



实时日志处理F l u m e和k a f k a技术分享

邵位

2016.5

1

Flume 框架、应用以及优化策略

2

kafka框架介绍和应用

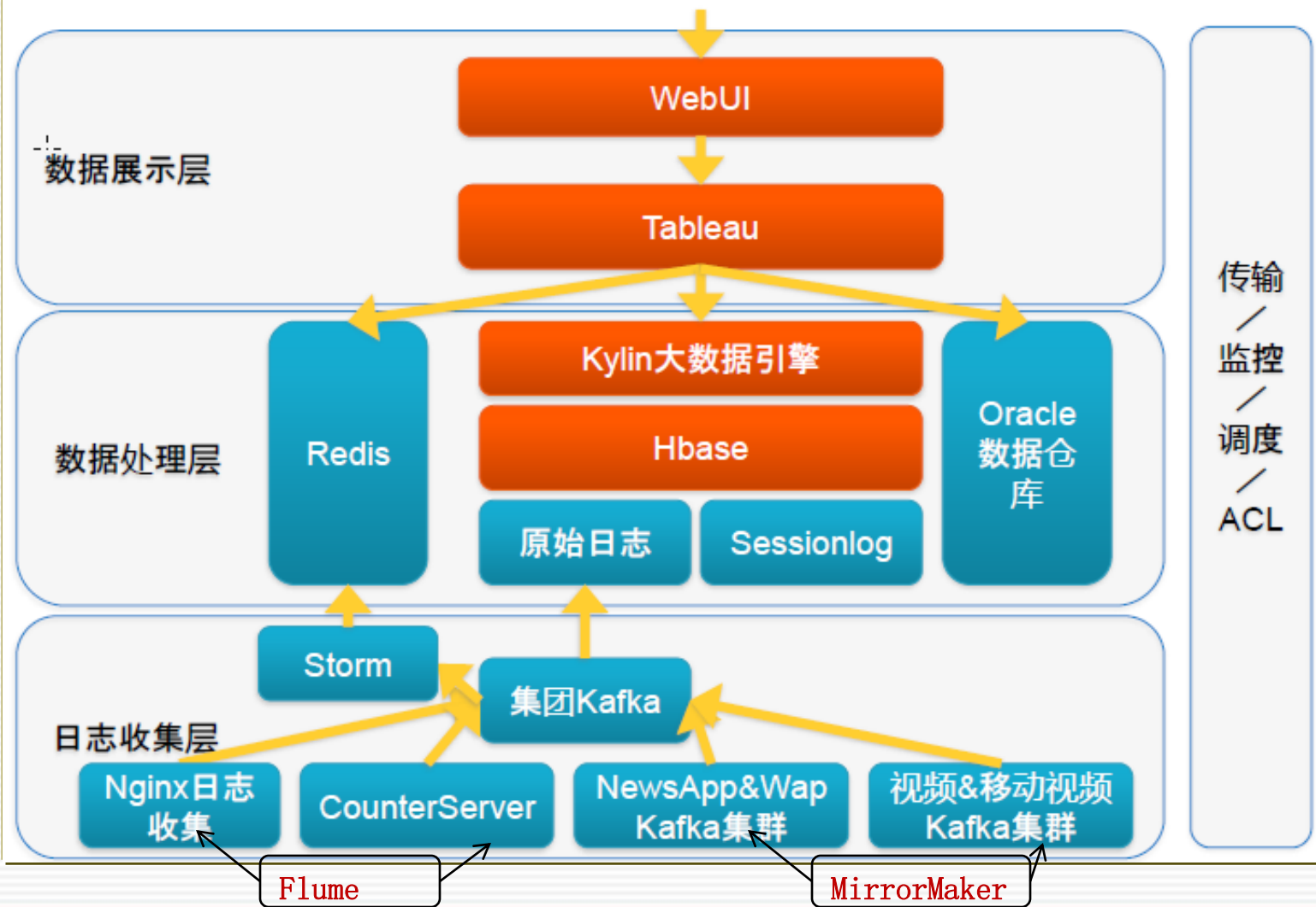
总体架构和Flume简介

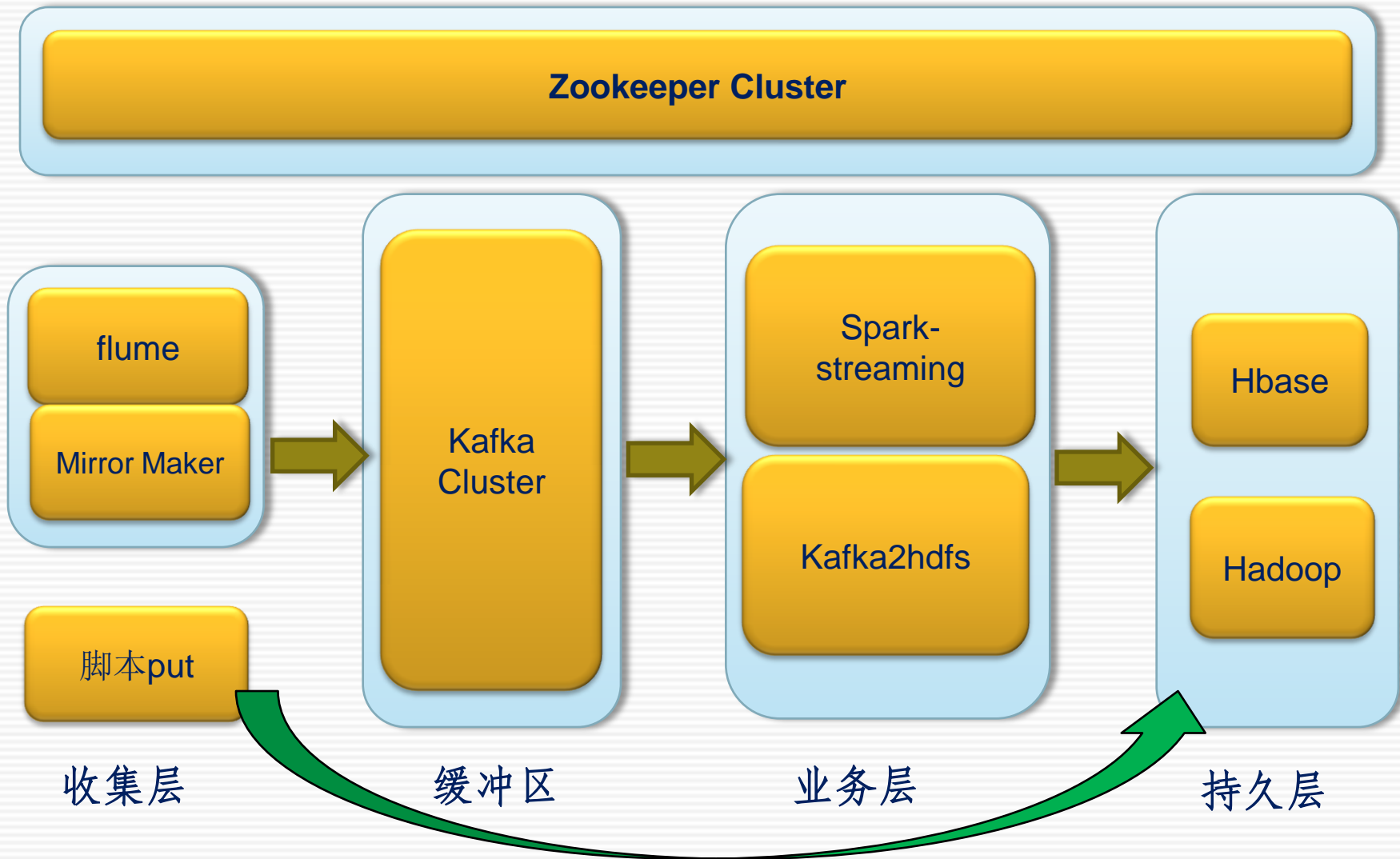
Memorychannel的通信机制和事务

Kafkachannel机制和调优

Flume的缺点和未来努力

大数据架构(2015) - 流程图





FLUME

012345678910111213141

0123456789101112131415161718191324654612316879841651

01234567891011121

314151617

Agent

01

Source

数据收集

02

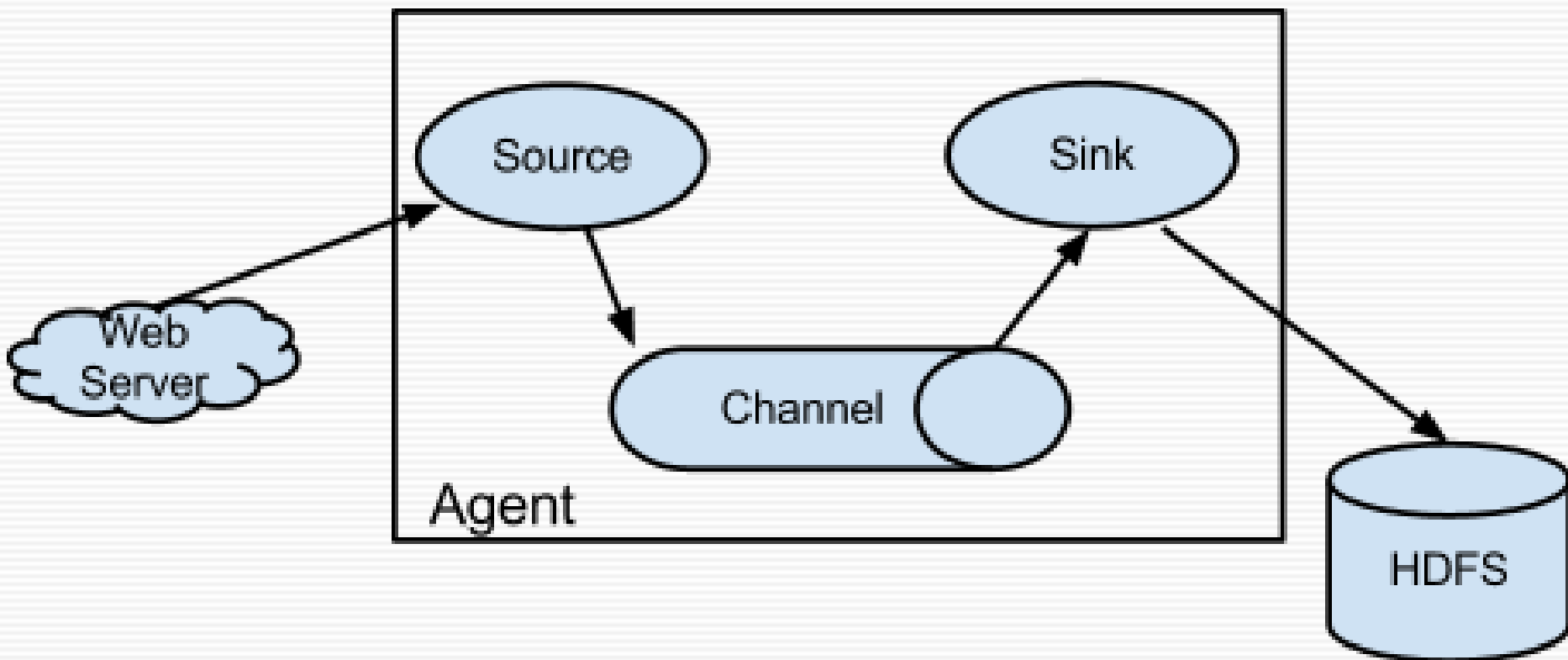
Channel

提供队列
缓存

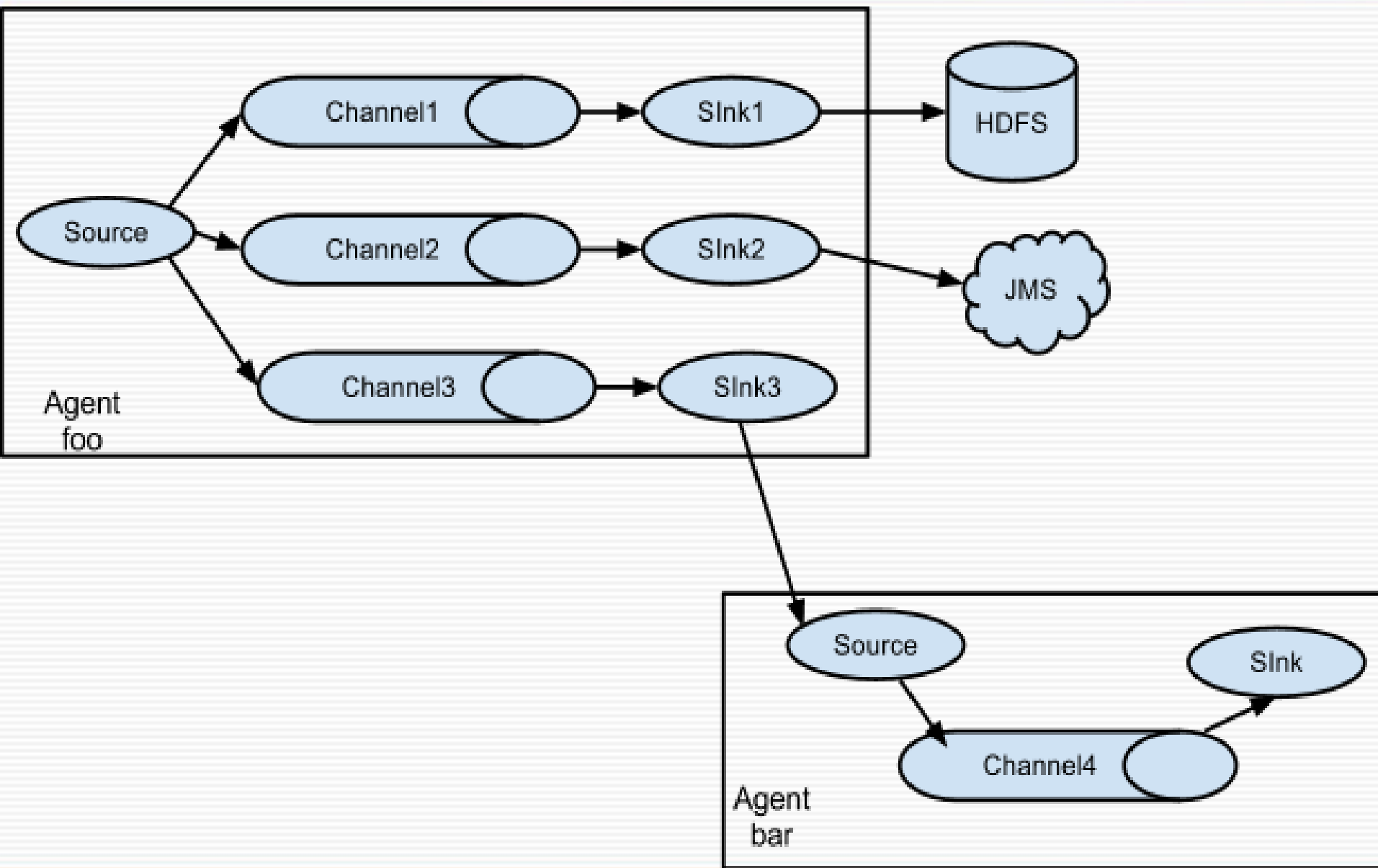
03

Sink

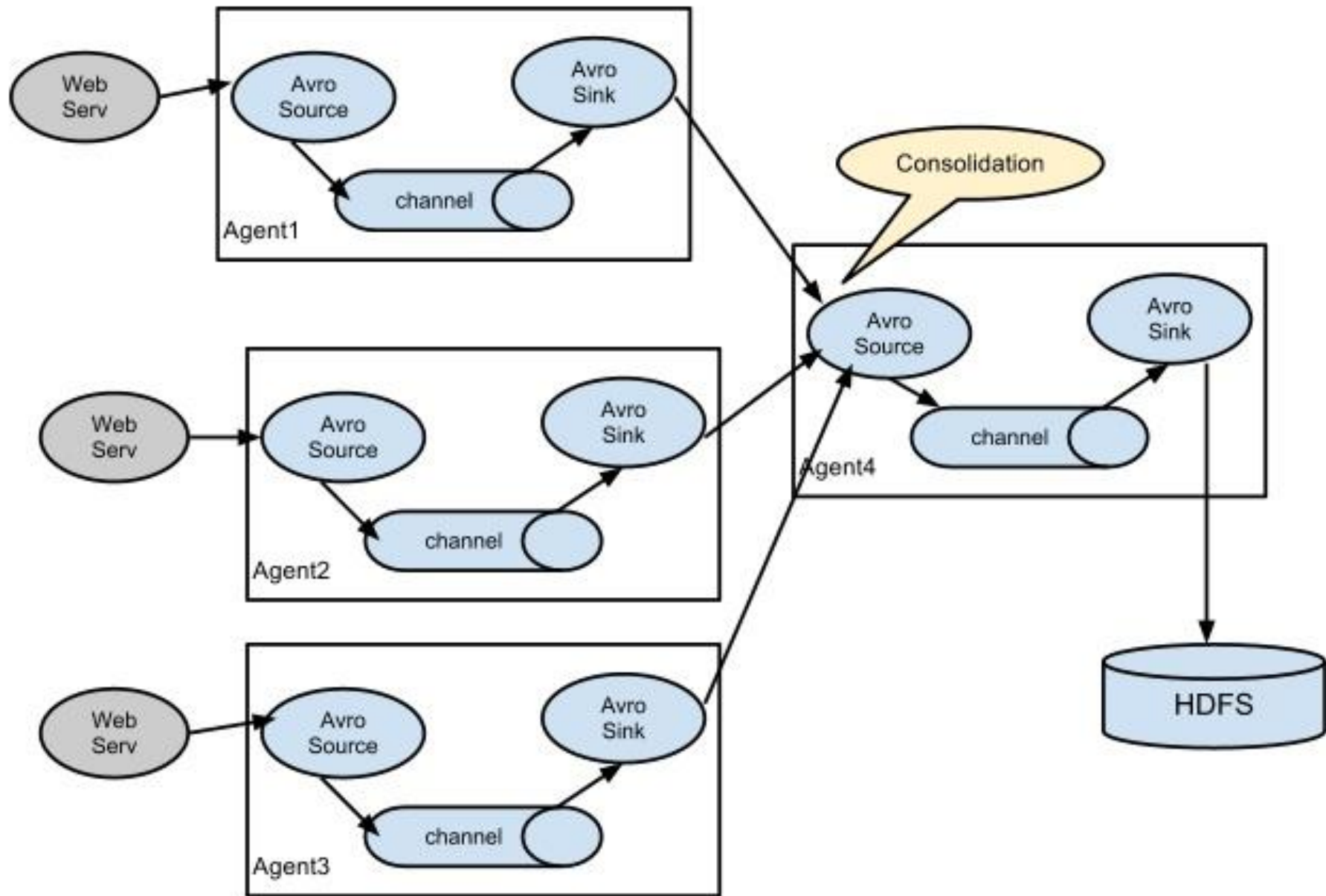
取数
kafkasink
hdfssink



Flume架构



Flume架构





高可靠



可扩展性



高效



可维护性

ExecSource:以运行Linux命令的方式，持续的输出最新的数据,实时传输，如**tail -F** 文件名指令

```
agent.sources.gather.type =
com.sohu.flume.source.ExecSource
agent.sources.gather.command = tail -F
/opt/countslog/countinfo.log
agent.sources.gather.channels = sendout
```

数据存储在内存队列中
速度最快

```
agent.channels.sendout.type =  
memory
```

```
agent.channels.sendout.capacity  
= 2000000
```

```
agent.channels.sendout.transacti  
onCapacity = 10000
```

capacity: Channel可以容纳events数量的最大容量, 就是queue的大小

