

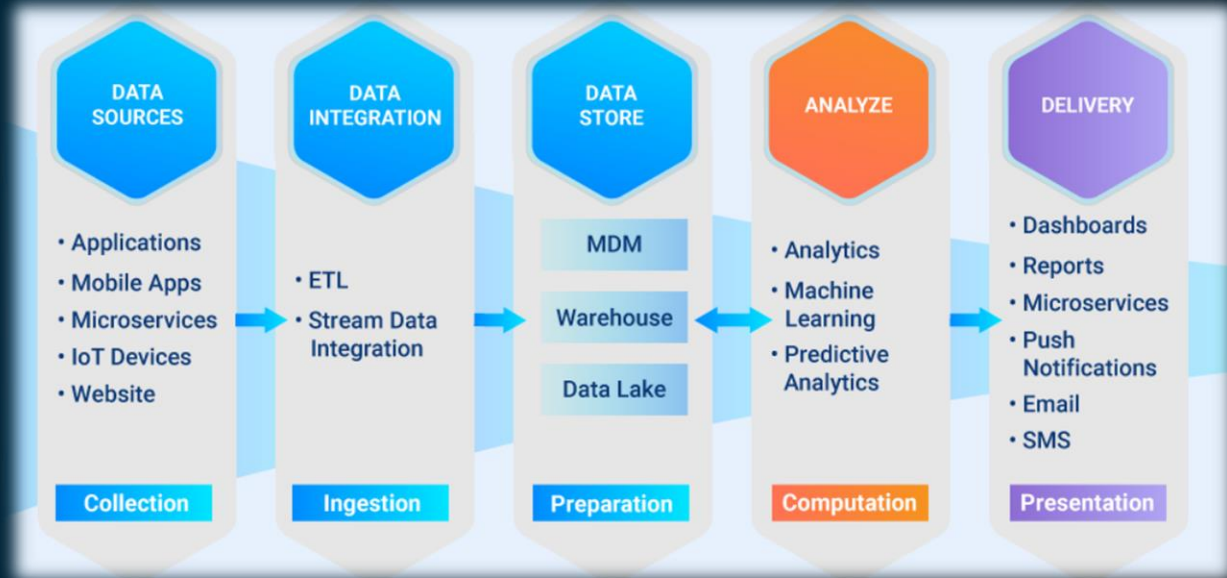
Apache Flume & ZooKeeper



Parham Pishro

Dr. Nasiri

Pipeline of Big Data





Apache Flume



Jul 26, 2012

Jul 2, 2013

May 20, 2015

Oct 4, 2017

Aug 16, 2022

1.2.0

1.3.1

1.4.0

1.5.0.1

1.6.0

1.7.0

1.8.0

1.9.0

1.10.1

1.11.0

Jan 2, 2013

Jul 16, 2014

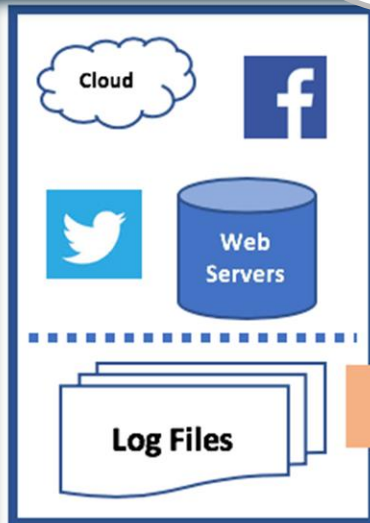
Oct 17, 2016

Jan 8, 2019

Oct 24, 2022

Introduction to Apache Flume

flume.apache.org

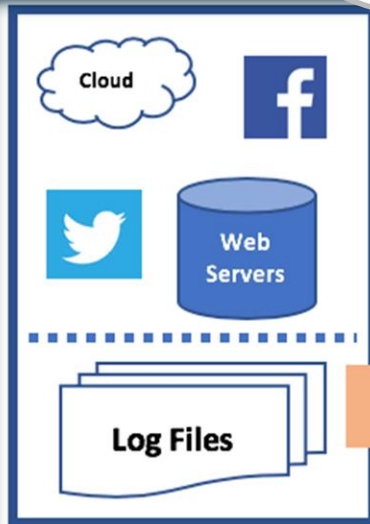


Flume is a distributed, reliable, and available service for efficiently collecting, aggregating, and moving large amounts of streaming event data.



