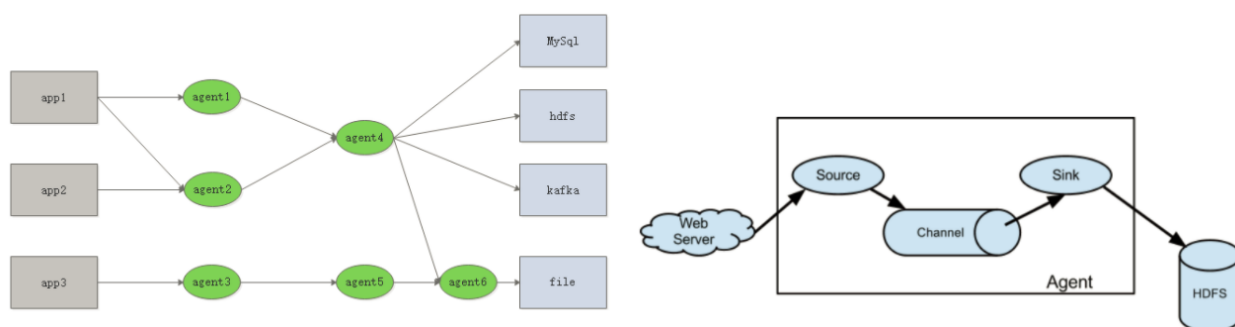


实时点击流日志处理



flume的核心是agent, 而agent包含source、channel、sink三个组件。

- source: source组件是专门用来收集数据的, 可以处理各种类型、各种格式的日志数据,包括 avro、thrift、exec、jms、spooling directory、netcat、sequence generator、syslog、http、legacy、自定义。
- channel: source组件把数据收集来以后, 临时存放在channel中, 即channel组件在agent中是专门用来存放临时数据的——对采集到的数据进行简单的缓存, 可以存放在memory、jdbc、file等等。
- sink: sink组件是用于把数据发送到目的地的组件, 目的地包括hdfs、logger、avro、thrift、ipc、file、null、hbase、solr、自定义

工作流程: flume把数据从数据源(source)收集过来, 再将数据送到指定的目的地(sink)。为保证传输的过程一定成功, 在送到目的地(sink)之前, 会先缓存数据(channel),待数据真正到达目的地(sink)后, flume再删除自己缓存的数据。在整个数据的传输的过程中, 流动的是event (基本单位), 因此事务保证是在event级别进行的。

event: event将传输的数据进行封装, 是flume传输数据的基本单位, 如果是文本文件, 通常是一行记录。event也是事务的基本单位。event在单个agent中经历source—channel—sink过程, 后面可能输出到下一个agent或者flume外的系统中。event本身为一个字节数组, 其携带headers(头信息)信息, 消息体, 消息内容

从上面图中可以看出flume支持多级的成网状数据流动, 非常的灵活好用, 这应该就是flume广泛使用原因吧。比如数据扇入到同一个agent或者扇出到多个agent。

flume配置

编辑: /root/bigdata/flume/conf/click_trace_log_hdfs.properties

```
# Name the components on this agent
al.sources = r1
al.sinks = k1 k2
al.channels = c1

# Describe/configure the source
al.sources.r1.type = exec
al.sources.r1.command = tail -F /root/meiduoSourceCode/logs/click_trace.log
al.sources.r1.channels = c1

al.sources.r1.interceptors = t1
al.sources.r1.interceptors.t1.type = timestamp

# Use a channel which buffers events in memory
al.channels.c1.type = memory
al.channels.c1.capacity = 1000
al.channels.c1.transactionCapacity = 100

al.sinks.k1.type = hdfs
al.sinks.k1.channel = c1
al.sinks.k1.hdfs.path = hdfs://localhost:9000/project2-meiduo-rs/logs/click-trace/%y-%m-%d
al.sinks.k1.hdfs.userLocalTimeStamp = true
al.sinks.k1.hdfs.filePrefix = click-trace-
al.sinks.k1.hdfs.fileType = DataStream
al.sinks.k1.hdfs.writeFormat = Text
al.sinks.k1.hdfs.round = true
al.sinks.k1.hdfs.roundValue = 10
al.sinks.k1.hdfs.roundUnit = minute

al.sinks.k2.channel = c1
al.sinks.k2.type = org.apache.flume.sink.kafka.KafkaSink
al.sinks.k2.kafka.topic = meiduo_click_trace
al.sinks.k2.kafka.bootstrap.servers = localhost:9092
al.sinks.k2.kafka.flumeBatchSize = 20
al.sinks.k2.kafka.producer.acks = 1
al.sinks.k2.kafka.producer.linger.ms = 1
al.sinks.k2.kafka.producer.compression.type = snappy
```

