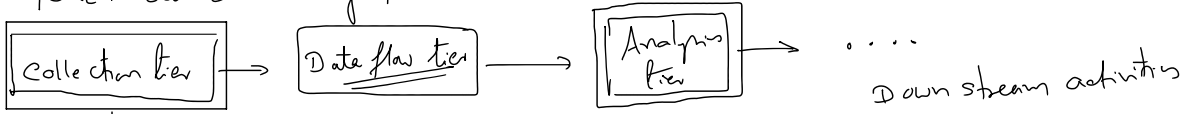


Agenda

Data flow layer

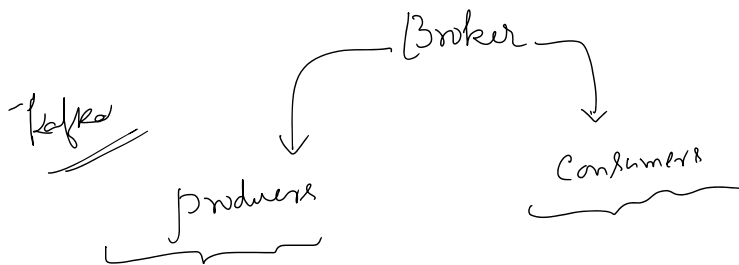
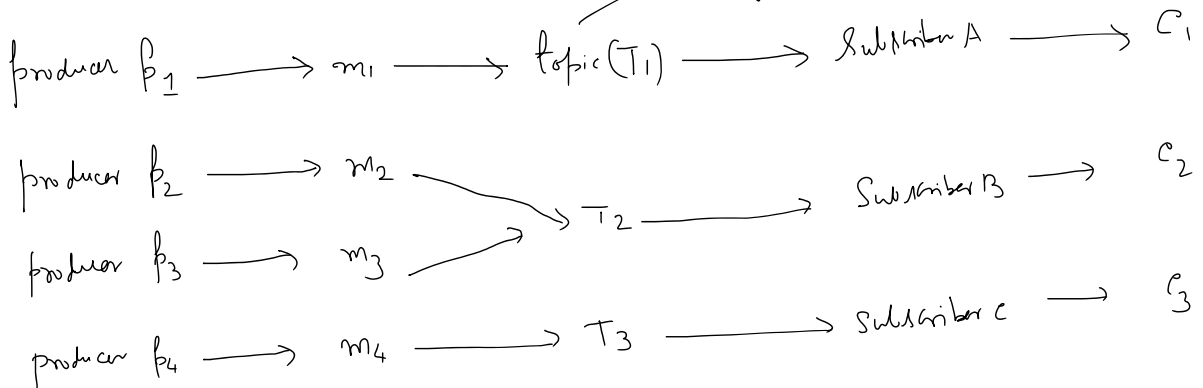
Generalised streaming Architecture



Common-interaction patterns

- client-server (Request-Response pattern). [When client must have an immediate response]
- pub-sub (producer-consumer pattern)

Complexity in terms of variety of data sources / formats
 Logical grouping of events



Data flow tier

message delivery semantics

→ Analysis tier

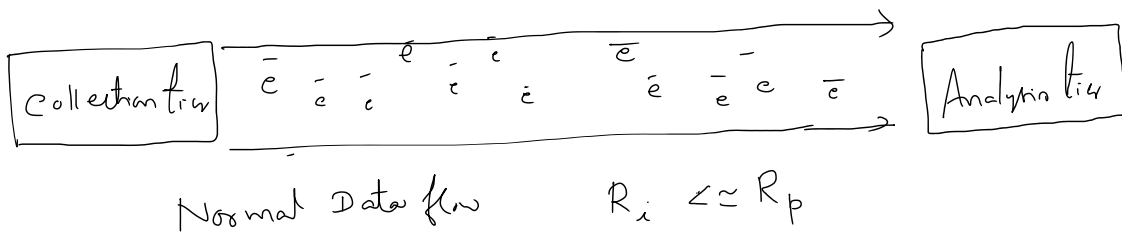
To decouple collection tier and Analysis tier

