



# ch1.

ch1

Hadoop. - 다수의 컴퓨터에서 대량의 데이터 처리

NoSQL  $\Rightarrow$  Hadoop  
DB기록 분산 처리

Redshift.

Hadoop  $\approx$  Data warehouse.

하드웨어 통합 정도가, 확장성이 쉽지만

Hadoop에서 대부분 처리

주로 Data 만

시각화.

BI.

Data Discovery

빅데이터의 특성은 DW와 다르게, 다수의 분산 시스템을 혼합하여  
확장성이 쉬운 데이터 처리 구조를 만든다.

③ 데이터 수평.

① stream  $\rightarrow$  stream 처리  $\rightarrow$  실시간 DB 등

$\downarrow$

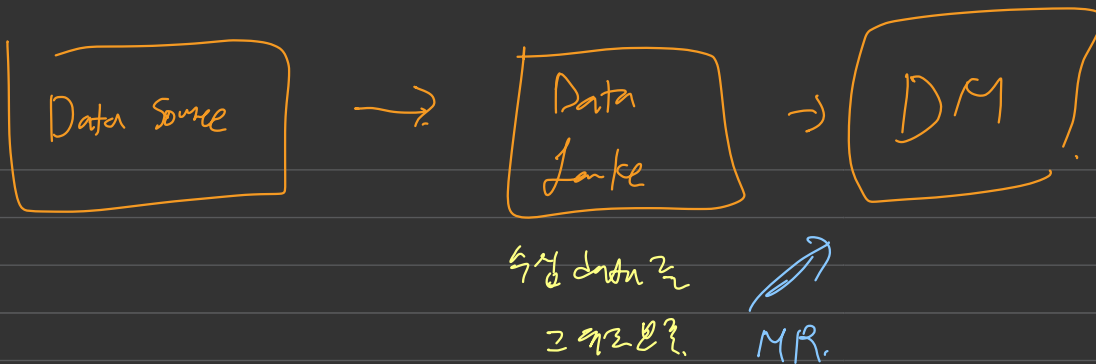
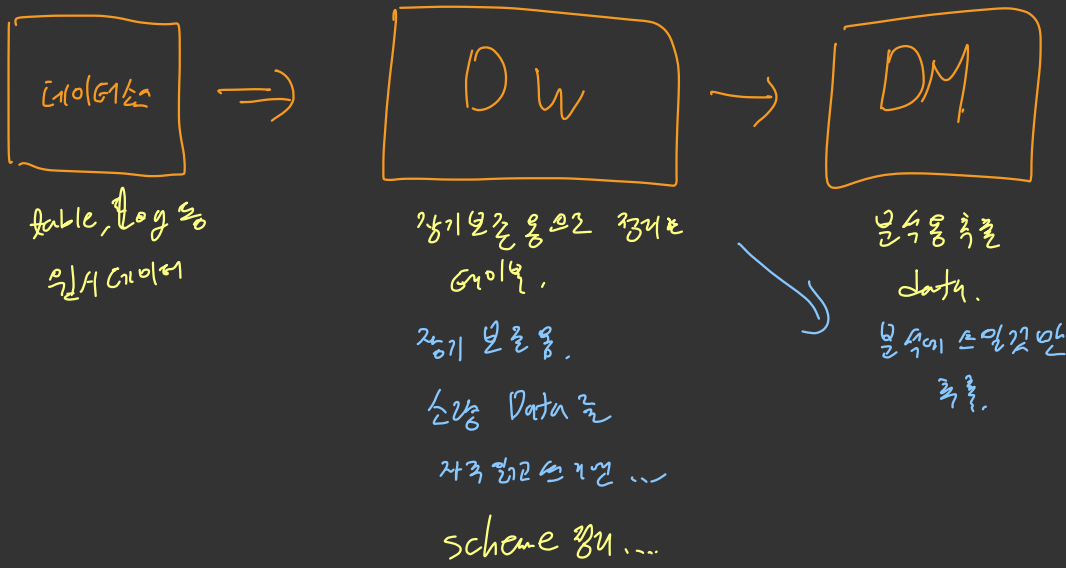
$\uparrow$

② Bulk.  $\rightarrow$  분산 storage  $\rightarrow$  분산 데이터 처리

SS = 객체 저장 구조

Map Reduce.

