

STREAM PROCESSING AT LINKEDIN: APACHE KAFKA & APACHE SAMZA

Processing billions of events every day

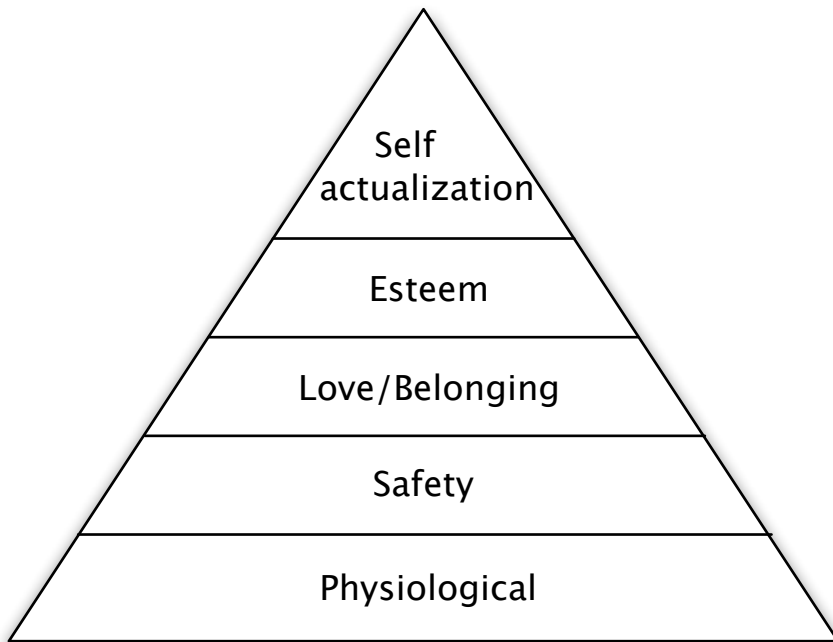
Neha Narkhede

- Co-founder and Head of Engineering @ Stealth Startup
- Prior to this...
 - ▣ Lead, Streams Infrastructure @ LinkedIn (Kafka & Samza)
 - ▣ One of the initial authors of Apache Kafka, committer and PMC member
- Reach out at @nehanarkhede

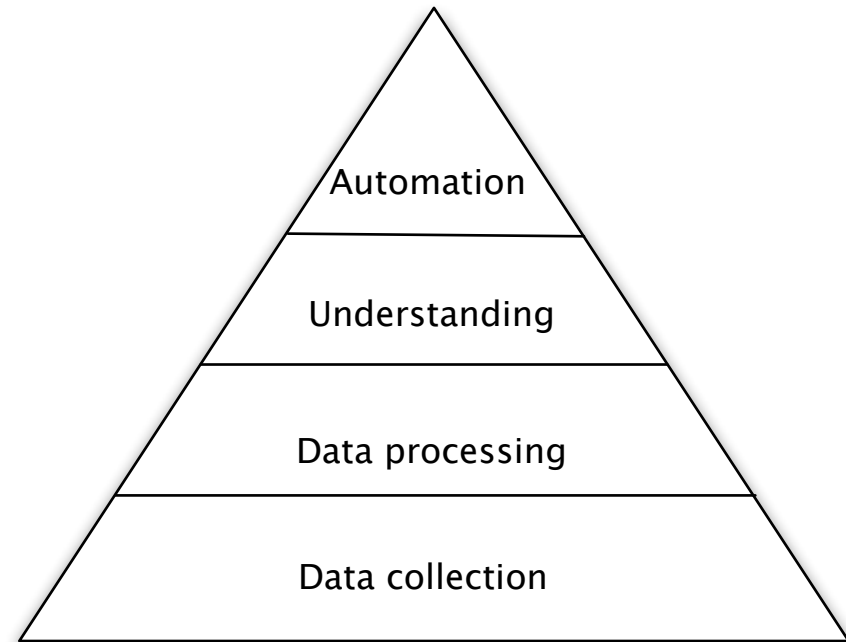
Agenda

- Real-time Data Integration
- Introduction to Logs & Apache Kafka
- Logs & Stream processing
- Apache Samza
- Stateful stream processing

The Data Needs Pyramid



Maslow's hierarchy of needs

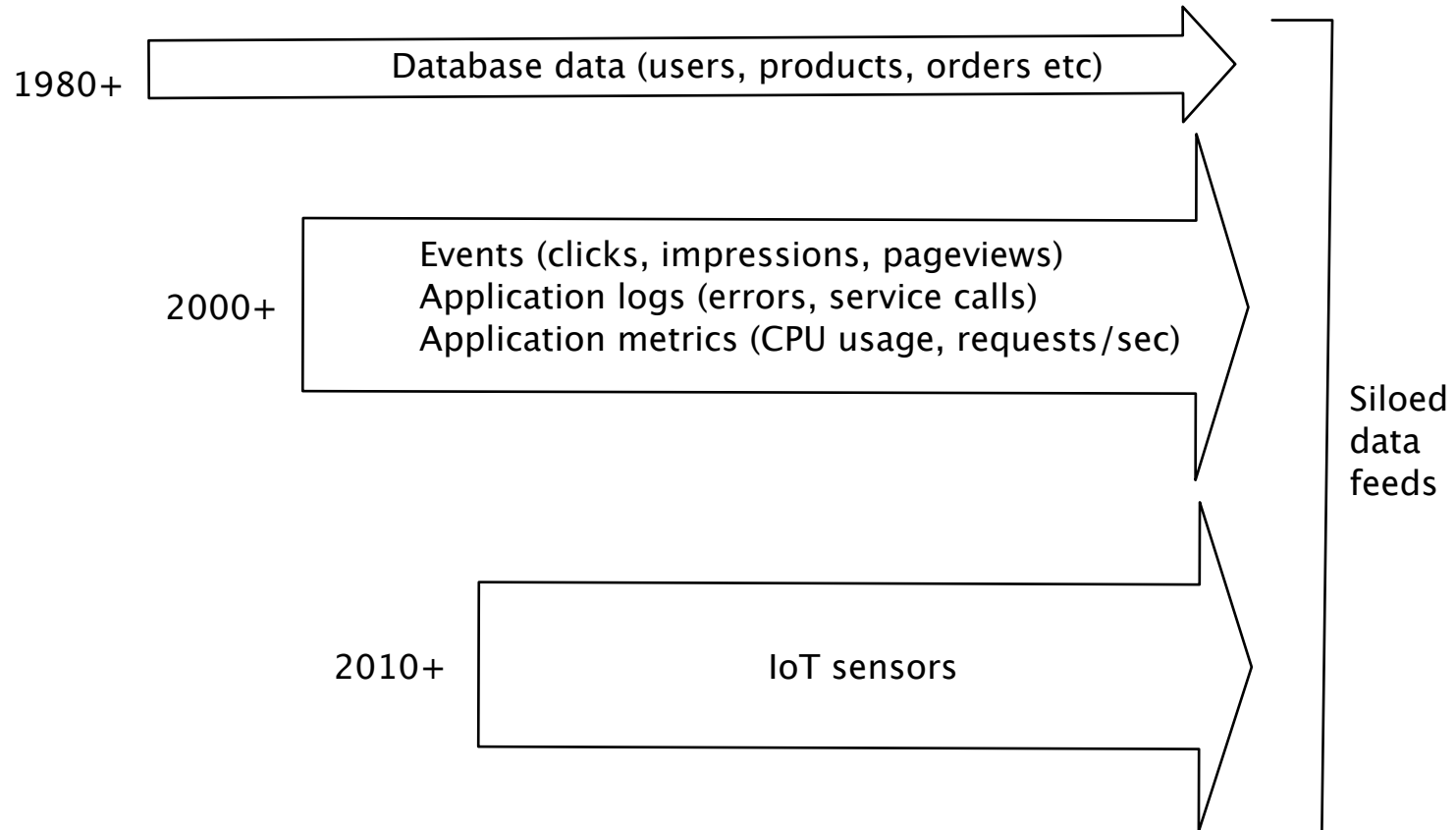


Data needs

Agenda

- **Real-time Data Integration**
- Introduction to Logs & Apache Kafka
- Logs & Stream processing
- Apache Samza
- Stateful stream processing

