

DeadBeef

博客园 首页 新随笔 联系 订阅 管理

随笔 - 69 文章 - 0 评论 - 0

Flume采集处理日志文件

1. Flume简介

Flume是Cloudera提供的一个高可用的，高可靠的，分布式的海量日志采集、聚合和传输的系统，Flume支持在日志系统中定制各类数据发送方，用于收集数据；同时，Flume提供对数据进行简单处理，并写到各种数据接受方（可定制）的能力。

1. 系统功能

1. 日志收集

Flume最早是Cloudera提供的日志收集系统，目前是Apache下的一个孵化项目，Flume支持在日志系统中定制各类数据发送方，用于收集数据。

1. 数据处理

Flume提供对数据进行简单处理，并写到各种数据接受方（可定制）的能力 Flume提供了从 console（控制台）、RPC（Thrift-RPC）、text（文件）、tail（UNIX tail）、syslog（syslog日志系统，支持TCP和UDP等2种模式），exec（命令执行）等数据源上收集数据的能力。

1. 工作方式

Flume采用了多Master的方式。为了保证配置数据的一致性，Flume[1] 引入了ZooKeeper，用于保存配置数据，ZooKeeper本身可保证配置数据的一致性和高可用，另外，在配置数据发生变化时，ZooKeeper可以通知Flume Master节点。Flume Master间使用gossip协议同步数据。

1. 流程结构

Flume的结构主要分为三部分：source、channel以及sink.其中source为源头，负责采集日志；channel为通道，负责传输和暂时储存；sink为目的地，将采集到的日志保存起来。在真正日志采集的过程中，根据待采集日志的类型以及存储需求，选择相应的类型的source、channel和sink进行配置，从而将日志采集并且保存起来。

1. Flume采集日志方案

1. 需求分析

1. 日志分类

操作系统：linux

日志更新类型：产生新日志，原日志结尾处追加

1. 采集时间需求

采集周期：短周期（一天之内）

1. 采集方案

1. 采集构架

使用flume采集日志文件的过程较简洁，只需选择恰当的source、channel和sink并且配置起来即可，若有特殊需求也可自己进行二次开发实现个人需求。

公告

昵称：deadbeef
园龄：2年
粉丝：3
关注：1
[+加关注](#)

< 2016年7月 >						
日	一	二	三	四	五	六
26	27	28	29	30	1	2
3	4	5	6	7	8	9
10	11	12	13	14	15	16
17	18	19	20	21	22	23
24	25	26	27	28	29	30
31	1	2	3	4	5	6

搜索

常用链接

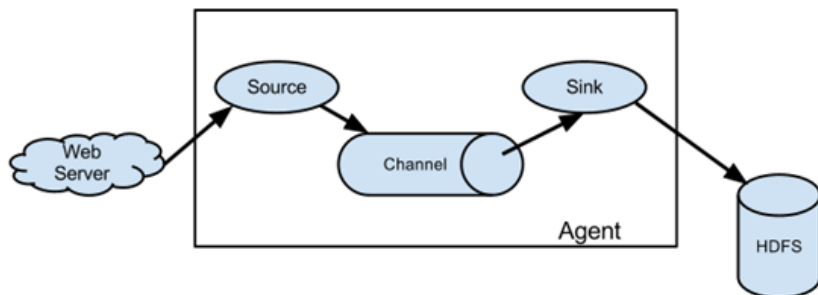
[我的随笔](#)
[我的评论](#)
[我的参与](#)
[最新评论](#)
[我的标签](#)
[更多链接](#)

随笔档案

2015年11月 (2)
2015年10月 (2)
2015年9月 (5)
2015年8月 (9)
2015年7月 (15)
2015年5月 (2)
2015年4月 (15)
2015年2月 (1)
2014年11月 (1)
2014年9月 (5)
2014年8月 (4)
2014年7月 (8)

阅读排行榜

- 2014高教社杯全国大学生...
- Flume采集处理日志文件(1...
- 【LINUX】pwnable.kr c...
- 【LINUX】pwnable.kr c...



5. 【PWN】pwnable.kr ech...

推荐排行榜

1. 【Security】常用的渗透测...

