data

Complexity → Collection → Dataflow → Processing tier → In memory Storage → . . . .

Merge ?

→ transformation ⎫
→ aggr          ⎬ Stream processing
→ filtering     ⎭
— Computationally intensive

Collection
— Interaction pattern

**Why cant Merge**

① Different data Sources
② Different Sources Connects  ⎫ Collection
③ Different formats           ⎭ — Interaction pattern

**Why do we need Dataflow tier?**

Modify Collection
Adding different
source connectors

Strong Coupling between
Collection tier
Processing tier

Collection

Processing tier

Dataflow        Regiment
chr

② **Impedance missmatch**

Collection tier → Data flow tier ? → Processing tier

Rate of ingestion
$R_i$

Rate of Processing
$R_p$

Manage The Throughput w.r.t Collection

Manage The throughput w.r.t Collation Priority

$$R_i \ll R_p \quad \text{(rare : monitoring use cases : atmost once - In general)}$$

$$R_i \geqslant R_p \quad \text{( Back pressure problem)}$$

$$R_i \gg R_p \quad \text{( most common)}$$

Amount     TxnID     timestmp
                     YYYY: MM: DD: H: m: s —

| | | | | | |

Define a Scheme

Amount  $t_{xnID}$  timestmp  ← Nothing to do with BL

In memory Table

→ data frame [ 'timestmp'] = data frame [ 'timestmp'].
                                    astype ("timestmp")

$$R_i \gg R_p$$

Collection          $P_2$
tier
                    $P_1$
Producer

Collection → Data flow → Processing tier        More Consumers
                                         $P_3$

Improve processing capacity
        $6 \times 10^4 \, kB$
        $6000 \times 10 \, kB$  6

$$R_i - R_p = 100 \text{ msgs/min}$$

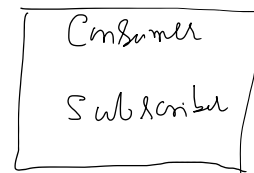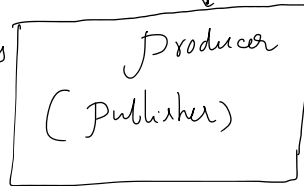partitions

