① Kafka — Partitions

② Full Text Search

③ NLP Techniques

④ Elastic Search

---

Swiggy

order
Service

Kafka Topic = Placed Orders

order-id

Notification
Service

I

F  E  D  C  B  A

I   →   C

II  →   ⊗D E

III →   ✗ B

Logistics
Service

II

get

E

offset

reset offset (B)

Restaurant
Service

IV

Producer
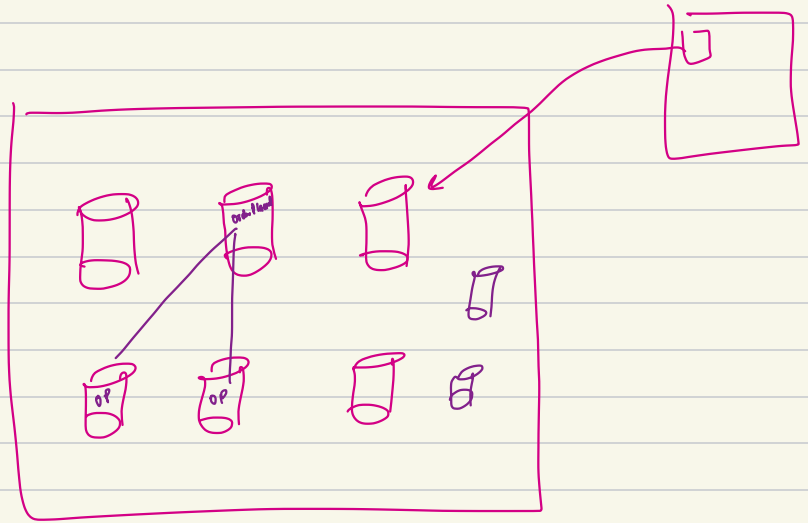
Consumers

Given that I want to build Kafka as a dependable / fault tolerant system,

the fundamental concept of replication should be employed :)

There may be a situation that 1 entire machine may NOT be sufficient to store all msgs of 1 queue also.