transactionCapacity每次事务处理的events数量, 就是putList和takeList的容量大小

从channel拿数据，发送到kafka

```
agent.sinks.pool.type =  
com.sohu.group.flume.sink.KafkaSink  
agent.sinks.pool.channel = Mychannel  
agent.sinks.pool.topic = countinfo  
agent.sinks.pool.broker =  
10.16.10.196:8092  
agent.sinks.pool.zookeeper =  
10.16.10.76:2181/kafka-0.8.1  
agent.sinks.pool.num_partitions = 72
```

总体架构和F l u m e 简介

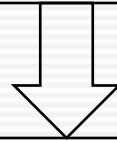
Memorychannel 的通信机制和事务

Kafkachannel 机制和调优

F l u m e 的缺点和未来努力

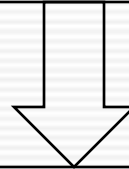
ExecSource:

```
flushEventBatch(List<Event> eventList)
```



ChannelProcessor (事务):

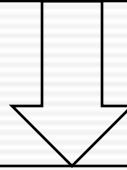
```
processEventBatch(List<Event> events)
```



Channel:

```
doPut(Event event)
```

Sink :
process ()



Channel:
doTake ()

doPut

把event写入数据结构putList

doTake

从queue得到event，放入takeList，返回event，然后发送到kafka

takeList的作用仅仅是commit失败时进行rollback!

Put事务提交

Event从putList放到queue中，
queue容量是capacity
然后清空putList

Take事务提交

清空takeList

Put事务回滚

清空putList，下次事务重新put

Take事务回滚

takeList回吐到queue

总体架构和Flume简介