Introduction to Apache Flume

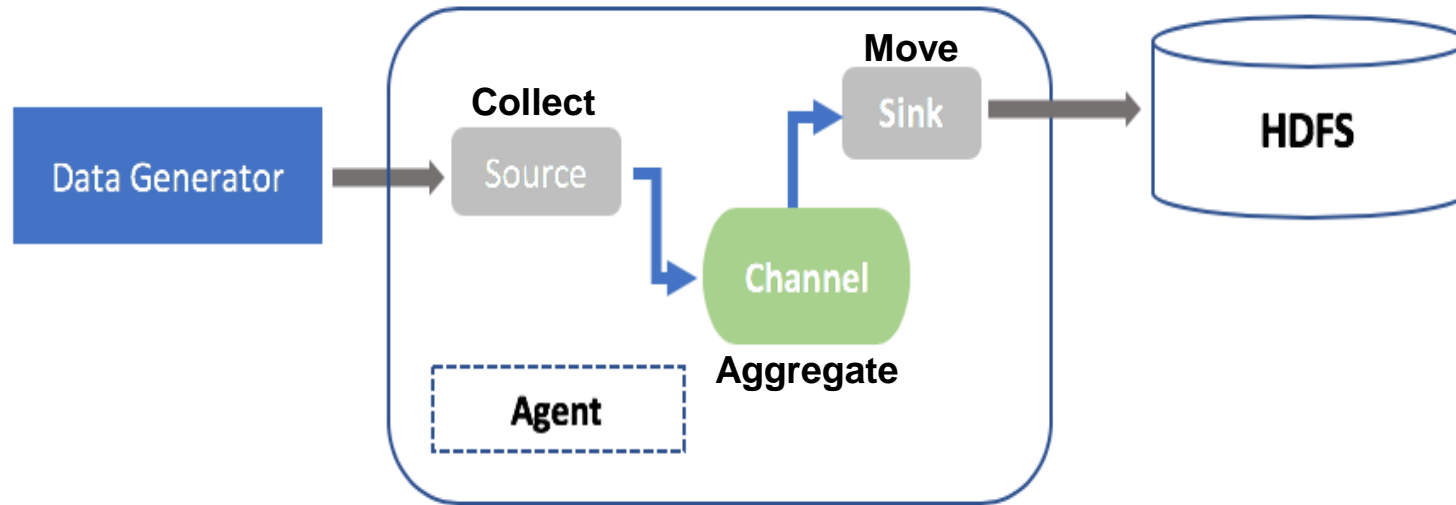
flume.apache.org



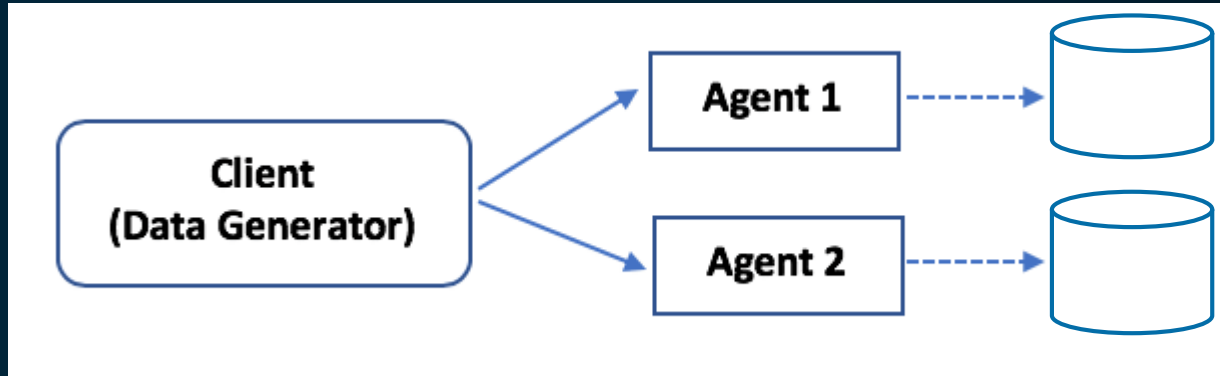
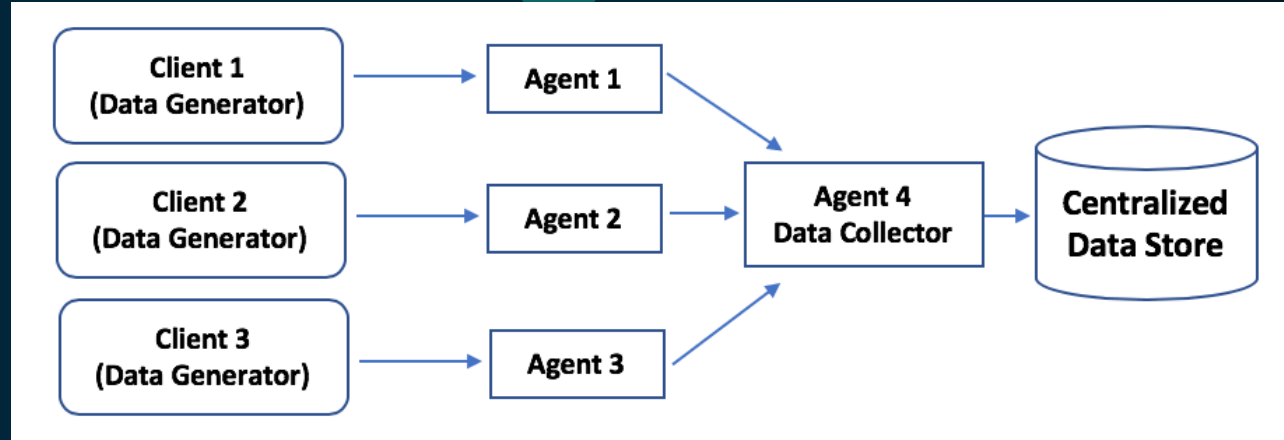
Flume is a distributed, reliable, and available service for efficiently collecting, aggregating, and moving large amounts of streaming data.



