



네이버 최대의 데이터 저장소 운영기 (HBase Locality기반 운영기)



정한룡 NAVER

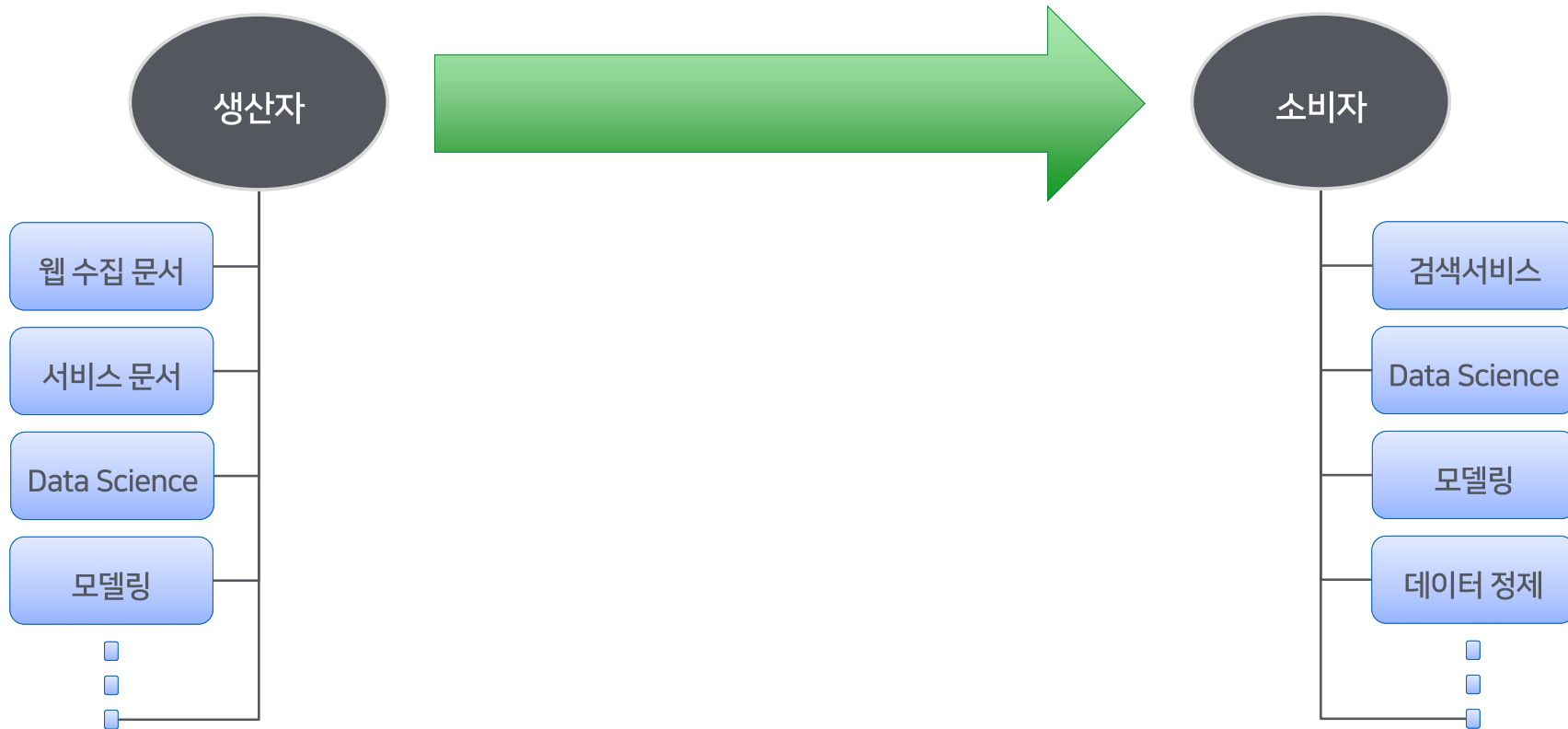
CONTENTS

1. HBase기반 데이터 저장소
2. HBase Locality란
3. HBase Locality를 통한 읽기 성능 높이기
4. 네이버 데이터 저장소에서 HBase Locality를 지키기 위한 운영
5. 네이버 데이터 저장소에서 효과



HBase기반 데이터 저장소

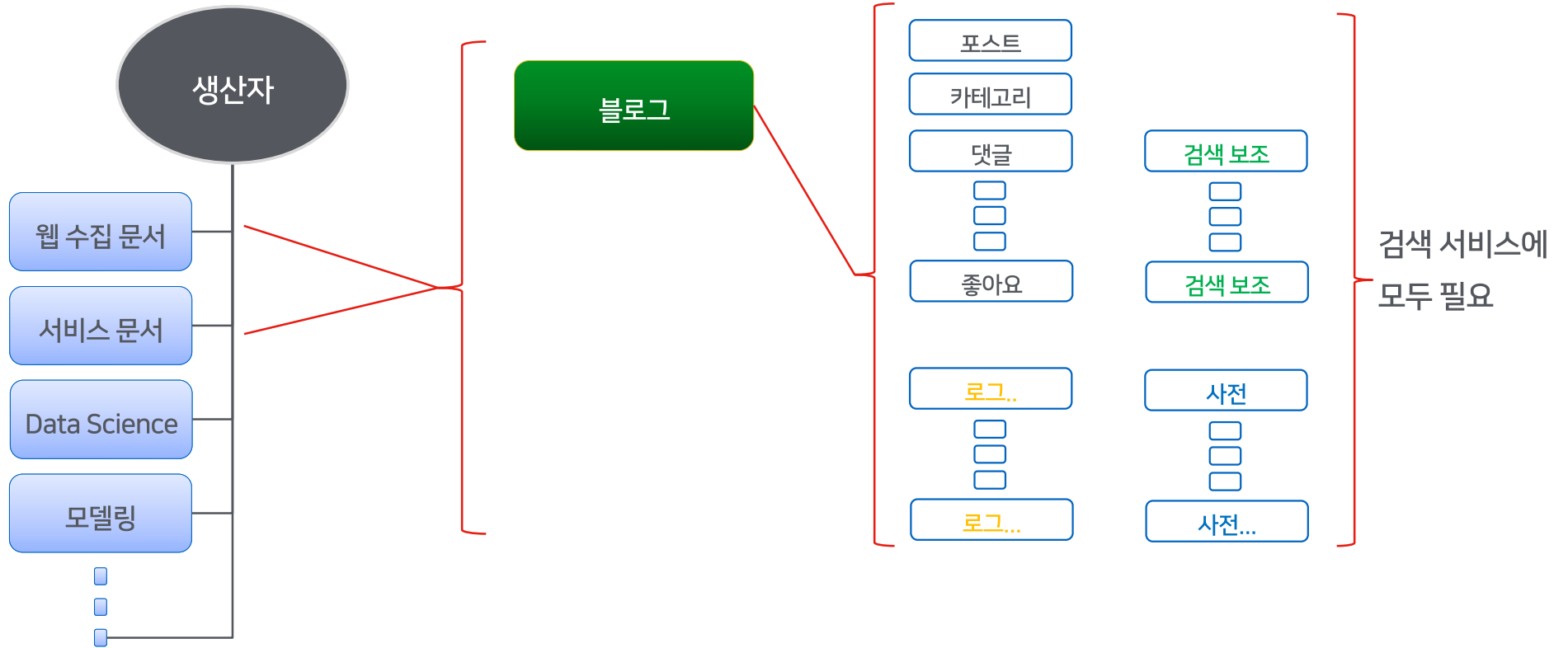
데이터 저장소가 있기 전



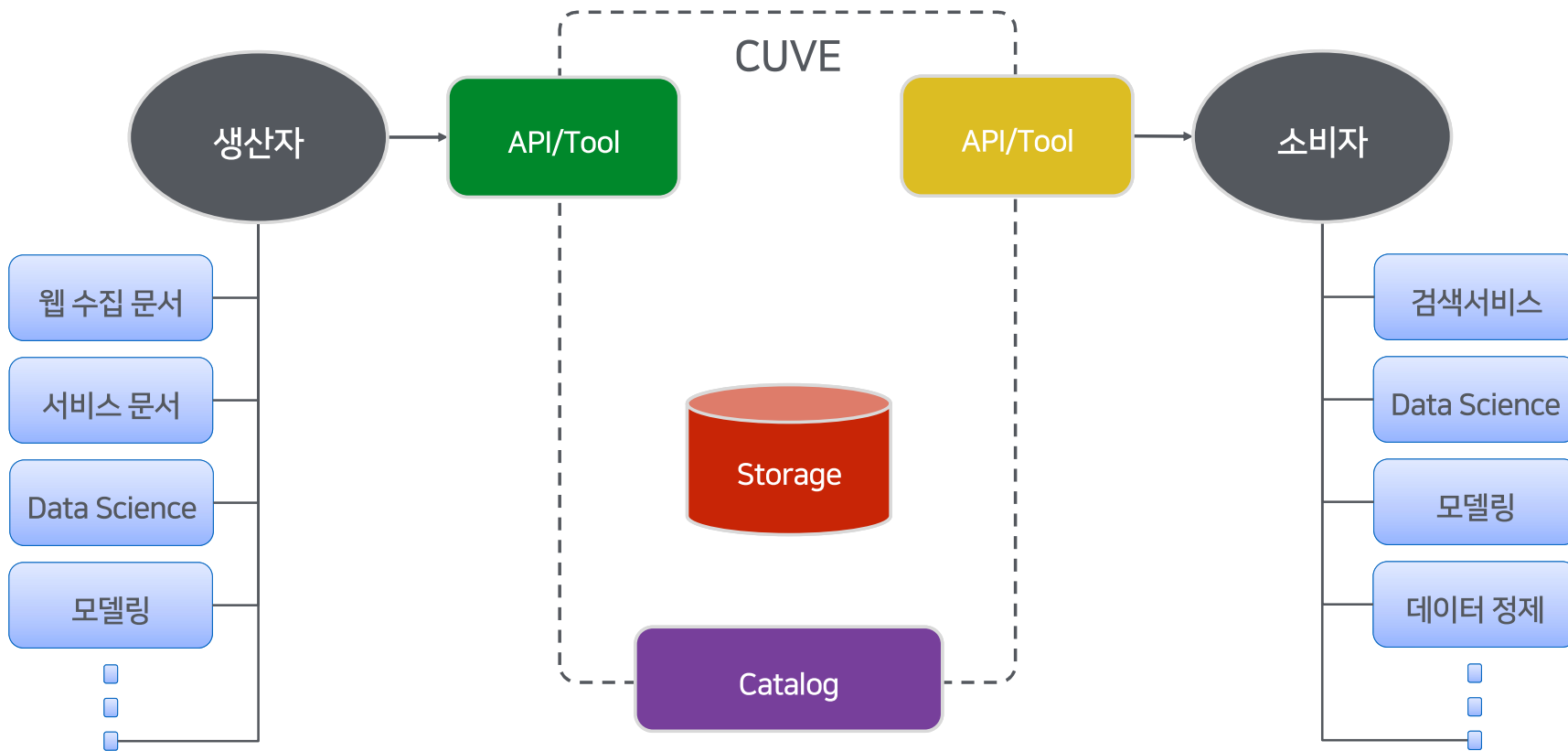
다수의 생산자



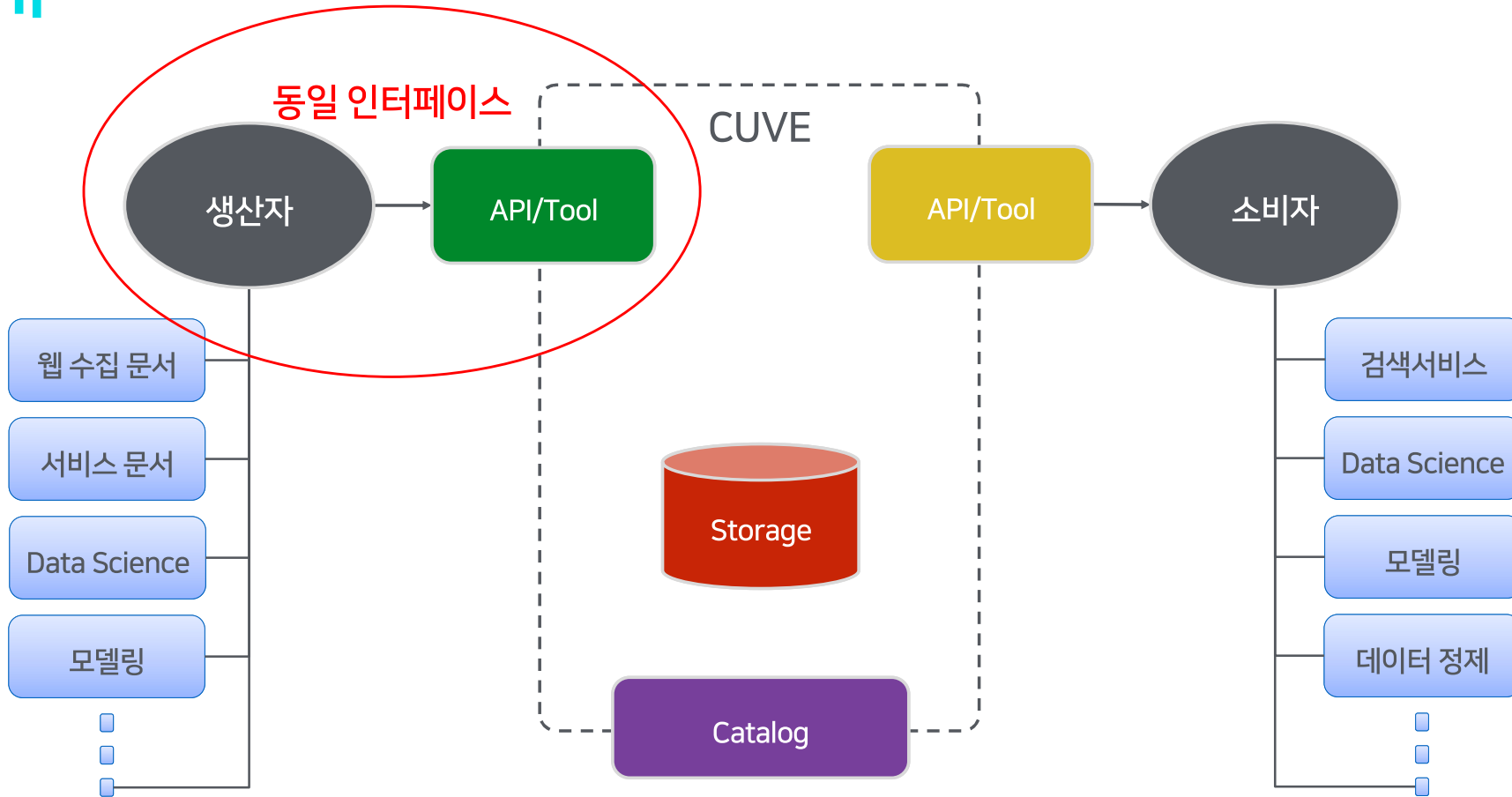
다양한 데이터



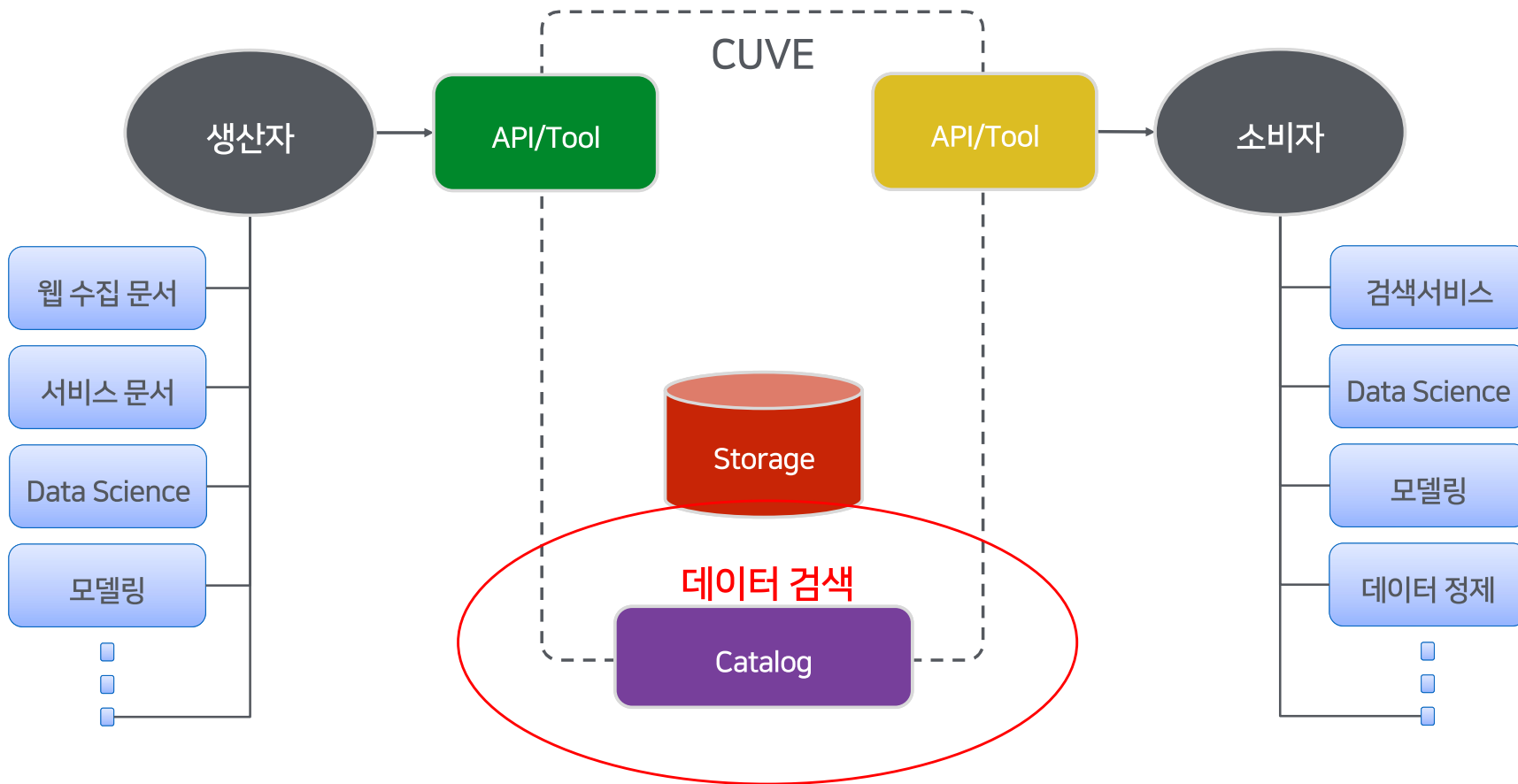
현재



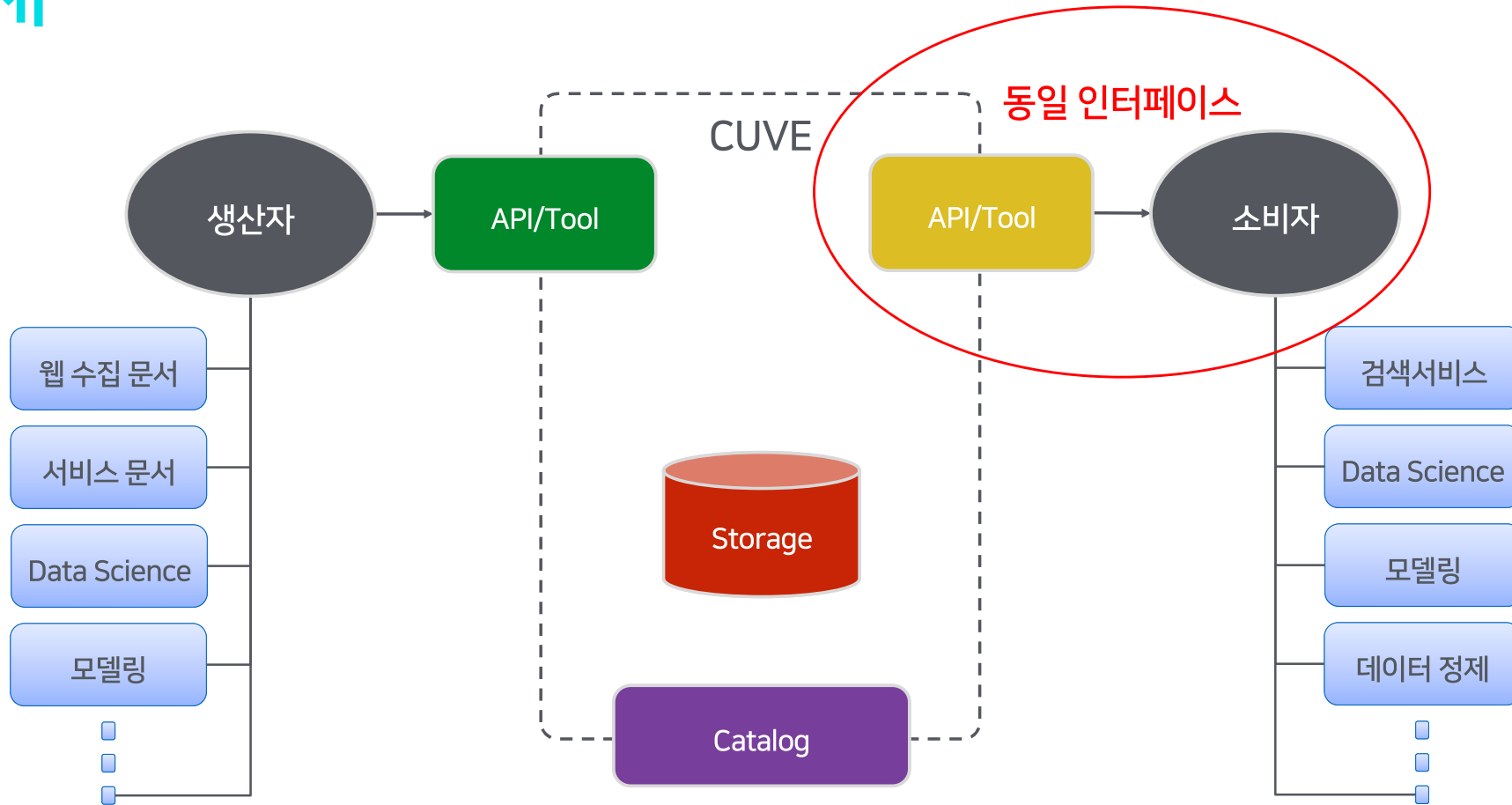
현재



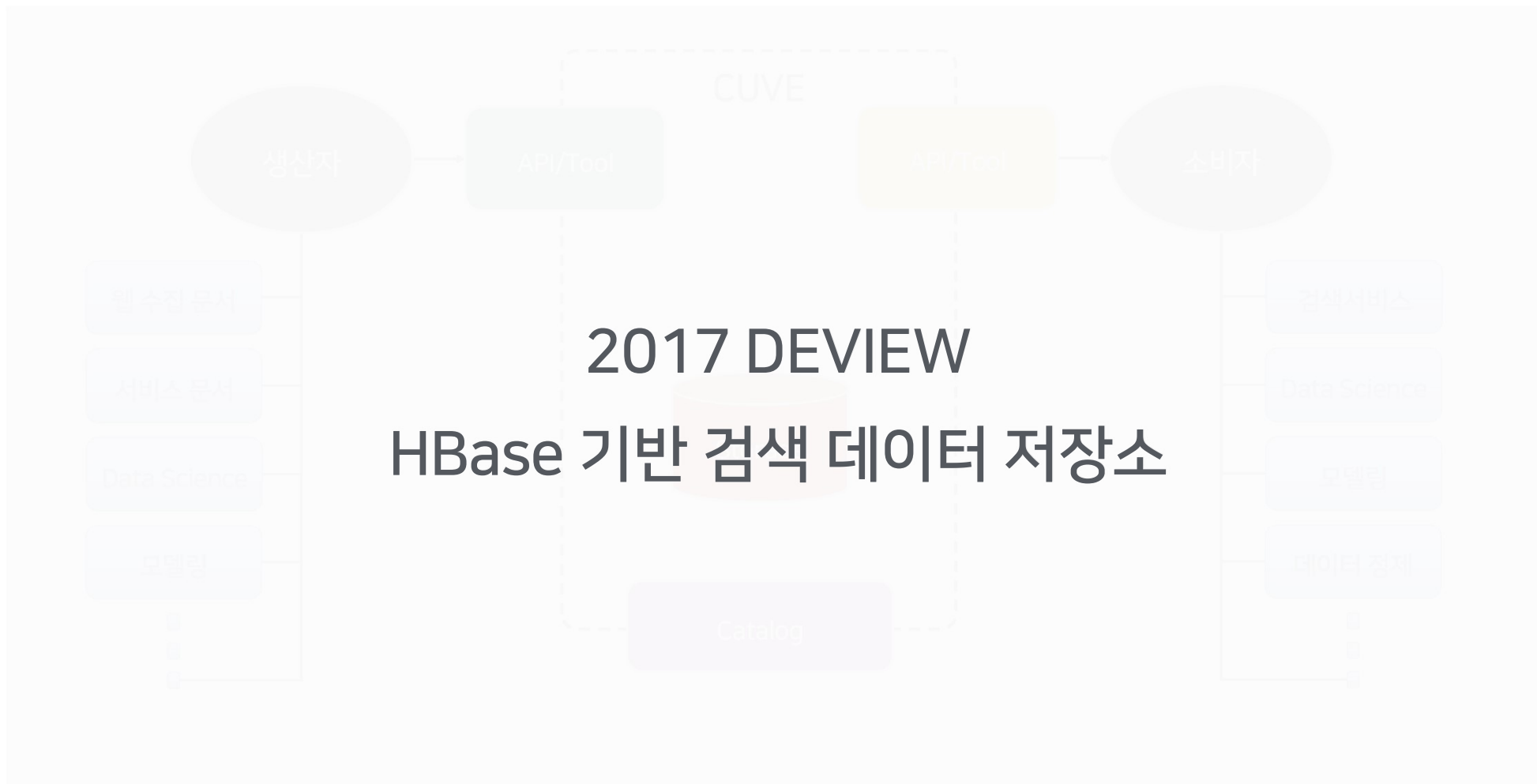
현재



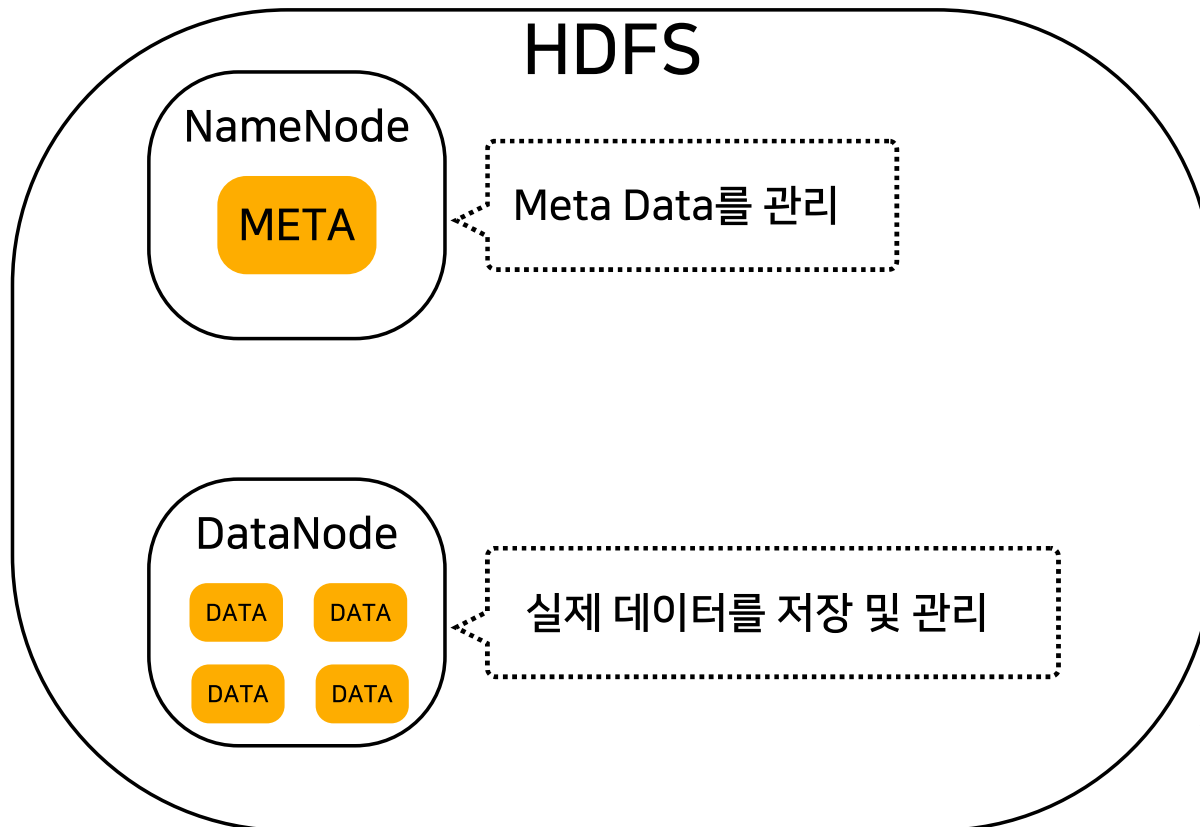
현재



HBase Key구조 및 자세한 건...



