



POLITECNICO
MILANO 1863

Kafka

Scaling continuous ingestion

Emanuele Della Valle

emanuele.dellavalle@polimi.it



POLITECNICO
MILANO 1863

Position of this class in the course



Taming
Continuous numerous flows that
can turn into a torrent
with
Complex Event Processing

Covered in the lectures on EPL and Esper

The background image shows a waterfall cascading down a steep, mossy cliff into a pool of water at the bottom. The scene is lush and green.

Taming
*Myriads of tiny flows
that you can collect*
with
Event-based Systems

The Apache Kafka logo, which consists of three interconnected circles forming a triangular shape.

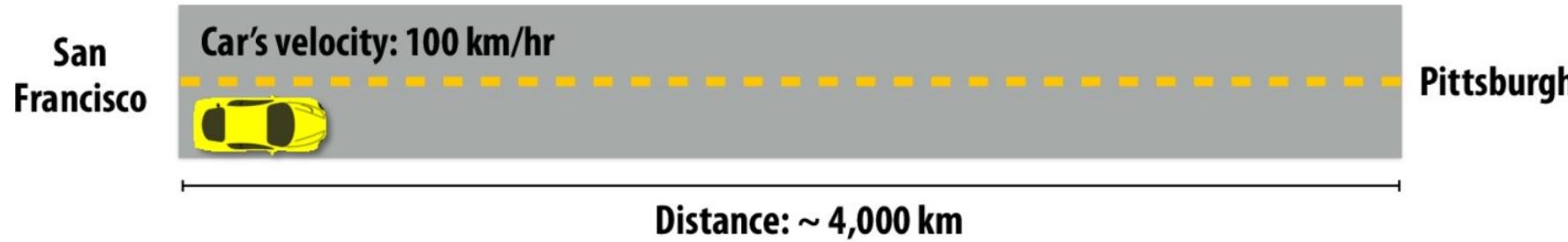
kafka

Throughput Latency Data/Message

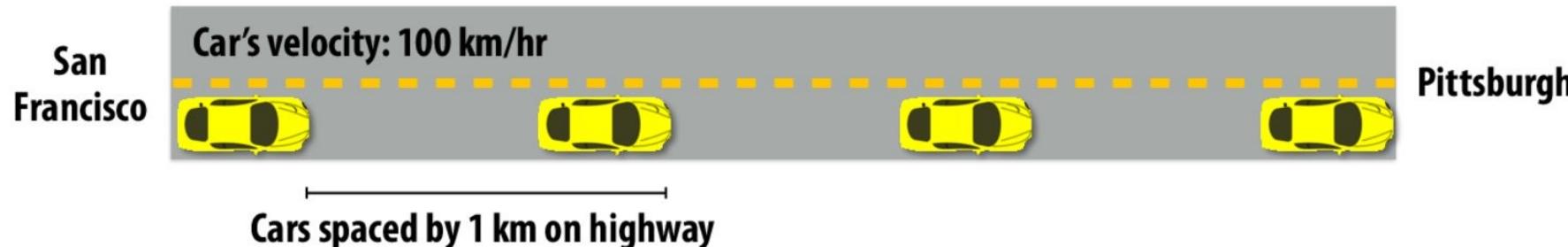
Few important concepts

Everyone wants to get to Pittsburgh!

(Latency vs. throughput review)