Architecture of Apache Flume

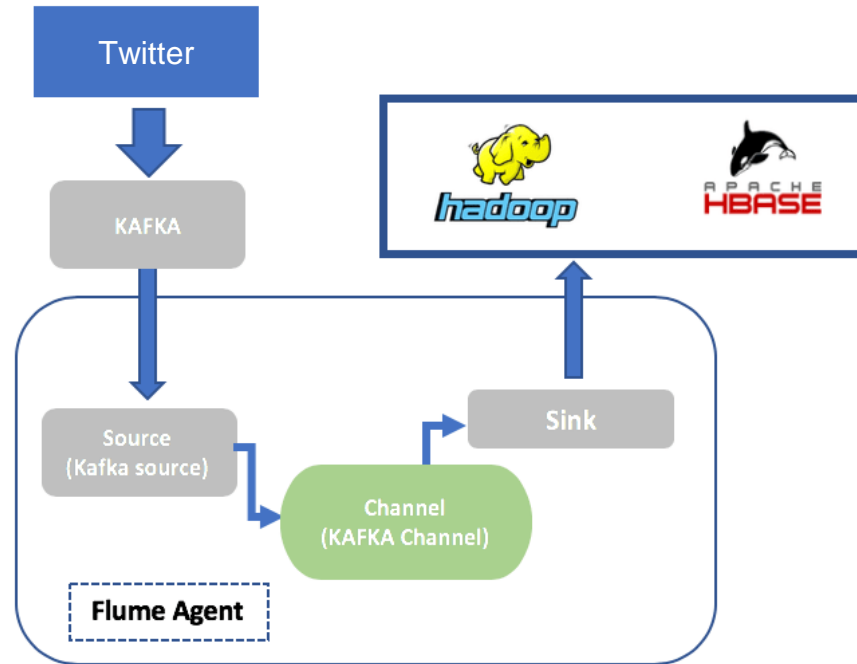


Two Types of Architecture



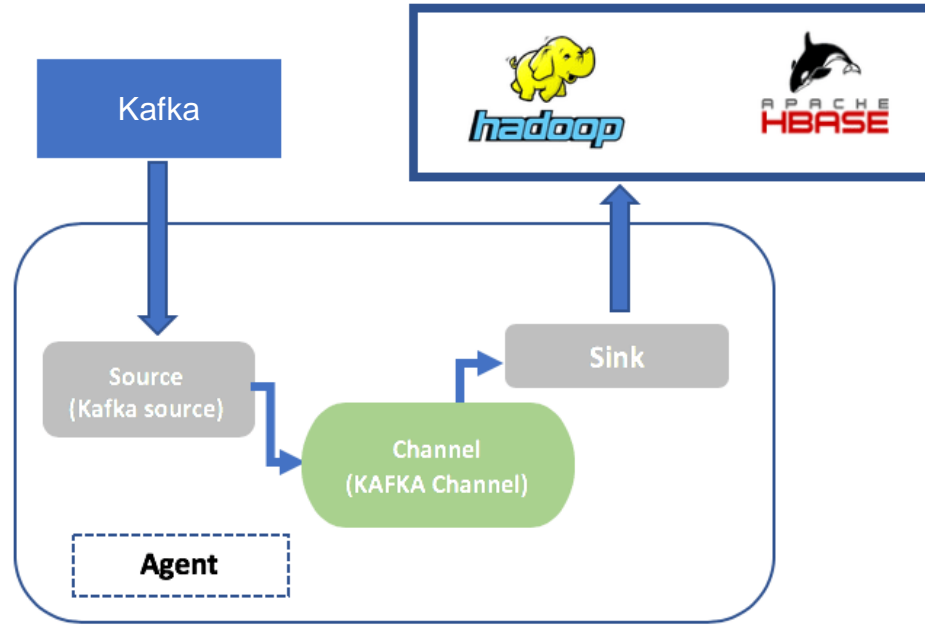
Example

Streaming Log Data to HDFS from Twitter



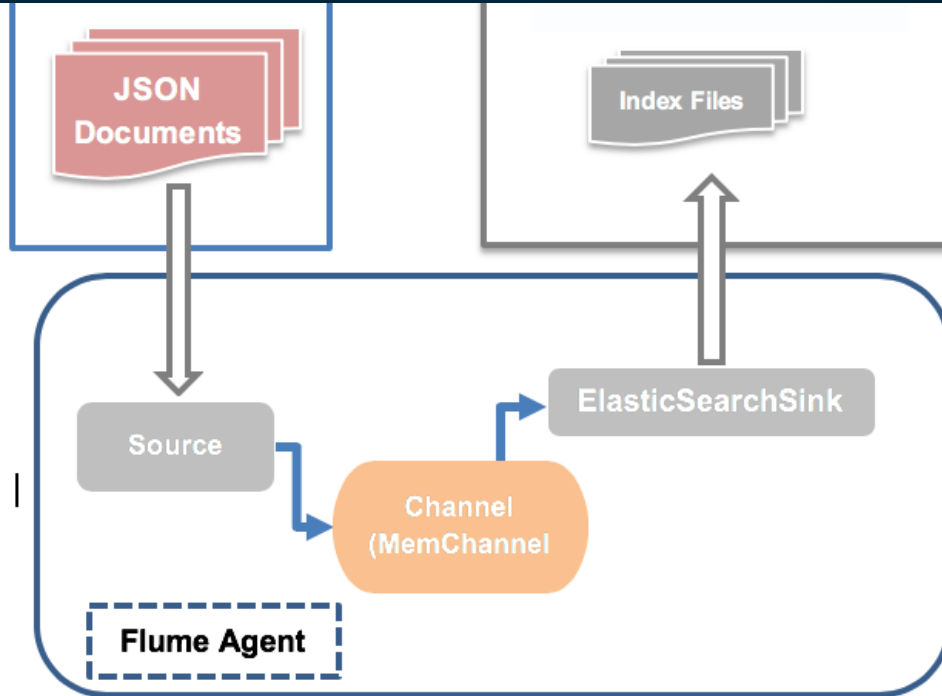
Example

Streaming Log Data from Kafka to HDFS



Example

Streaming Log Data to Elasticsearch



The Limitations of Apache Flume

- Complexity and difficulty of architecture to manage and maintain (From multiple sources to multiple destinations)
- Streaming is not 100% real-time
- Weakness to identify duplicate data



Apache ZooKeeper

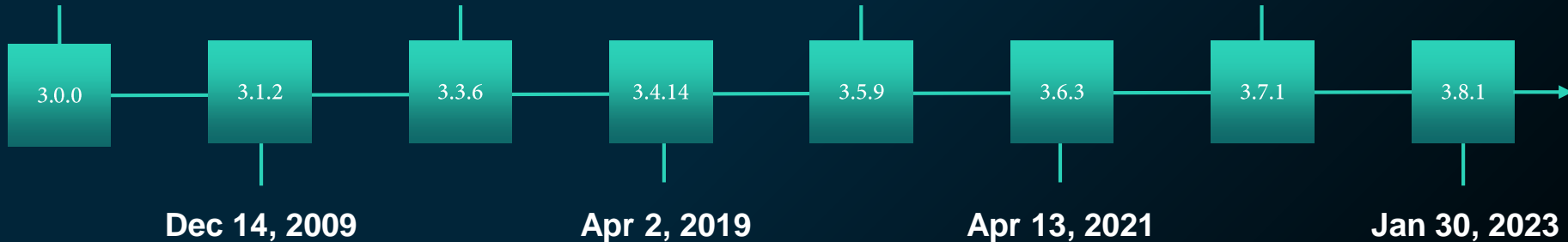


Oct 27, 2008

Aug 2, 2012

Jan 15, 2021

May 12, 2022



Why Do We Need Apache ZooKeeper?