HDFS에 대해서

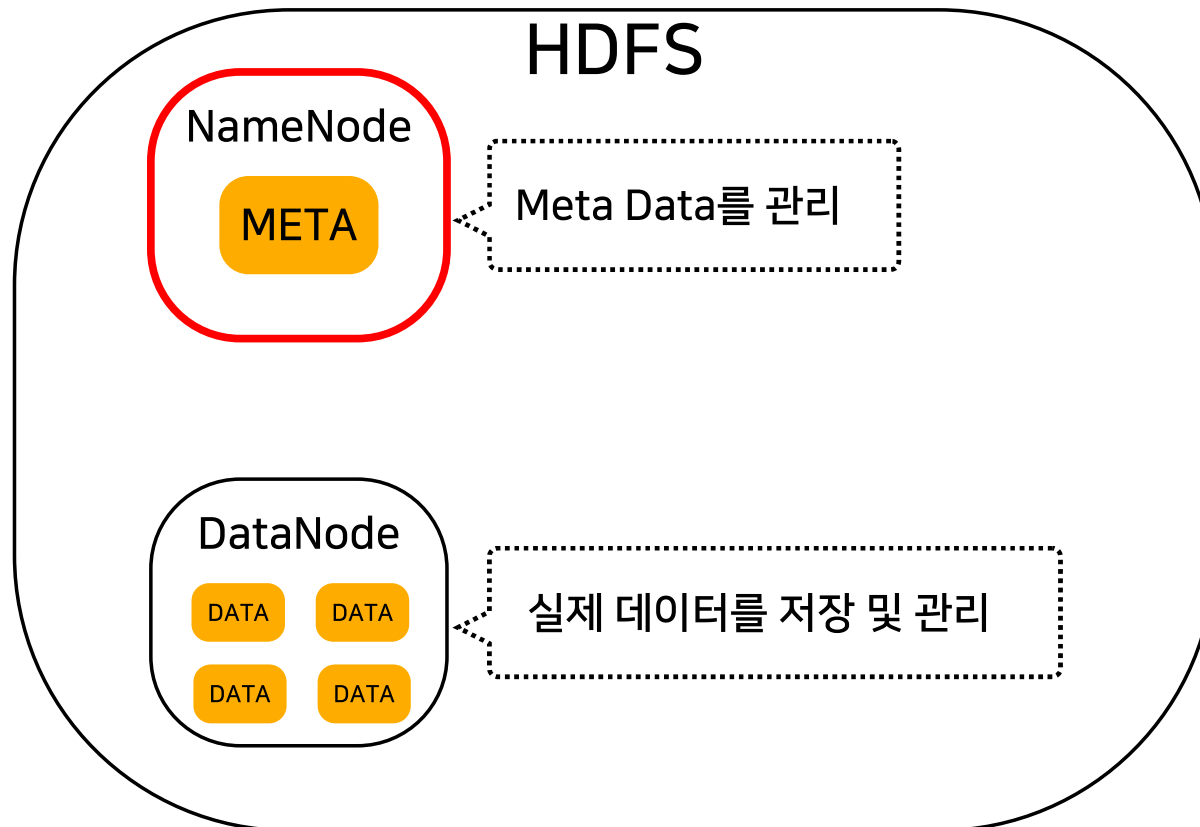


Hadoop Distributed File System

대용량 데이터를 분산된 서버에 저장하기 위하여 개발된 시스템

장애복구 (Replication)
스트리밍의 방식 데이터 접근
대용량 데이터 저장
데이터 무결성

HDFS에 대해서

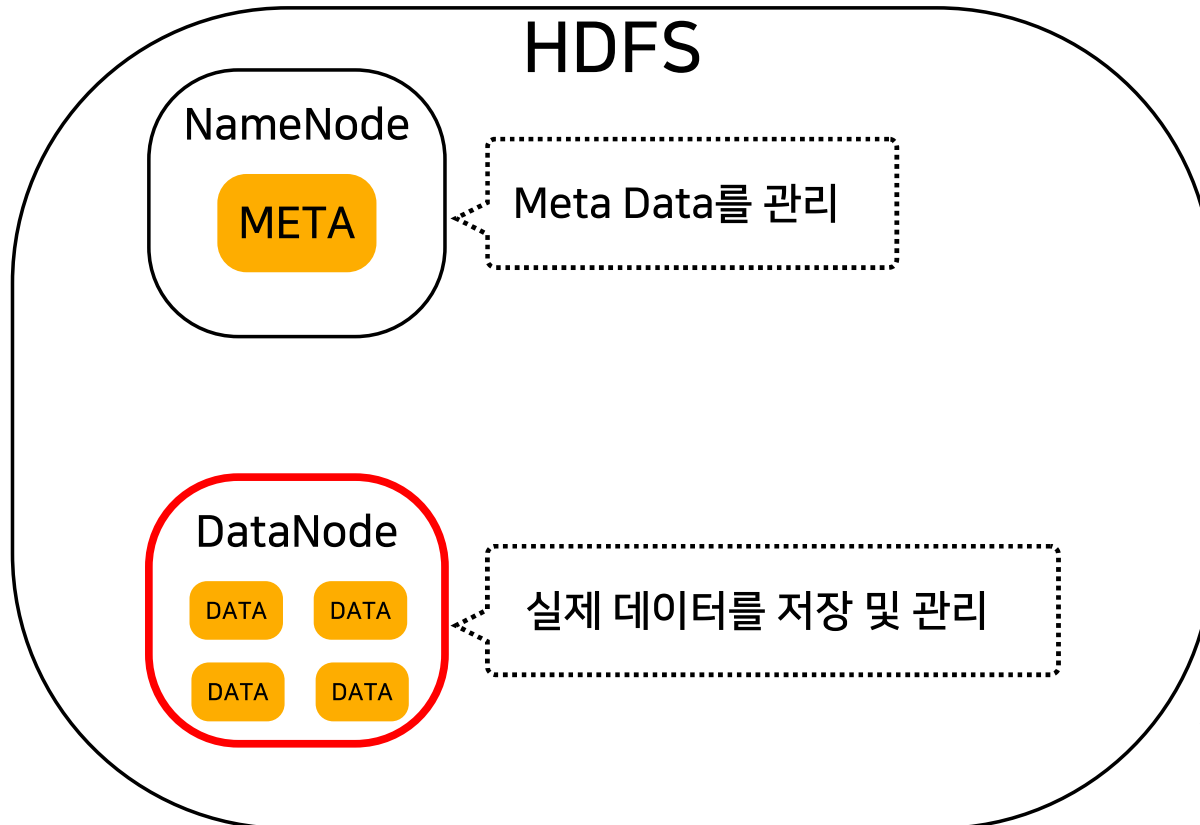


Hadoop Distributed File System

대용량 데이터를 분산된 서버에 저장하기 위하여 개발된 시스템

- 장애복구 (Replication)
- 스트리밍의 방식 데이터 접근
- 대용량 데이터 저장
- 데이터 무결성

HDFS에 대해서

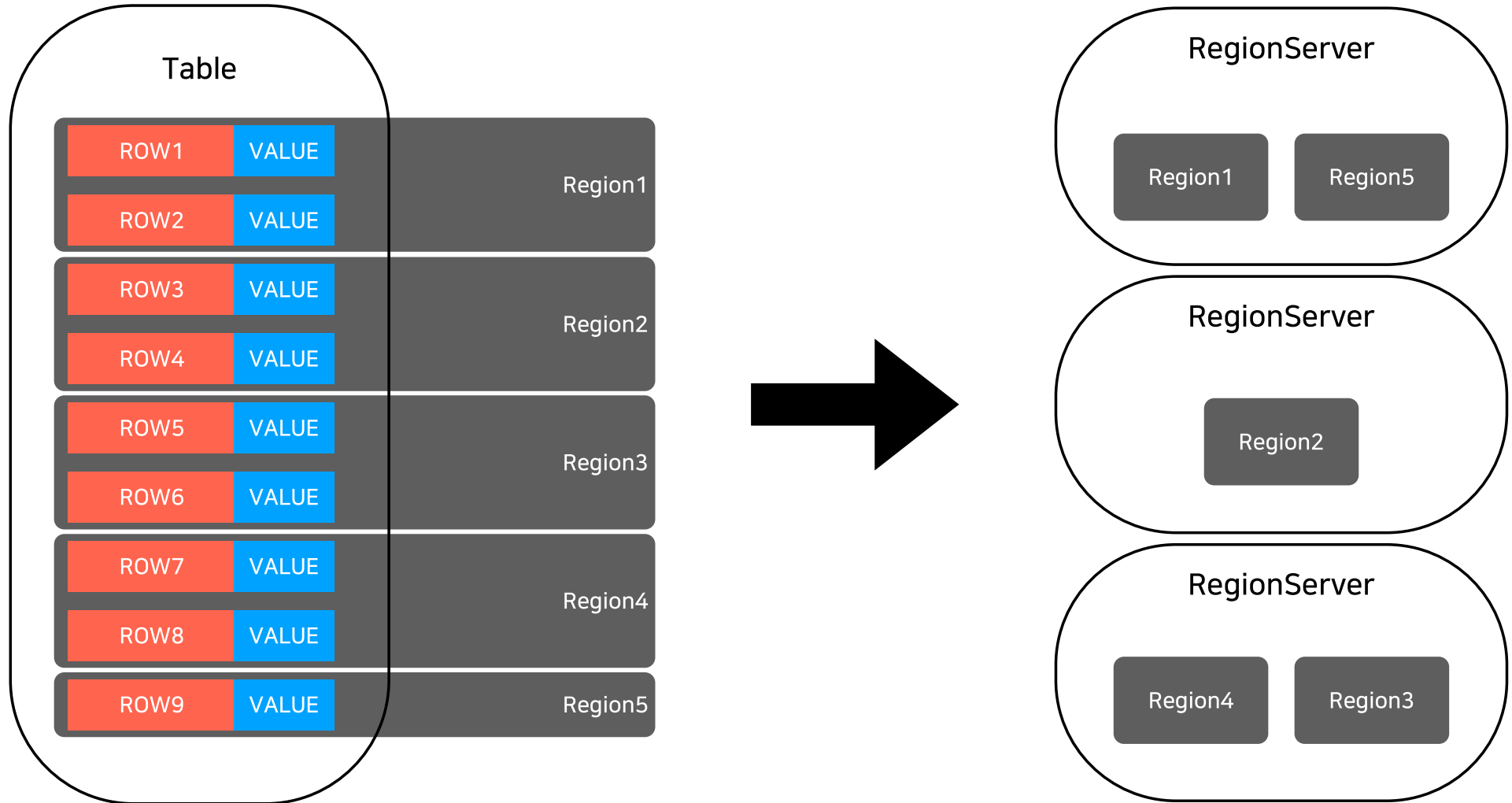


Hadoop Distributed File System

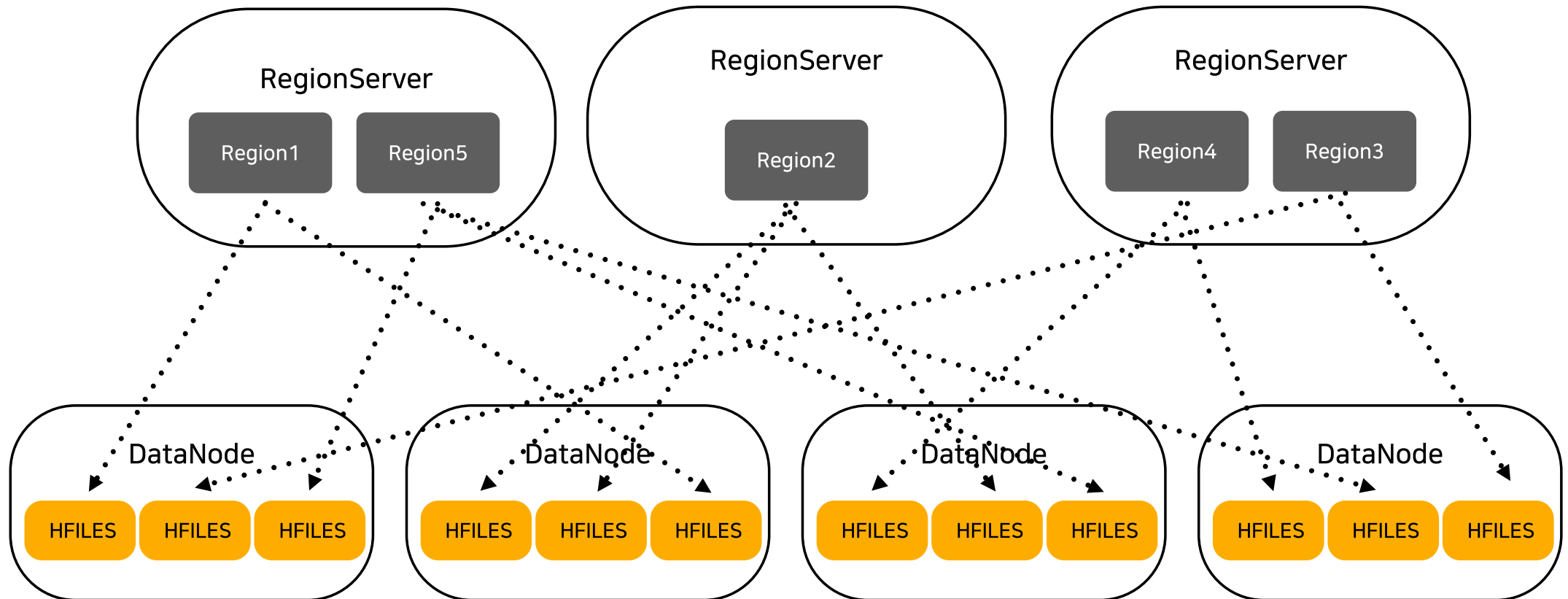
대용량 데이터를 분산된 서버에 저장하기 위하여 개발된 시스템

- 장애복구 (Replication)
- 스트리밍의 방식 데이터 접근
- 대용량 데이터 저장
- 데이터 무결성

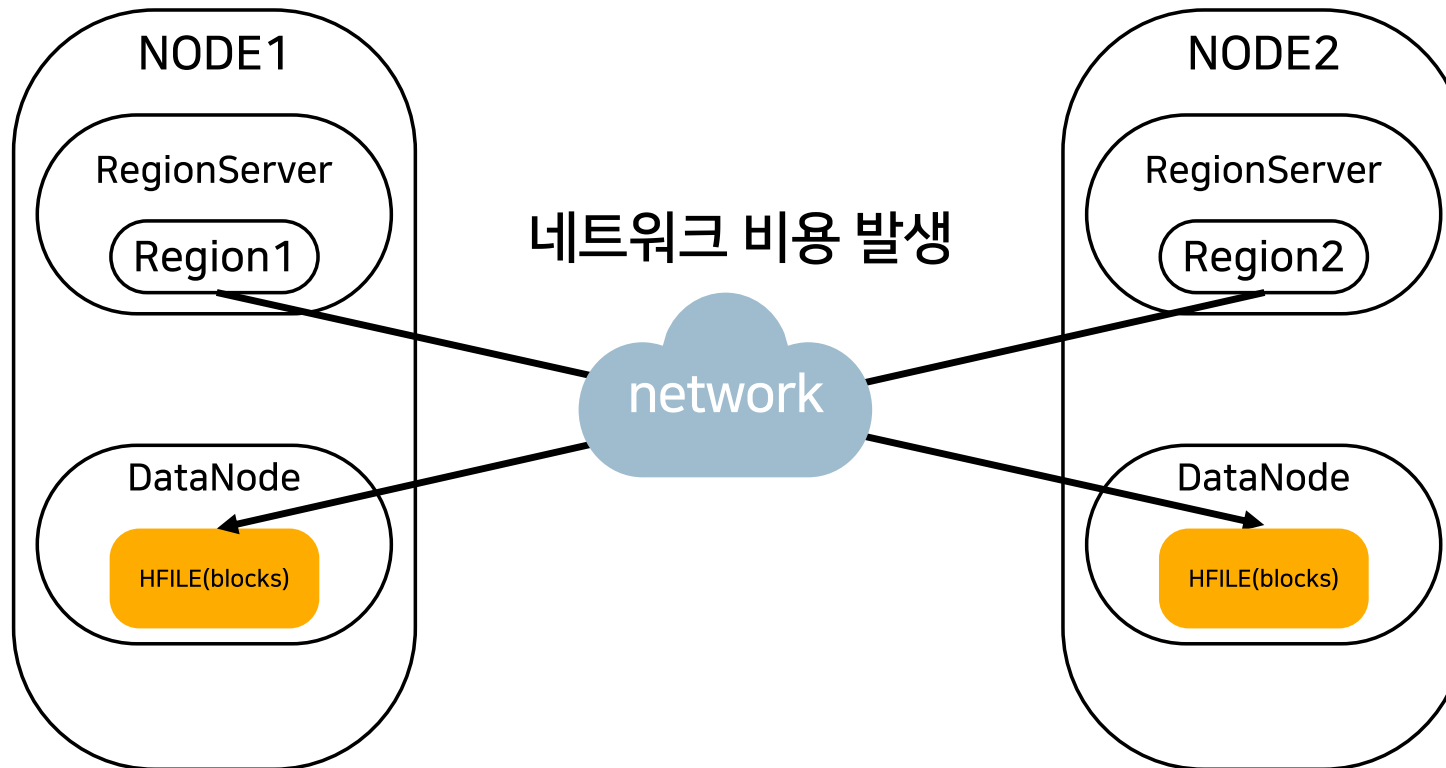
HBase에 대해서



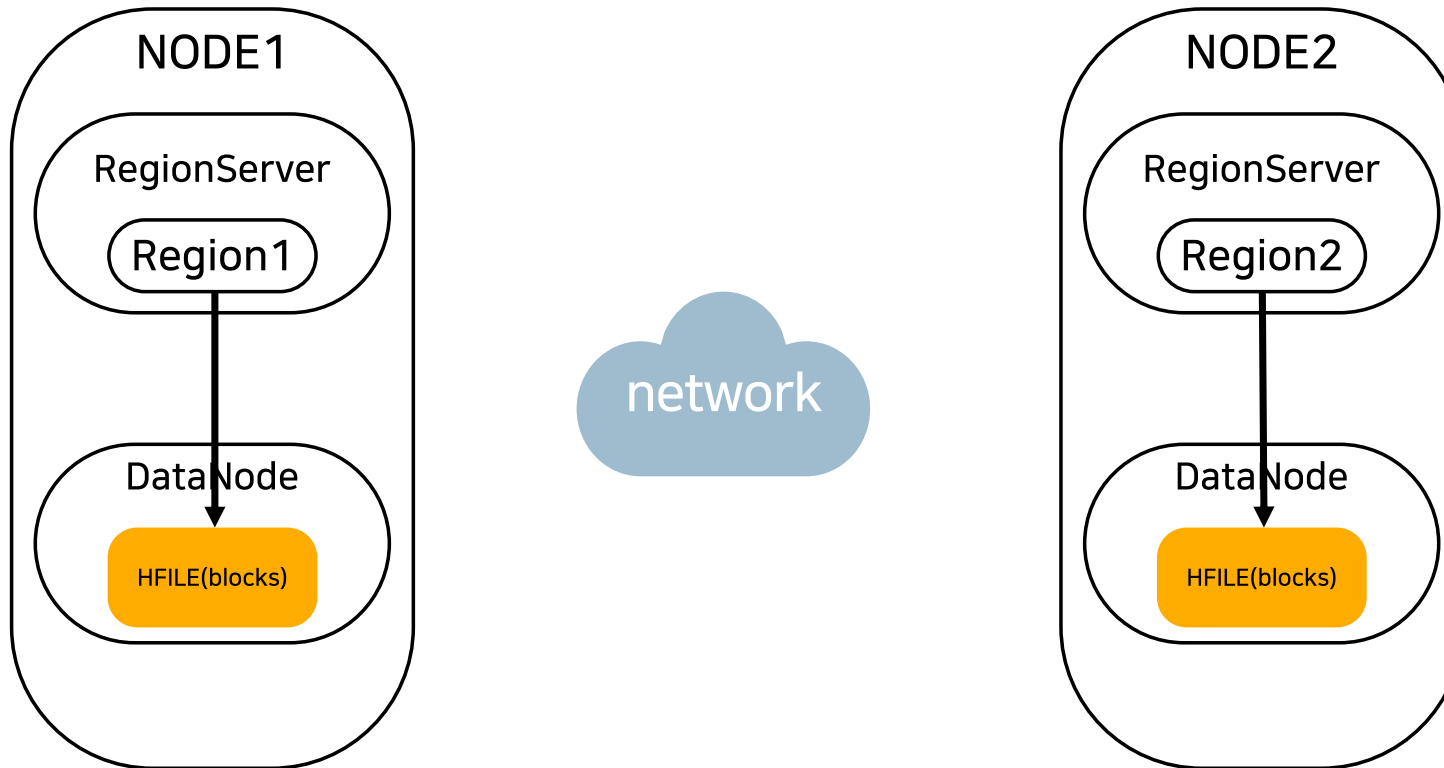
HBase에 대해서



Region - HFILES가 같은 장비에 있다면?



Region - HFILES가 같은 장비에 있다면?



Region - HFILES가 같은 장비에 있다면?