Considered as message Queue

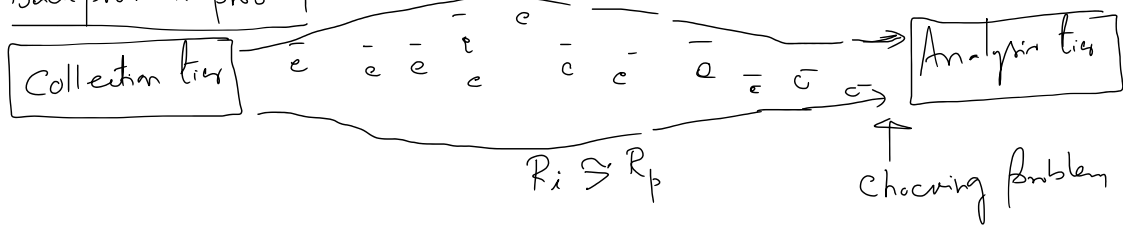
The collection tier produces events at much faster pace compared to the analysis tier

Rate of ingestion (R_i)

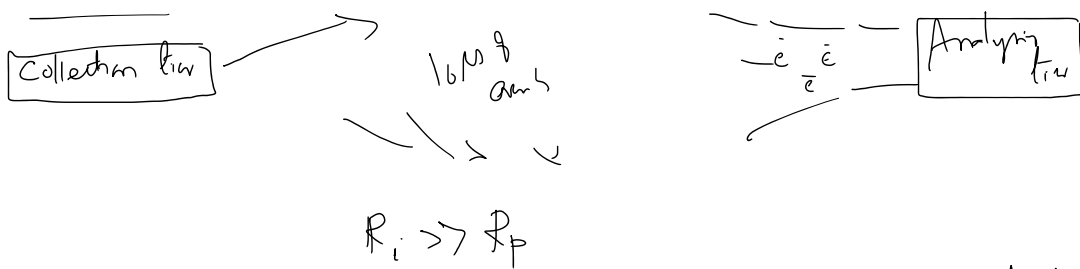
Rate of processing (R_p)



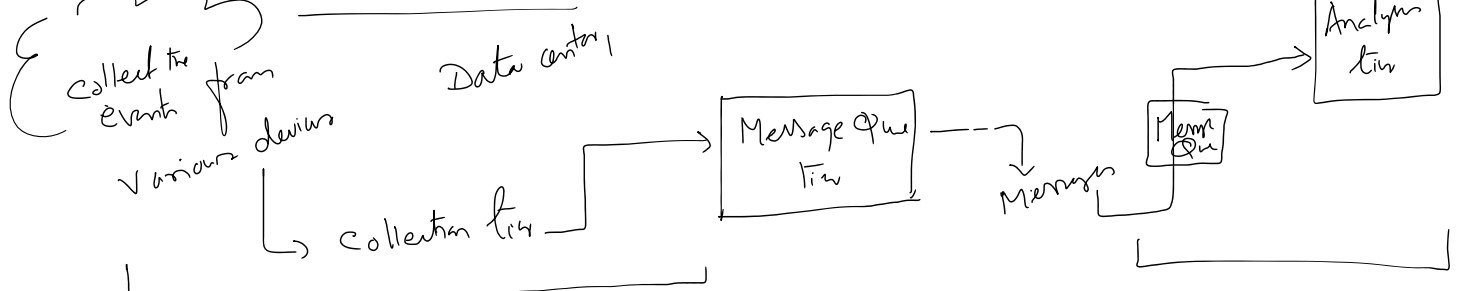
Backpressure Problem



Data loss



Durable messaging:



Durable messages provide a degree of fault tolerance allow off line consumption of events

Use case

What if the events are lost?