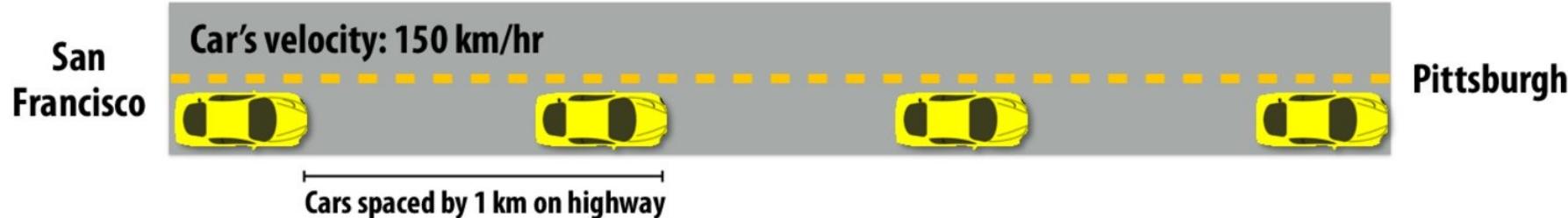


Latency of moving a person from San Francisco to Pittsburgh: 40 hours



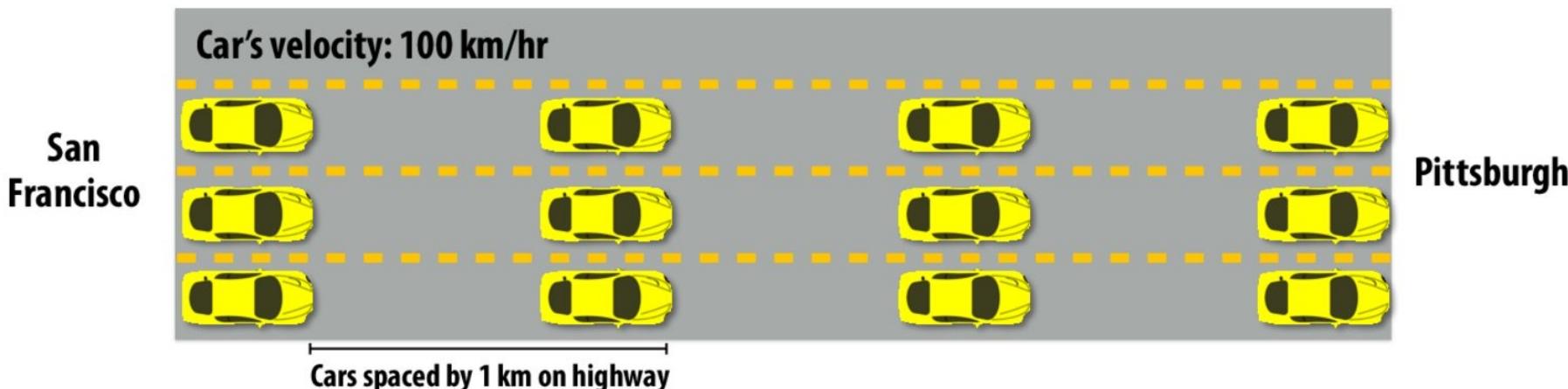
Throughput: 100 people per hour (1 car every 1/100 of an hour)

Improving throughput



Approach 1: drive faster!

Throughput = 150 people per hour (1 car every 1/150 of an hour)



Approach 2: build more lanes!

Throughput: 300 people per hour (3 cars every 1/100 of an hour)

Review: latency vs throughput

Latency

The amount of time needed for an operation to complete.

A memory load that misses the cache has a latency of 200 cycles

A packet takes 20 ms to be sent from my computer to Google

Asking a question on Piazza gets response in 10 minutes

throughput

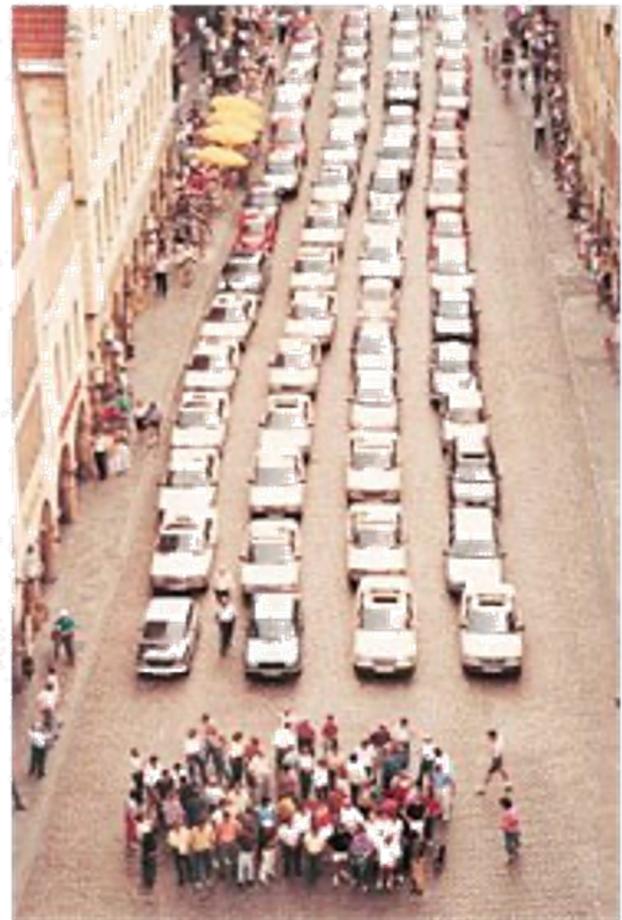
The rate at which operations are performed.

