# Why Are We Near Real-Time?

Kristo Raun
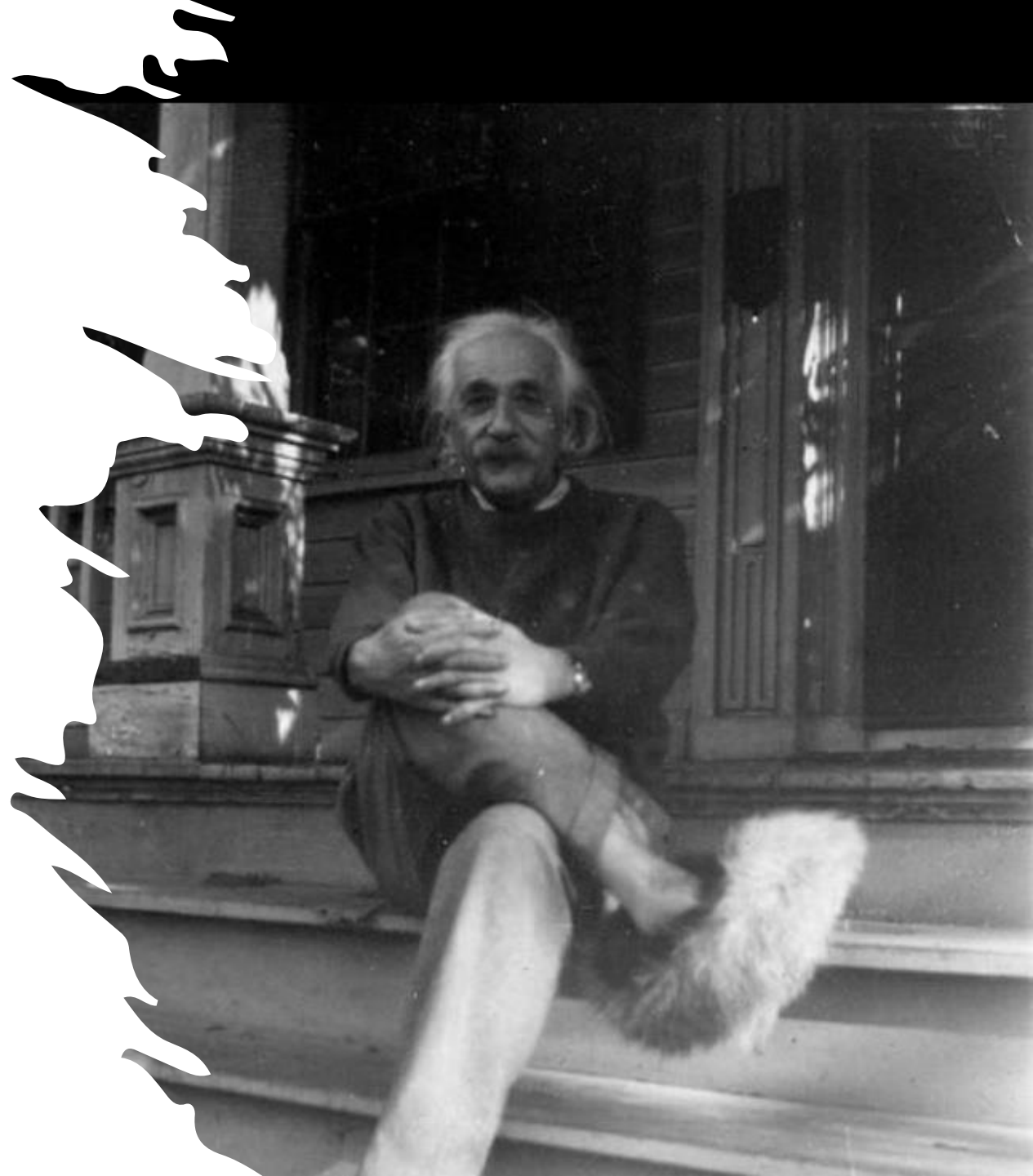
Introduction to Near Real-Time Data Analytics