具体过程为：按照需求配置一个agent，选取适当的source和sink，然后启动该agent，开始采集日志。

1. source

flume提供多种source供用户进行选择，尽可能多的满足大部分日志采集的需求，常用的source的类型包括avro、exec、netcat、spooling-directory和syslog等。具体的使用范围和配置方法详见[source](#)。

1. channel

flume中的channel不如source和sink那么重要，但却是不可忽视的组成部分。常用的channel为memory-channel，同时也有其他类型的channel，如JDBC、file-channel、custom-channel等，详情见[channel](#)。

1. sink

flume的sink也有很多种，常用的包括avro、logger、HDFS、hbase以及file-roll等，除此之外还有其他类型的sink，如thrift、IRC、custom等。具体的使用范围和使用方法详见[sink](#)。

1. Flume处理日志

Flume不止可以采集日志，还可以对日志进行简单的处理，在source处可以通过interceptor对日志正文处的重要内容进行过滤提取，在channel处可以通过header进行分类，将不同类型的日志投入不同的通道中，在sink处可以通过正则序列化来将正文内容进行进一步的过滤和分类。

1. Flume Source Interceptors

Flume可以通过interceptor将重要信息提取出来并且加入到header中，常用的interceptor有时间戳、主机名和UUID等，用户也可以根据个人需求编写正则过滤器，将某些特定格式的日志内容过滤出来，以满足特殊需求。

1. Flume Channel Selectors

Flume可以根据需求将不同的日志传输进不同的channel，具体方式有两种：复制和多路传输。复制就是不对日志进行分组，而是将所有日志都传输到每个通道中，对所有通道不做区别对待；多路传输就是根据指定的header将日志进行分类，根据分类规则将不同的日志投入到不同的channel中，从而将日志进行人为的初步分类。

1. Flume Sink Processors

Flume在sink处也可以对日志进行处理，常见的sink处理器包括custom、failover、load balancing和default等，和interceptor一样，用户也可以根据特殊需求使用正则过滤处理器，将日志内容过滤出来，但和interceptor不同的是在sink处使用正则序列化过滤出的内容不会加入到header中，从而不会使日志的header显得过于臃肿。

1. 附录

1. 常见的source

1. avro source

avro可以监听和收集指定端口的日志，使用avro的source需要说明被监听的主机ip和端口号，下面给出一个具体的例子：

```
al.sources = r1
```

```
al.channels = cl

al.sources.r1.type = avro

al.sources.r1.channels = cl

al.sources.r1.bind = 0.0.0.0

al.sources.r1.port = 4141
```

1. exec source

exec可以通过指定的操作对日志进行读取，使用exec时需要指定shell命令，对日志进行读取，下面给出一个具体的例子：

```
al.sources = r1

al.channels = cl

al.sources.r1.type = exec

al.sources.r1.command = tail -F /var/log/secure

al.sources.r1.channels = cl
```

1. spooling-directory source

spo_dir可以读取文件夹里的日志，使用时指定一个文件夹，可以读取该文件夹中的所有文件，需要注意的是该文件夹中的文件在读取过程中不能修改，同时文件名也不能修改。下面给出一个具体的例子：

```
agent-1.channels = ch-1

agent-1.sources = src-1

agent-1.sources.src-1.type = spooldir

agent-1.sources.src-1.channels = ch-1

agent-1.sources.src-1.spoolDir = /var/log/apache/flumeSpool

agent-1.sources.src-1.fileHeader = true
```

1. syslog source

syslog可以通过syslog协议读取系统日志，分为tcp和udp两种，使用时需指定ip和端口，下面给出一个udp的例子：

```
al.sources = r1

al.channels = cl

al.sources.r1.type = syslogudp

al.sources.r1.port = 5140

al.sources.r1.host = localhost

al.sources.r1.channels = cl
```

1. 常见的channel

Flume的channel种类并不多，最常用的是memory channel，下面给出例子：

```
al.channels = cl

al.channels.cl.type = memory

al.channels.cl.capacity = 10000

al.channels.cl.transactionCapacity = 10000

al.channels.cl.byteCapacityBufferPercentage = 20

al.channels.cl.byteCapacity = 800000
```