Increase in diversity of data



Explosion in diversity of systems

- Live Systems

- Voldemort

- Espresso

- GraphDB

- Search

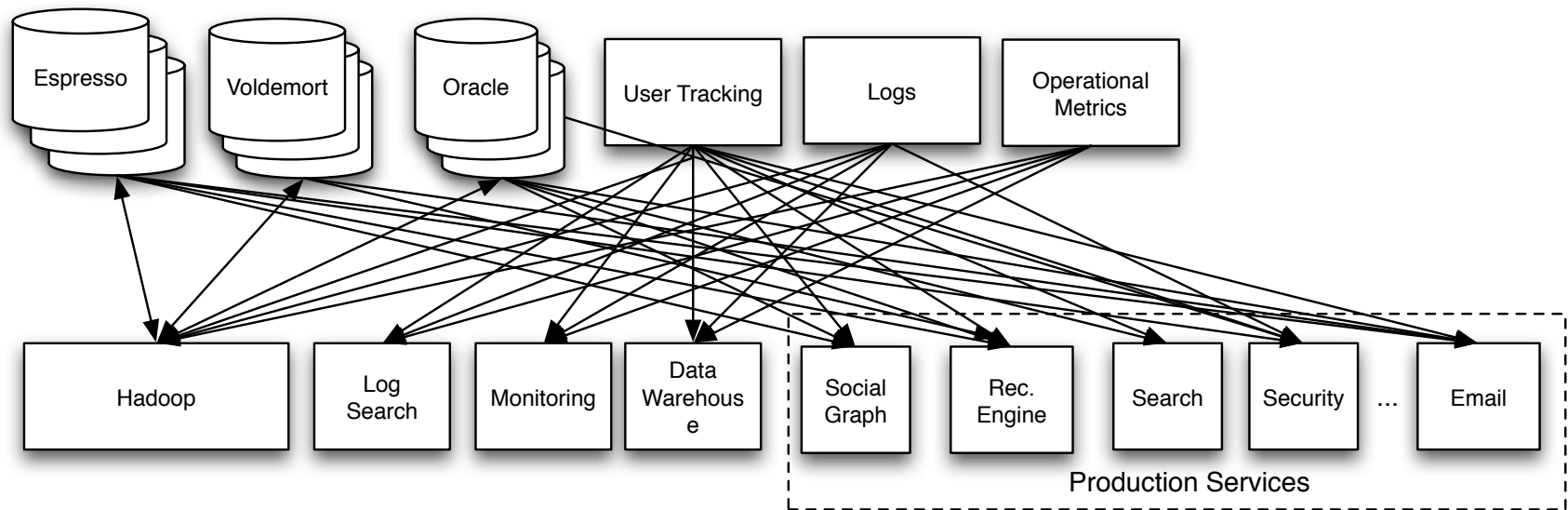
- Samza

- Batch

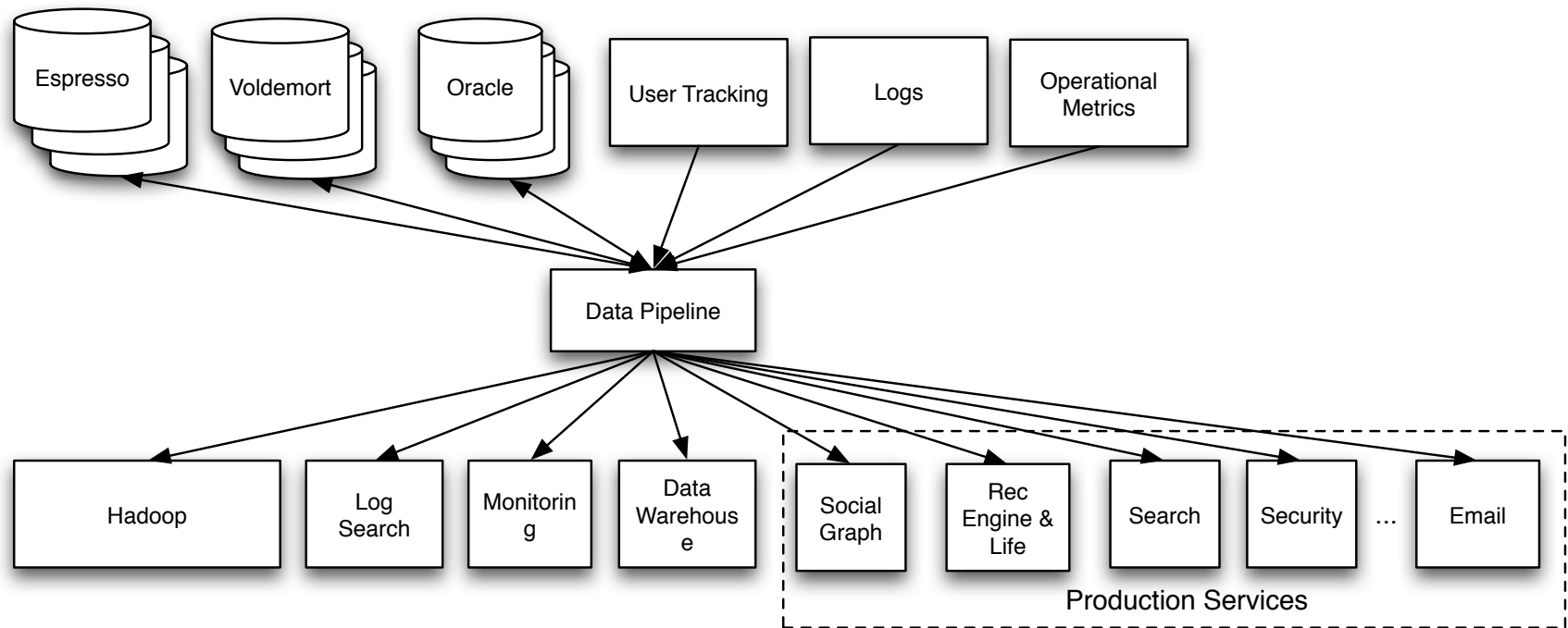
- Hadoop

- Teradata

Data integration disaster



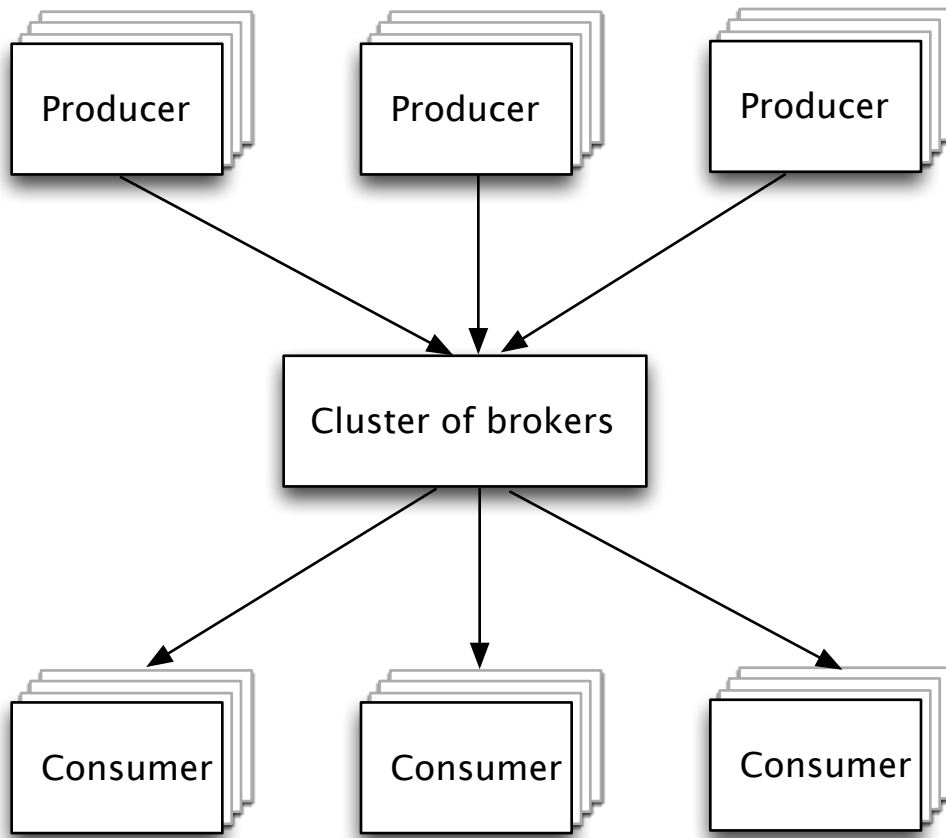
Centralized service



Agenda

- Real-time Data Integration
- **Introduction to Logs & Apache Kafka**
- Logs & Stream processing
- Apache Samza
- Stateful stream processing

Kafka at 10,000 ft



- Distributed from ground up
- Persistent
- Multi-subscriber

Key design principles

- Scalability of a file system
 - ▣ Hundreds of MB/sec/server throughput
 - ▣ Many TBs per server
- Guarantees of a database
 - ▣ Messages strictly ordered
 - ▣ All data persistent
- Distributed by default
 - ▣ Replication model
 - ▣ Partitioning model

Kafka adoption

intuit

Bloomberg

airbnb



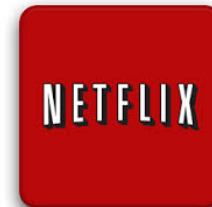
Pinterest



Microsoft



Adobe



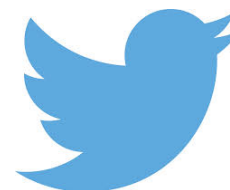
ebay

PayPal



box

Yandex



Square



Apache Kafka @ LinkedIn

- 175 TB of in-flight log data per colo
- Low-latency: ~1.5ms
- Replicated to each datacenter
- Tens of thousands of data producers
- Thousands of consumers
- 7 million messages written/sec
- 35 million messages read/sec
- Hadoop integration

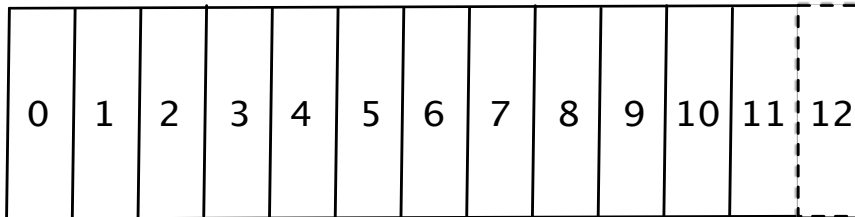
Logs

The data structure every systems engineer should know

The Log

1st record

next record written



- Ordered
- Append only
- Immutable

The Log: Partitioning

Partition 0

0	1	2	3	4	5	6	7	8	9	10	11	12
---	---	---	---	---	---	---	---	---	---	----	----	----

