

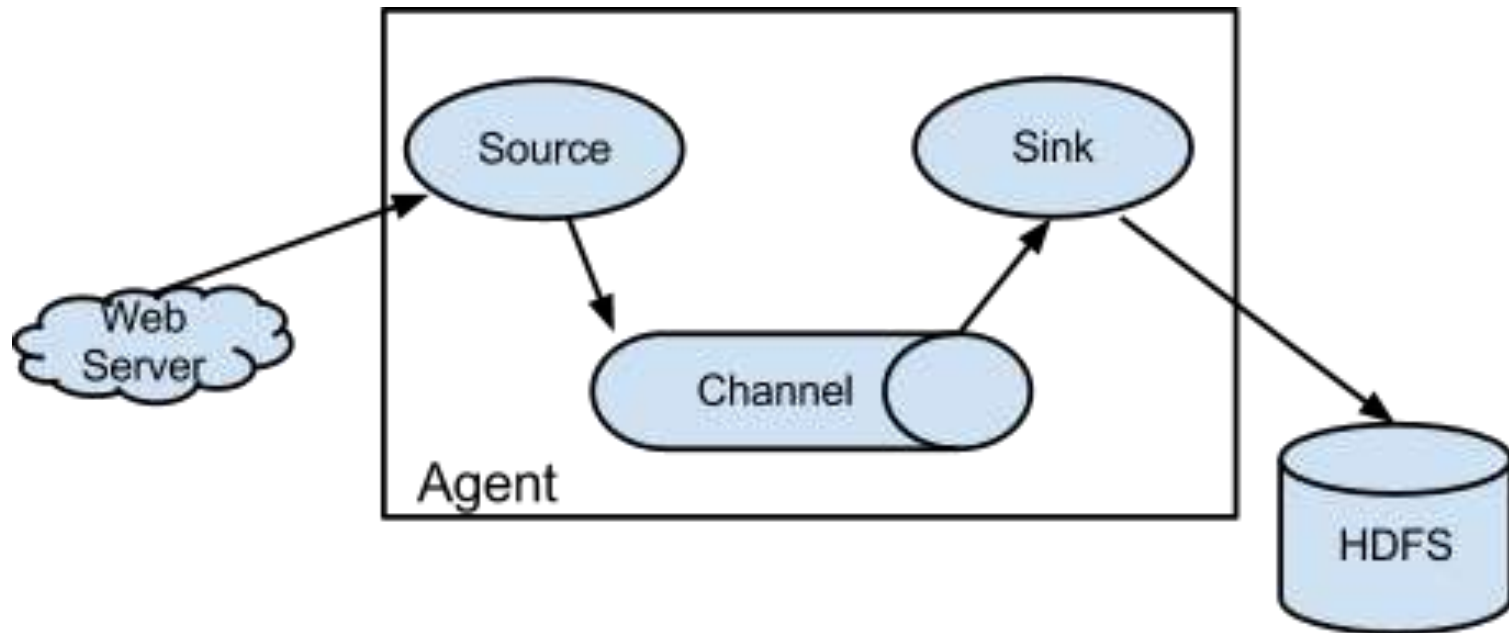


Hadoop 추가자료

Flume



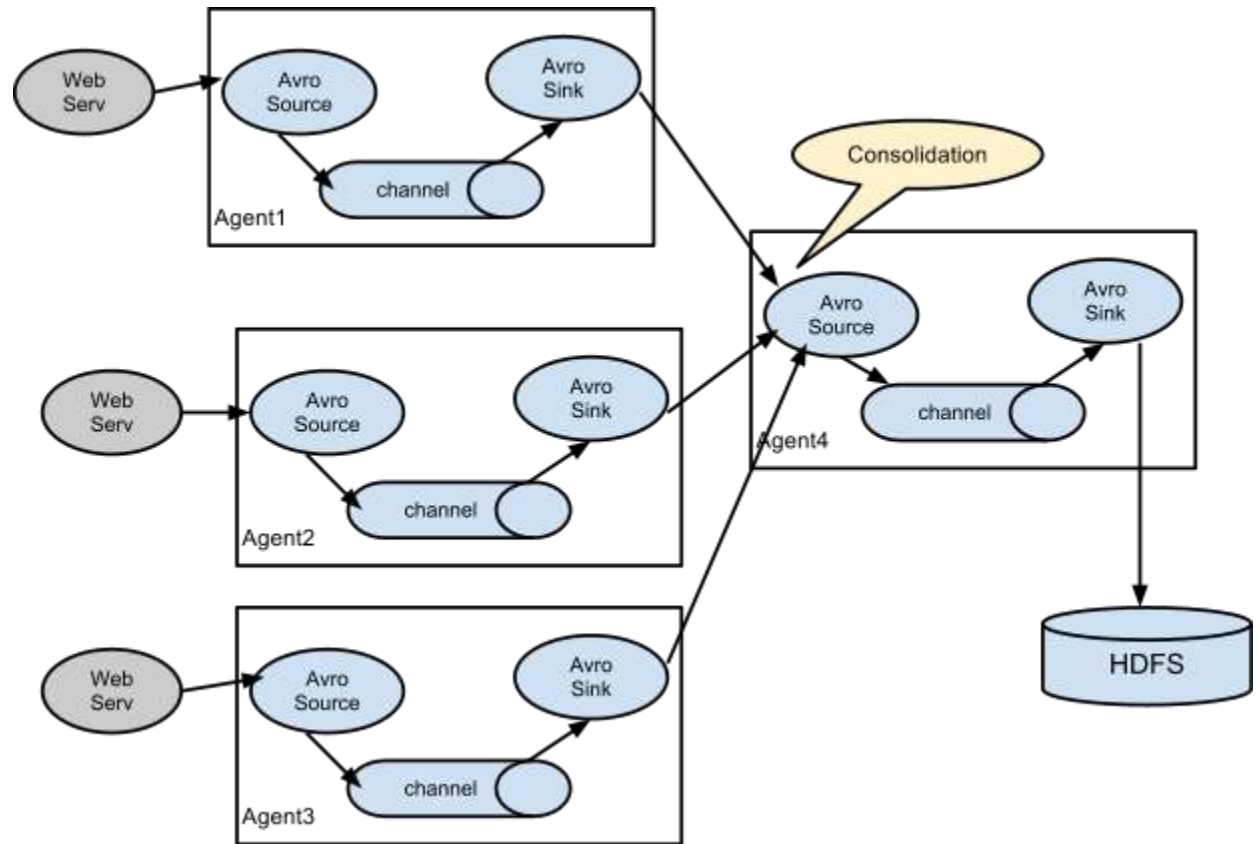
Apache Flume™



Flume-NG



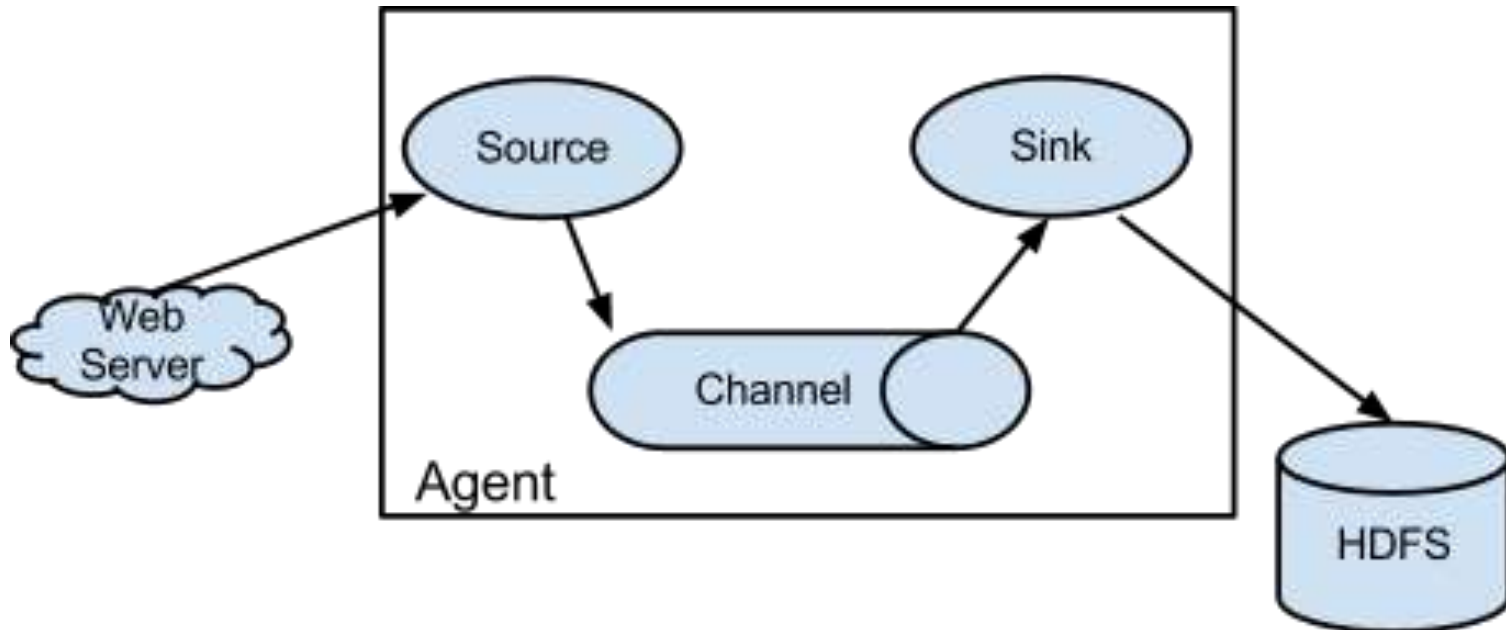
- 분산 데이터 수집/전송시스템
- 최초 설계 목적은 이벤트나 로그 구조의 데이터를 지속적으로 하둡 HDFS에 저장
- 에이전트
- 확장하여 다양한 분야에 활용 가능



Flume-NG



- 노드 : Flume이 구동되는 머신
- 모든 노드에는 “source”와 “sink”가 있음
 - source 예 : `tail -F /var/log/httpd/access_log`
 - sink 예 : `dfs(“hdfs://namenode/log/{host}%/%y%m%d”)`
- 데이터플로우 : 노드들의 체인



Flume-NG 설치



■ Flume 설치

```
$ tar -zxvf apache-flume-1.4.0-bin.tar.gz  
$ ln -s apache-flume-1.4.0-bin apache-flume  
cp flume-conf.properties.template flume-conf
```

■ conf/flume-env.sh 에 다음의 내용 추가

```
JAVA_HOME=/usr/java/java  
FLUME_CLASSPATH="/home/hadoop/apache-flume/lib/"  
