



第12期

《新一代互联网行为定向广告技术的挑
战与优化-》 -

品友互动专场

www.LAMPER.cn

QQ群：83304912

<http://weibo.com/lampercn>

Hadoop的ETL任务

—Flume使用及其优化

汪浩

目录

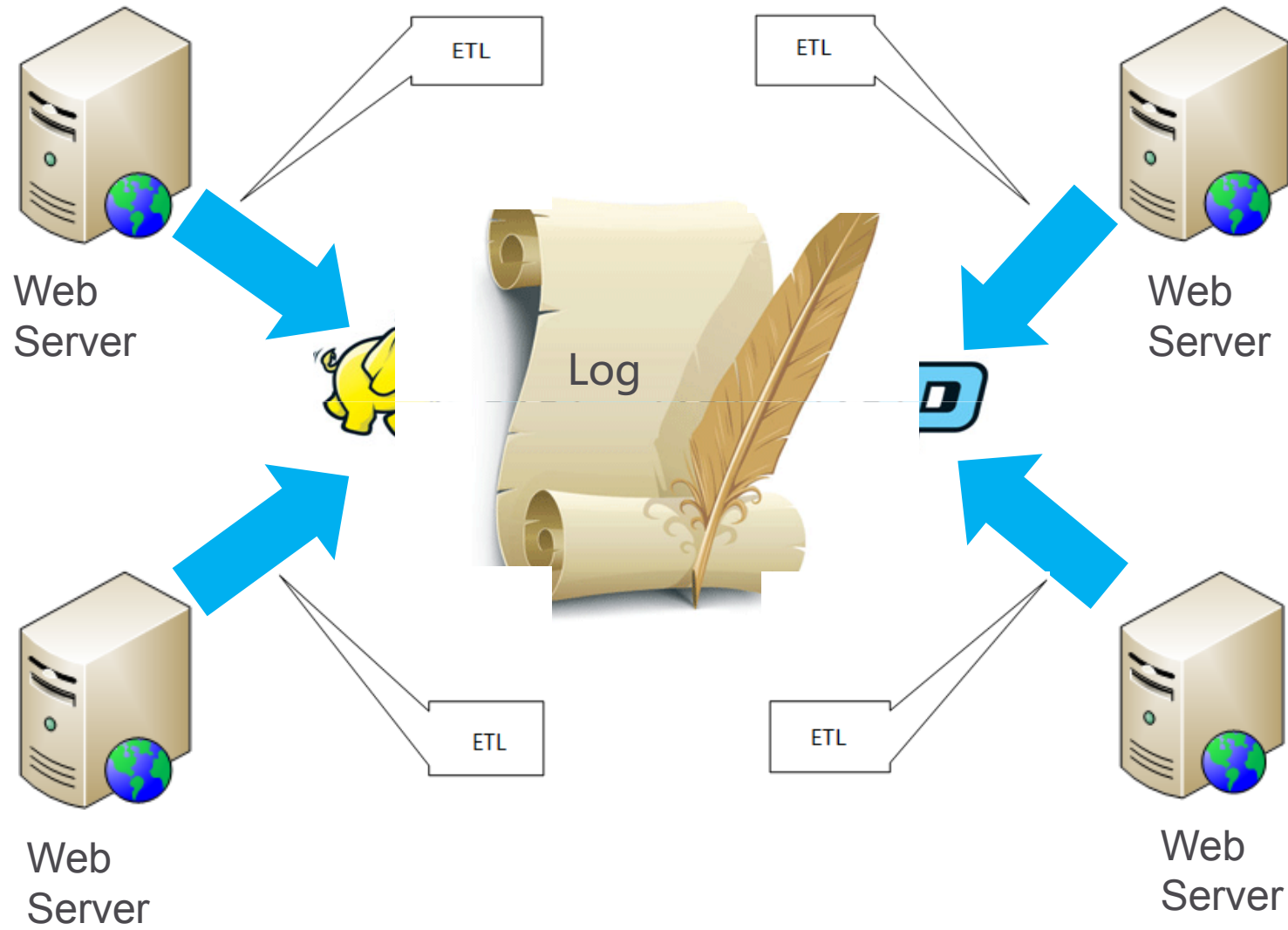
背景介绍

日志收集系统介绍

日志收集系统比较-Why Flume ?