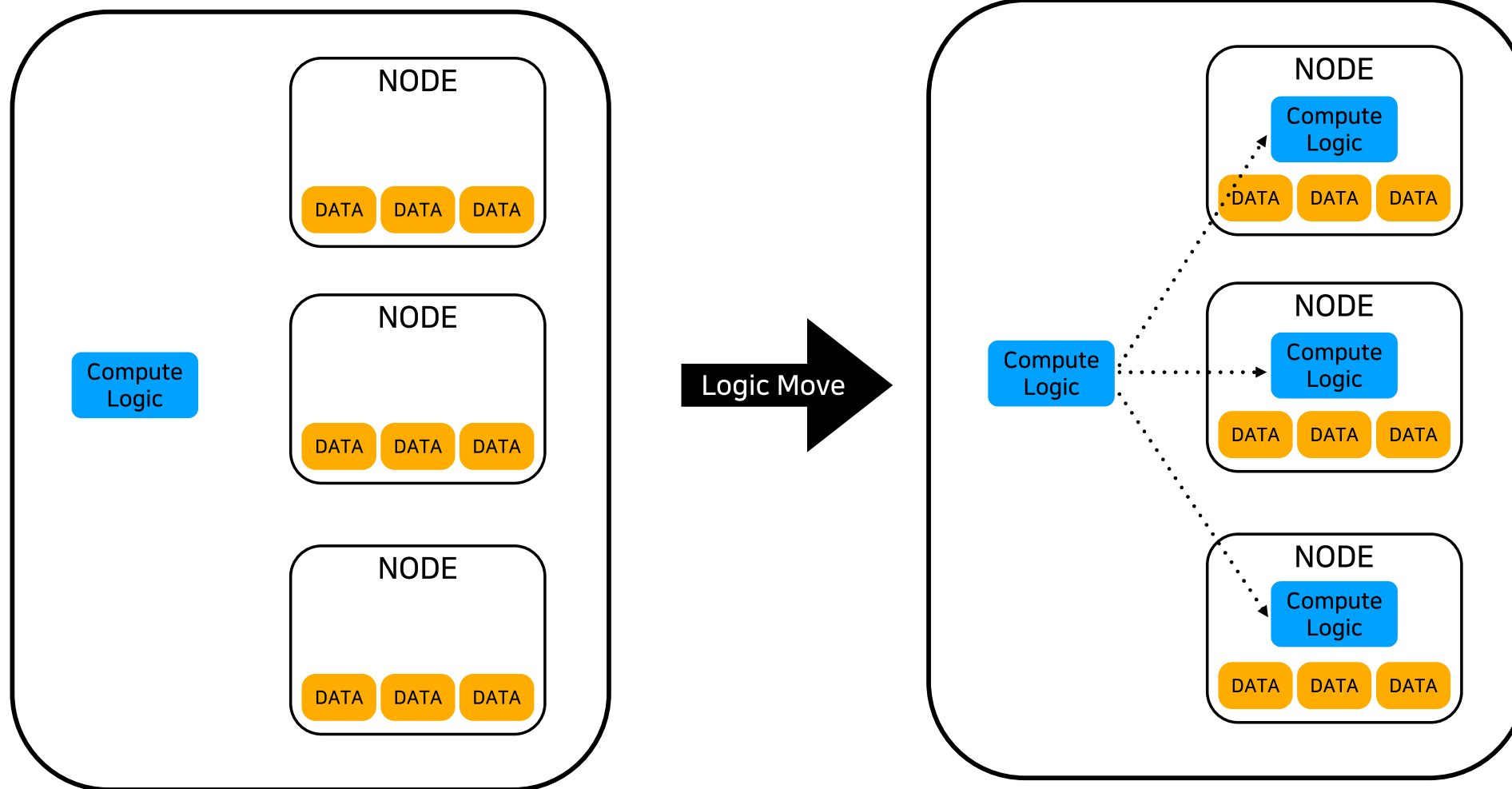


네트워크 비용 절약
읽기 속도 향상



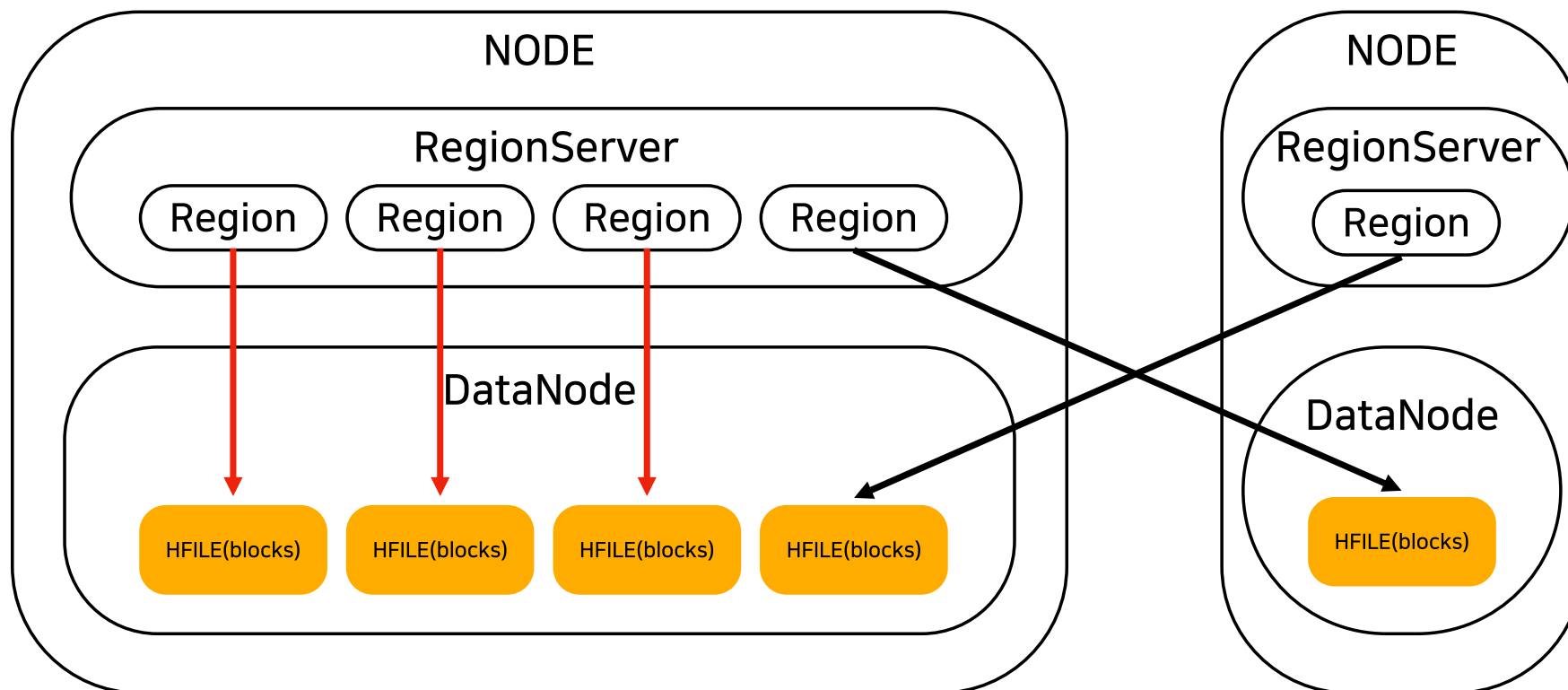
HBase Locality

Data Locality

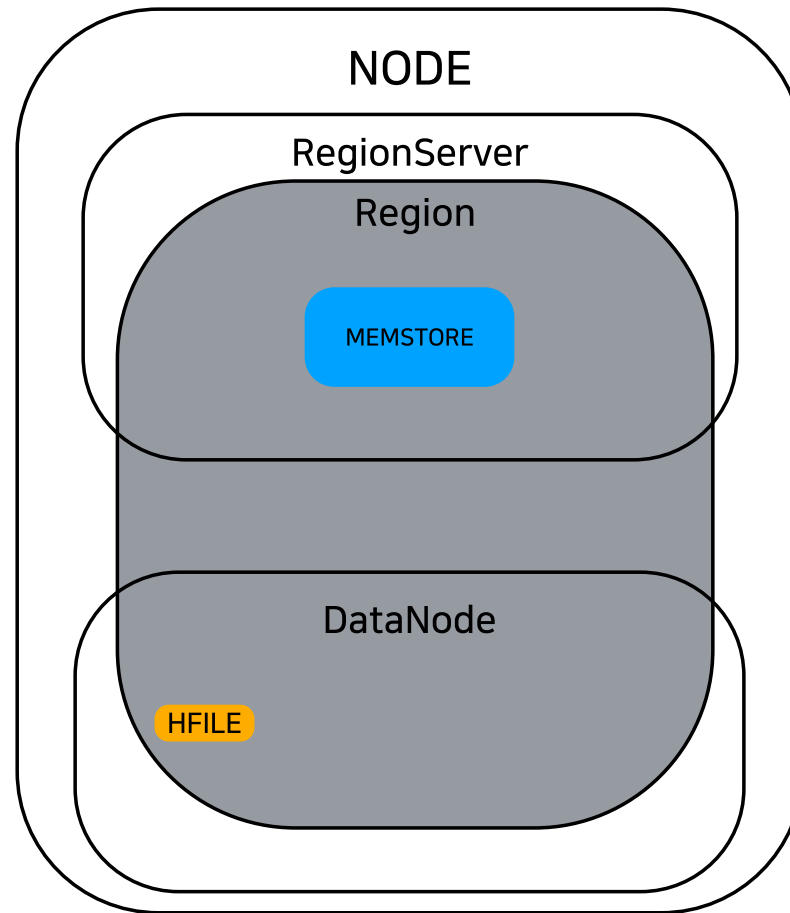


HBase Locality

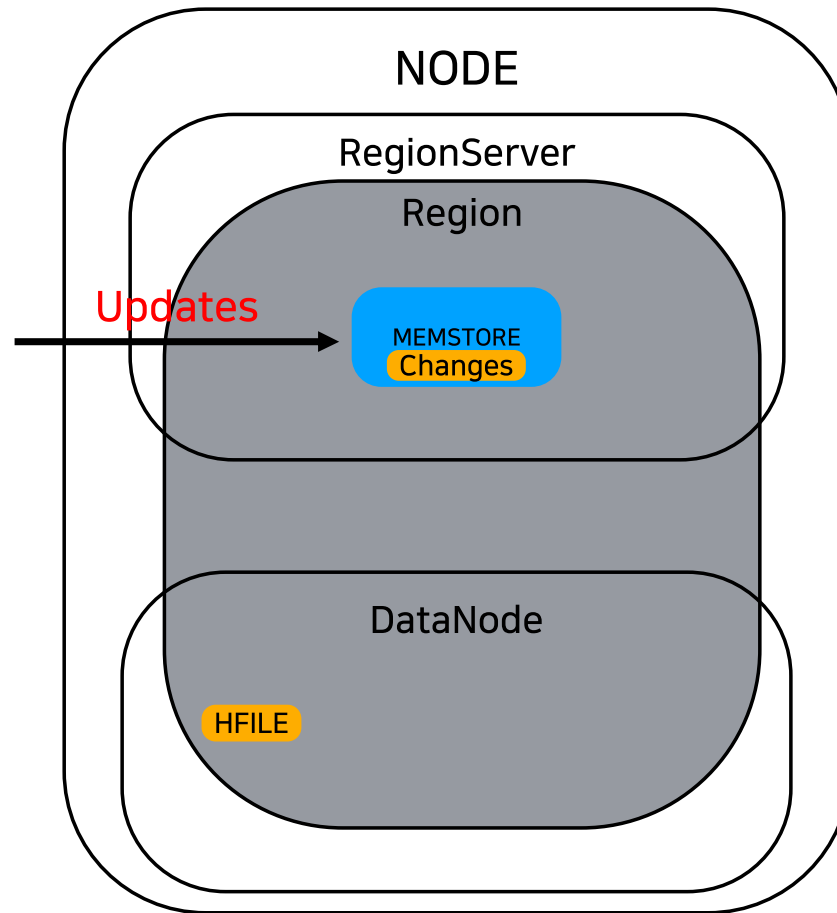
Region의 물리 데이터가 Local DataNode에 있는 것



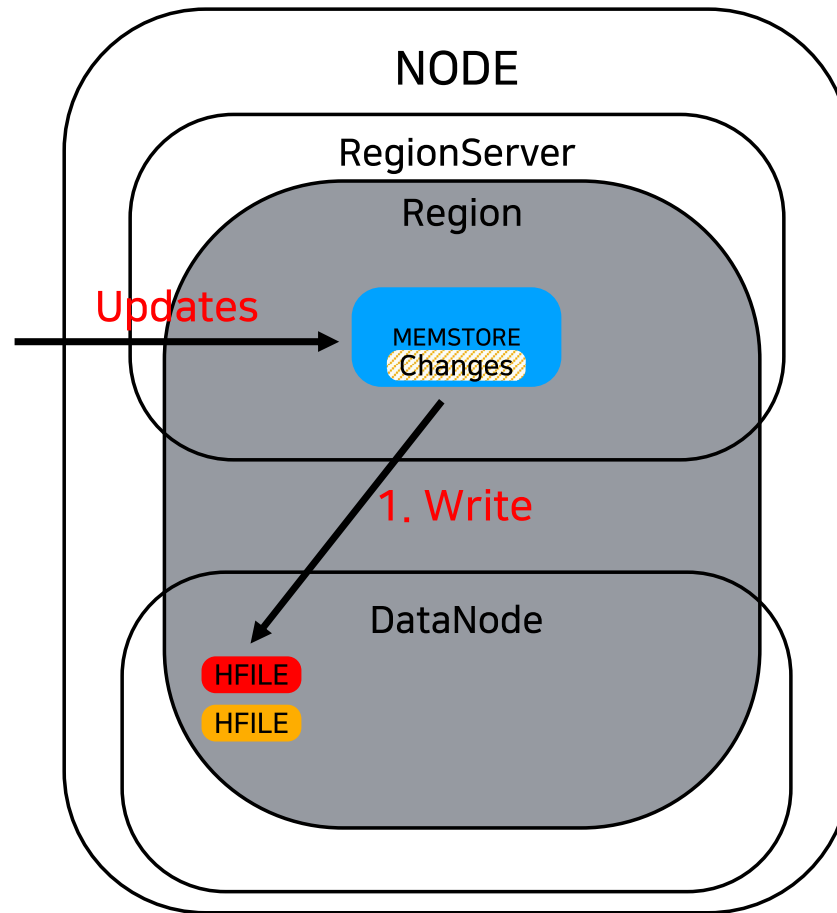
HBase에서 Locality를 올리는 방법



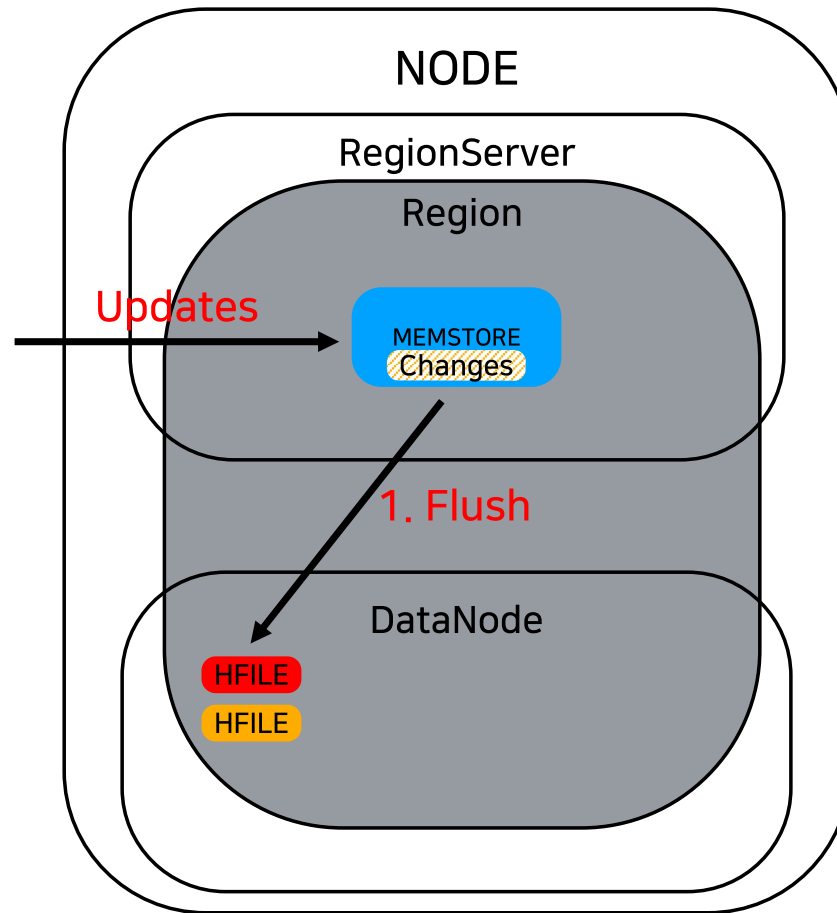
HBase에서 Locality를 올리는 방법



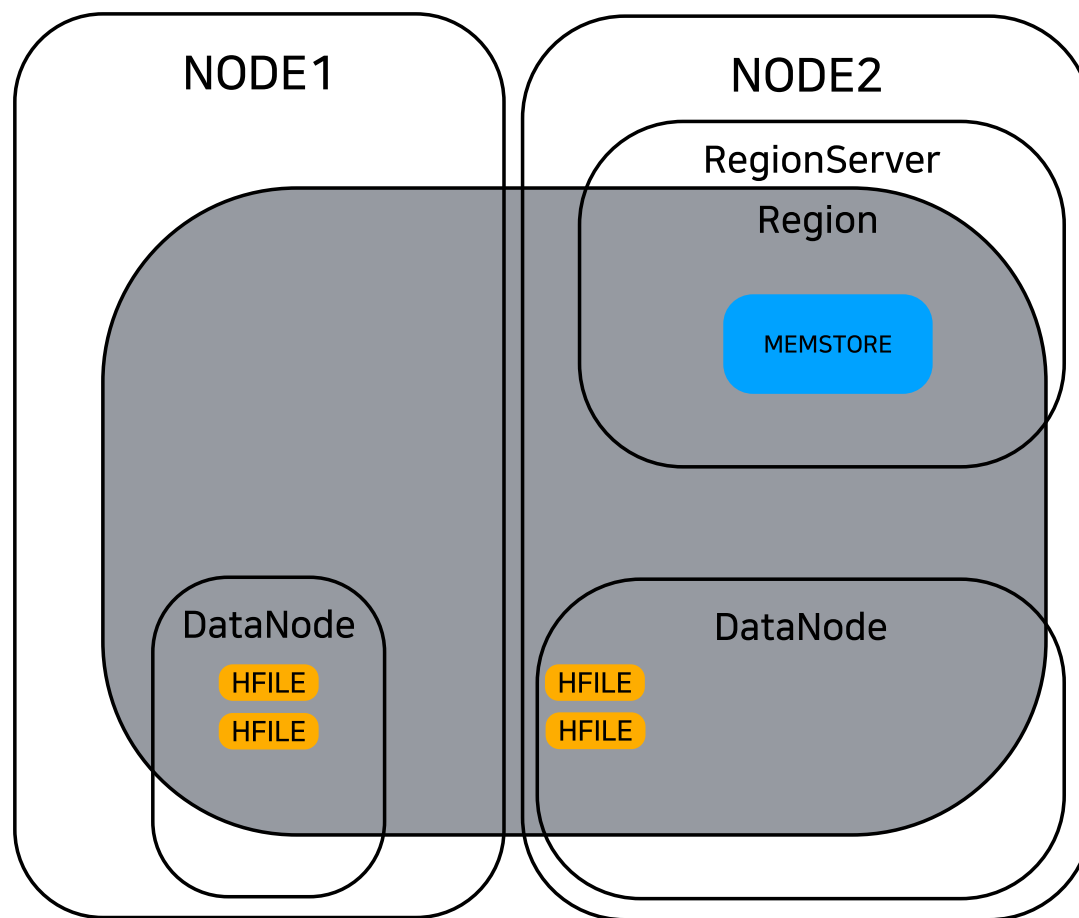
HBase에서 Locality를 올리는 방법



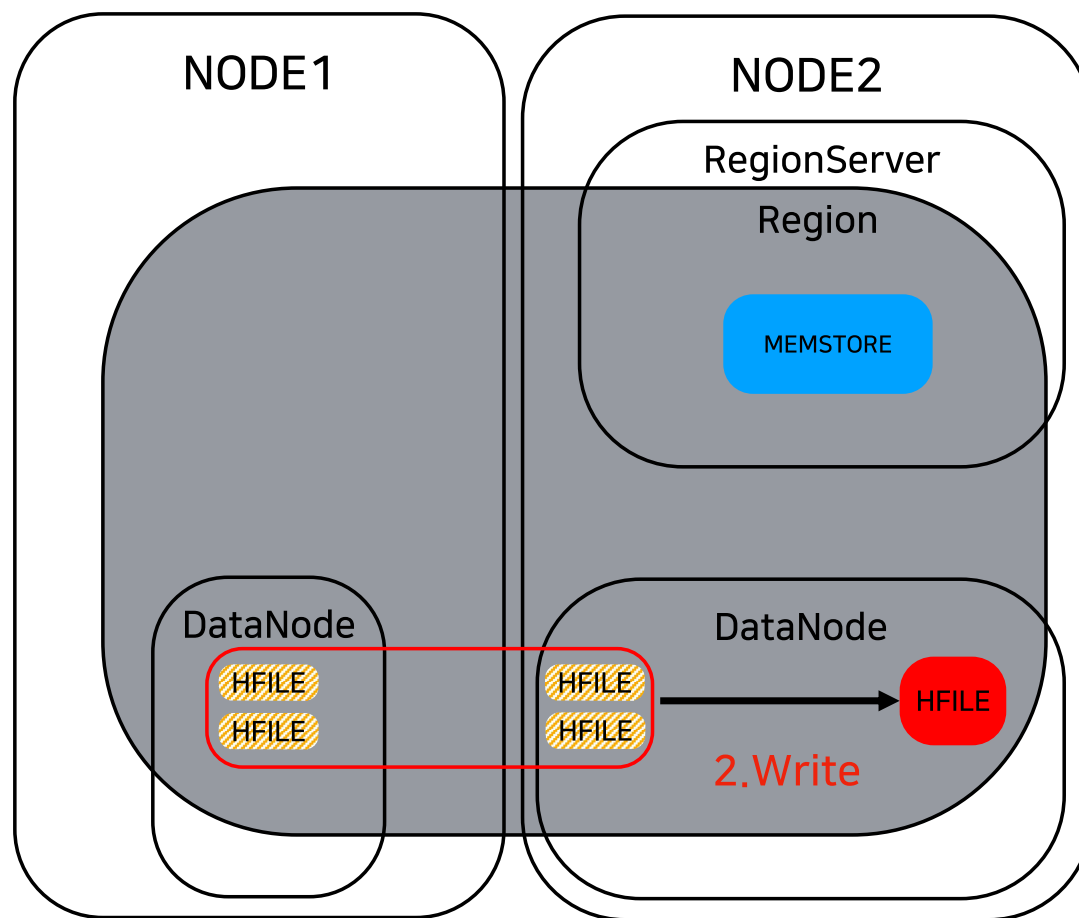
HBase에서 Locality를 올리는 방법



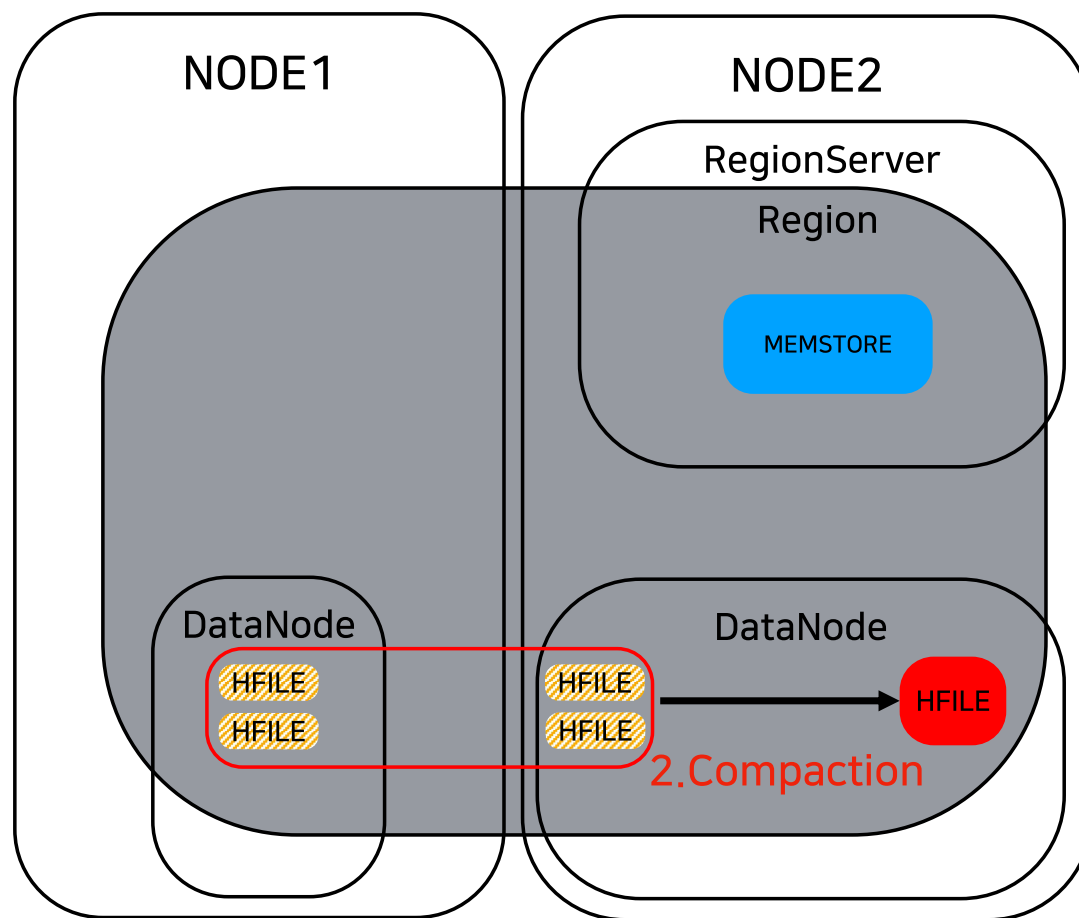
HBase에서 Locality를 올리는 방법



HBase에서 Locality를 올리는 방법

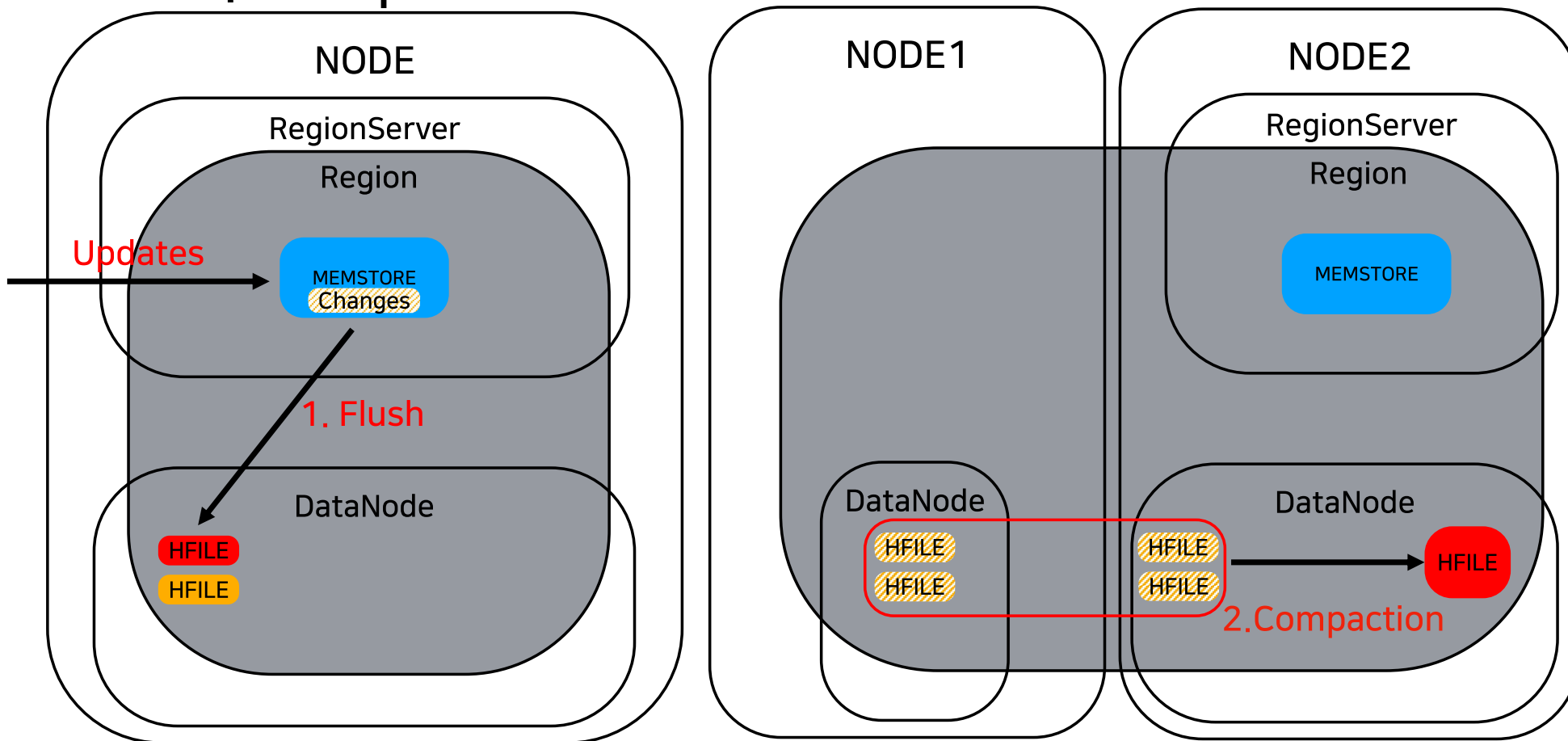


HBase에서 Locality를 올리는 방법



HBase에서 Locality를 올리는 방법

Flush와 Compaction

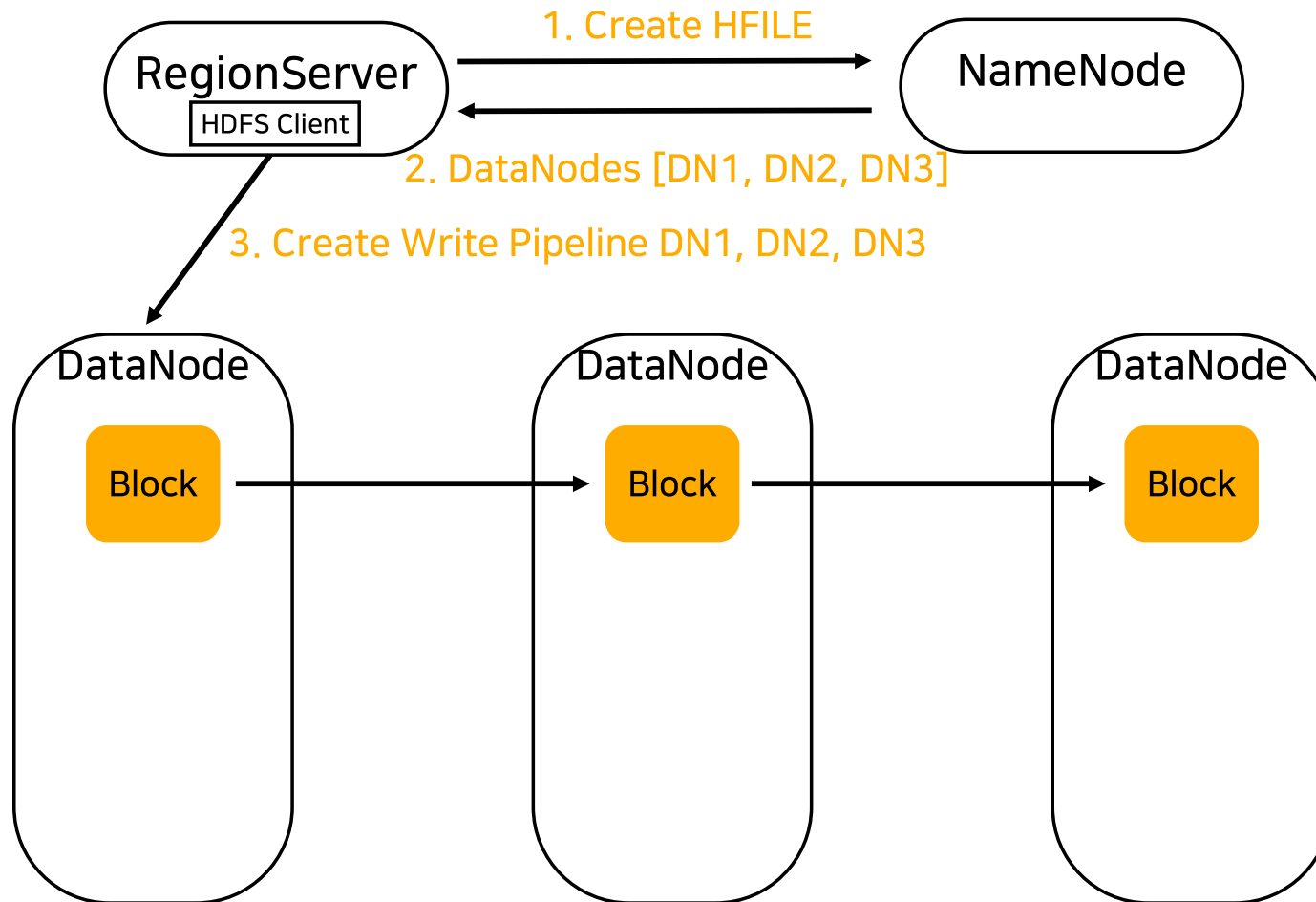


HBase에서 Locality를 올리는 방법

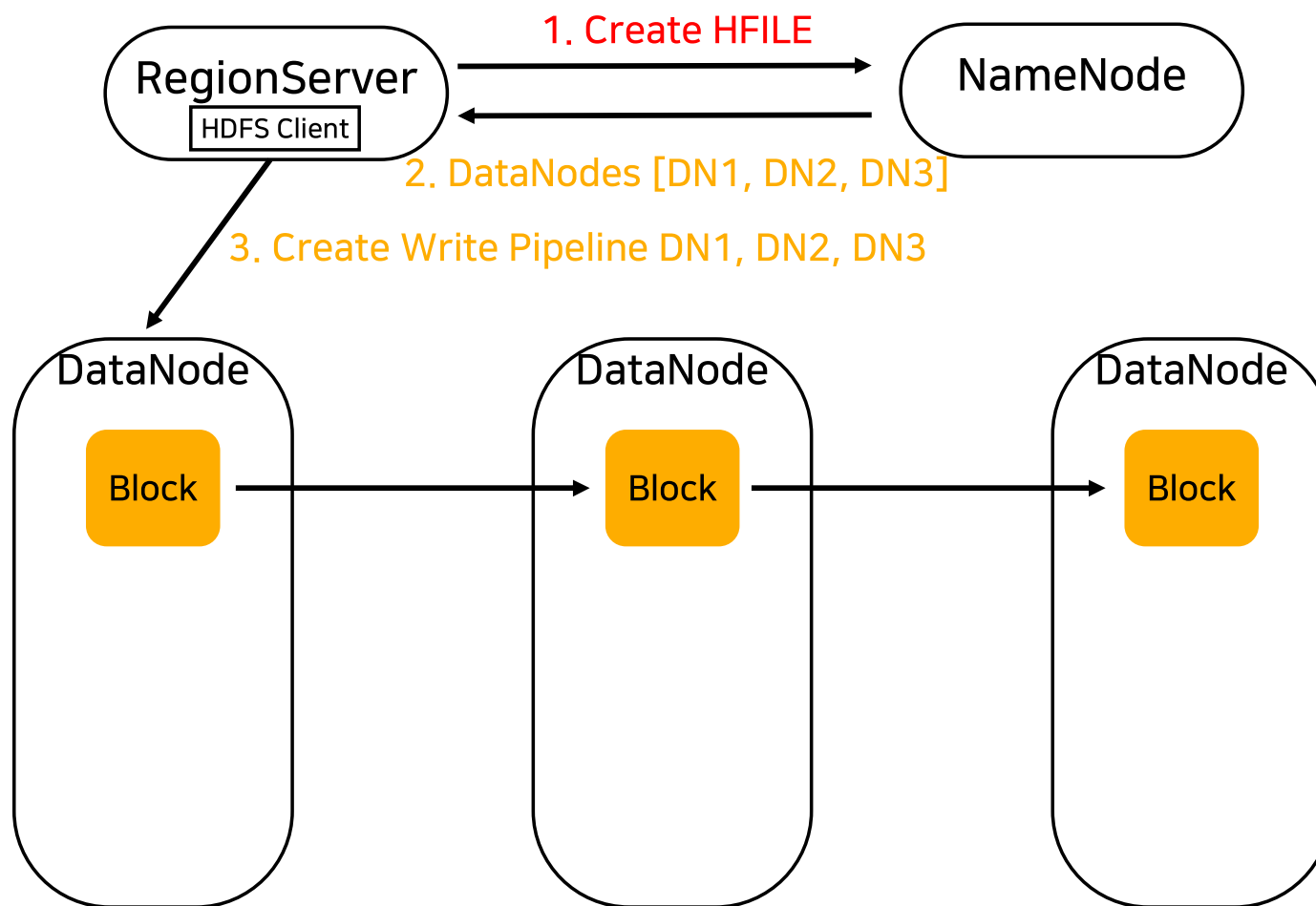
Flush와 Compaction



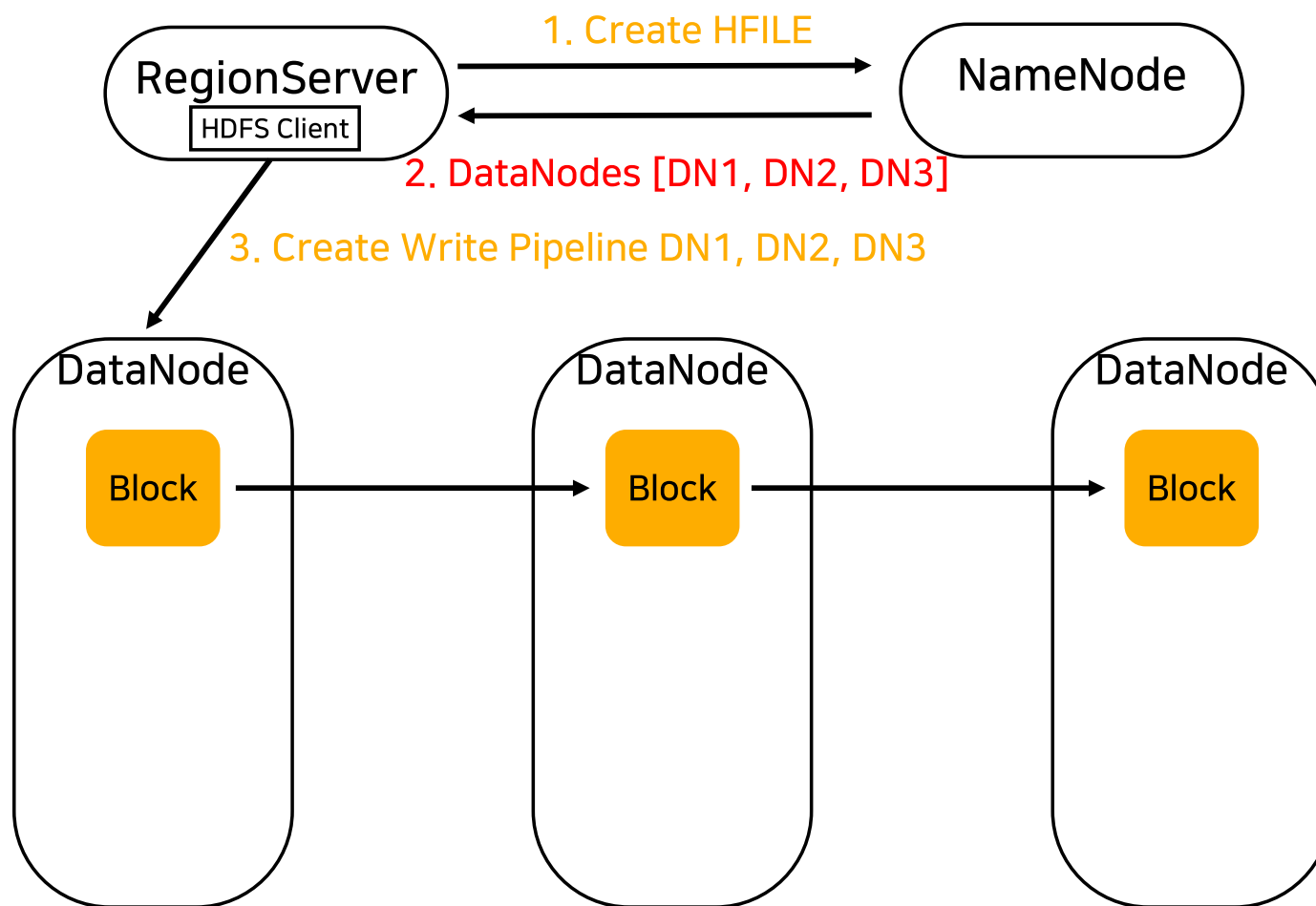
HDFS Write 통해서 Locality를 어떻게?



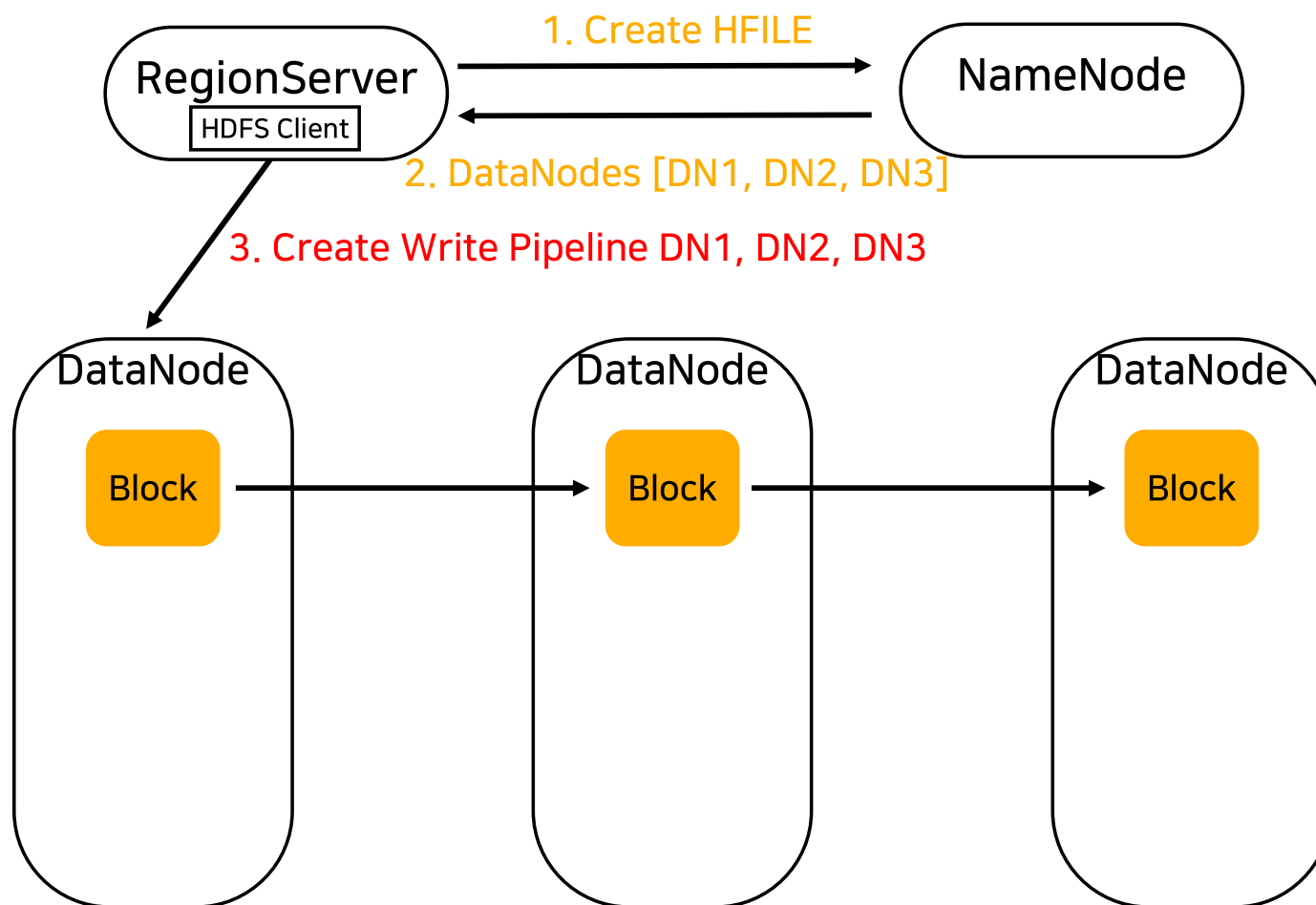
HDFS Write 통해서 Locality를 어떻게?



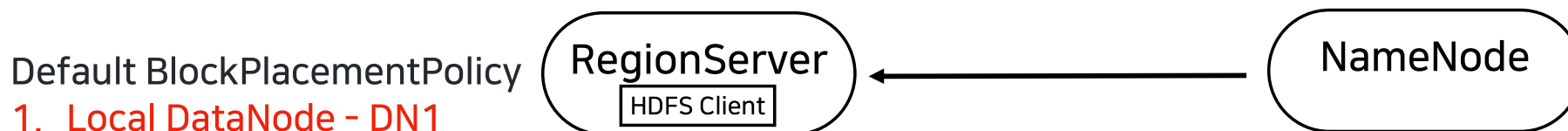
HDFS Write 통해서 Locality를 어떻게?



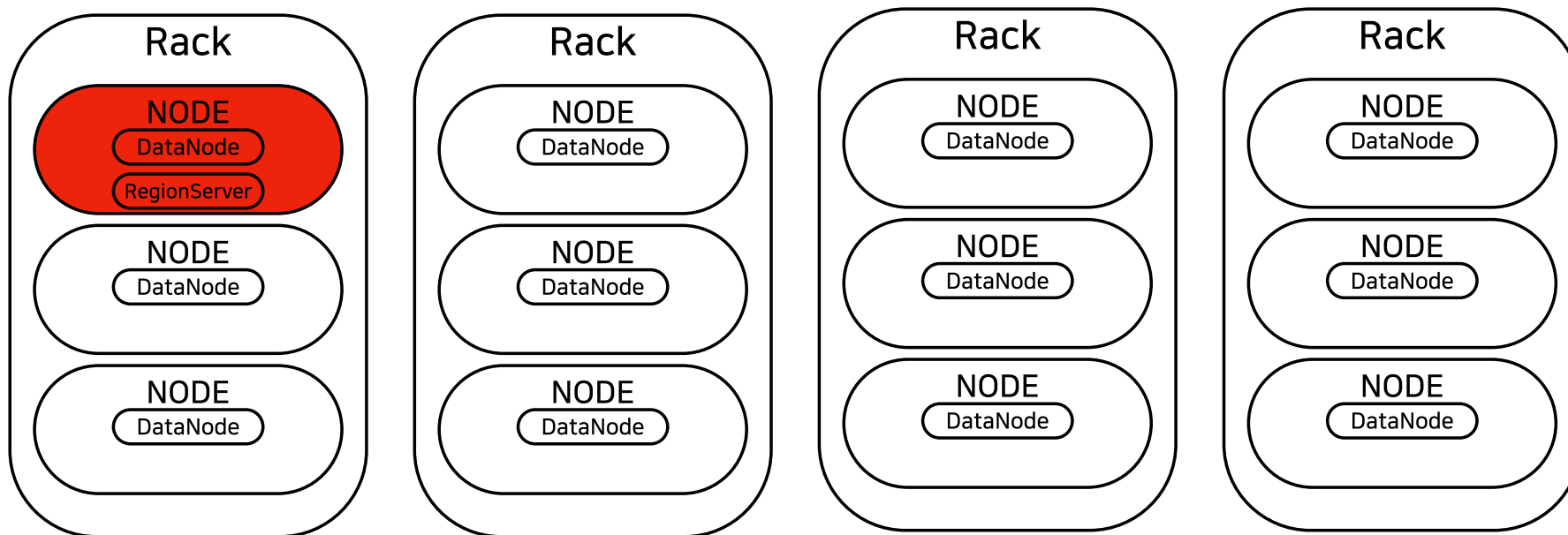
HDFS Write 통해서 Locality를 어떻게?



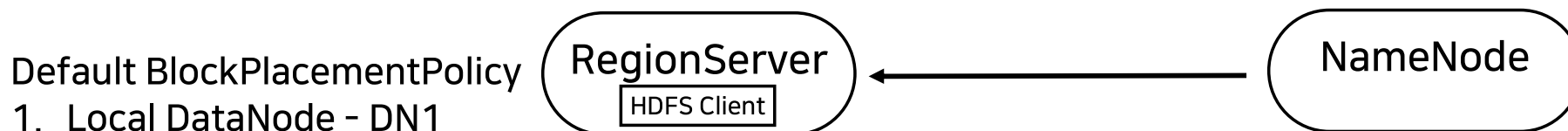
정답은 DataNode 선정 방법!