Architecture Design Perspective

- ① Impact on Business if the communication between collector tier and analysis tier is interrupted significantly
- ② How many days (Tolerance factor - message loss) of data the business can tolerate
- ③ Need for storage of historical data

Credit card txns data

Criteria	Inference
(i) Business Impact	
(ii) Tolerance for event loss	
(iii) Need for storing Historical Data	

Semantic :-

$$\text{Number of partitions Needed} = \max \left\{ \frac{T_t}{T_p}, \frac{T_k}{T_c} \right\}$$

T_t : Throughput of the system

T_p : Max. throughput of producer writing events onto single partition

T_c : Max throughput of a consumer reading events from a single partition

$$\text{System Throughput} = 6 \text{ GB/min} = 6 \times 10^3 \text{ MB/min} = \frac{6 \times 10^3}{60 \text{ sec}} = 100 \text{ MB/sec}$$

$$T_p = 5 \text{ MB/sec}$$

$$\frac{T_t}{T_p} = \frac{100 \text{ MB/sec}}{5 \text{ MB/sec}} = 20$$

estimated
no. of producers

$$20 = \max \left\{ 20, \frac{T_t}{T_c} \right\}$$

estimated number of
consumers

$$\frac{T_t}{T_c} \leq 20$$

$$\frac{T_t}{T_c} \leq t_c$$

$$t_c \geq \frac{T_t}{20} \Rightarrow 5 \text{ MB/sec}$$

Ex²

$$R_i = 24 \text{ MB/hr}$$

$$R_p = 2 \text{ MB/min}$$

$$R_i = \frac{24 \text{ MB}}{60 \text{ min}}$$

$$R_p = 2 \text{ MB/min}$$

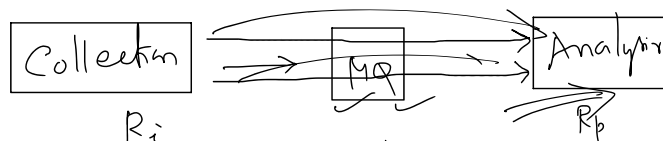
$$= 0.4 \text{ MB/min}$$

$$\frac{R_i}{R_p} = \frac{0.4}{2} = 0.2 \Rightarrow \frac{R_i}{R_p} = \frac{8}{10} = \frac{4}{5}$$

$$R_i = \frac{4}{5} R_p$$

$$\text{Avg msg size} = \underline{\underline{64 \text{ kB}}}$$

$$R_i \leq R_p$$



Exactly once, at least once using auxiliary message queue

$$|R_i - R_p| = \frac{1}{5} \times 2 \text{ MB/min}$$

1 min

0.4 MB

$$\frac{400 \text{ kB}}{25} \approx 6 \text{ kB}$$

$$1 \text{ min} \quad \frac{0.4 \text{ MB}}{64 \text{ kB}} = \frac{400 \text{ kB} \approx 6 \text{ kB}}{64 \times 4}$$

$$1 \text{ min} \rightarrow \underline{6 \text{ kB}} \text{ dr}$$

$$1 \text{ hr} : \underline{\underline{360 \text{ kB}}}$$

(k)

$$R_i = 120 \text{ MB/min}$$

~~120~~

$$R_p = 24 \text{ MB/hr}$$

$$24 / 60 \text{ min}$$

$$= \frac{4}{10} = 0.4 \text{ MB/min}$$

$$\frac{R_i}{R_p} = \frac{120 \text{ MB/min}}{0.4 \text{ MB/min}} = \frac{1200}{4} = 300$$

$$\boxed{R_i = 300 R_p}$$

$$R_i \gg R_p$$

$$n_i = \frac{120}{64 \text{ kB}} / \text{min} = \frac{\approx 2}{30} \text{ (MB}^3) \quad \frac{2000}{16} \checkmark$$

$$n_p = \frac{24 \text{ MB}}{64 \text{ kB}} / 60 \text{ min} =$$

$$= \frac{4 \times 10^8}{164 \times 10^3} = \frac{1000}{164}$$

$$R_p \approx 7 \text{ MB/min} \Rightarrow 7 \times 64 \text{ kB/min}$$

$$448 \text{ kB/min}$$

$$R_i = 2000$$

$$2000 - 7 = (1093) \times 64 \text{ kB} = \underline{\underline{1.1 \times 64 \text{ MB/min}}}$$

$$\text{Auxiliary storage} \approx 70 \text{ MB/min}$$

$$\text{Refresh hrly: } \underline{\underline{70 \times 60 = 4200 \text{ MB}}}$$