



셀프서비스 Stream Designer로 나의 시스템을 더욱 강력하게

(ksqlDB기반 사용자 Self Stream Processing 플랫폼 서비스)

SK Hynix **박재승**

박재승

SK Hynix



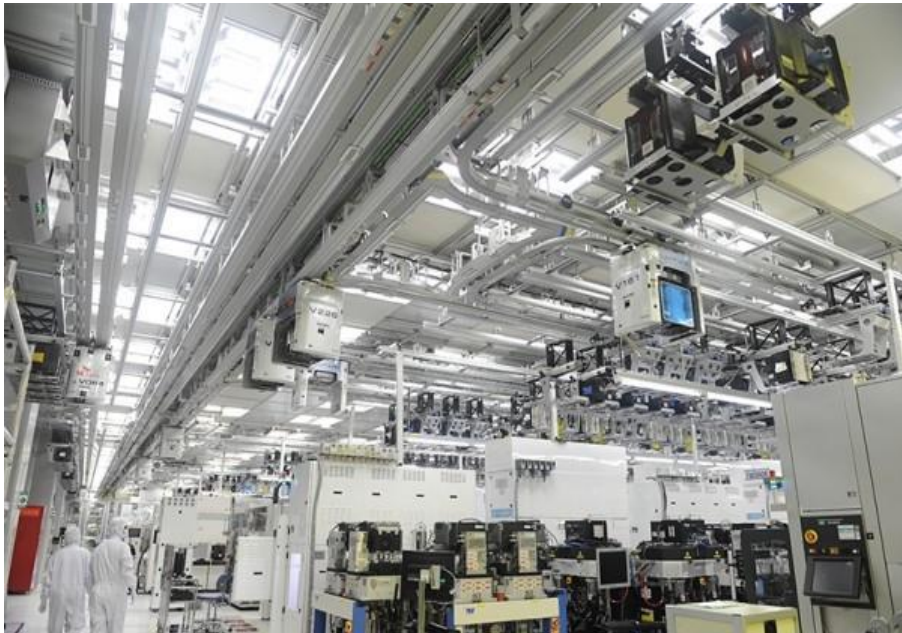
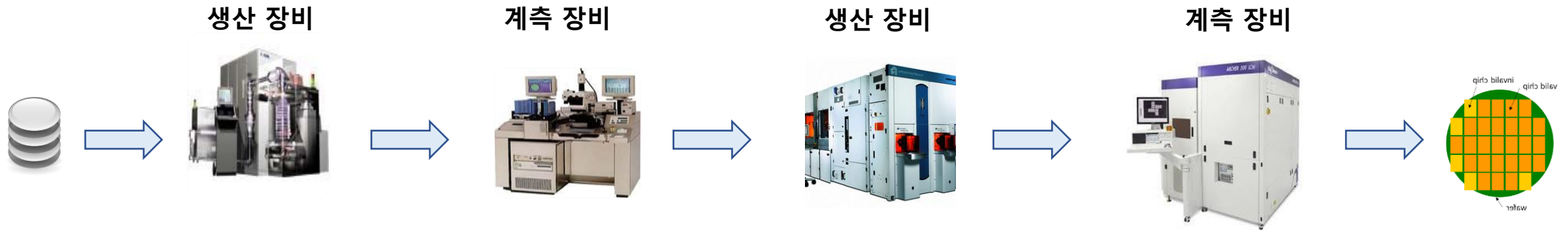
실시간 Data Platform / Data Catalog System 개발 리딩

기술 영역 | Data Engineer

Contents

- 01** 반도체 실시간 데이터
- 02** What's Realtime Platform
- 03** Data Hub
- 04** h-STREAM