배치 처리 - 어느정도 정제된 데이터를 흐름으로써 제공



## BI 도구

① BI 도구를 사용 특정 Data 분석 가능

간편.

Data 시각화 가능

② DM을 통해, BI 개발하기

편리함.

어떤 Data도

가용 가능

③ 시스템의 BI 도구를

도입, CSU 도입 가능

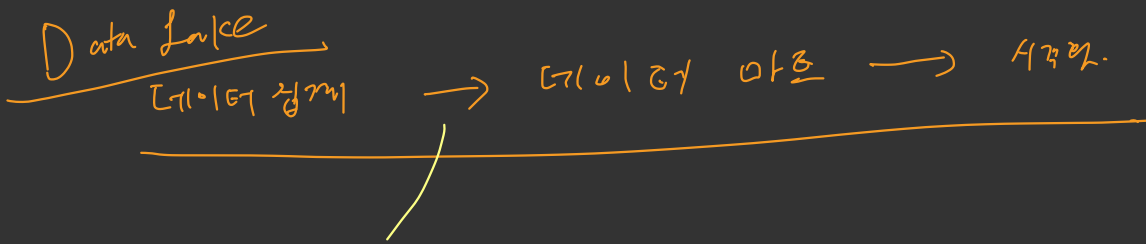
# CH 2.

Q 크로스정제:

크로스 table = 4값들이 보이기 쉬움

transaction table = 데이터가 탐색 가능

크로스  
정제



시각이 편리하므로  
해결하기 쉬움.

1) 모든 데이터를 메모리에 올림

2) 양쪽 라운드 MPP

가능하면 Data를 자주 읽기, 메모리 사용 최소화.

분산 데이터를 읽기 위해서는 매퍼 (병렬 처리) MPP.

Redshift, Google BigQuery

Column Oriented.

• 열지향 (column-oriented)

↳ 관련양 → Disk I/O 줄이기.



• 행지향 → 데이터 추가가 효율적.

↳ 검색을 고속화하기 위해 index를 만든다.

↳ index가 많다면 저장되는 모든 데이터를 load해서 레코드를 찾아야 함.

↳ 많은 I/O.

↳ Data 분석에는 어떤 column 사용될지 몰라 index는 필요치 않음...

☆ index에 의존하지 않는 고속 기술이 필요.

• 열지향 - 미리 column 단위로 정렬되어서, 필요한 column 만을  
로드하여 디스크 I/O를 한다.

↳ 접해는 고속이지만, 저장하는 시간이 많음...

↳ 압축률을 높임!

Redshift (하드웨어 최적화)  
(하드웨어 최적화)

• MPP 방식,

↳ 병렬로 여러 작은 디스크를 분할해서 가능

↳ 하나의 쿼리를 다수의 작은 task로 분할, 병렬로 실행

↳ 병목 현상 발생하지 않게끔, 데이터 크게 분할.

대리행-

↳ 분산 시스템에 모은.

시각화 도구

리뷰 DB에서 가져옴.

• Redash

• Superset

① 데이터 스토리지

• 마인스 2 종류.

② 쿼리 실행, 데이터베이스 연결

• 시계열은 Druid.

③ 그라운드 데이터 스토리지

④ 시각화 도구 어떤 것을 사용할까?