## Kafka Architecture

Where you are publishing
those emit : **topic**

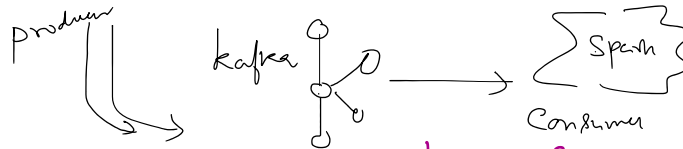Publisher — Subscriber



1. Producer sends <u>messages</u> (events)
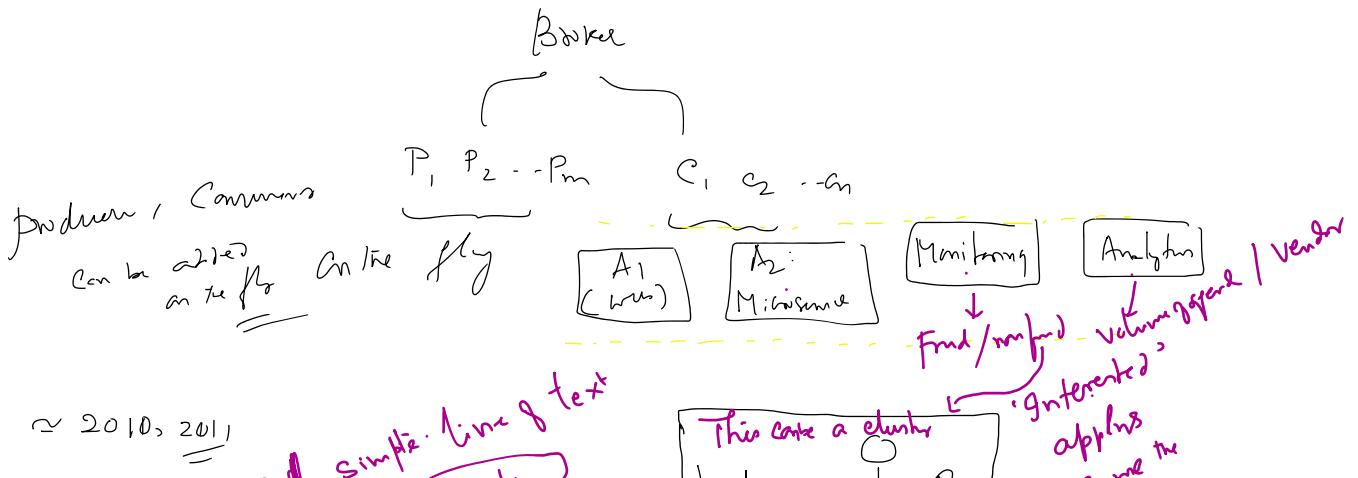
   primitive data records

2. Broker: responsible for recieving messages from producer and store them in a local storage

3. Consumer: Reads messages from broker and procers them
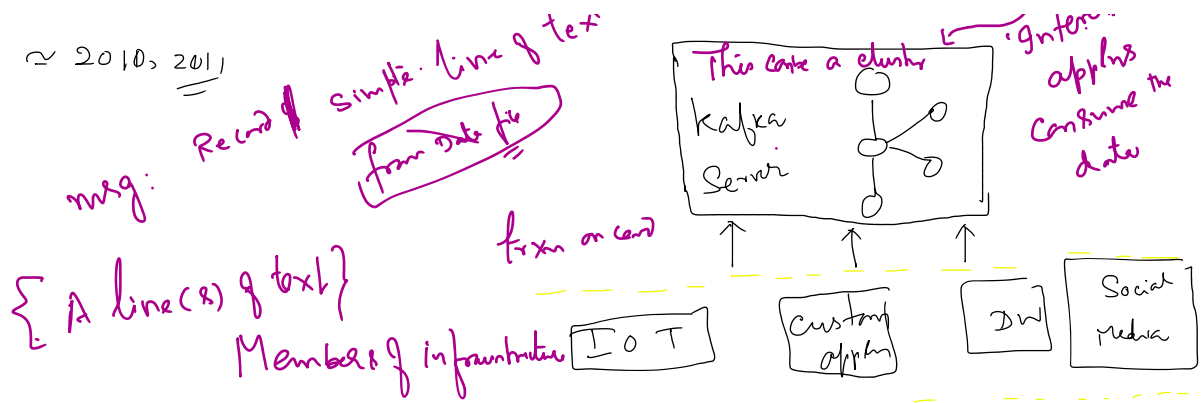
   Broker: A Name for Kafka Server

=> Broker acts as proxy between producers and consumers

Broker

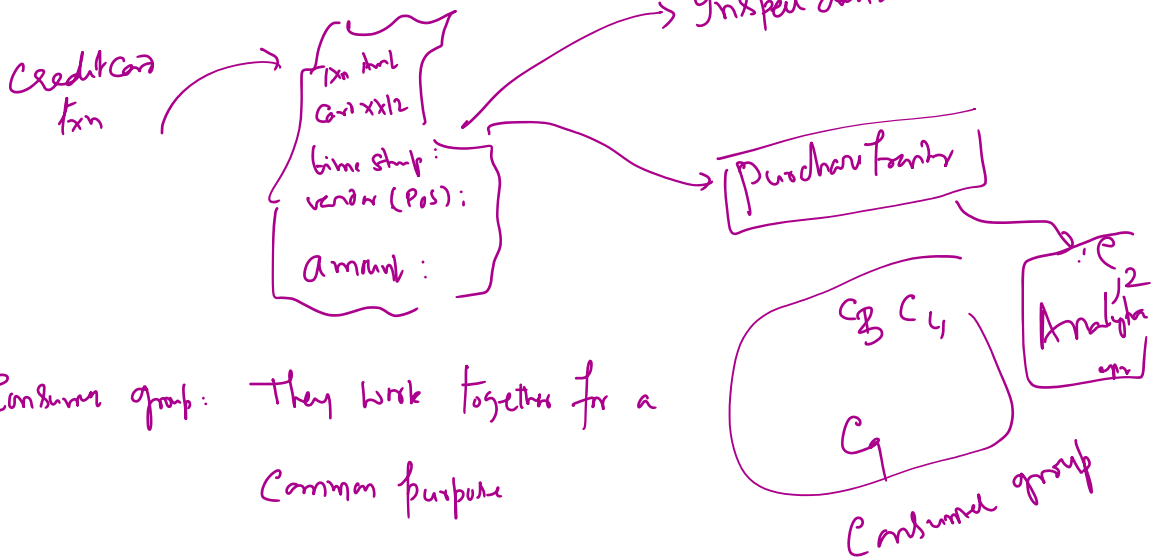$P_1, P_2 \cdots P_m$    $C_1, C_2 \cdots C_n$

Producers, Consumers
can be added on the fly

A1 (web)    A2 Microservice    Monitoring    Analytics

Frud/raufird — volume depend / vendor
· Interested applns

~ 2010, 2011

a simple line of text

This care a cluster

≈ 2010, 2011

msg: Record of simple line of text
                                    from Data file

{ A line(s) of text }

                    txn on card

Members of infrastructure      [ I O T ]     [Custom appln]    [DW]   [Social Media]

topic: Unique name for a Data stream
        ↳ Creating a topic is design-time decision                    [C₁]