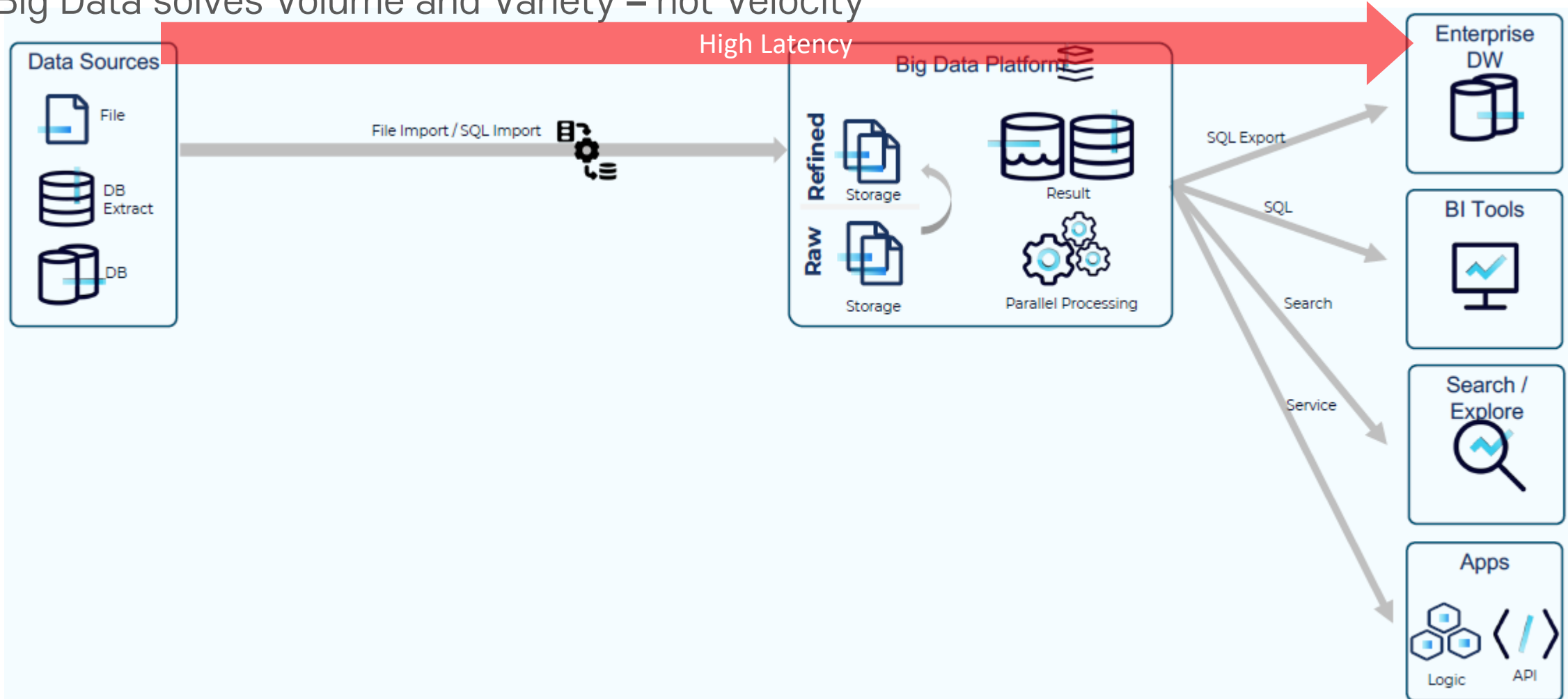
반도체 실시간 데이터



27,000 EPS
54 MB / sec

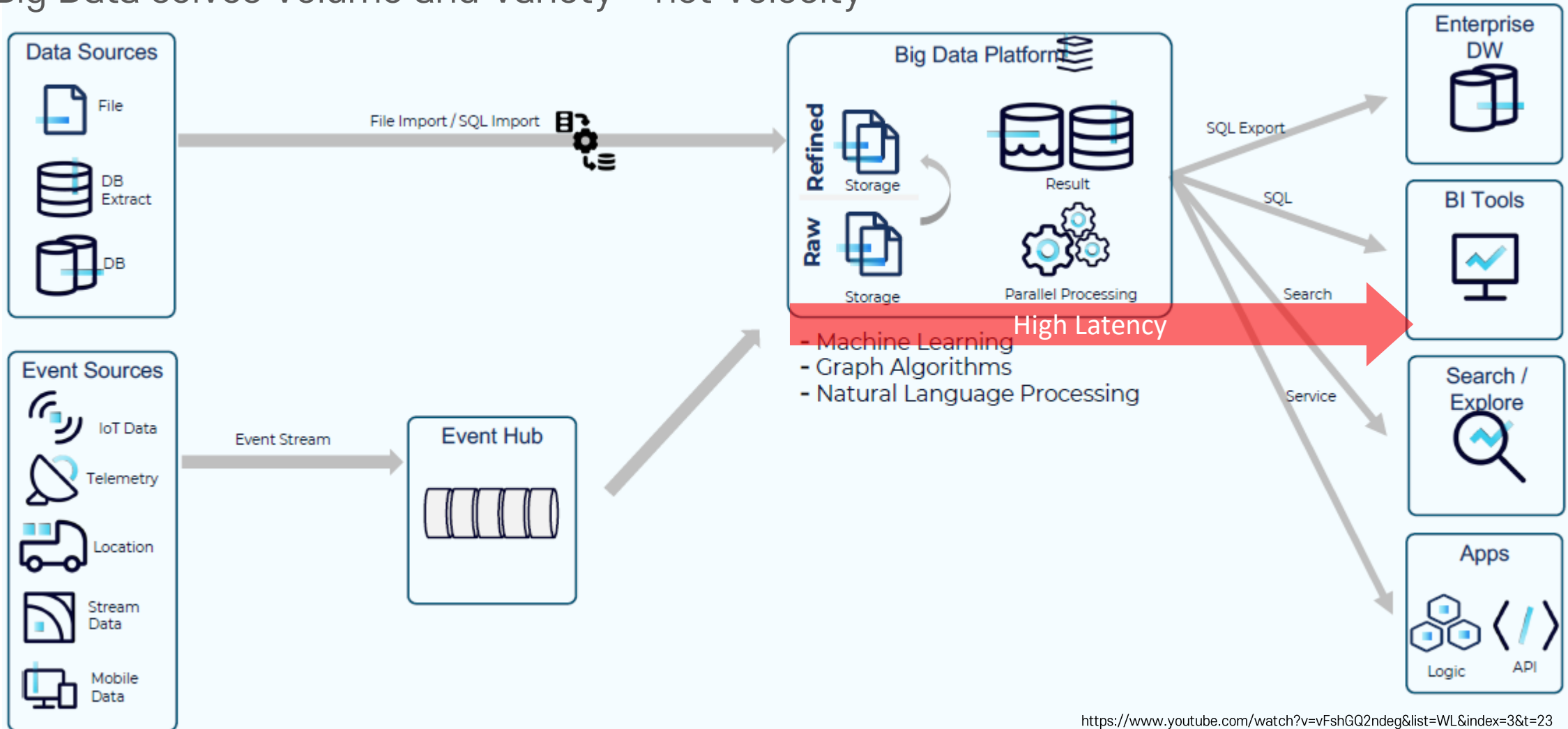
What's Realtime Platform?

Big Data solves Volume and Variety – not Velocity



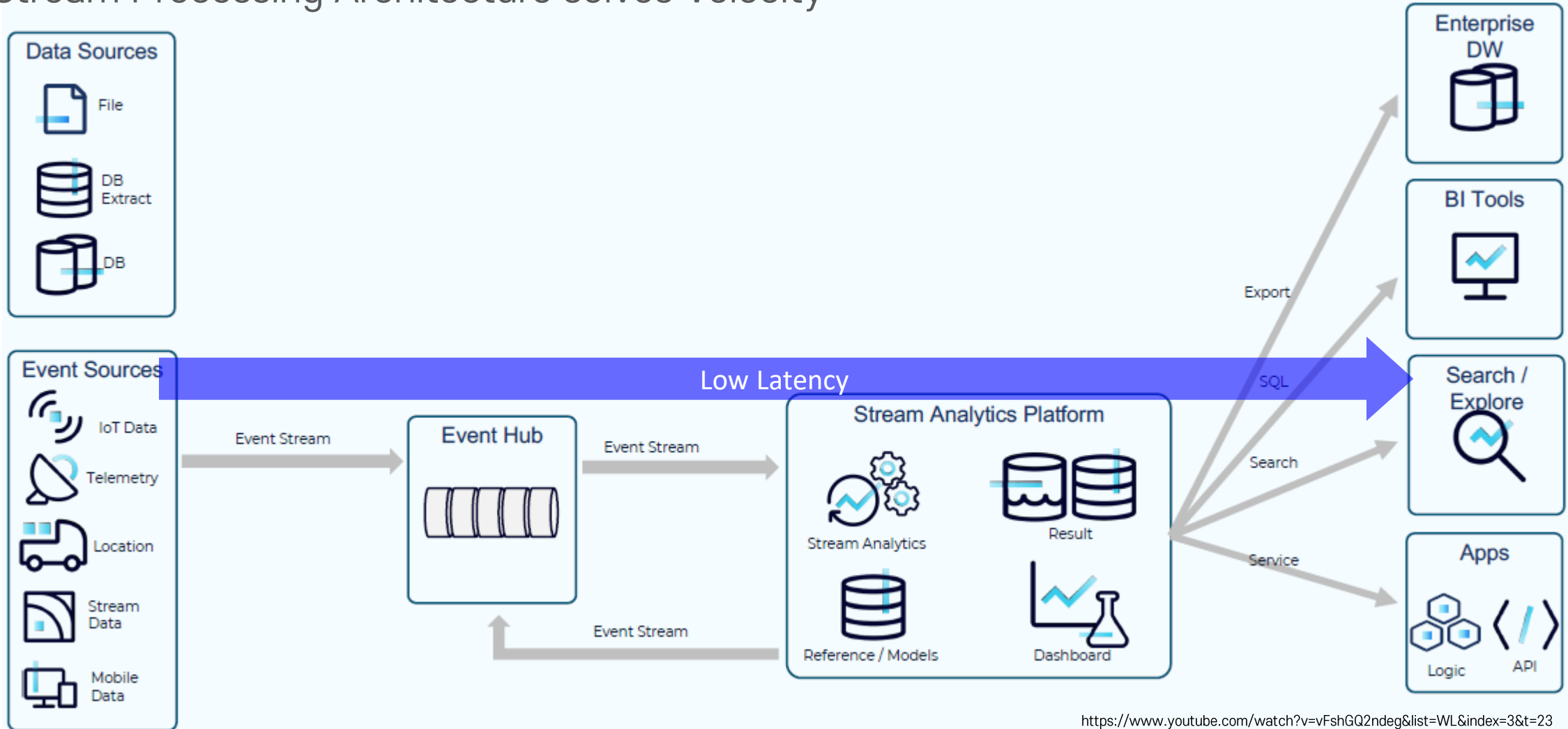
What's Realtime Platform?

Big Data solves Volume and Variety – not Velocity



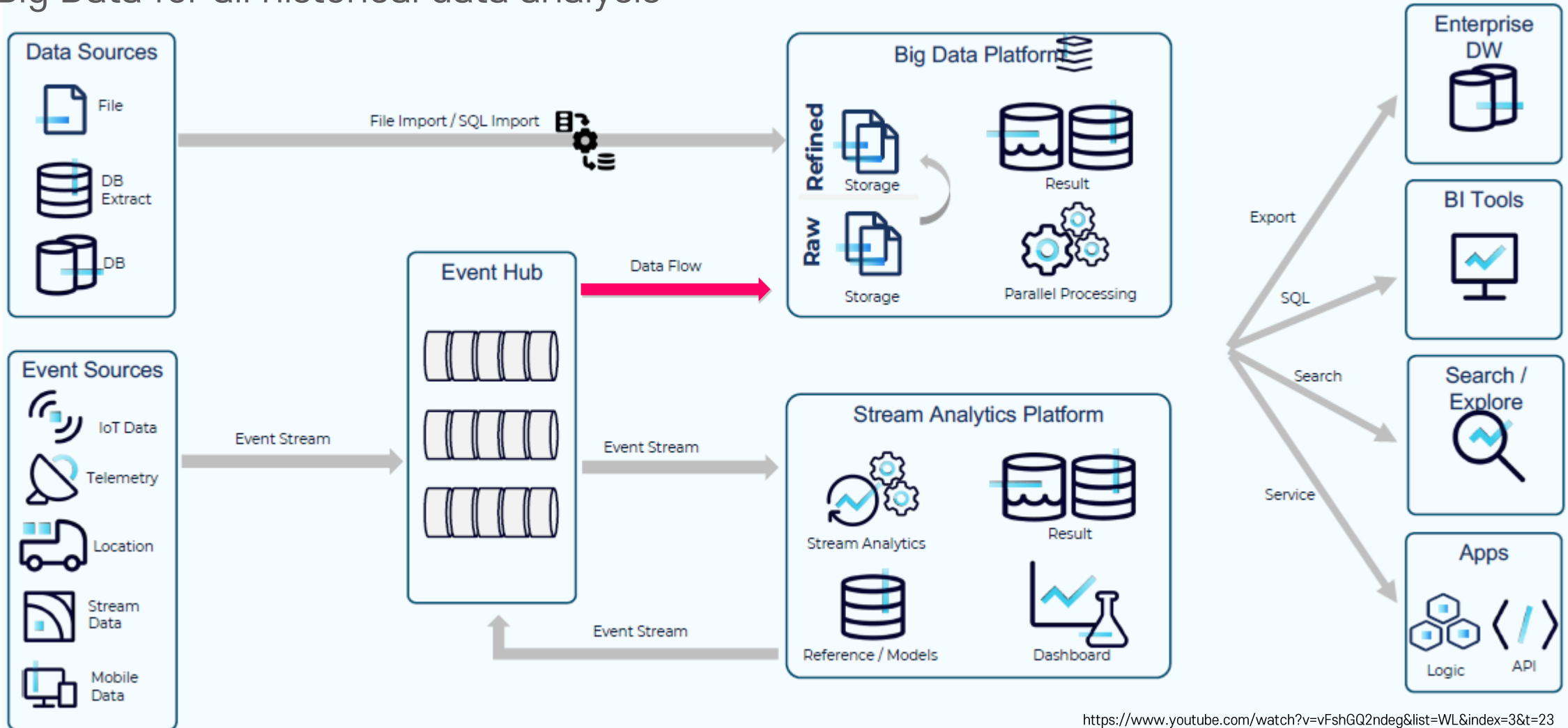
What's Realtime Platform?