1) 스프레드시트

2) 데이터 스토리지 (Redash)

3) Jupiter notebook

SQL

동적인 쿼리, 대량 데이터

Table report

☆

기이데이터

① 인터페이스

② 쿼리 실행을 위한 높은 수준의 데이터베이스

③ 데이터베이스로 연결된 데이터

바탕 화면으로 동적 쿼리, 대량 쿼리 실행 X 실행 ↑

☆ 데이터 마스터 기반 구조 - 시각화에 필요한 데이터만 받은 데이터베이스

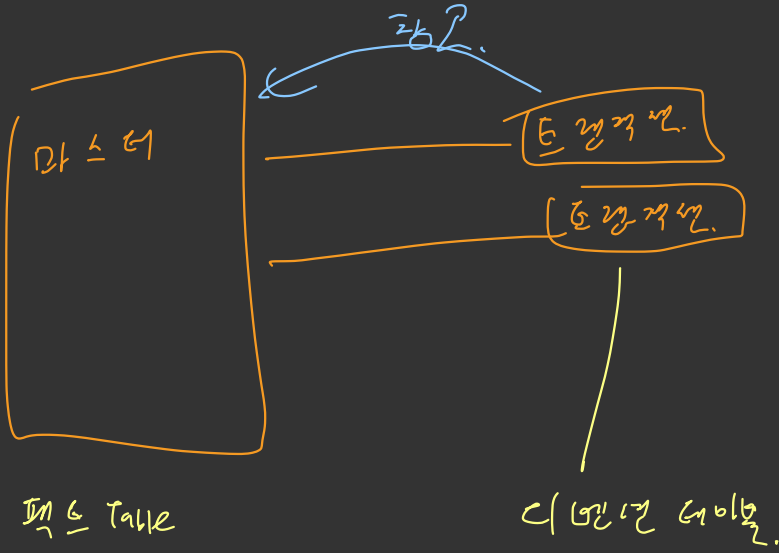
시각화에 적합한 데이터 만들기 (OLAP)

다차원 데이터.

모든 OLAP 데이터

BI 도구 + MPP DB를 위한 쿼리

o 테이블 비정규화.



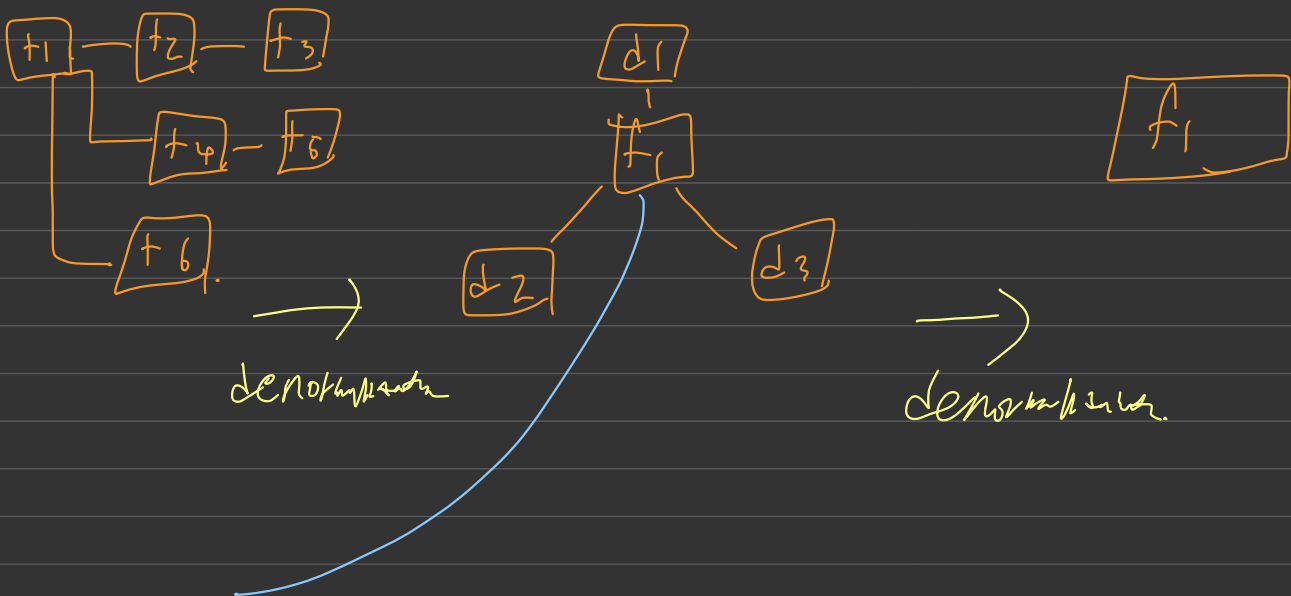
테이블 비정규화.

