

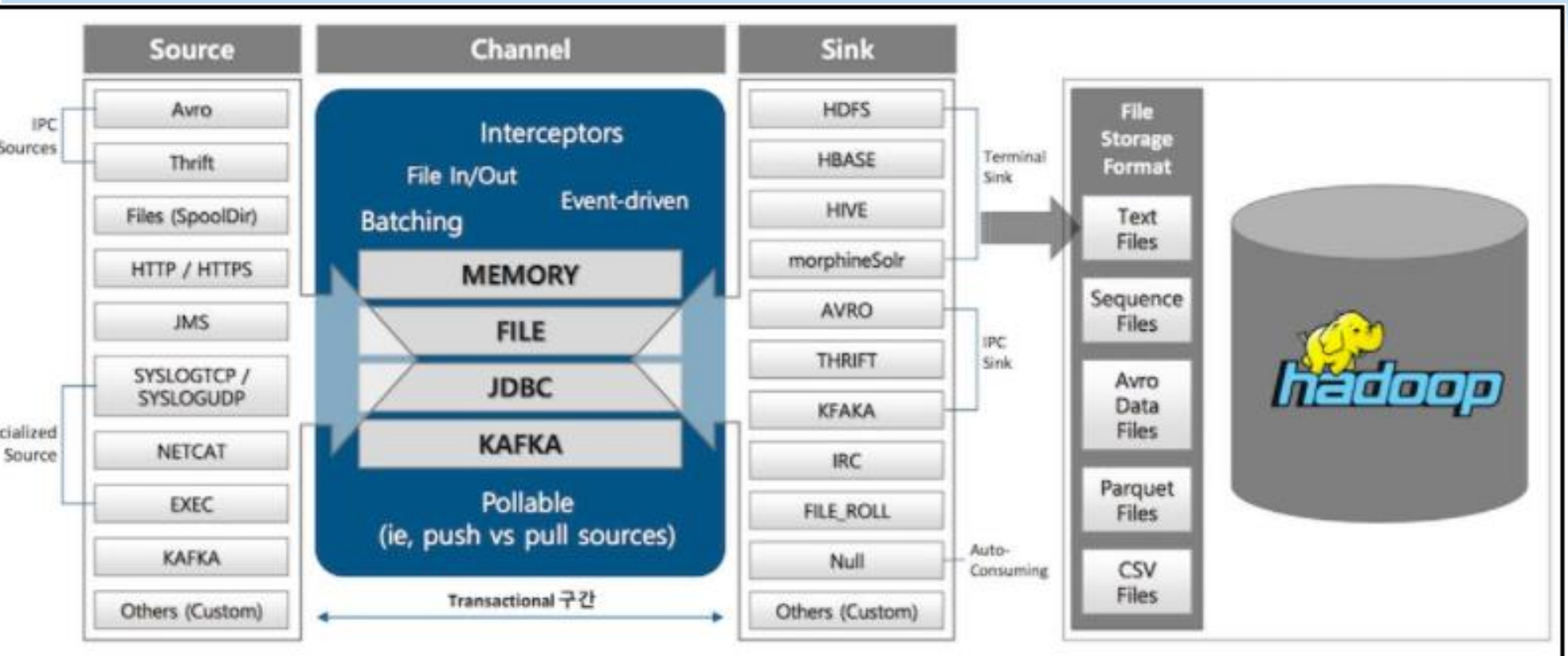
**3강**

# **빅데이터 수집**

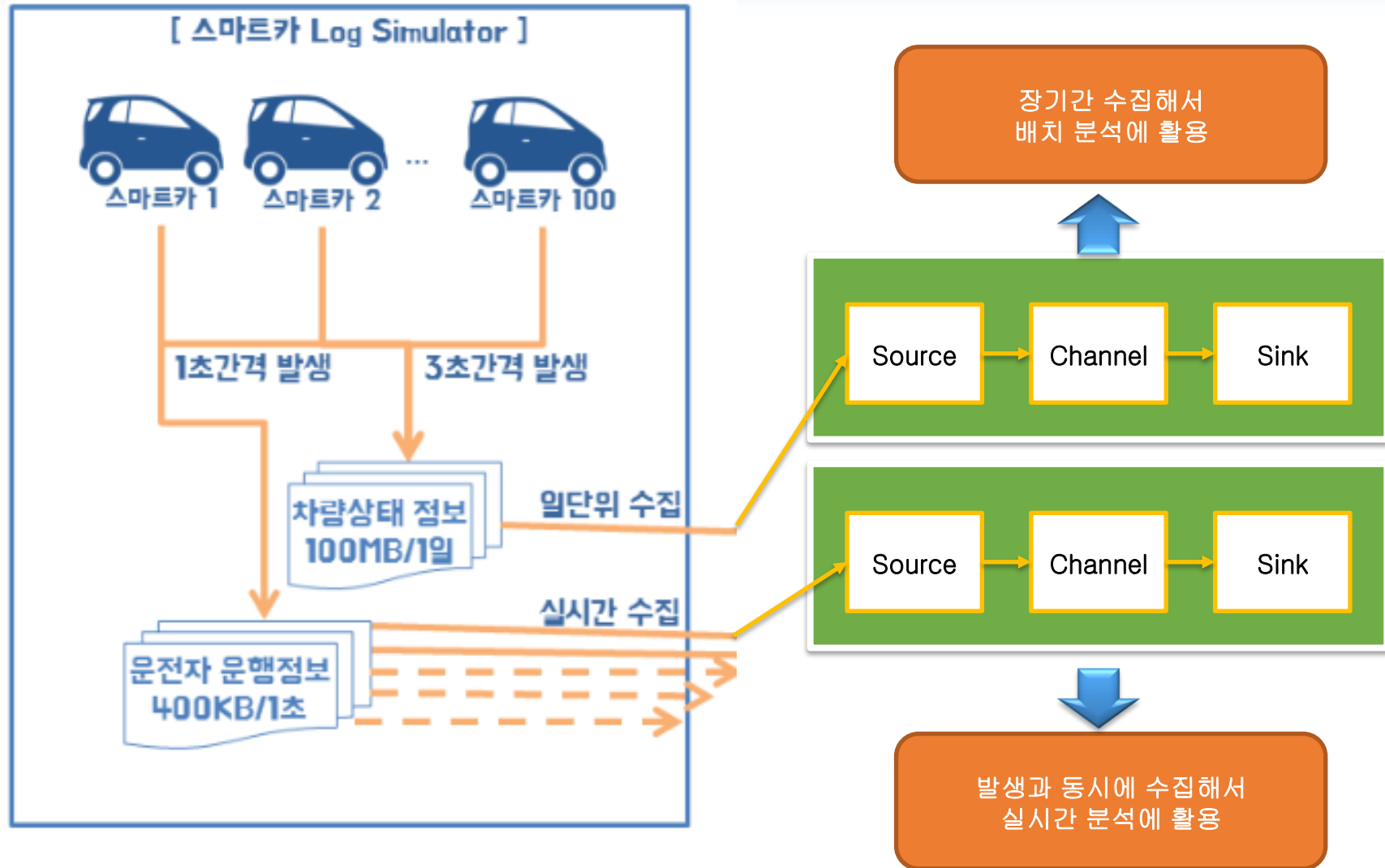


# Flume

# Flume이란



# Agent 생성



# Flume 설치

Cloudera Manager **클러스터** 호

**Cluster 1** **작업** 상태

**서비스 추가**

Add Hosts

상태

추가할 서비스 유형을 선택합니다.

서비스 유형	설명
<input type="radio"/> ADLS Connector	The ADLS Connect
<input type="radio"/> Accumulo	The Apache Accumulo CDH6.
<input checked="" type="radio"/> Flume	Flume은 대부분의 클러스터에서 사용할 수 있는 데이터 수집 서비스입니다.

클러스터 클릭하여 홈으로 이동  
작업 -> 서비스 추가 클릭  
Flume 선택  
Agent에서 server2 선택

☒ Select Dependencies

☒ 역할 할당 사용자 지정

☐ 변경 내용 검토

## 역할 할당 사용자 지정

여기에서 새 서비스에 대한 역할 할당을 사용자 지정합니다.

역할 할당을 호스트별로 볼 수도 있습니다.

**호스트별로 보기**

Agent x 1 새로 만들기

server02.hadoop.com ▼

# 설치 완료

## Cluster 1에 Flume 서비스 추가

- ✓ Select Dependencies
- ✓ 역할 할당 사용자 지정
- ✓ 변경 내용 검토
- 요약

### 요약

✓ 새 서비스가 클러스터에 설치 및 구성되었습니다.

