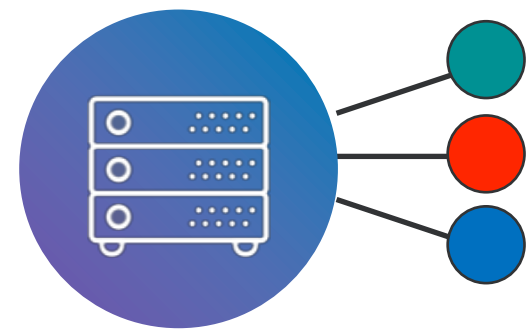
: Fix offline replicas

# Dynamically Balance the Cluster Load

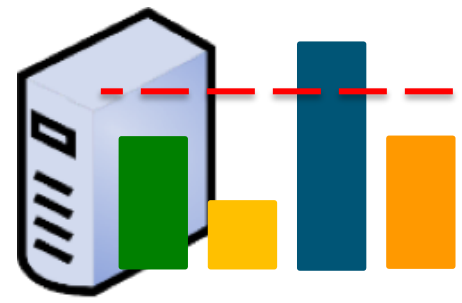Must satisfy *hard goals*, including:

: Guarantee rack-aware distribution of replicas

: Never exceed the capacity of broker resources – e.g. disk, CPU, network bandwidth
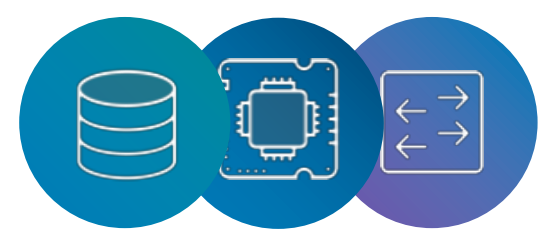
: Enforce operational requirements – e.g. maximum replica count per broker

# Dynamically Balance the Cluster Load

Satisfy *soft goals* as much as possible – i.e. best effort

: Balance disk, CPU, inbound/outbound network traffic utilization of brokers

: Balance replica distribution

: Balance potential outbound network load

: Balance distribution of partitions from the same topic

# **2** **Anomaly Detection with Self-Healing**

: Goal violation – rebalance cluster

: Broker failure – decommission broker(s)

: Metric anomaly – demote broker(s)

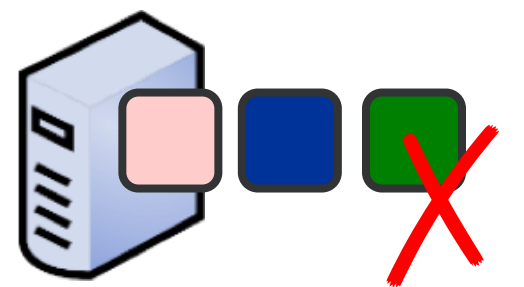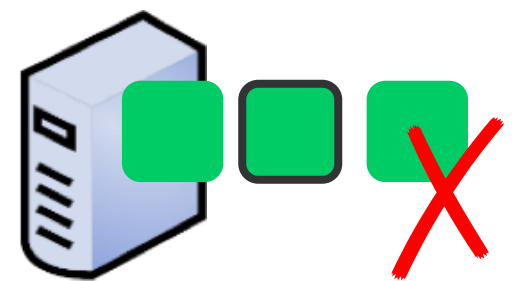## 3 Real-Time Monitoring of Kafka Clusters

: Examine the replica, leader, and load distribution

: Identify *under-replicated*, *under-min-ISR*, and *offline* partitions

: Check the health of brokers, disks, and user tasks

# Building Blocks of Management: Moving Replicas

# Building Blocks of Management: Moving Replicas

: Broker-0

1 2 1

: Broker-1

1 2

Replica Move

Broader impact, but expensive
• Requires data transfer*

* Replica swap: Bidirectional reassignments of distinct partition replicas among brokers

# Building Blocks of Management: Moving Leadership

**: Broker-0**

| 2 | 1 |

Leadership Move

**: Broker-1**

| 1 | 2 | 1 |

# Building Blocks of Management: Moving Leadership

: Broker-0

2  1

Leadership Move

└→ Cheap, but has limited impact
  • Affects network bytes out and CPU

: Broker-1

1  2  1

# A Multi-Objective Optimization Problem

Achieve conflicting cluster management goals while minimizing the impact of required operations on user traffic

# ARCHITECTURE

# Cruise Control Architecture



**REST API**

**Monitor**
- Sample Store
- Metric Sampler
- Capacity Resolver

**Analyzer**
- Goal(s)

**Anomaly Detector**
- Goal Violation
- Metric Anomaly
- Broker Failure
- Finder(s)
- Anomaly Notifier

**Executor**

Backup and Recovery

Load History T.

Metrics Reporter T.

Reported Metrics

Metrics Reporter

Broker Failures

**Kafka Cluster**

Throttled Proposal Execution

**Pluggable Component**
- Implements a public interface
- Accepts custom user code

**Internal Topic**
- Created and used by Cruise Control and its metrics reporter

# Metrics Reporter

Produces selected Kafka cluster metrics to the configured metrics reporter topic with the configured frequency