Memory can provide data to the processor at 25 GB/sec.

A communication link can send 10 million messages per second

The TAs answer 50 questions per day on Piazza

space required to transport 60 people



car



bus



bicycle

(Poster in city of Muenster Planning Office, August 2001) Credit: PressOffice City of Munster, Germany

@manudellavalle - <http://emanueledellavalle.org>



POLITECNICO
MILANO 1863



Data/Message matters!

- Fixed the size of a message (i.e., the space between two cars)
- The more people fit in the car, the higher the throughput
- The larger the data points per message, the higher the throughput
- **Compression** and **binary encoding** are the typical ways to increase data per message



Taming
*Myriads of tiny flows
that you can collect*
with
Event-based Systems



oo kafka



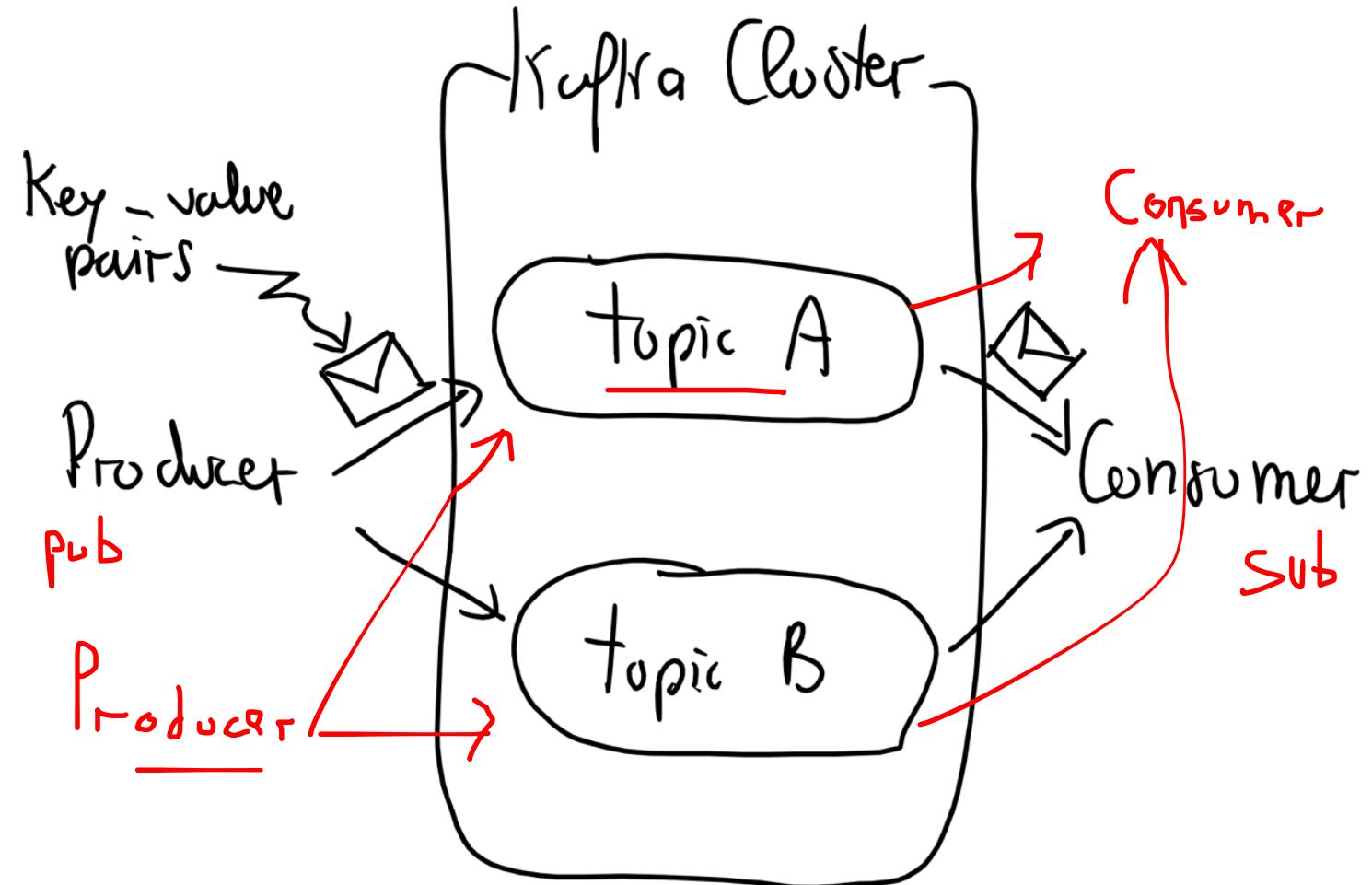