export PATH=$PATH:/home/hadoop/hadoop/bin/
```

Flume-NG 설치



■ Starting an agent

```
$ bin/flume-ng agent -n $agent_name -c conf -f conf/flume-  
conf.properties.template
```

Flume-NG 설치



■ conf/flume.conf 에 설정

agent의 각 요소에 이름을 부여

```
a1.sources = r1
```

```
a1.sinks = k1
```

```
a1.channels = c1
```

source 설정

```
a1.sources.r1.type = netcat
```

```
a1.sources.r1.bind = localhost
```

```
a1.sources.r1.port = 44444
```

sink 설정

```
a1.sinks.k1.type = logger
```

채널 설정

```
a1.channels.c1.type = memory
```

```
a1.channels.c1.capacity = 1000
```

```
a1.channels.c1.transactionCapacity = 100
```

source 와 sink 를 채널에 연결

```
a1.sources.r1.channels = c1
```

```
a1.sinks.k1.channel = c1
```

Flume-NG 설치



■ conf/flume.conf 에 설정

```
# Name the components on this agent
a1.sources = r1
a1.sinks = k1
a1.channels = c1
### source 설정
a1.sources.r1.type = exec
a1.sources.r1.command = tail -F /home/hadoop/syslog/a.txt
a1.sources.r1.channels = c1
### sink 를 hdfs로 설정
a1.sinks.k1.type = hdfs
a1.sinks.k1.channel = c1
a1.sinks.k1.hdfs.path = hdfs://hadoop01:9000/user/hadoop/logdata/a.txt
a1.sinks.k1.hdfs.filePrefix = events-
a1.sinks.k1.hdfs.round = true
a1.sinks.k1.hdfs.roundValue = 10
a1.sinks.k1.hdfs.roundUnit = minute
### 채널 설정
a1.channels.c1.type = memory
a1.channels.c1.capacity = 1000
a1.channels.c1.transactionCapacity = 100
### Bind the source and sink to the channel
a1.sources.r1.channels = c1
a1.sinks.k1.channel = c1
```