Stream Processing Architecture solves Velocity



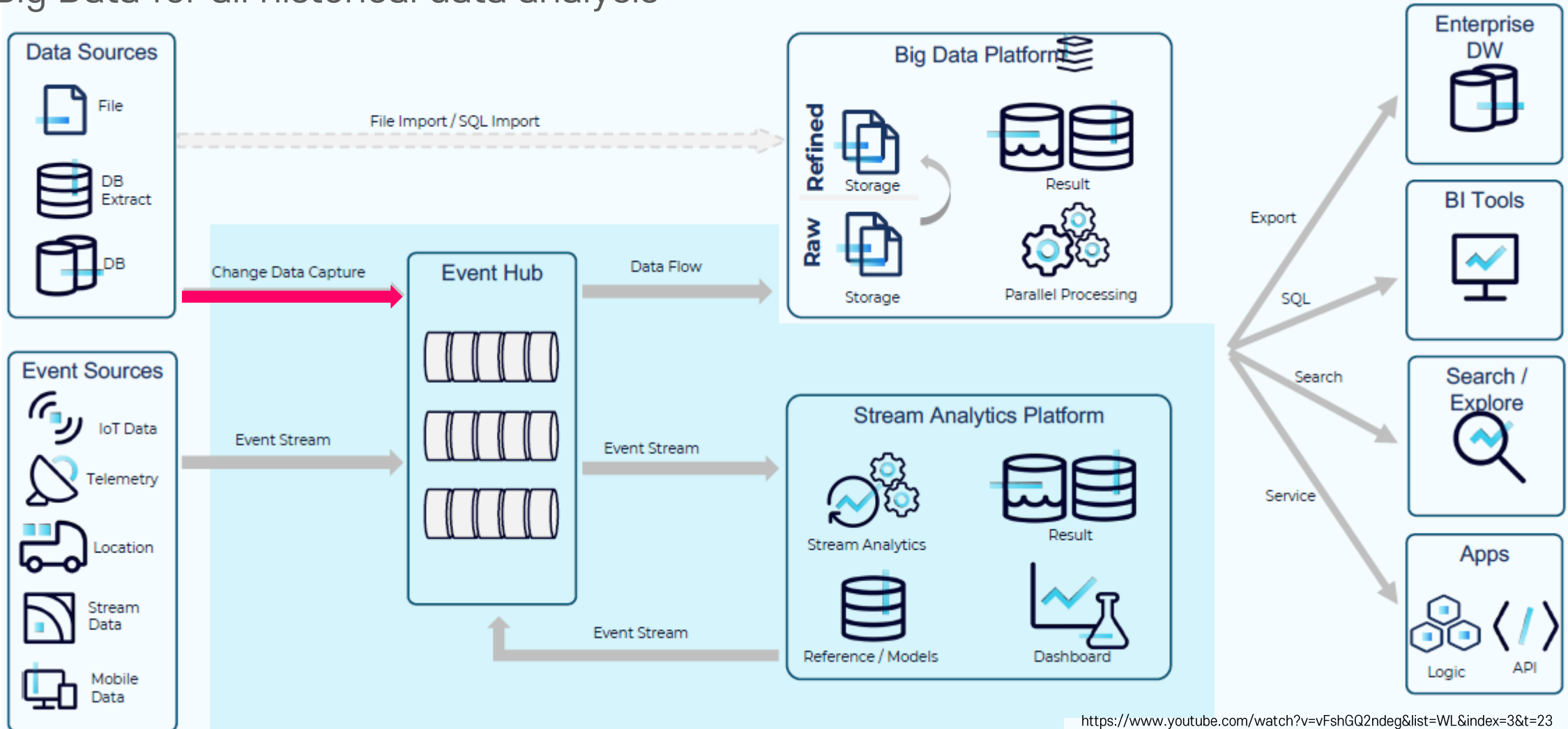
What's Realtime Platform?

Big Data for all historical data analysis



What's Realtime Platform?

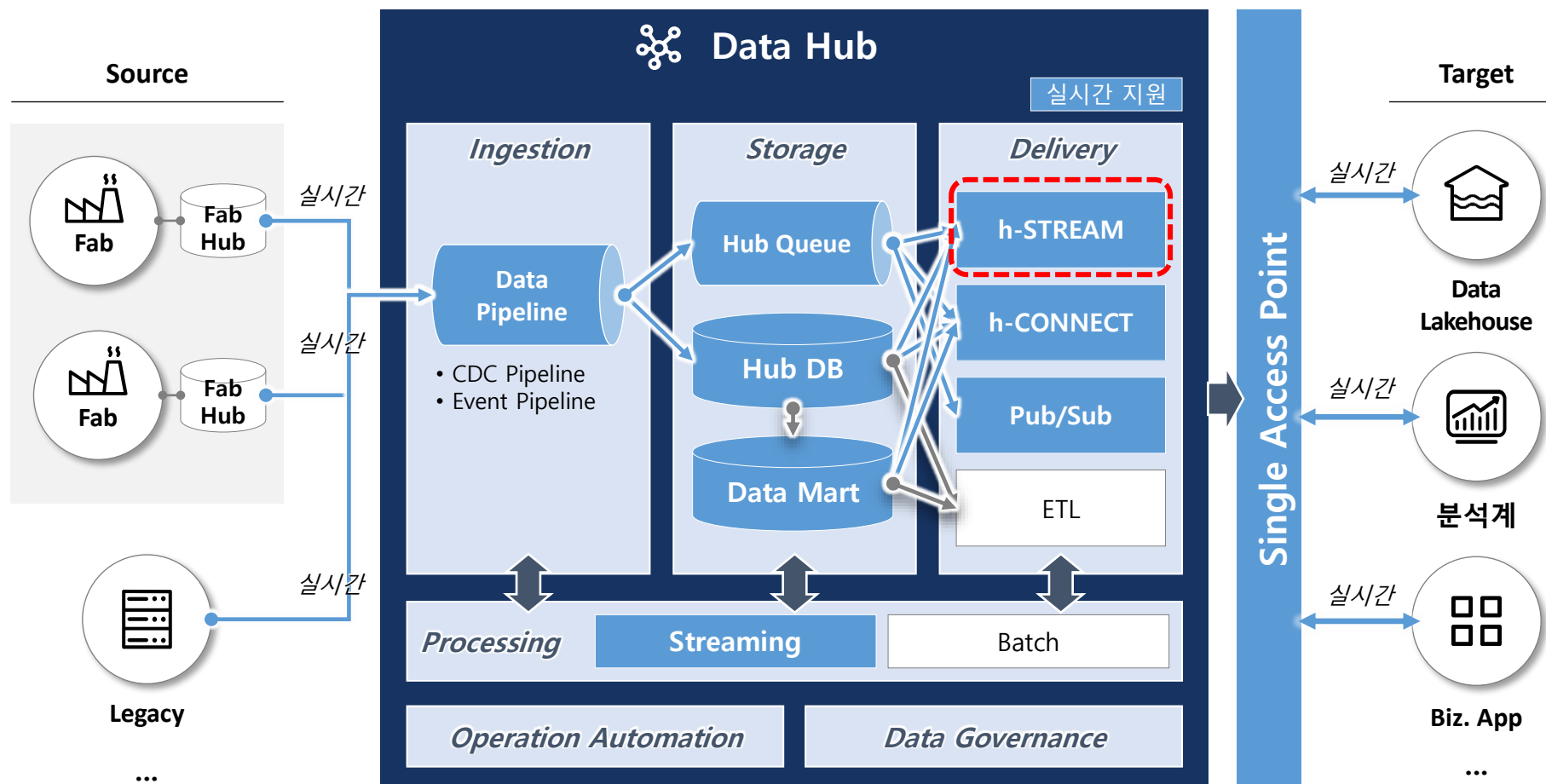
Big Data for all historical data analysis



Data Hub는 전사 데이터 연계를 위한 Hub 역할 수행
Source와 Target 간 실시간 데이터 Pipeline으로서 데이터 활용의 신속/적시성을 지원

Data Hub

Data Hub는 전사 데이터 연계를 위한 Hub 역할 수행
Source와 Target 간 실시간 데이터 Pipeline으로서 데이터 활용의 신속/적시성을 지원