启动flume对点击流日志进行采集，分别发送到kafka和hdfs: `flume-ng agent -f /root/bigdata/flume/conf/click_trace_log_hdfs.properties -n al`

启动Kafka(如果还未启动的话): `cd /root/bigdata/kafka && bin/zookeeper-server-start.sh -daemon config/zookeeper.properties && bin/kafka-server-start.sh config/server.properties`

```

In [1]: 1 import os
2 # 配置pyspark和spark driver运行时 使用的python解释器
3 JAVA_HOME = '/root/bigdata/jdk'
4 PYSARK_PYTHON = '/miniconda2/envs/py365/bin/python'
5 # 当存在多个版本时, 不指定很可能会导致出错
6 os.environ['PYSARK_PYTHON'] = PYSARK_PYTHON
7 os.environ['PYSARK_DRIVER_PYTHON'] = PYSARK_PYTHON
8 os.environ['JAVA_HOME'] = JAVA_HOME
9 # 注意, 如果是使用jupyter或ipython中, 利用spark streaming链接kafka的话, 必须加上下面语
10 # 同时注意: spark version>2.2.2的话, pyspark中的kafka对应模块已被遗弃, 因此这里暂时只
11 os.environ["PYSARK_SUBMIT_ARGS"] = "--packages org.apache.spark:spark-streaming-kafka
12 # 配置spark信息
13 from pyspark import SparkConf
14 import pyspark
15
16 SPARK_APP_NAME = "meiduo_logs"
17 SPARK_URL = "spark://192.168.58.100:7077"
18
19 conf = SparkConf() # 创建spark config对象
20 config = (
21     ("spark.app.name", SPARK_APP_NAME), # 设置启动的spark的app名称, 没有提供, 将随
22     ("spark.executor.memory", "2g"), # 设置该app启动时占用的内存用量, 默认1g, 指一
23     ("spark.master", SPARK_URL), # spark master的地址
24     ("spark.executor.cores", "2"), # 设置spark executor使用的CPU核心数, 指一台虚拟
25     # ("hive.metastore.uris", "thrift://localhost:9083"), # 配置hive元数据的访问,
26
27     # 以下三项配置, 可以控制执行器数量
28     # ("spark.dynamicAllocation.enabled", True),
29     # ("spark.dynamicAllocation.initialExecutors", 1), # 1个执行器
30     # ("spark.shuffle.service.enabled", True)
31     # ('spark.sql.pivotMaxValues', '99999'), # 当需要pivot DF, 且值很多时, 需要修改, 默
32 )
33 # 查看更详细配置及说明: https://spark.apache.org/docs/latest/configuration.html
34
35 conf.setAll(config)
36
37 # 利用config对象, 创建spark session
38 sc = pyspark.SparkContext(master=SPARK_URL, conf=conf)

```

```

In [2]: 1 # 注意: 初次安装并运行时, 由于使用了kafka, 所以会自动下载一系列的依赖jar包, 会耗费一定
2
3 from pyspark.streaming.kafka import KafkaUtils
4 from pyspark.streaming import StreamingContext
5
6 # 第2个参数表示 the time interval (in seconds) at which streaming data will be divided
7 ssc = StreamingContext(sc, 0.5)
8
9 kafkaParams = {"metadata.broker.list": "192.168.58.100:9092"}
10 dstream = KafkaUtils.createDirectStream(ssc, ["meiduo_click_trace"], kafkaParams)