Flume使用心得和优化

背景介绍





日志收集系统介绍

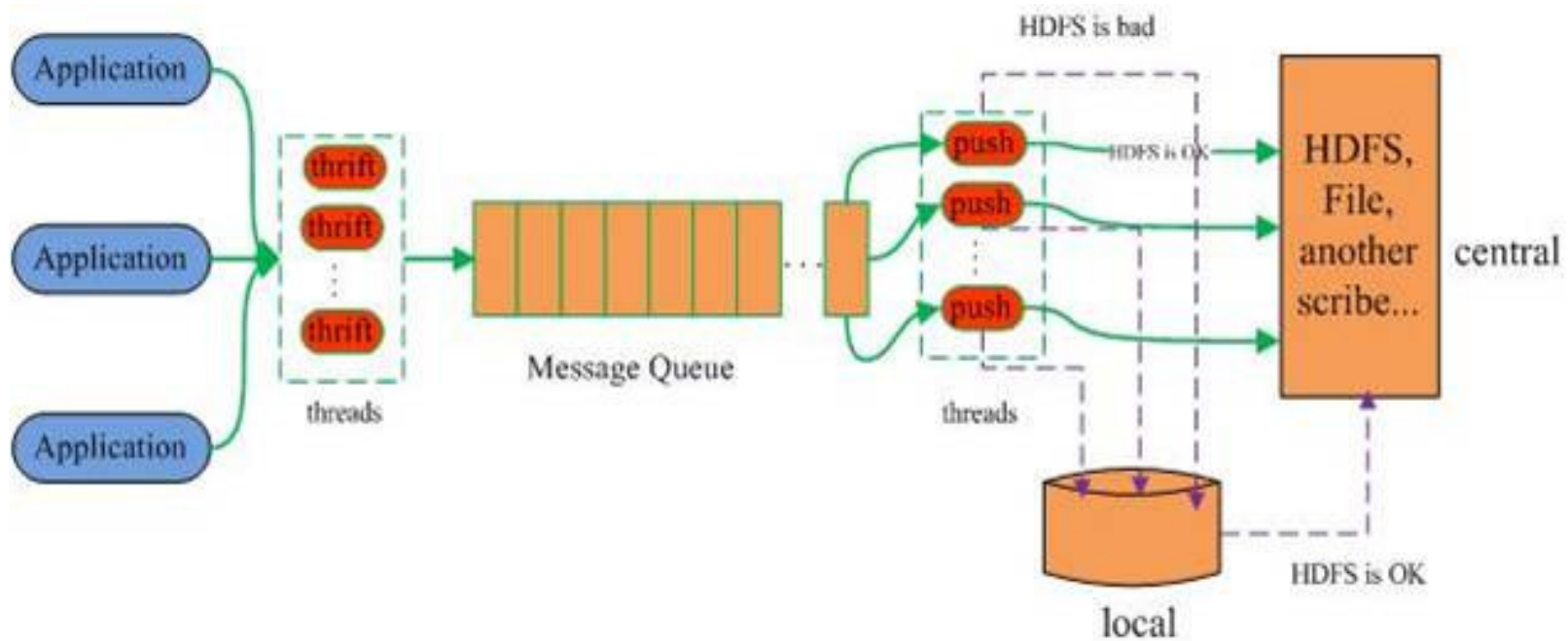


FLUME

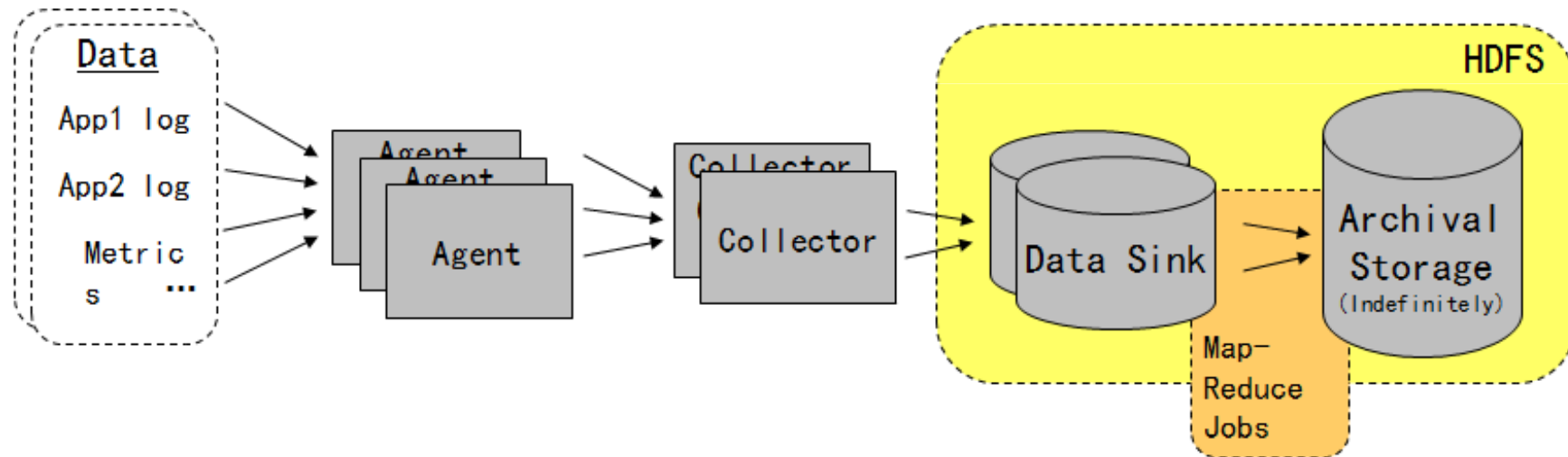


日志收集系统介绍——Scribe

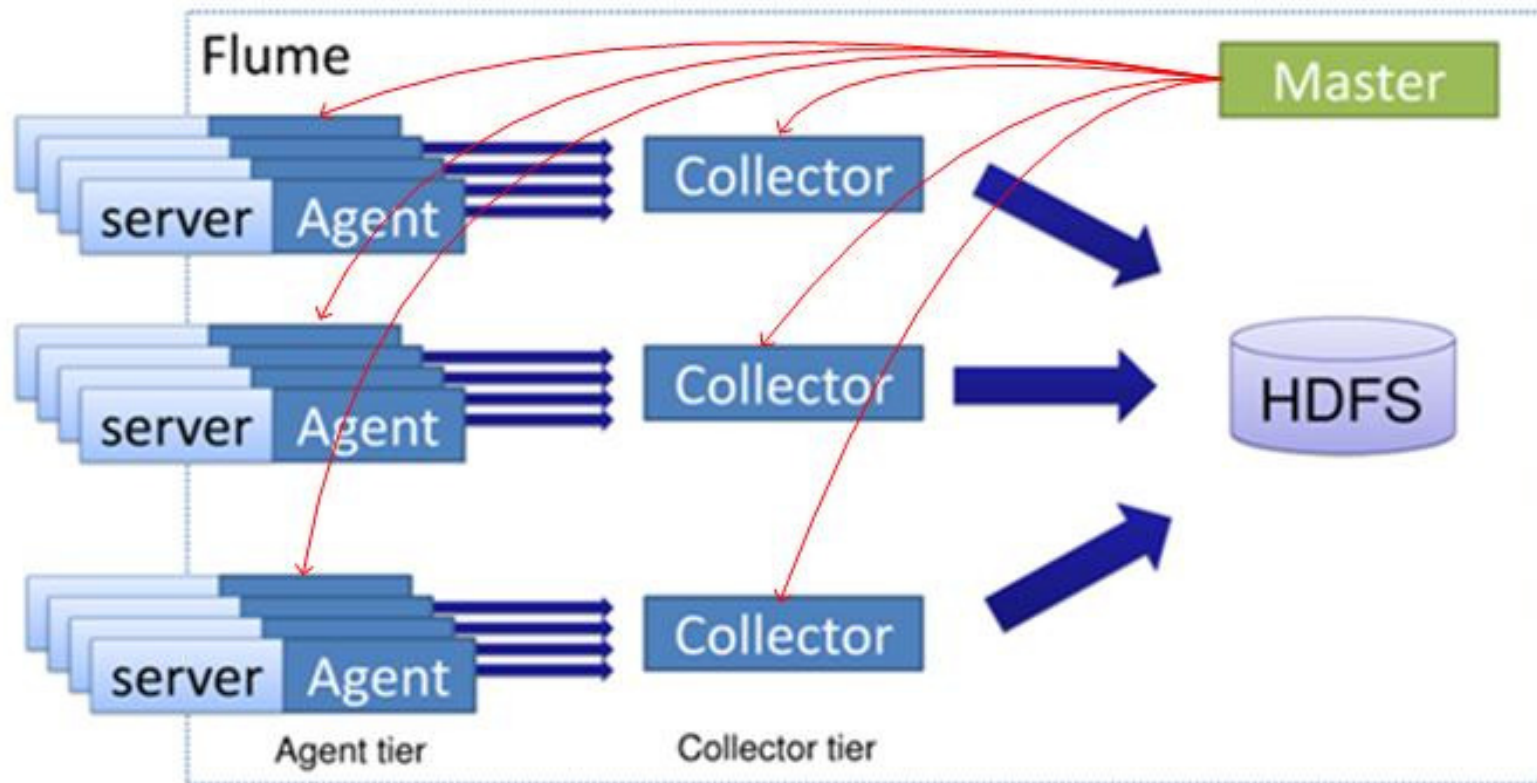
facebook
Scribe