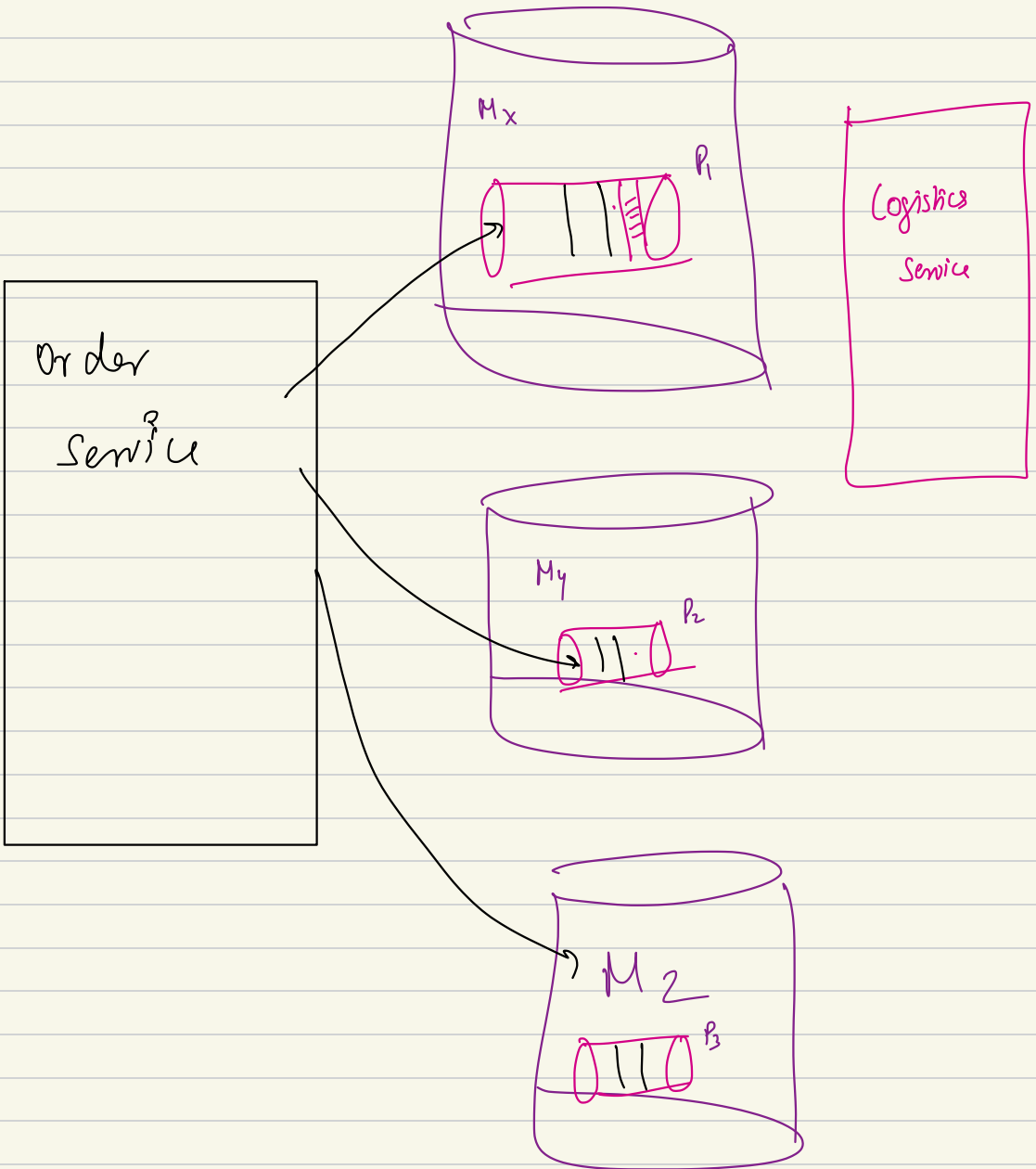The usual way to solve this → SHARDING

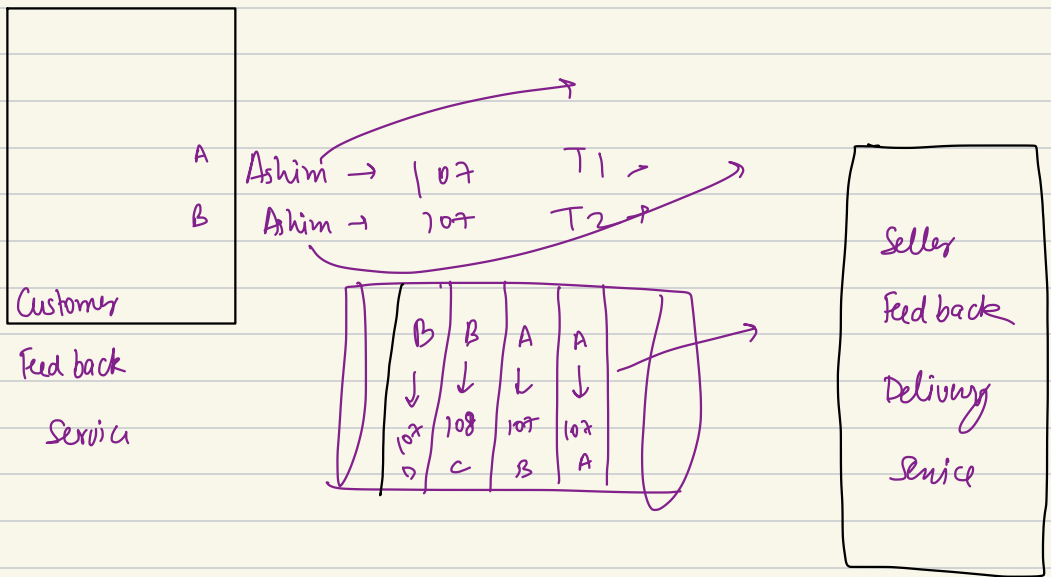But we can't use sharding directly before understanding it...