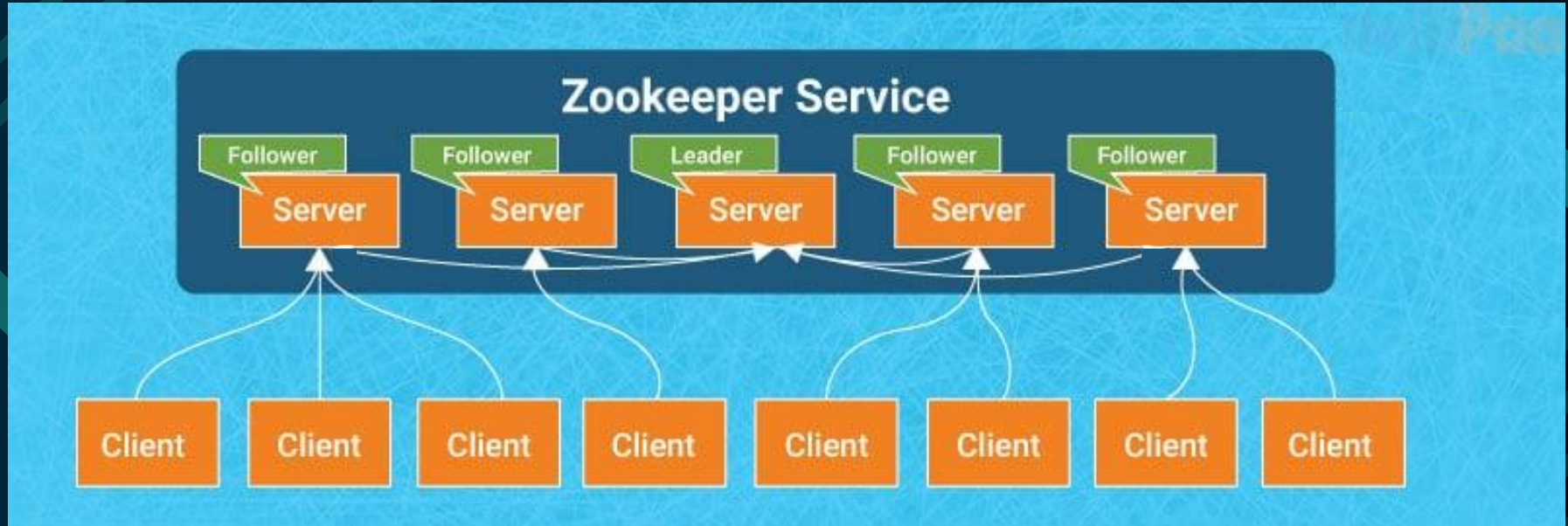
Coordinating Different Services in Hadoop Ecosystem
(Systems Configuration with Data Synchronization)



Why Do We Need Apache ZooKeeper?