日志收集系统介绍——Chukwa



日志收集系统介绍——Flume



日志收集系统介绍——Flume

- Flume基本概念
 - 数据路径
 - Nodes在数据路径上
 - Nodes上存在Source和Sink
 - Nodes可以设置为不同的角色
 - 控制路径
 - 心跳检测
 - 指定Sources和Sinks
 - 控制Nodes间的数据流

Agent

Collector

Master

日志收集系统介绍——Flume

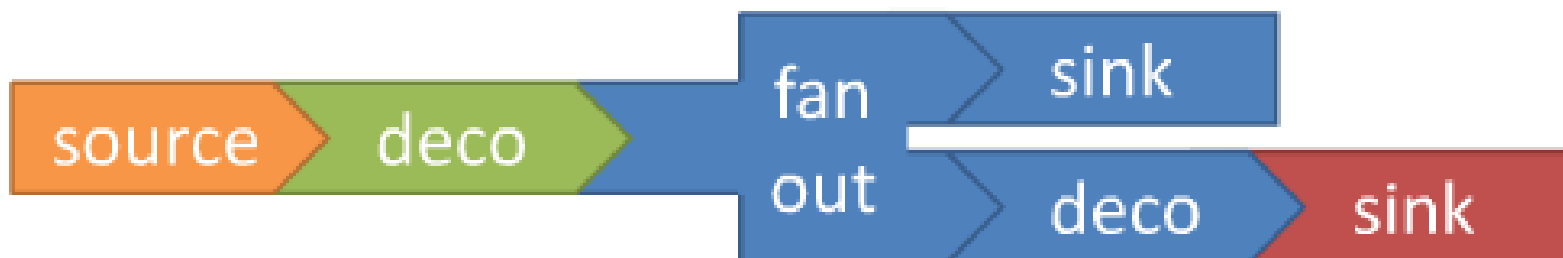


日志收集系统介绍——Flume

- Flume可扩展性

1. Flume易于扩展的原因

- 简单的Source和Sink APIs
- 基于事件流的设计易于把简单的操作组合成复杂的操作
- 插件式的架构易于用户使用自己编写的Sources、Sinks和Decorators



