

④ 빅데이터 스트리밍 형태의 data 수집.

빅데이터의 대부분은 저장성이 높은 분산 storage에 저장.

기본은 (객체 storage 이다)
 HDFS (Hadoop)
 S3 (Amazon)

⑤ 객체 storage는 다수의 컴퓨터를 사용하며 파일은 여러 디스크에 분산해.

↳ 중복해, 복제본만.

⑥ 데이터 수집.

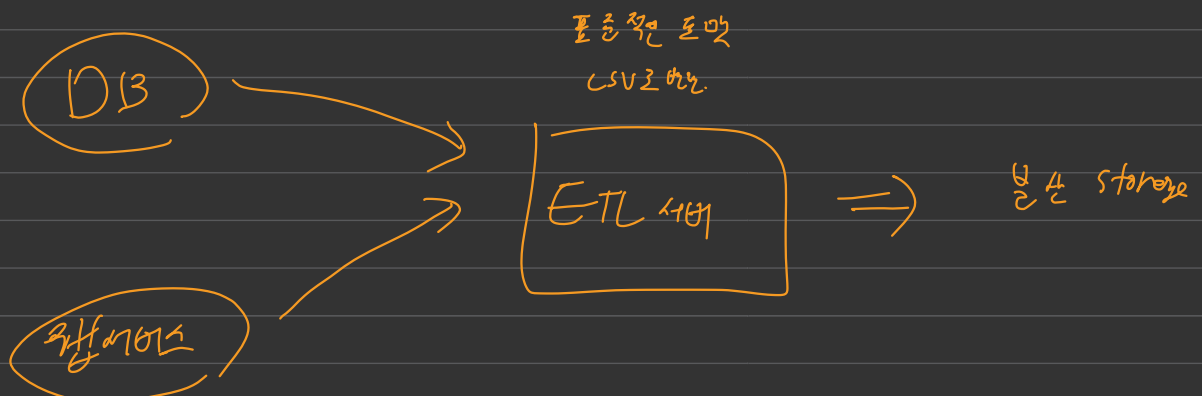
↳ time series 데이터, 작은 데이터는 즉방식 | 모아서 | 여러 큰 파일로 만들.

↳ 너무 큰 data는 적당히 나눈다.

객체 storage에서 효율적으로 처리할 수 있는 파일 크기는

대략 (1MB ~ 1GB)

⑦ 빅데이터 data 전송.



⑥ 성능 x 신뢰성 메시지 배송의 tradeoff.

메시지 배송에 대해 보내진 데이터를 분산시스템에 저장할수록 코어가 클수록.

대량의 메시지를 안정적으로 받기 위해서는 번번한 크기에 그 성능이 낮아질 수 있고
성능이 매우 높고 규모에 따라 성능을 얼마든지 늘릴 수 있는 다양한
필요하다.

★ 빅데이터에서는 종종 Data를 실시간으로 기록할 수 있는 클러스터가 필요하다.
(message Broker)

빅데이터를 위한 분산형 메시지 브로커는 오픈소스의 각

Apache Kafka, Amazon Kinesis 등 다양하다.

