# Getting to Know Apache Kafka's Architecture

**Ryan Plant**
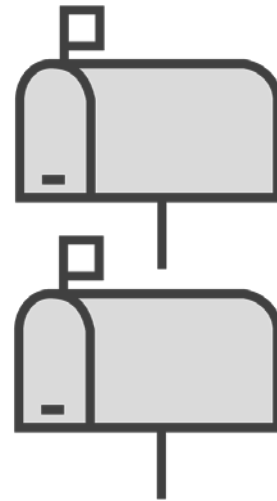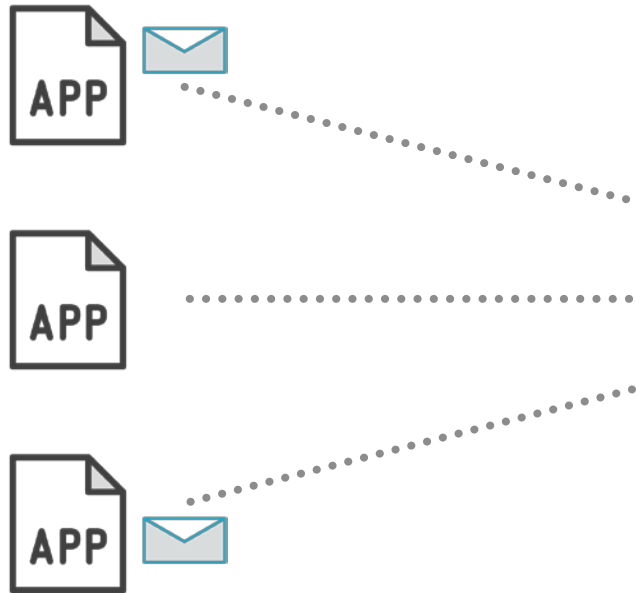COURSE AUTHOR

@ryan_plant    blog.ryanplant.com
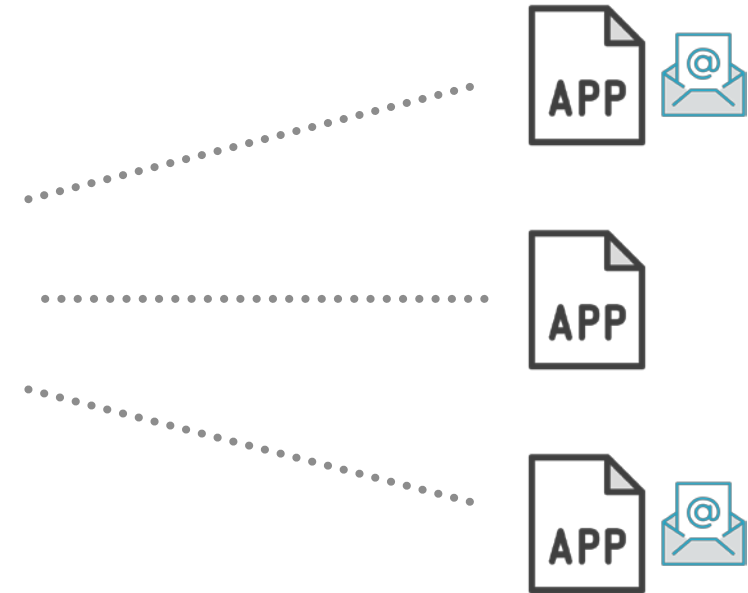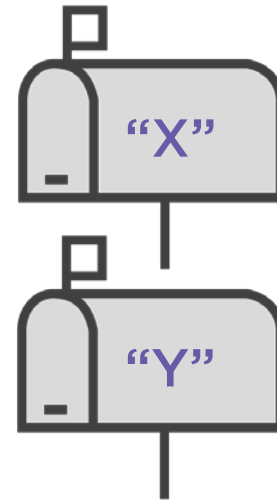
# Apache Kafka as a Messaging System

**Producers**

**Consumers**

# Apache Kafka as a Messaging System

**Producers**

**Topics**

**Consumers**

To: "X"

Retrieve: "X"

APP

APP

APP

APP

"X"

"Y"

APP

APP

To: "Y"

Retrieve: "Y"