**참고:** 여전히 새 서비스를 시작해야 할 수 있습니다. 시작하기 전에 오래된 구성이 포함된 모든 종속 서비스를 재시작하는 것이 좋습니다. 이러한 작업은 아래에서 **완료**를 클릭하면 주 페이지에서 수행할 수 있습니다.

변경내용검토 후 완료

뒤로

완료

# 메모리 설정

Cluster 1

**Flume** 작업 ▾

상태    인스턴스    **구성**    명령    메트릭 세부 정보    차트 라이브러리    감사    쿼리 링크 ▾

**heap**

**필터**

▼ 범위

Flume (서비스 차원)	0
Agent	5

▼ 범주

Flume-NG Solr 싱크	0
고급	2
기본	0
로그	0
리소스 관리	1

메모리가 부족하면 힙 덤프 ☒ Agent Default Group

힙 덤프 디렉토리  
oom\_heap\_dump\_dir

Flume 이동  
구성에서 heap 검색  
100으로 변경

Agent의 Java 힙 크기(단위: Agent Default Group ↻  
바이트)

**100** MiB ▾

# 서비스 시작

Cloudera Manager 클러스터 ▾ 호스

Cluster 1

Flume

작업 ▾

시작

중지

재시작

상태 테스트

작업 시작 클릭

컨텍스트 Flume Oct 5, 9:48:05 PM 중단

0/1단계가 완료되었습니다.

☒ Show All Steps ☐ Show Only Failed Steps ☐ Show Only Running Steps

> 서비스의 1 역할 시작

0/1개의 시작 명령이 완료되었습니다.