Architecture of Apache ZooKeeper



Zookeeper Service



Server

Leader

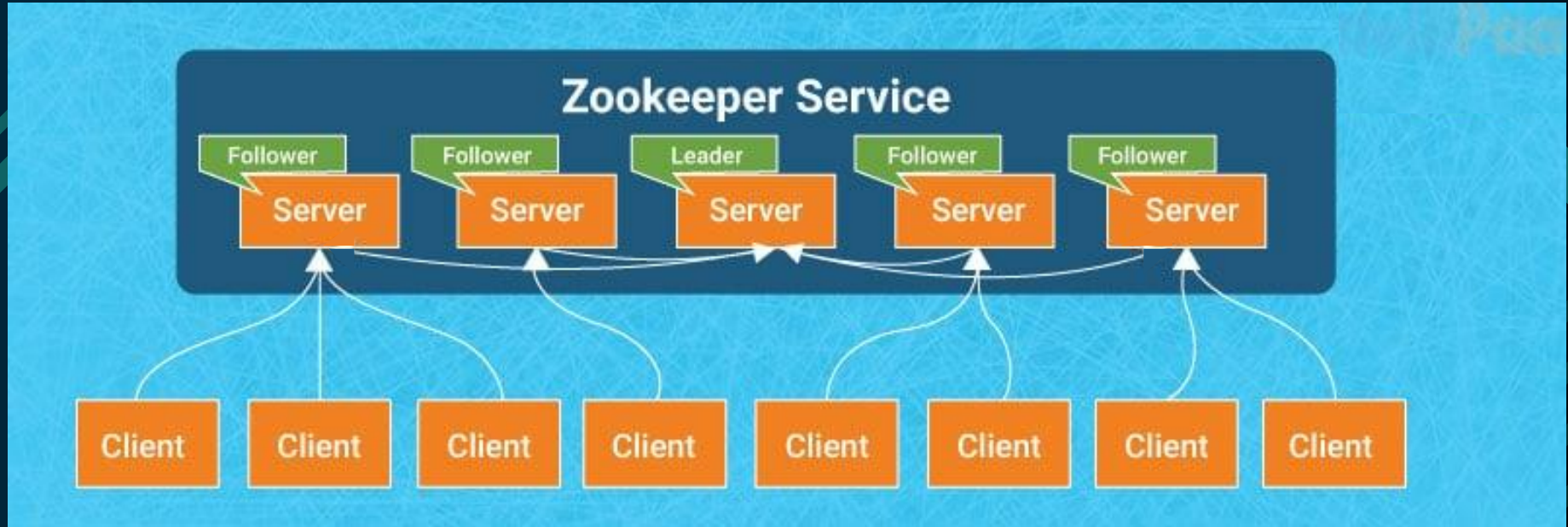
Follower

Client

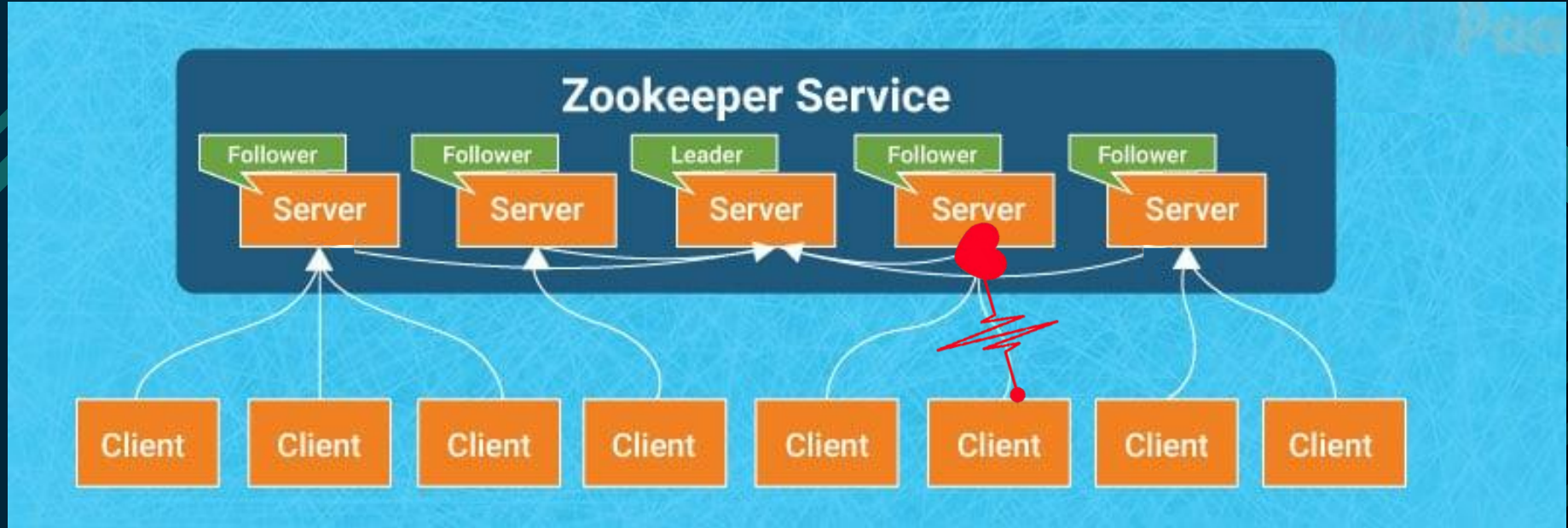
Working of Apache Zookeeper



Working of Apache Zookeeper



Working of Apache Zookeeper