```

...

```
In [29]: 1 import re
2 def map(row):
3     match = re.search("\
4 exposure_timesteamp<(P<exposure_timesteamp>.*?)> \
5 exposure_loc<(P<exposure_loc>.*?)> \
6 timesteamp<(P<timesteamp>.*?)> \
7 behavior<(P<behavior>.*?)> \
8 uid<(P<uid>.*?)> \
9 sku_id<(P<sku_id>.*?)> \
10 cate_id<(P<cate_id>.*?)> \
11 stay_time<(P<stay_time>.*?)>", row[1])
12
13     result = []
14     if match:
15         result.append(("exposure_timesteamp", match.group("exposure_timesteamp")))
16         result.append(("exposure_loc", match.group("exposure_loc")))
17         result.append(("timesteamp", match.group("timesteamp")))
18         result.append(("behavior", match.group("behavior")))
19         result.append(("uid", match.group("uid")))
20         result.append(("sku_id", match.group("sku_id")))
21         result.append(("cate_id", match.group("cate_id")))
22         result.append(("stay_time", match.group("stay_time")))
23     return result
24
25 def foreachRDD(rdd):
26     print("foreachRDD", rdd.collect())
```

```
In [30]: 1 dstream.map(map).foreachRDD(foreachRDD)
```

```
In [31]: 1 ssc.start()
```

...

```

In [32]: 1 # 生成日志
2 import logging#log: 记录
3 import time
4
5 def get_logger(logger_name, path, level):
6
7     # 创建logger
8     logger = logging.getLogger(logger_name)
9     # level: OFF、FATAL、ERROR、WARN、INFO、DEBUG、ALL或者自己定义的级别
10    logger.setLevel(level)
11
12    # 创建formatter
13    # %(asctime)s: 打印日志的时间
14    # %(message)s: 打印日志信息
15    fmt = '%(asctime)s: %(message)s'
16    datefmt = '%Y/%m/%d %H:%M:%S'
17    formatter = logging.Formatter(fmt, datefmt)
18
19    # 创建handler
20    # FileHandler: writes formatted logging records to disk files
21    handler = logging.FileHandler(path)
22    handler.setLevel(level)
23
24    # 添加handler和formatter 到 logger
25    handler.setFormatter(formatter)
26    logger.addHandler(handler)
27
28    return logger
29
30 click_trace_logger = get_logger('click_trace', '/root/workspace/3.rs_project/project2/
31                                logging.DEBUG)
32
33 # 点击流日志
34 exposure_timesteamp = time.time()
35 exposure_loc = 'detail'
36 timesteamp = time.time()
37 behavior = 'pv' # pv fav cart buy
38 uid = 1
39 sku_id = 1
40 cate_id = 1
41 stay_time = 60
42 ## 假设某点击流日志记录格式如下:
43 click_trace_logger.info("exposure_timesteamp<%d> exposure_loc<%s> timesteamp<%d> behav
44                        %(exposure_timesteamp, exposure_loc, timesteamp, behavior, uid
45

```

```

foreachRDD []
foreachRDD [(('exposure_timesteamp', '1608781198'), ('exposure_loc', 'detail'), ('tim
esteamp', '1608782037'), ('behavior', 'pv'), ('uid', '1'), ('sku_id', '1'), ('cate_i
d', '1'), ('stay_time', '60'))]
foreachRDD []
foreachRDD []
foreachRDD []
foreachRDD []
foreachRDD []

```

```
foreachRDD []
foreachRDD []
foreachRDD []
foreachRDD []
foreachRDD []
foreachRDD [('exposure_timesteamp', '1608781198'), ('exposure_loc', 'detail'), ('tim
esteamp', '1608782037'), ('behavior', 'pv'), ('uid', '1'), ('sku_id', '1'), ('cate_i
d', '1'), ('stay_time', '60')], [('exposure_timesteamp', '1608781198'), ('exposure_lo
c', 'detail'), ('timesteamp', '1608782037'), ('behavior', 'pv'), ('uid', '1'), ('sku_
id', '1'), ('cate_id', '1'), ('stay_time', '60')], [('exposure_timesteamp', '16087811
```

In []: 1 ssc.stop()