# Apache Kafka as a Messaging System
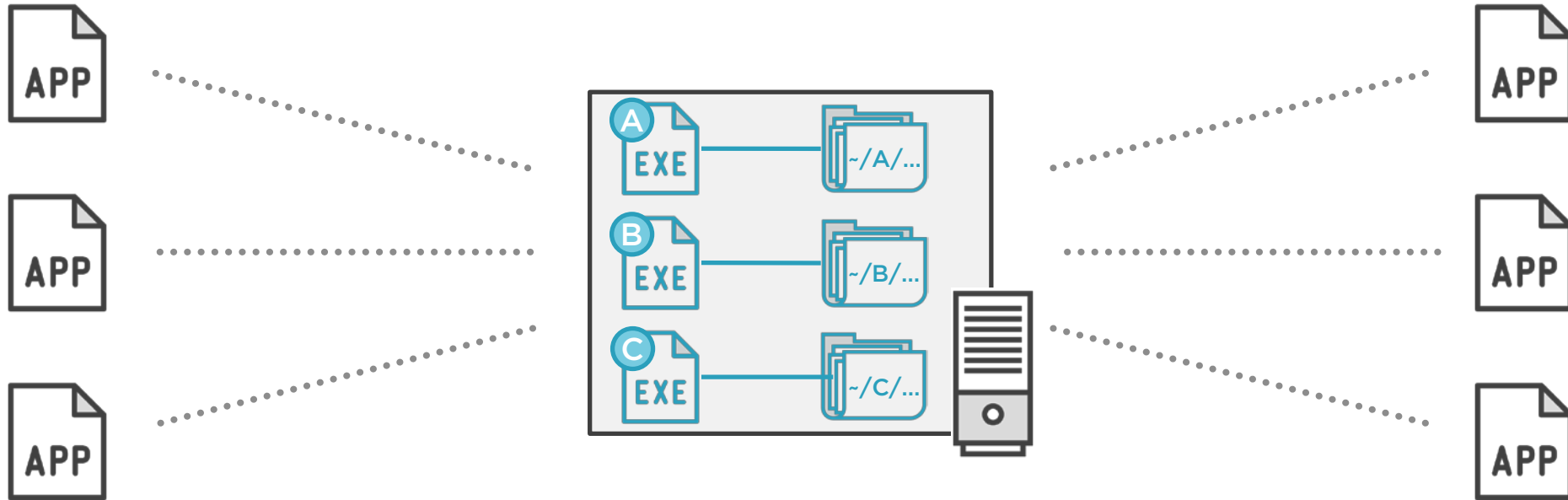
APP

APP

APP

**Broker**

APP

APP

APP

# Apache Kafka as a Messaging System

**Producers**

**Broker**

**Consumers**

# How Apache Kafka Starts to Differentiate

**Producers**

**Consumers**

LinkedIn: 1,400 brokers => 2 petabytes per week

| Broker | Broker |
|--------|--------|
| Broker | Broker |

"A high-throughput distributed messaging system."