日志收集系统介绍——Flume

- Flume可扩展性

2. Connector的种类

-Sources (产生数据)

source

Console、Exec、Syslog、Scribe、IRC、Twitter

-Sinks (发送数据)

sink

Console、Local files、HDFS、S3

-Decorators

deco

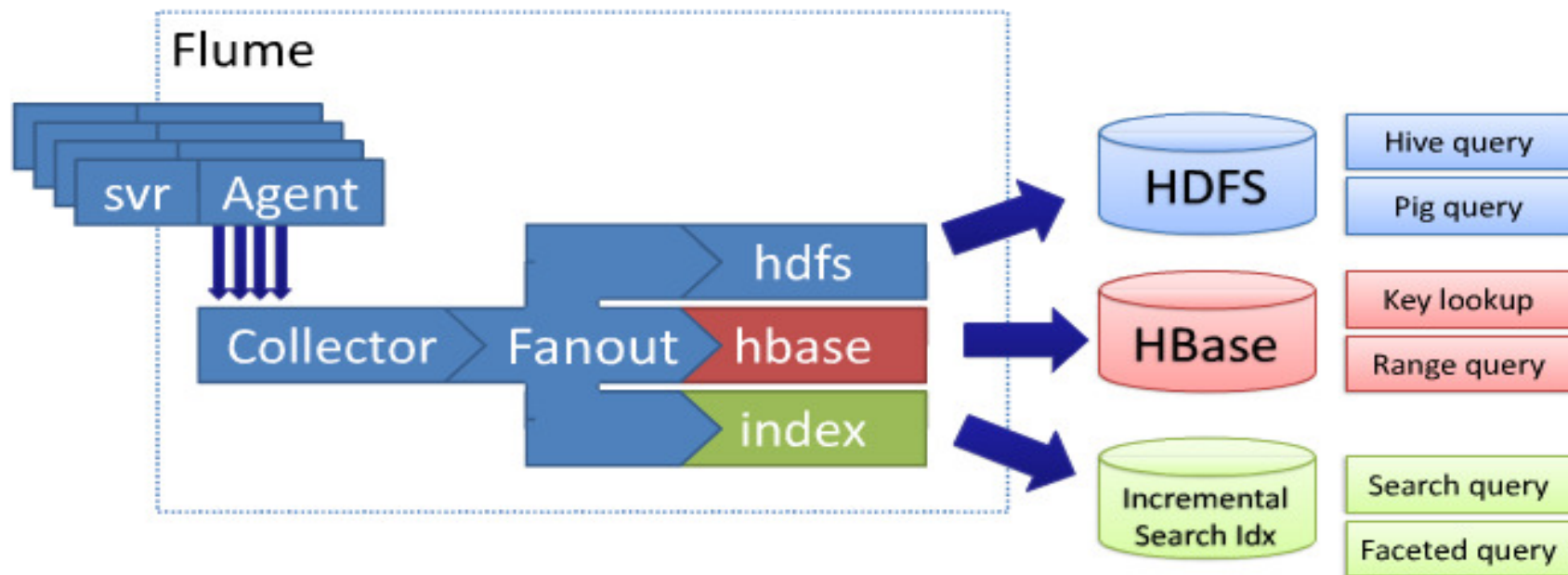
(在数据发送给Sinks之前对数据进行处理)