⑦ Push

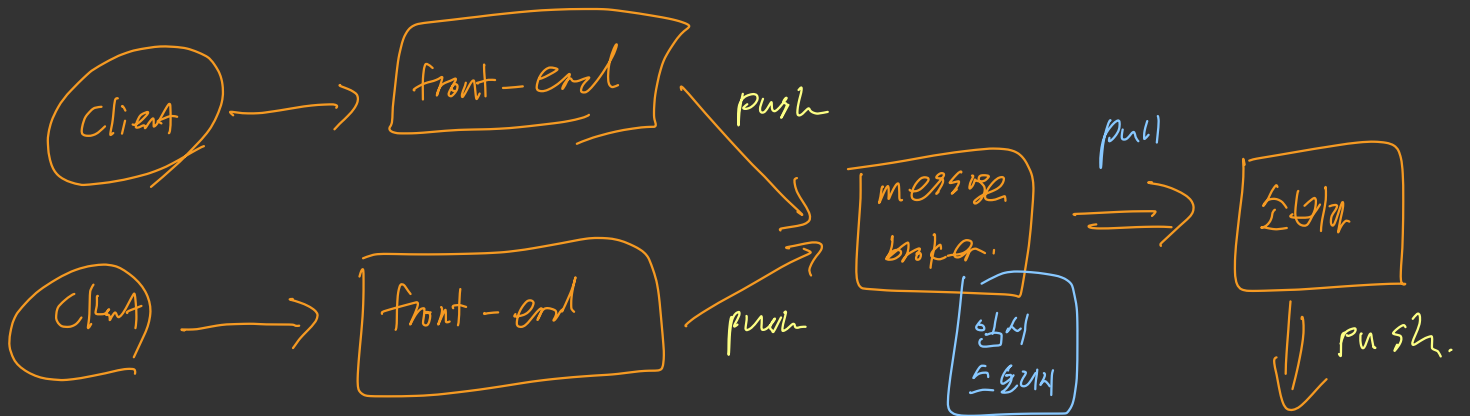
송신측의 메시지로 데이터를
보내는 방식

vs

Pull

수신측이 스스로 데이터를 가져오는 방식

② 풀링의 메시지 배분



분산 처리

메시지 브로커는 높은 번으로 데이터를 쓰는 것이 최적화되어 있고,

여기에 노드에 부하 분산함으로써 성능을 끌어 올릴 수 있는 덕분에 확장성을 가지고 있다.

↳ 푸싱 형의 메시지 배분은 모든 메시지 브로커에 적용시키고, 일정한 번으로 커닝 데이터를 분산 스토리하여 기록하여 성능 문제를 피한다.

★ 풀링 형은 메시지 화일 사이즈 적당 하면 도움이 된다.

③ 메시지 낙하당 - 메시지 브로커에 저장된 데이터는 복수의 다른 노드에

있어 있을 수 있다. 이를 통해 메시지가 복수 노드에 여러개로 분산

신 리딩 문제 ... 3가지 신게 병행.

- 모바일 리딩과 같은 신뢰성이 낮은 네트워크에서는 반드시 중복, 누락이 발생함.

1) at most once → 메시지는 한번만 전송, 실패해서 4차원 가능성 있음

2) exactly once → 메시지는 수신되거나 중복 없이 1번만 전달됨.

3) at least once → 메시지는 확실하게 전달됨. 여러번 전달될 가능성이 있음.

중복 제거는 높은 비용의 문제임.

↳ TCP에서는 메시지에 시퀀스 번호를 붙이나, 분산 시스템에서는 시퀀스 번호
그냥... 그냥 X

모든 메시지에 일련번호 넣기는 힘들, 성능 향상 제한됨...

→ 이를 해결하기 위해!?

↳ 1) 중복을 허용한 중복 제거. (bulk 제거 가능)

- 전송 과정에서 파일명만 바뀌어, 동일함.