손쉽게 실시간 데이터를 **확인** 하고 싶어요

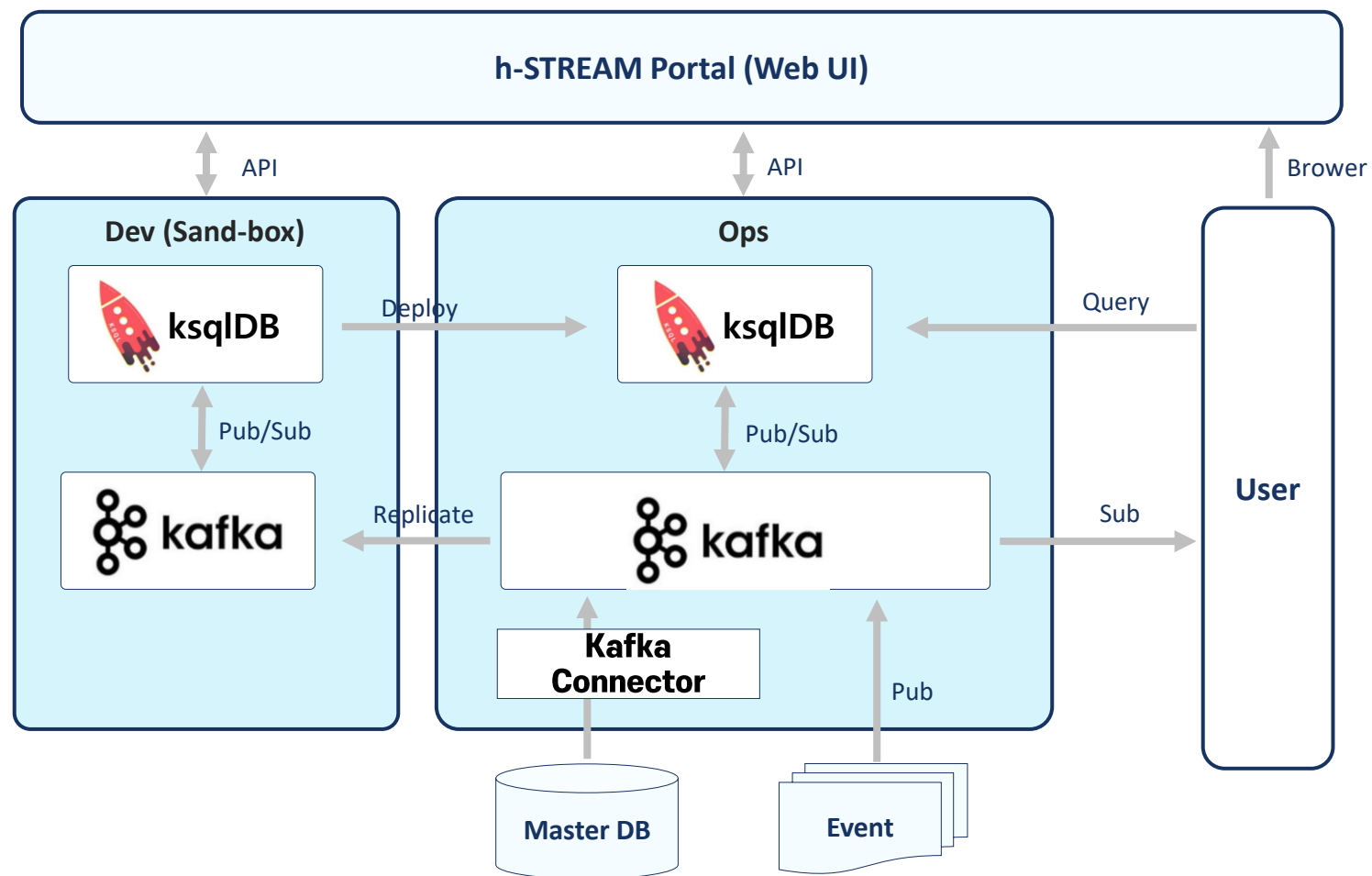
실시간 데이터를 **조작** 해서 **나만의 실시간 데이터**를 만들고 싶어요

컴퓨팅 자원이 따로 없어 Stream Processing을 못하고 있어요 $\pi\pi$

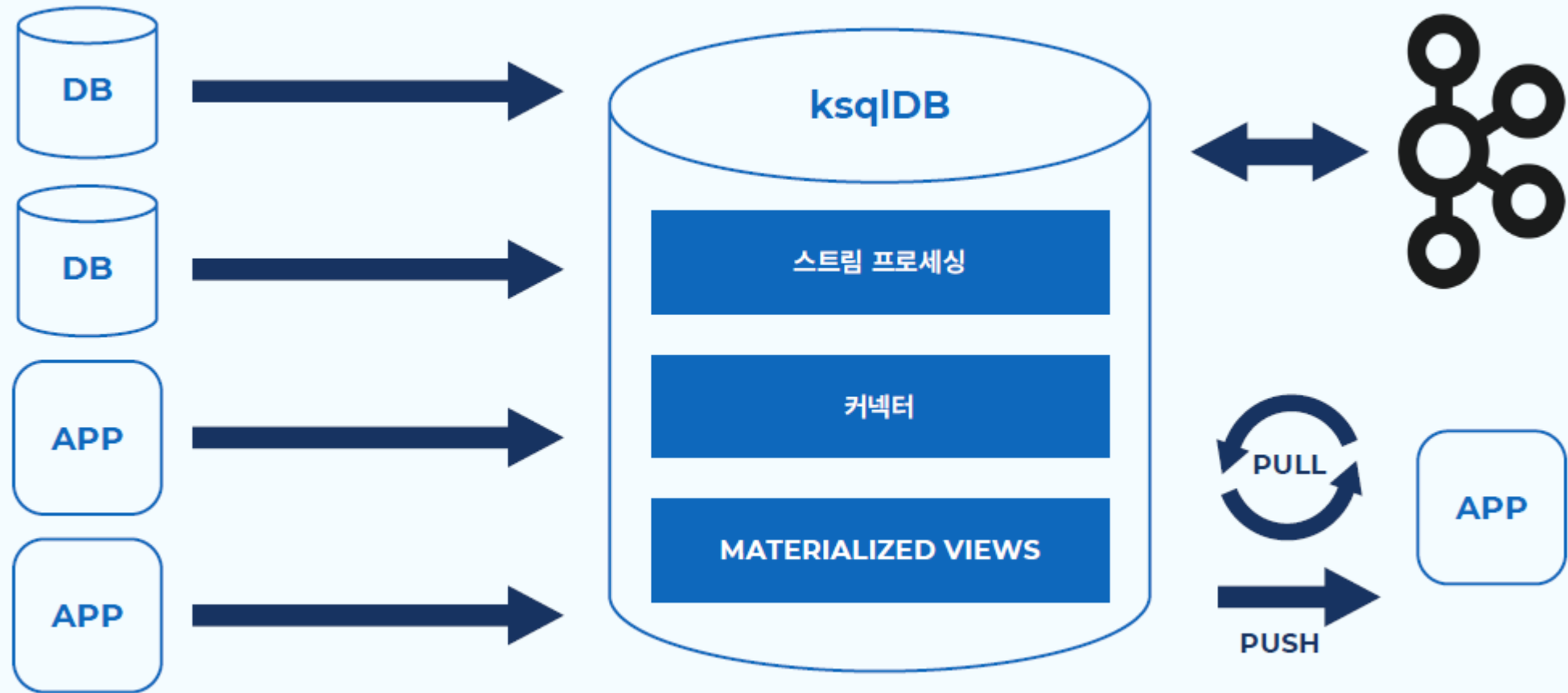
안정적인 Stream Processing 로직을 운영 시스템에 연계하고 싶어요

Stream Processing **노하우**를 **공유** 받고 싶어요

Architecture



ksqlDB



Stream Processing Tools

Kafka Producer & Consumer

```
ConsumerRecords<String, String> records =
consumer.poll(100);
Map<String, Integer> counts = new DefaultMap<String,
Integer>();
for (ConsumerRecord<String, Integer> record : records) {
    String key = record.key();
    int c = counts.get(key)
    c += record.value()
    counts.put(key, c)
}
for (Map.Entry<String, Integer> entry : counts.entrySet()) {
    int stateCount;
    int attempts;
    while (attempts++ < MAX_RETRIES) {
        try {
            stateCount = stateStore.getValue(entry.getKey())
            stateStore.setValue(entry.getKey(), entry.getValue() +
stateCount)
            break;
        } catch (StateStoreException e) {
            RetryUtils.backoff(attempts);
        }
    }
}
```