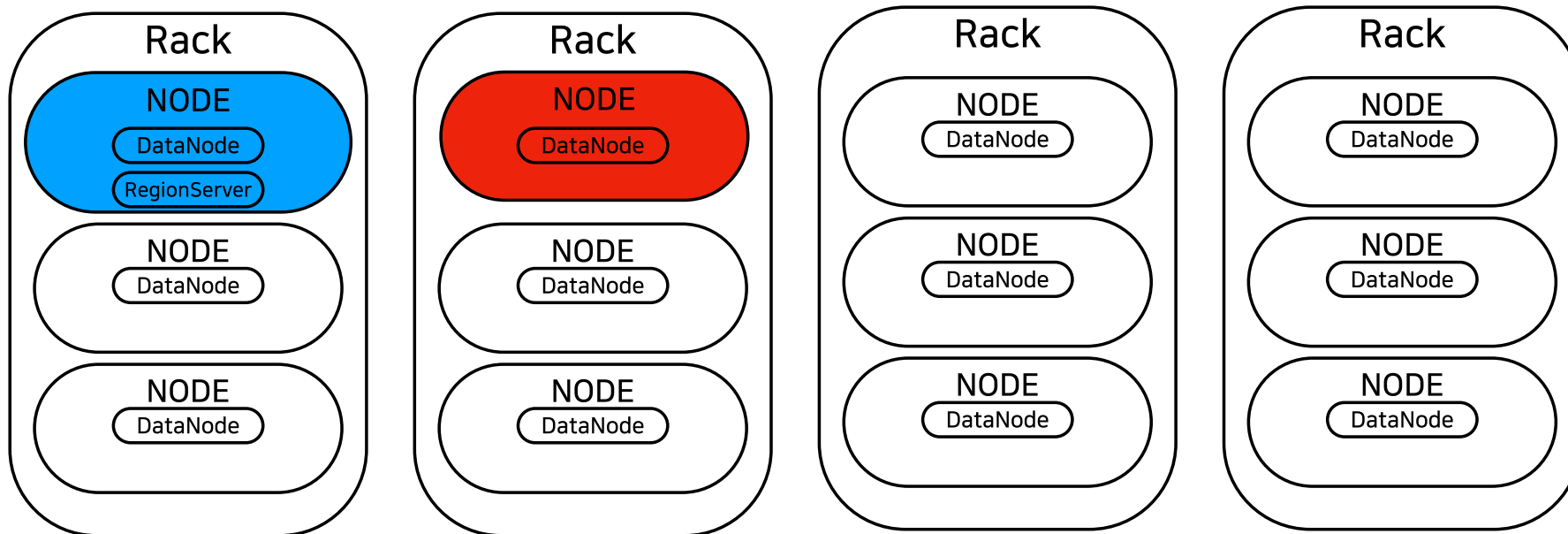
1. Local DataNode - DN1
2. 다른 Rack DataNode - DN2
3. DN2와 동일한 Rack DataNode - DN3
4. Random - DNn



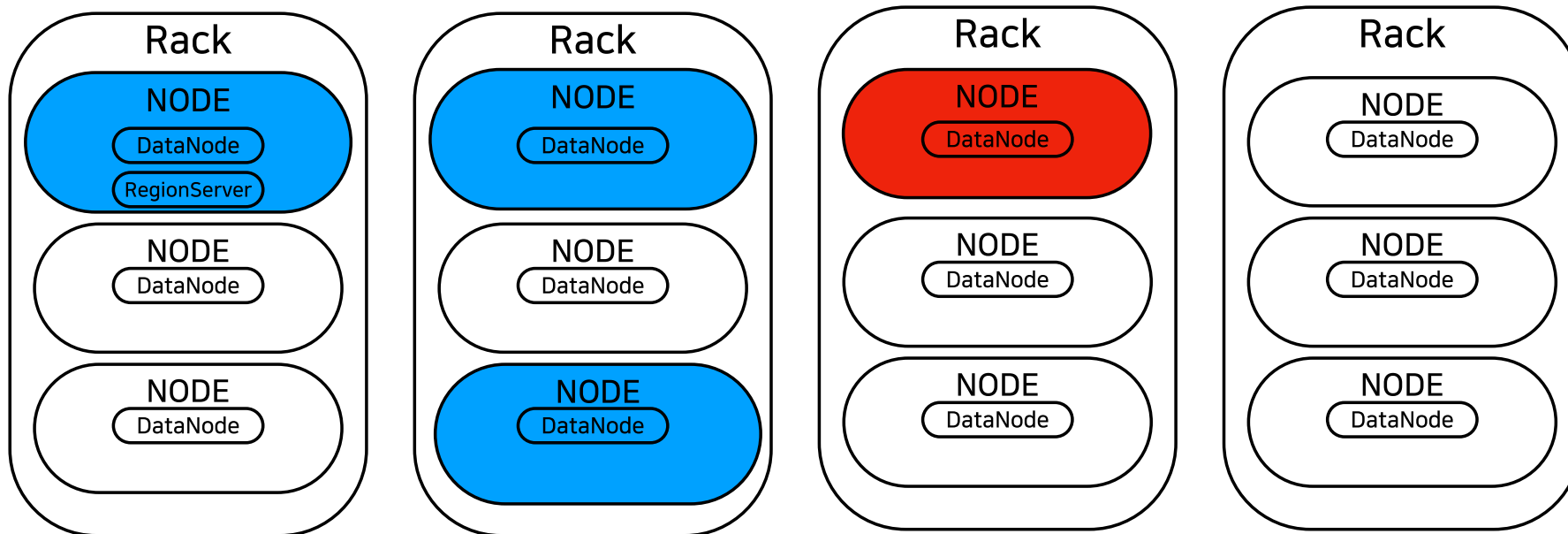
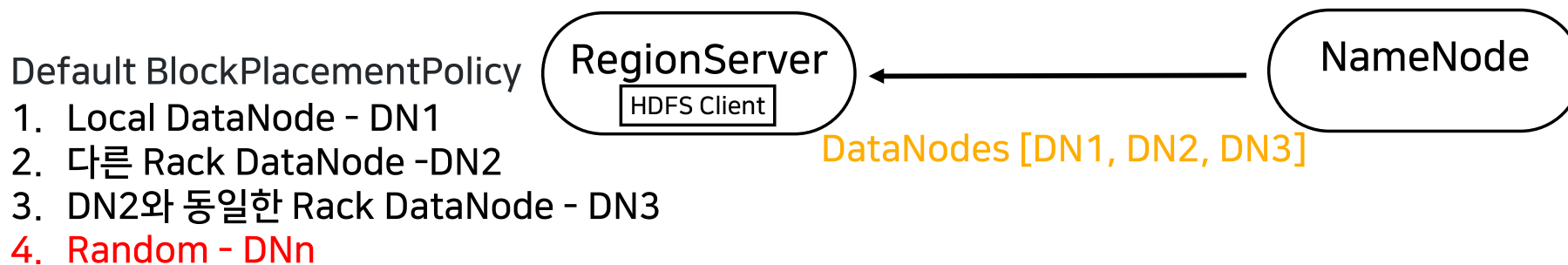
정답은 DataNode 선정 방법!



1. Local DataNode - DN1
2. 다른 Rack DataNode - DN2
3. DN2와 동일한 Rack DataNode - DN3
4. Random - DNn

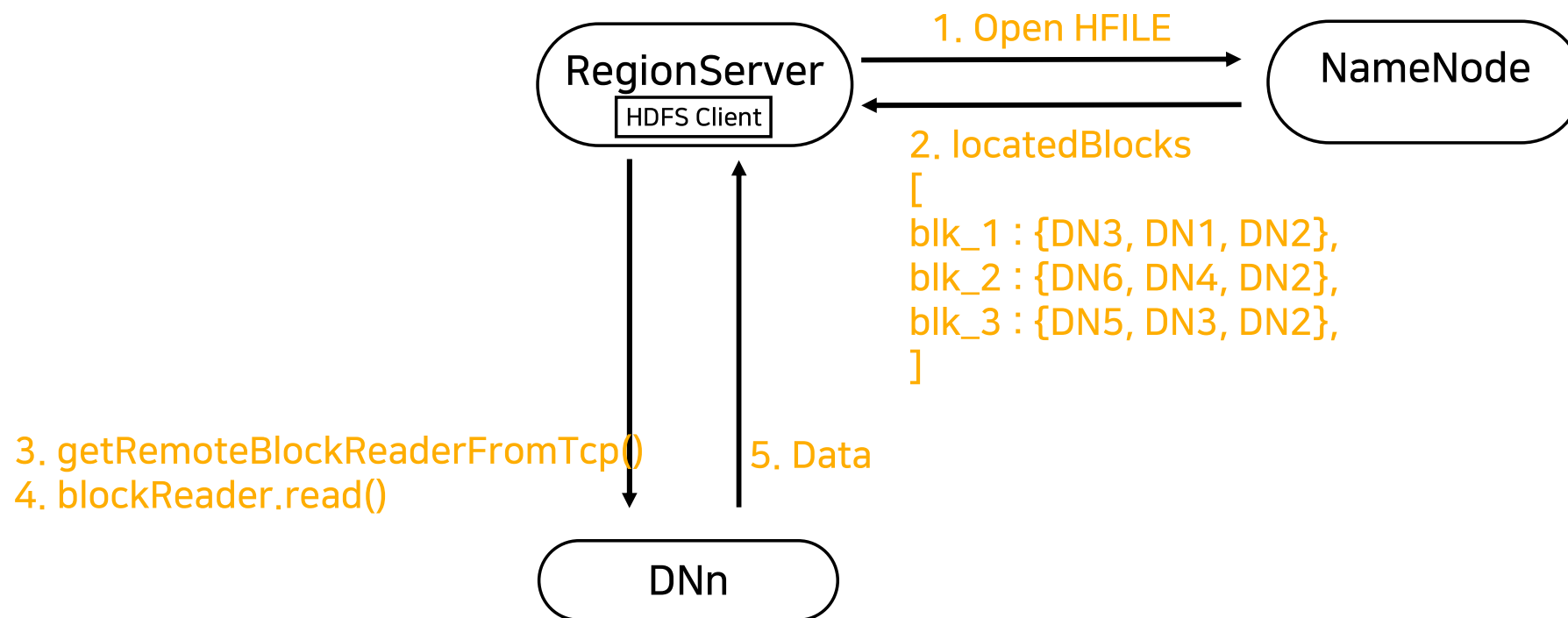


정답은 DataNode 선정 방법!

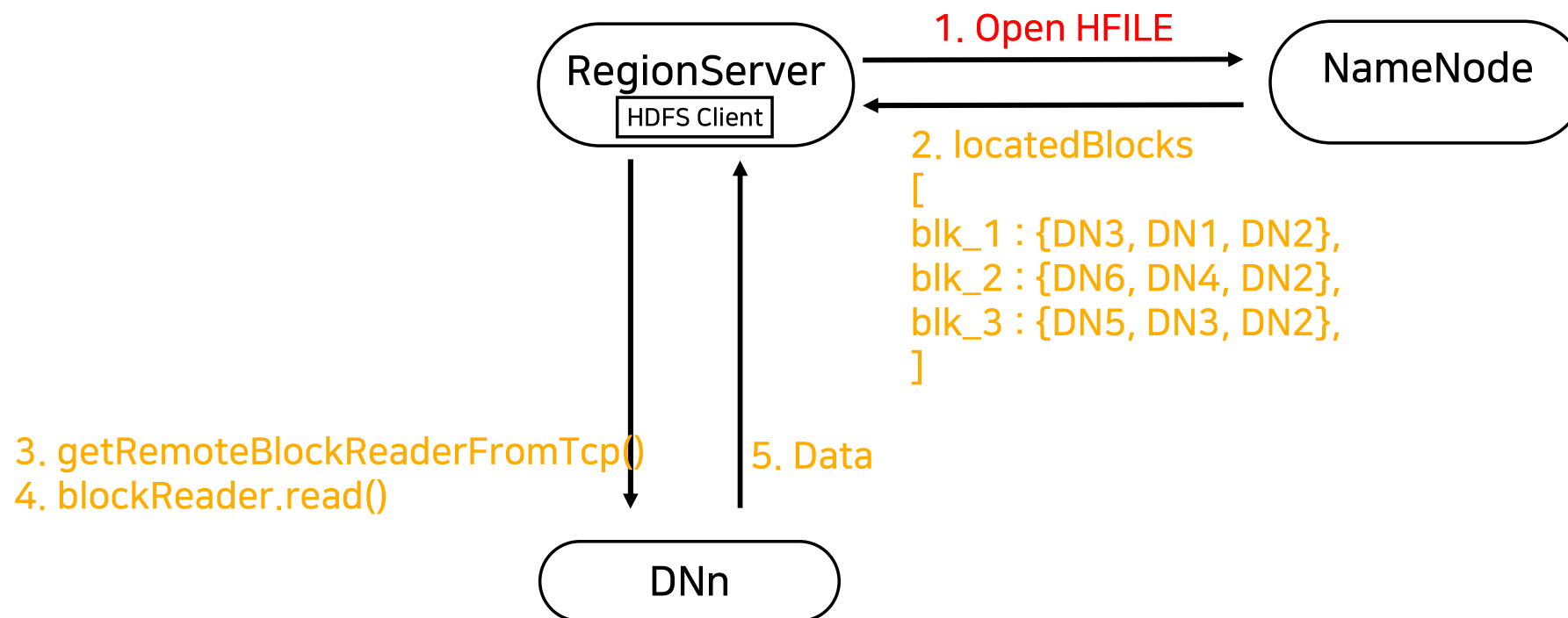


HBase Locality로 읽기 성능 높이는 방법

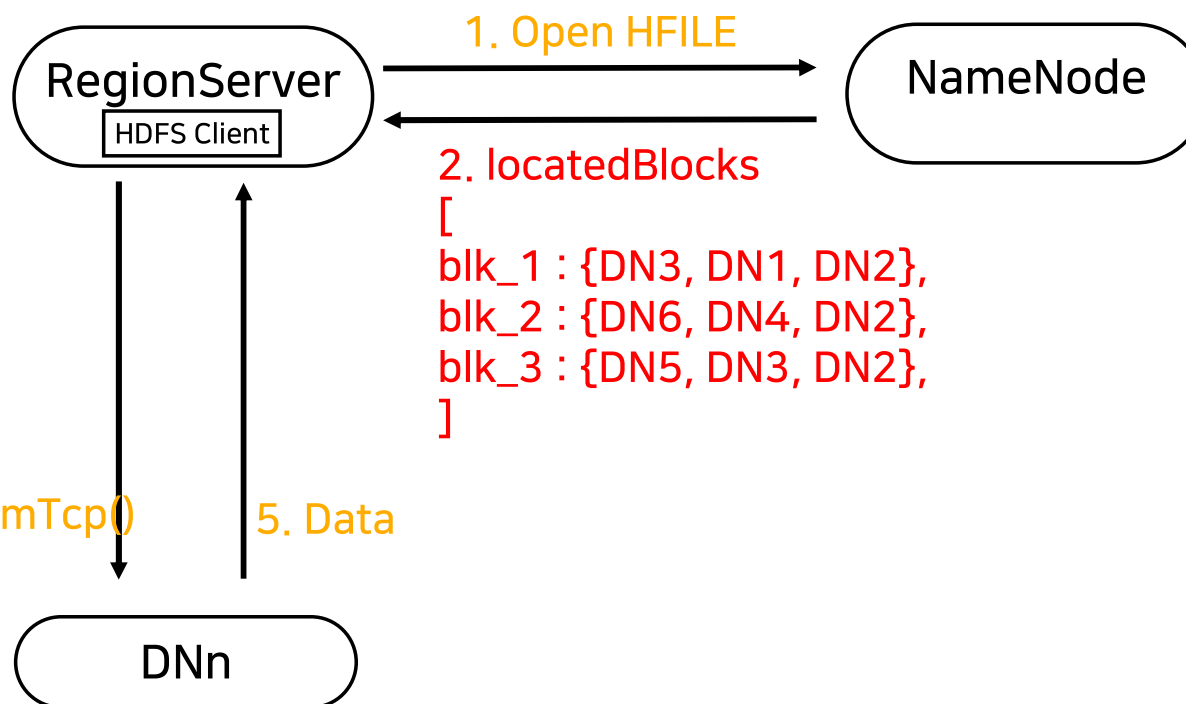
1. LocatedBlocks 순서



1. LocatedBlocks 순서

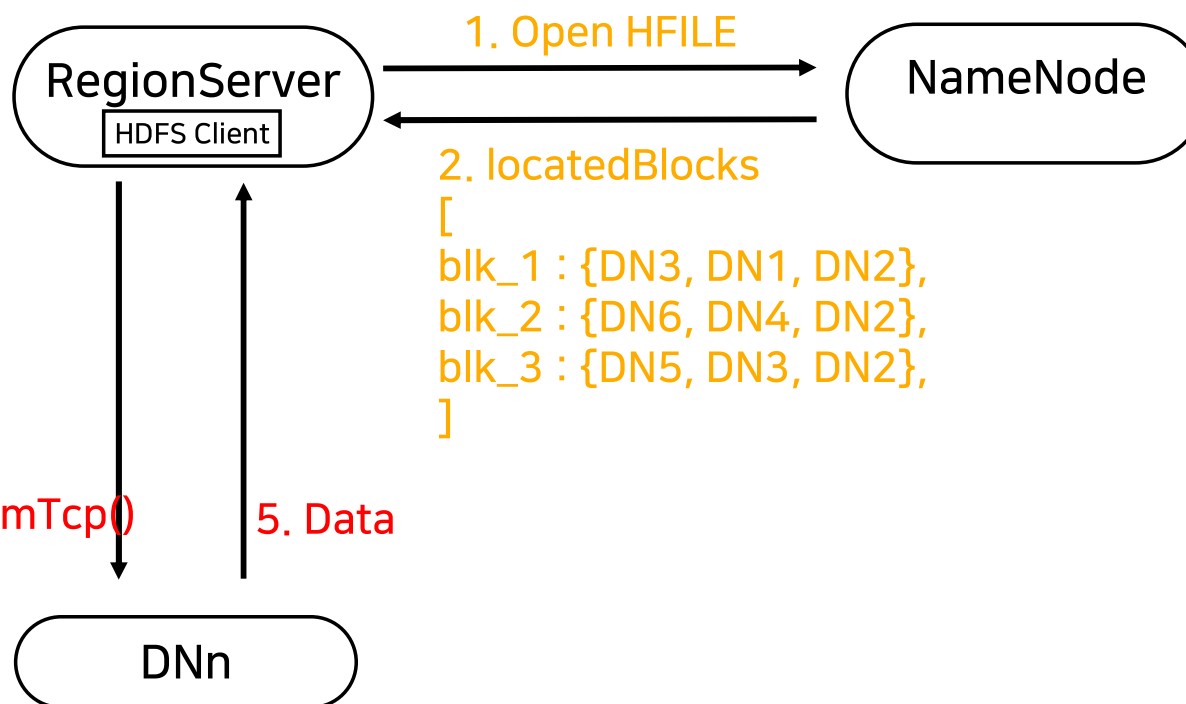


1. LocatedBlocks 순서



3. getRemoteBlockReaderFromTcp()
4. blockReader.read()

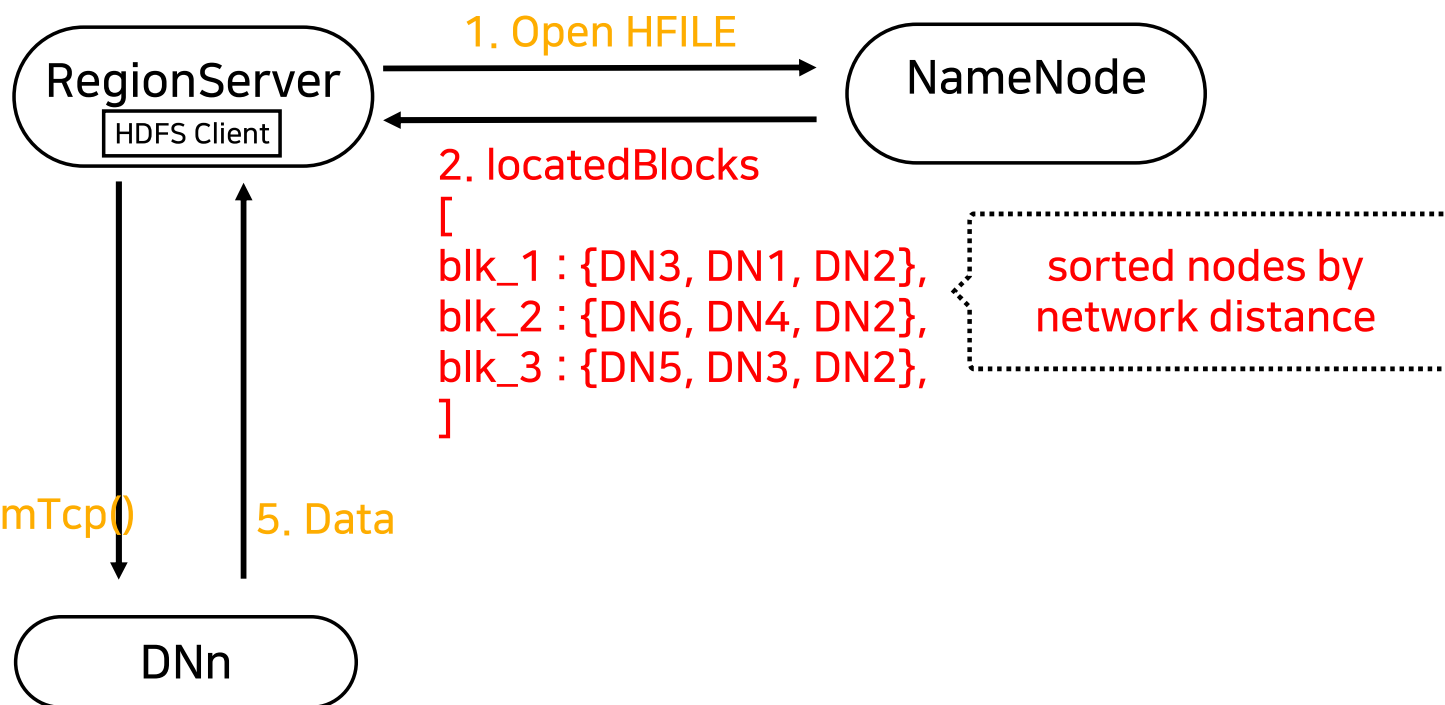
1. LocatedBlocks 순서



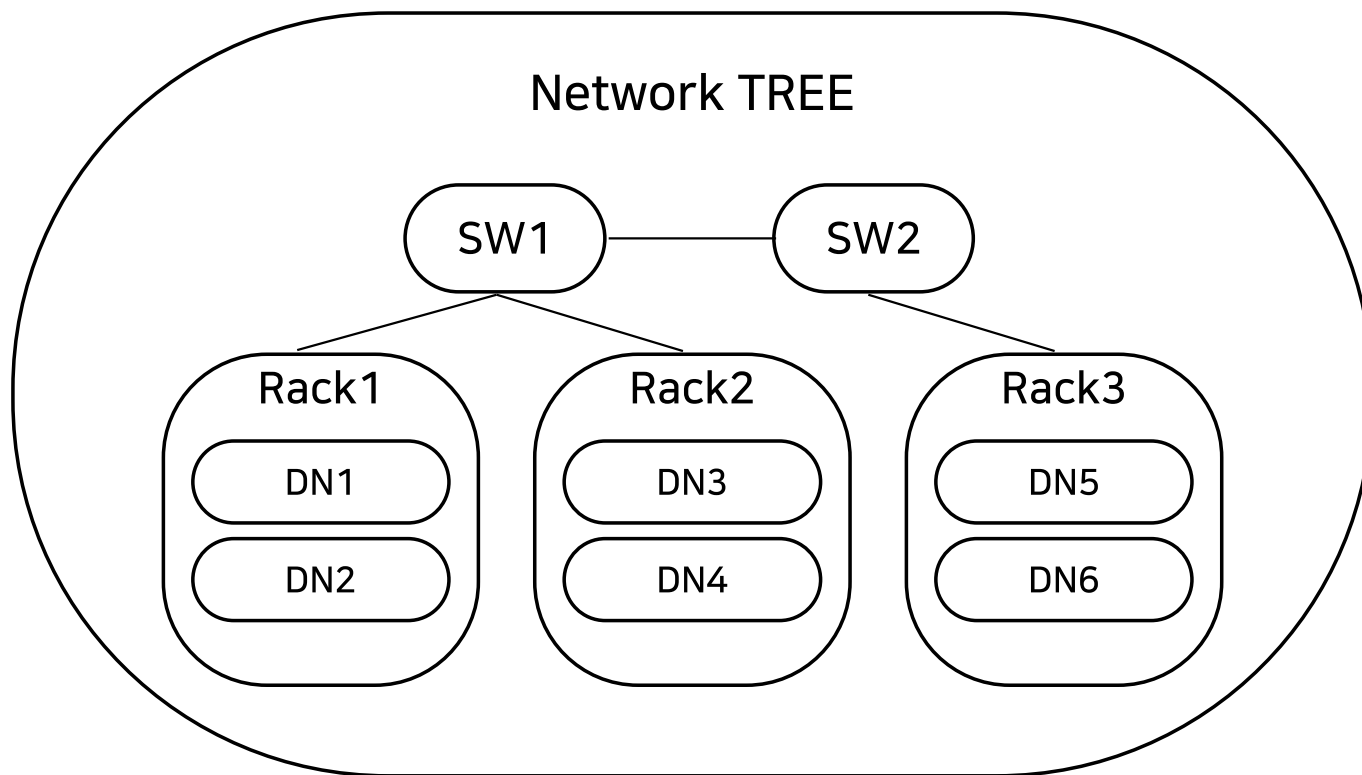
3. getRemoteBlockReaderFromTcp()
4. blockReader.read()