# Monitor
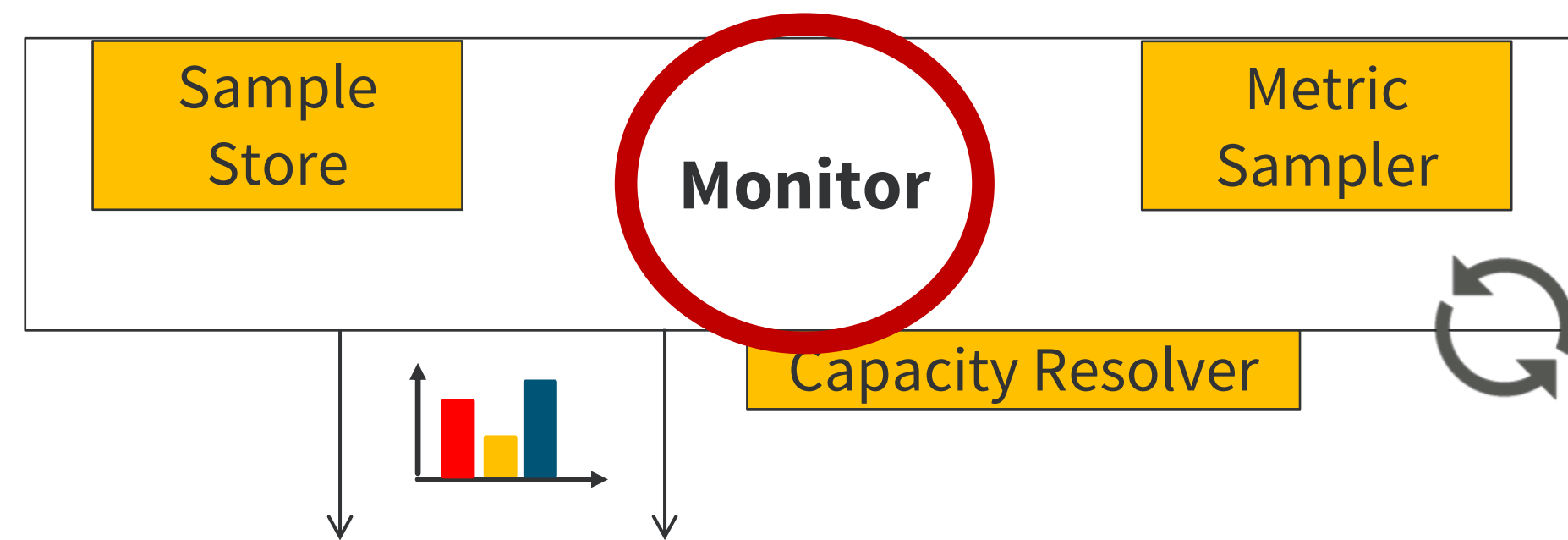


Generates a model ( ) to describe the cluster

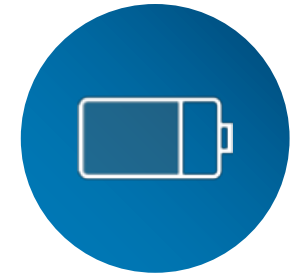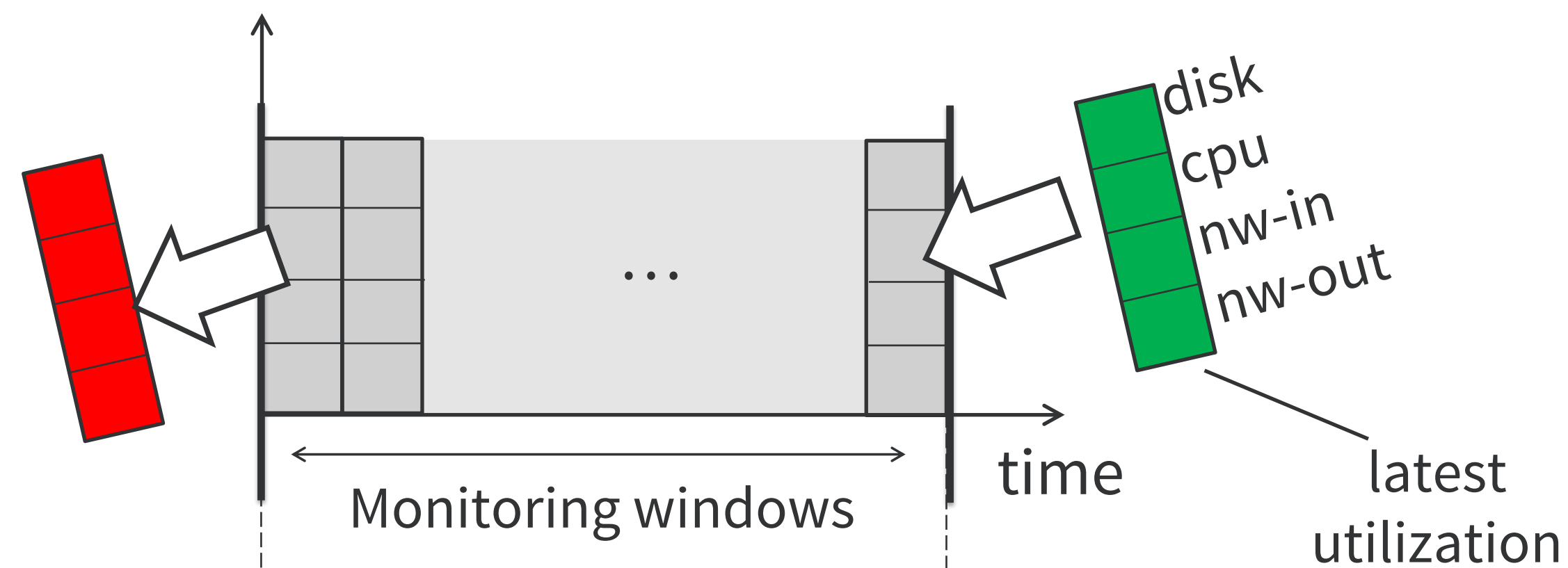# Monitor: Cluster Model (📊)

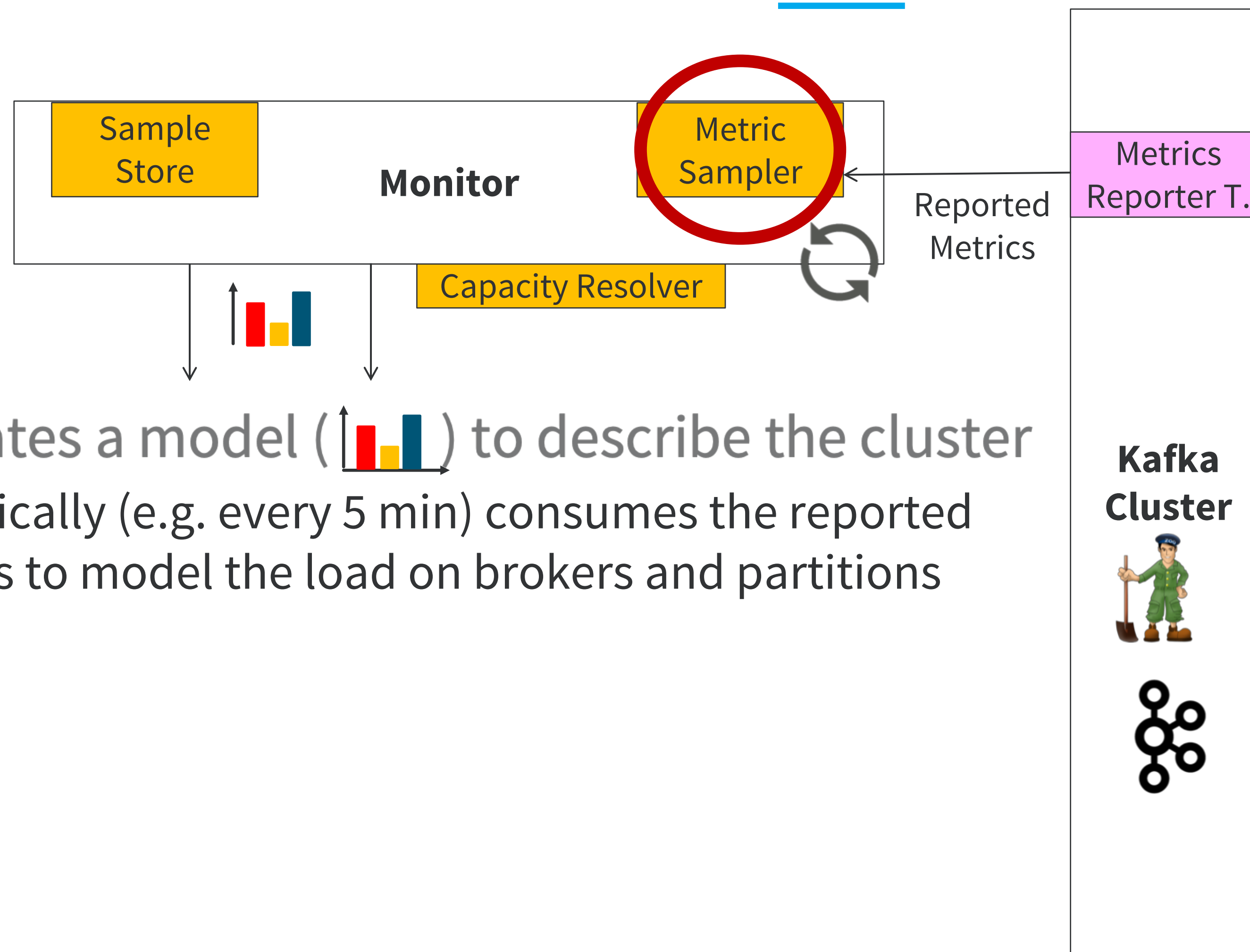🔵 : *Topology* – rack, host, and broker distribution

🔵 : *Placement* – replica, leadership, and partition distribution

🔵 : *Load* – current and historical utilization of brokers and replicas

# Monitor: Metric Sampler



Generates a model ( ) to describe the cluster

- Periodically (e.g. every 5 min) consumes the reported metrics to model the load on brokers and partitions

# Monitor: Sample Store

Backup and Recovery

Load History T.

Sample Store

Monitor

Metric Sampler

Capacity Resolver

Generates a model (📊) to describe the cluster

- Periodically (e.g. every 5 min) consumes the reported metrics to model the load on brokers and partitions
- Produces broker and partition models to load history topic, and uses the stored data to recover upon failure