Wire batching、compression、sampling、throughput
throttling

日志收集系统介绍——Flume

- Flume可扩展性

3. 示例



日志收集系统介绍——Flume

- Flume可靠性
三种级别的故障恢复模式

1. agentBESink



日志收集系统介绍——Flume

- Flume可靠性
三种级别的故障恢复模式

2. agentDFOSink



日志收集系统介绍——Flume

- Flume可靠性
三种级别的故障恢复模式

3. agentE2ESink

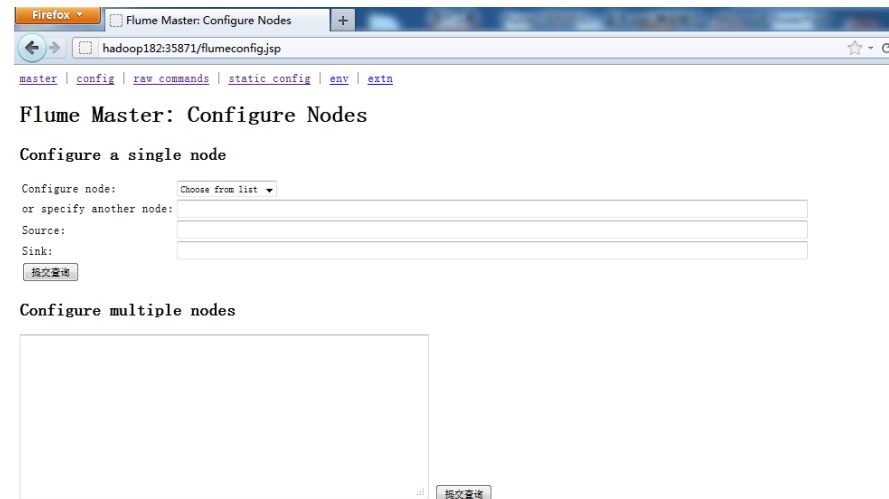


日志收集系统介绍——Flume

- Flume可管理性

1. 管理接口

-Web Page



-Flume Shell

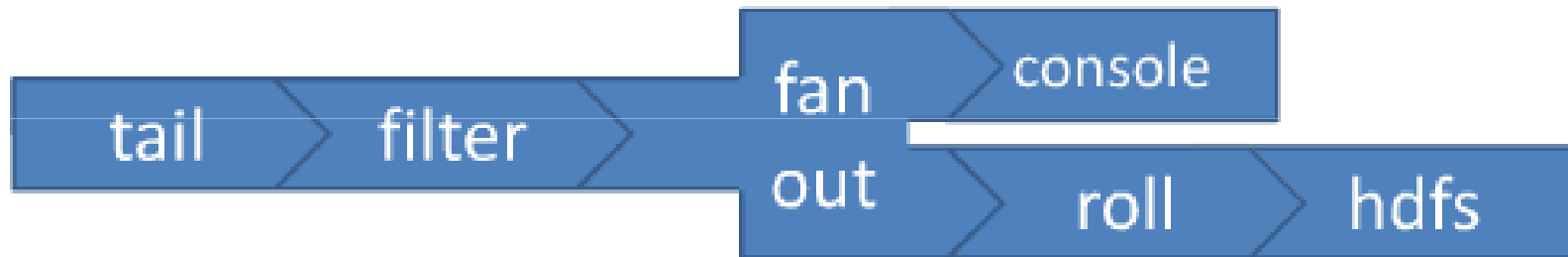
```
-bash-3.2$ flume shell
2012-01-31 15:50:00,932 [main] INFO conf.FlumeConfiguration: Loading configurations from /etc/flume/conf
=====
FlumeShell v0.9.4-cdh3u1
Copyright (c) Cloudera 2010, All Rights Reserved
=====
Type a command to execute (hint: many commands
only work when you are connected to a master node)

You may connect to a master node by typing:
  connect host[:adminport=35873[:reportport=45678]]

[flume (disconnected)] help
```

日志收集系统介绍——Flume

- Flume可管理性
2.示例



Configuring FlumeNode:

```
tail( "file" ) | filter [ console, roll(1000)  
{ dfs( "hdfs://namenode/user/flume" ) } ] ;
```

日志收集系统介绍——Flume



IPINYOU
品友互动



日志收集系统比较- Why Flume?