블록 파티셔닝을 사용함.

① 정규화 (정규화)

② 스키마 카탈로그

③ 비정규화



① 스키마 카탈로그 : 데이터를 테이블 이름과 속성으로

② fact table을 복제하여 구축함.

↳ K, ZD 마스킹 기술.

★ 데이터 웨어하우스 구축은 스키마가 중요.

① 다차원 모델 시각화에 데이터 테이블은 중요하지 않음.

↳ 표정만, 클러스터 구분.  
↳ X

## CH3 대규모 분산처리 FRAMEWORK

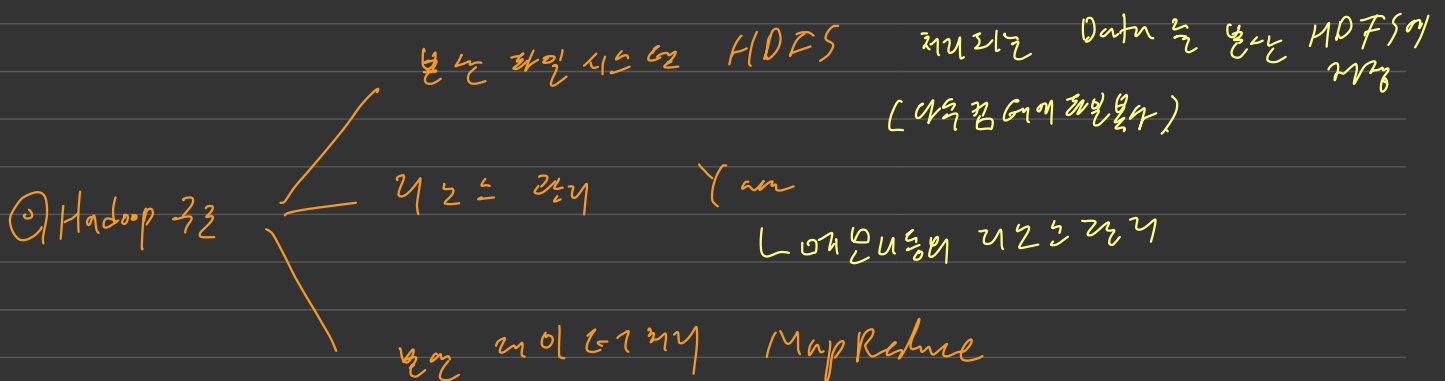
구현 방법.

- ① 쿼리 데이터 - schema table
- ② 비구조화. - text, binary, image
- ③ schemaless (비구조화) - CSV, JSON, XML.

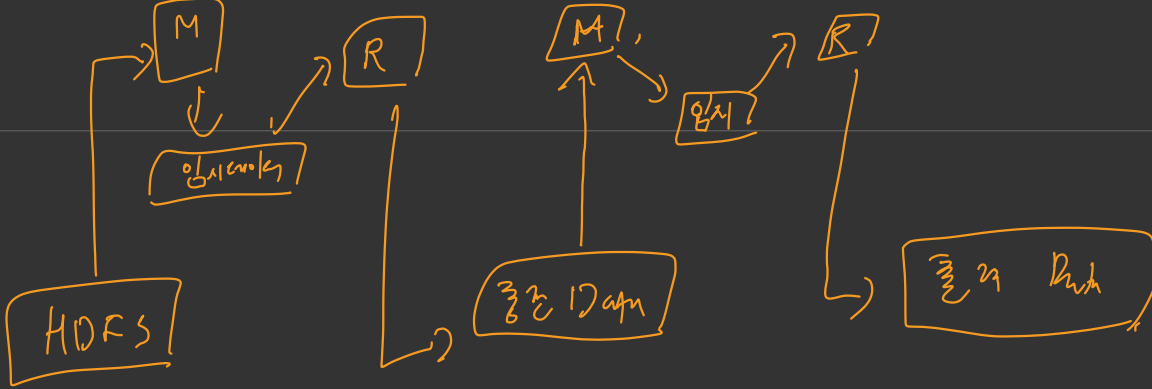
★ 비구조화 데이터를 읽어들이고 처리할 storage로 변환할 것 필요함

(데이터 가공 / 압축을 위해 많은 Resource...)