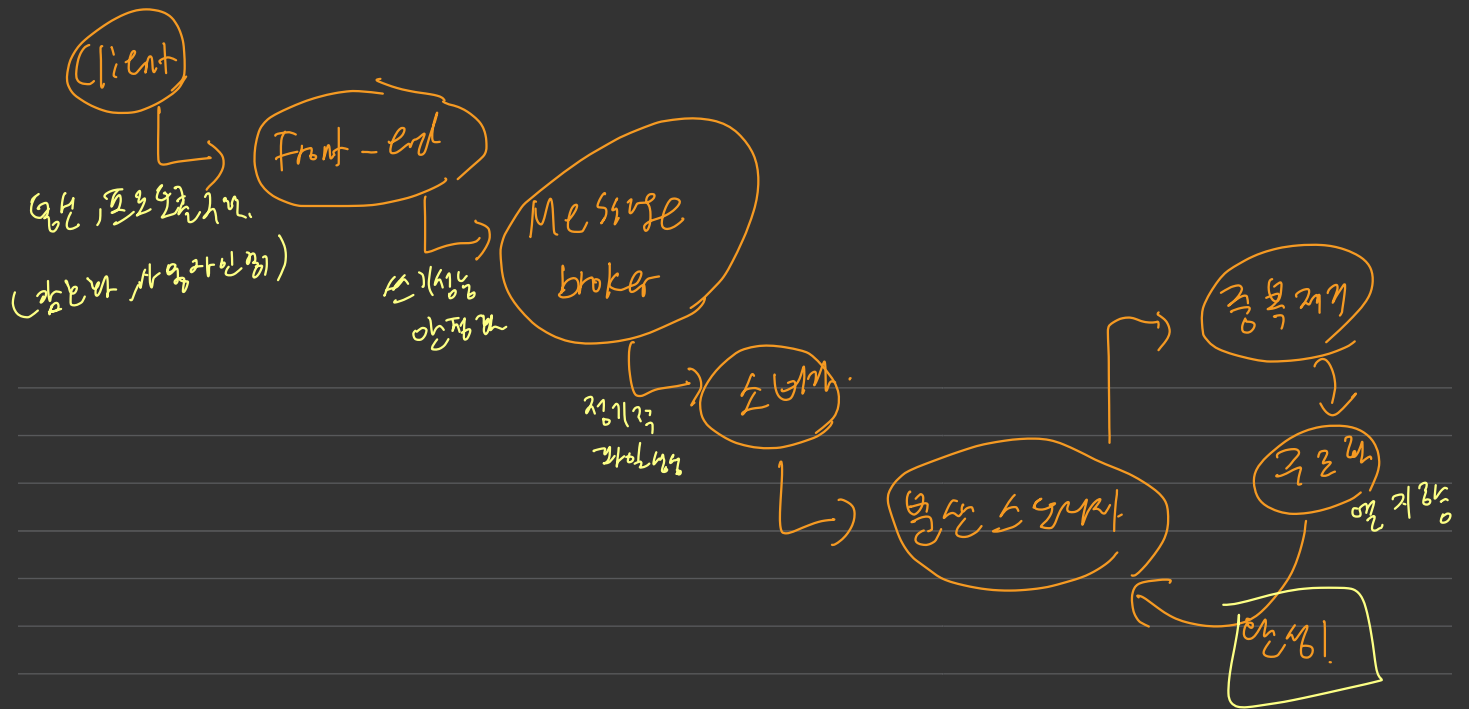
2) 고유 ID에 의한 중복 제거

- 리전에 따른 ID 부여 기구, 그 보다는 늦게는 메시지에 중복을 허용.

3) 종단적인 신뢰성 (End to End)

↳ 중복 전송으로 보지 않 *Idempotent* 하고, 프로그램 내에서
중복 제거

② 데이터 수집의 파이프라인 (재능기각된 데이터 분류에 적합한 storage)



① 시계열 데이터의 최적화

↳ 스트리밍 형태의 메시지 전송은 " 도착 시간직접"이 불가능.

이벤트시간 - 클라이언트 상에서 메시지가 생성된 시간

프로세스시간 - 서버가 처리하는 시간.

데이터 분석과 대량버피는 시간

② 프로세스 시간에 의한 불합과 문제점.

블래스드 네거미 데이터는 낮은 단계에서는 이벤트 시간이 아니라, "프로세스 시간"은 사용량이
보통이다.

EX) 2017년 월 1만 미 동작한 데이터만, 다수의 데이터가 포함되어 있다. ^{이벤트 기록}

이 경우 매우 많은 파일을 열어야
원하는 이벤트를 찾을 수 있다

다수의 파일을 모두 검색하는 것은 풀스캔 (full scan), 시스템의 부하를 크게
높이는 요인이 된다.

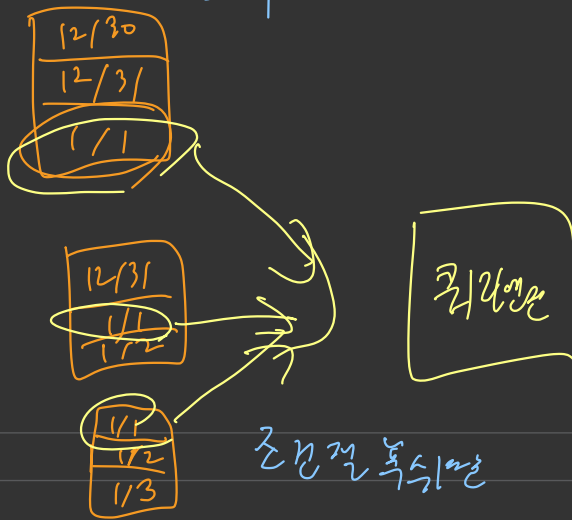
이벤트 인덱스 - 이벤트 시간에 의한 정제의 효율화

LRDB 메시 index를 만드는 것과 마찬가지로, 이벤트 시간에 인덱스를 만든다.