Queue is FIFO

---

PROBLEM: when you divide a queue across multiple machines, it No longer remains 1 queue; it becomes multiple pieces of 1 queue.
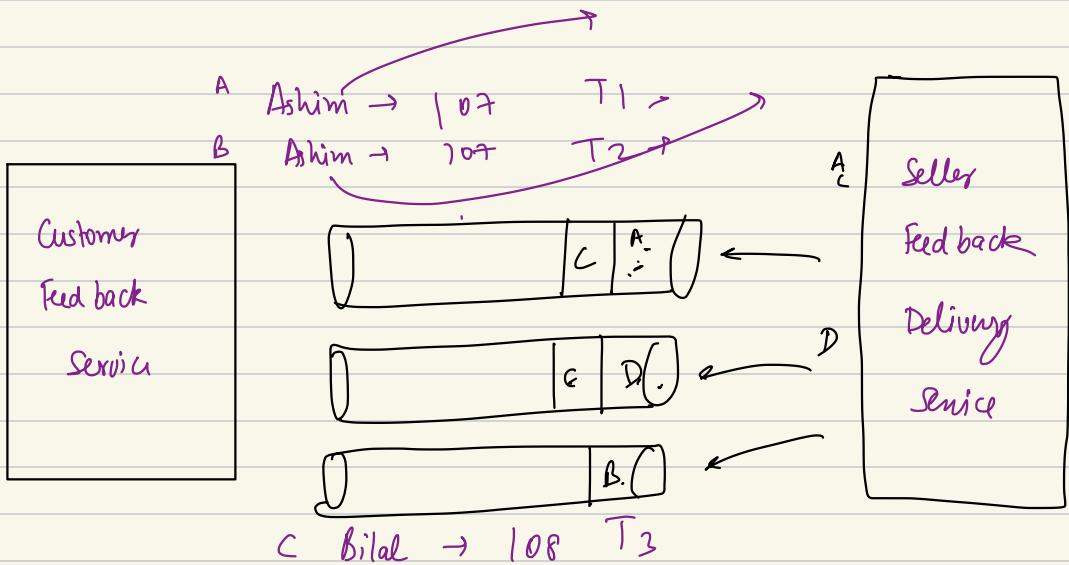
Kafka Topic | Queue  $\equiv$  Placed Orders



$M_x$

$P_1$

Logistics Service

Order Service?

$M_y$

$P_2$

$M_2$

$P_3$

The Kafka Topic when partitioned, no longer gives FIFO gaurentees across all partitions.

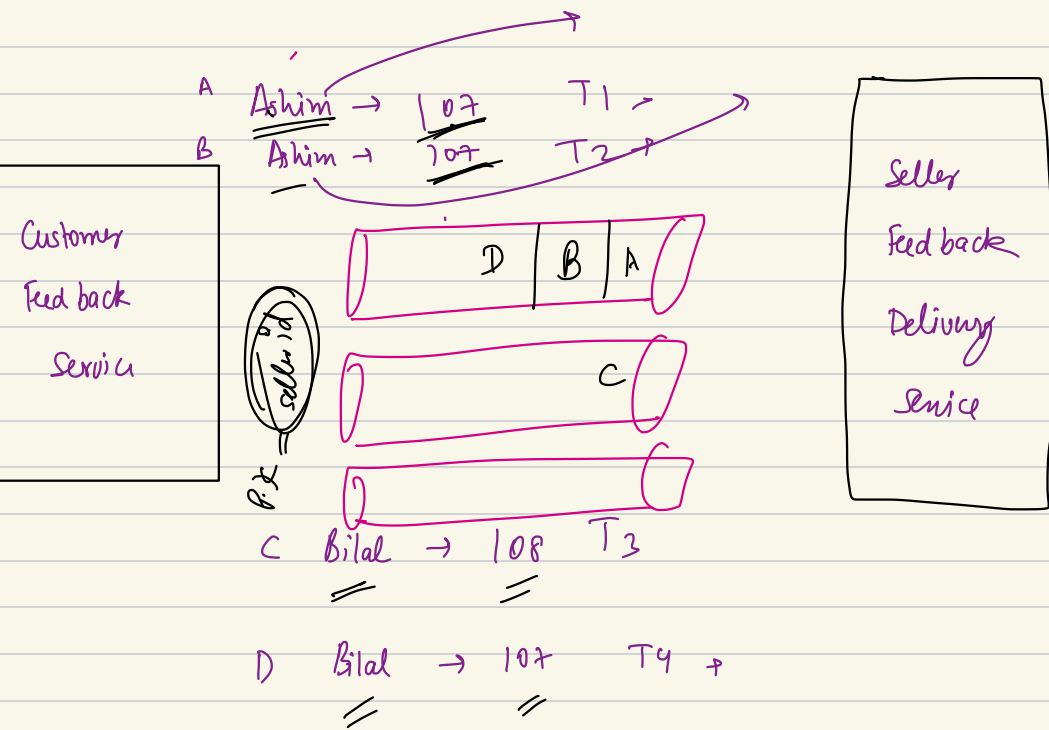Only for msgs of the same partition, you can have FIFO gaurentees.