August 2022
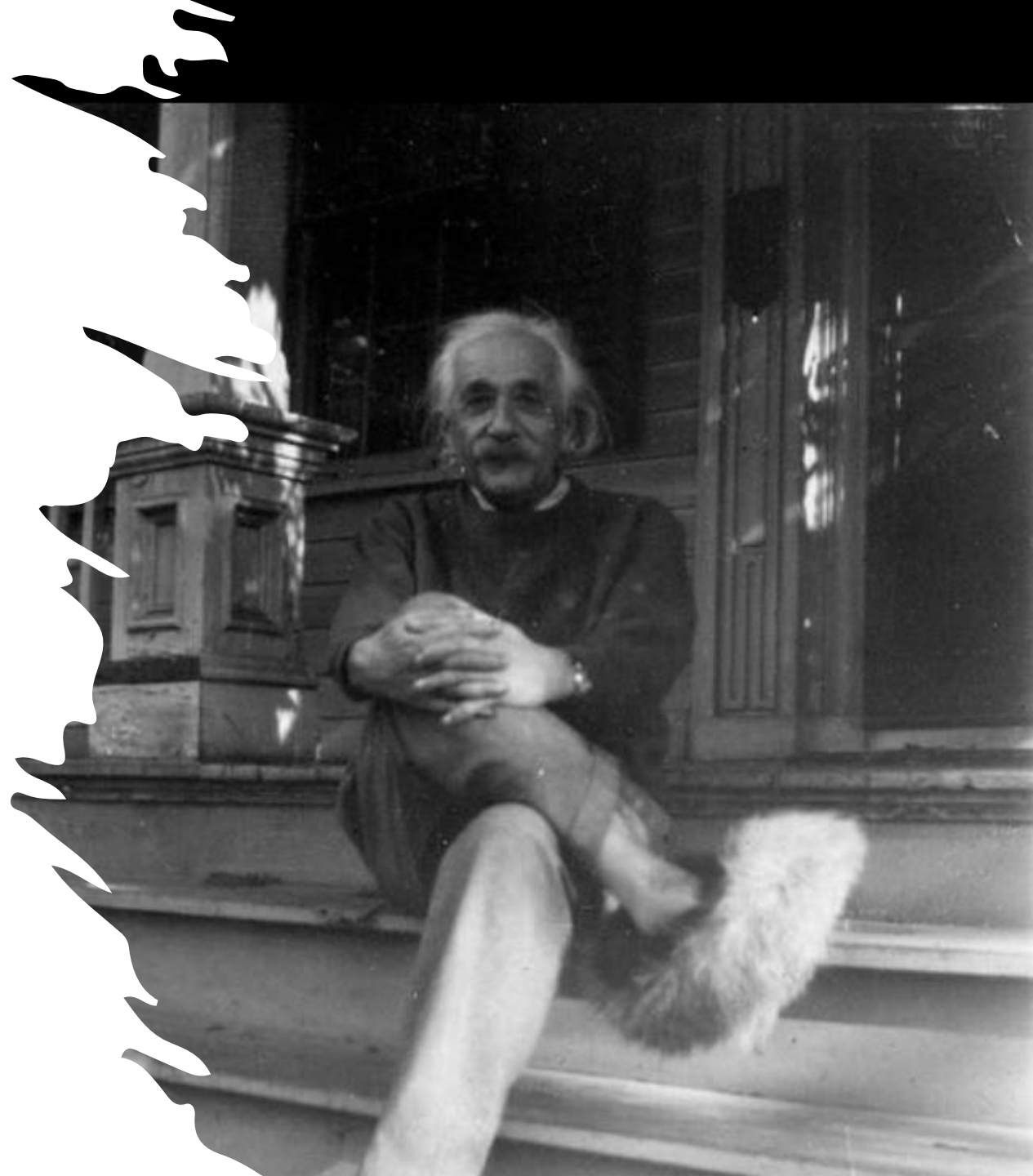
UNIVERSITY OF TARTU