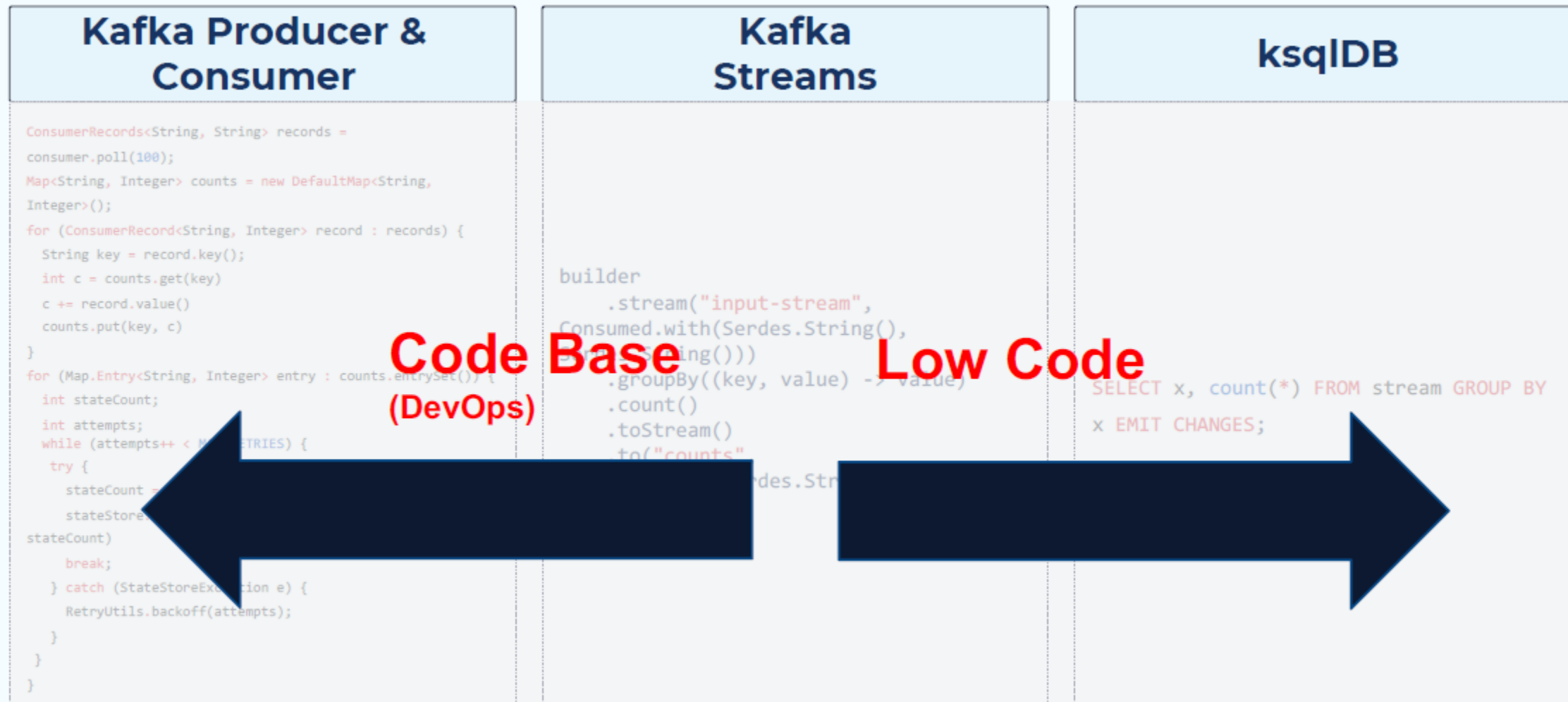
Kafka Streams

```
builder
    .stream("input-stream",
Consumed.with(Serdes.String(),
Serdes.String()))
    .groupBy((key, value) -> value)
    .count()
    .toStream()
    .to("counts",
Produced.with(Serdes.String(),
Serdes.Long()));
```

ksqlDB

```
SELECT x, count(*) FROM stream GROUP BY
x EMIT CHANGES;
```

Stream Processing Tools



h-STREAM

Topic / Stream / Table

Topic



Stream



=



Topic

+



Structure

Table



=



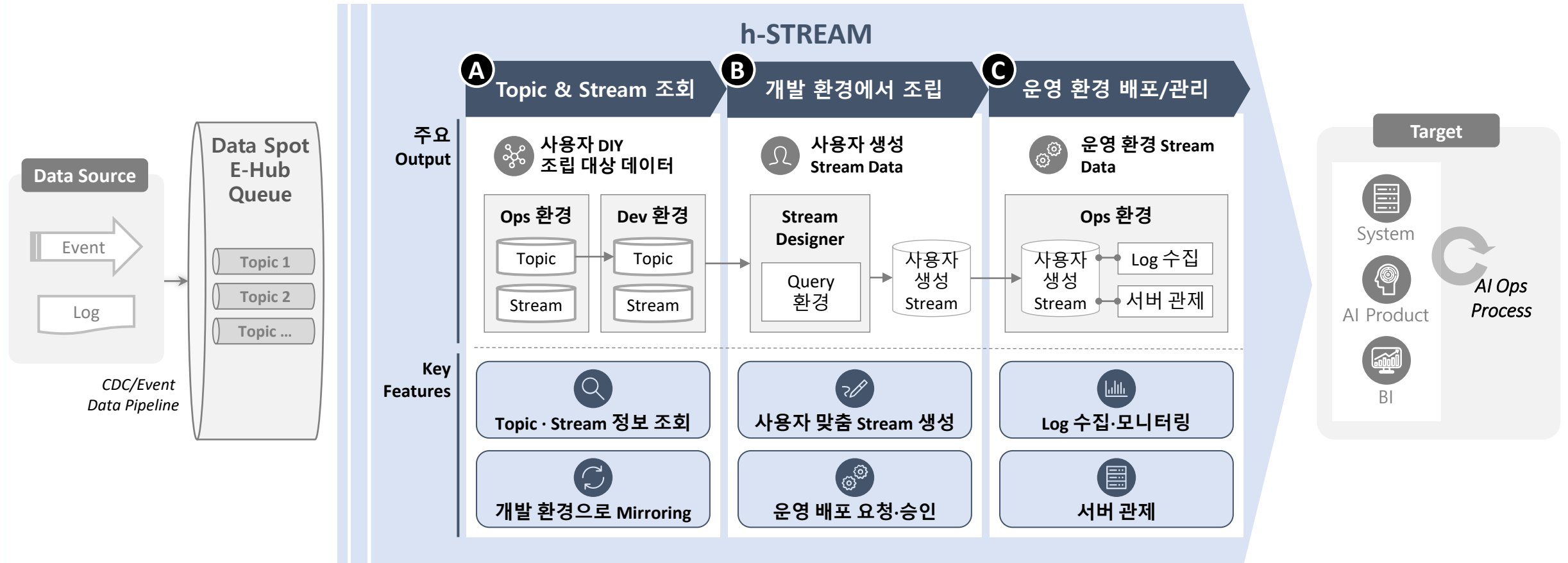
Stream

+



RocksDB

h-STREAM Features



Topic & Stream Search

Topic · Stream 정보 조회

Topic · Stream 목록 조회

- 개발(Sandbox) 및 운영 환경에서 생성/운영 중인 Topic 및 Stream 목록 조회

Topic

Total: 568 | 운영 개발 *개발/운영 중 선택*

<input type="checkbox"/>	토픽명
<input type="checkbox"/>	connect-configs
<input type="checkbox"/>	connect-offsets
<input type="checkbox"/>	ftl.adapter.cj.m11.edp.edes.hub.cmdcol
<input type="checkbox"/>	ftl.adapter.cj.m11.edp.edes.hub.alarhistory

Topic · Stream 상세 정보 조회

- Topic/Stream 명, 배포 상태, 생성자, 생성일, 최종 수정일 등 해당 Topic/Stream에 대한 상세 정보 조회

Stream 상세

Stream

스트림명	test1name3		
KSQL Query	test_query		
생성자	시스템(system)	생성일	2022-09-08 10:07:43

Topic

토픽명	cdc.dcp.m14.dcp_dcoldervitem_det	파라미터 갯수
복제수	2	Sync 중인 복제수
메세지 갯수	0	메세지 사이즈

개발 환경으로 Mirroring

Topic 복사 (운영 환경 → Sandbox)

요청

제목	요청 구분	상태	생성일
Topic 1	토픽 복사	접수	2022-09-02 11:24:07
Topic 2	토픽 복사	요청	2022-09-02 11:24:07

상세
조회

요청 상세

요청 구분	토픽 복사	상태	요청
제목	subject1	요청자	시스템(system)
요청 내용	content1	생성일자	2022-09-02 11:24:07

- 운영 Topic 조회 후, 사용자가 직접 DIY 조립을 수행할 항목들을 개발(Sandbox) 환경으로 복사 요청
- 관리자는 해당 요청 사항을 확인한 후, 이를 승인하여 개발 환경으로 복제

Stream Development

사용자 맞춤 Stream 생성

KSQL 기반 Query 작성 (Stream Designer)