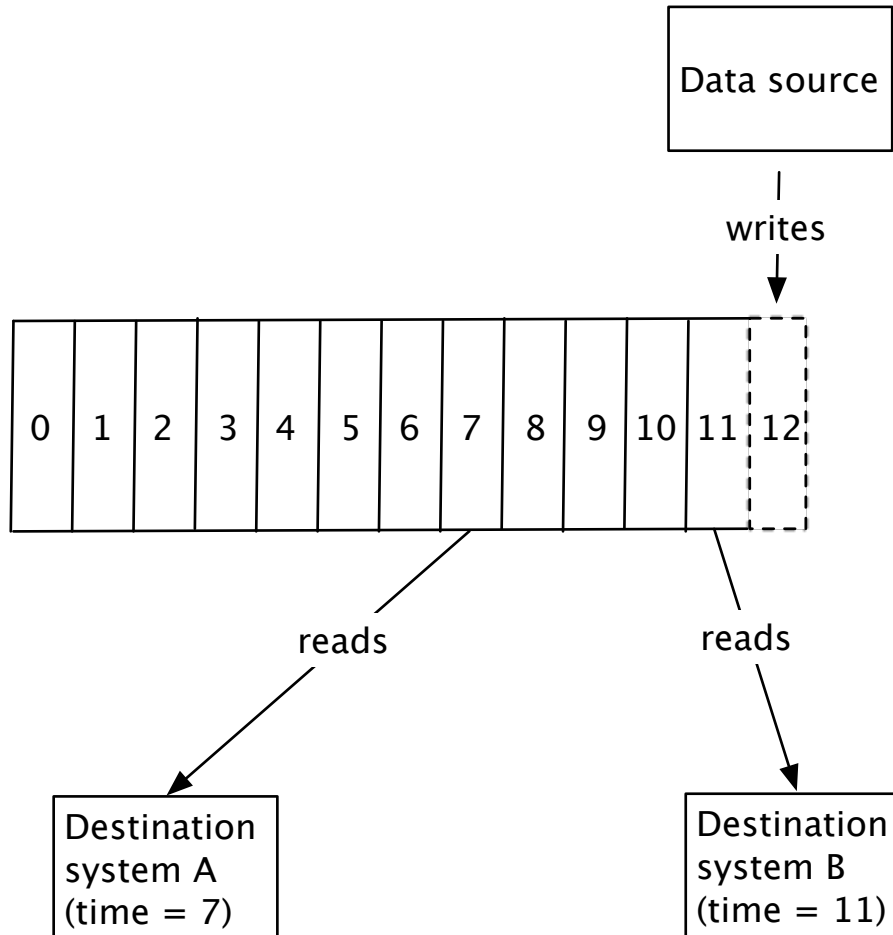
Partition 1

0	1	2	3	4	5	6	7	8	9
---	---	---	---	---	---	---	---	---	---

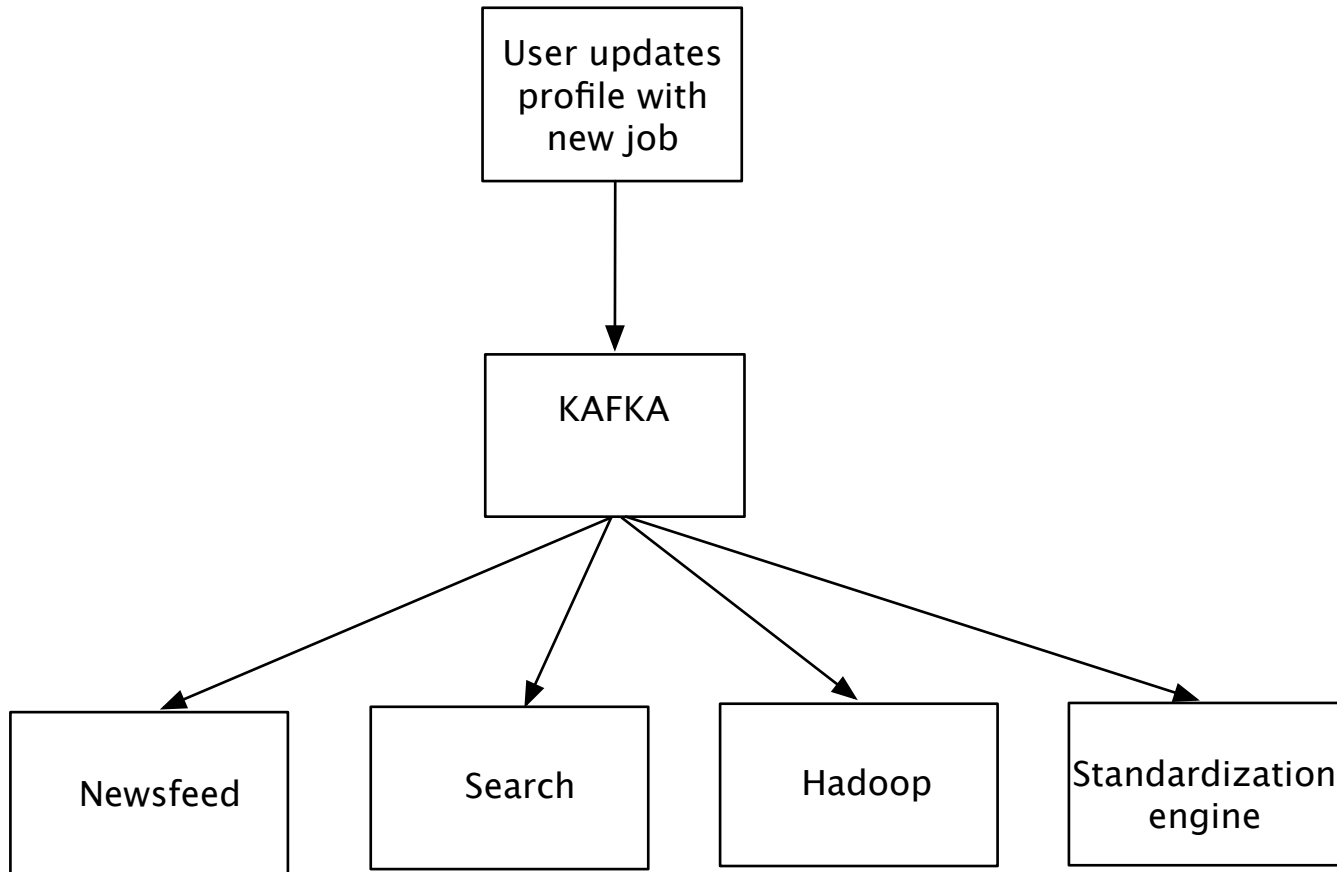
Partition 2

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
---	---	---	---	---	---	---	---	---	---	----	----	----	----	----	----	----

Logs: pub/sub done right



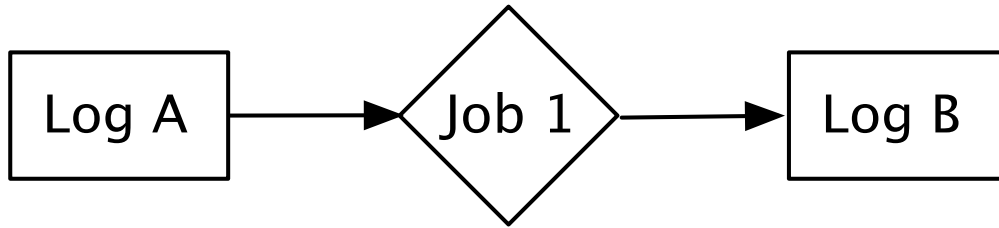
Logs for data integration



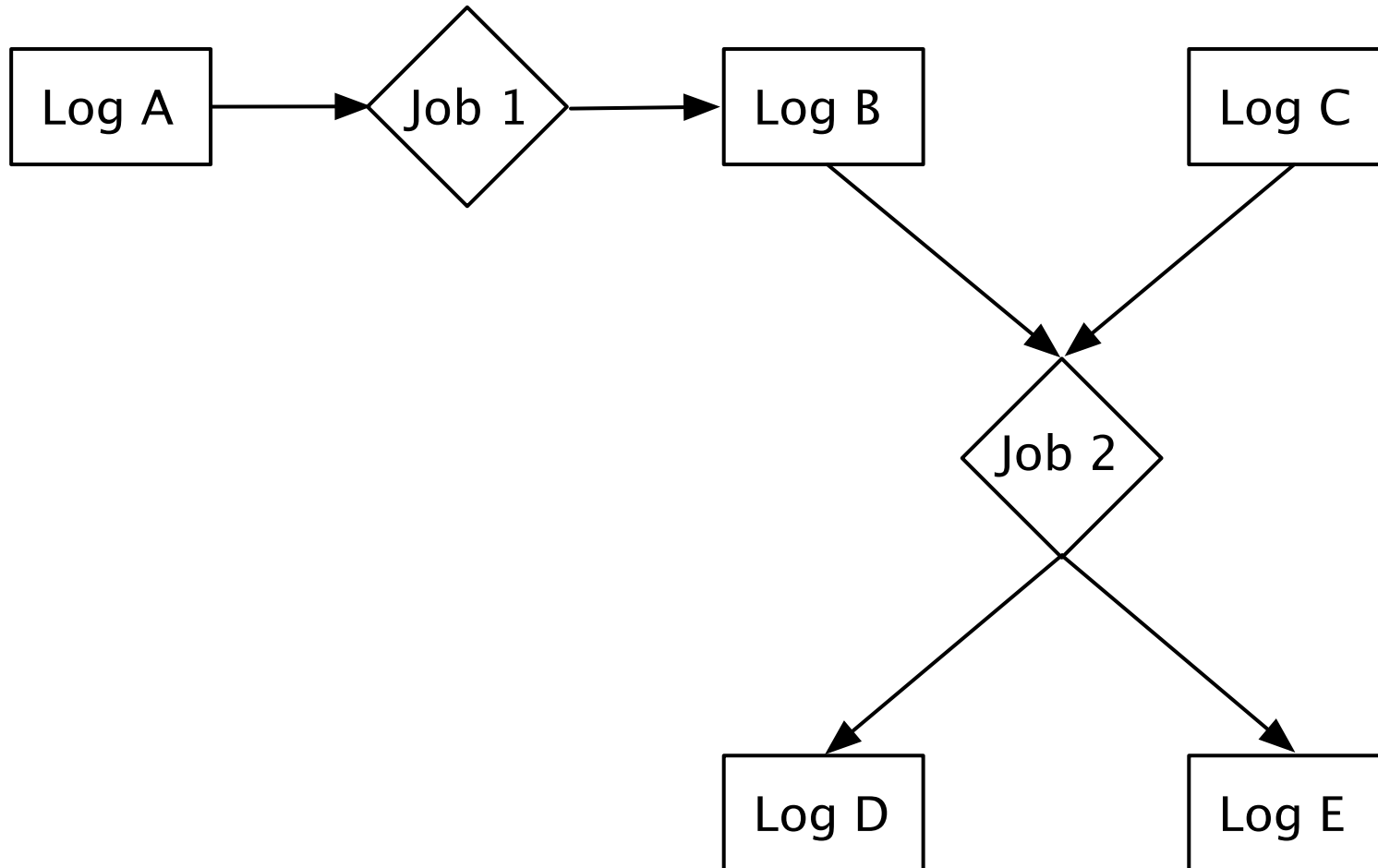
Agenda

- Real-time Data Integration
- Introduction to Logs & Apache Kafka
- **Logs & Stream processing**
- Apache Samza
- Stateful stream processing