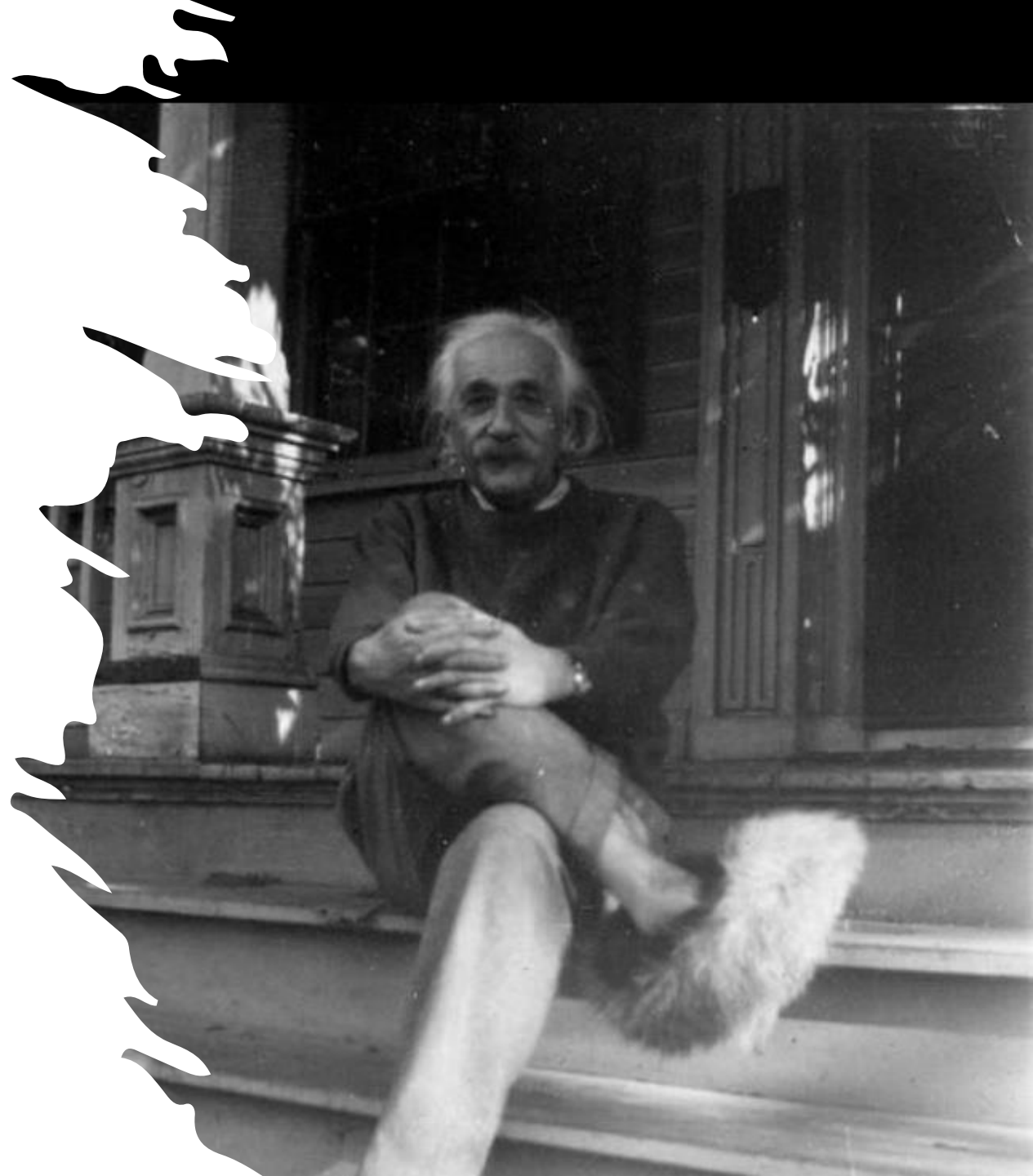
# Everything is relative