Stream Designer

Stream Designer Query Editor

```
1 CREATE STREAM TEST_MINWOCK_20220824_427 (FAB STRING, ENV STRING, MODE STRING, RCV STRING, MSG_TYPE STRING, MSG_TYPE2 STRING, MSG_SRC STRING)
```

Stream Data 생성을 위한 Query 작성

References

- Table 1
- Table 2
- Stream 1

Result

```
{
  "commandId": "stream/10101834_TEST_MINWOCK_20220824_427/create",
  "commandStatus": {
    "status": "SUCCESS",
    "message": "Stream created",
    "queryId": null
  },
  "commandSequenceNumber": 222,
  "warnings": []
}
```

결과 조회

- Sandbox 환경에서 KSQL 기반 Query 작성하여 신규 Stream Data 생성
- Query 작성 참고를 위한 Reference Table 및 타 Stream Data 조회 가능

운영 배포 요청

Stream 배포 요청 (Sandbox → 운영 환경)

요청

제목	요청 구분	상태	생성일
Stream 1	스트림 배포	접수	2022-09-02 11:24:07
Stream 2	스트림 배포	요청	2022-09-02 11:24:07

요청 상세

요청 구분	스트림 배포	상태	요청
제목	subject1	요청자	시스템(system)
요청 내용	content1	생성일자	2022-09-02 11:24:07

- 개발 환경에서 최종 생성된 Stream Data를 운영 환경에 배포하기 위해서, 관리자에게 '스트림 배포' 요청
- 관리자는 요청 건 확인 후 승인하여 운영 환경 배포 수행

Stream Deploy

Log 수집·모니터링

Log 확인

Audit Logs

Total: 113

IP Address	요청 API	성공여부	발생 시간	Created At	Creator
123.123.123.123	/test	성공	2022-07-06 17:00:13		홍길동(h1231234)
123.123.123.123	/test	실패	2022-07-06 17:00:13		홍길동(h1231234)
123.123.123.123	/test	성공	2022-07-06 17:00:13		홍길동(h1231234)
123.123.123.123	/test	실패	2022-07-06 17:00:13		홍길동(h1231234)
123.123.123.123	/test	성공	2022-07-06 17:00:13		홍길동(h1231234)
123.123.123.123	/test	성공	2022-07-06 17:00:13		홍길동(h1231234)

< 1 >

- h-STREAM에서 발생한 Log 수집
 - IP Address, 요청 API, 성공 여부, 발생 시간, 생성 시간, 생성자
- Search 기능 활용하여 키워드 기반 Log 검색/조회 가능

서버 관제

서버 목록 조회

- h-STREAM 서버 및 상태 정보 조회

호스트명	IP Address	상태
hostName1	123.123.123.121	성공
hostName2	123.123.123.122	성공
hostName3	123.123.123.123	성공
hostName4	123.123.123.124	실패