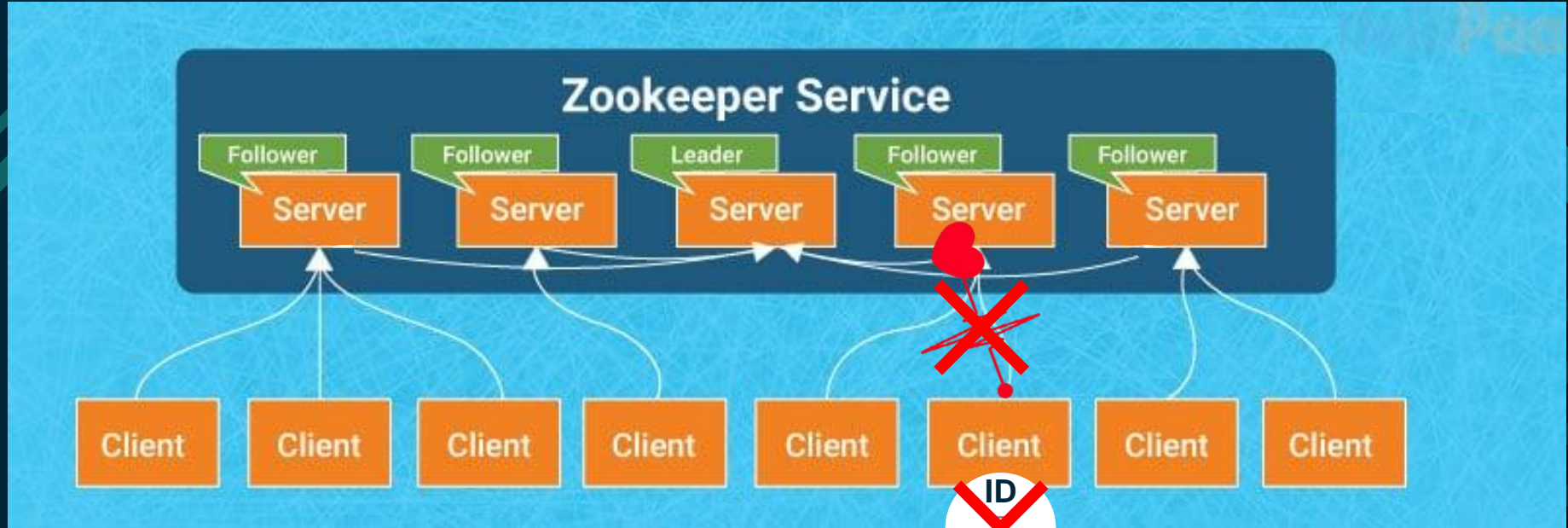
Working of Apache Zookeeper



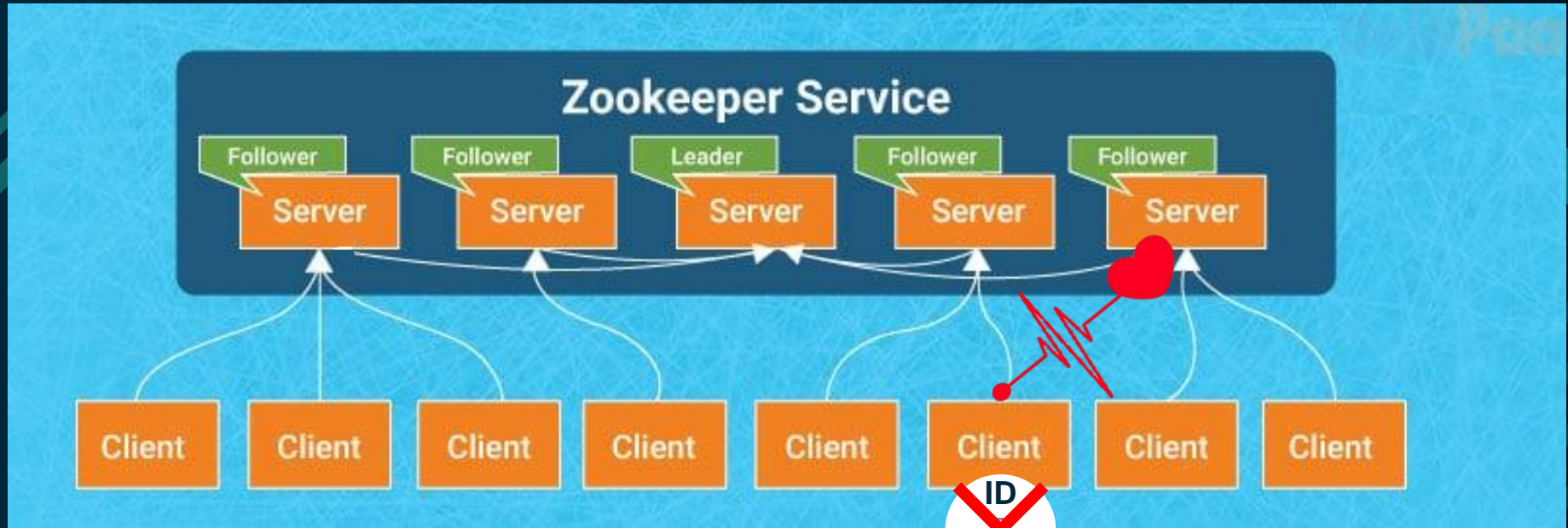
Working of Apache Zookeeper



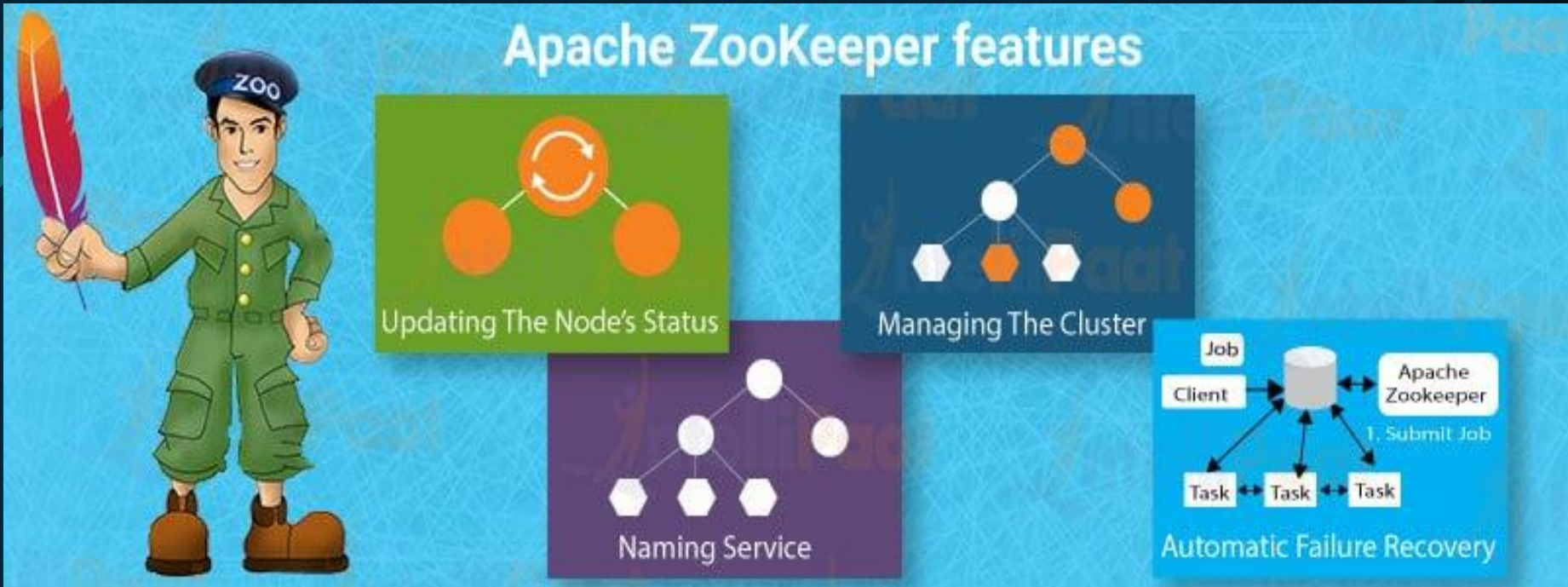
Working of Apache Zookeeper



Working of Apache Zookeeper



Features of Apache ZooKeeper