**Kafka Cluster**

# Monitor: Capacity Resolver
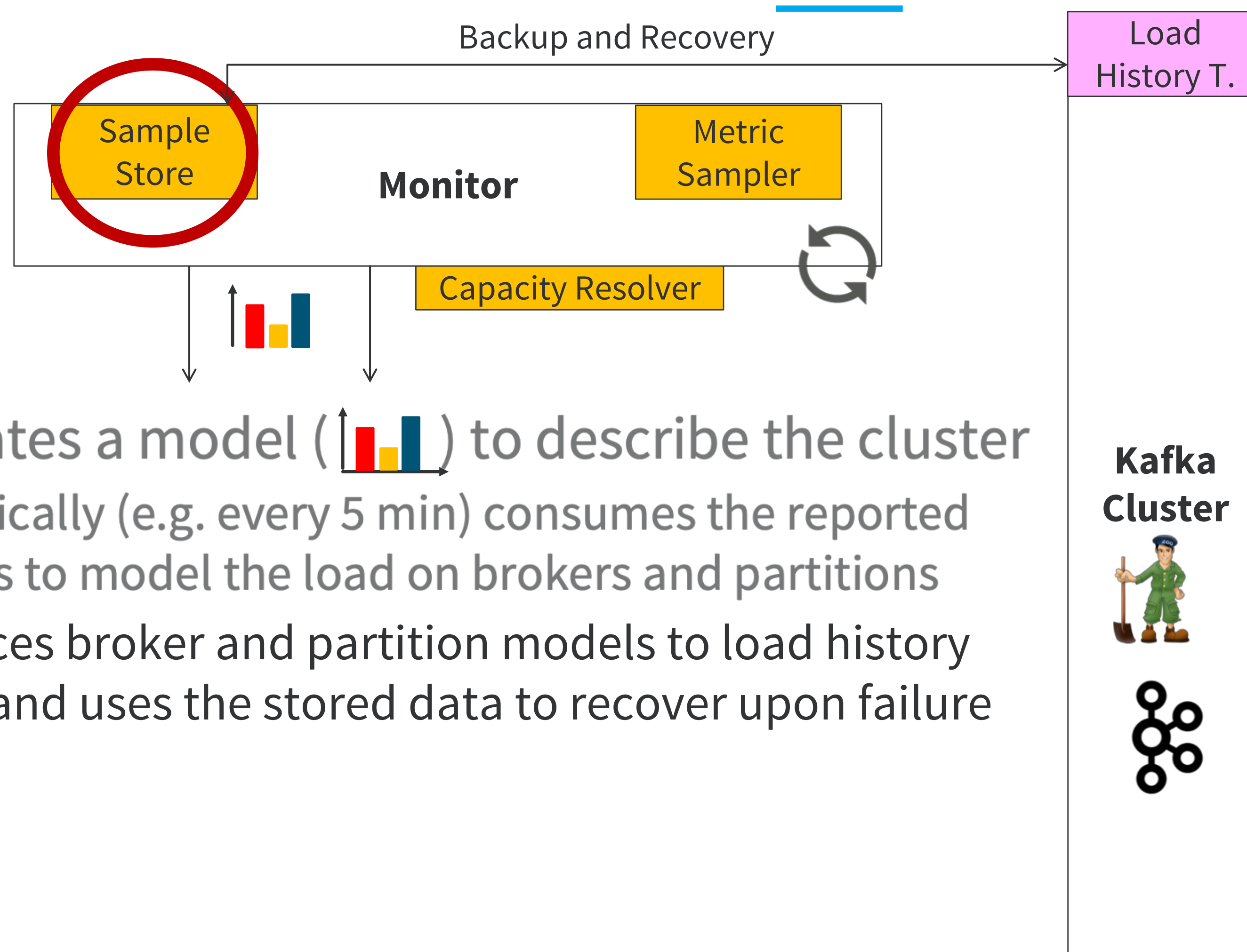


Generates a model ( ) to describe the cluster

- Periodically (e.g. every 5 min) consumes the reported metrics to model the load on brokers and partitions
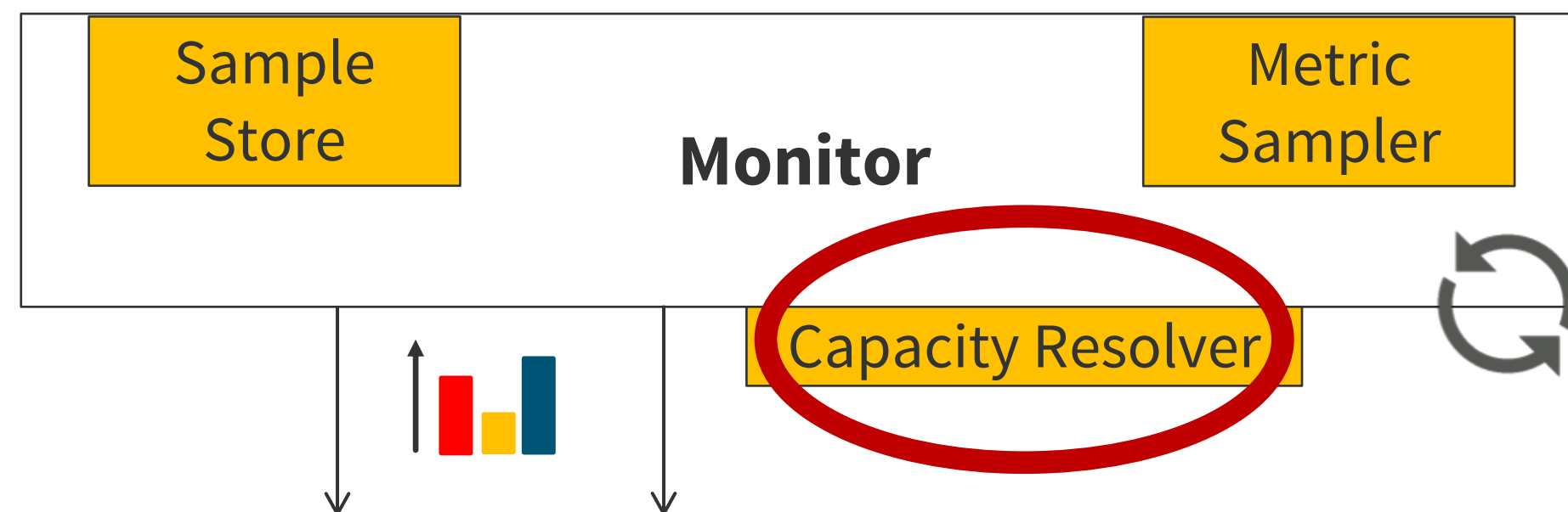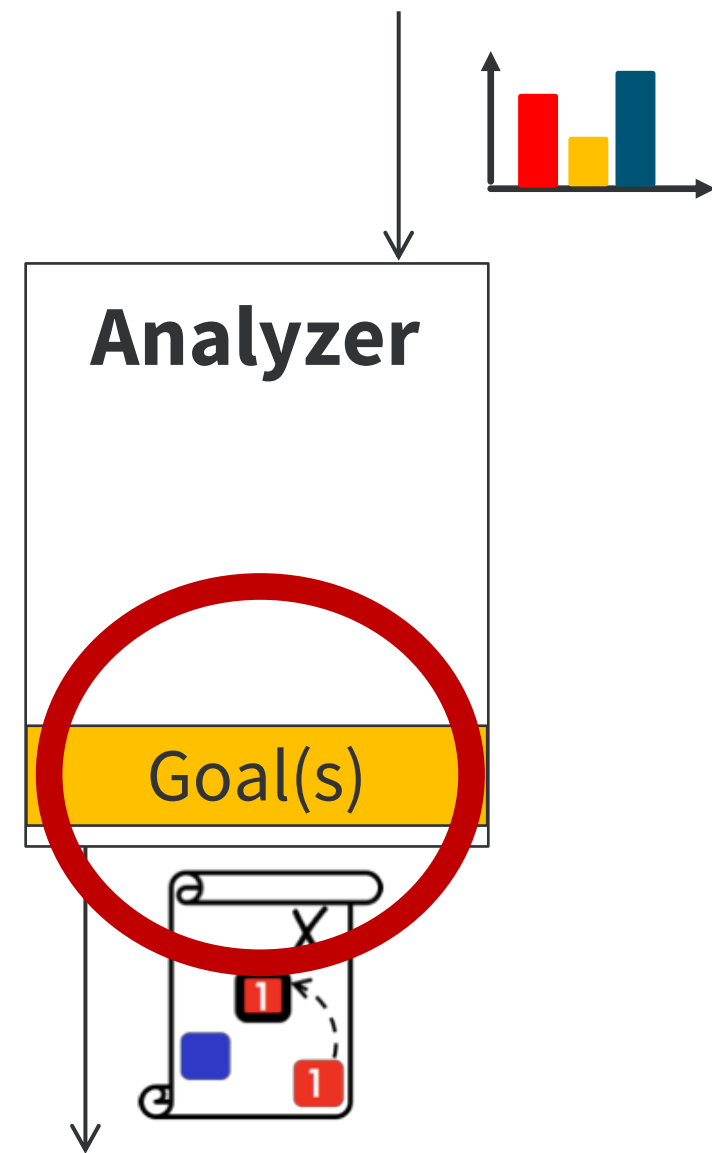
- Produces broker and partition models to load history topic, and uses the stored data to recover upon failure

- Gathers the broker capacities from a pluggable resolver

# Analyzer



Generates proposals to achieve goals via a fast and near-optimal heuristic solution
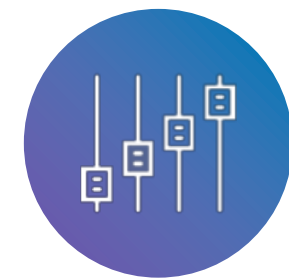
# Analyzer: Goals



Generates proposals to achieve goals via a fast and near-optimal heuristic solution