Kafka History

- In 2010, it was created at LinkedIn to
 - Simplifying data pipelines
 - Processing streaming data for batch and real-time analytics
 - In 2012, it became a top-level Apache project
 - Now, it is now at the core of LinkedIn's architecture
 - Performs extremely well at very large scale
 - Produces over 2 *trillion* messages per day
- and it is in use at many organizations
- Twitter, Netflix, Goldman Sachs, Hotels.com, IBM, Spotify, Uber, Square, Cisco...



A Conceptual View of Kafka

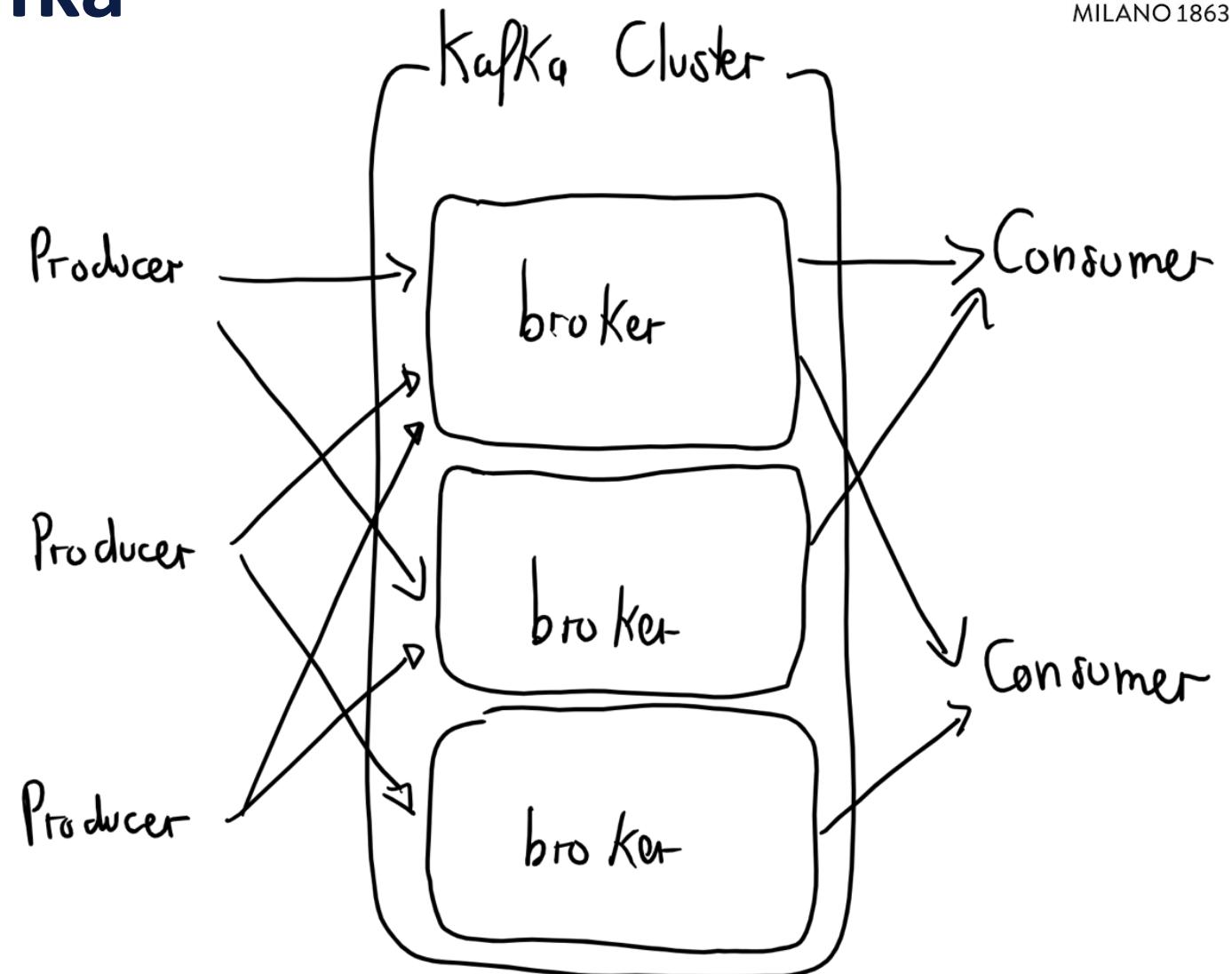
- **Producers** send messages on topics
 - Publishers in pub-sub terms
- **Consumers** read messages from topics
 - Subscribers in pub-sub terms
- **Topics** are streams of messages
- A **message** is a key-value pair + metadata
- A Kafka cluster manages topics
 - No content-based routing