# Everything is relative

When you sit with a nice girl for two hours you think it's only a minute.
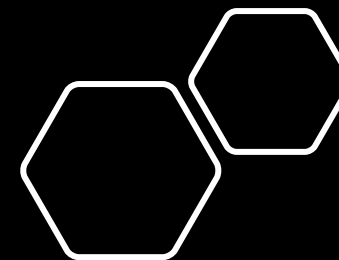
# Everything is relative

When you sit with a nice girl for two hours
you think it's only a minute.

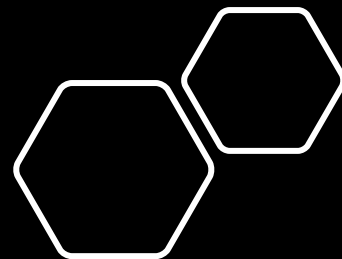When you receive streaming data for a minute
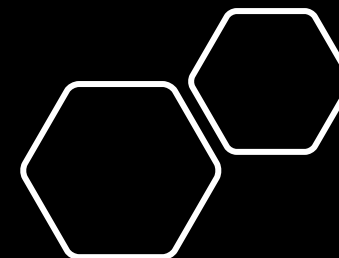you think it's two hours.

# About me

# About me
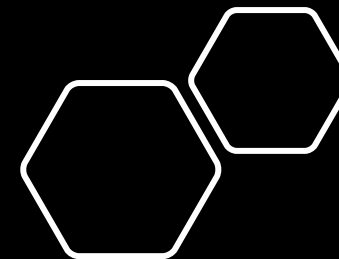
2008 – BA start

# About me

2008 – BA start

2011 – working w data
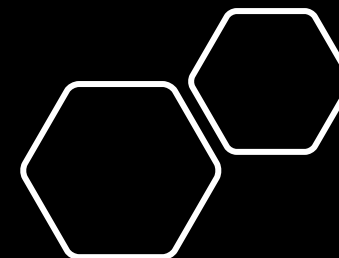
# About me

- 2008 – BA start
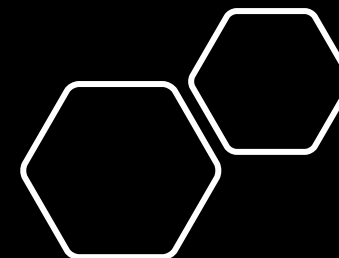
- 2011 – working w data

- 2015 – MSc start, 1st IT job

# About me

- 2008 – BA start

- 2011 – working w data

- 2015 – MSc start, 1st IT job

- 2018-19
    - – 1st child
    - – live in farm
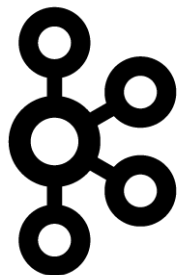    - – data engineer

# About me

- 2008 – BA start

- 2011 – working w data

- 2015 – MSc start, 1st IT job

- 2018-19
  - 1st child
  - live in farm
  - data engineer
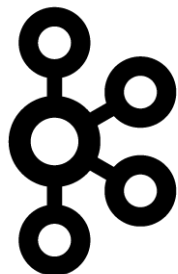
- 2021 – PhD start
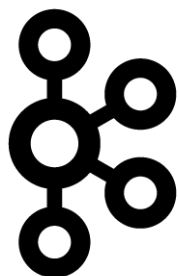
# Agenda for the week

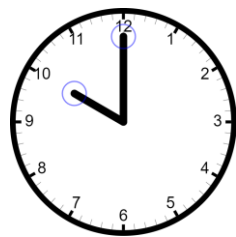# Agenda for the week

# Agenda for the week



Kafka



Flink

# Agenda for the week


Kafka


Flink


Visuals

# Agenda for the week
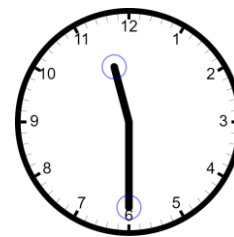


Kafka



Flink



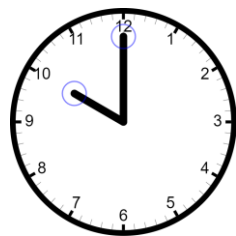Visuals

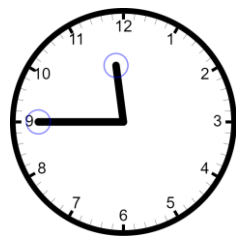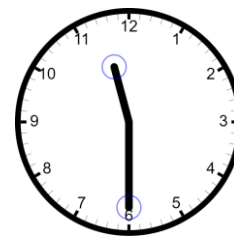

Project



Project
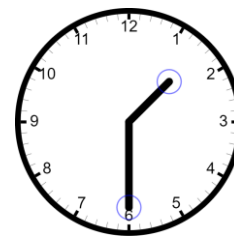
# Agenda for today

# Agenda for today

Why are we near real-time?