# The Apache Kafka Cluster

# The Apache Kafka Cluster

**Producers**

**Cluster**
**Size:** 1

**Consumers**

Broker

Broker

Broker

Broker

# The Apache Kafka Cluster

**Producers**

**Cluster Size:** **2**

**Consumers**

# The Apache Kafka Cluster



**Producers**

**Cluster
Size:  2**

**Consumers**

# Distributed Systems

**KAFKA BROKERS**

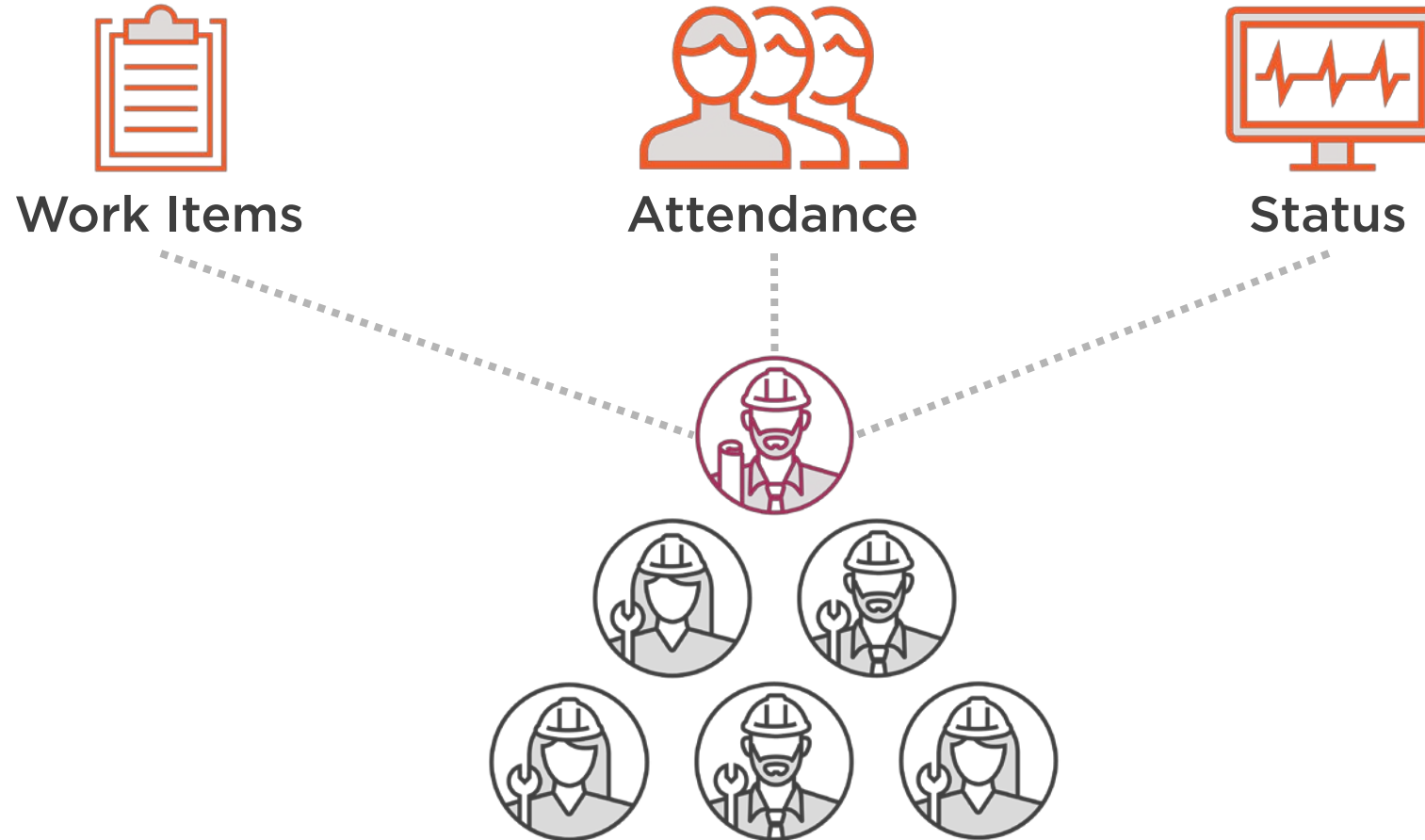Collection of resources that are instructed to achieve a specific goal or function

Consist of multiple workers or nodes

The system of nodes require coordination to ensure consistency and progress towards a common goal