Memorychannel的通信机制和事务

Kafkachannel机制和调优

Flume的缺点和未来努力

doPut

把event写入List<byte[]>
serializedEvents

doTake

没有使用

doCommit

把List serializedEvents发送
到kafka

发送成功才会事务Commit

doRollback

清空serializedEvents
下次事务再次发送

agent.channels.ctestfk.kafka.

batch.num.messages = 10000

设置为跟transactionCapacity一致

日志每秒条数	kafka.batch.num.messages	延时	每发送20万条耗时
3万	10000	无延时	7.1秒
3万	2000	延时15分钟	12.1秒
3万	不设置(200)	延时1小时多	34-54秒

默认是20

```
agent.sources.s.batchSize = 10000
```

```
agent.sinks.k.hdfs.batchSize =  
10000
```

batchSize一定不能大于
transactionCapacity

原生的是随机写入

增加了负载均衡功能，轮询
写入kafka

自带序列化方式avro, 有乱码问题

目前可通过配置序列化方式来解决乱码问题

Fai lover



以插件的形式开发

后台监控进程

自动进行**channel**切换

通过配置可启动或关闭**failover**

监控时间间隔可配置

日志每秒条数	channel	延时	每发送20万条耗时
10万	memorychannel	报错后无响应	8.6秒
10万	kafkachannel	9到25分钟	8.2秒
5万	memorychannel	报错后无响应	8.1秒
3万	memorychannel	无延时	7.1秒
3万	kafkachannel	无延时	7秒

总体架构和F lume简介

Memorychannel 的通信机制和事务

Kafkachannel 机制和调优

F lume的缺点和未来努力

◆ OOM 问题

Flume 启动时的最大堆内存大小默认是 20M，线上环境很容易 OOM，因此需要在 `flume-env.sh` 中添加 JVM 启动参数：

```
JAVA_OPTS=“-Xms8192m -Xmx8192m ...”
```