# Agenda for today

Why are we near real-time?

Apache Kafka setup
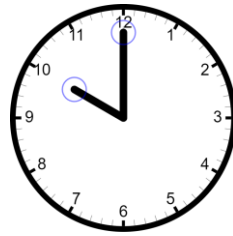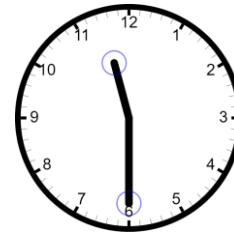
# Agenda for today

Why are we near real-time?

Apache Kafka setup

Apache Kafka practice

# Why?

# Why?

# Why?

# Why?

# Why?

# Why?

# Why?

# Timeline

1992

Bill Inmon
*EDW*

Batch processing, data warehouses

# Timeline

# Timeline

1992        1996        2000

Bill Inmon
*EDW*

Ralph Kimball
*Facts & Dimensions*

Daniel Linstedt
*Data Vault*

Batch processing, data warehouses

# Timeline

| 1992 | 1996 | 2000 | 2006 |
|------|------|------|------|

Bill Inmon
*EDW*

Ralph Kimball
*Facts & Dimensions*

Daniel Linstedt
*Data Vault*

Hadoop

Batch processing, data warehouses  ………  data lakes

# Timeline



1992 — Bill Inmon — *EDW*
1996 — Ralph Kimball — *Facts & Dimensions*
2000 — Daniel Linstedt — *Data Vault*
2006 — Hadoop
2011 — Kafka, Flink

Batch processing, data warehouses ……… data lakes ……… stream processing

# Timeline

# Timeline



World Wide Web      Google Search      Facebook      iPhone

1992      1996      2000      2006      2011      2014

Bill Inmon
*EDW*

Ralph Kimball
*Facts & Dimensions*

Daniel Linstedt
*Data Vault*

Hadoop

Kafka   Flink

Spark

Batch processing, data warehouses  ………  data lakes  ………  stream processing

# Role

# Role

# Role

# Role

# Role

# Role

# Role

# Role

# Role

# Data architecture

## Ingest

- Iot Devices
- Logs, Files
- Customer data, Financial transactions
- Weather data
- Business Apps

Event Hubs

Azure blob storage

IoT Hub

## Analyze

Continuous Intelligence/Real-time analytics

Stream Analytics

Reference Data
SQL DB, Blob store

Real-time scoring
Azure ML service

## Deliver

**Alerts and actions**
Event Hubs, Service Bus, Azure Functions etc

**Dynamic Dashboarding**
Power BI

**Data Warehousing**
Azure Synapse Analytics

**Storage/ Archival**
SQL DB, Azure Data Lake Gen 1 & Gen 2, Cosmos DB, Blob storage, etc

# Accepted Latency

## Time Dimension v/s SLA

SLA is order of Hours / Day

**Batch**



- Pre generated reports
- Cross grain resolution - trends

# Accepted Latency

## Time Dimension v/s SLA

SLA is order of Hours / Day

**Batch**



- Pre generated reports
- Cross grain resolution - trends

SLA is of order of Mins / Hour

**Near Real Time**



- Adhoc queries
- Mid resolution – aggregated counters

*Source: Practical Real-time Data Processing and Analytics (2017)*

# Accepted Latency

## Time Dimension v/s SLA

SLA is order of Hours / Day

**Batch**



- Pre generated reports
- Cross grain resolution - trends

SLA is of order of Mins / Hour

**Near Real Time**



- Adhoc queries
- Mid resolution – aggregated counters

SLA is of order of Msec / Secs

**Real Time**



- Event Driven
- High resolution – each event counts

# Analytics

# Analytics

# Analytics

# Analytics

# Analytics

# Analytics

Scalability

# Challenges



Ordering

Consistency

Fault Tolerance

# Recap



- Data generation boom
- Data processing shift
  - From batch to streaming
- We need a live view
  - Operational analytics
  - Near real-time (seconds/minutes latency)
- Challenges
  - Scalability – Ordering – Consistency – Fault Tolerance

# Not discussed but important

- Reverse ETL
- Machine Learning
- DataOps
- Security
- Data Management



It's Super Important

# Kafka

# Messaging

# Messaging

# Messaging

- Hello, would you like to hear a TCP joke?