↳ Distributed Hadoop.







↳ 해결하기 위해 사용된 Tez (Hive 2.4.4) 중간 Data 없앴.

↳ 사용 가능한 Resource 사용

Task X

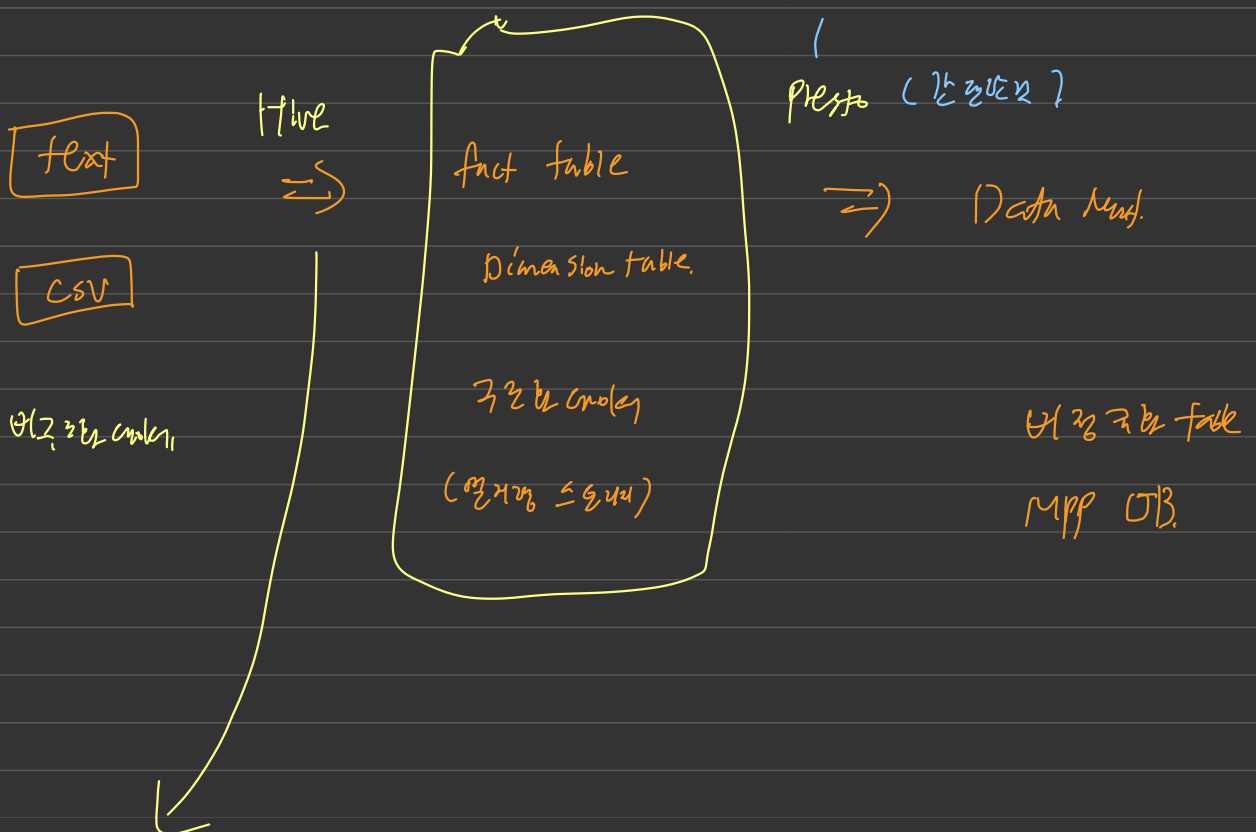
② Spill

↳ 사용 가능한 많은 Data를 RAM에 저장함.

↳ MR을 대역.