5. Data

1. LocatedBlocks 순서



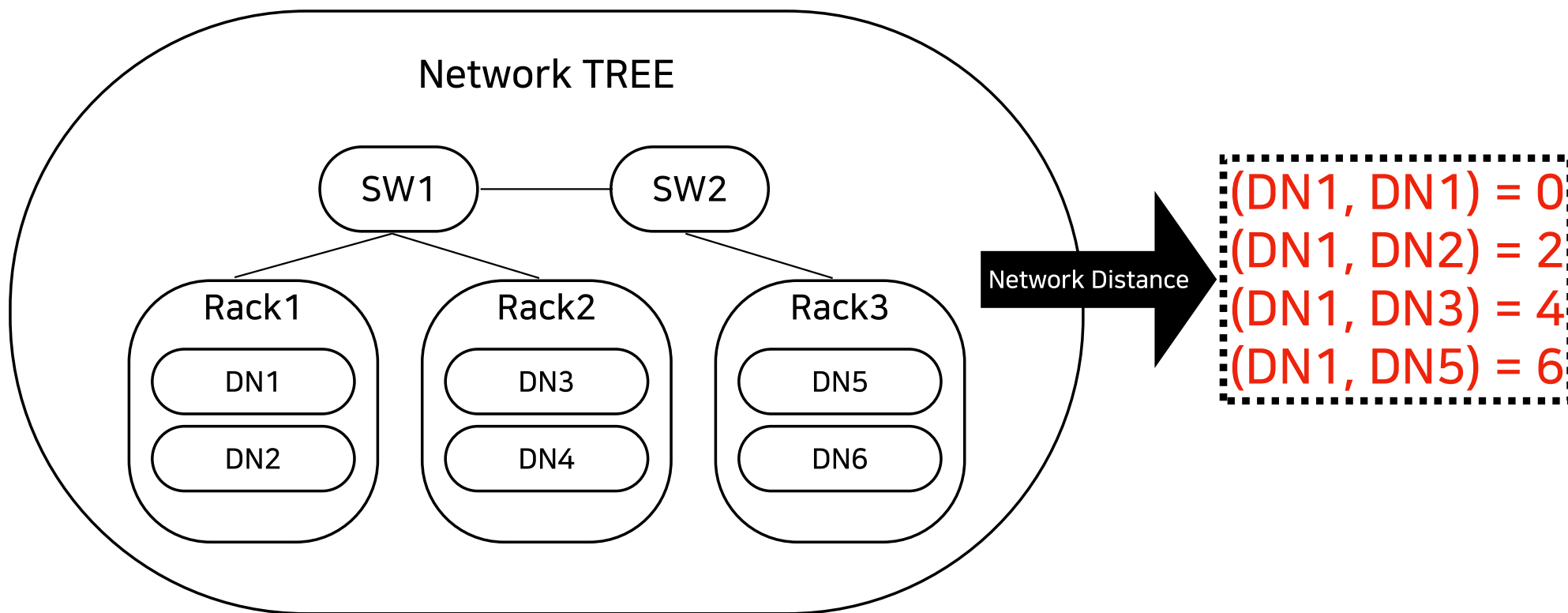
1. LocatedBlocks 순서



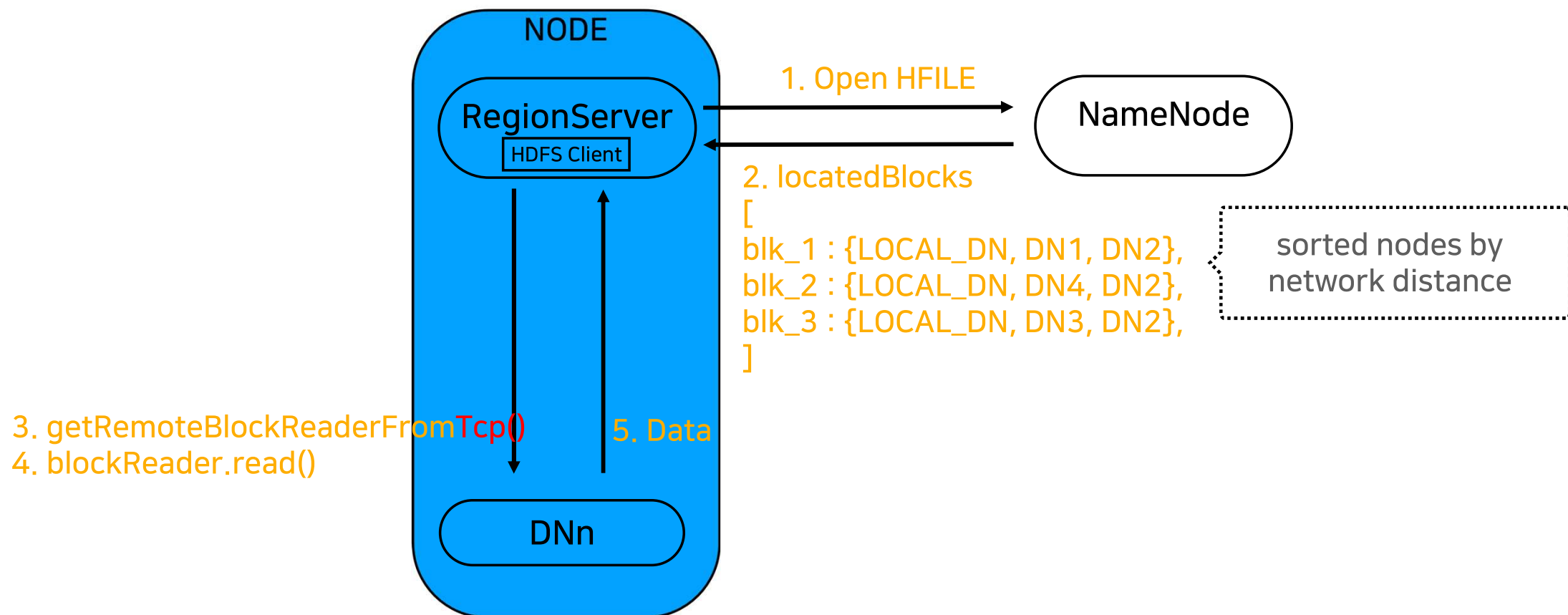
RACK INFO

- DN1 /SW1/RACK1
- DN2 /SW1/RACK1
- DN3 /SW1/RACK2
- DN4 /SW1/RACK2
- DN5 /SW2/RACK3
- DN6 /SW2/RACK3

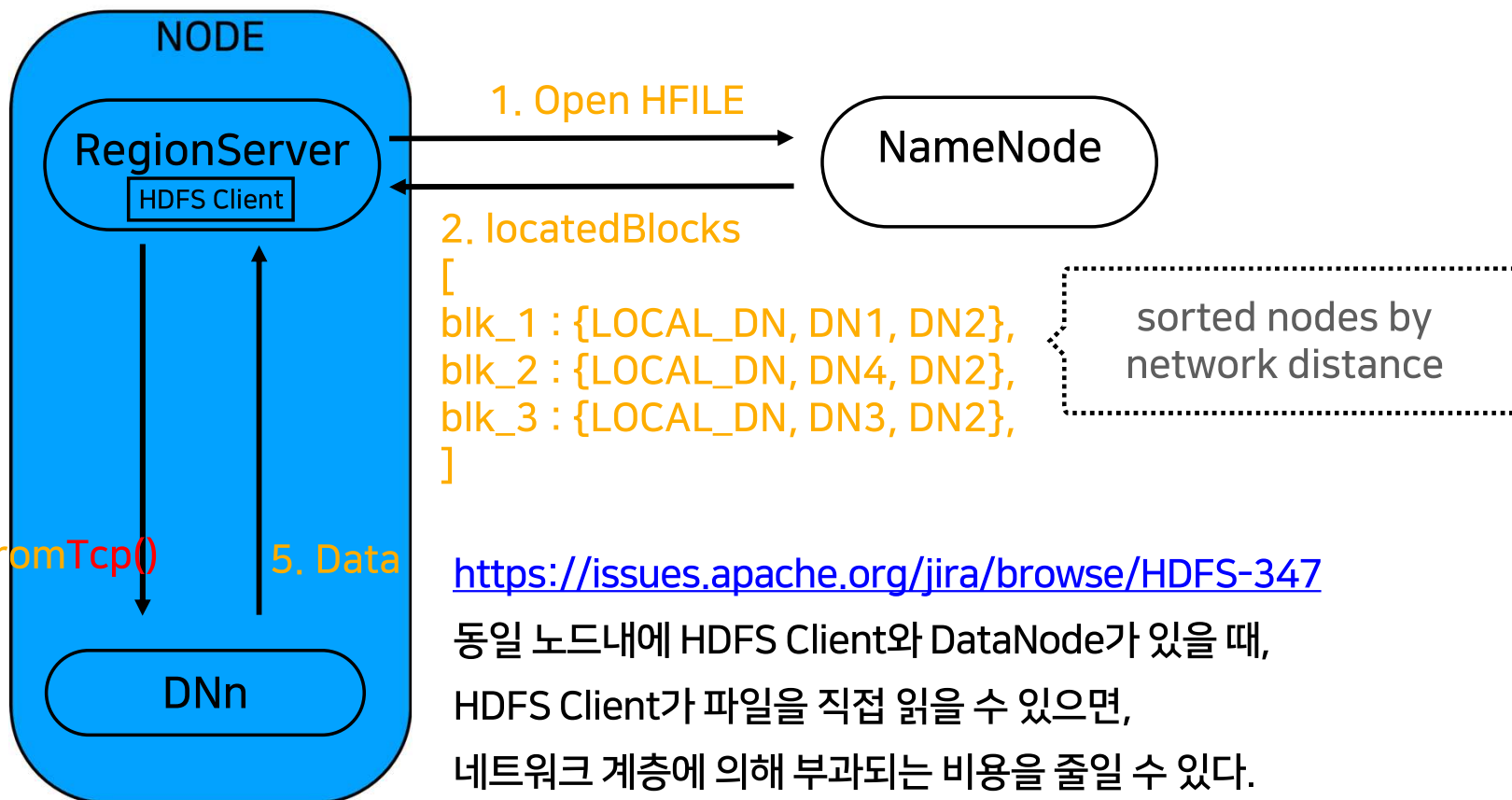
1. LocatedBlocks 순서



1. LocatedBlocks 순서

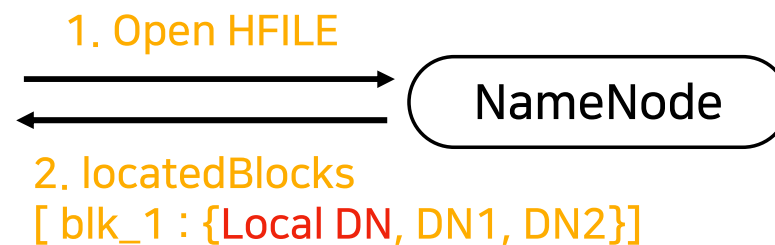
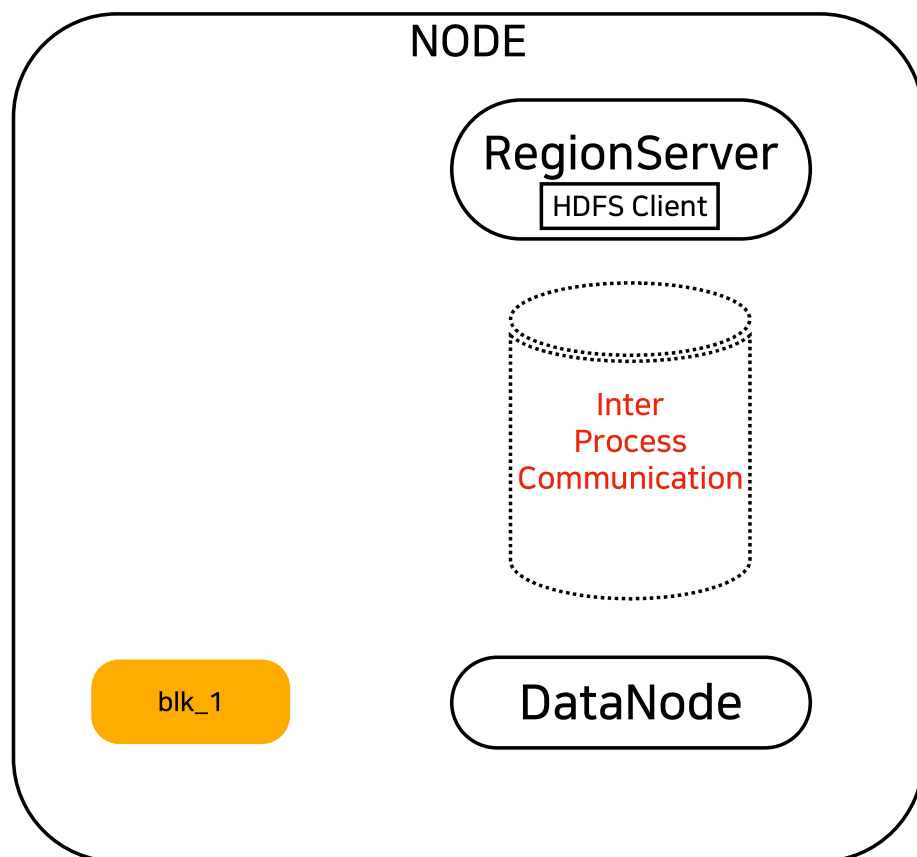


TCP 통신은 왜?

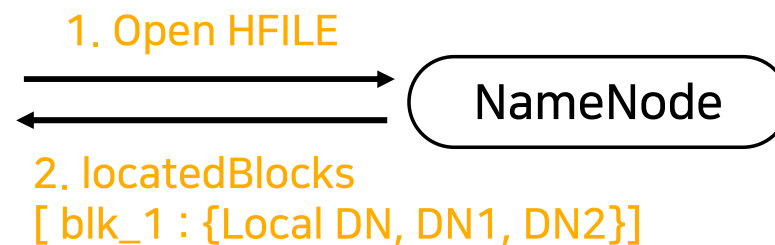
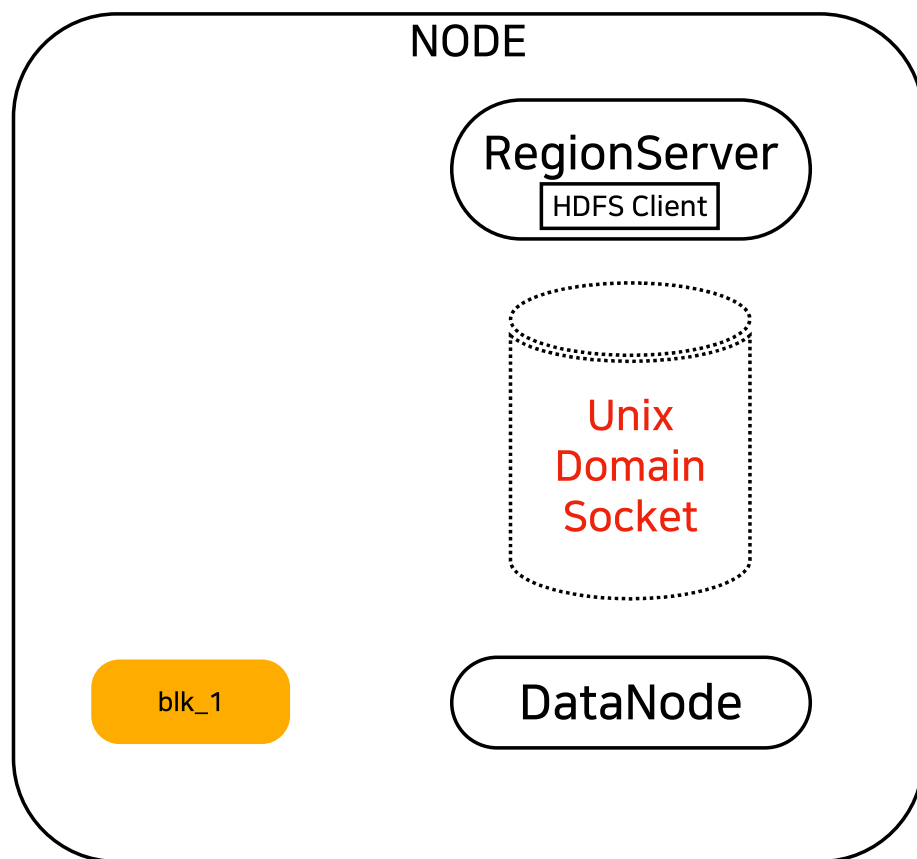


3. getRemoteBlockReaderFromTcp()
4. blockReader.read()

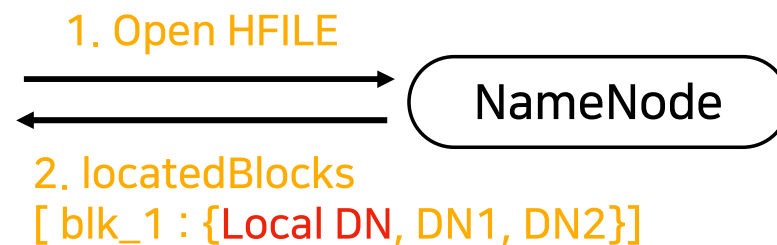
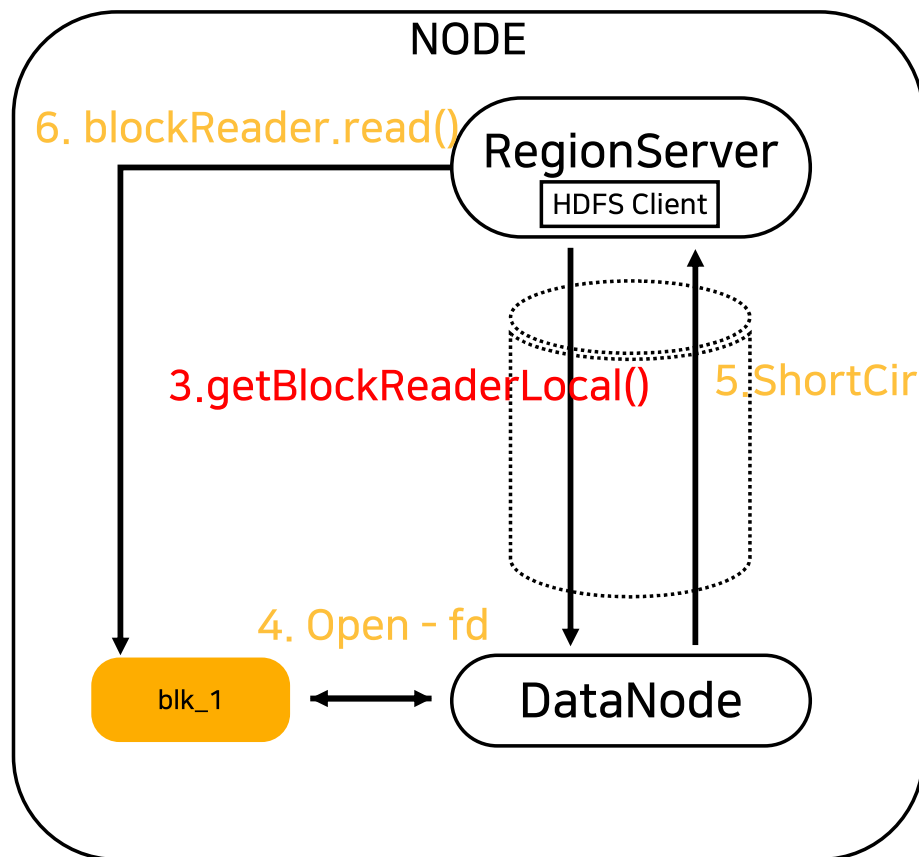
2. HDFS Short-Circuit Local Reads



2. HDFS Short-Circuit Local Reads



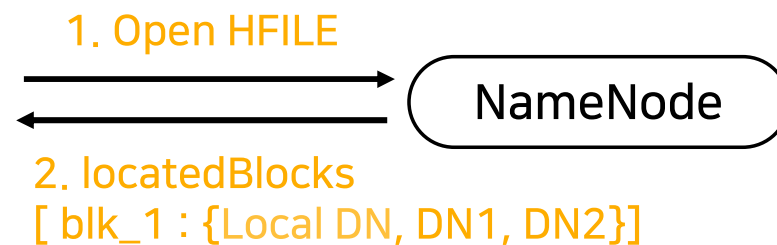
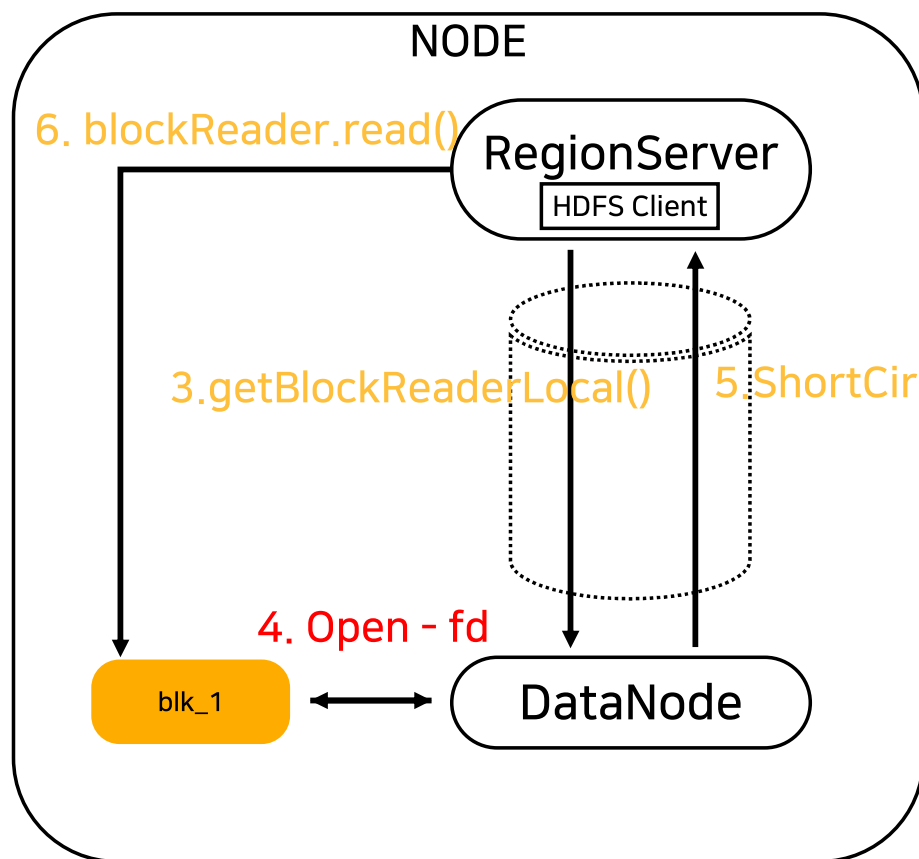
2. HDFS Short-Circuit Local Reads



```
public class ShortCircuitReplica{
    private final FileInputStream dataStream
    private final Slot slot
    ...
}

public FileInputStream(FileDescriptor fd) {
}
```

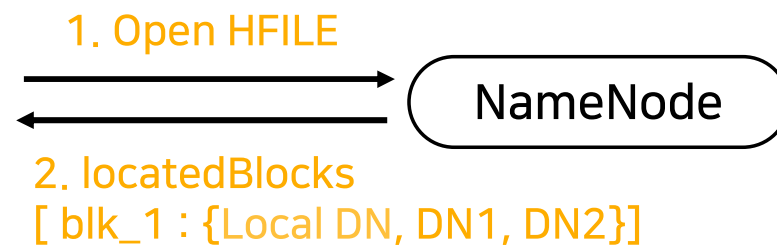
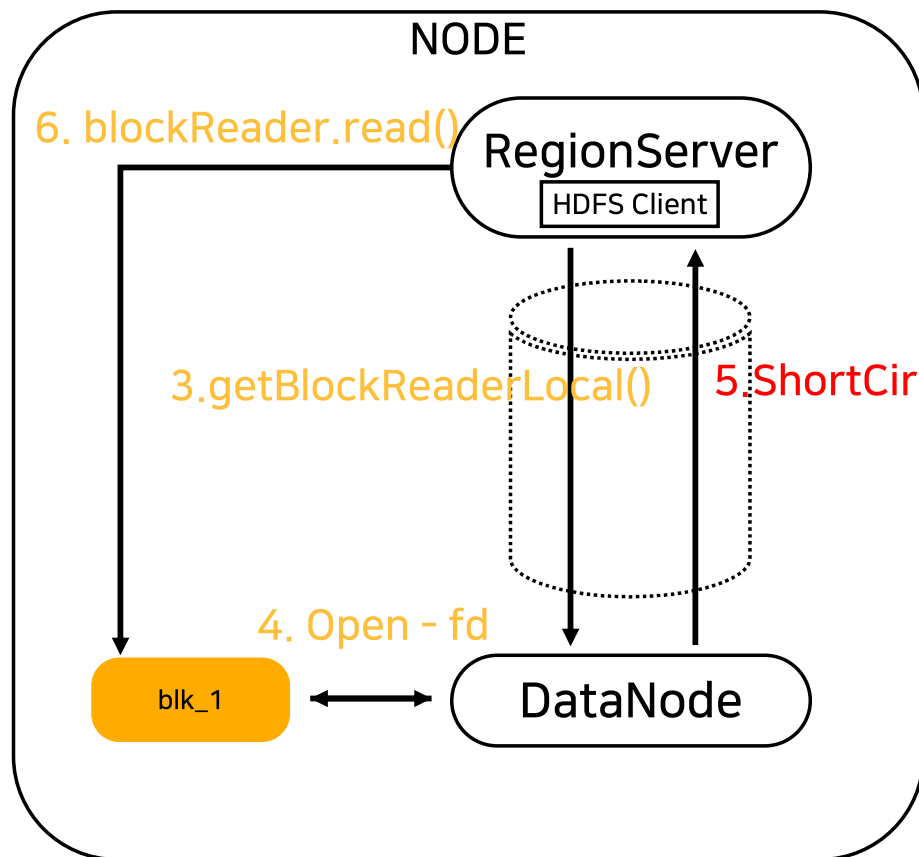
2. HDFS Short-Circuit Local Reads



```
public class ShortCircuitReplica{
    private final FileInputStream dataStream
    private final Slot slot
    ...
}

public FileInputStream(FileDescriptor fd) {
}
```

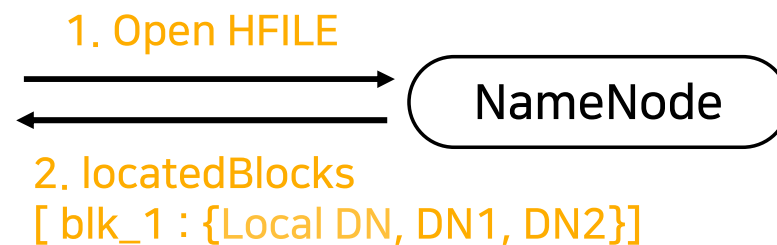
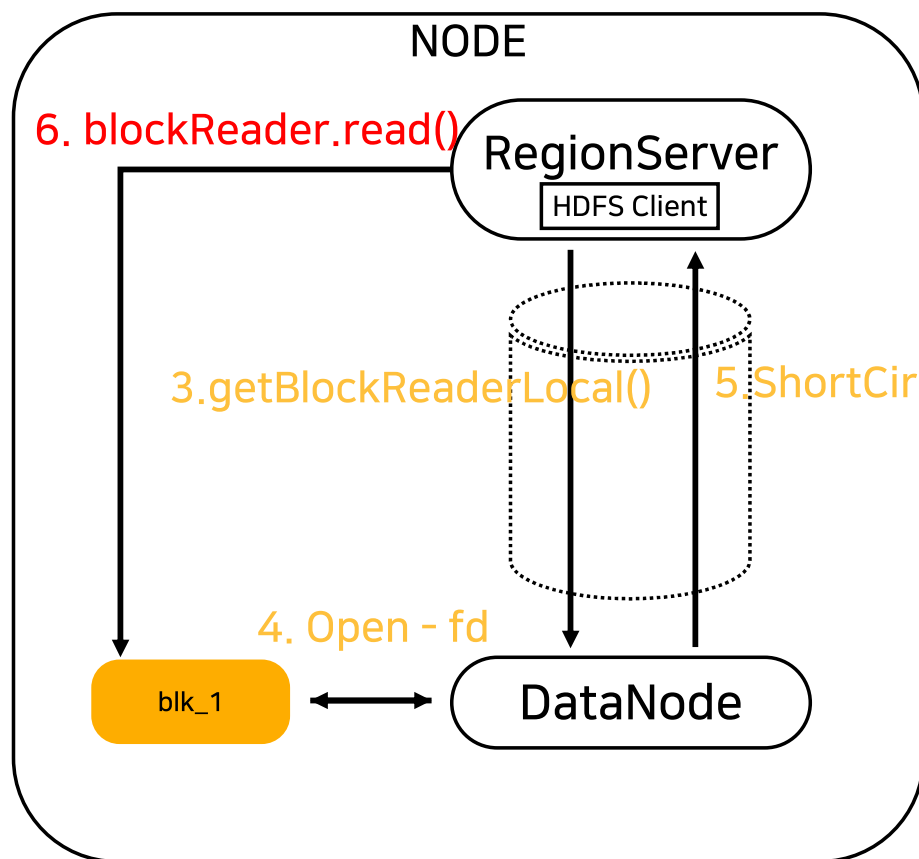
2. HDFS Short-Circuit Local Reads



```
public class ShortCircuitReplica{
    private final FileInputStream dataStream
    private final Slot slot
    ...
}

public FileInputStream(FileDescriptor fd) {
}
```

2. HDFS Short-Circuit Local Reads

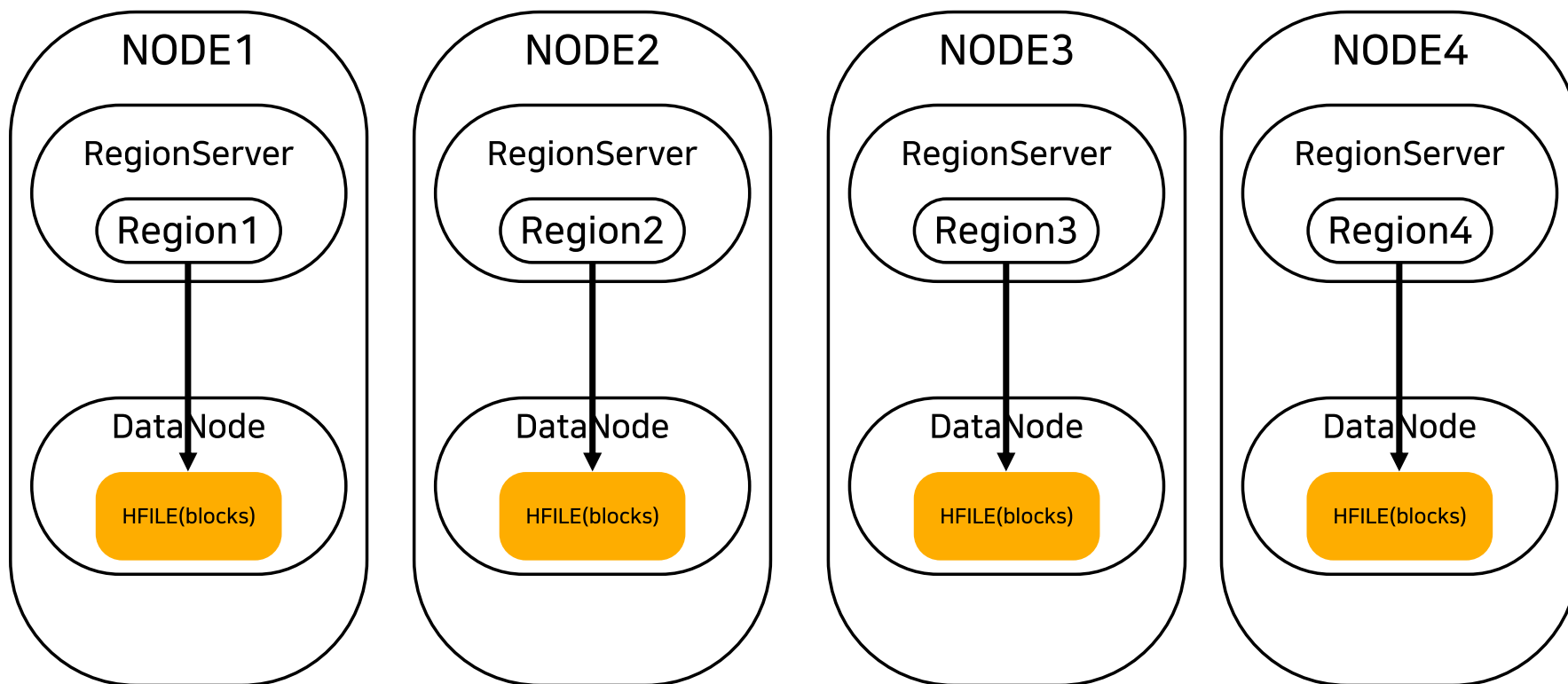


```
public class ShortCircuitReplica{
    private final FileInputStream dataStream
    private final Slot slot
    ...
}

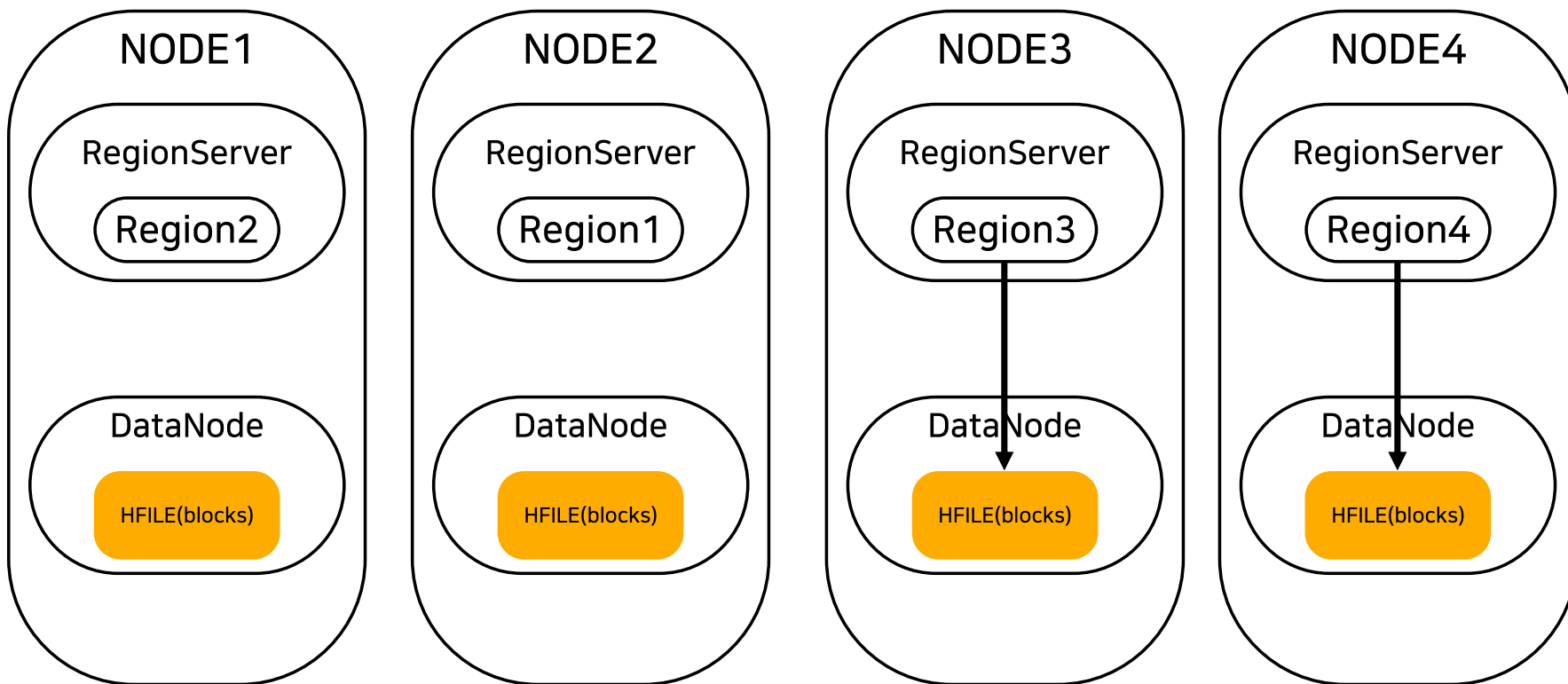
public FileInputStream(FileDescriptor fd) {
}
```

네이버 데이터 저장소에서 HBase Locality를 지키기 위한 운영

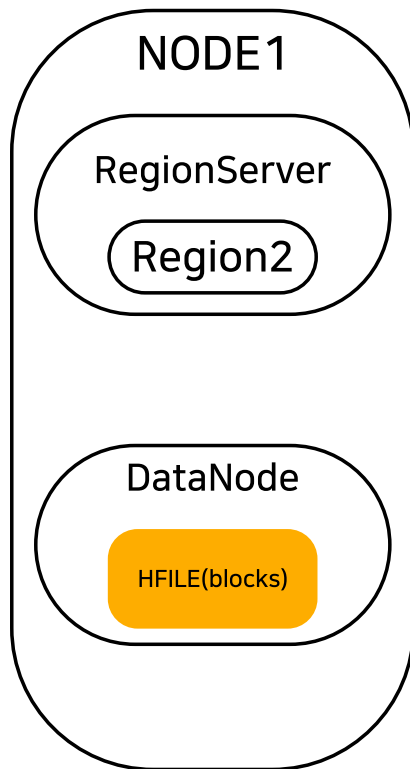
1. 롤링 리스타트 후



1. 롤링 리스타트 후



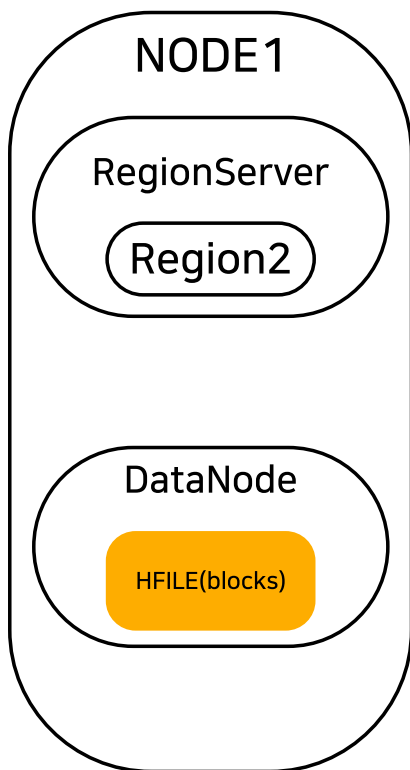
1. 롤링 리스타트 후



HDFS://{HDFS_NAMESPACE}/HBase/data
/{NameSpace}/{Table}/{Region}/{CF}/HFILES

Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name
-rw-r--r--	hbase	hbase	202.33 MB	2021. 4. 14. 오후 2:08:39	3	256 MB	50bb7714f50f4d4b9c04ea3c3c6e8162
-rw-r--r--	hbase	hbase	103.79 GB	2021. 4. 14. 오전 1:10:45	3	256 MB	bda77ddd53ac46b6b6362c71d3efb944
-rw-r--r--	hbase	hbase	1.41 GB	2021. 4. 14. 오후 12:18:45	3	256 MB	be58c4896dde44389d3ee587ee2e53b6
-rw-r--r--	hbase	hbase	67.37 MB	2021. 4. 14. 오후 2:46:09	3	256 MB	e29b8ef79e1a40288a005bd3a154d644

1. 롤링 리스타트 후



File information - bda77ddd53ac46b6b6362c71d3efb944

Download

Block information -- Block 0

Block ID: 1124248829

Block Pool ID: [REDACTED]

Generation Stamp: 50508946

Size: 268435456

Availability:

- NODE1
- NODE2
- NODE3

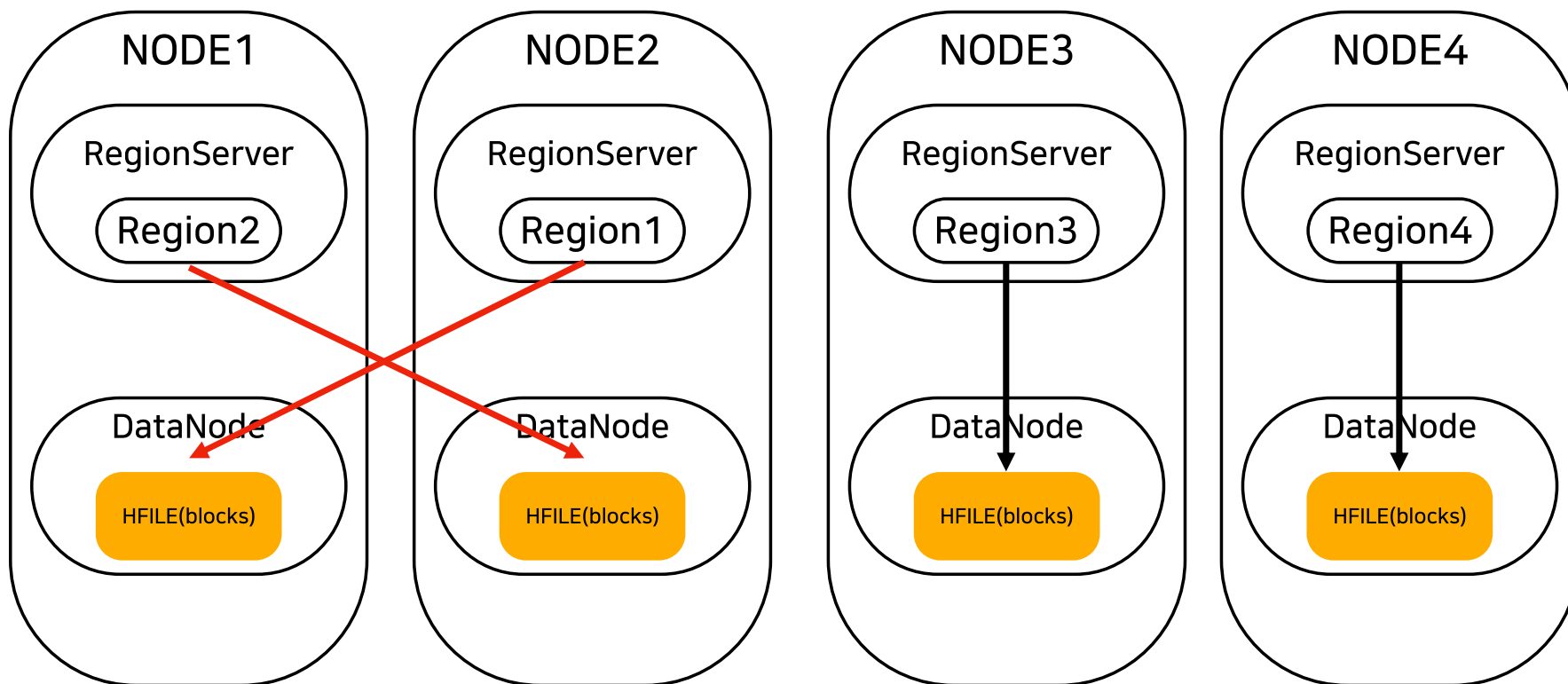
Permissions: -rw-r--r--

Size	Name
VB	50bb7714f50f4d4b9c04ea3c3c6e8162
VB	bda77ddd53ac46b6b6362c71d3efb944
VB	be58c4896dde44389d3ee587ee2e53b6
VB	e29b8ef79e1a40288a005bd3a154d644

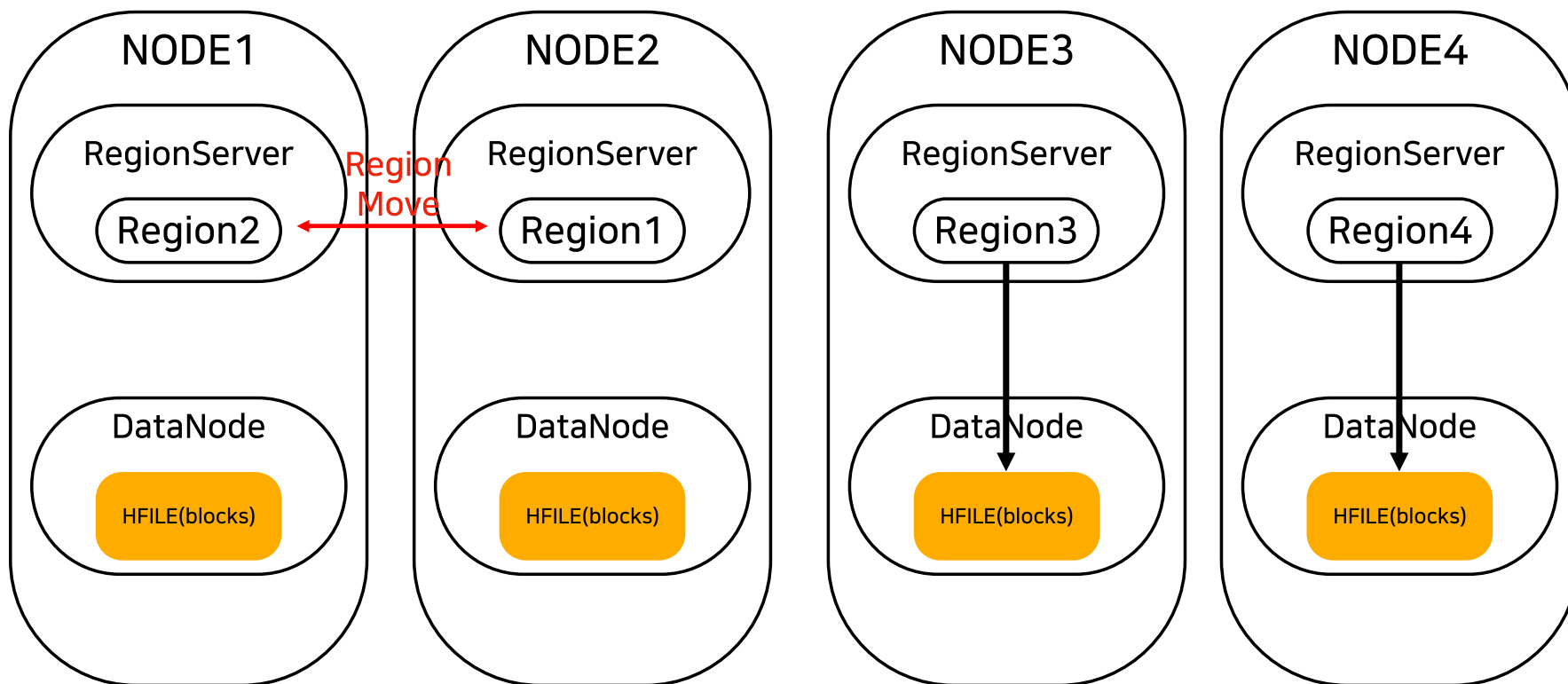
e/data
:F}/HFILES

Close

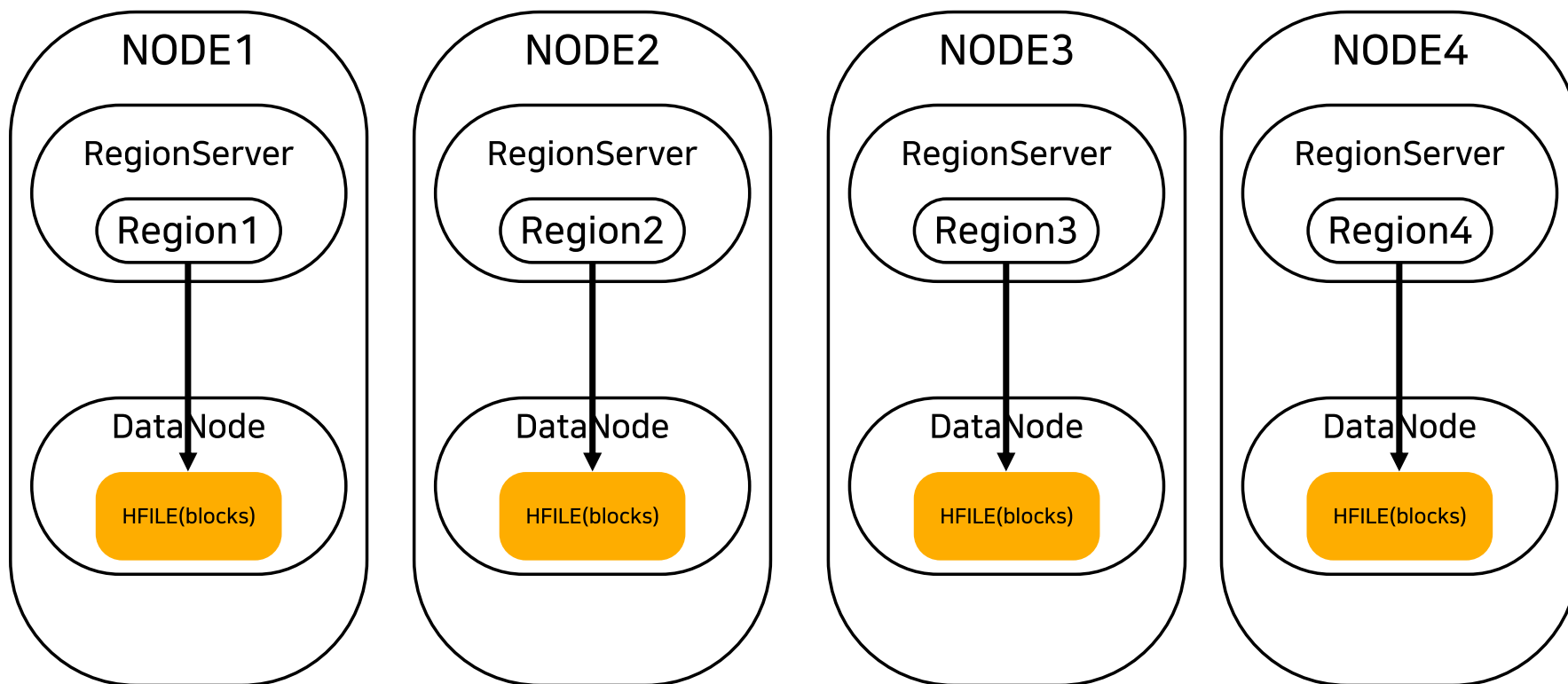
1. 롤링 리스타트 후



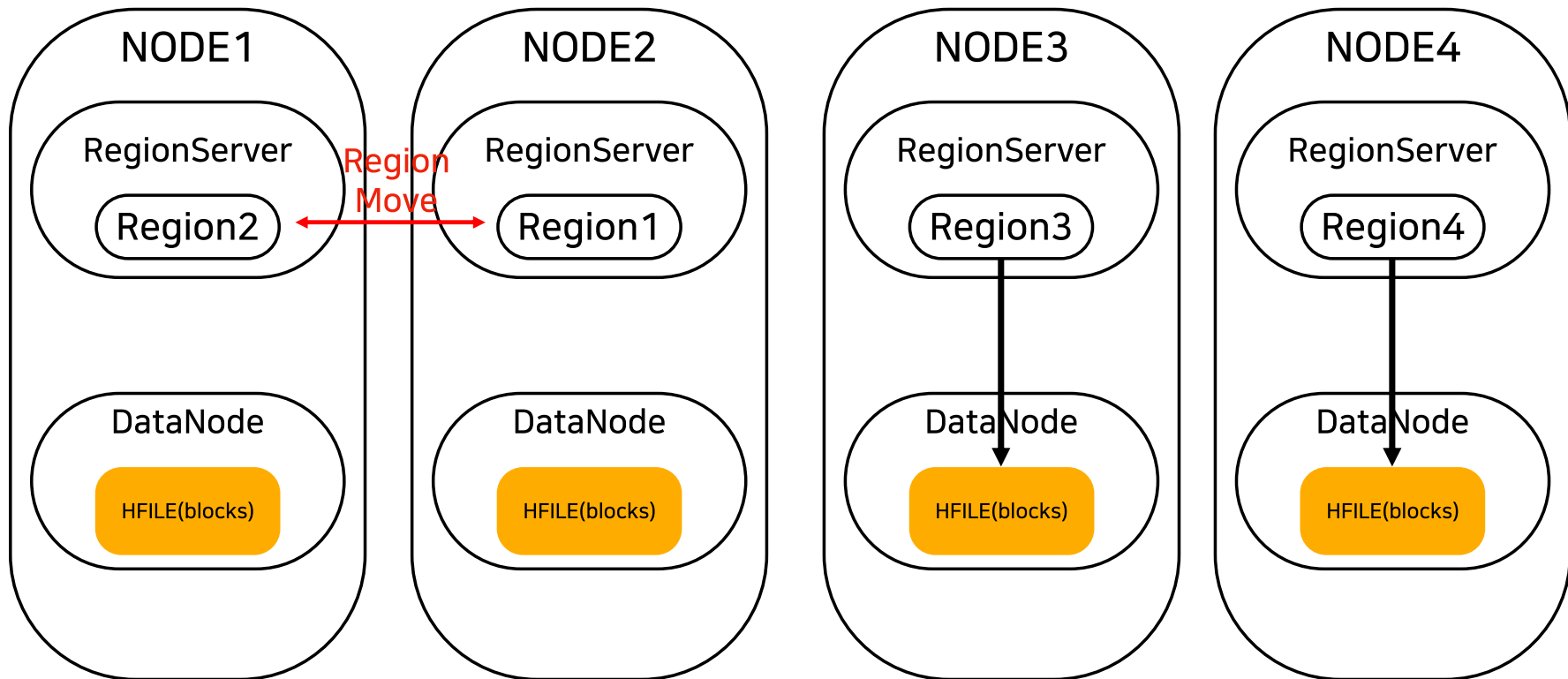
1. 롤링 리스타트 후



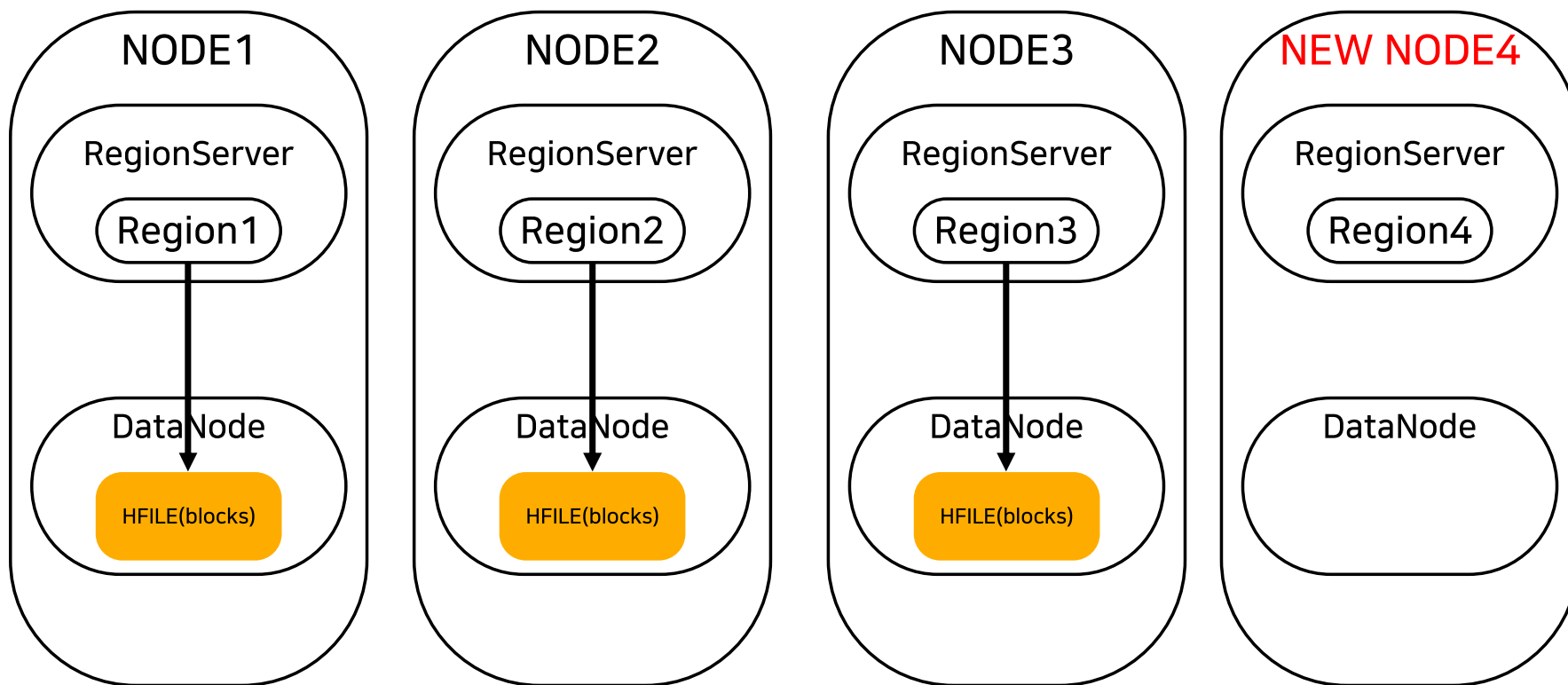
1. 롤링 리스타트 후



2. 장비 수리 후

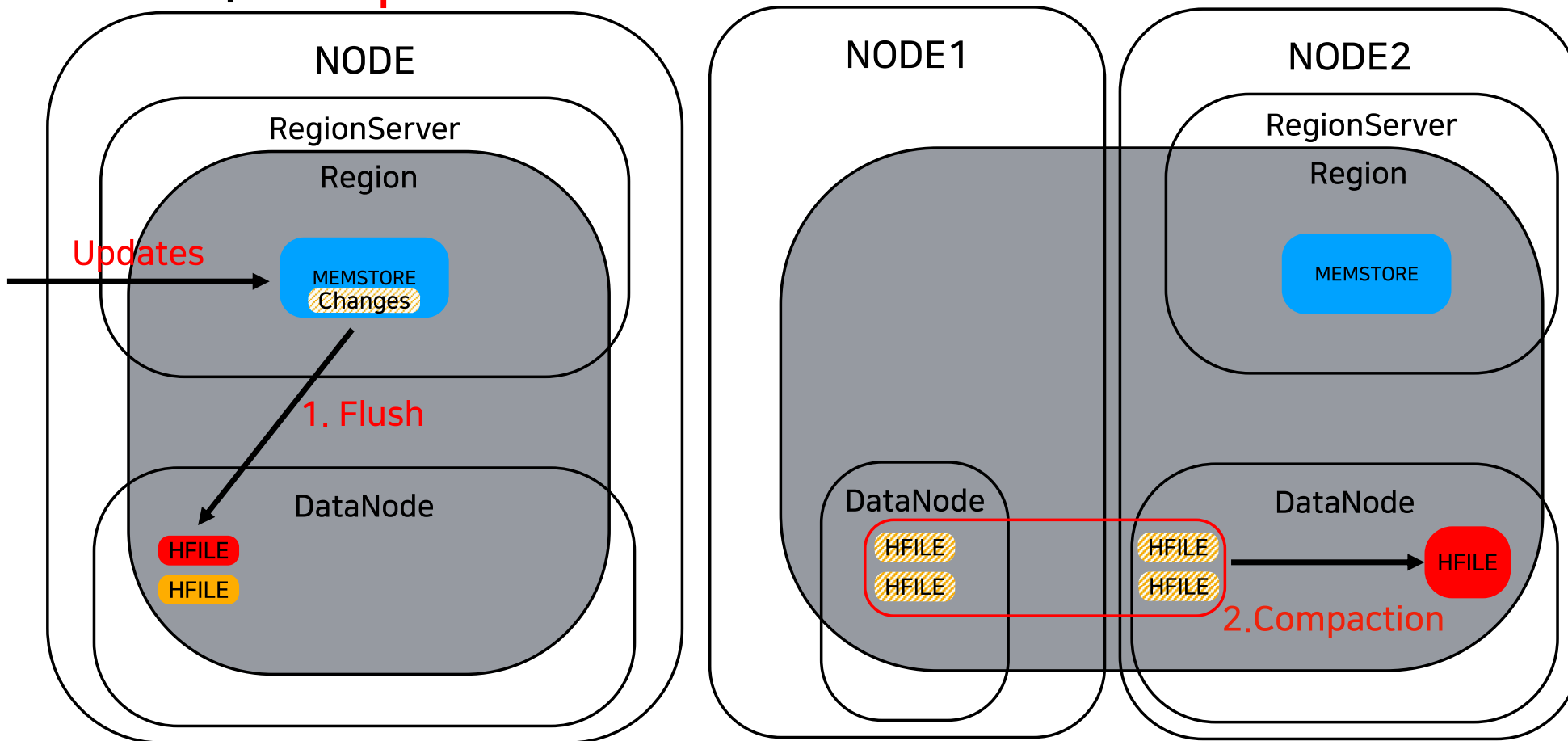


3. 새로운 장비 추가 후

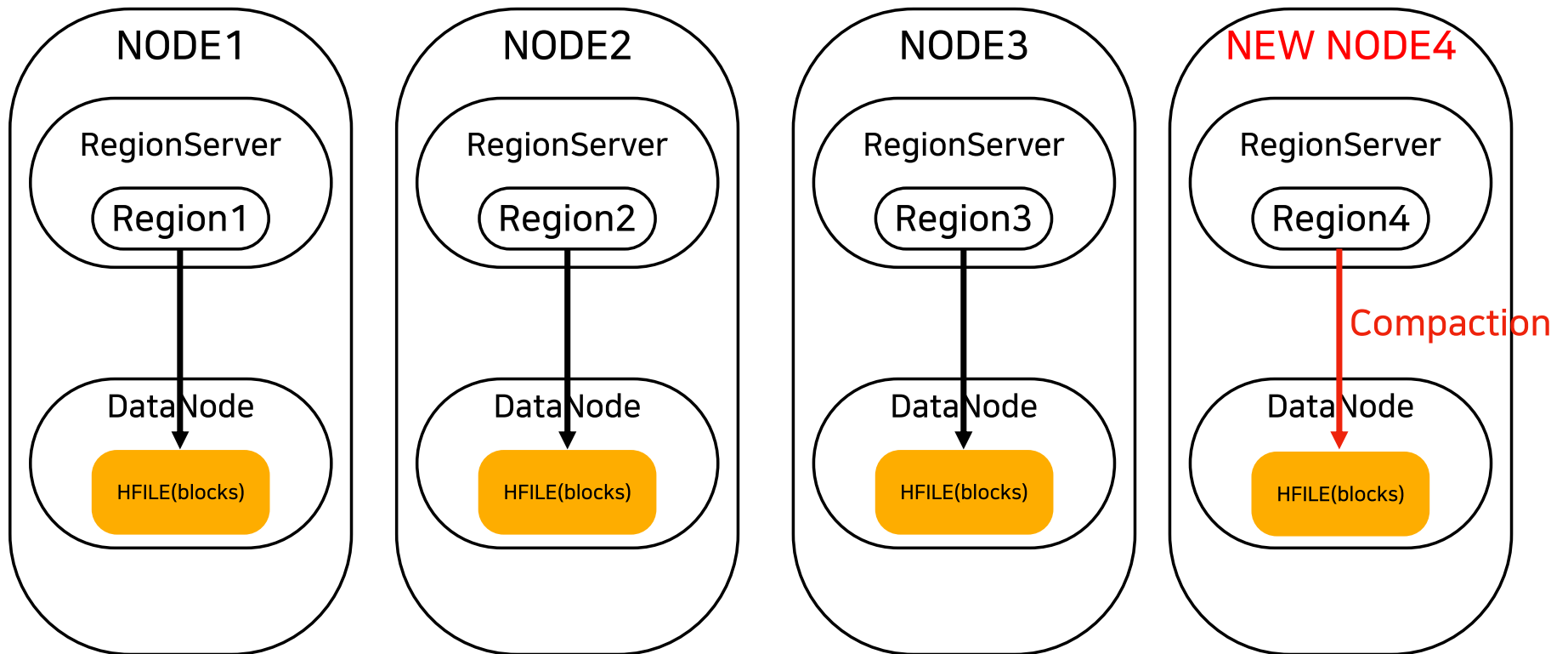


HBase에서 Locality를 올리는 방법

Flush와 Compaction



3. 새로운 장비 추가 후



클러스터 규모와 비례하는 디스크 장애

- 수천대의 클러스터
- 장비 한대에서 사용하는 DISK는 12개
- 클러스터 규모와 증가하는 DISK 장애율

Cluster	장비 수 (대)
Cluster A	xxxx
Cluster B	xxx
Cluster C	xxx
...	...
총	xxxx대

해당 년. 월	장애 발생 횟수 (건)
2021. 4	151
2021. 5	210
2021. 6	102
2021. 7	126
2021. 8	249
2021. 9	168

클러스터 규모와 비례하는 디스크 장애

- 한국 클러스터만 약 6000대
- 장비 한대에서 사용하는 DISK는 12개
- 클러스터 규모와 증가하는 DISK 장애율

클러스터 규모와 반비례 하는 HBase Locality!

클러스터명	클러스터 규모	클러스터 수	클러스터당 장애율
Cluster A	약 1500	2021.4	151
Cluster B	약 500	2021.5	210
Cluster C	약 500	2021.6	102
...	...	2021.7	126
...	...	2021.8	249
총	약 6000대	2021.9	168

DataNode Hot Swap Drive

- DataNode Damon을 종료하지 않고 디스크 추가/제거 가능
- Hadoop 2.7.0 부터 지원
- 12개의 디스크 중 고장난 디스크만 교체 & 투입 가능
- HBase Locality를 최대한 유지하면서 교체 가능

DataNode Hot Swap Drive

DataNode 설정

```
<property>  
<name>dfs.datanode.data.dir</name>  
  <value>/disk1/datanode,/disk2/datanode,/disk3/datanode</value>  
</property>
```

NODE

DataNode

DISK1

DISK2

DISK3

DataNode Hot Swap Drive

DataNode 설정

```
<property>  
<name>dfs.datanode.data.dir</name>  
  <value>/disk1/datanode,/disk3/datanode</value>  
</property>
```

NODE

DataNode

DISK1

DISK2

DISK3

DataNode Hot Swap Drive

DataNode 설정

```
quartz.scheduler.class=org.quartz.simpl.SimpleThreadPoolScheduler
dfs.nfs.address=${dfs.nfs.address}
dfs.nfs.url=${dfs.nfs.url}
dfs.nfs.dir=${dfs.nfs.dir}
dfs.nfs.permissions=${dfs.nfs.permissions}
```

Reconfig

```
hdfs dfsadmin -reconfig datanode ${DATANODE}:${PORT} start
```

DataNode

1001

2001

DataNode Hot Swap Drive

DataNode 설정

```
<property>  
<name>dfs.datanode.data.dir</name>  
  <value>/disk1/datanode,/disk3/datanode</value>  
</property>
```

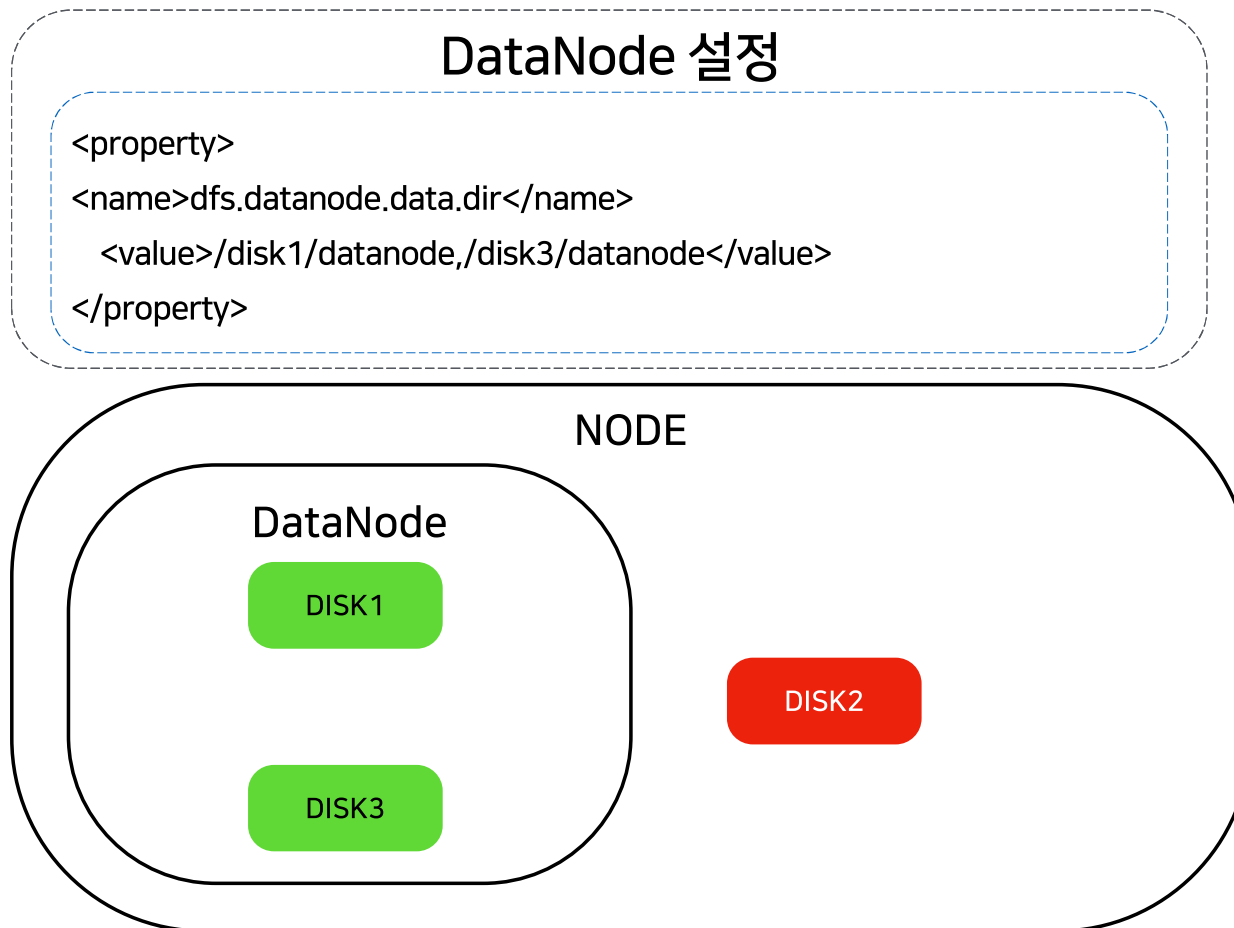
NODE

DataNode

DISK1

DISK3

DISK2



DataNode Hot Swap Drive

DataNode 설정

```
<property>  
<name>dfs.datanode.data.dir</name>  
  <value>/disk1/datanode,/disk2/datanode,/disk3/datanode</value>  
</property>
```

NODE

DataNode

DISK1

DISK3

NEW_DISK2

DataNode Hot Swap Drive

DataNode 설정

```
quorum=quorum
dfs.namenode.secondary.dfs.name.dir=${dfs.namenode.secondary.dfs.name.dir}
dfs.replicate=${dfs.replicate}
dfs.replication=${dfs.replication}
dfs.replicate=${dfs.replicate}
dfs.replication=${dfs.replication}
```

Reconfig

```
hdfs dfsadmin -reconfig datanode ${DATANODE}:${PORT} start
```

DataNode

1001

1001

1001

DataNode Hot Swap Drive

DataNode 설정

```
<property>  
<name>dfs.datanode.data.dir</name>  
  <value>/disk1/datanode,/disk2/datanode,/disk3/datanode</value>  
</property>
```

NODE

DataNode

DISK1

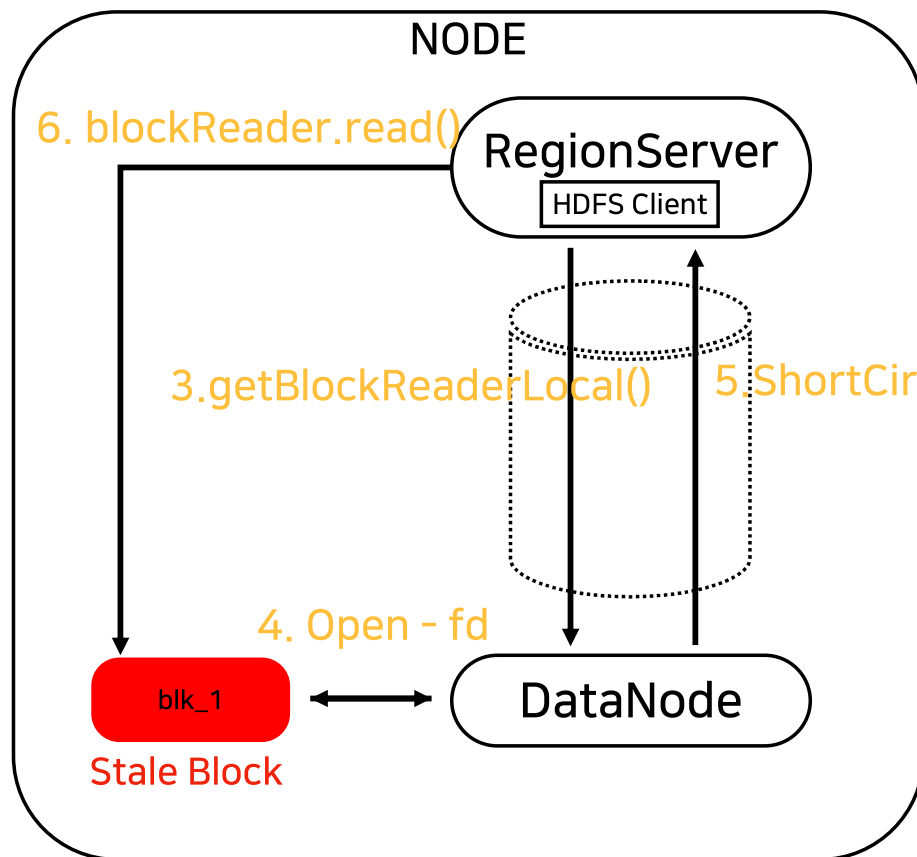
NEW_DISK2

DISK3

언마운트 불능 현상1

- 간헐적으로 디스크가 언마운트 되지 않는 현상 발생
 - -f 옵션을 주고 강제로 언마운트함
 - 이 경우, OS에서는 계속 오류가 발생하고 있는 것으로 인지
 - 혹은 FD를 물고있는 프로세스 종료 (ex: RegionServer)
- 원인은 HDFS의 short-circuit이 HBase에서 예기치 못한 동작을 초래한 것
 - 디스크 장애 시, HDFS에서는 즉시 해당 디스크의 FD를 닫음
 - 하지만, lazy 방식으로 FD를 관리하는 HBase에서는 열어두고 있음
 - 결과적으로 해당 디스크가 사용 중(busy)이므로 언마운트가 제대로 되지 않는 현상이 발생