◆数据重复写入

回滚后重复写入

可能的解决方法：

1

关闭批量写入

2

消费端进行去重

◆数据丢失:

MemoryChannel

Flume进程挂掉

补救:

hdfs put来补数

预防:

监控: 进程和日志监控

◆ 日志量比原来大

1

flume.log 日志量约60MB/天

2

监控进程调kafka api产生日志，可做日志滚动定期删除。

1

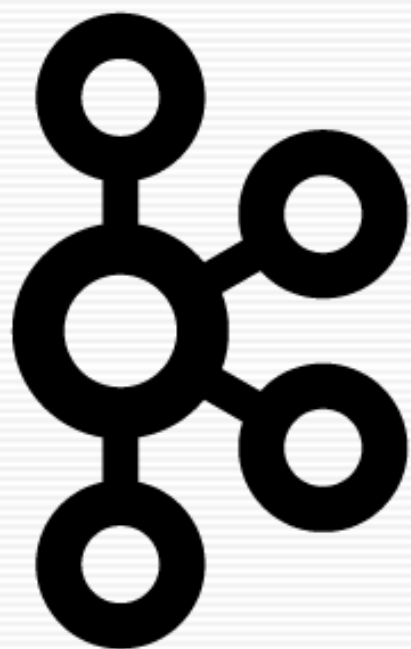
高效稳定的**Source**

2

Channel参数调优，**avro**序列化

3

Kafkachannel+其他**sink**



kafka

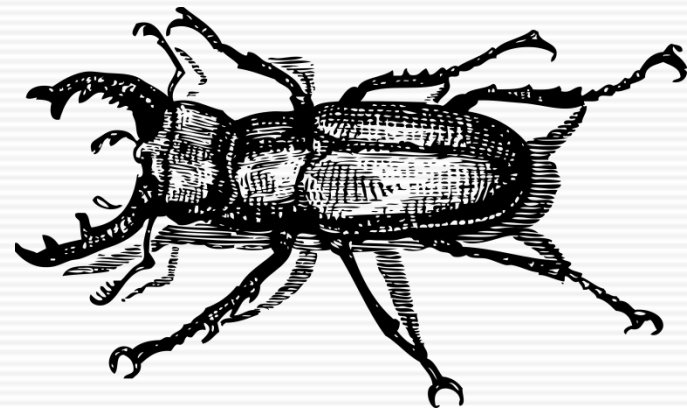
Kafka介绍

架构设计实现

应用现状

未来规划

- LinkedIn 开发
- 发布-订阅消息系统



可扩展性

高吞吐

消息持久化

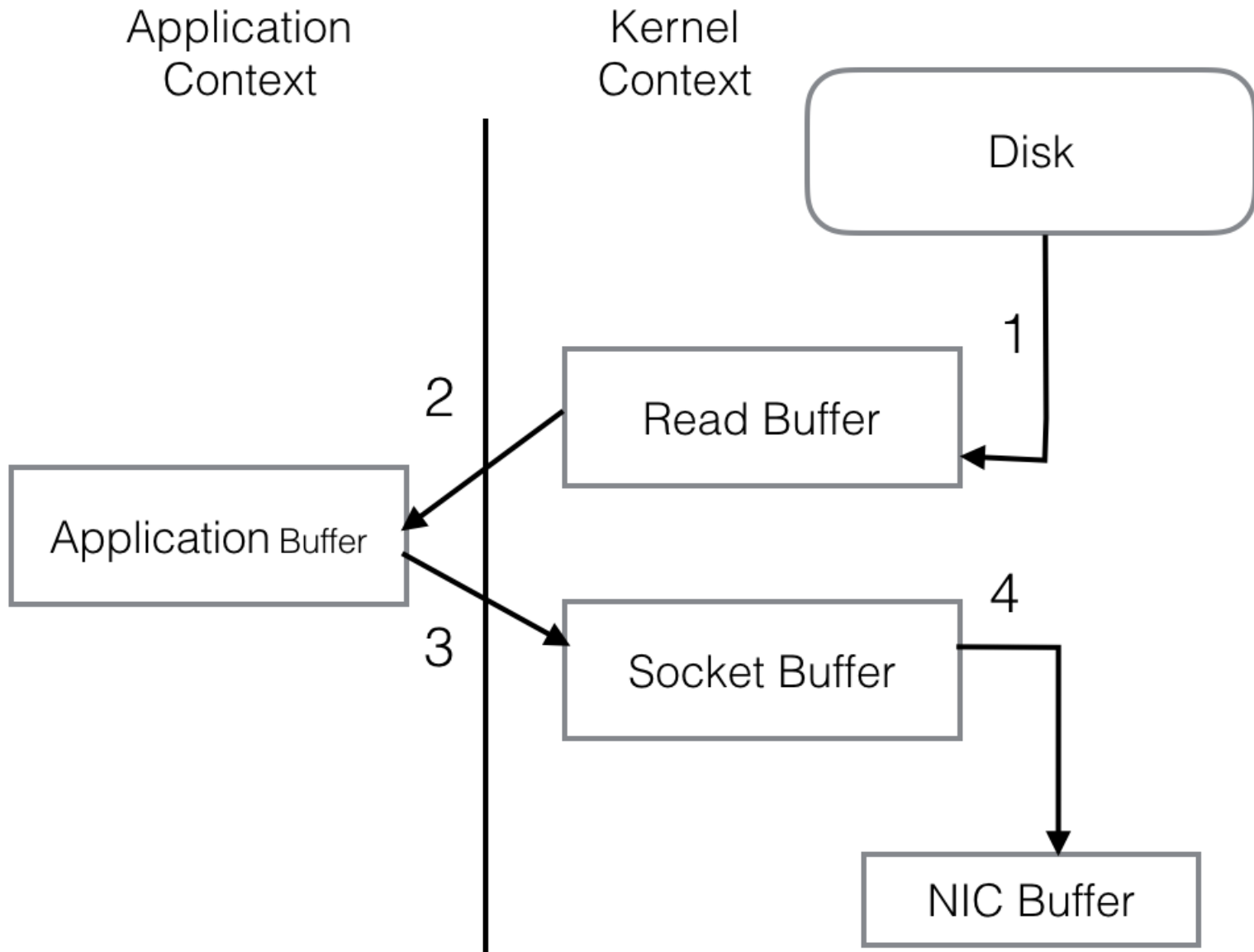
分布式

Kafka介绍

架构设计实现

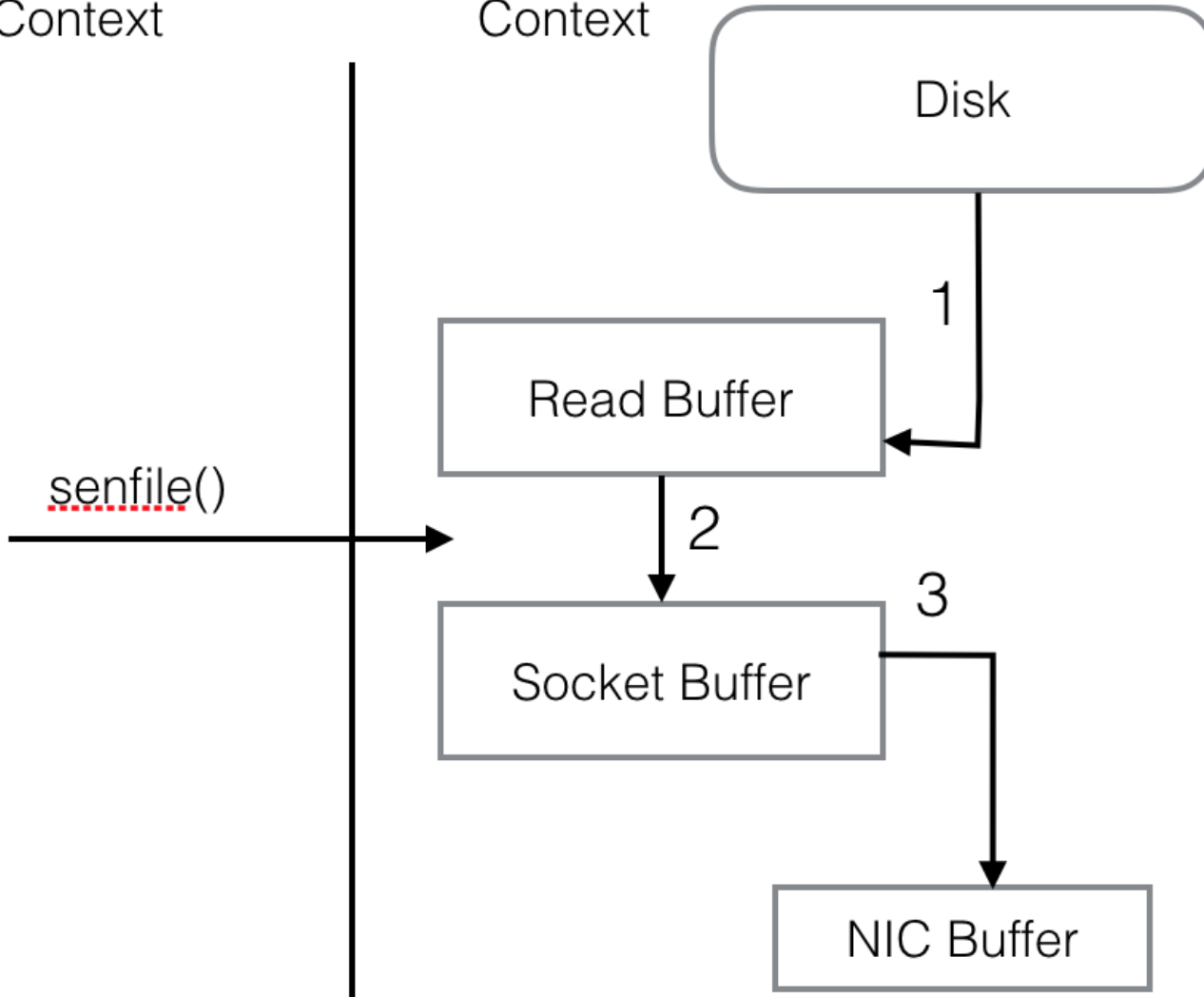