Customer
Feed back
Service

A  Ashim → 107    T1
B  Ashim → 107    T2

B  B  A  A
↓  ↓  ↓  ↓
107 108 107 107
D  C  B  A

Seller
Feed back

Delivery
Service

C  Bilal → 108  T3

D  Bilal → 107   T4

A   Ashim → 107    T1

B   Ashim → 107    T2

Customer
Feed back
Service

A
C

Seller
Feed back
Delivery
Service

C   Bilal → 108   T3

D   Bilal → 107    T4

E
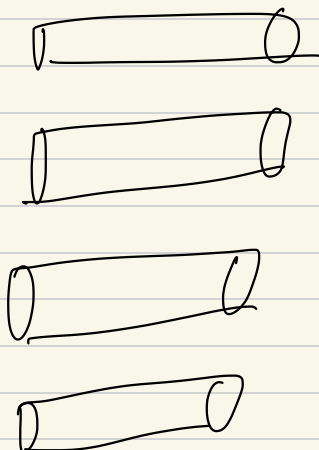
Order within the same partition is NOT lost;
but global ordering is LOST

Kafka has to live with it
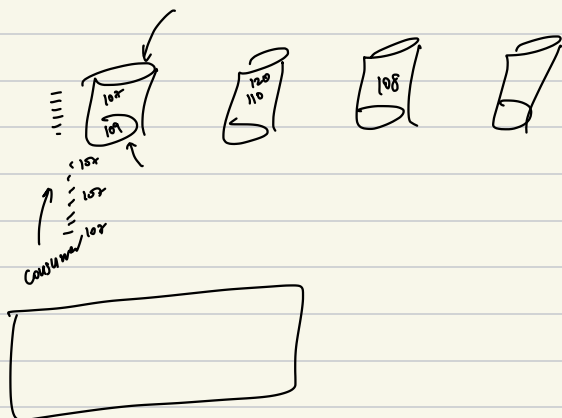
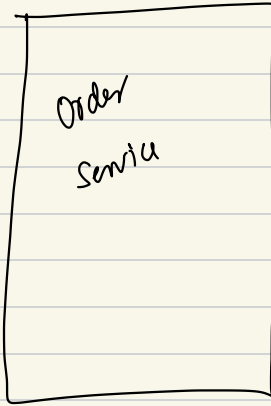A  Ashim → 107    T1
B  Ashim → 107    T2

Customer
Feed back
Service

seller id

| | D | B | A |

| | | C |

C  Bilal → 108   T3

D  Bilal → 107   T4

Seller
Feed back

Delivery

Service

lp

Order

Restaurants

$T_1$
$T_2$
$T_3$

$Ts_{III}$
$Ts_{I'}$
$T_4$

$Ts_{II}$

$10x$
$10x$

$120$
$110$

$108$

$< 15x$
$10x$
columna $10x$

SWIGGY

Topic ≡ Placed Orders → 4 partitions

Order Senvia

poll (PO, I)
q

I

| O. | 2 | E. | A |

poll (PO, II)