중단 닫기



# Kafka

# Kafka

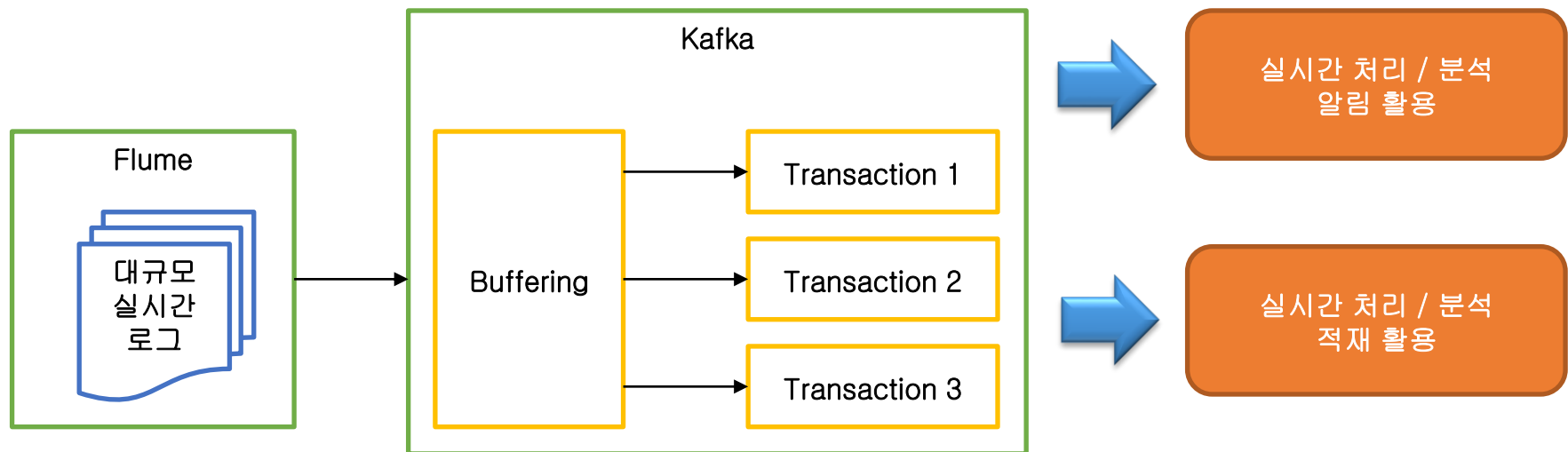
## ❖ 정의

- Message Oriented Middleware
- 대규모 메시지 데이터를 비동기 방식으로 중계
- 버퍼링을 통한 시스템의 안정적 전송

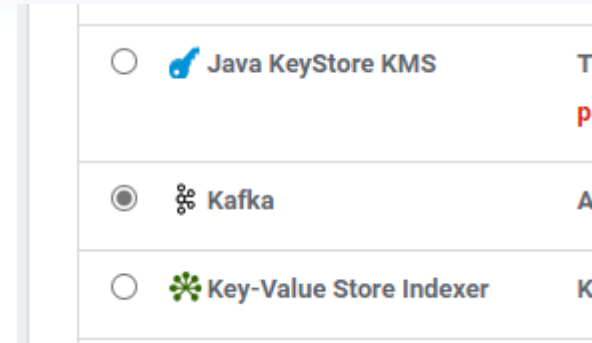
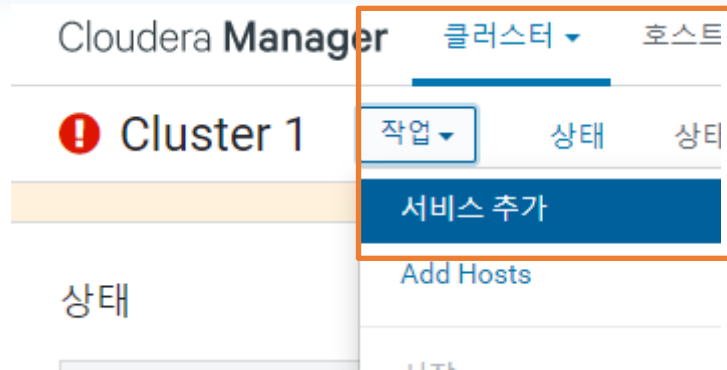
## ❖ 주요 기능

- Broker : 서비스 인스턴스
- Topic : Broker에서 데이터의 발행/소비 처리를 위한 저장소
- Provider : Broker의 특정 Topic에 데이터 전송
- Consumer : Broker의 특정 Topic에서 데이터 수신

# 카프카 활용



# 설정



클러스터 클릭  
작업 -> 서비스 추가  
Kafka 선택  
Broker를 server2로 변경

## 역할 할당 사용자 지정

여기에서 새 서비스에 대한 역할 할당을 사용자 지정할 수 있지만 단일 호스트에 너무 많은 수의 역할을 할당하는 등 올바르지 않게 할당할

역할 할당을 호스트별로 볼 수도 있습니다. [호스트별로 보기](#)

🌀 Kafka Broker × 1 새로 만들기

server02.hadoop.com ▼

🌀 Kafka MirrorMaker

호스트 선택

🌀 Gateway

호스트 선택

# 설치 완료

✓ Select Dependencies

✓ 역할 할당 사용자 지정

✓ 변경 내용 검토

● 명령 세부 정보

○ 요약

## 첫 번째 실행 명령

상태 ○ 실행 중 📅 Oct 5, 10:13:38 PM

중단

✓ 0/1 단계가 완료되었습니다.

☒ Show All Steps

☐ Show Only Failed Steps

☐ Show Only Running Steps

➤ ○ **Run a set of services for the first time**

0/1 단계가 완료되었습니다.

나머지 기본값으로 설정하고  
완료

Cluster 1

✓ Kafka

작업 ▼

상태

인스턴스

구성

명령

차트 라이브러리

감사

쿼리 링크 ▼

data reten

필터

▼ 범위

**Data Retention Time**

log.retention.ms

Kafka Broker Default Group ↺

10

분 ▼

# Flume 설정

# 구성 파일 분석

## ❖ Sources, channels, sinks 이름 지정

- `SmartCar_Agent.sources = SmartCarInfo_SpoolSource`
- `SmartCar_Agent.channels = SmartCarInfo_Channel`
- `SmartCar_Agent.sinks = SmartCarInfo_LoggerSink`

## ❖ 소스 타입 설정

- `SmartCar_Agent.sources.SmartCarInfo_SpoolSource.type = spooldir`

## ❖ 스푼 위치 지정

- `SmartCar_Agent.sources.SmartCarInfo_SpoolSource.spoolDir = /home/pilot-pjt/working/car-batch-log`