Benefits of Apache ZooKeeper



Simplicity: Coordination is done

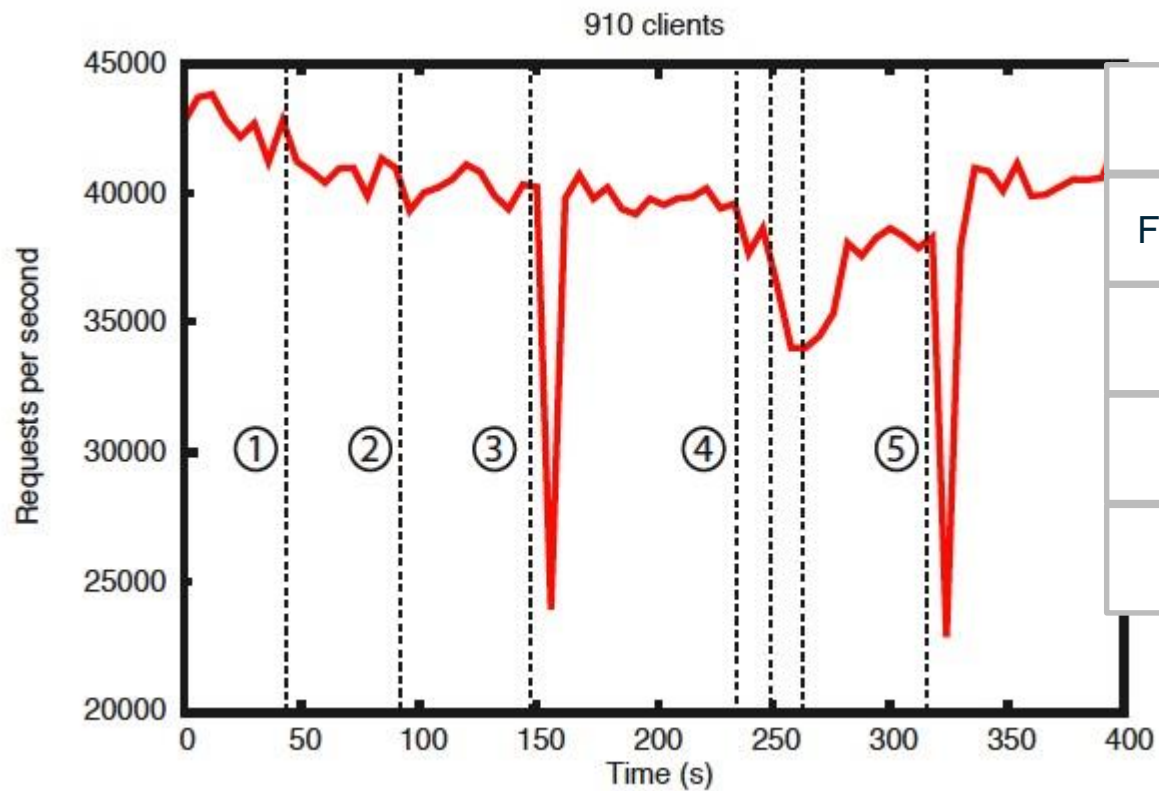
Reliability: $\left\lceil \frac{n}{2} \right\rceil + 1$

Order: Stamping each update with a number

Speed: Runs with a ratio of 10:1 (when 'reads' are more common)

