

Stream processing and analytics

Quiz₁ → ~ 6th Contact Session (7th Contact Session)
 → ~ 30min { No. of questions = 10 ; each 0.5 mark }
 → Window { Couple of days including weekends }

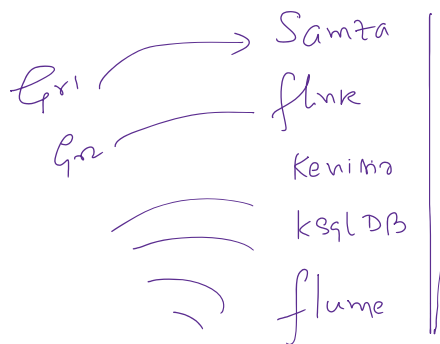
Quiz₂ → ~ 13th Contact Session
 → 30min { No. of questions = 20 ; each $\frac{1}{4} = 0.25$ mark }
 → Window (3 days { Fri, Mon })

Content: Apache Kafka, Apache Spark open access documentation

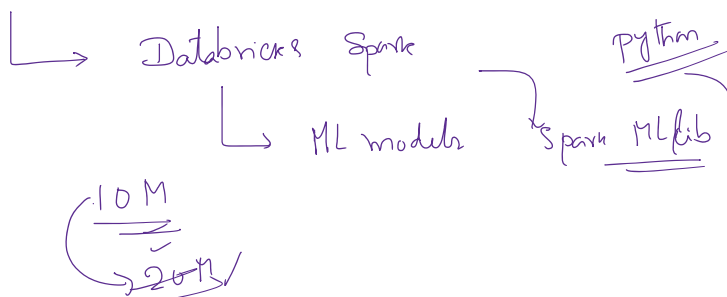
⊕
 Slides

Assignment 1 (Exponential learning)

↳ Explore streaming platforms
 ↳ mapped to individual groups



Assignment II



2011

Quiz I
5

Quiz II
5

A1
10

A2
10

Total
30

⇒ group allocation (as per your choice) 2-3 weeks time: groups finalization
~ 5th Contact Session

30 min/day {Excluding Assignment questions}

BF / Coffee time

Data Engineering / Data Architect

CS: 16

CS₁ CS₂ CS₃ CS₄ CS₅ CS₆ CS₇ CS₈ CS₉ CS₁₀ CS₁₁ CS₁₂ CS₁₃ CS₁₄ CS₁₅ CS₁₆

① Understanding reqs

② Proposing Data Arch.

③ Evaluation of Data Arch.

④ Aspects — Non-functional Requirements

⑤ Kafka Arch.

Kafka Graph aspects
streaming algorithms
and time/space complexity issues

PySpark

~~Stream~~

~~Orbital
Depends~~

Spark structured Storage

Architecture and execution of jobs on Spark

Stream processing:

↳ Computations on potentially endless and constantly evolving

Source of data

Tired
Computation model

Create
Retrieve

changing...

Functions
Data representation

Datatypes

String
Float
i1

Create
Retrieve
Update
Delete

changing contents/data
original contents/data

Data types

String
float
int

Moving Data towards code

Data at rest
(persisted data)

def computeSum(lst):

Sum = 0

for element in lst:

Sum += element

return Sum

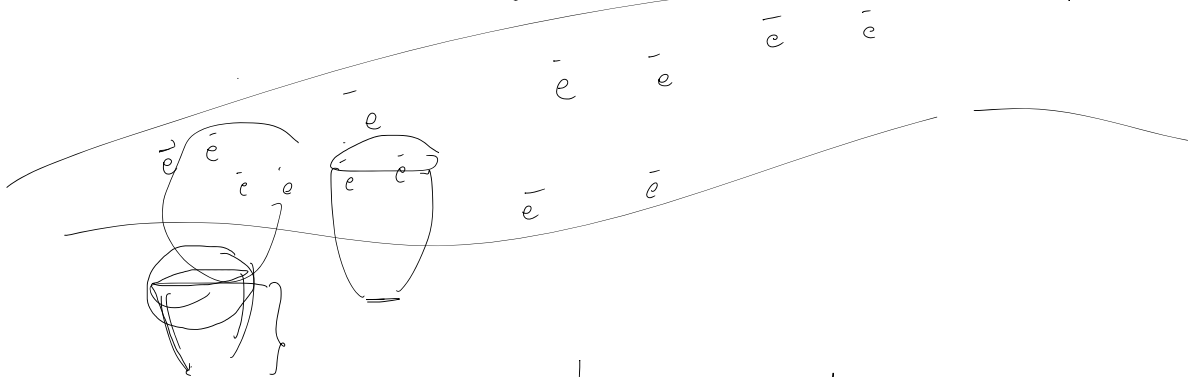
Query

aggreg / func. etc

Data fetched either from file/data store

streaming as computational model

Continuous query processing



Code is executed on fast moving Data

Data stream:

It comprises of endless stream of immutable data (events)

Size of event ~ KB (< 1MB)

Event

is an immutable fact regarding something that occurred within the System under consideration