★ 데이터 마트 구축의 Pipeline

Hive 쿼리 엔진도 가능 (연산 시간 절약)



① 지름만 fact table로 만들기.

[17] 2023 04 27 2023 fact table 2023

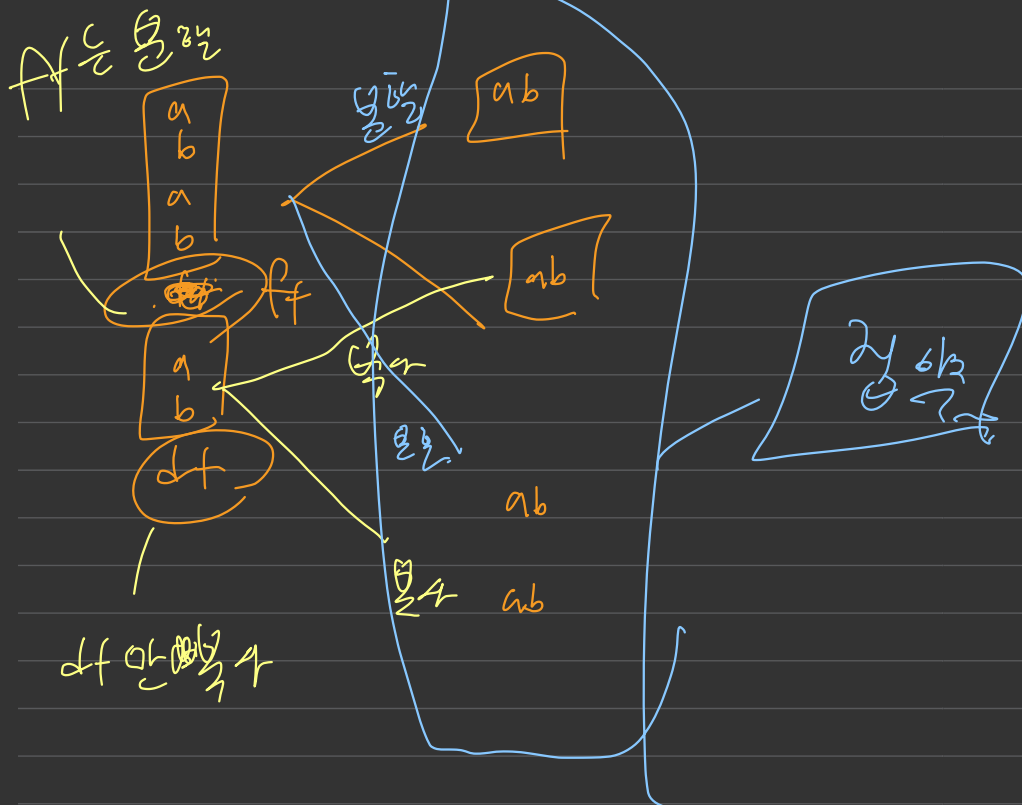
② Data 2720 크기 22  $\frac{2720}{22}$ 로 나누어 분석.

② Preto.

↳ Hive Meta store에 table & metadata!

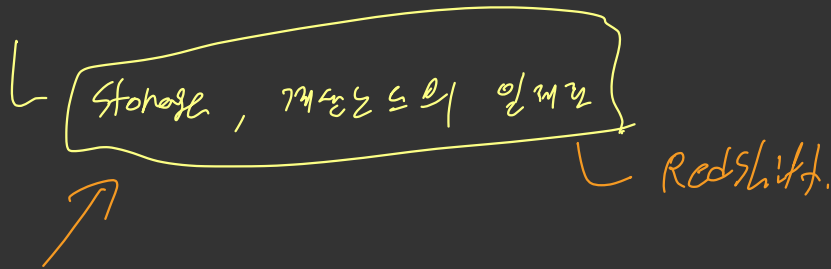
우리분처럼 리얼의 실링 계약을 맺는

— 불순 비율 (MR) 및  $\rightarrow$  3D 케이스의 결함도 식별.