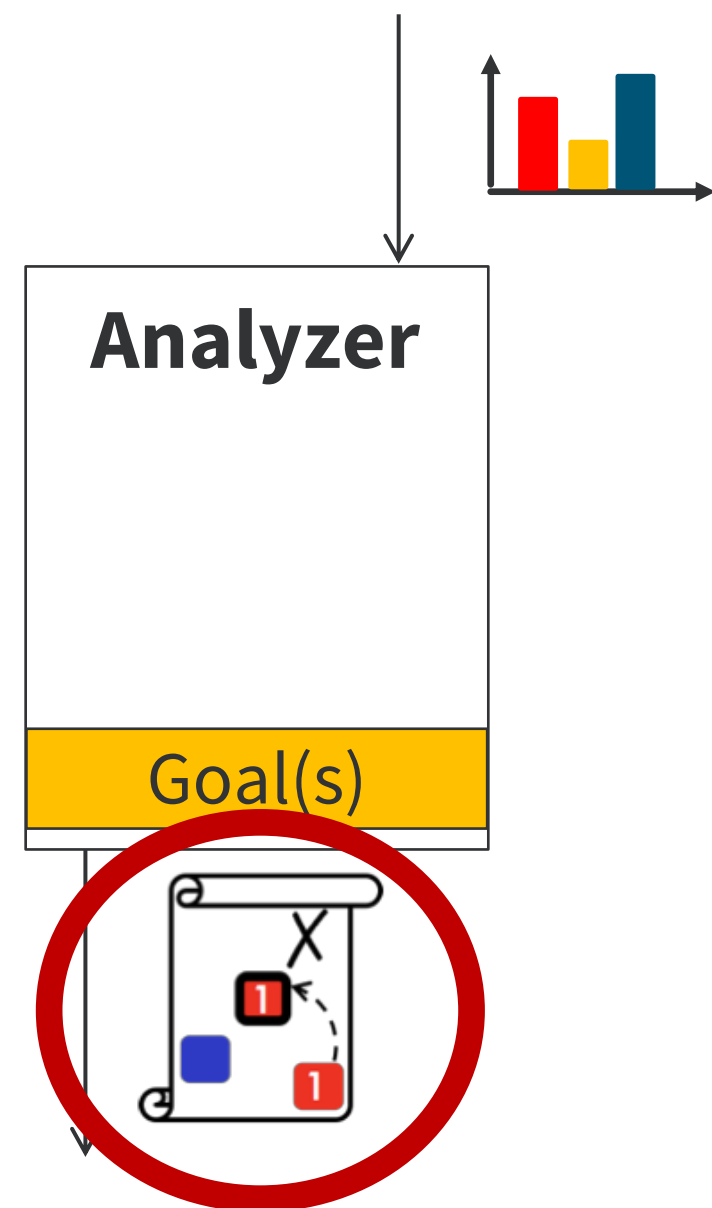: *Priorities* – custom order of optimization

: *Strictness* – hard (e.g. rack awareness) or soft (e.g. resource utilization balance) optimization demands
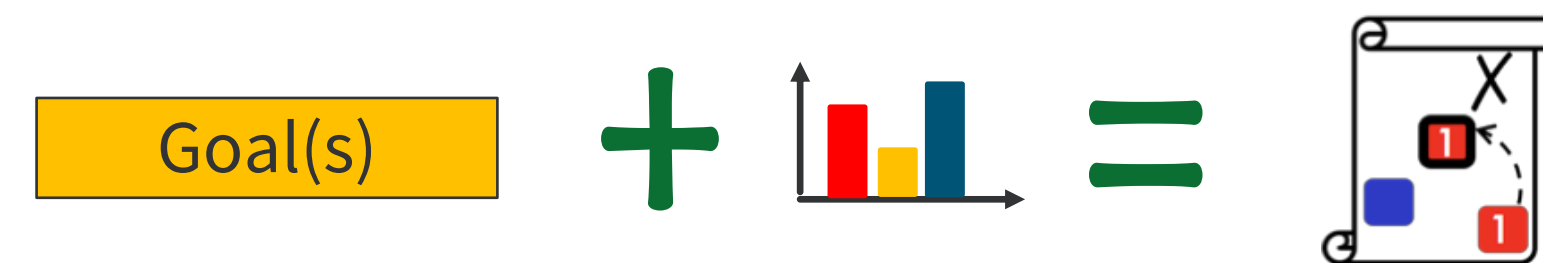
: *Modes* – e.g. kafka-assigner (https://github.com/linkedin/kafka-tools)

# Analyzer: Proposals

Generates proposals to achieve goals via a fast and near-optimal heuristic solution

Goal(s) **+** = 

**Proposals** – in order of priority:

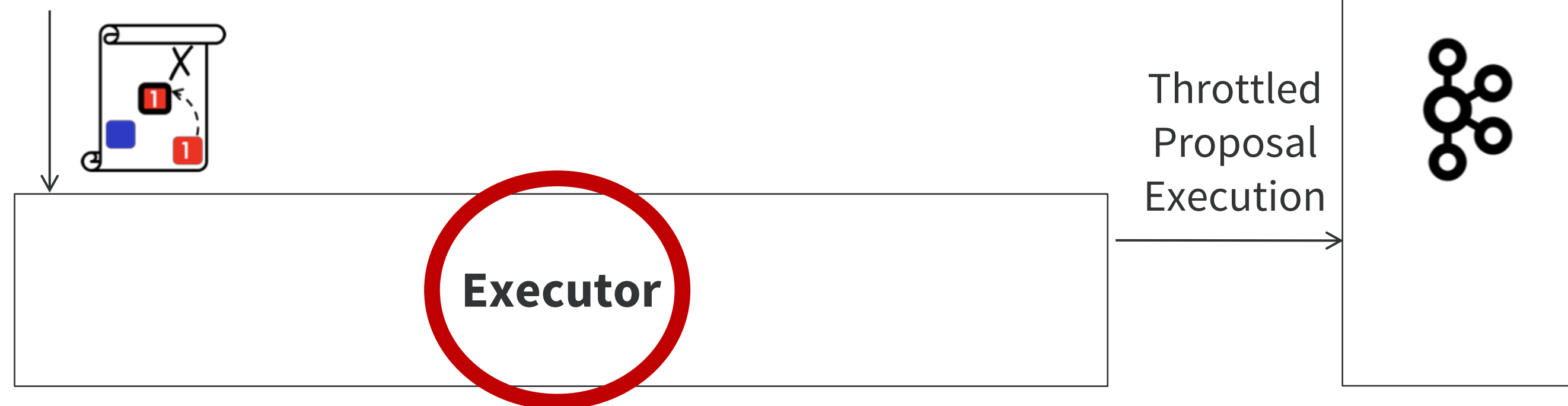• Leadership move > Replica move > Replica swap

# Executor

Proposal execution:

- Dynamically controls the maximum number of concurrent leadership / replica reassignments
- Ensures only one execution at a time
- Enables graceful cancellation of ongoing executions
- Integration with replication quotas (KIP-73)

**Kafka Cluster**

Executor

Throttled Proposal Execution

# Anomaly Detector

Identifies, notifies, and fixes (self-healing):

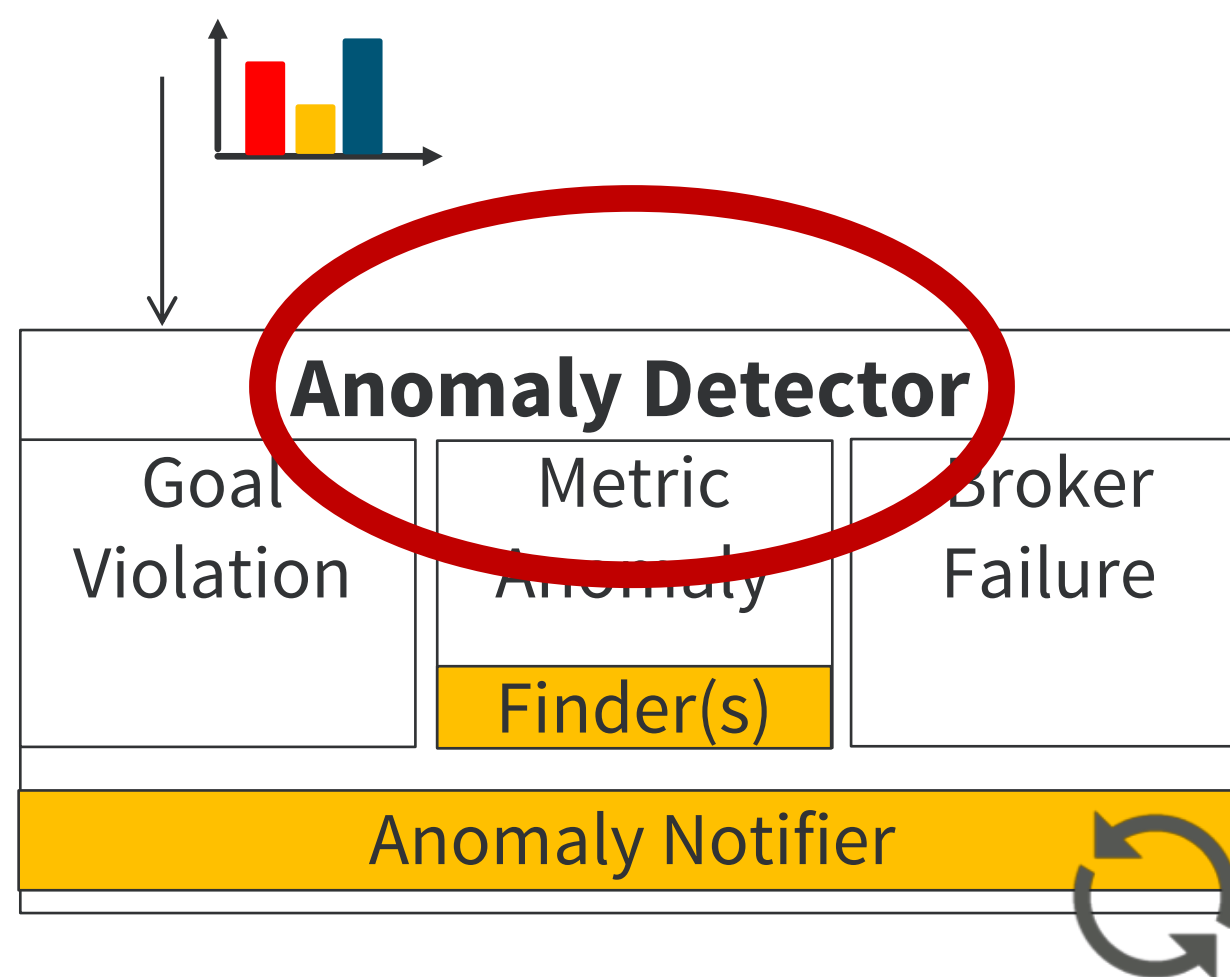- Violation of anomaly detection goals
- Broker failures
- Metric anomalies