Data records in the context of data streaming are called events



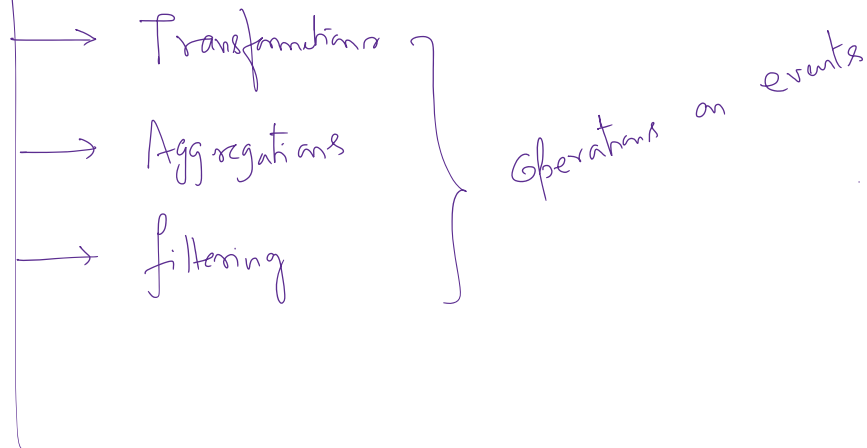
card swipe at pos



Card ends in x x 4593
Spent 20,000 at Shopmart

time stamp

⇒ Un-mutable { Cannot modify Content of the event }



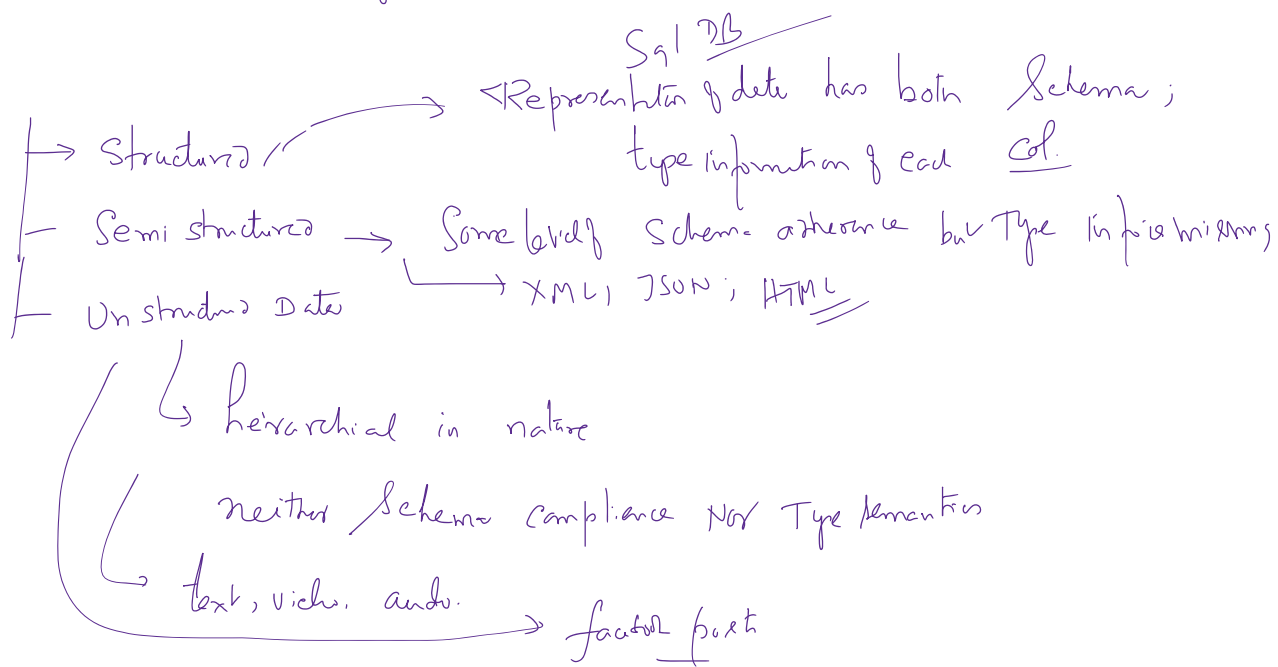
Examples of stream processing

① App logs : time stamp : <context> | status
 { 2025-07-27:18:00:05:10 : exception / failure / warning / long transaction / error }

② Web Analytics:

③ Real time pricing of a stock

Data Types



Size of the data

Size of The data

$$1 \text{ Byte} = 8 \text{ bits}$$

$$kB \approx 1024 \text{ Bytes} \approx 10^3 \text{ Bytes}$$

kB

MB

GB

TB

$$1024 \text{ kB}$$

$$1024 \text{ MB}$$

$$1024 \text{ GB}$$

$$\approx 10^3 \text{ kB}$$

$$10^3 \text{ MB}$$

$$10^3 \text{ GB}$$

$$1 \text{ TB} \approx 10^3 \text{ GB}$$

$$1 \text{ TB} \approx 10^3 \text{ GB} = 10^3 (10^3 \text{ MB}) = \underline{\underline{10^6 \text{ MB}}}$$