Flume-NG 설치



■ 예제 실행

```
$ bin/flume-ng agent --conf ./conf/ -f conf/flume.conf W  
-Dflume.root.logger=DEBUG,console -n a1
```

■ 다음과 같이 실행 로그 출력

```
2013-06-18 14:00:49,784 (hdfs-hdfs-sink-call-runner-0) [INFO -  
org.apache.flume.sink.hdfs.BucketWriter.doOpen(BucketWriter.java:189)] Creating  
hdfs://localhost:54310/tmp/system.log//FlumeData.1371589249458.tmp
```



[Go to parent directory](#)

Name	Type	Size	Replication	Block Size	Modification Time	Permission	Owner	Group
events-.1382271326590	file	0.41 KB	3	64 MB	2013-10-20 21:15	rw-r--r--	hadoop	supergroup
events-.1382271326591	file	0.31 KB	3	64 MB	2013-10-20 21:19	rw-r--r--	hadoop	supergroup

Flume-NG 설치



■ Flume Sources

Flume Sources	Avro Source
	Thrift Source
	Exec Source
	JMS Source
	Spooling Directory Source
	NetCat Source
	Syslog Sources
	Syslog UDP Source
	HTTP Source
	Legacy Sources
	Custom Source

Flume-NG 설치



■ HDFS Sink

Name	Default	Description
channel	–	
type	–	The component type name, needs to be <code>hdfs</code>
hdfs.path	–	HDFS directory path (eg <code>hdfs://namenode/flume/webdata/</code>)
hdfs.filePrefix	FlumeData	Name prefixed to files created by Flume in hdfs directory
hdfs.fileSuffix	–	Suffix to append to file (eg <code>.avro</code> - <i>NOTE: period is not automatically added</i>)
hdfs.inUsePrefix	–	Prefix that is used for temporal files that flume actively writes into
hdfs.inUseSuffix	<code>.tmp</code>	Suffix that is used for temporal files that flume actively writes into
hdfs.rollInterval	30	Number of seconds to wait before rolling current file (0 = never roll based on time interval)
hdfs.rollSize	1024	File size to trigger roll, in bytes (0: never roll based on file size)
hdfs.rollCount	10	Number of events written to file before it rolled (0 = never roll based on number of events)
hdfs.idleTimeout	0	Timeout after which inactive files get closed (0 = disable automatic closing of idle files)
hdfs.batchSize	100	number of events written to file before it is flushed to HDFS
hdfs.codeC	–	Compression codec. one of following : <code>gzip</code> , <code>bzip2</code> , <code>lzo</code> , <code>lzop</code> , <code>snappy</code>
hdfs.fileType	SequenceFile	File format: currently <code>SequenceFile</code> , <code>DataStream</code> or <code>CompressedStream</code> (1) <code>DataStream</code> will not compress output file and please don't set <code>codeC</code> (2) <code>CompressedStream</code> requires set <code>hdfs.codeC</code> with an available <code>codeC</code>
hdfs.maxOpenFiles	5000	Allow only this number of open files. If this number is exceeded, the oldest file is closed.
hdfs.minBlockReplicas	–	Specify minimum number of replicas per HDFS block. If not specified, it comes from the default Hadoop config in the classpath.
hdfs.writeFormat	–	Format for sequence file records. One of “Text” or “Writable” (the default).
hdfs.callTimeout	10000	Number of milliseconds allowed for HDFS operations, such as open, write, flush, close. This number should be increased if many HDFS timeout operations are occurring.
hdfs.threadsPoolSize	10	Number of threads per HDFS sink for HDFS IO ops (open, write, etc.)
hdfs.rollTimerPoolSize	1	Number of threads per HDFS sink for scheduling timed file rolling