可靠性	
Scribe	Scribe Server和中央Scribe Server之间、中央Scribe Server和Store之间都有容错机制。但是Scribe server发生故障时，内存中少量的数据会丢失，磁盘上的数据不会丢失。
Chukwa	Agent定期记录已发送给Collector的数据偏移量，一旦出现故障后，可根据偏移量继续发送数据。
Flume	Agent和Collector，Collector和Store之间均有容错机制，且提供3种级别的可靠性保证。

日志收集系统比较- Why Flume?

可扩展性	
Scribe	Agent是一个Thrift Client，新的功能需要自己实现；Collector也就是一个Thrift server。
Chukwa	本身自带少量的Agents和Collectors。
Flume	本身提供丰富的Agents、Collector和Decorators，插件式的结构十分方便扩展。

日志收集系统比较- Why Flume?

可管理性	
Scribe	修改配置文件。
Chukwa	修改配置文件。
Flume	通过Web或者Flume Shell。

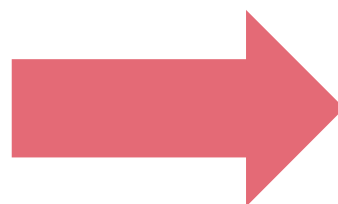
日志收集系统比较- Why Flume?

综合考虑以下因素：

良好的可靠性

良好的可扩展性

良好的可管理性



Flume使用心得及优化

1.Master单点故障

2.Collector负载均衡

3.HDFS Small Files

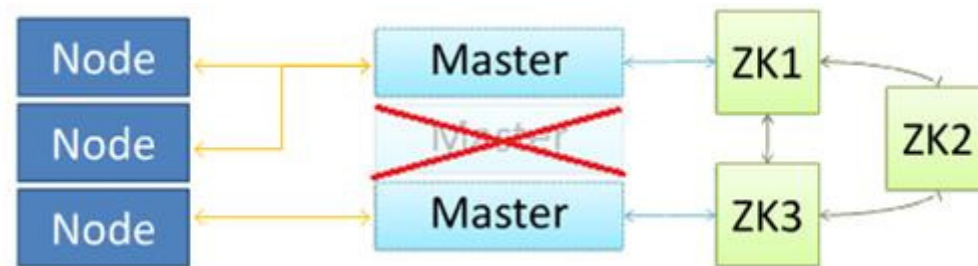
4.数据传输速率

5.CPU使用率

6.内存使用率

Flume使用心得及优化

- Master单点故障



Flume使用心得及优化

- Collector负载均衡

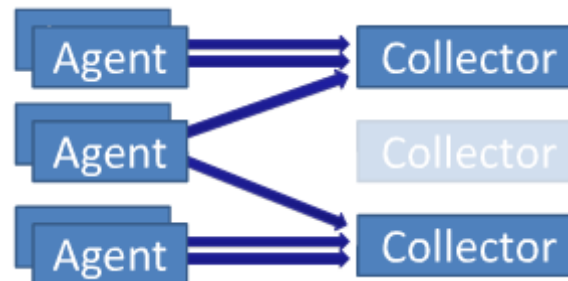
1. Agents的数据可以在逻辑上划分，发往不同的Collectors



Flume使用心得及优化

- Collector负载均衡

*2.通过预先设置自动failover达到分流
当有Collector无法工作时
当加入新的Collector时*



Flume使用心得及优化

- HDFS small files

日志传输到HDFS上，存在大量的小文件

- NameNode造成压力

- 产生大量的Map任务

```
Namenode (Filename, numReplicas, block-ids, ...)  
/users/sameerp/data/part-0, r:2, {1,3}, ...  
/users/sameerp/data/part-1, r:3, {2,4,5}, ...
```