⧗ Disk failures (JBOD)

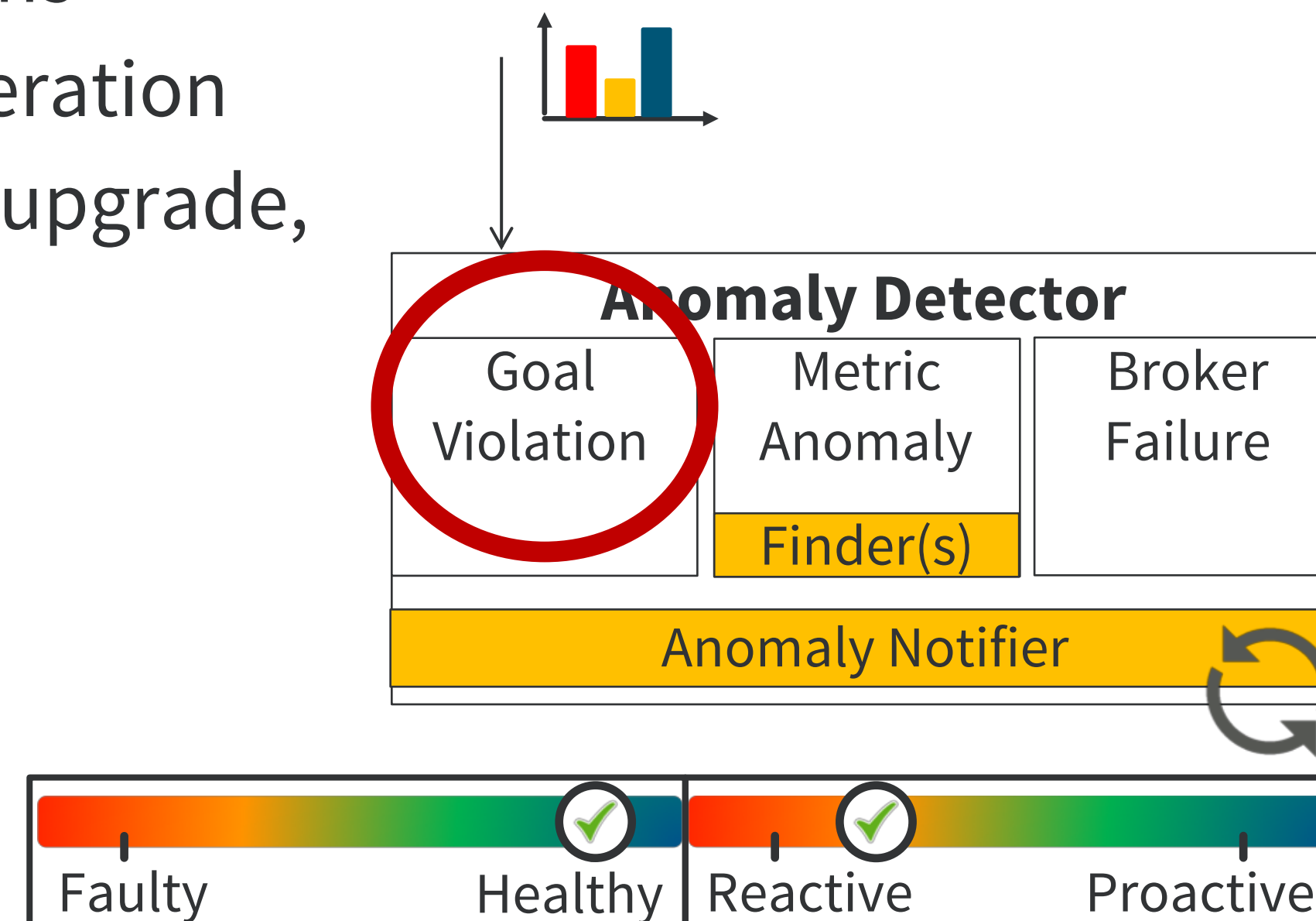👍👎 : Faulty vs Healthy Cluster

🛡⚡ : Reactive vs. Proactive Mitigation

# Anomaly Detector: Goal Violations and Self-Healing

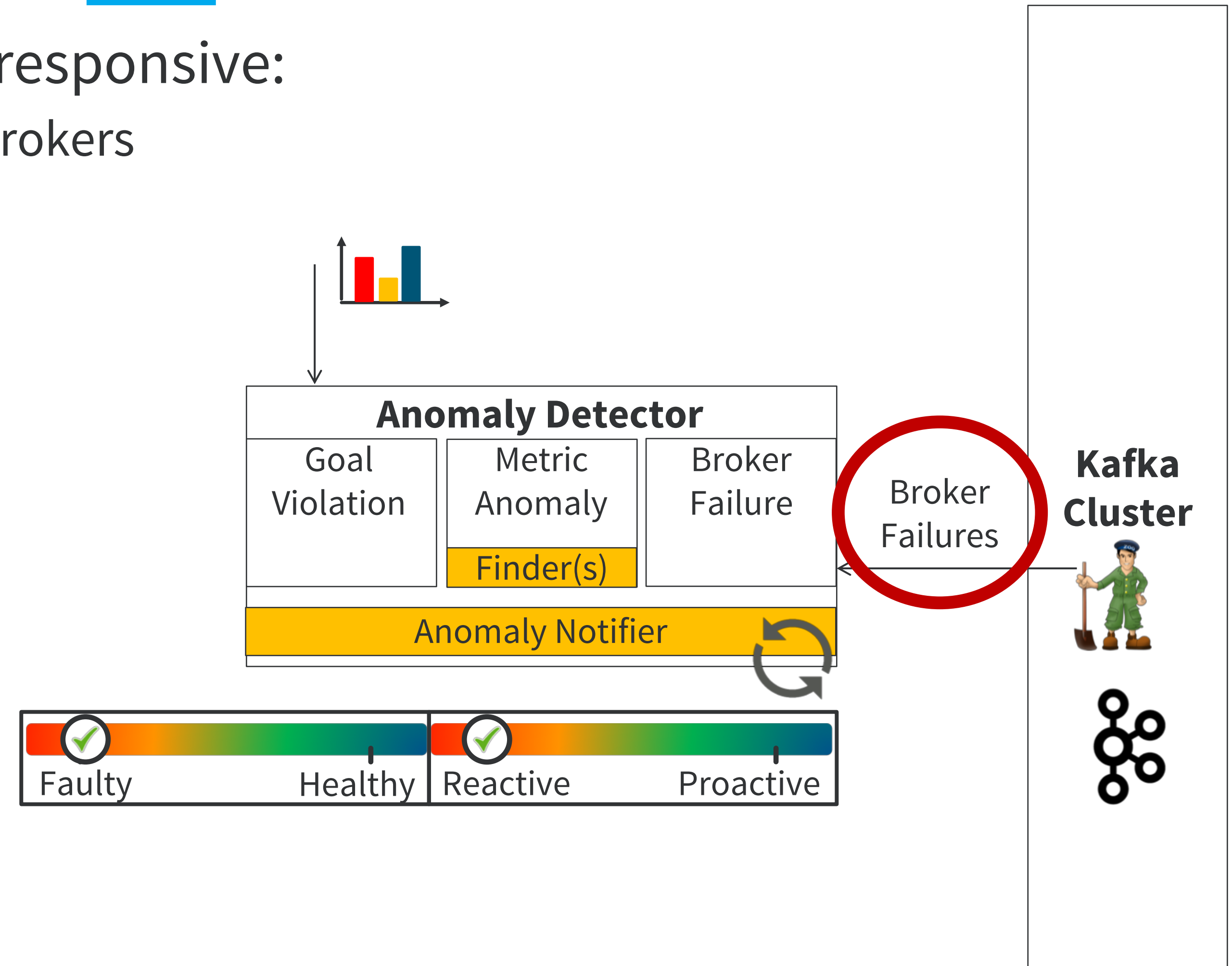Checks for the violation of the anomaly detection goals