$$L \cap \sum_{i=1}^n \mathbb{C} e_i = \sum_{i=1}^n \mathbb{C} e_i \quad \left( \sum_{i=1}^n \mathbb{C} e_i \right)$$


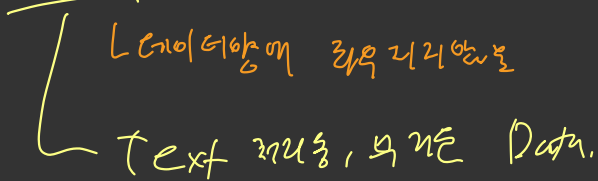
## ① A 프레임워크 선택하기

① MPP DB - 많은 수의 데이터의 고속 집계.

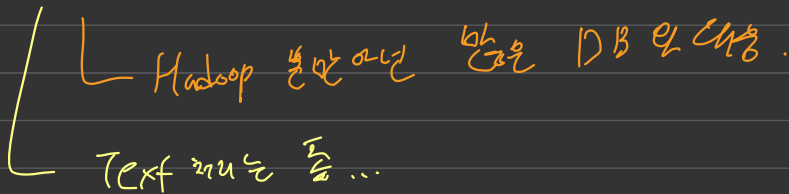


MPP - BI 연결.

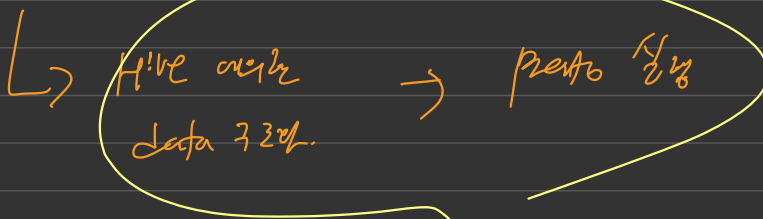
② Hive - <sup>TEZ</sup> 높은 확장성.



③ Presto - 속도 중요.



④ Spark!!!!



하둡의 script로 가능!!!!

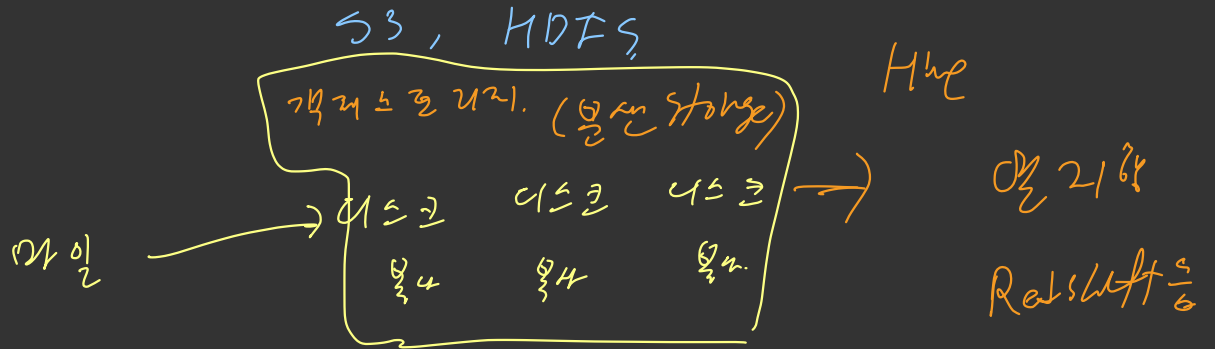
즉! ETL도 하둡에서 바로 ④ SQL로 집계 결과로 보내는 등 가능!