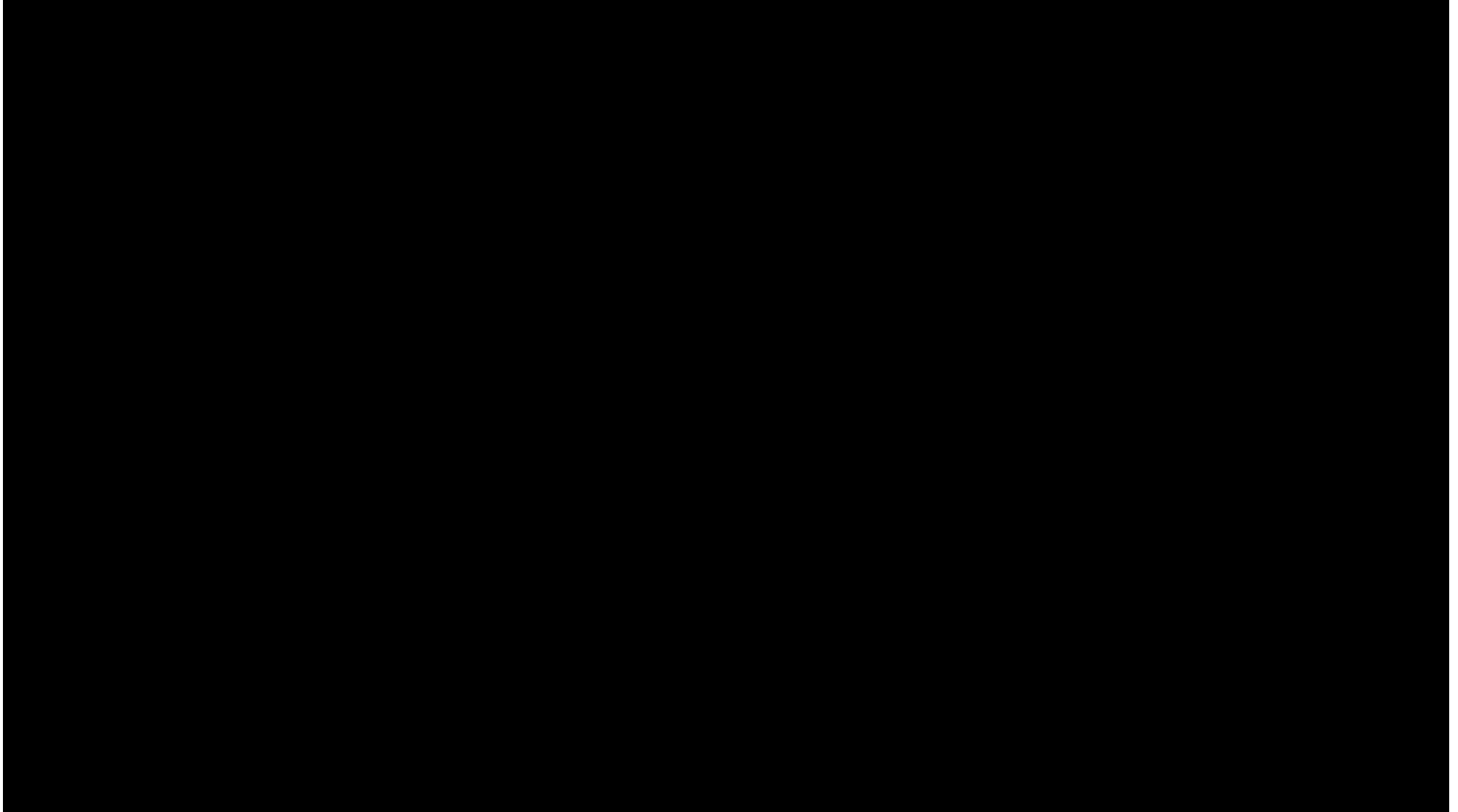
운영 관리

- 개별 서버 목록 및 상태 확인 후 각 서버에 '중지', '재가동', 'LOG보기' 명령 수행

관리 수행

관리
중지 재가동 LOG보기
중지 재가동 LOG보기
중지 재가동 LOG보기

Demo



Expected Effect



실시간 데이터의
Self-Service 분석 구현

**Self-Service
선별 및 조립**

ksql 쿼리 기반 손쉬운
topic 조합 및 재구성 지원

**Code 작성
최소화**

별도의 API 개발 필요 없이,
Web UI 기반으로 topic 접근



실시간 데이터 기반
비즈니스 Use Case 확산

**실시간 데이터
가치 입증**

실시간 데이터 적용 사례
발굴 기반 Biz. Impact 실현

**적시적
의사결정 지원**

핵심 의사결정을 위한 분석
모델에 적시적 Input 제공



올바른 의사결정을 위한
데이터 투명성 증대

**Data
Governance**

사용자 배포 Stream이
전사 Catalog와 연동

**Data
Timeliness**

의사결정 시점에 상응하며
시차로 인한 오류 제거



전사 데이터 자산
활용 활성화

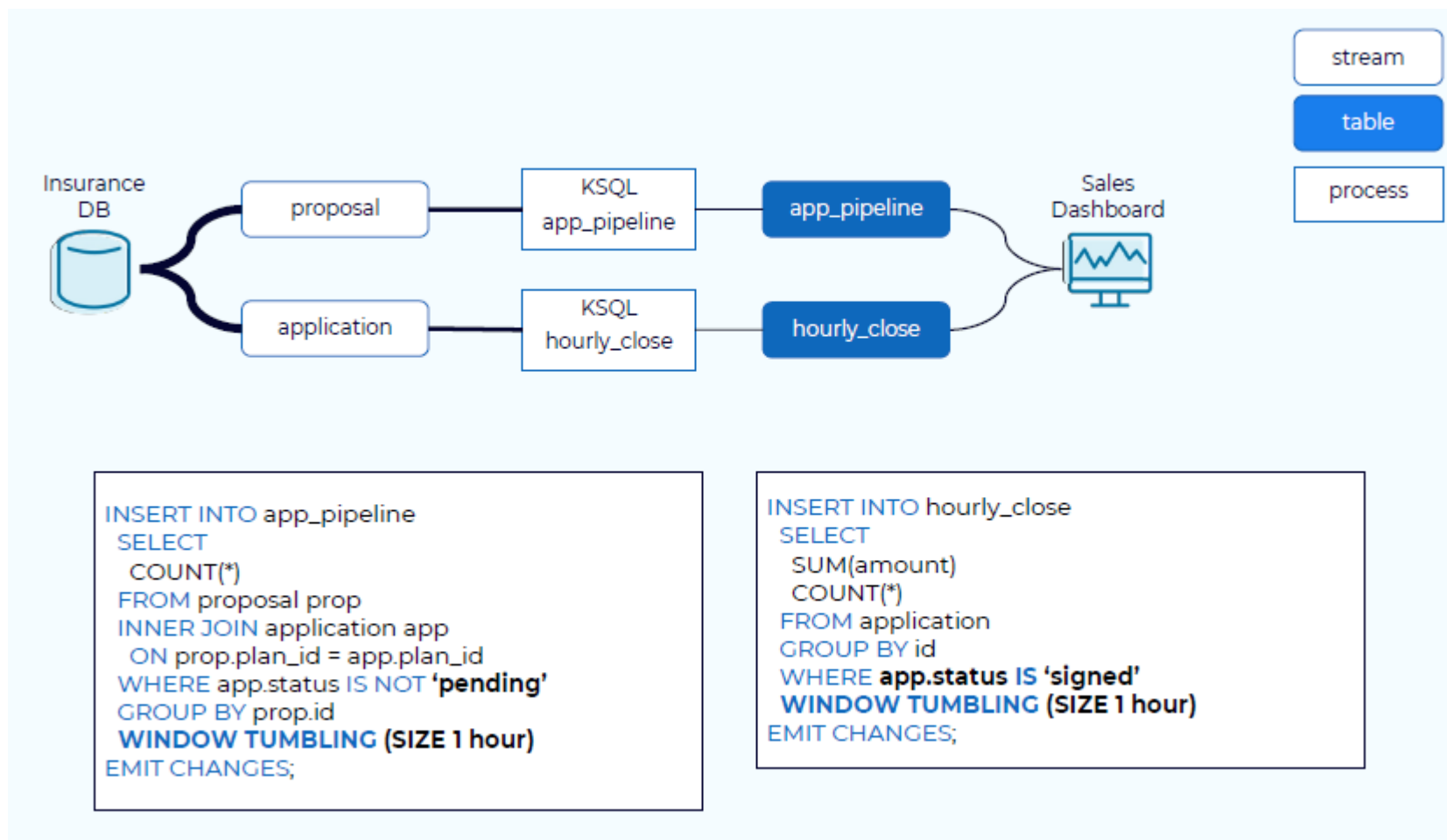
**CDC 데이터
활용 확산**

활용률이 상대적으로 낮은
CDC 데이터의 활용도 강화

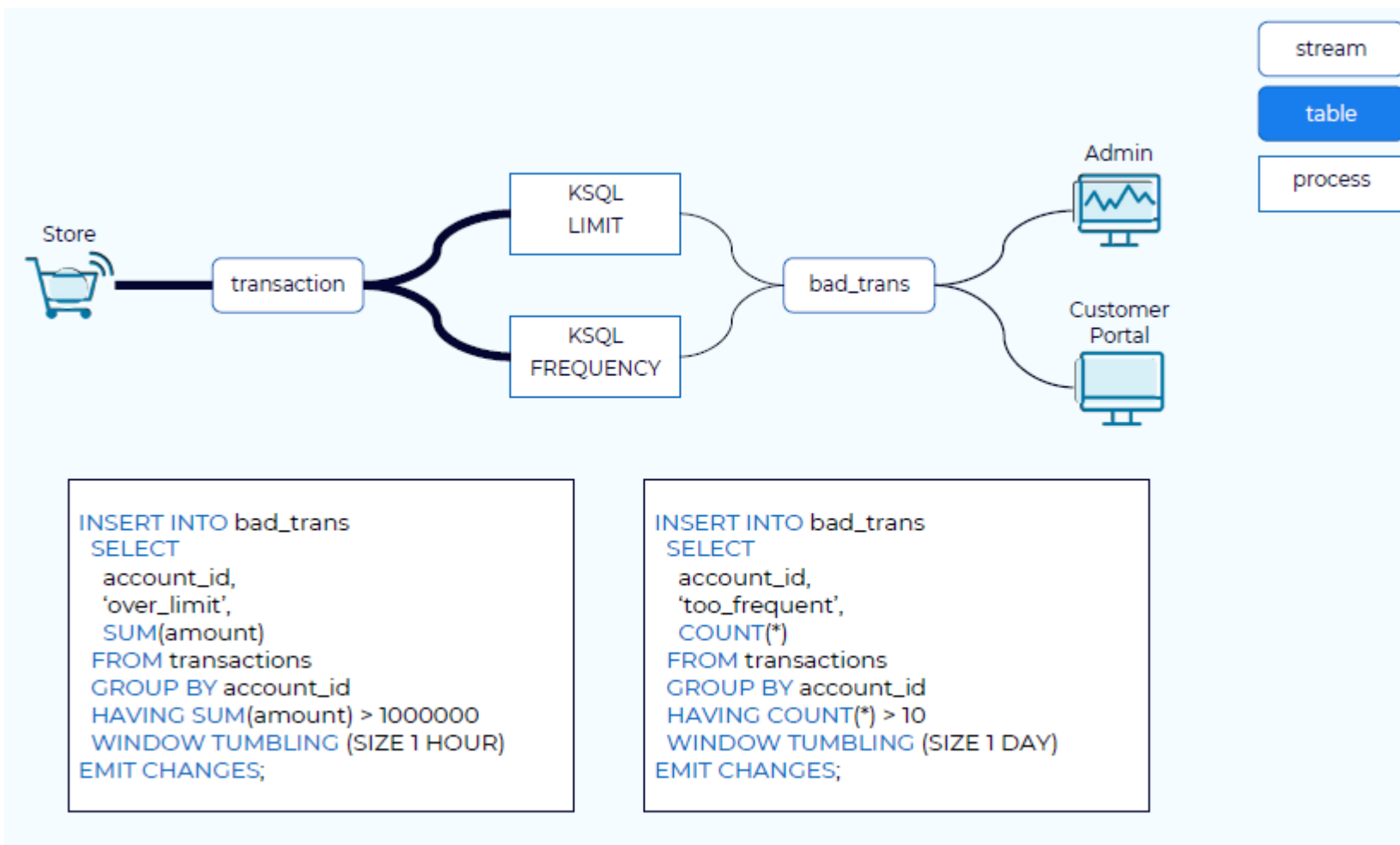
**분석 문화
고도화에 기여**

In-stream Analytics 문화
확산을 위한 기반 제공

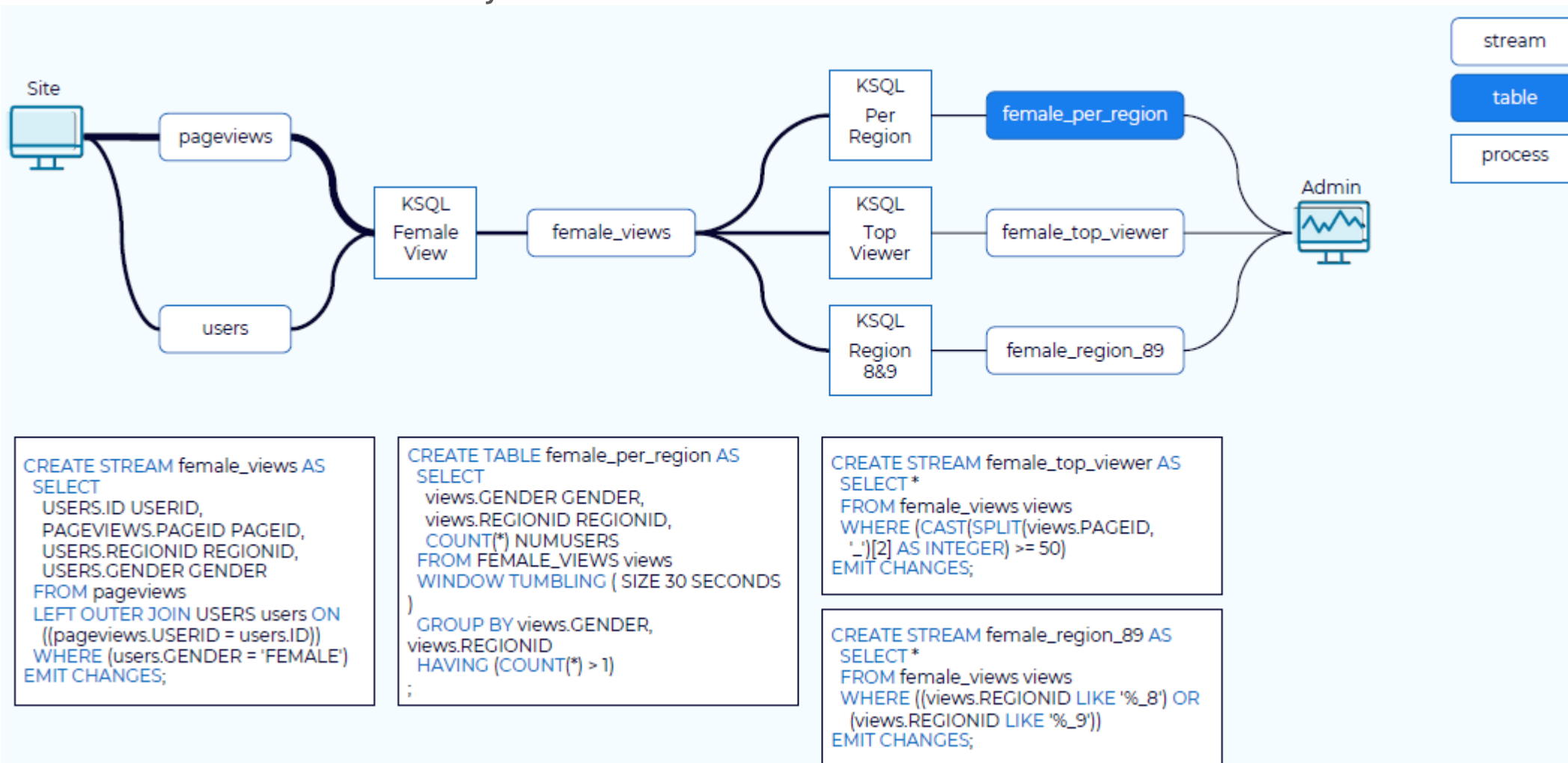
Use Case > Updating Dashboards