Each node communicates with each other though messages

Distributed Systems: Controller Election

Work Items    Attendance    Status

# Distributed Systems: The Cluster



KAFKA CLUSTER

# Distributed Systems: Getting Work Done



PRODUCER

**Worker availability and health**

**Task redundancy**

# Distributed Systems: Getting Work Done (Reliably)

# Distributed Systems: Getting Work Done (Reliably)

Sources of Work in Apache Kafka

PRODUCER

KAFKA CLUSTER

CONSUMER

# Distributed Systems: Communication and Consensus

**Worker node membership and naming**

**Configuration management**

**Leader election**

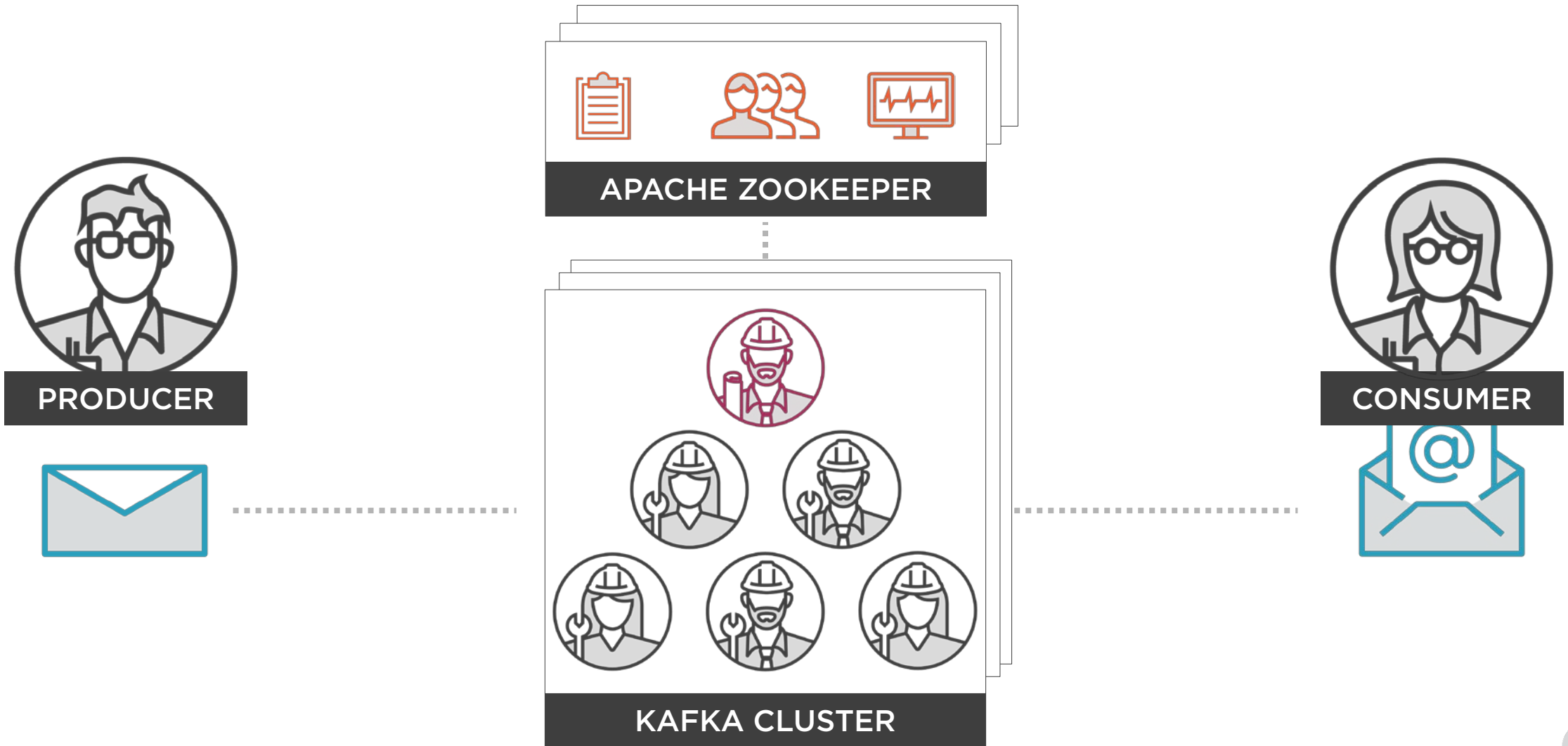**Health status**

# Apache Zookeeper

**Centralized service for maintaining metadata about a cluster of distributed nodes**

- Configuration information
- Heath status
- Group membership

**Hadoop, HBase, Mesos, Solr, Redis, and Neo4j**

**Distributed system consisting of multiple nodes in an "ensemble"**

# Apache Kafka's Distributed Architecture

# Summary

**Apache Kafka is a Pub-Sub messaging system, consisting of:**

- Producers and Consumers
- Brokers within a Cluster

**Characteristics of distributed systems**

- Worker node roles: Controllers, Leaders, and Followers
- Reliability through replication
- Consensus-based communication

**Role of Apache Zookeeper**