Credit Card  →  { Txn Amt           →  Inspect Event
Txn              Card XX12
                 time stmp:
                 vendor (PoS);  →  [Purchase Frontier]

                 amount: }

                                                    [Analytics appn]²

Consumer group: They work together for a            C₃ C₄

               Common purpose                        Cₐ

                                            Consumed group

Kafka cluster
_____



Kafka broker B₁          B₂

                B₃

        Kafka cluster
        _____

                Broker
               { Subscribe (Topic T)                    C₁

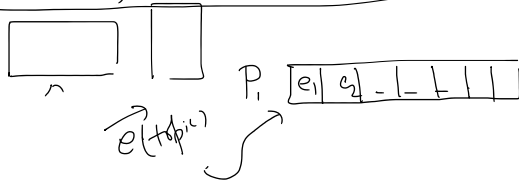P₁    P₂    P̶₃̶

                                            C₁

Broker is responsible for storing in local Stage

↳ Broker will be facing Storage issues

Break The topic into multiple smaller parts and distribute it over multiple Computers in kafka cluster → partition
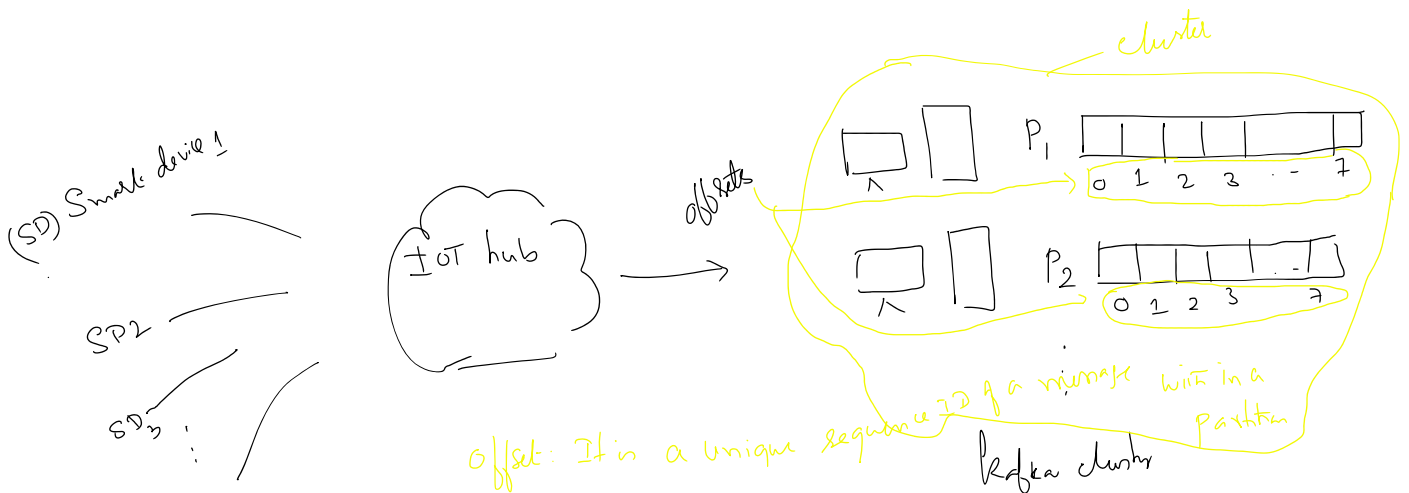
partition : It is smallest unit sitting on a single machine

P₁ | e₁ | e₂ | _ | _ | | | |

e(topic)

A partition is small independent portion of a topic

No overlap of events across the partitions

No. of partition in a topic is a design decision

cluster

(SD) Small device 1

SP2

SD₃ :

IoT hub →

offsets

P₁ | | | | | |
0  1  2  3  -  7

P₂ | | | | | |
0  1  2  3  -  7

Kafka cluster

Offset: It is a unique sequence ID of a message within a partition

⊙ It will be automatically assigned by the broker to every message as it arrives in the partition

⇒ messages are stored in a partition as 'append'

Offset ID is arrival order number for msg

Topic Name $\longrightarrow$ partition $\longrightarrow$ partition number $\longrightarrow$ offset number

$\phantom{Topic Name} \smile$ Broker helper (c) $\smile$

$*$ Partition is orted for Scalability

· Estimate for number of partitions

Number of partition $(N)$ = max $\{N_p, N_e\}$

$$N_p = \frac{\text{Total System Throughput}}{\text{Max. Throughput Producer writing msg into single partition}}$$

$$N_p = \frac{T_t}{T_p}$$

$$N_c = \frac{T_t}{T_c} = \frac{\text{Total System Throughput}}{\text{Max Throughput Consumer reading a msg from a single partition}}$$