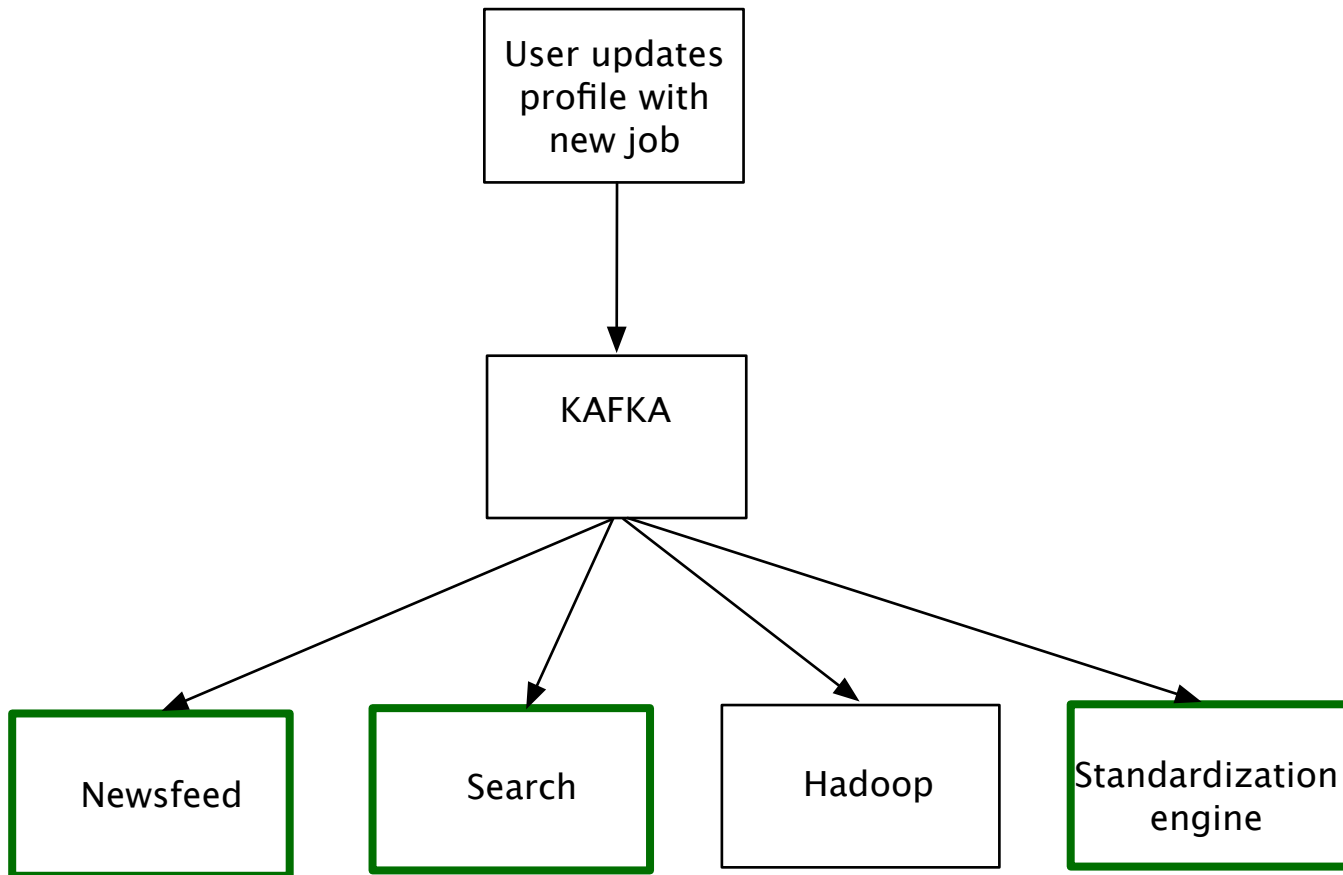
Stream processing = $f(\log)$



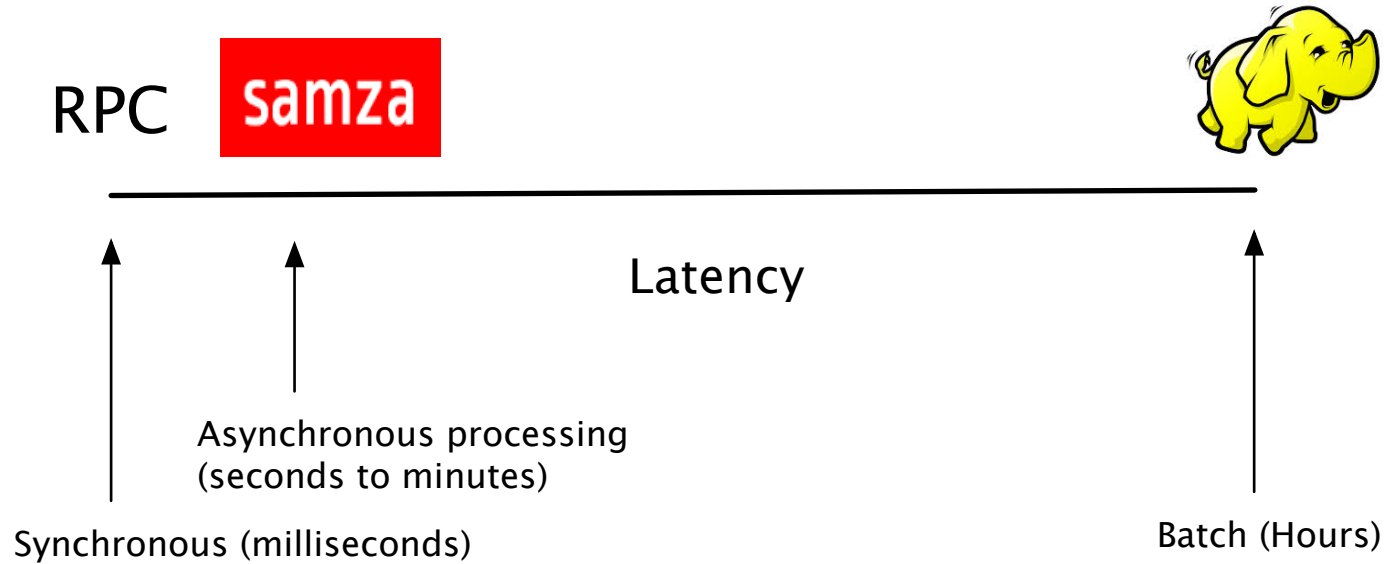
Stream processing = $f(\log)$



Apache Samza at LinkedIn



Latency spectrum of data systems



Agenda

- Real-time Data Integration
- Introduction to Logs & Apache Kafka
- Logs & Stream processing
- **Apache Samza**
- Stateful stream processing

Samza API

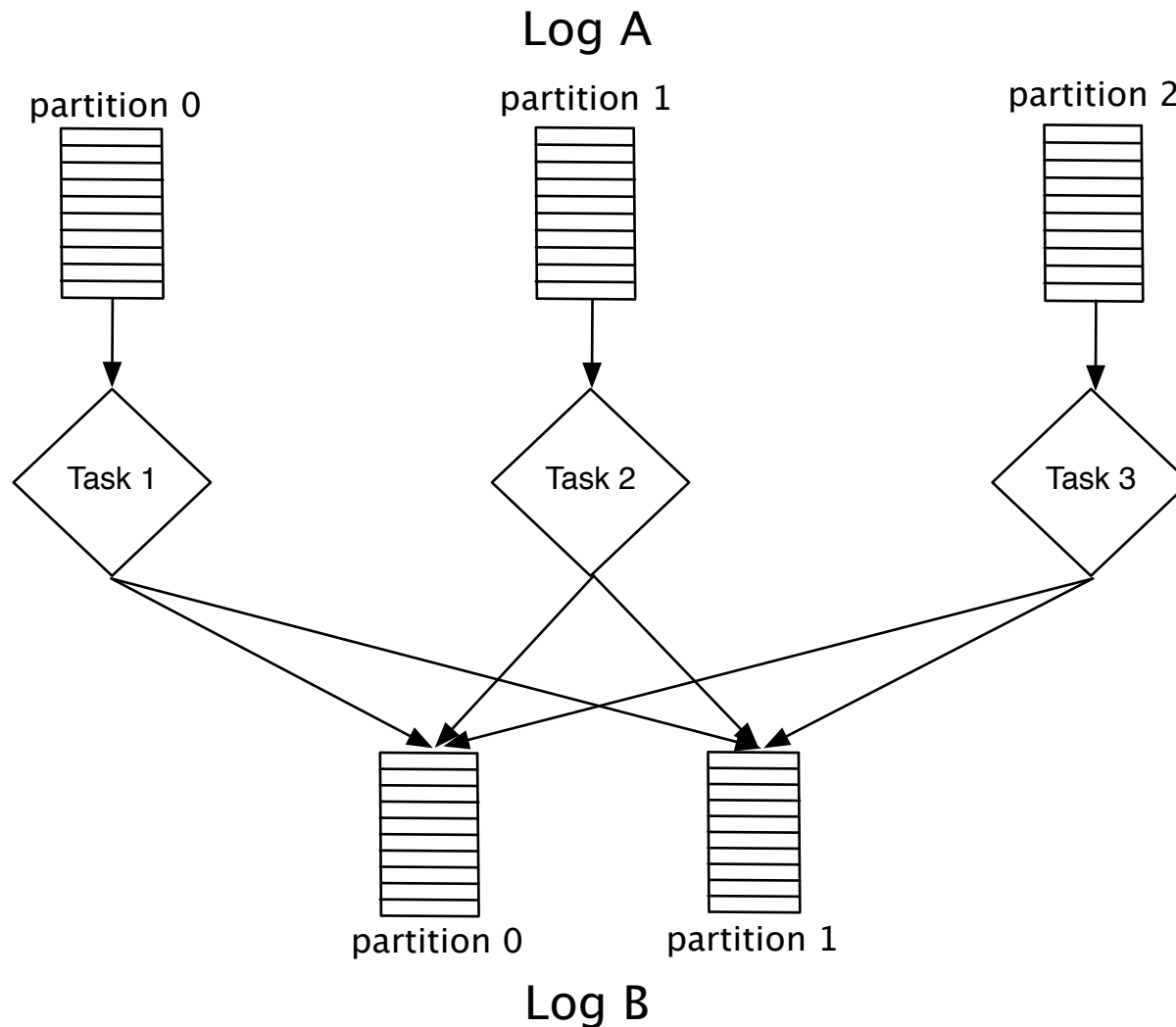
```
public interface StreamTask {  
    void process (IncomingMessageEnvelope envelope,  
                 MessageCollector collector,  
                 TaskCoordinator coordinator);  
}
```

getKey(), getMsg()