# Messaging

- Hello, would you like to hear a TCP joke?
- Hello, yes, I'd like to hear a TCP joke.

# Messaging

- Hello, would you like to hear a TCP joke?
- Hello, yes, I'd like to hear a TCP joke.
- OK, I'll tell you a TCP joke.

# Messaging

- Hello, would you like to hear a TCP joke?
- Hello, yes, I'd like to hear a TCP joke.
- OK, I'll tell you a TCP joke.
- OK, I'll hear a TCP joke.

# Kafka motivation

Traditional
Message-Queue

# Kafka motivation

**Traditional Message-Queue**

producer — 1 ... 2 ... 3 → | 1 | 2 | 3 | — 1 → consumer

— 2 → consumer

— 3 → consumer

**Kafka**

producer — 1 ... 2 ... 3 → | 1 | 2 | 3 | — 1 ... 2 ... 3 → consumer

— 1 ... 2 ... 3 → consumer

— 2 ... 3 → consumer

- Distributed streaming platform
  - Publish & subscribe
  - Store streams durably
  - Process streams as they occur

# Kafka cluster

Cluster

broker_id=1

broker_id=2

broker_id=3

# Kafka topic

The core abstraction Kafka provides for a stream of records is the **topic**.
A topic is a category or feed name to which records are published.

# Kafka topic

*producers*

E-shop

Retail

# Kafka topic

producers

kafka

E-shop

Retail

Cluster

broker_id=1

broker_id=2

broker_id=3

topics

Sales

# Kafka topic

producers

kafka

E-shop

Retail

Cluster

broker_id=1

broker_id=2

broker_id=3

topics

Sales

Returns

# Kafka topic

*producers*                                    *kafka*

```
E-shop
```

```
Retail
```