- Identifies *fixable* and *unfixable* goal violations
- Self-healing triggers a cluster rebalance operation
- Avoids false positives due to broker failure, upgrade, restart, or release certification

# Anomaly Detector: Broker Failures

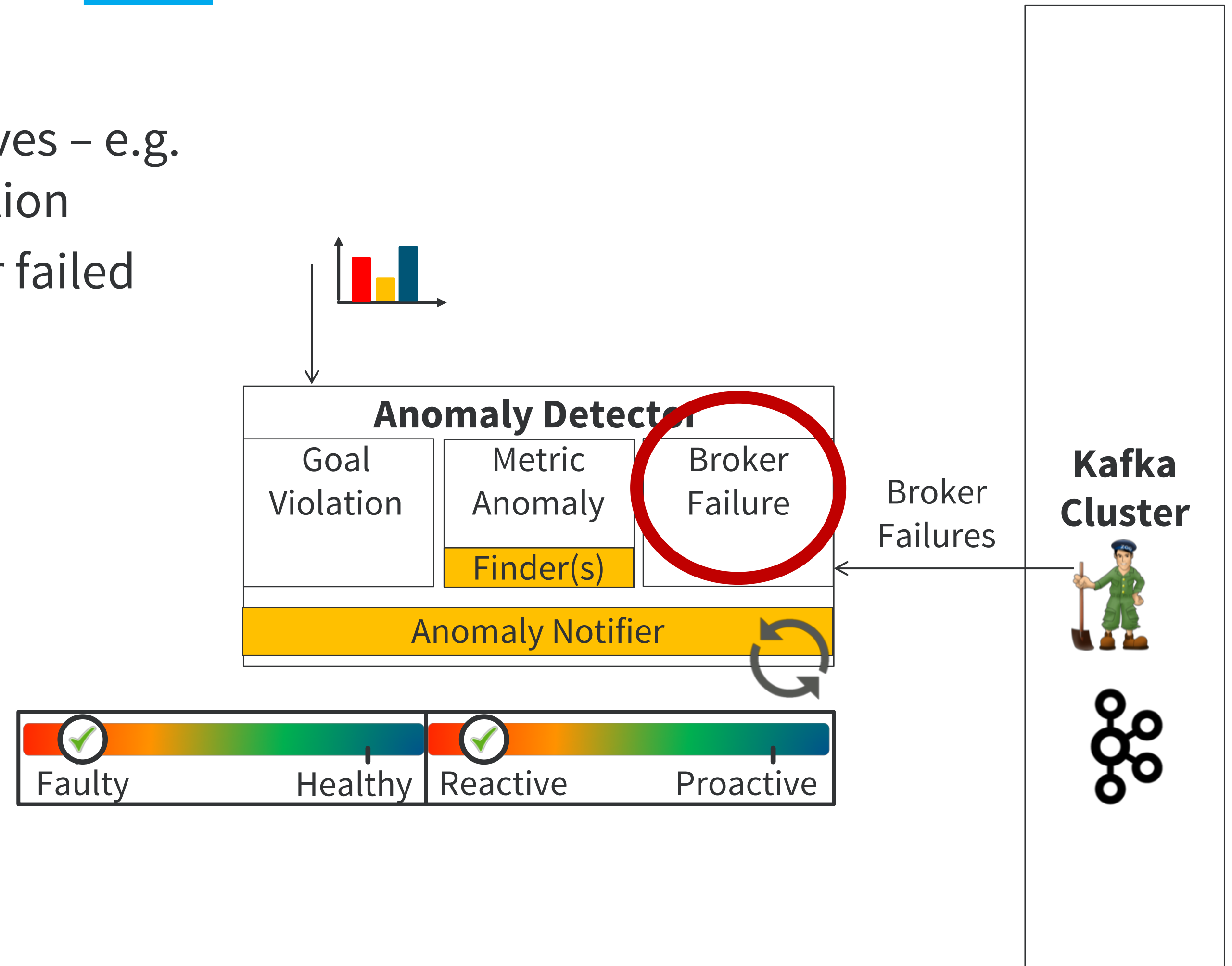Concerned with whether brokers are responsive:
- Ignores the internal state deterioration of brokers
- Identifies fail-stop failures

# Anomaly Detector: Broker Failures and Self-Healing

Checks for broker failures:

- Enables a grace period to lower false positives – e.g. due to upgrade, restart, or release certification
- Self-healing triggers a remove operation for failed brokers

37

# Anomaly Detector: Reactive Mitigation
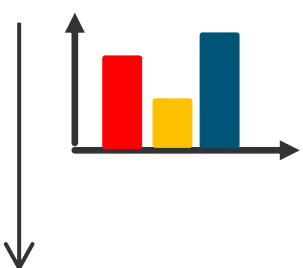
Cluster maintenance becomes costly

Requires immediate attention of affected services

Poor user experience due to frequent service interruptions

Server & network failures

$\sim$ 
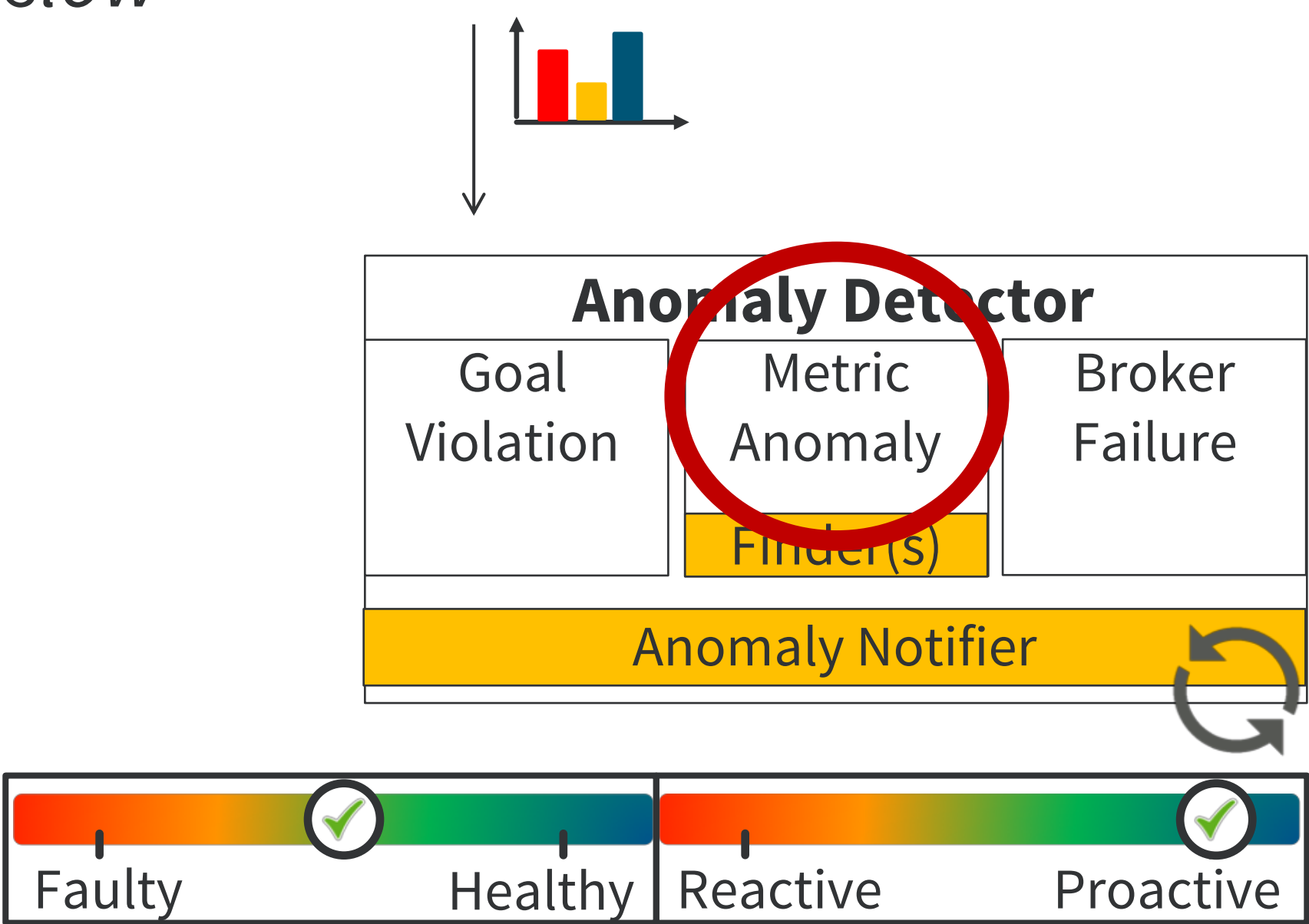- Size of clusters
- Volume of user traffic
- Hardware degradation