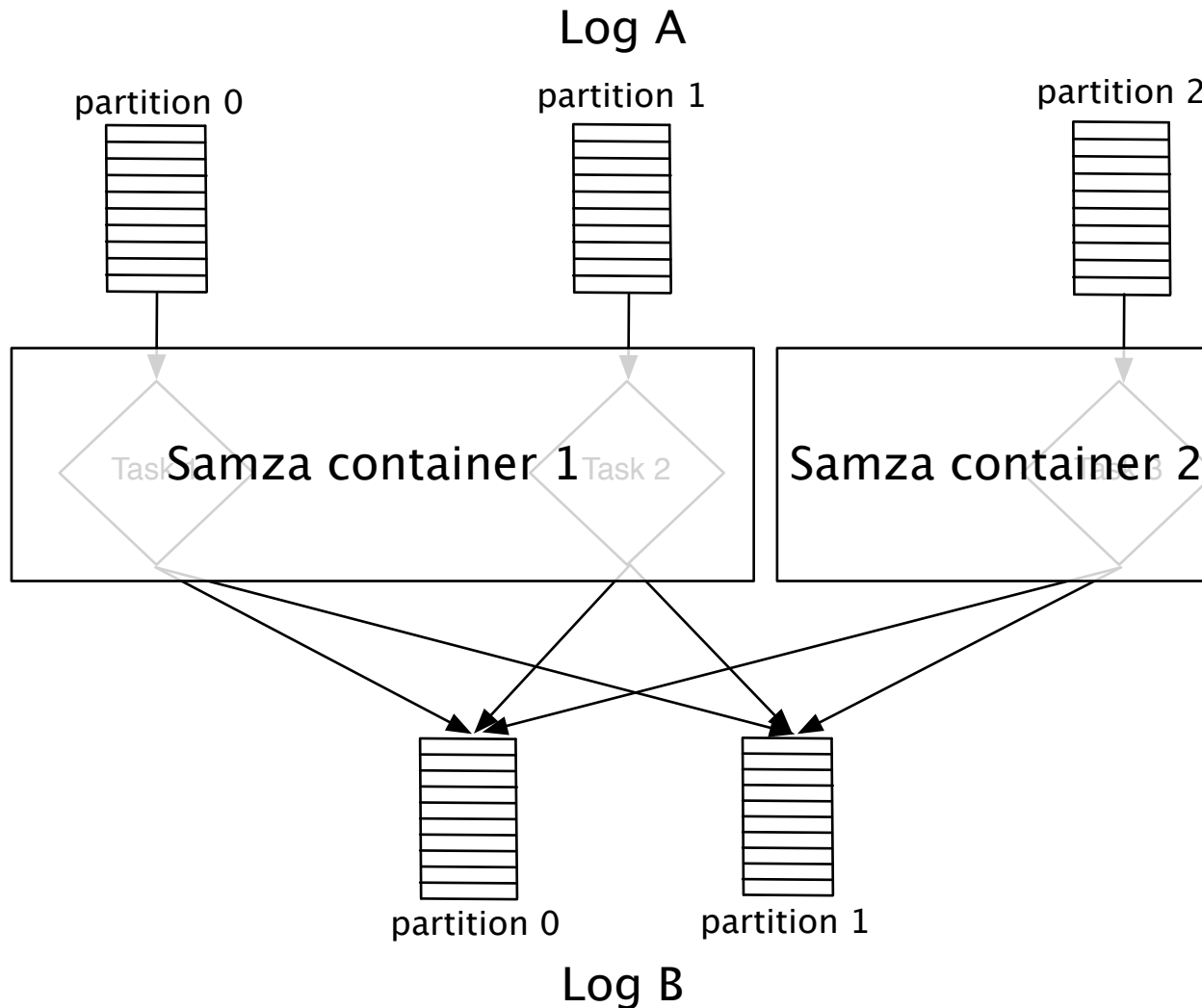
commit(), shutdown()

sendMsg(topic, key, value)

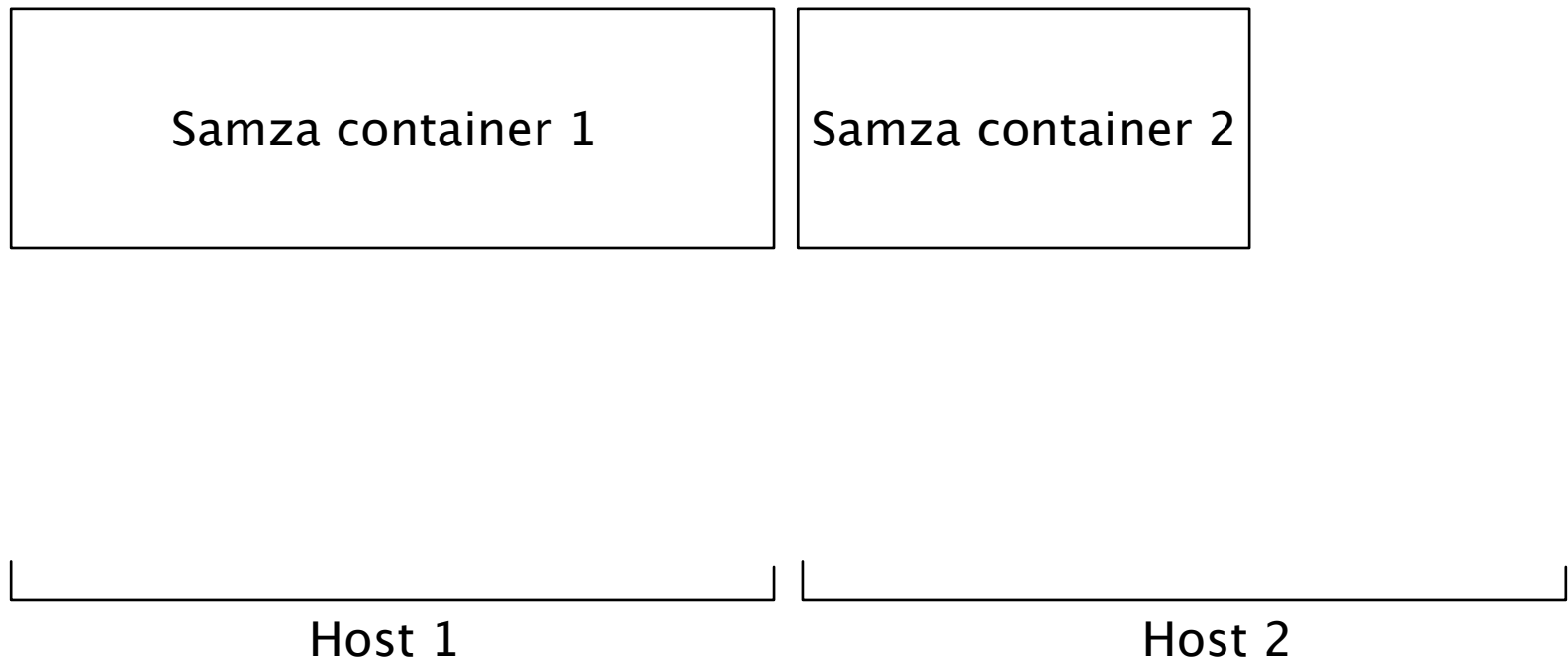
Samza Architecture (Logical view)



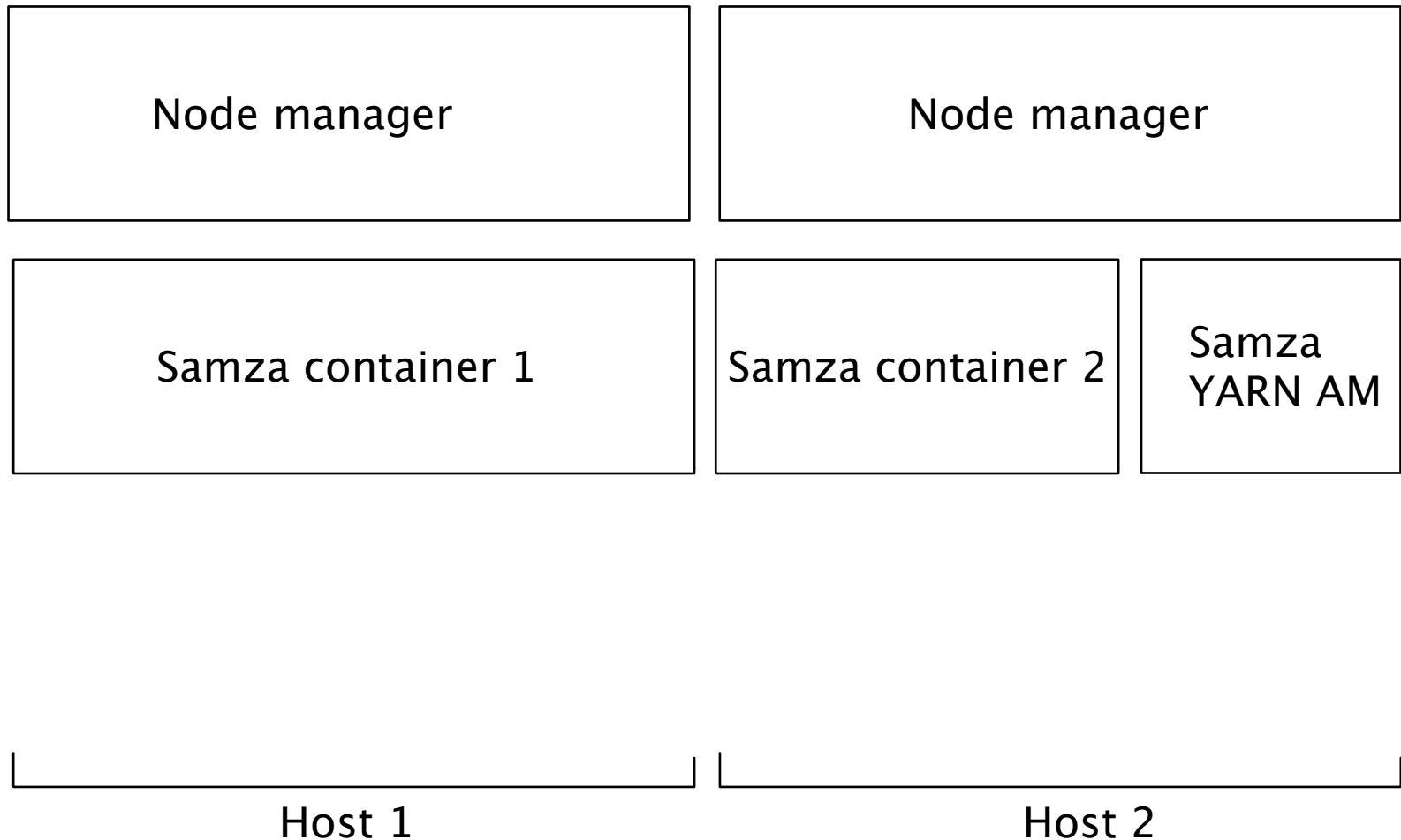
Samza Architecture (Logical view)