Block이 유효하지 않을 때

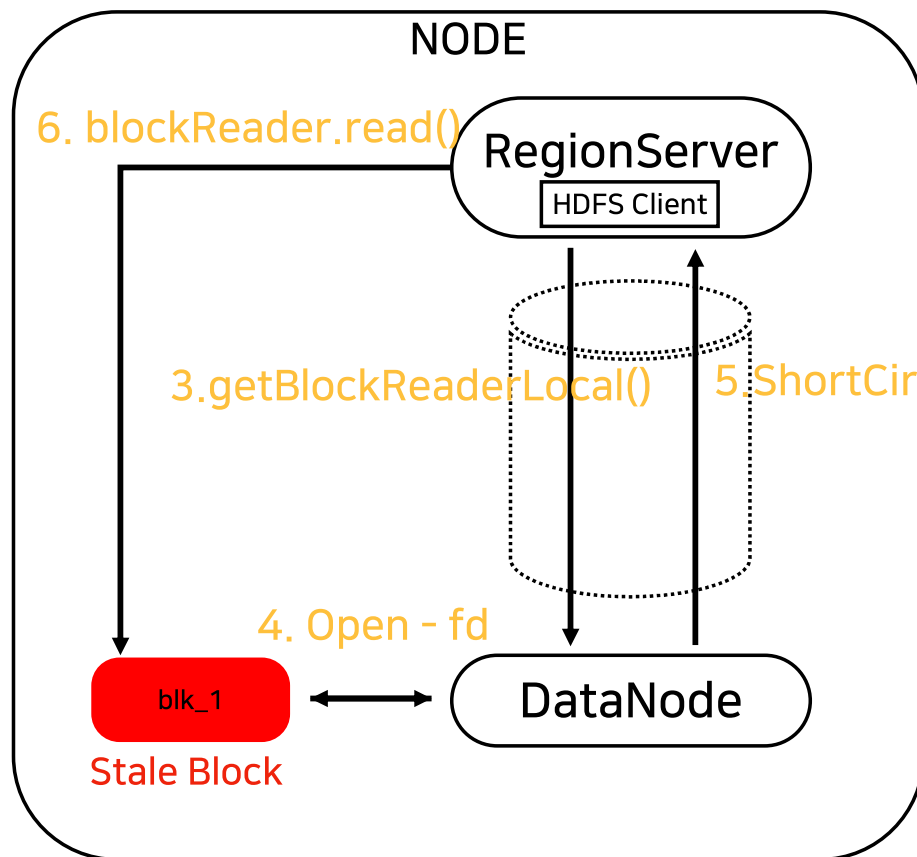


1. Open HFILE
2. locatedBlocks
[blk_1 : {Local DN, DN1, DN2}]

```
public class ShortCircuitReplica{
    private final FileInputStream dataStream
    private final Slot slot
    ...
}

public FileInputStream(FileDescriptor fd) {
}
```

Block이 유효하지 않을 때

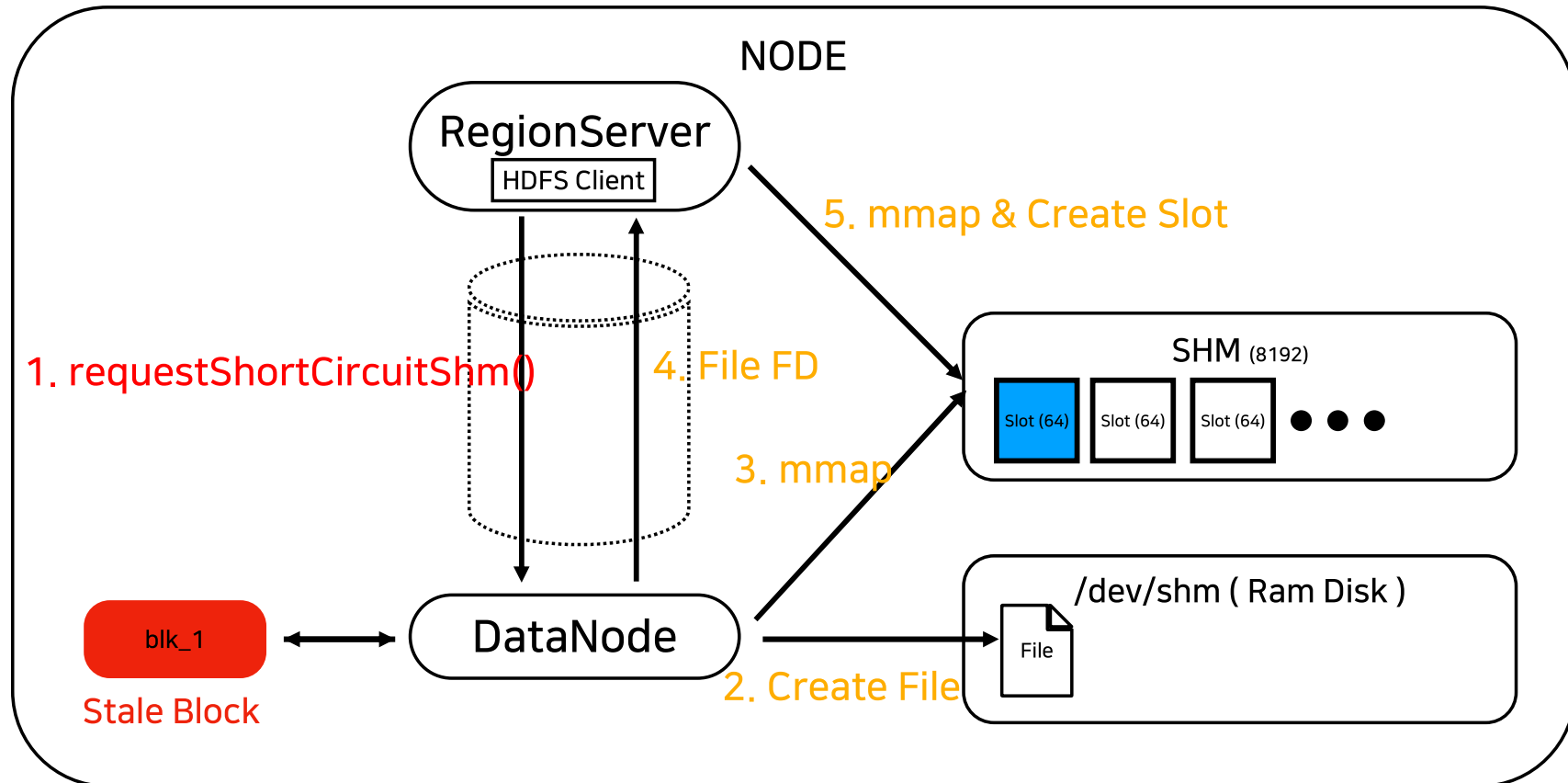


1. Open HFILE
2. locatedBlocks
[blk_1 : {Local DN, DN1, DN2}]

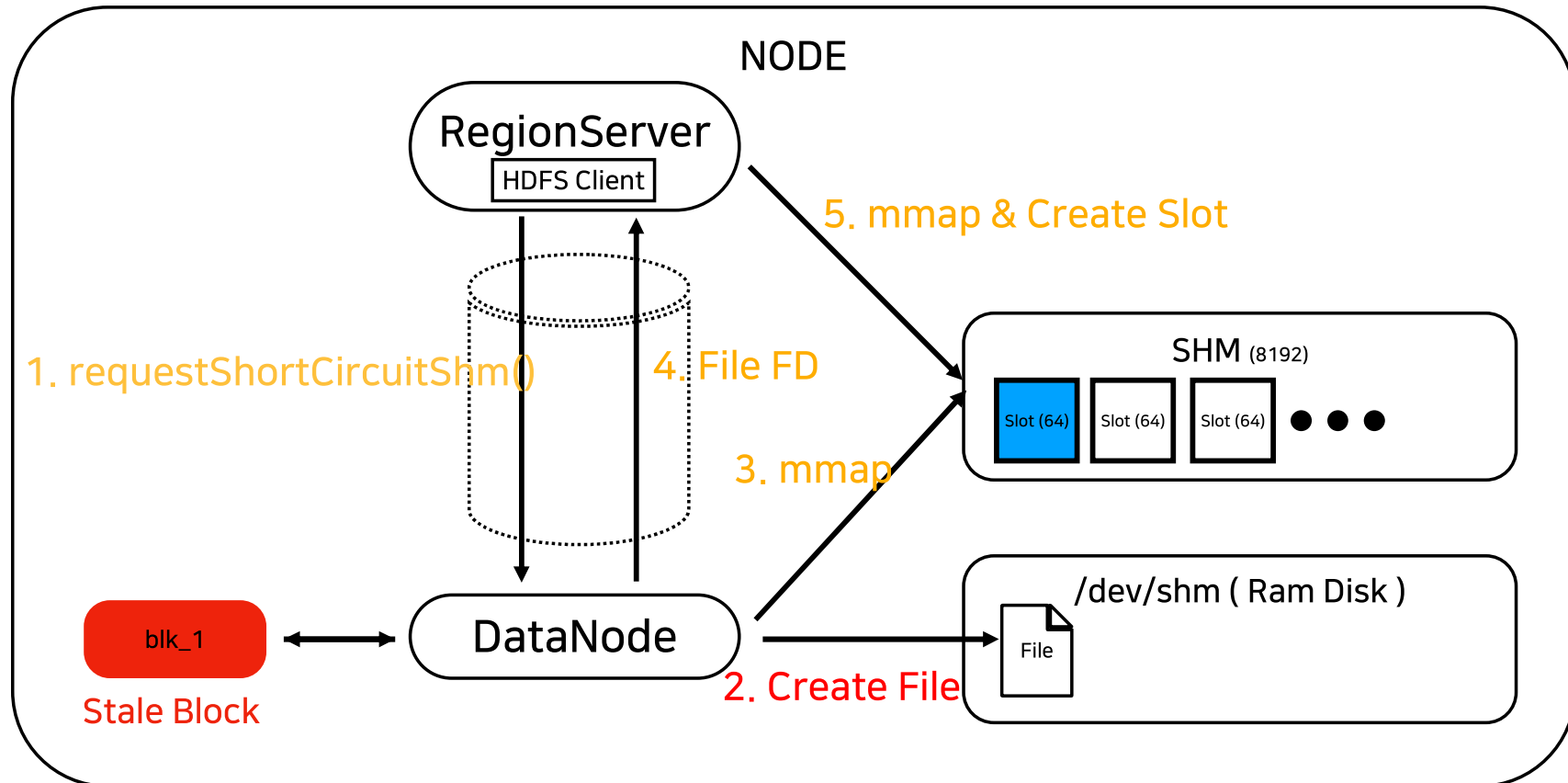
```
public class ShortCircuitReplica{
    private final FileInputStream dataStream
    private final Slot slot
    ...
}

public FileInputStream(FileDescriptor fd) {
}
```