Flume使用心得及优化

- HDFS small files

优化方案：

- 使用CollectorSink 时设置rollmillis参数
- 在flume-site.xml中配置flume.collector.roll.millis

Flume使用心得及优化

- 数据传输速率

1. Batch

-使用batch(n,maxlatency) Decorator对Event进行批处理，提高系统吞吐量和资源的利用率



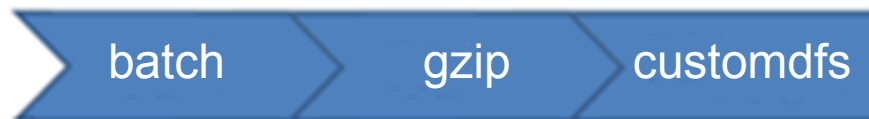
- 在Agent上使用batch时，也要相应的在Collector上使用unbatch，将集成的Event分解成原始的单个Event

Flume使用心得及优化

- 数据传输速率

2.Compression

- 使用gzip Decorator对Event进行压缩，降低网络传输的数据量，提高数据传输速率。
- 一般结合batch一起使用。



- 在Agent端使用gzip时，在Collector也要相应的使用gunzip，对Event进行解压。
- 我们使用gzip可以减少80%的数据量。

Flume使用心得及优化

- 数据传输速率

3. checksum

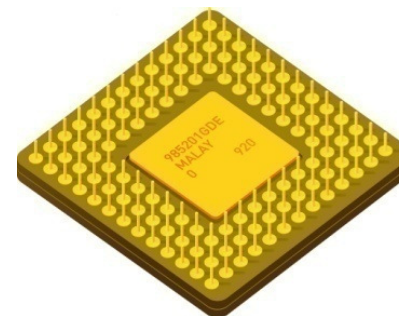
- Collector接受到数据后，对数据进行Checksum，验证数据的正确性。通过改进Checksum算法，缩短数据的验证的时间
- 我们通过改进Checksum，使数据验证的时间减少33%

Flume使用心得及优化

- CPU使用率

使用TailDirSource 时，发现消耗了大量的CPU

- 轮询文件的时间200ms
- 大量的文件消耗消耗很多的CPU



优化方案：

- 延长轮询文件的时间间隔
- 长时间内容没有改变的文件将不再轮询
- 减小Log所在的路径深度

优化后，使CPU的使用率降低了2/3

Flume使用心得及优化

- 内存使用率

使用TailDirSource消耗内存

-TailDirSource 里大量使用了Direct Buffer



优化方案：

调整MaxDirectMemorySize的大小，限制Direct Memory的使用

优化后，使内存的使用率降低了1/2

Flume在品友互动的使用 – 优化

与君共勉——

Citius, Altius, Fortius



我们的技术：

海量数据、云计算、分布式、数据挖掘、机器学习、
精准定向、用户行为分析

海量数据、云计算、分布式、数据挖掘、精准定向、
数据分析、Hadoop, Redis, Hbase, Hive, Pig,
Oozie, Ganglia, Flume, Lucene, LIBSVM,
Mahout, Zookeeper.....

尽管每一个词都名声显赫，热的发紫，但这的确就是
我们每天正在做的和使用着的；

我们的产品：Optimus  Folo8 

品友互动感谢您的关注，希望继续支持：

官方网站：<http://www.ipinyou.com.cn>

官方微博：<http://weibo.com/pinyouhudong>

招聘微博：<http://weibo.com/pinyouhudonghr>

我知道，你那一本正经的外表下面那颗躁动的心已经
蠢蠢欲动了!!!!

来吧！这儿的舞台无比广阔，这儿的技術绝对前沿，
这儿的事情会让你激动得颤抖!!!!

加入品友：hr@ipinyou.com



世界一流的研发团队

欢迎加入品友互动!
hr@ipinyou.com