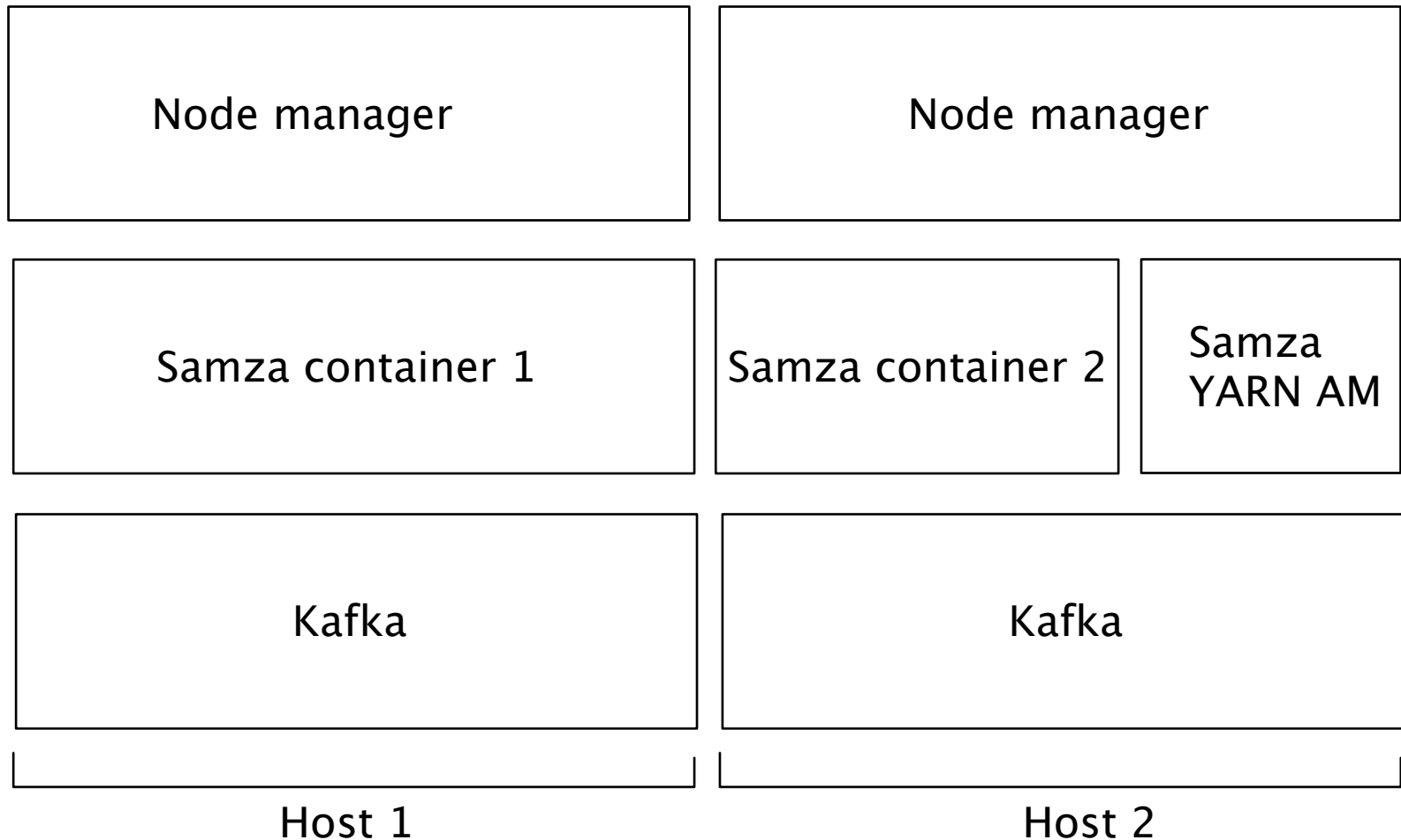
Samza Architecture (Physical view)



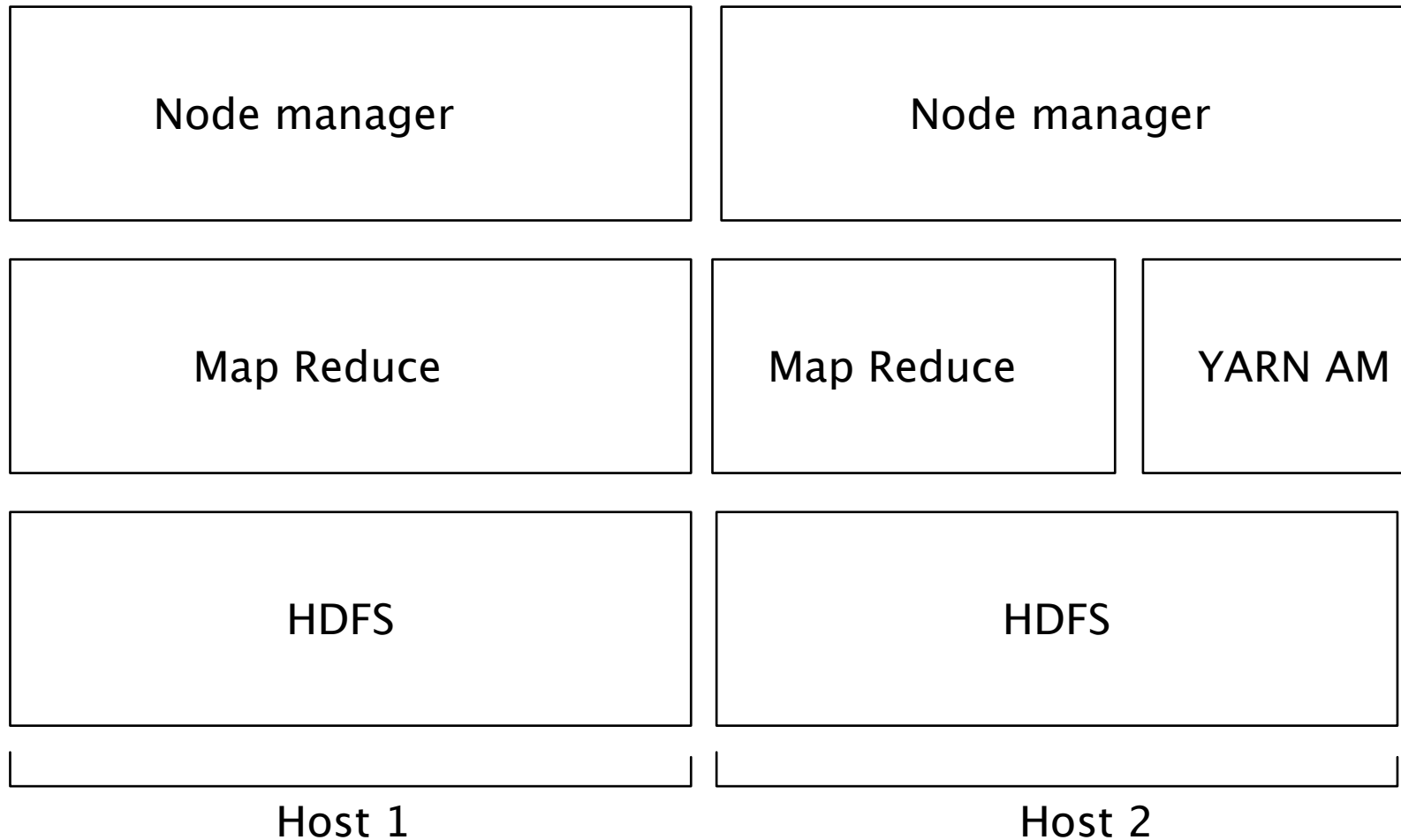
Samza Architecture (Physical view)



Samza Architecture (Physical view)



Samza Architecture: Equivalence to Map Reduce



M/R Operation Primitives

- Filter records matching some condition
- Map $\text{record} = f(\text{record})$
- Join Two/more datasets by key
- Group records with same key
- Aggregate $f(\text{records within the same group})$
- Pipe job 1's output \Rightarrow job 2's input

M/R Operation Primitives on streams

- Filter records matching some condition
- Map $\text{record} = f(\text{record})$
- Join Two/more datasets by key
- Group records with same key
- Aggregate $f(\text{records within the same group})$
- Pipe job 1's output \Rightarrow job 2's input

Requires state
maintenance

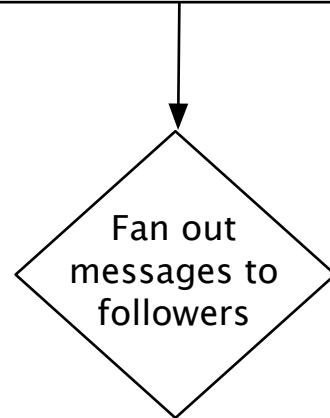


Agenda

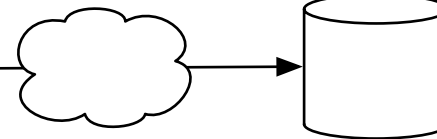
- Real-time Data Integration
- Introduction to Logs & Apache Kafka
- Logs & Stream processing
- Apache Samza
- **Stateful stream processing**