## ❖ 삭제 정책 설정

- `SmartCar_Agent.sources.SmartCarInfo_SpoolSource.deletePolicy = immediate`

## ❖ Channel에 전송할 batch 크기

- `SmartCar_Agent.sources.SmartCarInfo_SpoolSource.batchSize = 1000`



# 구성 파일 분석

## ❖ Channel 설정

- `SmartCar_Agent.channels.SmartCarInfo_Channel.type = memory`

## ❖ Channel에 저장할 최대 event 수

- `SmartCar_Agent.channels.SmartCarInfo_Channel.capacity = 100000`

## ❖ Source에서 가져오거나 sink에 전달할 최대 event 수

- `SmartCar_Agent.channels.SmartCarInfo_Channel.transactionCapacity = 10000`

## ❖ 정상 동작 확인 위해 flume으로 수집된 데이터를 logge로 전송

- `SmartCar_Agent.sinks.SmartCarInfo_LoggerSink.type = logger`

## ❖ SpoolSource와 Memory channel 연결

- `SmartCar_Agent.sources.SmartCarInfo_SpoolSource.channels = SmartCarInfo_Channel`

## ❖ Logger와 Memory channel 연결

- `SmartCar_Agent.sinks.SmartCarInfo_LoggerSink.channel = SmartCarInfo_Channel`

# interceptors

## ❖ 이해하기

- Source와 Channel 중간에서 데이터 가공하는 역할
- 플럼 Source에서 유입되는 데이터 중 일부를 수정/추가/가공/정제 등
- 플럼 데이터 전송 단위 : Event = Header + Body
- interceptors는 Header 특정값 추가, Body 데이터 가공

## ❖ 변수선언

- `SmartCar_Agent.sources.SmartCarInfo_SpoolSource.interceptors = filterInterceptor`

## ❖ 정규 표현식을 이용한 필터링

- `SmartCar_Agent.sources.SmartCarInfo_SpoolSource.interceptors.filterInterceptor.type = regex_filter`

## ❖ 14자리 날짜 형식으로 시작하는 데이터

- `SmartCar_Agent.sources.SmartCarInfo_SpoolSource.interceptors.filterInterceptor.regex = ^\d{14}`

## ❖ 제외된 이벤트 처리


- `SmartCar_Agent.sources.SmartCarInfo_SpoolSource.interceptors.filterInterceptor.excludeEvents = false`

# 카프카 연동

- ❖ 외부 수행 명령 결과를 flume으로 가져와 수집
  - SmartCar\_Agent.sources.DriverCarInfo\_TailSource.type = exec
  - SmartCar\_Agent.sources.DriverCarInfo\_TailSource.command = tail -F /home/pilot-pjt/working/driver-realtime-log/SmartCarDriverInfo.log
  - SmartCar\_Agent.sources.DriverCarInfo\_TailSource.restart = true
  - SmartCar\_Agent.sources.DriverCarInfo\_TailSource.batchSize = 1000
- ❖ 카프카 연동
  - SmartCar\_Agent.sinks.DriverCarInfo\_KafkaSink.type = org.apache.flume.sink.kafka.KafkaSink
  - SmartCar\_Agent.sinks.DriverCarInfo\_KafkaSink.topic = SmartCar-Topic
  - SmartCar\_Agent.sinks.DriverCarInfo\_KafkaSink.brokerList = server02.hadoop.com:9092
  - SmartCar\_Agent.sinks.DriverCarInfo\_KafkaSink.requiredAcks = 1
  - SmartCar\_Agent.sinks.DriverCarInfo\_KafkaSink.batchSize = 1000

# 구성 파일 등록

Cluster 1

 **Flume** 작업 ▾

상태    인스턴스    **구성**    명령

검색

시스템 그룹

Flume(서비스 전체)

flume

Agent 이름

Agent Default Group ↺

SmartCar\_Agent

구성 파일

Agent Default Group ↺

```
SmartCar_Agent.sources = SmartCarInfo_SpoolSource DriverCarInfo_TailSource
SmartCar_Agent.channels = SmartCarInfo_Channel DriverCarInfo_Channel
SmartCar_Agent.sinks = SmartCarInfo_LoggerSink DriverCarInfo_KafkaSink

SmartCar_Agent.sources.SmartCarInfo_SpoolSource.type = spooldir
```