Scalability: Enhancement the performance by deploying more machines

Benefits of Apache ZooKeeper



Simplicity

Failure and recovery of a follower

Failure and recovery of a different follower

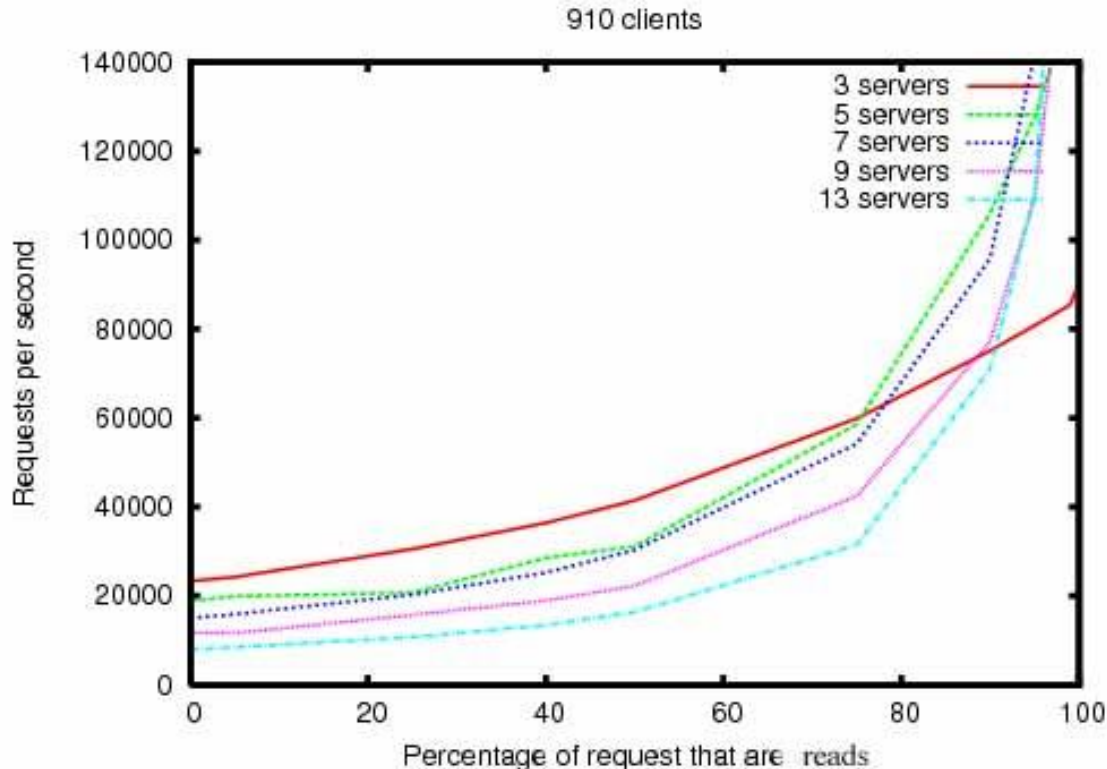
Failure and recovery of the leader

Failure and recovery of two followers

Failure and recovery of another leader

Scalability: Enhancement the performance by deploying more machines

Benefits of Apache ZooKeeper



Simpli

Reliab

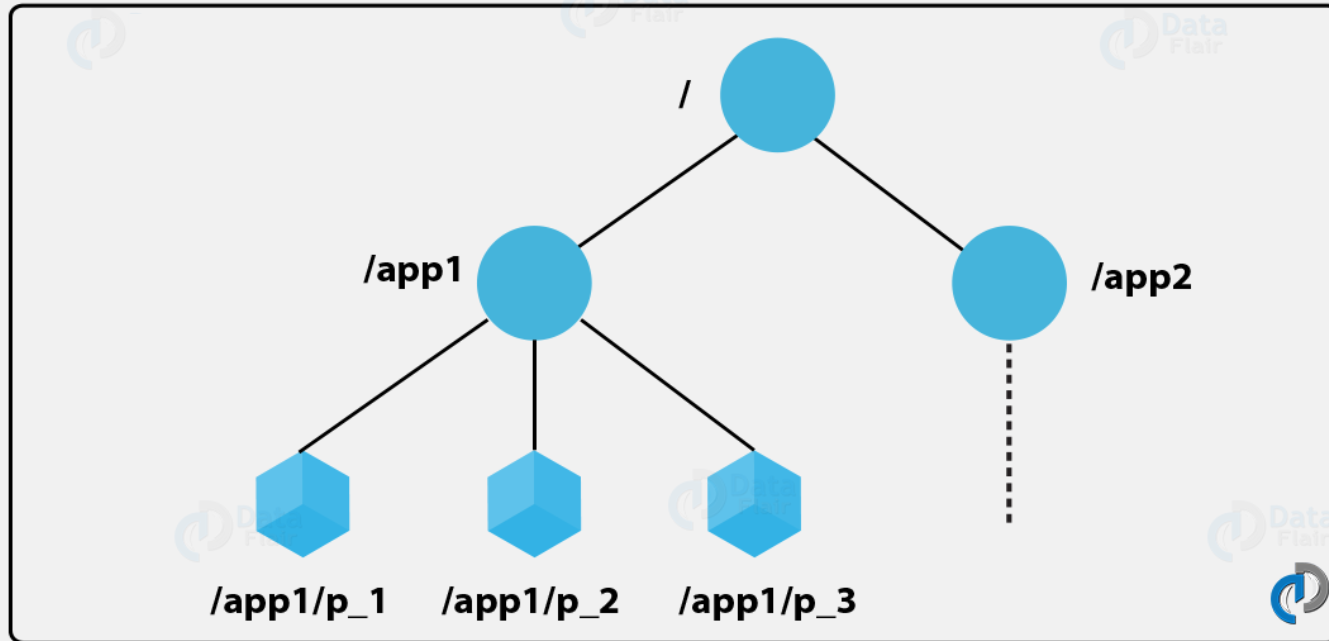
Order:

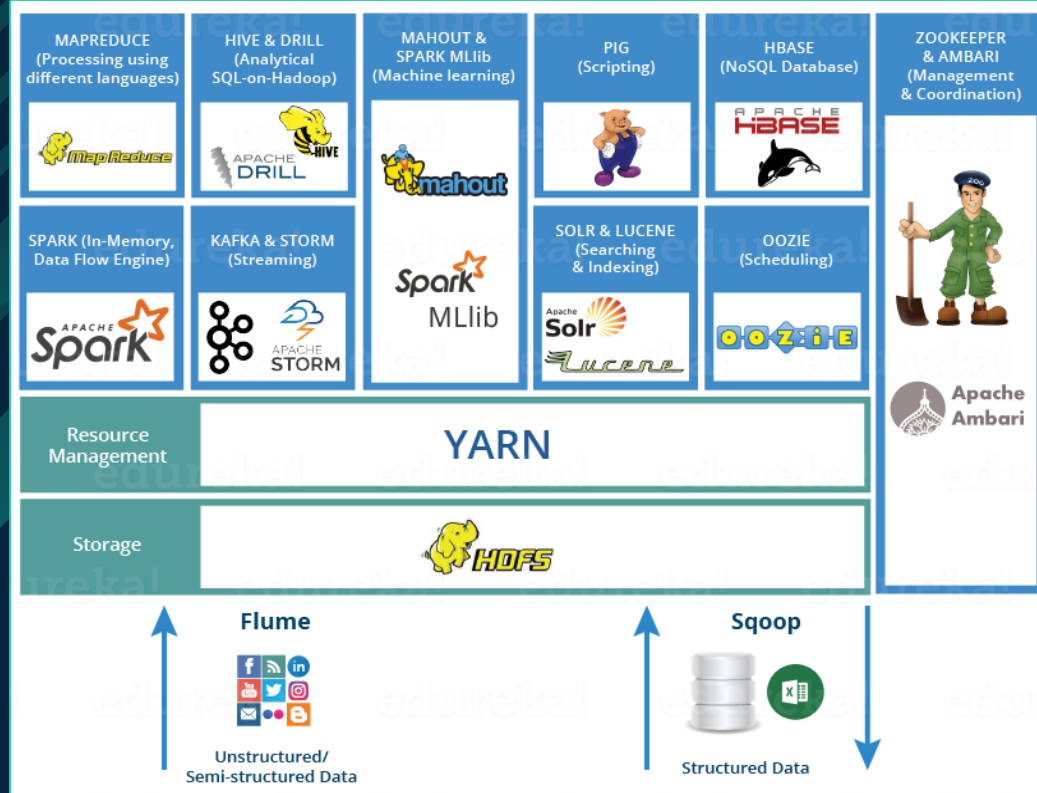
Speed:

Scalability: Enhancement the performance by deploying more machines

Data Model and The Hierarchical Namespace

Name = Sequence of path elements separated by a slash (/)
Every node is identified by a path.







**Thanks
for
Your
Attention**