应用现状

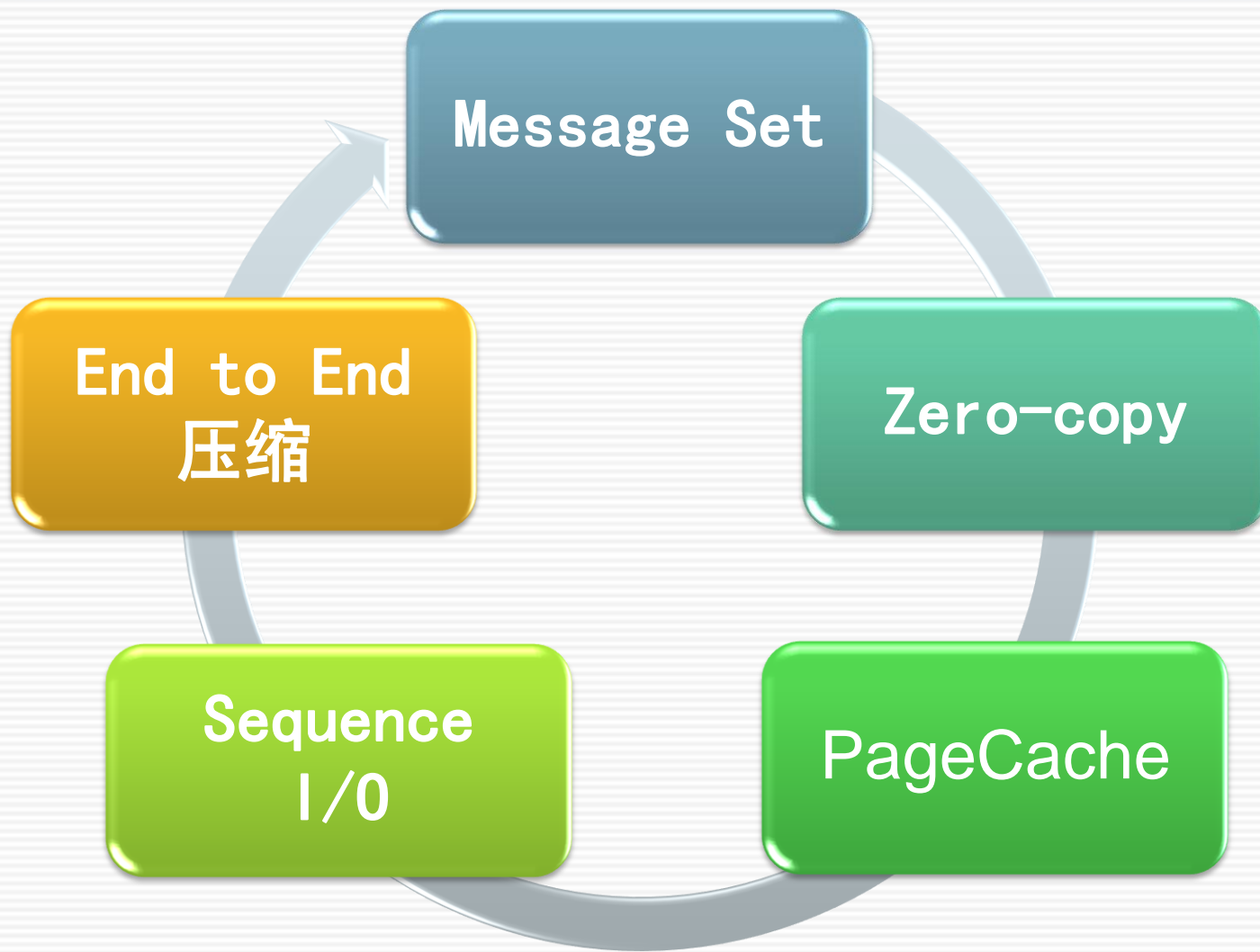
未来规划



Application Context

Kernel Context



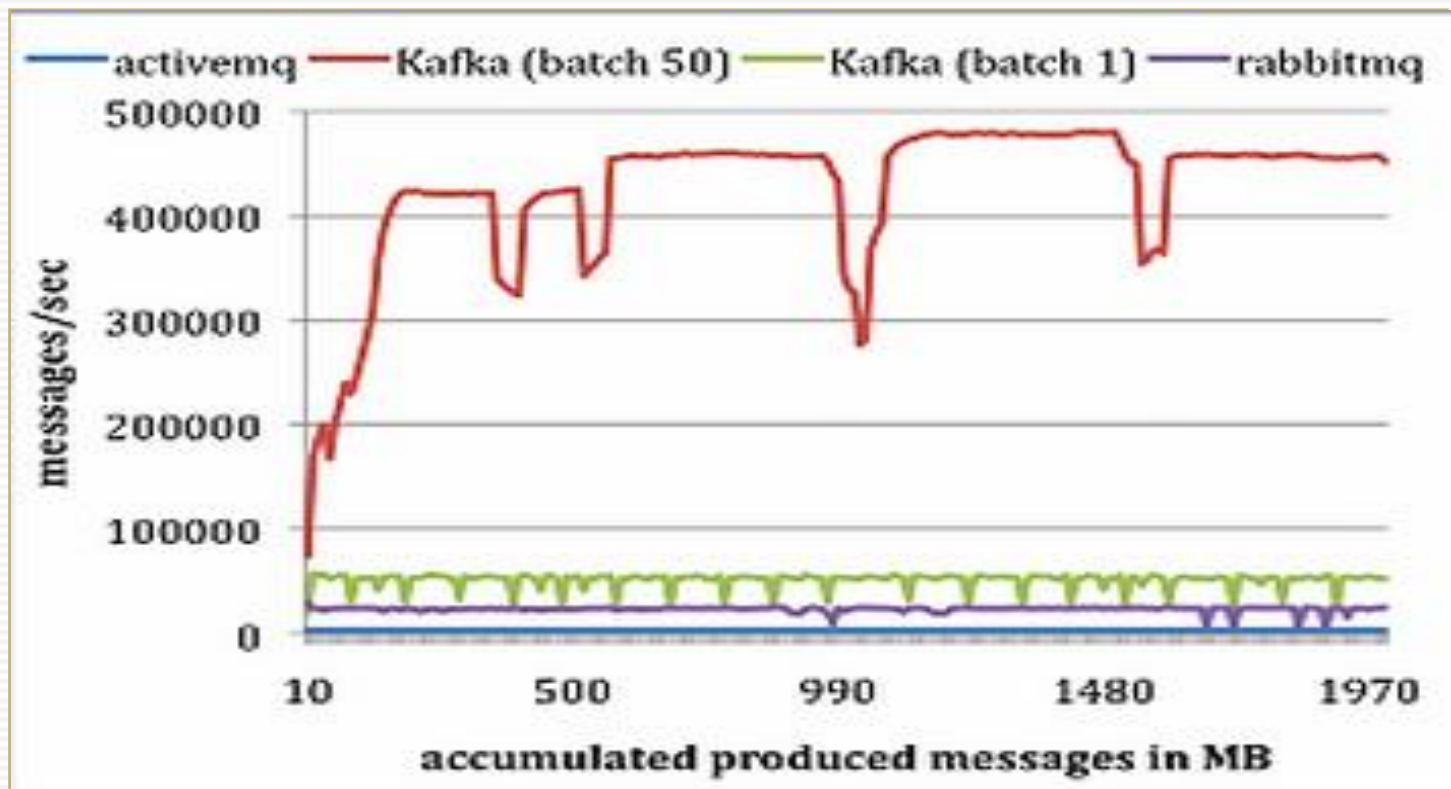


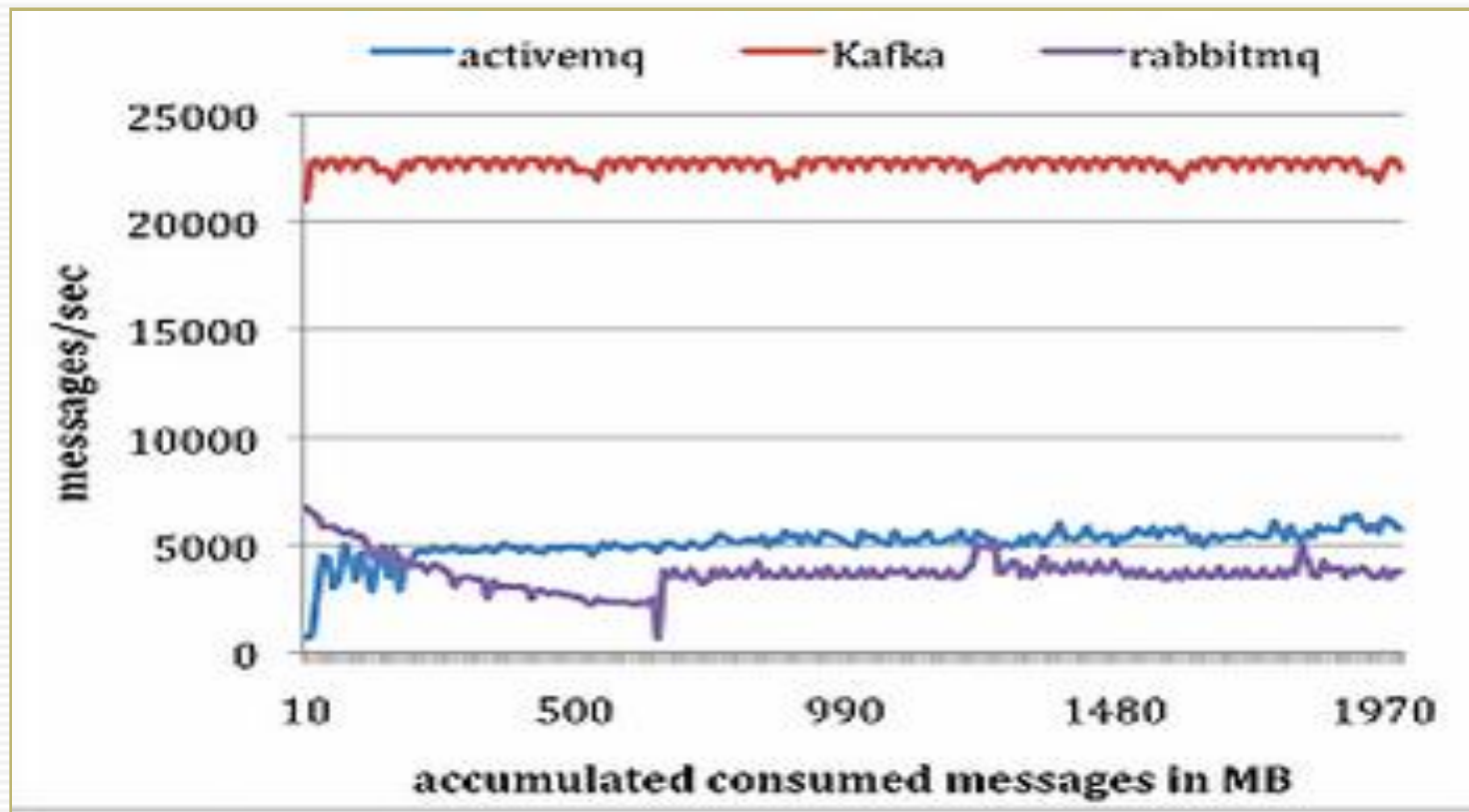
日志
收集

消息
队列

流式
处理

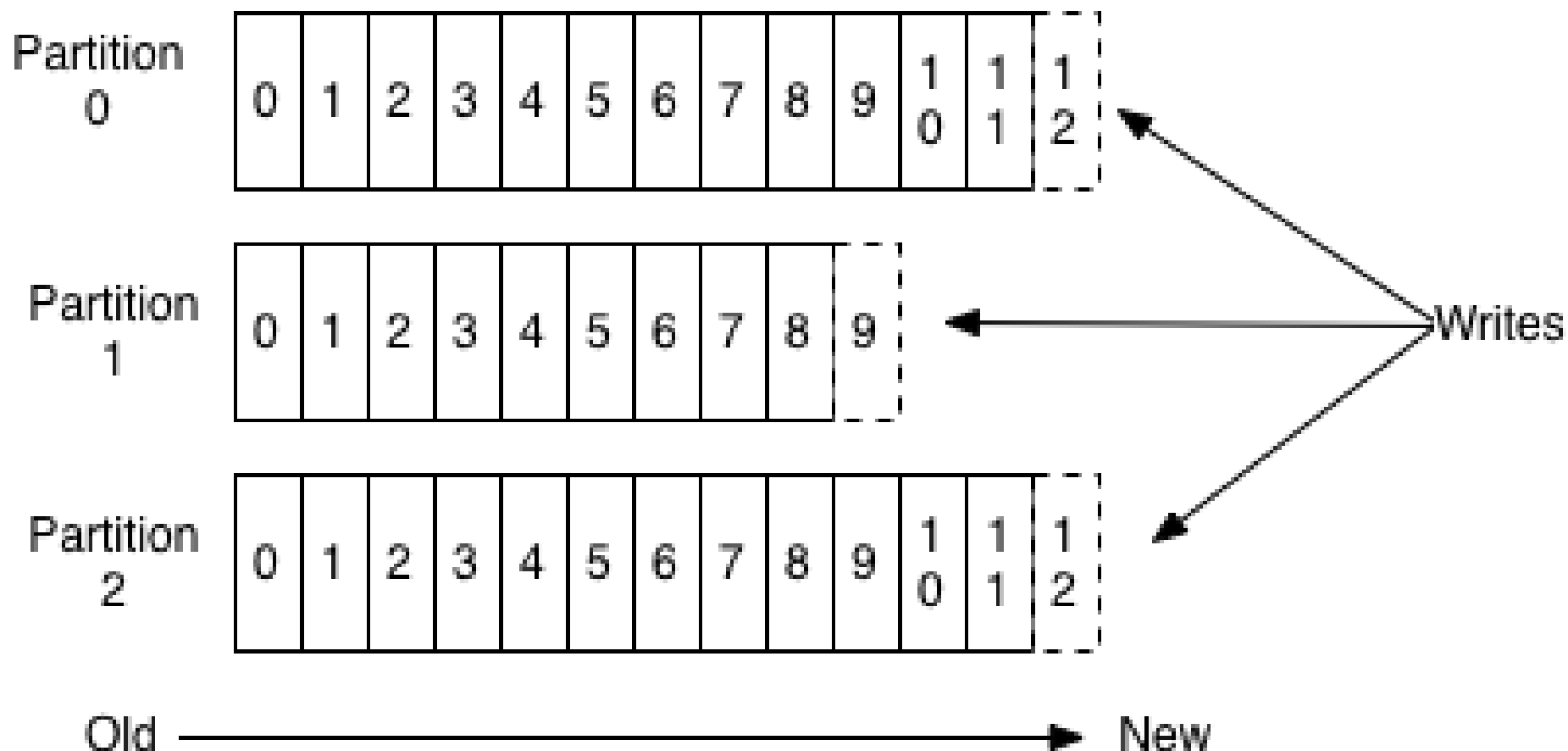
持久
化日
志

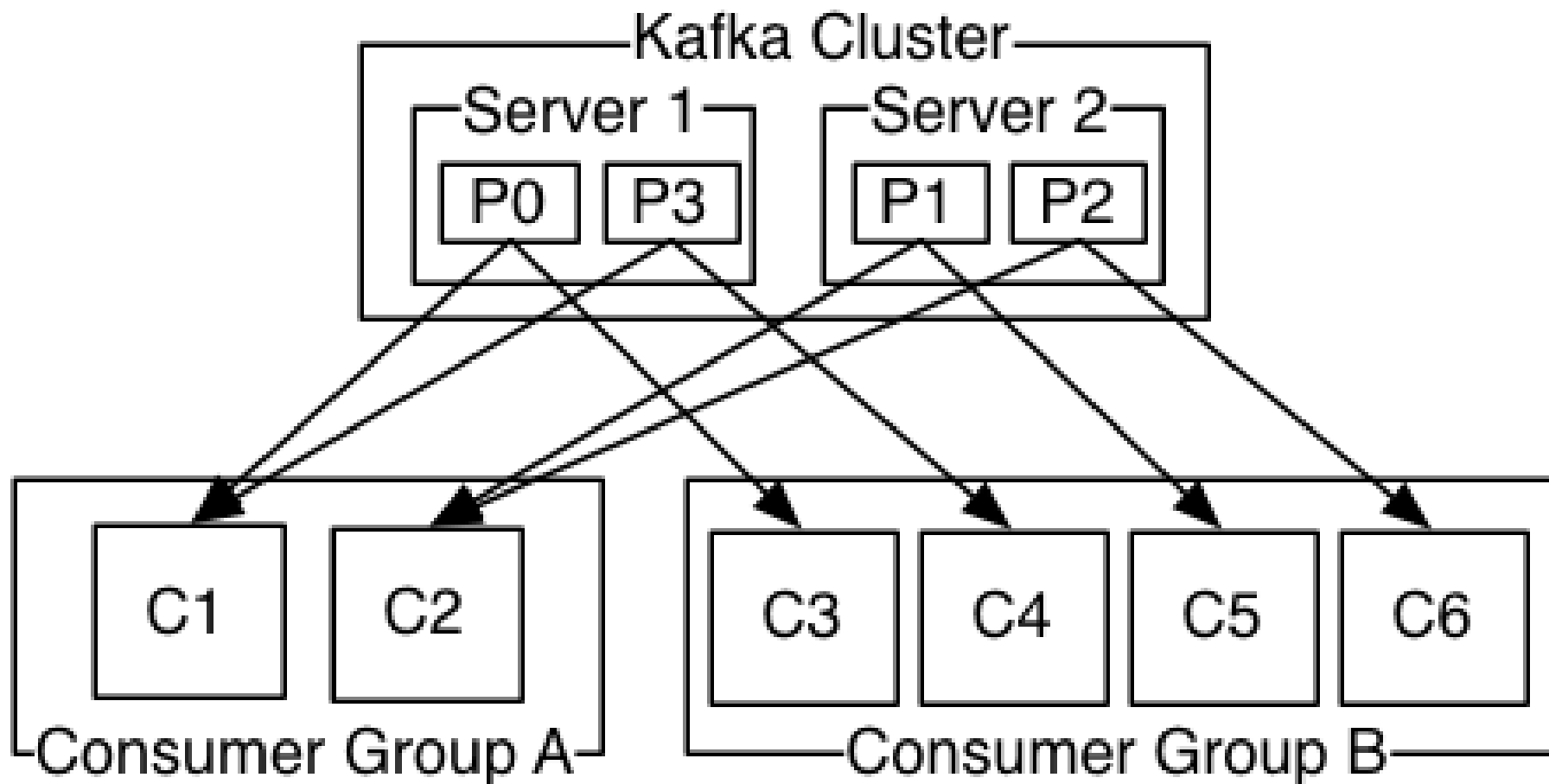


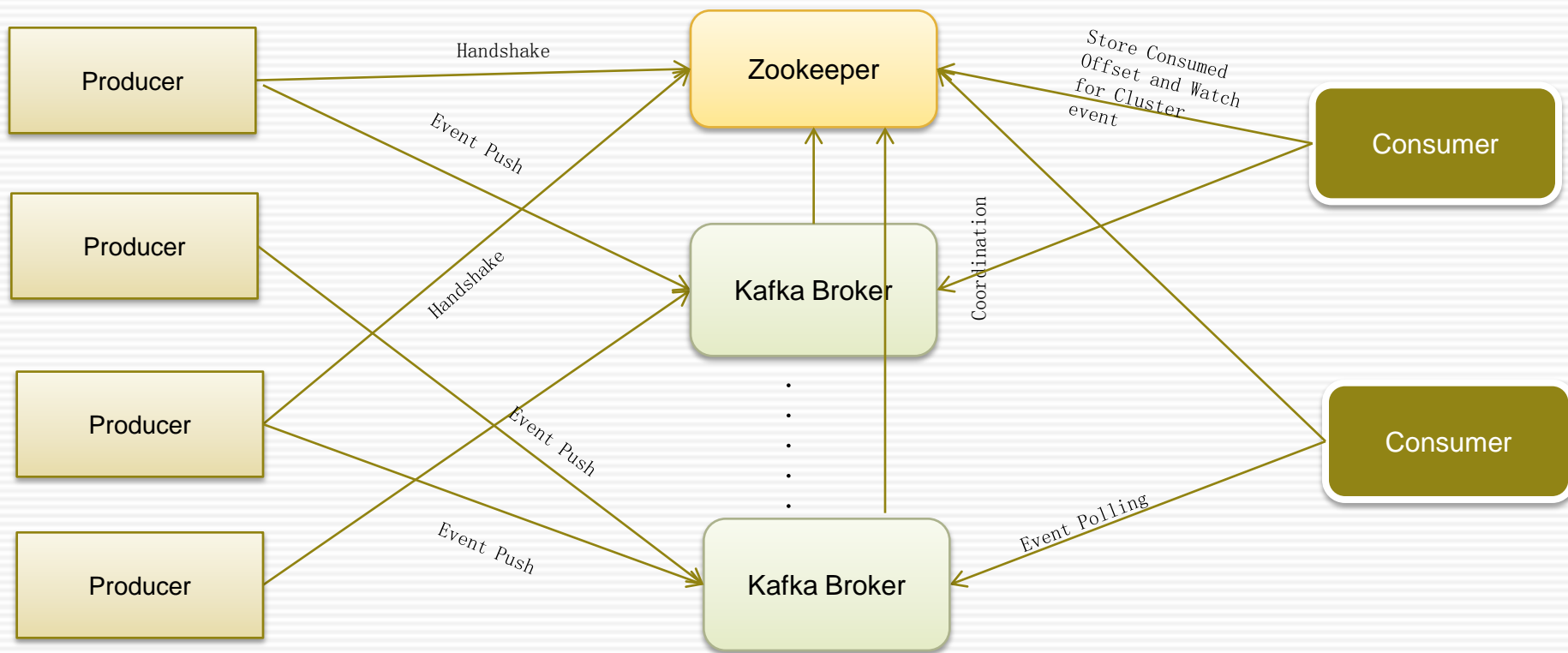


- **Broker** Kafka集群中的一个节点服务
- **Topic** 消息源的分类
- **Partition** 基于topic的分组。实现topic负载均衡的方法之一
- **Offset** 可以简单理解为消息在文件中的位置或者累计的量
- **Message** 通信的基本单位

Anatomy of a Topic







- 基于Partition主从设计