Example: Newsfeed

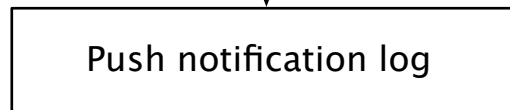
User ... posted "..."
User 989 posted "Blah Blah"
User 567 posted "Hello World"



External connection DB

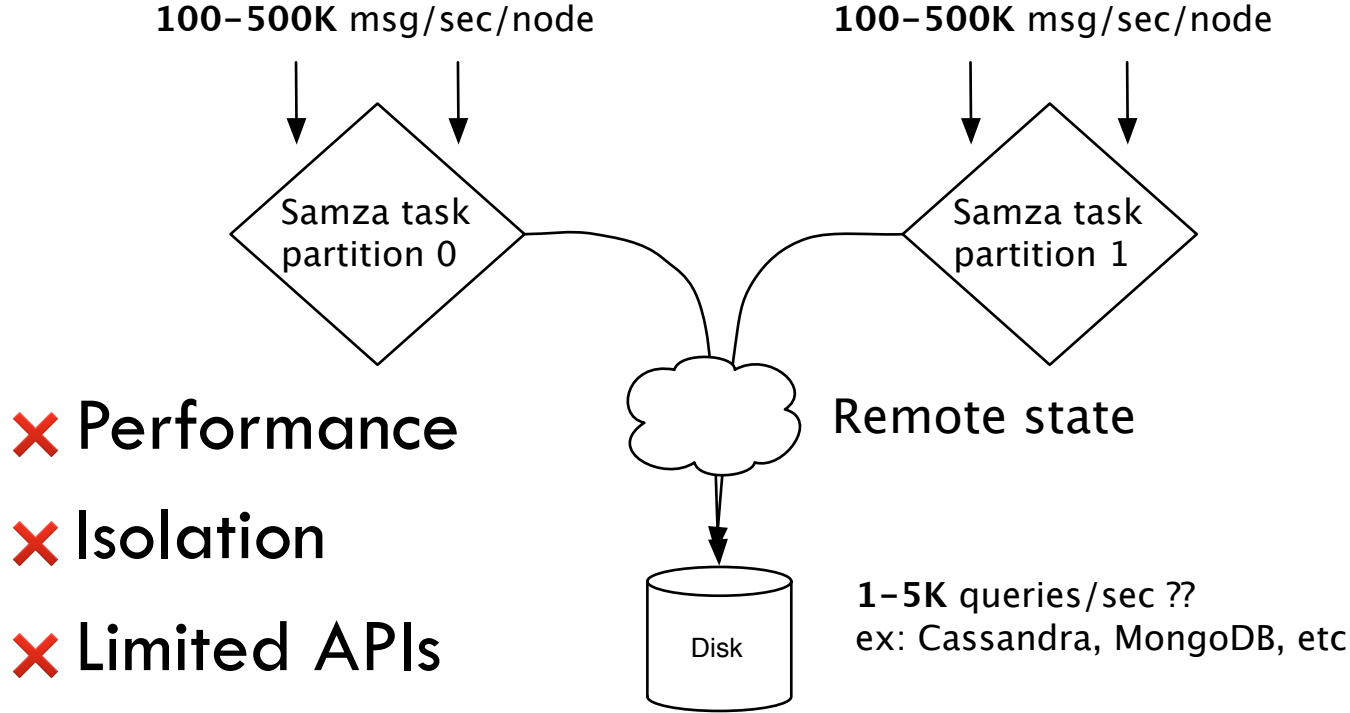


567 -> [123, 679, 789, ...]
999 -> [156, 343, ...]

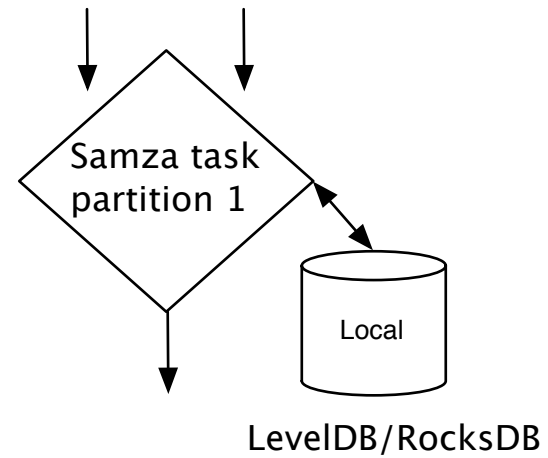
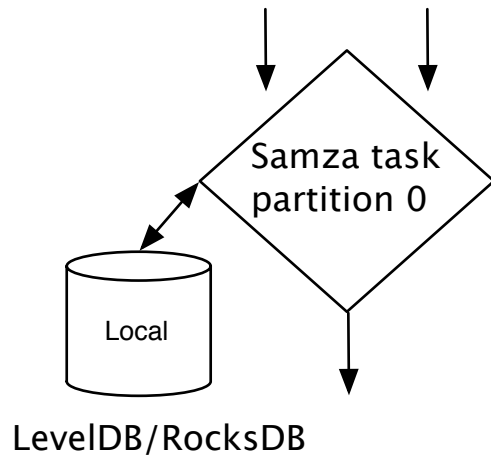


Refresh user 123's newsfeed
Refresh user 679's newsfeed
Refresh user ...'s newsfeed

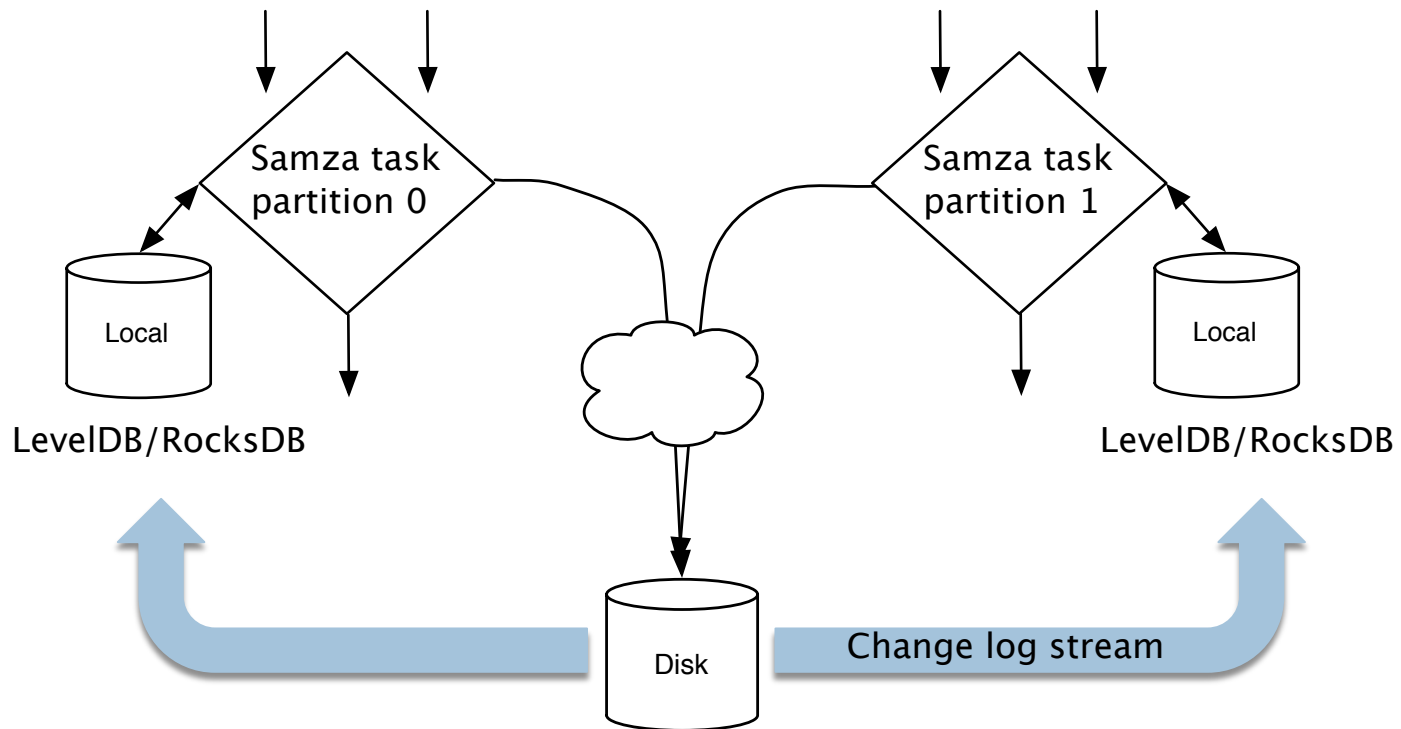
Local state vs Remote state: Remote



Local state: Bring data closer to computation



Local state: Bring data closer to computation



Example Revisited: Newsfeed

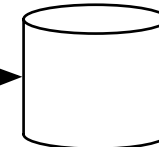
User ... posted "..."
User 989 posted "Blah Blah"
User 567 posted "Hello World"