1. 常见的sink

1. logger sink

logger顾名思义，就是将收集到的日志写到flume的log中，是个十分简单但非常实用的sink

1. avro sink

avro可以将接受到的日志发送到指定端口，供级联agent的下一跳收集和接受日志，使用时需要指定目的ip和端口：例子如下：

```
al.channels = c1

al.sinks = k1

al.sinks.k1.type = avro

al.sinks.k1.channel = c1

al.sinks.k1.hostname = 10.10.10.10

al.sinks.k1.port = 4545
```

1. file roll sink

file_roll可以将一定时间内收集到的日志写到一个指定的文件中，具体过程为用户指定一个文件夹和一个周期，然后启动agent，这时该文件夹会产生一个文件将该周期内收集到的日志全部写进该文件内，直到下一个周期再次产生一个新文件继续写入，以此类推，周而复始。下面给出一个具体的例子：

```
al.channels = c1

al.sinks = k1

al.sinks.k1.type = file_roll

al.sinks.k1.channel = c1

al.sinks.k1.sink.directory = /var/log/flume
```

1. hdfs sink

hdfs与file roll有些类似，都是将收集到的日志写入到新创建的文件中保存起来，但区别是file roll的文件存储路径为系统的本地路径，而hdfs的存储路径为分布式的文件系统hdfs的路径，同时hdfs创建新文件的周期可以是时间，也可以是文件的大小，还可以是采集日志的条数。具体实例如下：

```
al.channels = c1

al.sinks = k1

al.sinks.k1.type = hdfs

al.sinks.k1.channel = c1

al.sinks.k1.hdfs.path = /flume/events/%y-%m-%d/%H%M/%S

al.sinks.k1.hdfs.filePrefix = events-

al.sinks.k1.hdfs.round = true

al.sinks.k1.hdfs.roundValue = 10

al.sinks.k1.hdfs.roundUnit = minute
```

1. hbase sink

hbase是一种数据库，可以储存日志，使用时需要指定存储日志的表名和列族名，然后agent就可以将收集到的日志逐条插入到数据库中。例子如下：

```
al.channels = c1

al.sinks = k1

al.sinks.k1.type = hbase

al.sinks.k1.table = foo_table

al.sinks.k1.columnFamily = bar_cf
```

```
al.sinks.k1.serializer =
org.apache.flume.sink.hbase.RegexHbaseEventSerializer

al.sinks.k1.channel = c1
```

好文要顶

关注我

收藏该文



deadbeef
关注 - 1
粉丝 - 3
[+加关注](#)

0

0

(请您对文章做出评价)

» 下一篇：[某互联网创业公司人才需求书](#)

posted @ 2014-07-28 10:33 deadbeef 阅读(1130) 评论(0) 编辑 收藏

[刷新评论](#) [刷新页面](#) [返回顶部](#)

注册用户登录后才能发表评论，请 [登录](#) 或 [注册](#)，[访问网站首页](#)。

- 【推荐】50万行VC++源码: 大型组态工控、电力仿真CAD与GIS源码库
- 【推荐】融云即时通讯云 - 豆果美食、Faceu等亿级APP都在用
- 【福利】你是我的好朋友，我要送你个天猫红包
- 【活动】蚂蚁金服开放平台合作伙伴大会(北京8.10)

ActiveReports

企业级报表服务平台

单独部署、集成应用、报表制作、数据整合
权限管理、移动办公、二次集成开发

立即了解

最新IT新闻:

- 51岁郭富城创业了，为何选择做男士护肤品？
- Facebook市值赶超股神巴菲特公司 全美排名第五
- 为了研发无人驾驶系统 苹果又挖来了QNX创始人
- 亚马逊第二季度净利8.57亿美元 同比大涨832%
- Alphabet第二季度净利润48.8亿美元 同比增长24%

» 更多新闻...

消息推送领导品牌全面升级

[详情点击](#)

最新知识库文章:

- 可是姑娘，你为什么要编程呢？
- 知其所以然（以算法学习为例）
- 如何给变量取个简短且无歧义的名字
- 编程的智慧
- 写给初学前端工程师的一封信

» 更多知识库文章...