- 1 Leader; (n-1) follower
- 主分区可读写、从分区备份

Higher Level

- 自动维护**offset**
- 开发简单,省心

Lower Level

- 消费者单独维护**offset**
- 扩展灵活,开发复杂

Kafka介绍

架构模型

应用现状

未来规划

- Kafka 18台+6台
- Zookeeper 7台

- 12 核、24G、15T、千兆网卡
- kafka为0.8.1版本
- Zk为3.4.5

- Available Topic **46**
- Total Topic **81**个
- 入口带宽 30M左右
- 出口带宽 80M左右
- 磁盘空间使用 7T
- 备份3天，副本3

- kafka2hdfs 44个
- spark-streaming 8个
- 其他部门

集群规模

集群配置

现状

消费者

1

Partition分配算法简单

2

Topic无法删除

3

未提供安全机制

4

没有可视化工具

脚本监控

- 基于zabbix监控

可视化监
控工具

- kafka-web-console

addr-10.16.43.151: Outgoing traffic on interface eth2

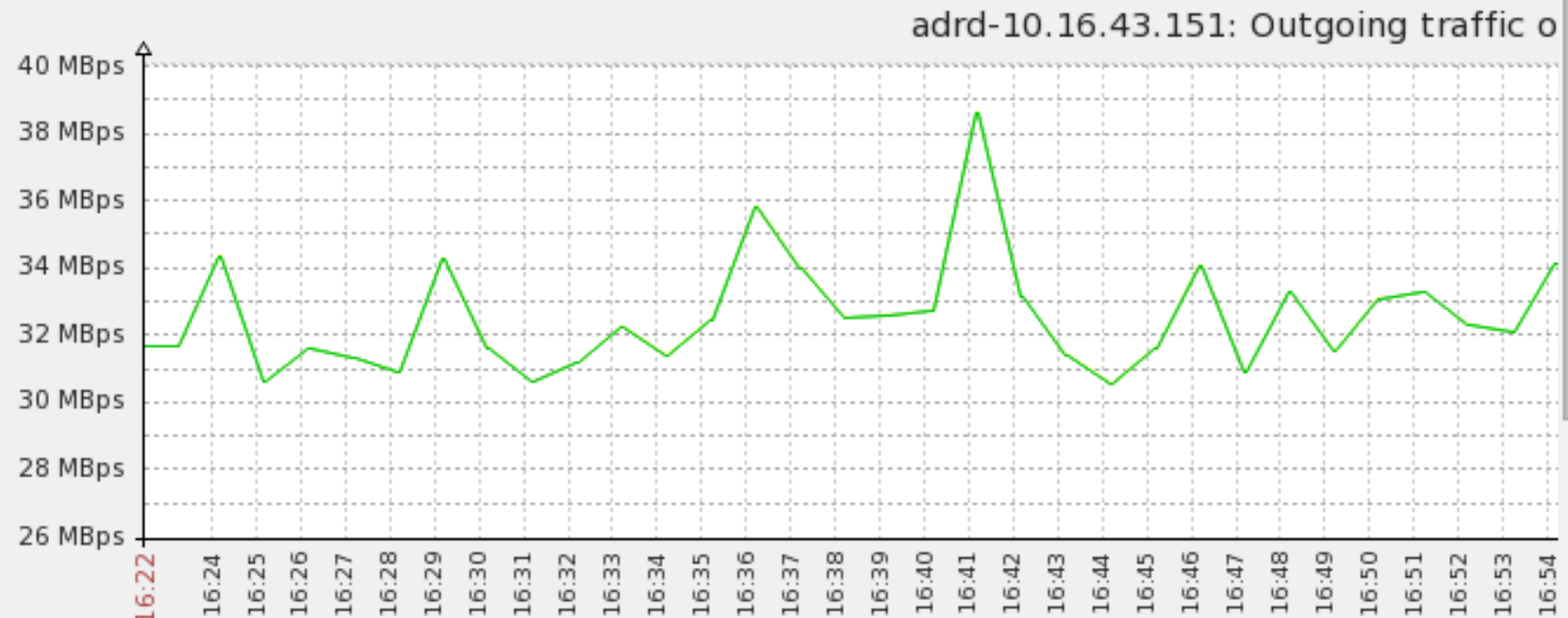
Graph

Hide filter

Zoom: [1h](#) [2h](#) [3h](#) [6h](#) [12h](#) [1d](#) [7d](#) [14d](#) [1m](#) [3m](#) [6m](#) [All](#)



Navigation controls: << [6m](#) [1m](#) [7d](#) [1d](#) [12h](#) [1h](#) | [1h](#) [12h](#) [1d](#) [7d](#) [1m](#) [6m](#) >>



Kafka Web Console

Zookeeper	Host	Port	MessagesIn	DataFlowIn	DataFlowOut
10.16.34.187	10.16.43.156	9997	600088051856	353435265670835	16367384891895
10.16.34.187	10.16.43.154	9997	457310381415	288479776774169	10890692693081
10.16.34.187	10.16.43.155	9997	551196801136	314870018675549	12743944346353
10.16.34.187	10.16.43.152	9997	509670525665	296860862562628	10739156515256
10.16.34.187	10.16.43.153	9997	425038570469	274962794689447	10047326830863
10.16.34.187	10.16.43.150	9997	476325154731	308158746674060	11460771913911

• 技术方案

Scala play、
zookeeper、
jmx、sbt

• 功能

Brokers、
Topic的监控
监控项配置
告警等功能

Kafka介绍

架构模型

应用现状

未来规划

0.9+

SSL或SASL进行验证

native offset storage

新的ConsumerAPI、NIO

Ip链接数限制

读写授权管理



Thank You !