Status update log

User ... followed ...
User 123 followed 567
User 890 followed 234

New connection log

Fan out
messages to
followers

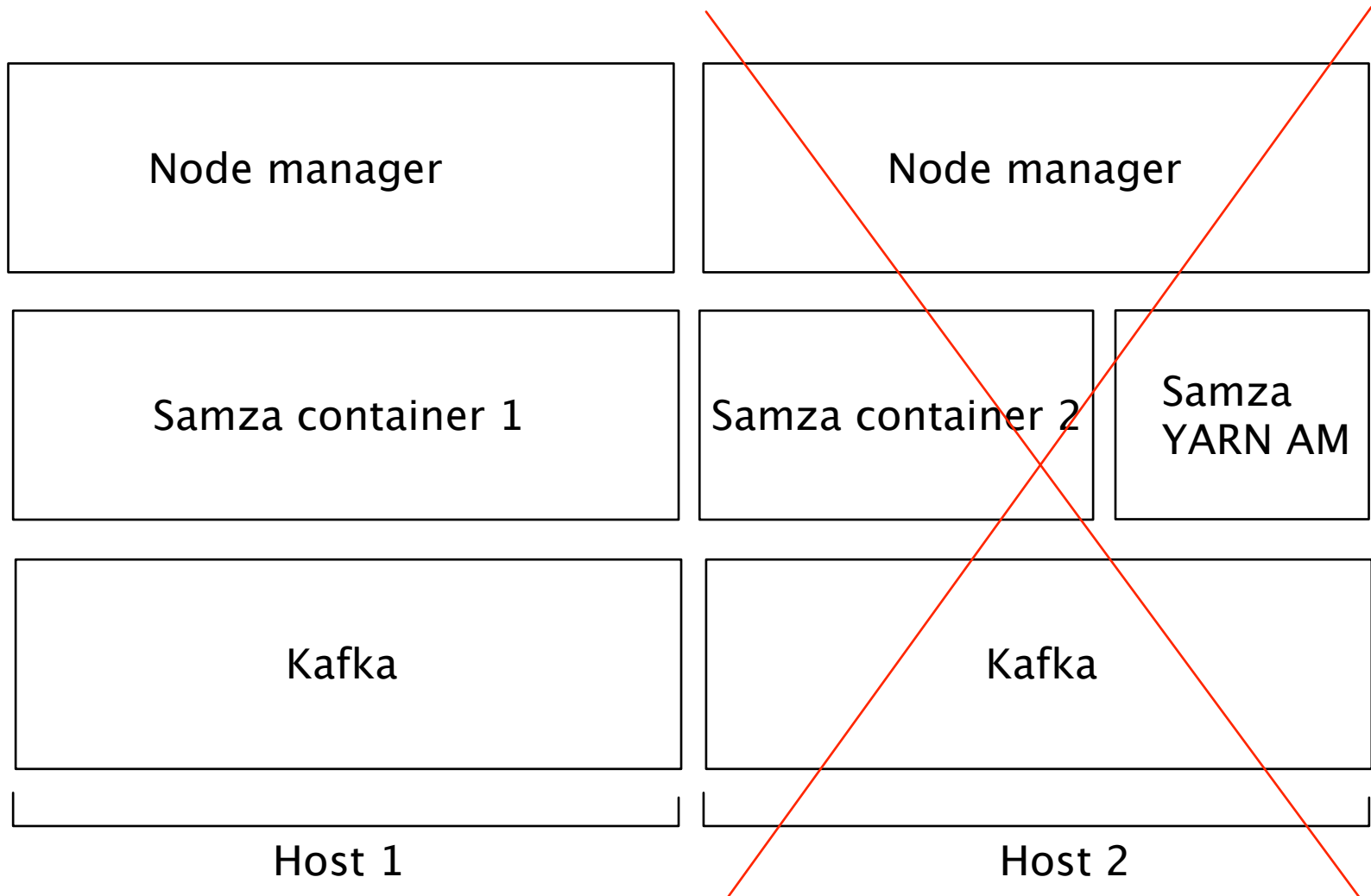


567 -> [123, 679, 789, ...]
999 -> [156, 343, ...]

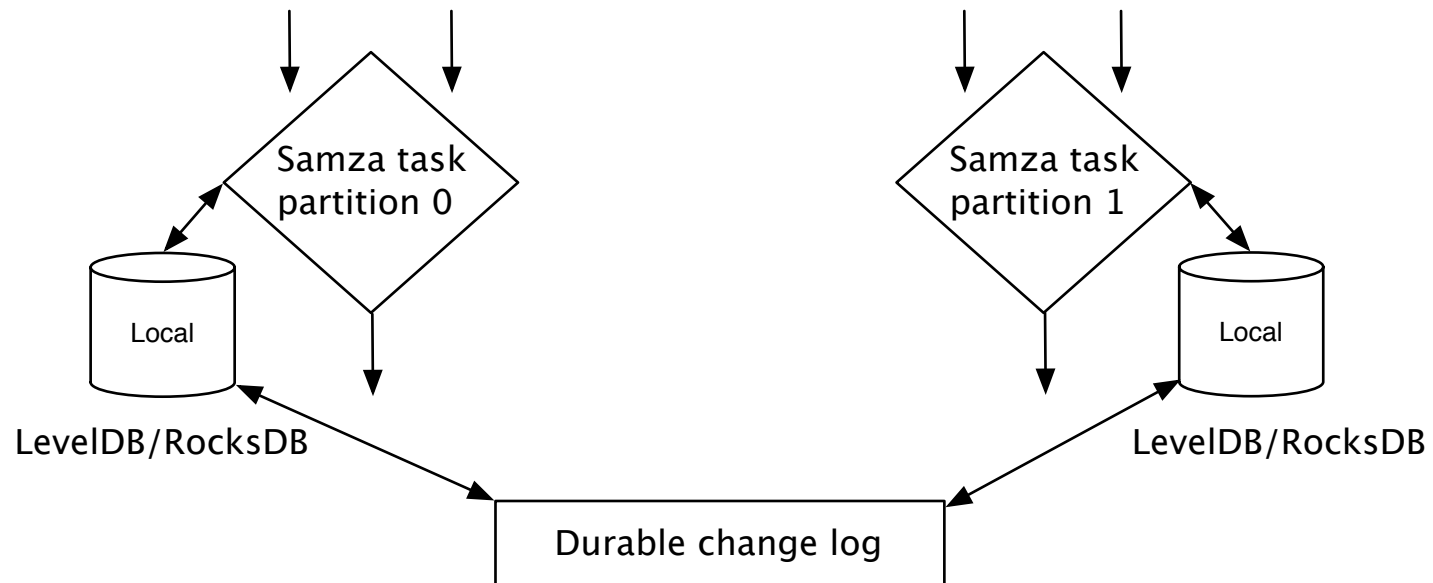
Push notification log

Refresh user 123's newsfeed
Refresh user 679's newsfeed
Refresh user ...'s newsfeed

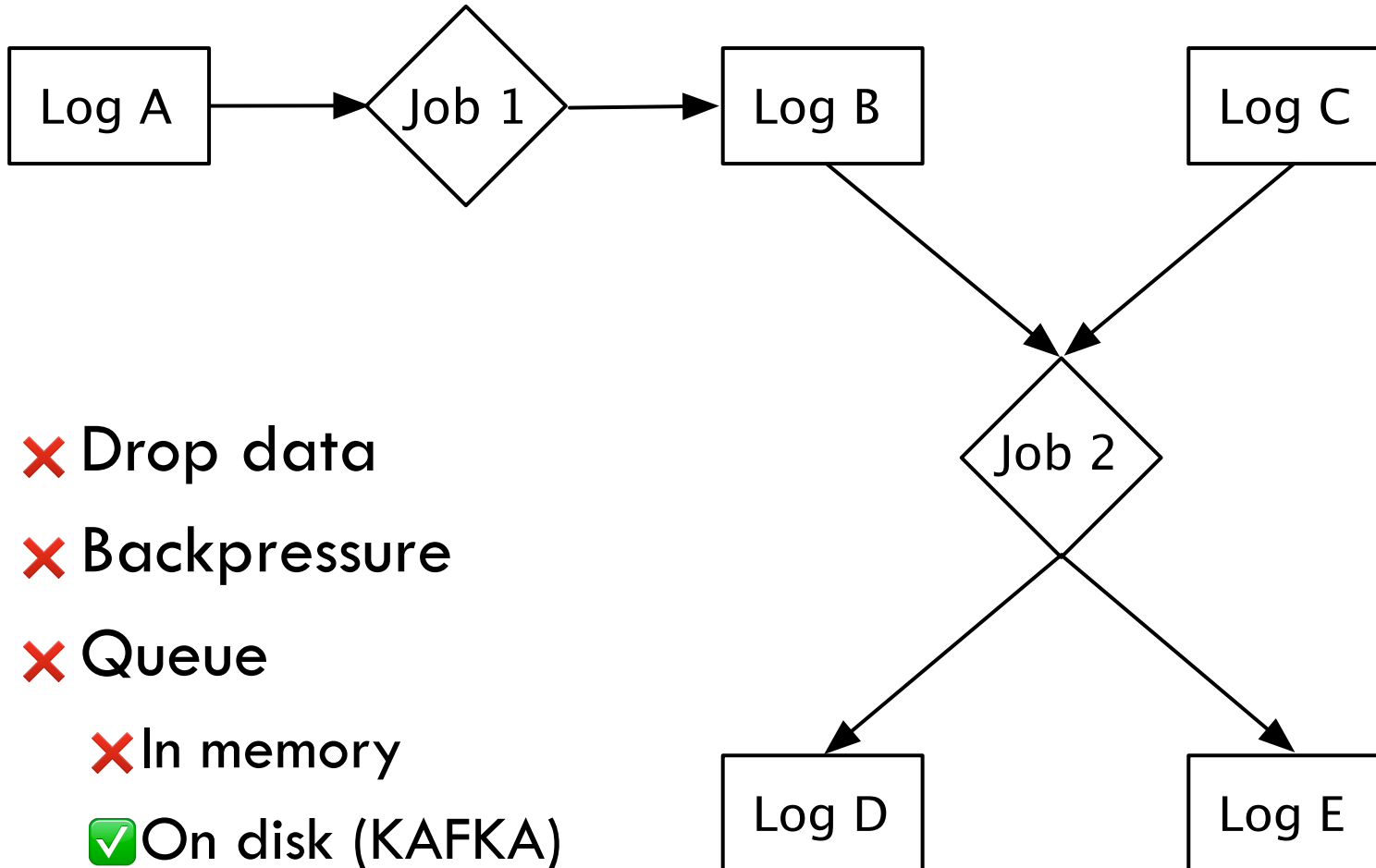
Fault tolerance?



Fault tolerance in Samza



Slow jobs



Summary

- Real time data integration is crucial for the success and adoption of stream processing
- Logs form the basis for real time data integration
- Stream processing = $f(\text{logs})$
- Samza is designed from ground-up for scalability and provides fault-tolerant, persistent state

Thank you!

- The Log
 - http://bit.ly/the_log
- Apache Kafka
 - <http://kafka.apache.org>
- Apache Samza
 - <http://samza.incubator.apache.org>
- Me
 - @nehanarkhede
 - <http://www.linkedin.com/in/nehanarkhede>