A System View of Kafka

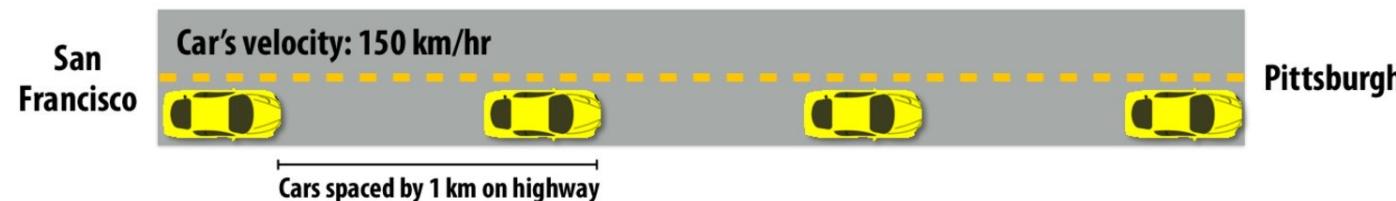
- **Brokers** are the main storage and messaging components of the Kafka cluster
- Producers publish messages on brokers
- Consumers subscribe to brokers to receive messages



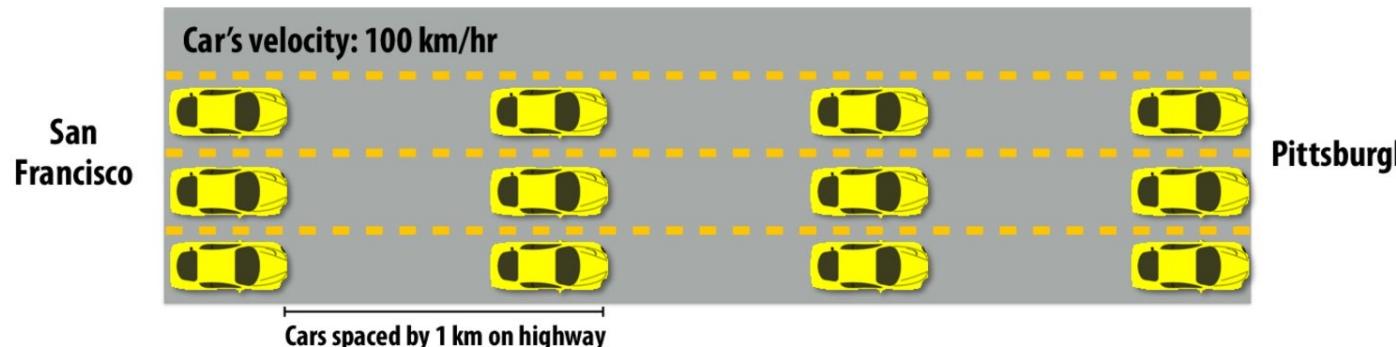