Cluster

broker_id=1

broker_id=2

broker_id=3

*topics*

———————→  Sales

———————→  Returns

———————→  Visits

# Kafka topic

*producers*

*kafka*

*consumers*

E-shop

Retail

Cluster

broker_id=1

broker_id=2

broker_id=3

Aftersales

*topics*

→ Sales

→ Returns

→ Visits

# Kafka topic



producers

kafka

consumers

E-shop

Retail

Cluster

broker_id=1

broker_id=2

broker_id=3

Aftersales

Partners

topics

Sales

Returns

Visits

# Kafka topic



producers

kafka

consumers

E-shop

Retail

Cluster

broker_id=1

broker_id=2

broker_id=3

Aftersales

Partners

Marketing

topics

Sales

Returns

Visits

# Kafka topic

producers

kafka

consumers

E-shop

Retail

Cluster

broker_id=1

broker_id=2

broker_id=3

Aftersales

Partners

Marketing

Finance

topics

Sales

Returns

Visits

# Kafka message?

*producers*

# Kafka message?

- (Key)
- Value

*producers*
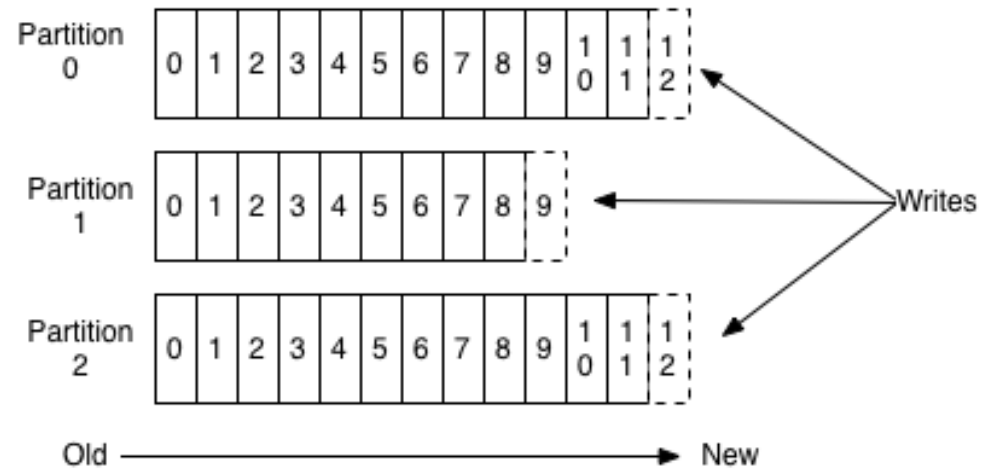


E-shop
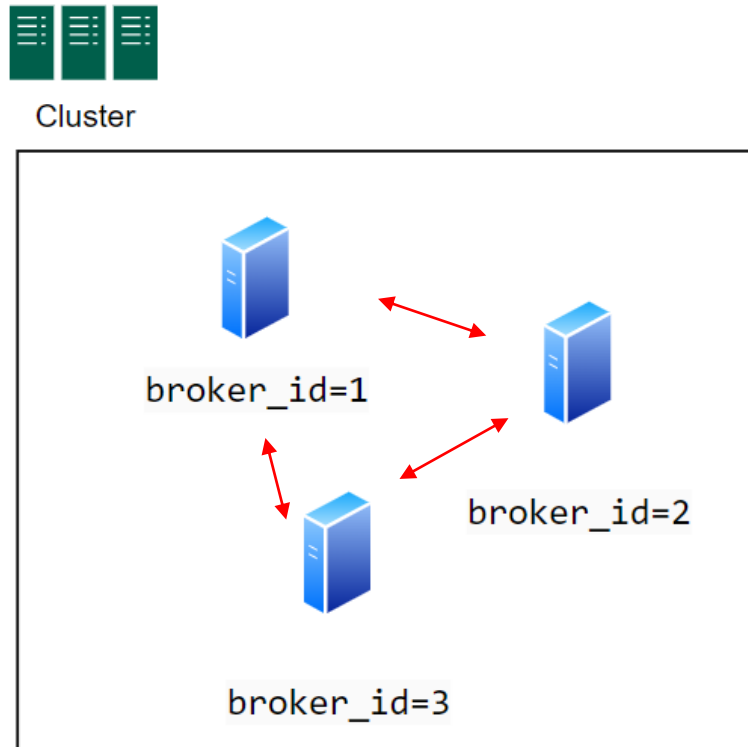
Retail

# Kafka partitions

- Scalability



Anatomy of a Topic

# Kafka replication

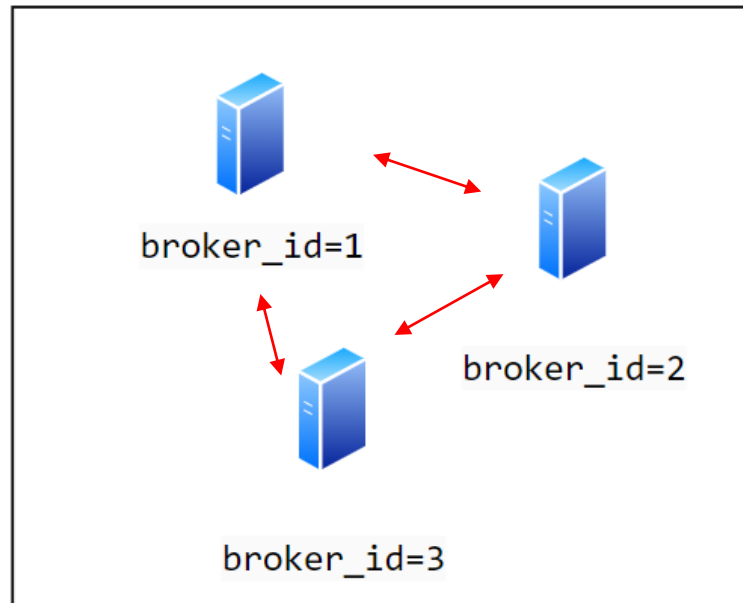- Fault-tolerance: if a broker is down, another broker can serve the data