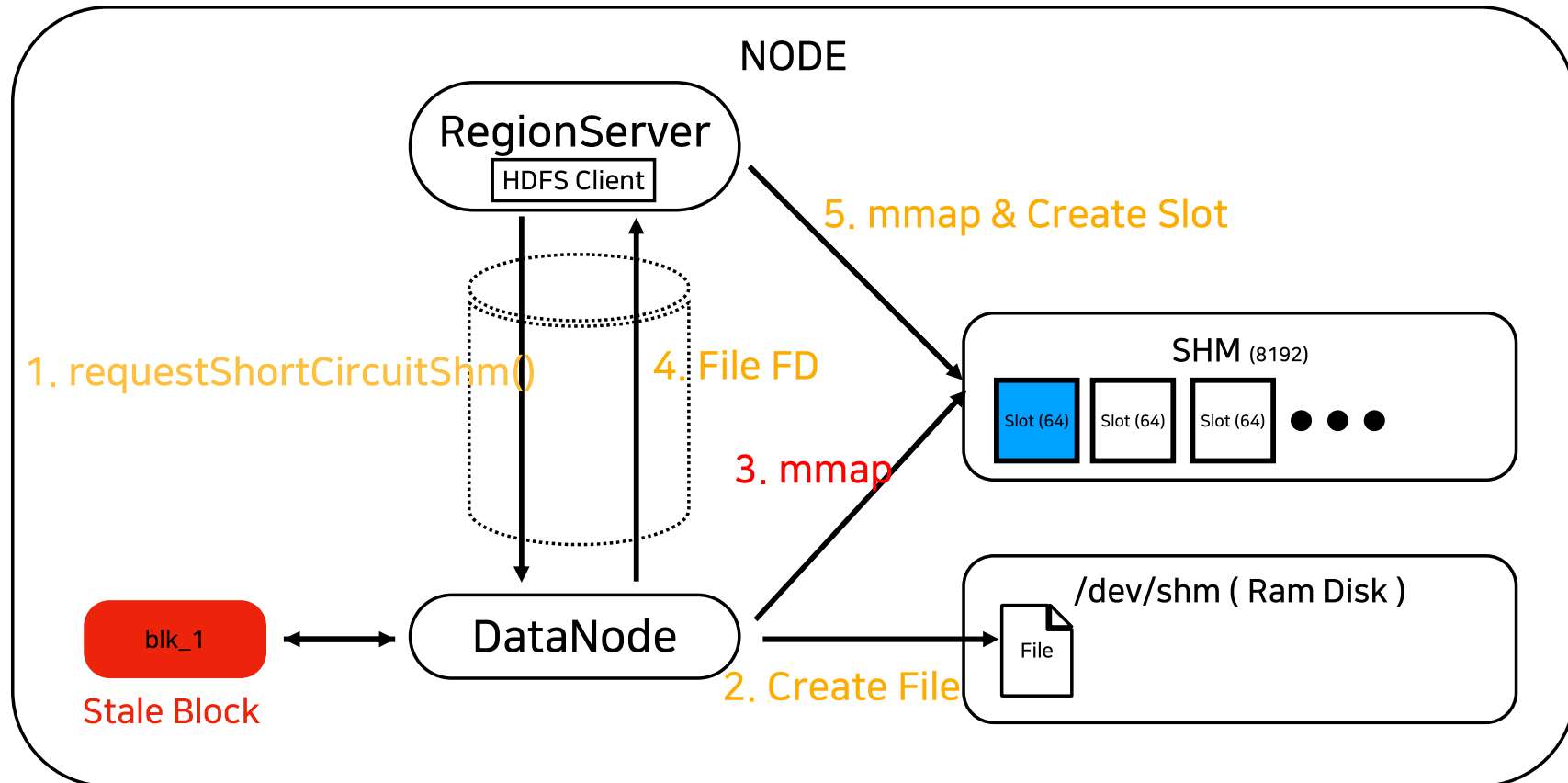

Block이 유효하지 않을 때



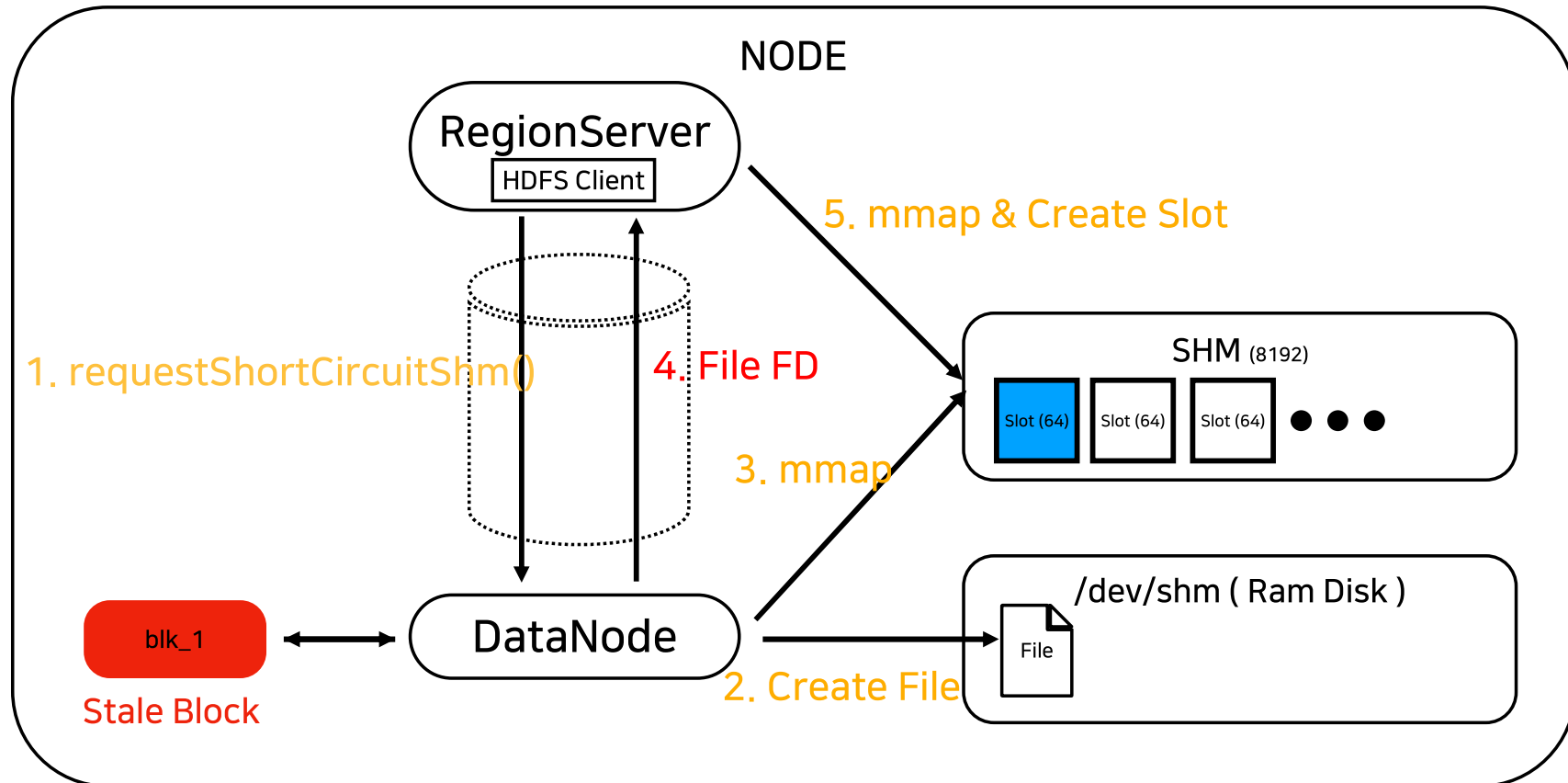
Block이 유효하지 않을 때



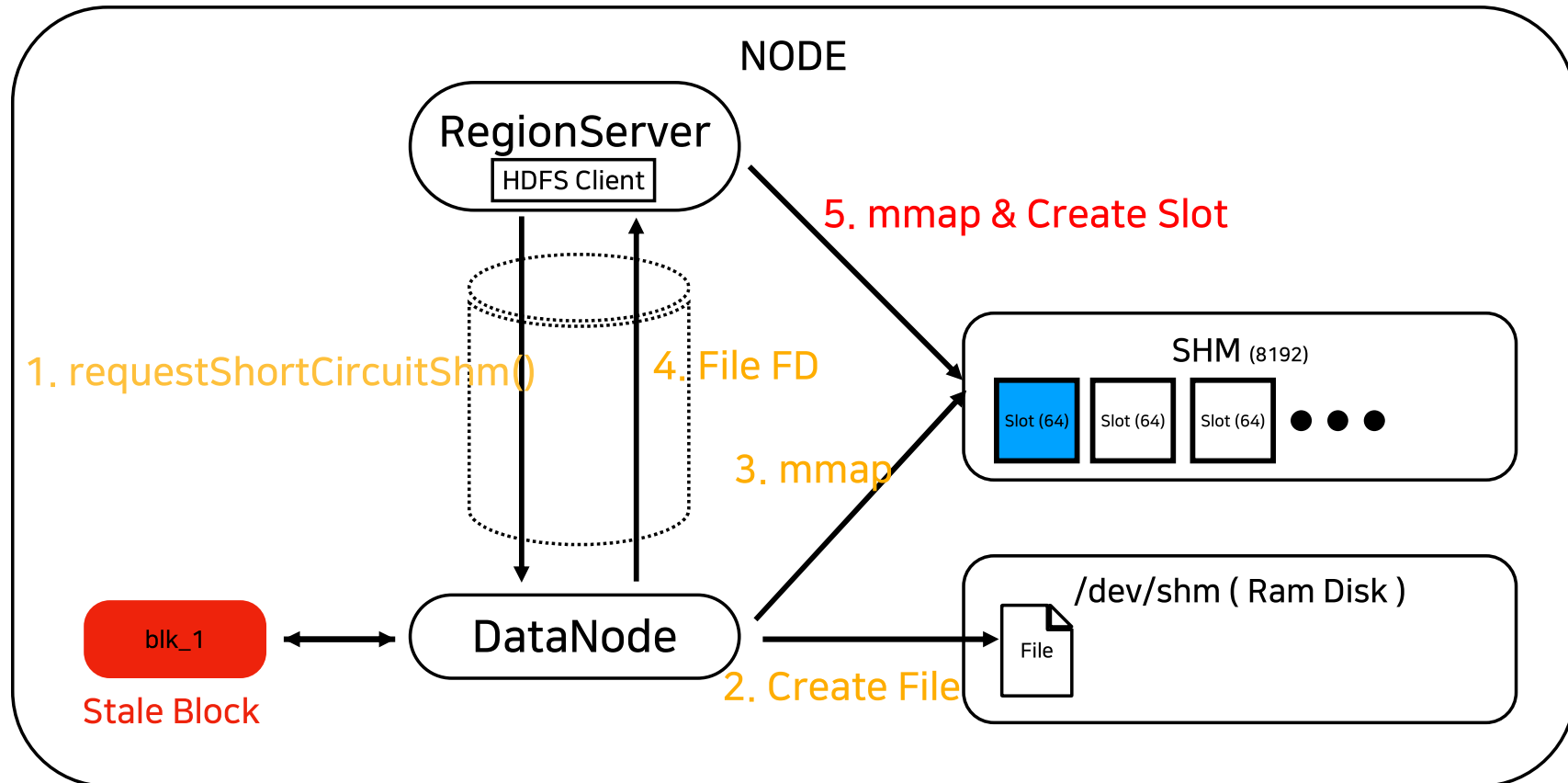
Block이 유효하지 않을 때



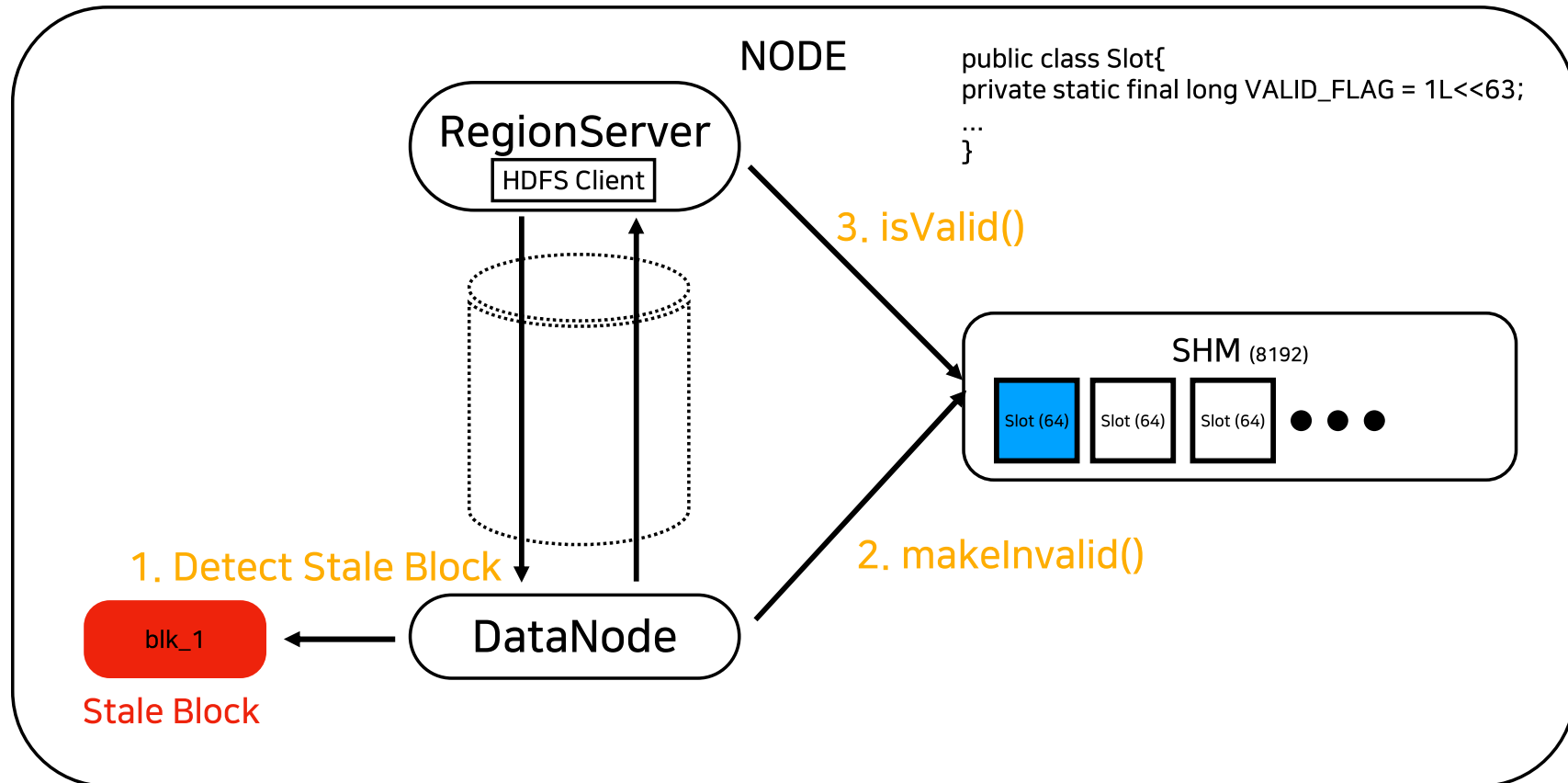
Block이 유효하지 않을 때



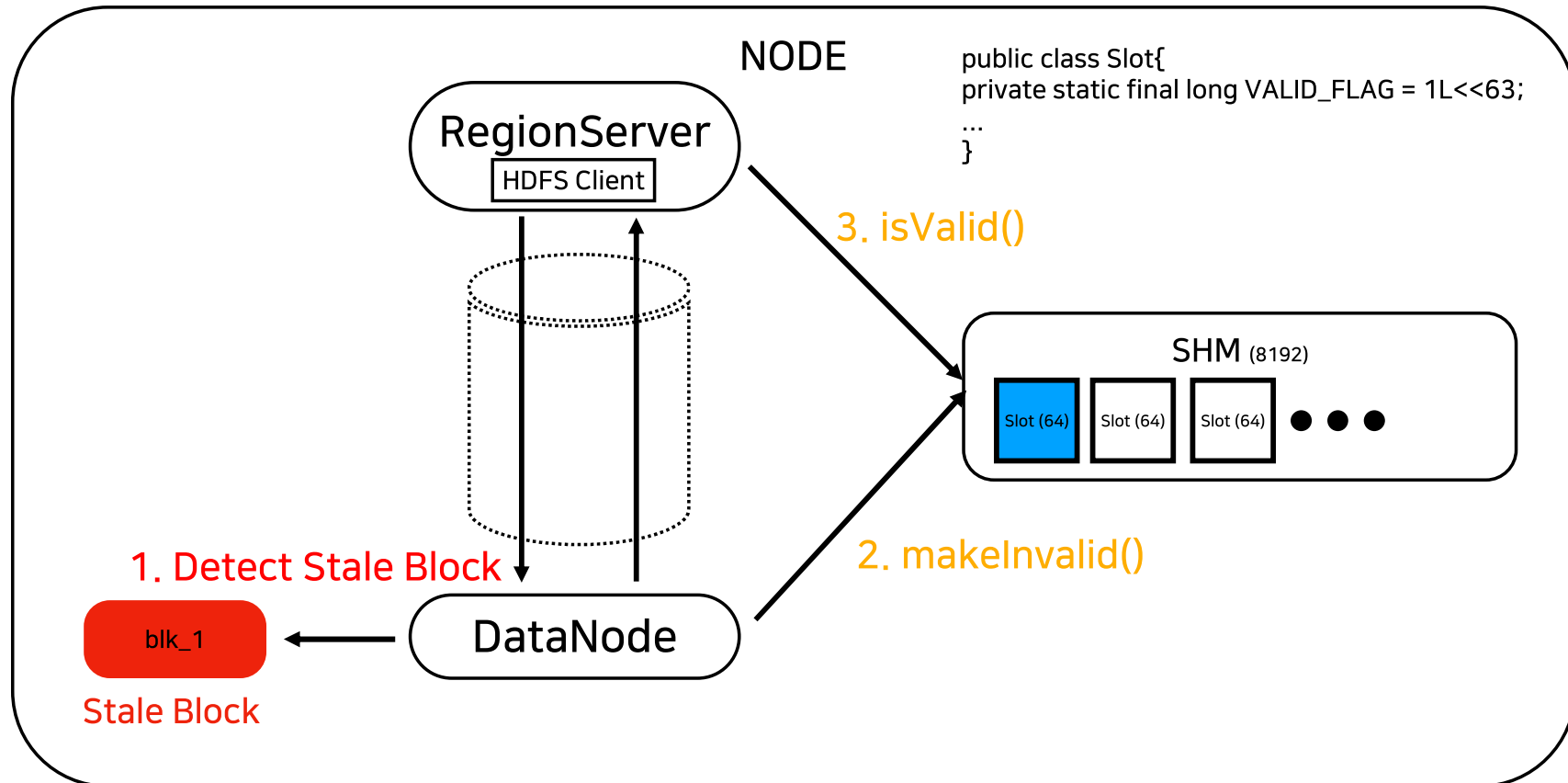
Block이 유효하지 않을 때



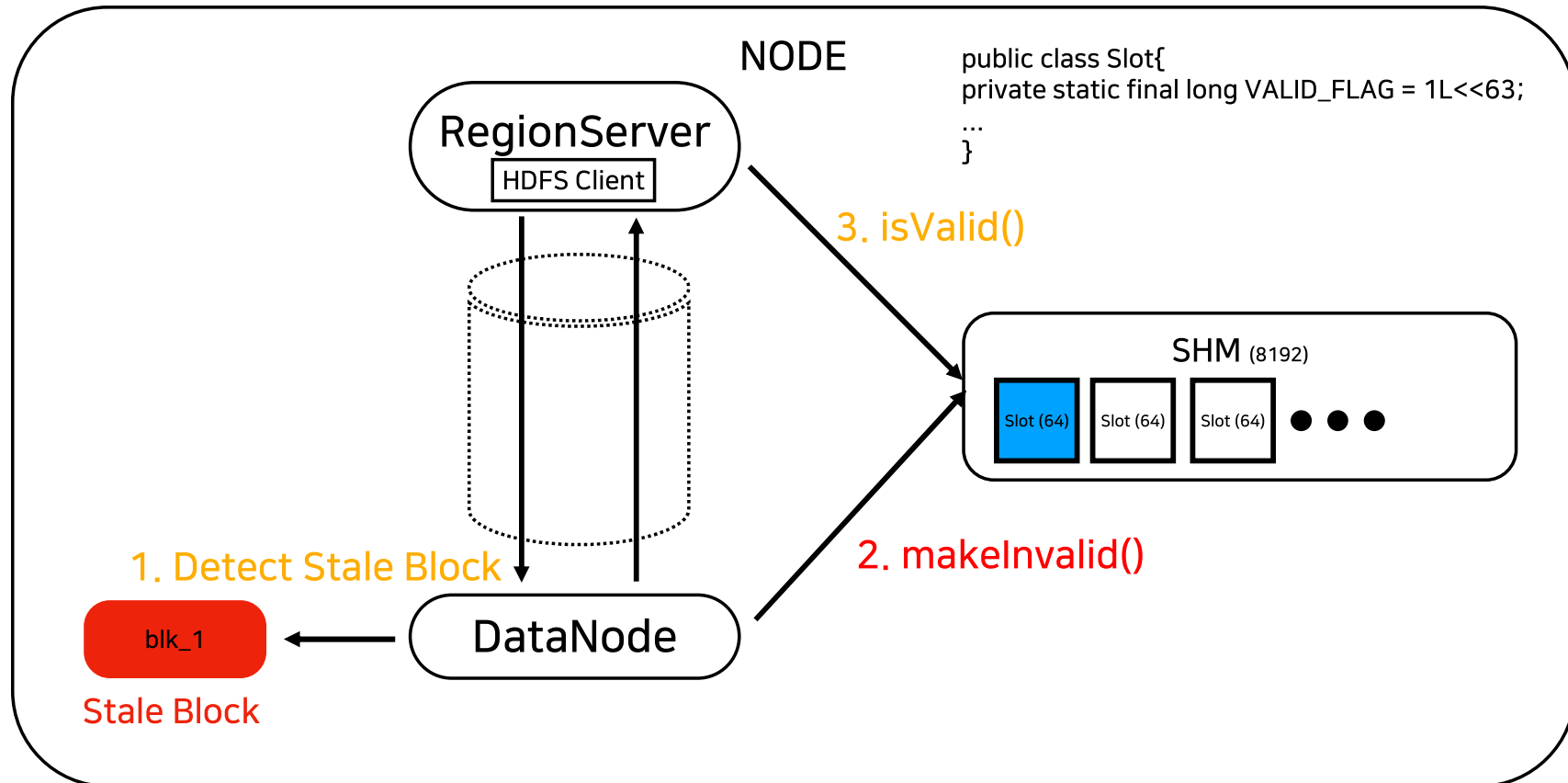
Block이 유효하지 않을 때



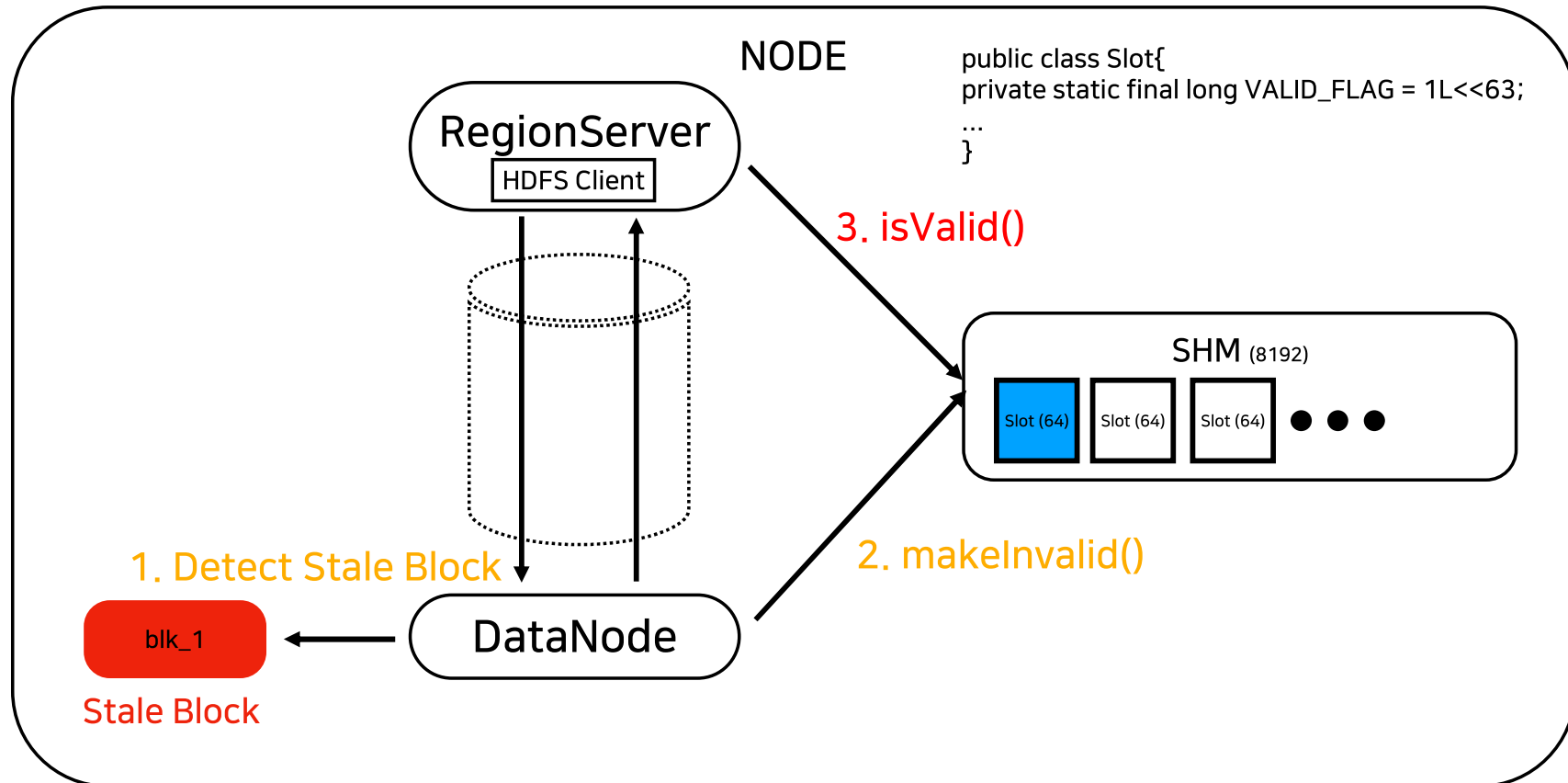
Block이 유효하지 않을 때



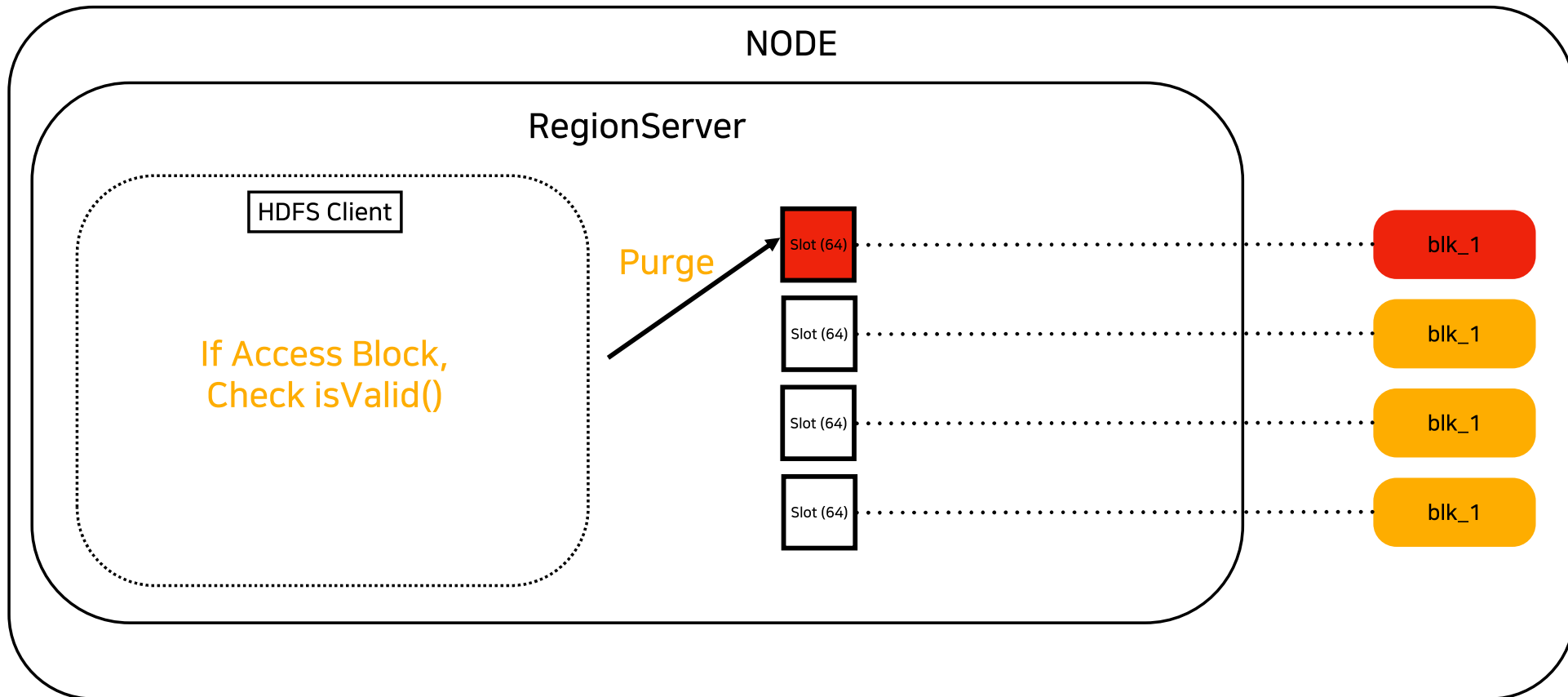
Block이 유효하지 않을 때



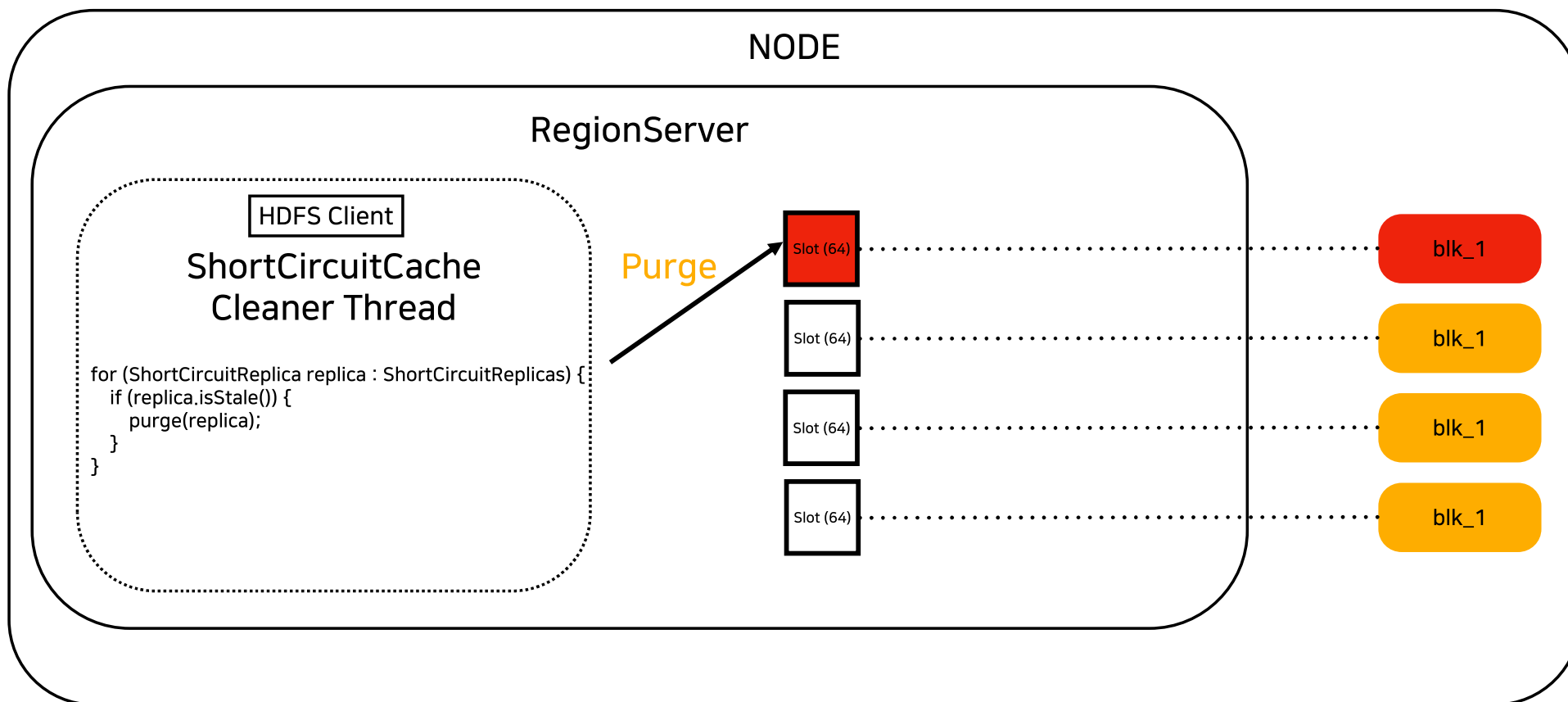
Block이 유효하지 않을 때



Lazy Purge



Aggressive Purge



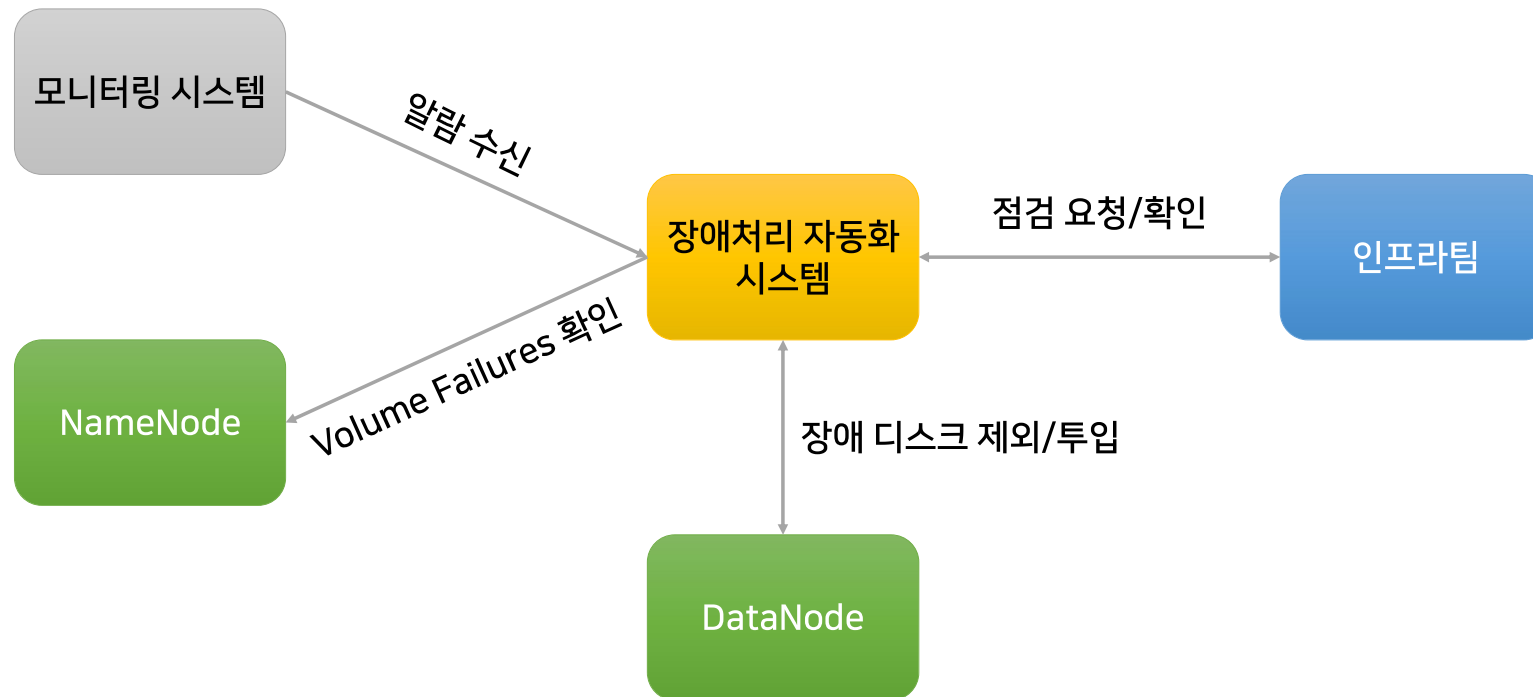
Aggressive Purge



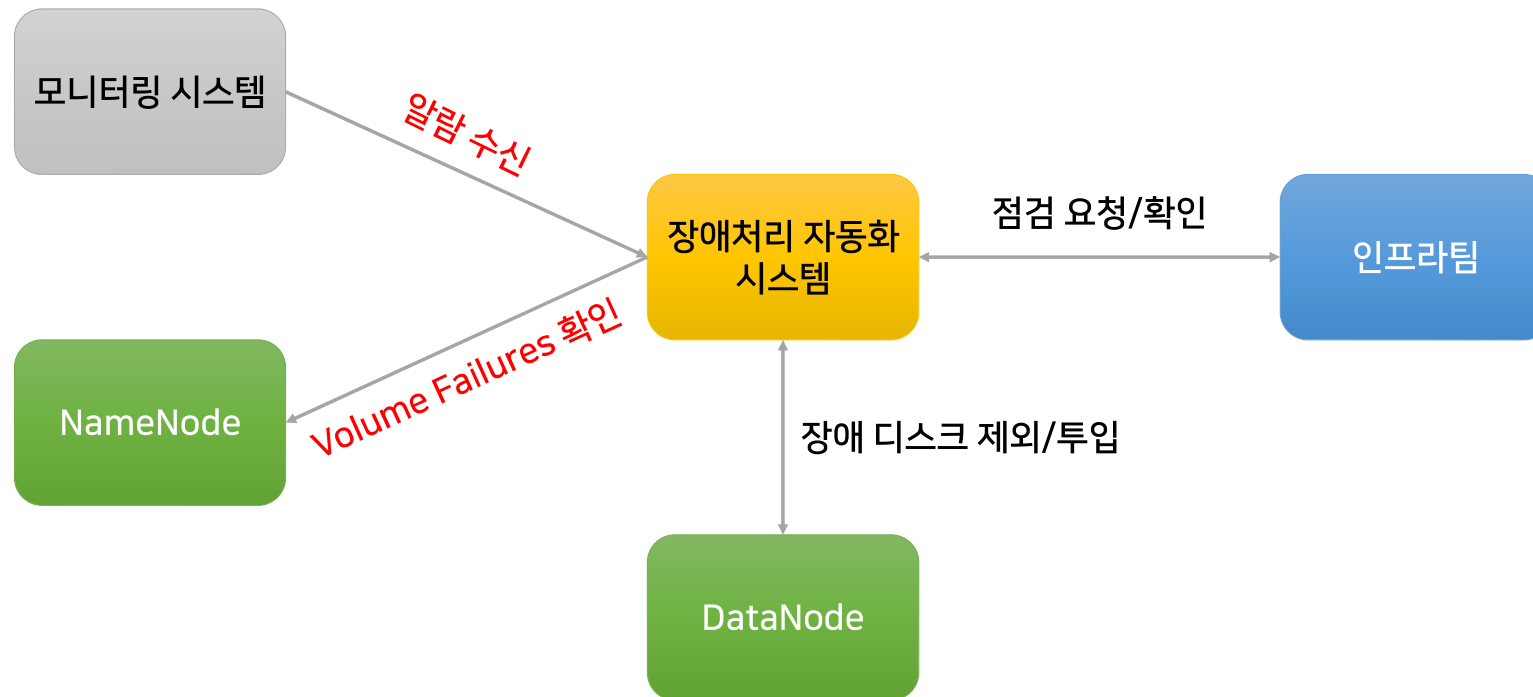
언마운트 불능 현상2

- HADOOP 3.0.0에서 발생
- Disk가 고장으로 메타데이터를 읽을 수 없을 때 볼륨 레퍼런스를 릴리즈 하지 않는 버그 (HDFS-15963)

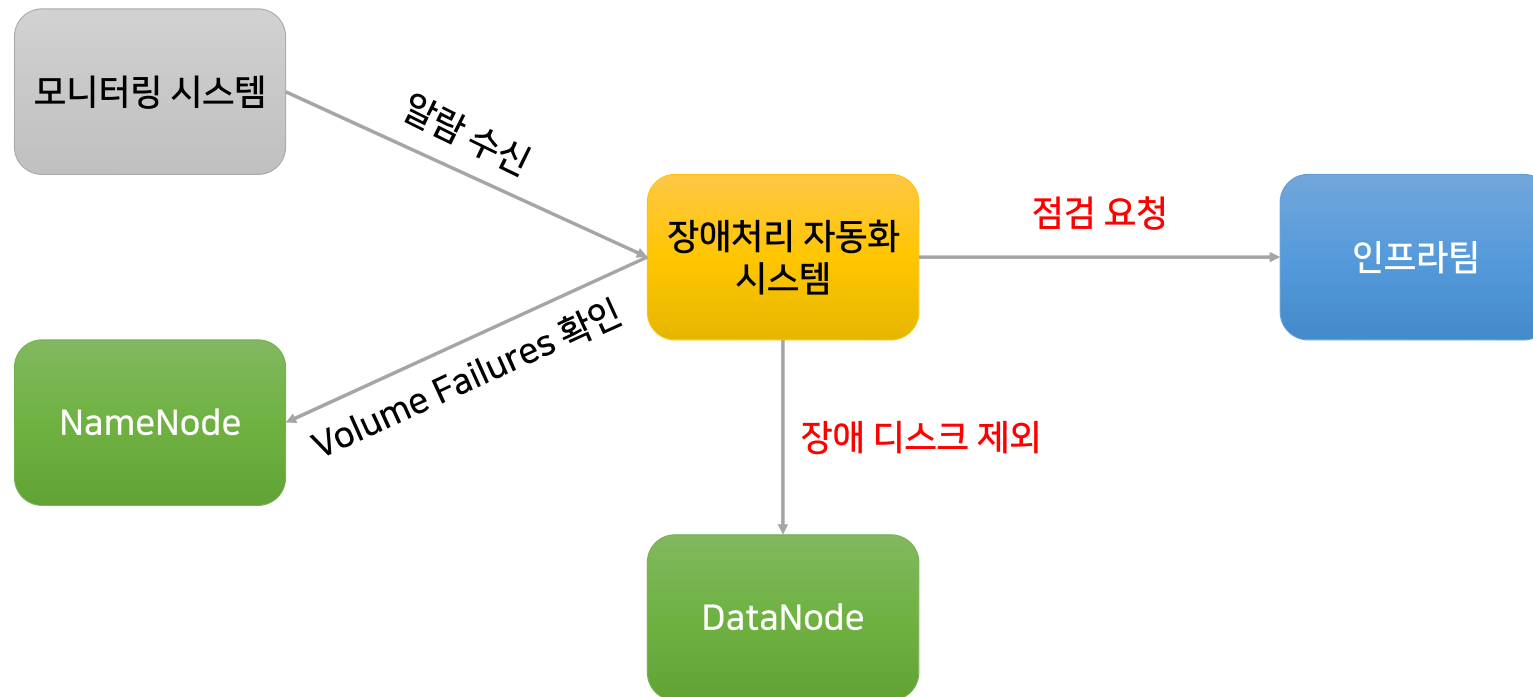
디스크 장애 처리 자동화



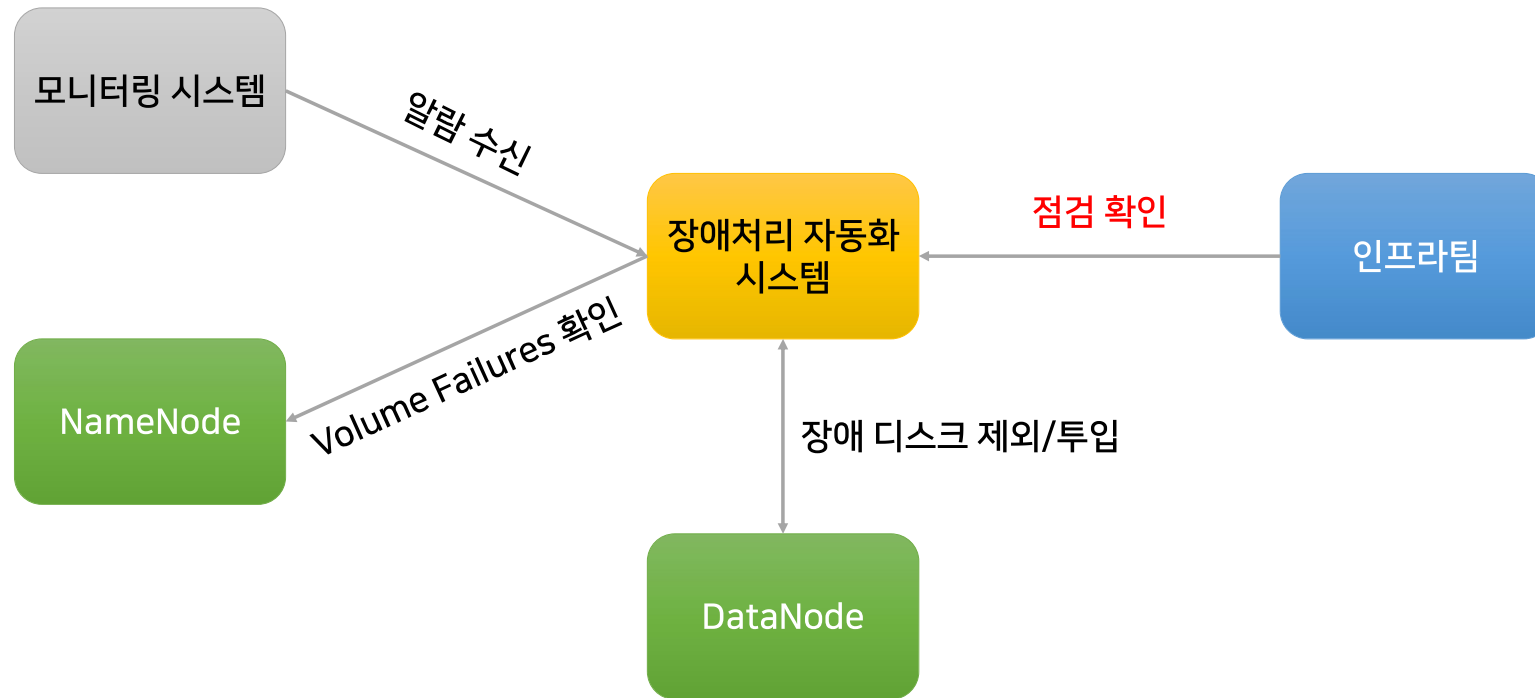
디스크 장애 처리 자동화



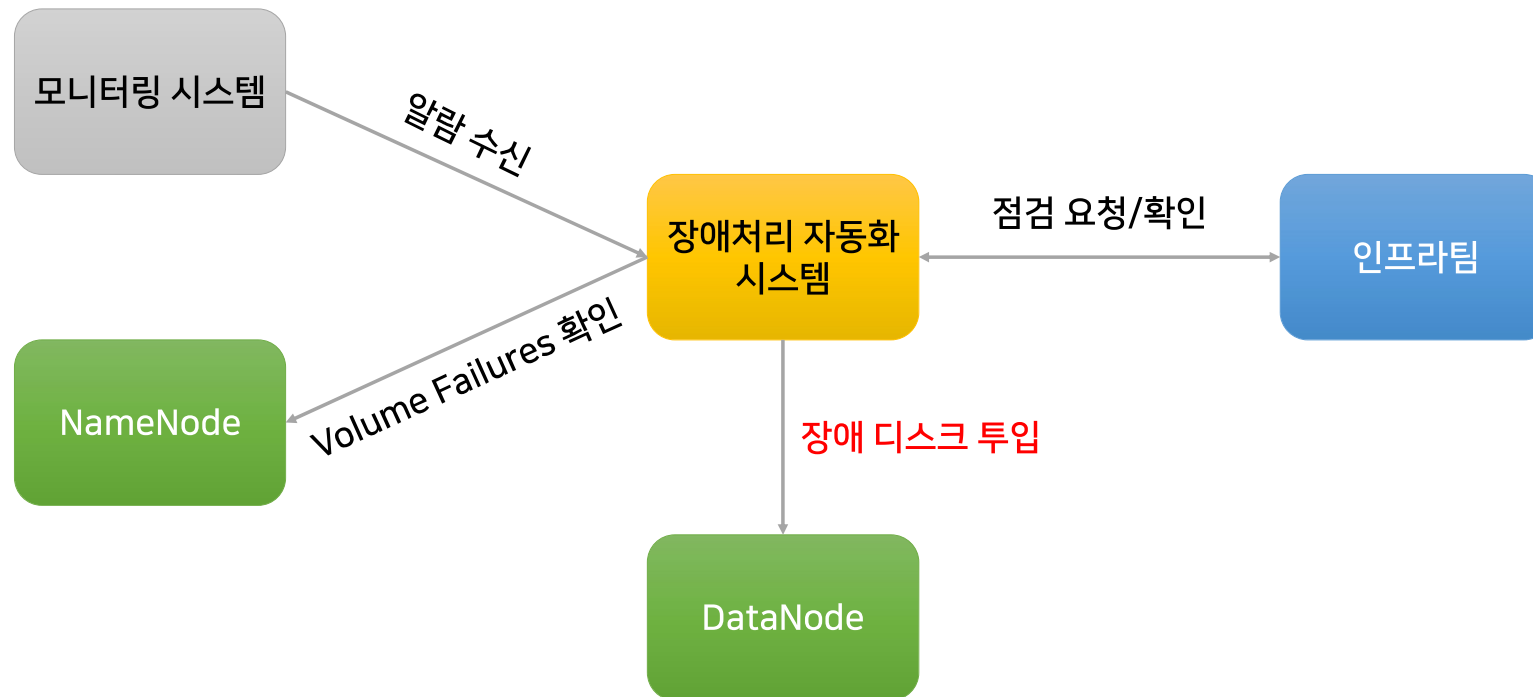
디스크 장애 처리 자동화



디스크 장애 처리 자동화



디스크 장애 처리 자동화

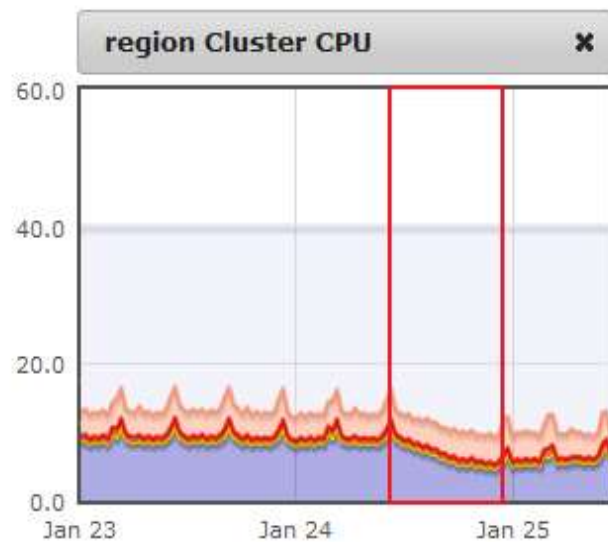


네이버 데이터 저장소에서 효과

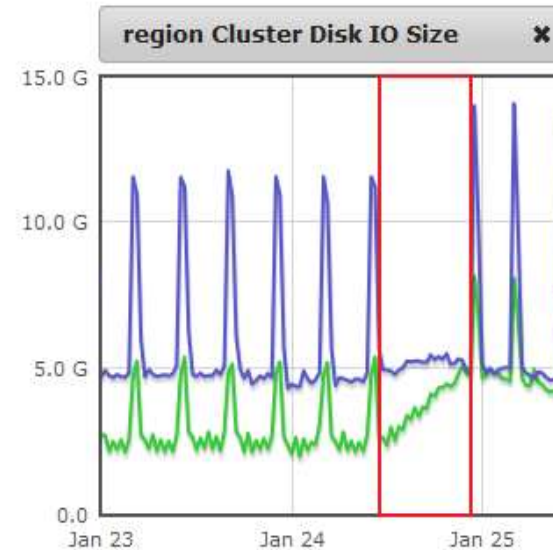
네이버 데이터 저장소에서 효과

HBase Regionserver 백여대


01/24 11:00 경 ShortCircuit 설정 적용 (Rolling restart)



파란색 user cpu가 40% 줄어듬



녹색 Disk Read가 2G/s 늘어남



We're hiring



AI & Data Platform 소개

<https://naver-career.gitbook.io/kr/service/search/ai-and-data-platform>

경력 등록

<https://d2.naver.com/news/7591059>



AI & Data Platform

Service

Meta

데이터 거버넌스

Content

저장

AI

Search Service

AI

Log

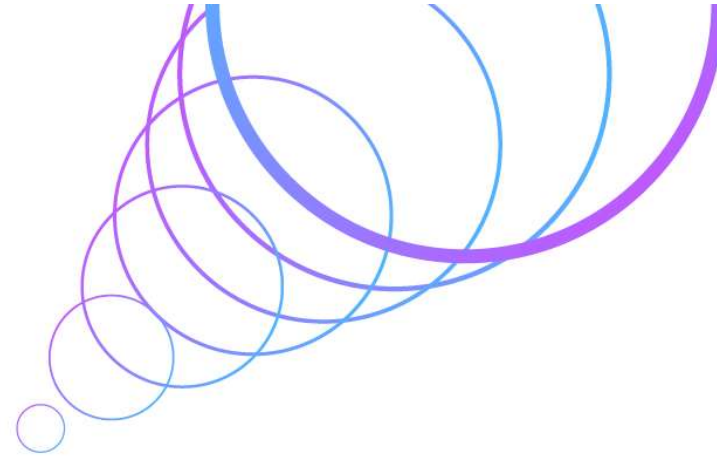
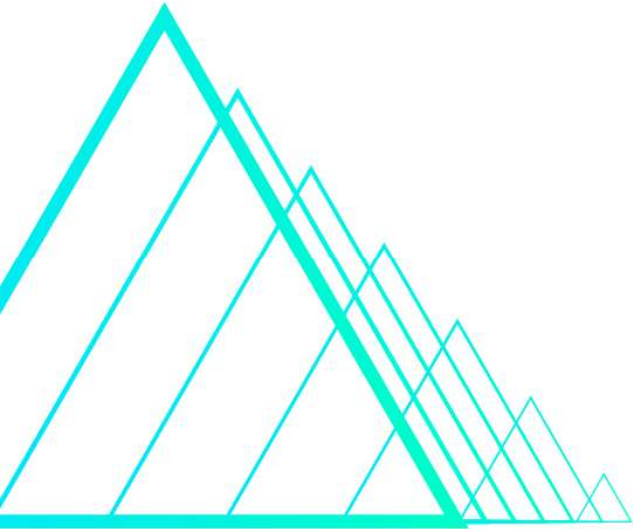
수집

연산/처리

Research

- Data catalog
- Data governance
- Data quality management
- Data store
- Data distribute
- ML training
- ML serving
- ML pipeline
- Processing resource
- GPU scheduling
- Log ingest
- Log real-time search
- Log analysis
- Log ETL





Thank You

