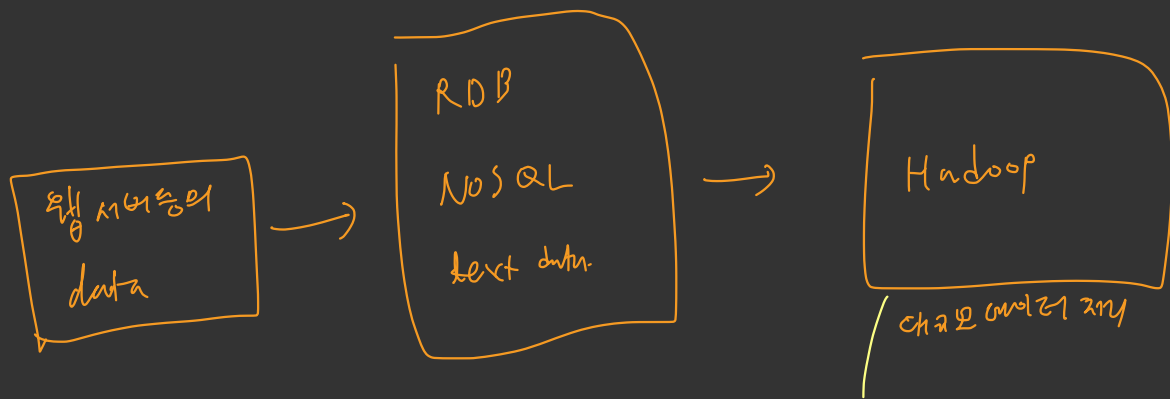


Raw
데이터 처리: 데이터 가공



다수의 컴퓨터에서 대량의 데이터
저장 시도함.

(*) 웹 서버 웹 페이지 등에서
모으는 방대한 data storage가
필요.

Hive - Hadoop에서 SQL 같은 쿼리 언어를
실행하기 위한 software

Hadoop - 모은 데이터를 다량에 걸쳐

DWH 환경 = 소프트웨어와 하드웨어가 결합...
하드웨어도 고려해야함.

★ 데이터 전송 방식.

메신 →. 버크홀 : 어떤가에 존재하는 데이터를 정리에 후로, DB중에서 가져옴

외출증가후서부 스트리밍 방. : 차례대로 생성되는 데이터를 끊임없이 보내는 방법.

모바일 앱등에 들어간 이벤트 데이터 및, embedded 장비
→ 시계열 (time series) 에 저장하므로 지금 무슨일이 일어나는지
알 수 있다

나름으로
→ 확장성이 높으므로 사용해야함

클라우드나지 → S3, NoSQL 등

객체 스트리밍

클라우드 Data 처리 (컴퓨팅, ETL 프로세스)

→ MapReduce, Data 처리량 ↑, 많은 컴퓨팅 자원

나중에 필요한 데이터가 있어서 리액티브 DB

방식

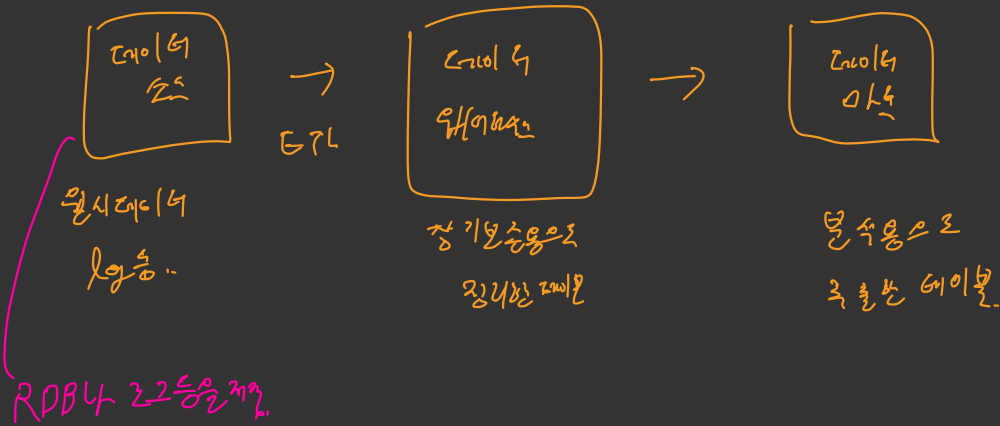
① SQL 사용

② ETL (리액티브 데이터 워크로드 제공)

→ 이 리액티브 DB에 맞게 파일 형식 변환

데이터 웨어하우스는 대량의 데이터를 장기 보존하는 것에 최적화되어 있다.

↳ 소량의 데이터를 자주 쓰고 읽는 데는 적합하지 않다



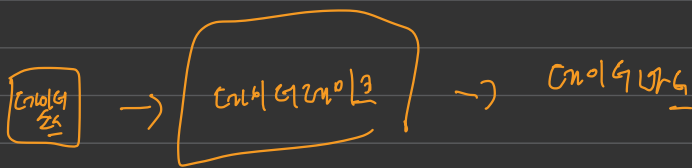
★ 데이터 레이크

모든 데이터가 데이터 웨어하우스를 거치지 않고도 만들어지지 않는다.

예) Binary data.

모든 데이터를 그대로 저장한다.

나중에 필요한 것들만 꺼내서 사용한다.



★ 데이터 레이크는 단순 스토리지이며, 그것만으로 데이터를 가공할 수 없다.

그래서 사용하려면 Map Reduce 이다.

↳ 분석에 필요한 데이터를 가공, 집계하고, 이것을 데이터 마스로
즉시 쓰기, 분석한다.

필수적으로 원본을 꼭 쓰는 Data

데이터
소스

→ Data Lake → Data Mart

← Data Engineer →

스작업으로 데이터는 집계 ⇒ Ad hoc 분석
[필수적인 데이터 원소]

Ad hoc 분석은, Data Mart 만들지 않고 데이터베이스나 웨어하우스에 직접 연결하는 경우가 많다.

★ 데이터 파이프라인의 큰 흐름은 변하지 않는다!!!

★ ① 저장할 수 있는 데이터 용량이 제한이 없도록

★ ② 데이터를 효율적으로 추출할 수단이 있어야 함.

★ DataFrame = 표 형식의 데이터를 다루는 객체.

↳ 스키마의 변경 안에서 가급적 집계를 본 줄로.

JSON, 텍스트 레이어들로 한번에 다룸