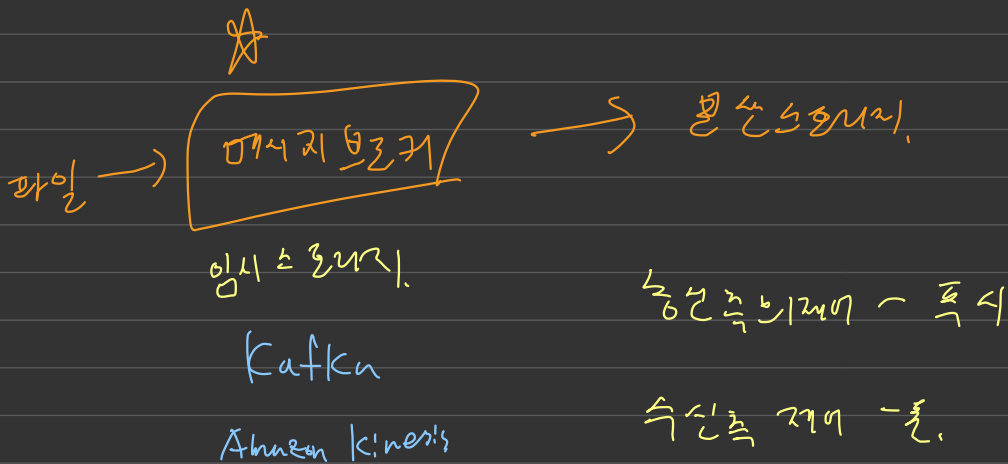
# Ch 4 빅 데이터의 특징

① 빅 데이터 스트리밍 형태의 Data 특성

• 객체 스트리밍 데이터 저장 - 분산도 처리에 대안되는 기술이다.



★ 스트리밍 형태의 Data 전송.



② 시계열 Data의 특징

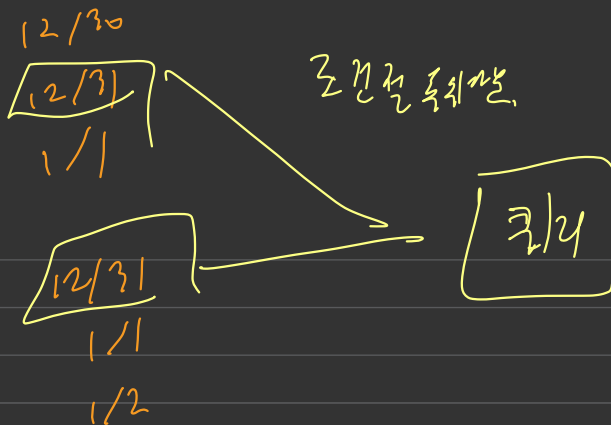
/

① 시계열 index.

↳ 이벤트 시계열 index

② 시간 점록치

시간 점으로 group by



1/1

1/2

1/3

여러 가지 시계열 table 사용 가능

③ 이벤트 시계열의 특징

✗

④ DM을 이벤트 시계열로 정렬. — DM 시계열 만들기

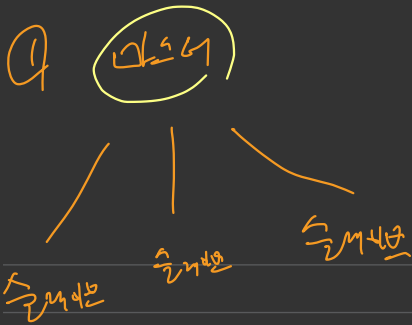
# NO SQL DB

① 분산 KVS.

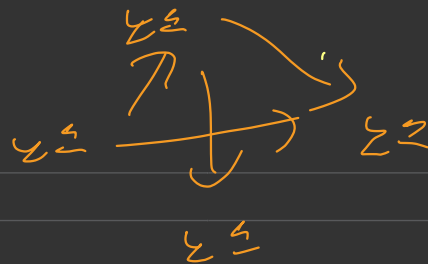
↳ 모든 데이터는 키-값 형태로 저장.

↳ 노드간의 복제된 쿼리 처리.

Client



② P2P



DynamoDB - P2P 형의 분산 아키텍처

↳ 사용자에 대해서 기록 변경 ↑, OK!

## ② Wide column store

↳ 모든 KVS를 보지 않음, 2개 이상의 컬럼의 데이터만 저장함.

|       |       |         |
|-------|-------|---------|
| Row 1 | col 1 | value 1 |
|       | col 2 | value 2 |
| Row 2 | col 1 | value 1 |
|       | col 2 | value 2 |

## ③ Document Store.

↳ 데이터 저장 형태

이름 데이터

22 저장 중 ...