# Kafka 동작 확인

- ❖ [root@server02 /]# kafka-console-producer --broker-list server02.hadoop.com:9092 -topic SmartCar-Topic
- ❖ >20/10/11 13:01:43 INFO clients.Metadata: Cluster ID: djul1ODDShmRdbl\_Ufhxhg
- ❖ test message
- ❖ >
  
- ❖ [root@server02 /]# kafka-console-consumer --bootstrap-server server02.hadoop.com:9092 --topic SmartCar-Topic --partition 0 --from-beginning
- ❖ >20/10/11 13:01:43 INFO clients.Metadata: Cluster ID: djul1ODDShmRdbl\_Ufhxhg
- ❖ test message

# 수집 기능 테스트

# 자동차 시뮬레이터 동작

- ❖ [root@server02 working]java -cp bigdata.smartcar.loggen-1.0.jar com.wikibook.bigdata.smartcar.loggen.CarLogMain 20200101 3&
- ❖ [root@server02 working]# java -cp bigdata.smartcar.loggen-1.0.jar com.wikibook.bigdata.smartcar.loggen.DriverLogMain 20200101 3 &
- ❖ [root@server02 working]# cd SmartCar/
- ❖ [root@server02 SmartCar]# cat SmartCarStatusInfo\_20200101.txt
- ❖ [root@server02 SmartCar]# cd ..
- ❖ [root@server02 working]# cd driver-realtime-log/
- ❖ [root@server02 driver-realtime-log]# tail -f SmartCarDriverInfo.log

## 분석 정보 이동


- ❖ [root@server02 driver-realtime-log]# cd ..
- ❖ [root@server02 working]# cd SmartCar/
- ❖ [root@server02 SmartCar]# ls
- ❖ SmartCarStatusInfo\_20200101.txt
- ❖ [root@server02 SmartCar]# mv SmartCarStatusInfo\_20200101.txt /home/pilot-pjt/working/car-batch-log/



# 서비스 재시작

Cloudera Manager 클러스터 ▾ 호스트 ▾

Cluster 1

✓ Flume 작업 ▾ 

상태 인스턴스

상태 테스트

✓ Agent 상태  
양호 상태의 Agent 1개.

시작  
중지  
**재시작**  
롤링 재시작  
역할 인스턴스 추가

Cloudera Manager 클러스터 ▾ 호스트 ▾

Cluster 1

! Kafka 작업 ▾

상태 인스턴스

상태 테스트

시작  
**재시작**  
롤링 재시작

## 수집 기능 점검

- ❖ [root@server02 ~]# tail -f /var/log/flume-ng/flume-cmf-flume-AGENT-server02.hadoop.com.log

```
2020-10-11 16:52:40,749 INFO org.apache.flume.sink.Lc
 31 30 31 32 33 35 39 32 34 2C 45 20190101235924,E }
2020-10-11 16:52:40,749 INFO org.apache.flume.sink.Lc
 31 30 31 32 33 35 39 32 38 2C 45 20190101235928,E }
```

- ❖ [root@server02 ~]# kafka-console-consumer --bootstrap-server server02.hadoop.com:9092 --topic SmartCar-Topic --partition 0

```
20190101003428,B0003,4,0,F,N,37,D04
20190101003426,C0001,2,0,R3,R,10,D06
20190101003438,H0002,1,0,F,N,5,D04
20190101003430,B0003,3,0,F,N,52,D05
```

# 자동차 시뮬레이터 종료

```
[root@server02 ~]# ps -ef | grep smartcar.log
```

```
root    4557    1  1 16:51 ?        00:00:04 java -cp bigdata.smartcar.loggen-  
1.0.jar com.wikibook.bigdata.smartcar.loggen.DriverLogMain 20190101 3
```

```
root    4584    1  3 16:51 ?        00:00:07 java -cp bigdata.smartcar.loggen-  
1.0.jar com.wikibook.bigdata.smartcar.loggen.CarLogMain 20190101 3
```

```
root    5467  5033  0 16:54 pts/1    00:00:00 grep smartcar.log
```

```
[root@server02 ~]# kill -9 4557
```

```
[root@server02 ~]# kill -9 4584
```

```
[root@server02 ~]# ps -ef | grep smartcar.log
```

```
root    5504  5033  0 16:55 pts/1    00:00:00 grep smartcar.log
```