**Anomaly Detector**

| Goal Violation | Metric Anomaly | Broker Failure |
|---|---|---|
| | Finder(s) | |

Anomaly Notifier

# Anomaly Detector: Metric Anomaly

Checks for abnormal changes in broker
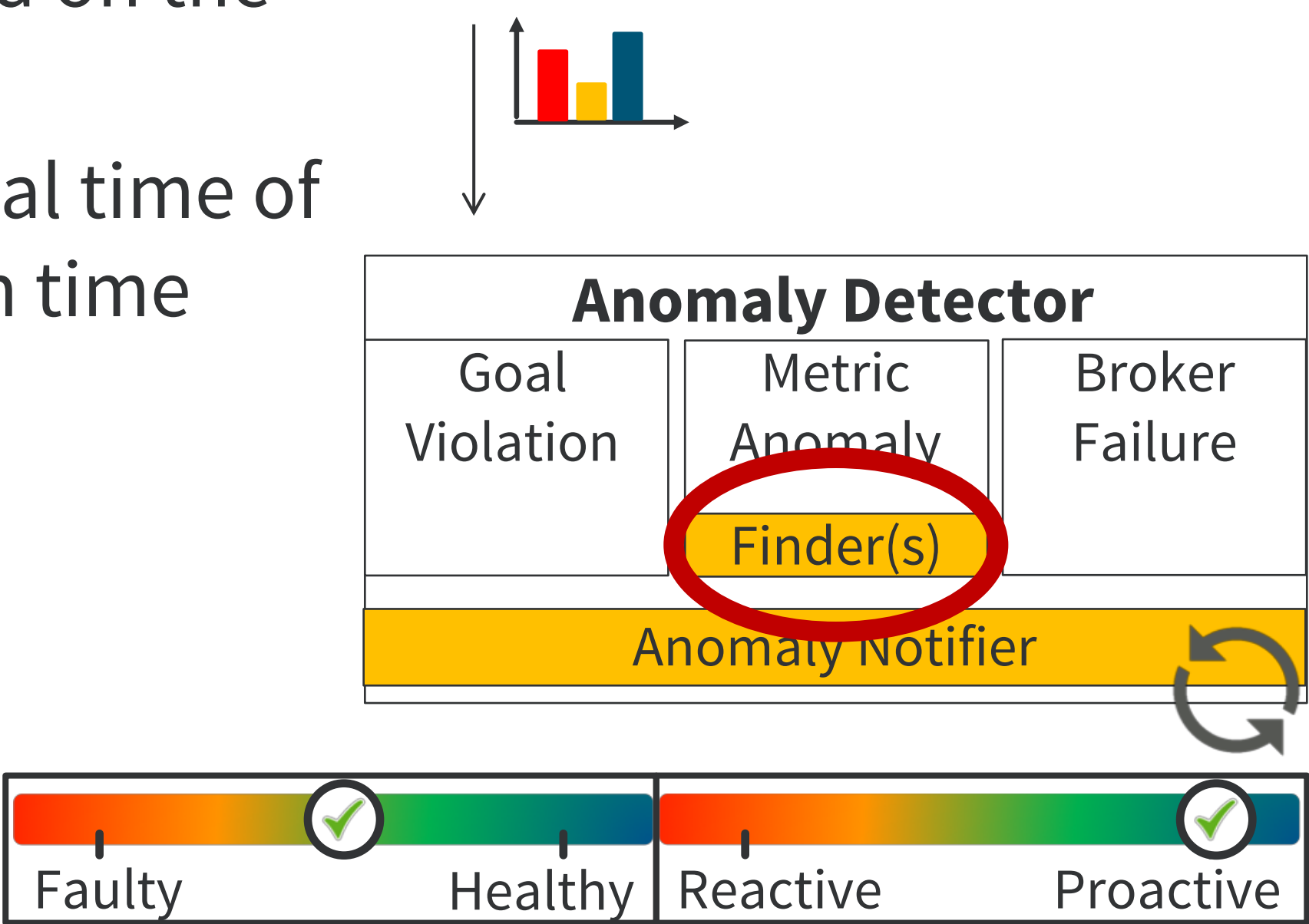metrics – e.g. a recent spike in log flush time:

- Self-healing triggers a demote operation for *slow* brokers



Anomaly Detector

| Goal Violation | Metric Anomaly | Broker Failure |

Finder(s)

Anomaly Notifier

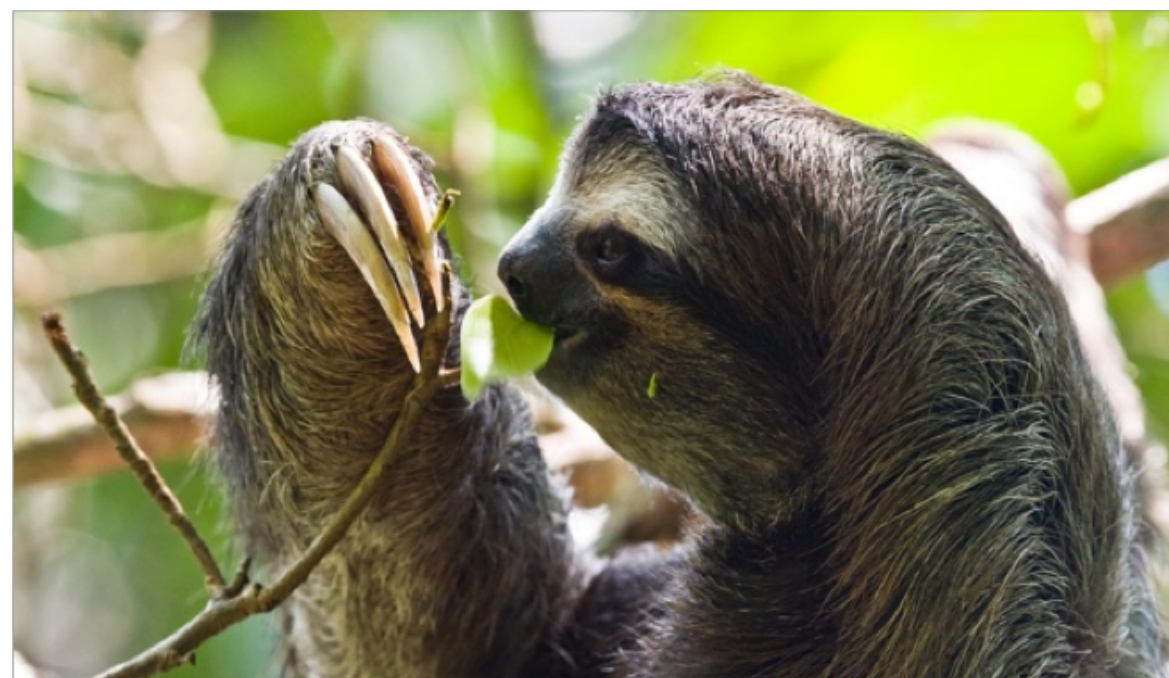Faulty    Healthy    Reactive    Proactive

# Anomaly Detector: Metric Anomaly

Compares current and historical metrics to detect slow brokers:

- The comparison in the default finder is based on the percentile rank of the latest metric value
- Metrics of interest are configurable – e.g. local time of produce / consume / follower fetch, log flush time
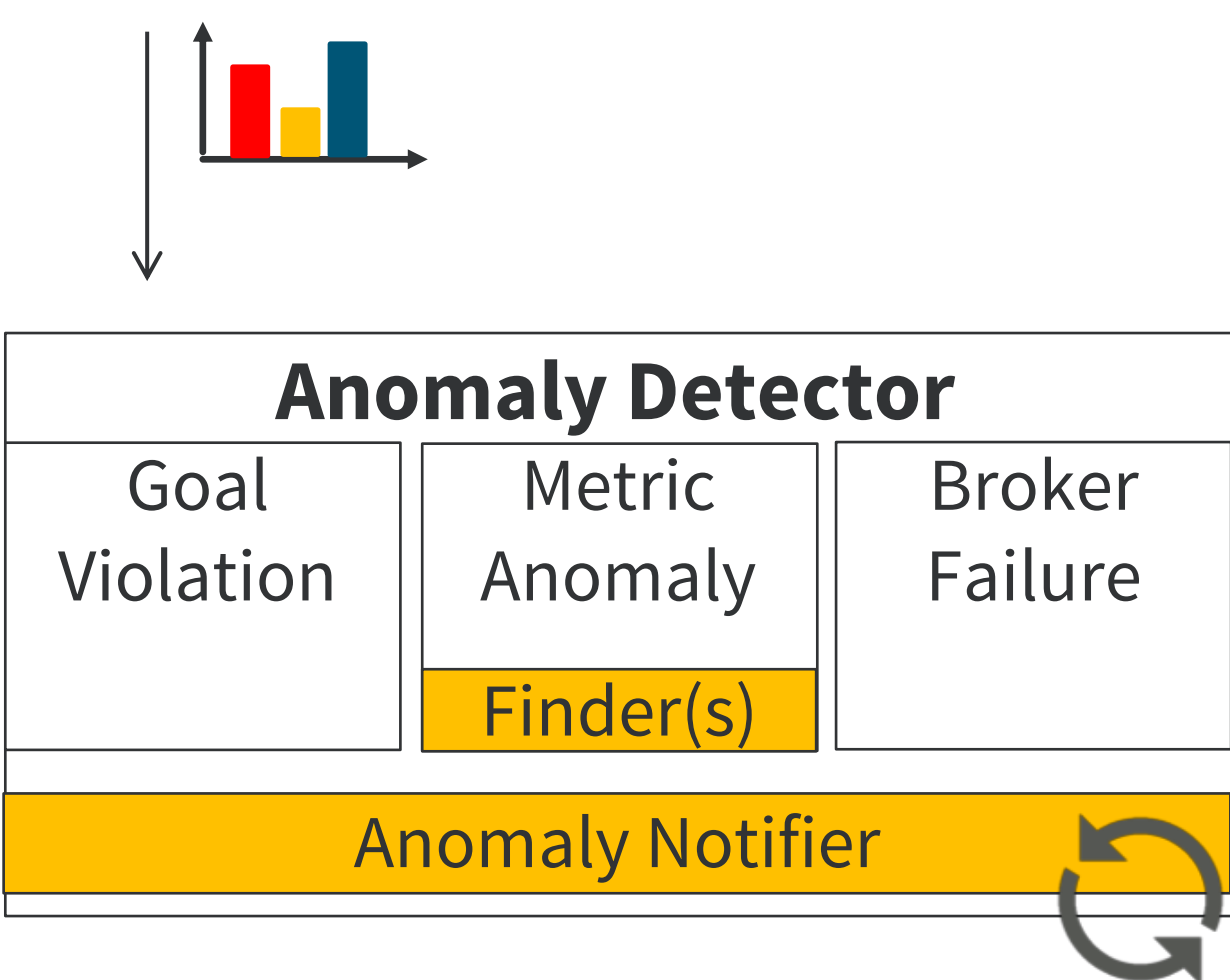- Supports multiple active finders

# Anomaly Detector: Proactive Mitigation

*In-place fix of slow / faulty brokers is non-trivial*

- The root cause could be a hardware issue (e.g. a misbehaving disk), a software glitch, or a traffic shift
- Hence, the mitigation strategies are agnostic of the particular issue with the broker

**Anomaly Detector**

| Goal Violation | Metric Anomaly | Broker Failure |
| --- | --- | --- |
| | Finder(s) | |

Anomaly Notifier

# REST API

Supports sync and async endpoints including:

**REST API**

GET
- Cluster Load
- Partition Load
- Proposals
- Kafka Cluster State
- Cruise Control State
- User Tasks

POST
- Add / Remove / Demote Broker
- Rebalance Cluster
- Fix Offline Replicas (JBOD)
- Stop Ongoing Execution
- Pause / Resume Sampling
- Admin – ongoing behavior changes

GUI & multi-cluster management

# Managing the Manager – *Monitoring Cruise Control*

Reported *JMX* metrics include:

| Executor | : Started, stopped, and ongoing executions in different modes, and the status of balancing tasks |

| Anomaly Detector | : Broker failure, goal violation, and metric anomaly rate |

| Monitor | : Cluster model and sampling performance |

| Analyzer | : Stats on proposal generation |

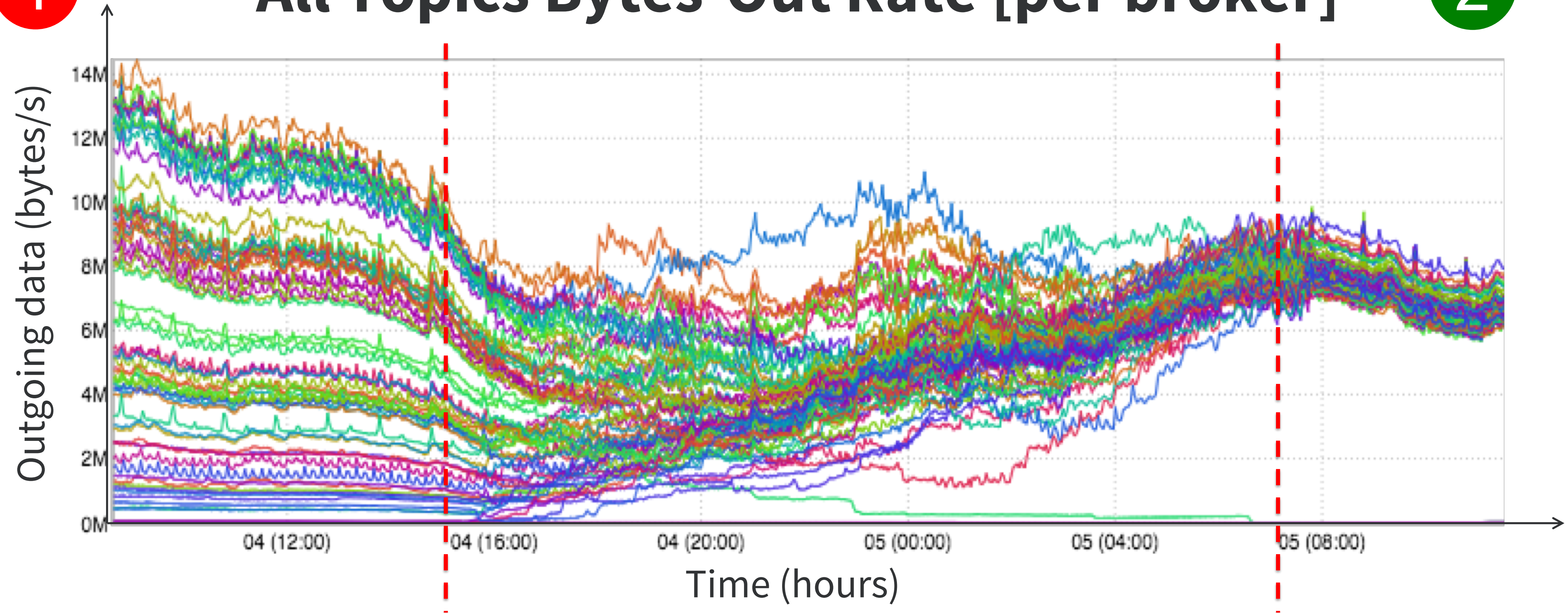# Evaluation: Remove Brokers and Rebalance

**1**

## All Topics Bytes-In Rate [per broker]

**2**

# Evaluation: Remove Brokers and Rebalance

**All Topics Bytes-Out Rate [per broker]**

# Evaluation: Remove Brokers and Rebalance

**1**

**Number of Partitions [per broker]**

**2**

# Summary

A system that provides effortless management of Kafka clusters

✔ Admin Operations for Cluster Maintenance

✔ Anomaly Detection with Self-healing

✔ Real-Time Monitoring of Kafka Clusters

⌛ Integration with Other Systems – e.g. Apache Helix

# More…

___

**Fork** 196
**Watch** 116
**Star** 1,052

: Open source repository
(**https://github.com/linkedin/cruise-control**)

: Gitter room (**https://gitter.im/kafka-cruise-control**)

: UI (**https://github.com/linkedin/cruise-control-ui**)

# Rate today's session



Session page on conference website



O'Reilly Events App