II | P | D. | C | B |

Read from head

III | OM | G | I | n |

IV | R | N | L | R |

(R.S)

II) Restaurant Senvite

Consumer group

Notification Senvia

Consumer group

| I | → ~~2~~ 2 |
| II | → ~~X~~ C |
| III | → G |
| IV | → N |

:)

NS
| I | — |
| II | — |
| III | — |
| IV | — |

Topic

Consume
5 services

107 → X, A, Z, Y, C, D

108 → F, G, H

P(107+108) → | X | F | G | A | Z | Y | H | C | D |

# ELASTIC SEARCH

Linked In

Parameterized
Searches

Full Text
Search

A POST:    India is a beabuhful country

B POST :   Canada India relabasship is at
           a new low

C POST :   Canada Engineering colleges might be

           a sham.

D POST:    India lost the CWC

E POST:    I love my job. I love my life.


query: ( India cricket )


SQL 😃 😃 😃

Select * from CONTENT_TABLE
        where content % india/cricket/."
              LIKE

# NoSQL

→ Document DB seems to be the most relevant DB out of all the DBs we studied to store above posts / web-pages

But still full text search is a challenge. it is unoptimized to perform it.

---

# FIRST PRINCIPLES =>

<span style="border:1px solid">Inverse Index</span>

Glossary

Ram: 8, 10, 20, 24, 26

Aniket: 8, 80, 98, 101

Akash : 1, 100, 101

A POST:     | India is a beabthful country |

B POST :     Canada India relatianship is at
                     a new Low

C POST :     Canada Engineering colleges might be
                     a sham.

D POST:     India lost the CWC

E POST :     I love my job. I love my life.

---

## Inverted Index

India        —    A , B , D
Canada       —    B , C
Country      —    A
Beautiful    —    A
relationship —    A
Cricket      =    (D)
CWC          =    D

query:

"Cricket India"

Union
+
Rankivg ☆

How to create an inverted index.

Apache Lucene

① STOP WORD REMOVAL

a, an, the

② STEMMING

→ run, running, ran
→ crash, crashed, crashes

→ US of America, United States of America, USA

happy happiness happier → happy

India, Indian, Indianess

③

③  Tokenization

POST B: Canada  India  relationship  is at
a new low

relation

Canada  India  relationship  ~~is~~  ~~at~~
~~a~~ new low

Input

  │  — Remove Stopwords
  ↓

┌──────────────┐
│              │
└──────────────┘
  │  Stemming
  ↓

┌──────────────┐
│  ≡≡≡≡≡≡     │
└──────────────┘
  │  Tokenization
  ↓

┌────┐
│ ⊥  │  vector
└────┘

Unigram

Canada

India

relations

new

low

Bi-grams

Canada    India

India     relation

Relation   new

new      low

Trigram

4-gram

5-gram

bi – skip – grams  ( 1 )

Tri – skip – grams  ( 1 )

---

Canada  India  relashsship  is at
        a new  low

after    Tokeninatias

Canada

India

relatias

new

low

Canada India

  India  relatias

relatias  new

  new  low

Content
↓
Vector of Tokens
☺

Canada relation
India new

relation low

# LinkedIn

## 1 M Posts

1 → Post ; Vector
2 → Post ; vector

3

4

1M → Post , vector

All Tokens created globally
$\downarrow$
Assign some importance score
to each Token

Tf-IDf

1M   Documents

raw document          Processed   vector

1M   P. vectors

↓ Filtration

a few million interesting Tokens

| | | | |
|---|---|---|---|
| Canada | : | [ List of doc ids | ] |
| Canada India | : | [ | ] |
| India | | [ | ] |
| ___ | | [ | ] |
| ___ | | [ | ] |
| ___ | | [ | ] |
| ___ | | [ | ] |
| ___ | | [ | ] |
| ___ | | [ | ] |

doc1    doc2    doc3

vector of Tokens

Set of Tokens

$— f_i$
$f_2$

$—— f_3$

$— f_5$
$f_4$

$\downarrow$ Filter

smaller set