Use Case > Fraud Detection



Use Case > Site Access Analysis



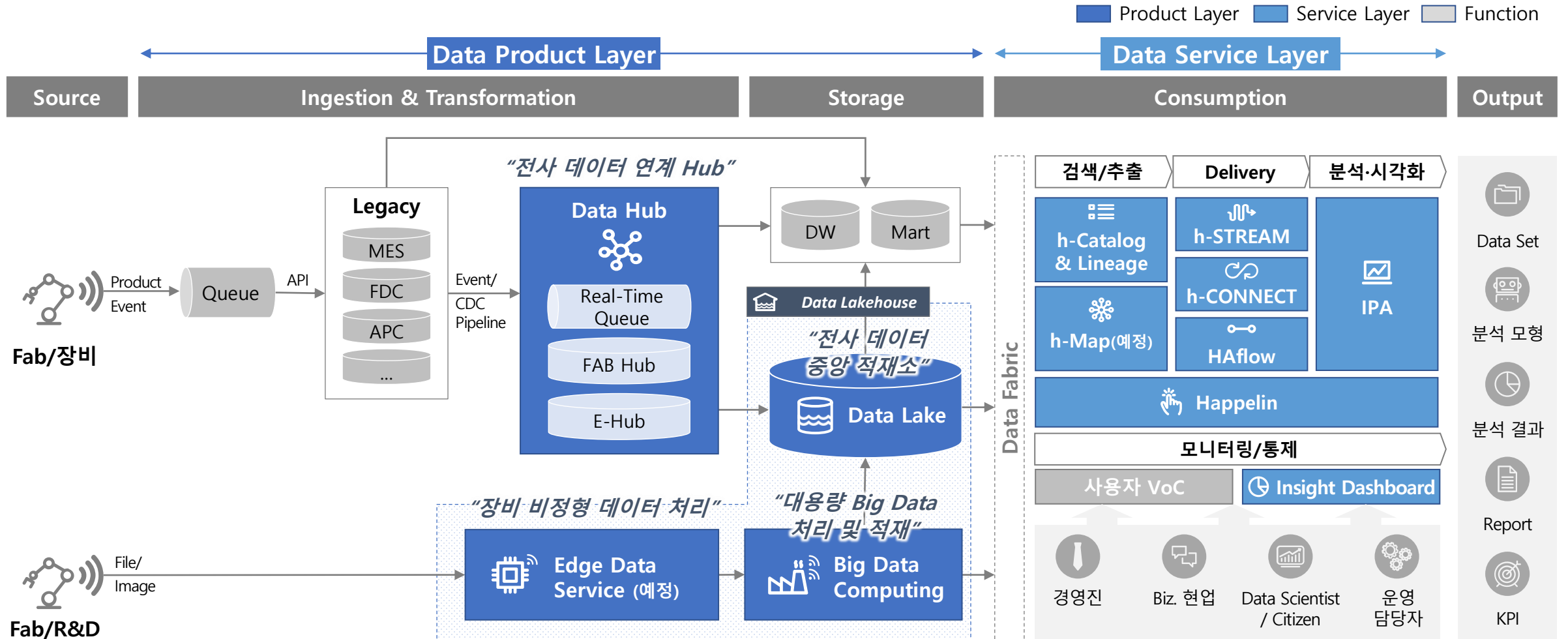
Q&A

E-mail jseung.park@sk.com

Thank you

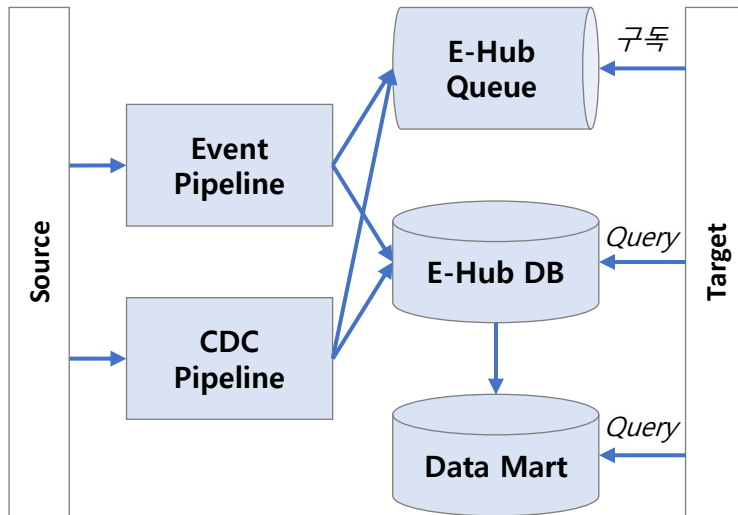
Data Platform

Data Platform Eco System



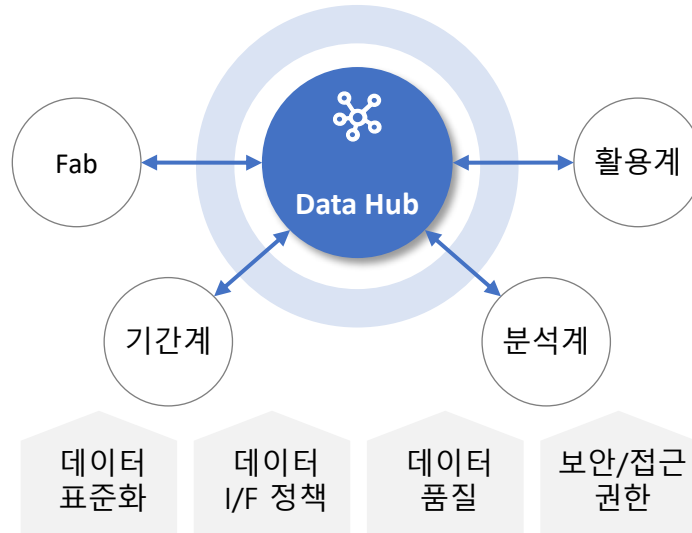
Real-time Data Platform

① 실시간 데이터 연계



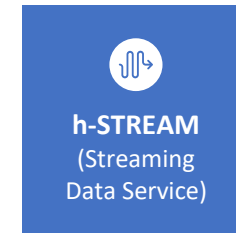
- End-to-End Near Real-Time 데이터 연계 (데이터 발생 시점부터 10초 이내)
- E-Hub Queue 적재 데이터 **실시간 구독**
- 실시간성 **Mart向 Ad-hoc Query** 실행 가능

② 전사 Interface Hub 역할 수행



- 전사 데이터 연계 **Single Access Point**
- I/F 일원화를 통한 **기간계 운영 부담 경감**
- 통일된 I/F 운영 기준 기반 **데이터 품질 및 일관성 확보**

③ 다양한 Delivery Pipeline



- Sandbox 환경에서 사용자가 Streaming 데이터를 DIY 조립/생성
- 신규 생성된 데이터를 사용자의 Notebook 환경에 전달



- 실시간 변경 데이터를 사용자가 지정한 Table과 지속적으로 자동 동기화



- E-Hub Queue 적재 데이터를 Target 시스템에서 구독

- 사용자가 **실시간 데이터를 목적에 맞게 즉시 활용**할 수 있도록 Delivery Pipeline 제공
- VM/Notebook으로 안정적인 적재 지원