LRDB 메시 index에 대응하는 카운트나 분산 DB 등 이용되면, 시간으로 인덱스된
 table 생성 가능
 ↳ 같은 메시는 중복 지양

이벤트 시간 축적 - 이벤트 시간에 의한 정제의 효율화

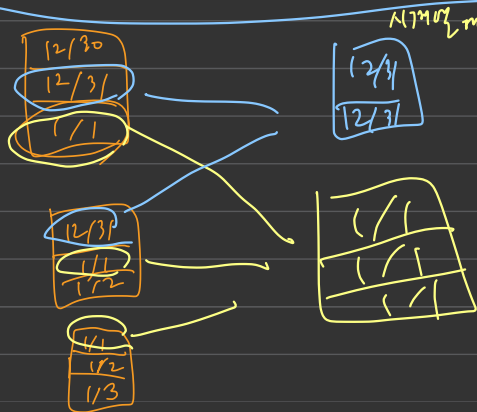
↳ 이벤트 발생 시간으로 정제할 수 있게
 하는 처리를 하고 있다



[데이터가 중복] 연속적으로 배치되어 있는 데이터를 리드하도록 하고 ↑

↳ 가능한 지기로 만들면 X (데이터가 적으면 만들 load가 많음)

이벤트 시간에 의한 정제 - 데이터의 효율성, 이벤트 시간 데이터



이벤트 시간 메시
 data는 효율성
 가 높음...

데이터마스트는 이벤트를 바탕으로 접근하기

↳ 데이터마스트 만드는 것 자체에 대한 활용도도 높음.

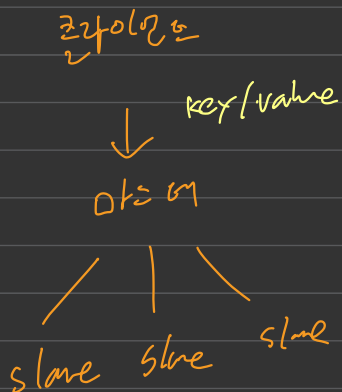
⑤ 비구조화 데이터의 분산스토리지

↳ NoSQL 사용시, 데이터는 단순히 보편적 저장뿐 아니라, application에
맞춘 이용, 실시간 검색 가능.
작성이 쉬움. (write)

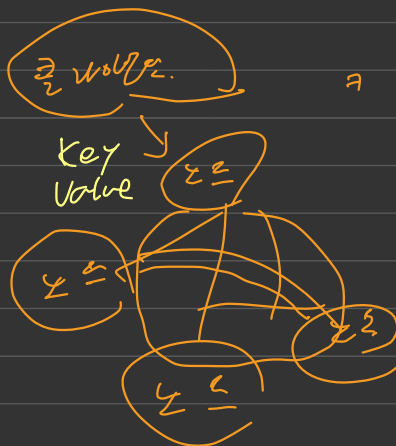
• 분산 KVS - 디스크로의 쓰기성능을 높이기.

↳ distributed key value store
모든 데이터를 키-value 쌍으로 저장하므로 NoSQL DB.

① Master/slave



② P2P형



• Amazon Dynamo DB

- ↳ 항상 안정된 읽기 쓰기 성능을 제공하도록 디자인된 분산형 NoSQL DB.
- ↳ P2P형.

• 데이터 관리 소프트웨어 - 구조화 데이터를 분석하기 위한 것.

- ↳ 분산 KVS를 발전시켜 2개 이상의 서버에 커미티 데이터는 저장하는 수 있도록 함
- ↳ Google Bigtable
- ↳ Apache HBase
- ↳ Apache Cassandra.
- ↳ Row에 column명과 Value를 같이 저장함.

• Apache Cassandra.

- ↳ 데이터 관리 소프트웨어, (NoSQL)
- ↳ 스키마를 결정할 필요가 없다,
- ↳ insert into (주요)으로 동작.
- ↳ 동적인 레코드를 읽어옴
- ↳ P2P형.

• 돈작민트스키 - 스키마리스 데이터 관리하기.

- ↳ 주로 데이터 처리 무관성 보장.

JSON 형식 복잡하게 구성된 스키마리스 데이터는 그 데이터의 형식에
저장된 쿼리를 실행할 수 있도록 함.