# Kafka replication

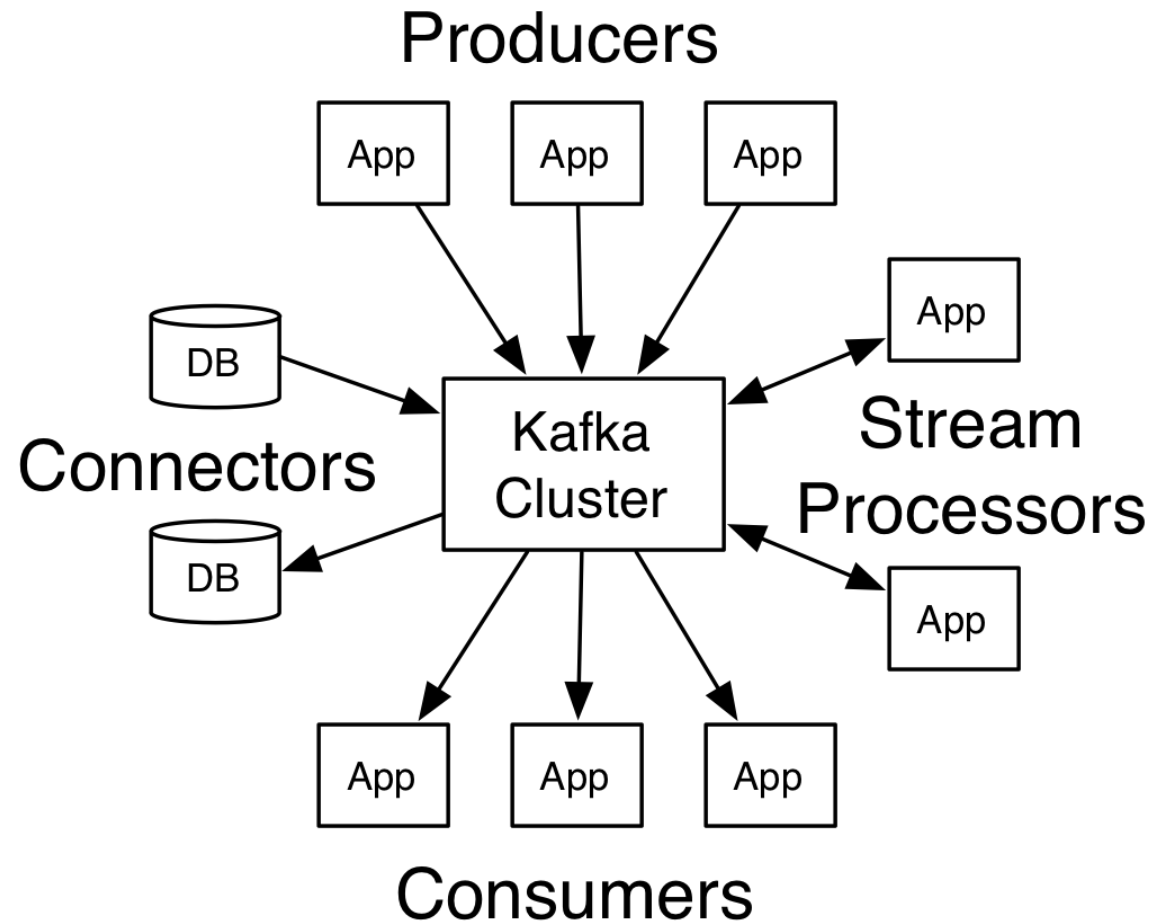- Fault-tolerance: if a broker is down, another broker can serve the data



ISR = In-Sync Replica

# Example

# Kafka core

# Kafka extended

- Schema registry
- Kafka Connect
  - CDC (Change Data Capture)
- Zookeeper, Kafka Raft (Kraft)
- ksqlDB