Notice the similarity

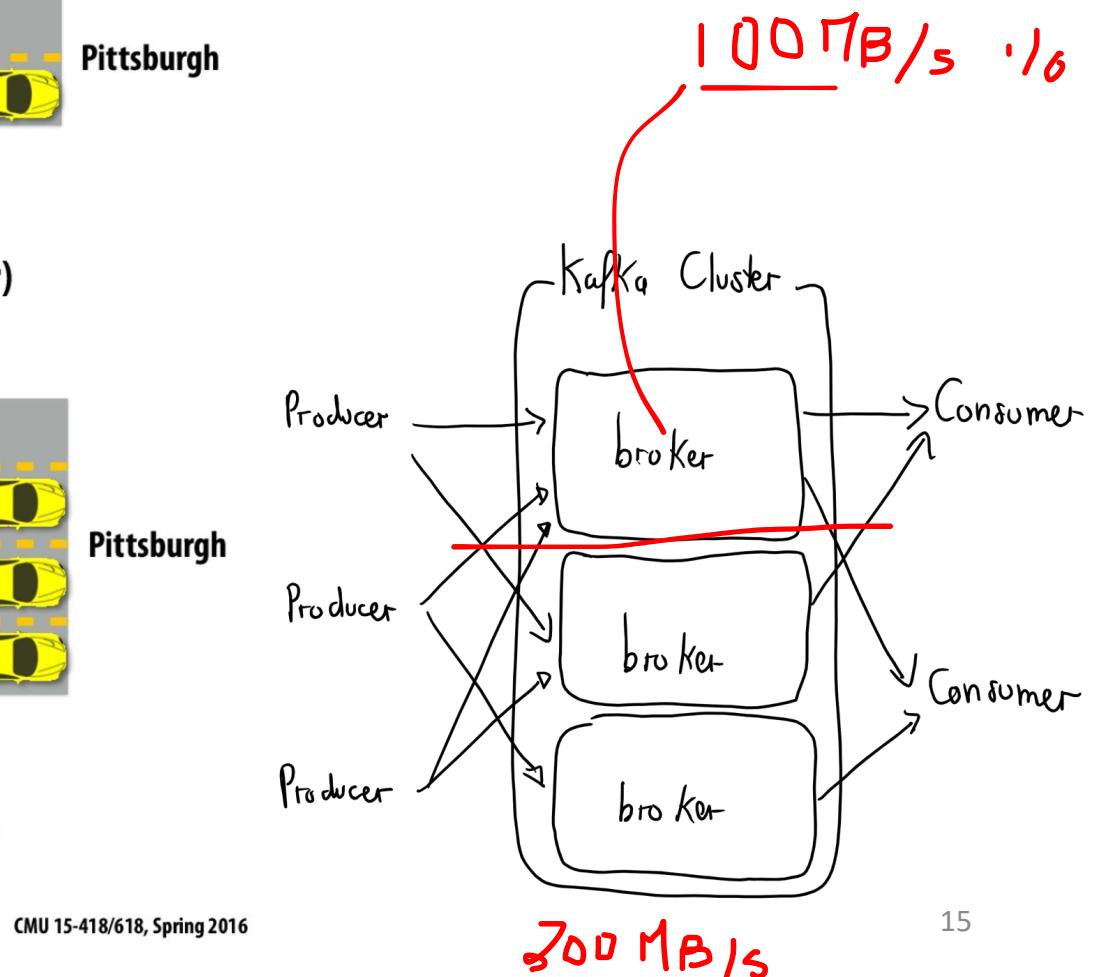
Improving throughput



Approach 1: drive faster!
Throughput = 150 people per hour (1 car every 1/150 of an hour)



Approach 2: build more lanes!
Throughput: 300 people per hour (3 cars every 1/100 of an hour)



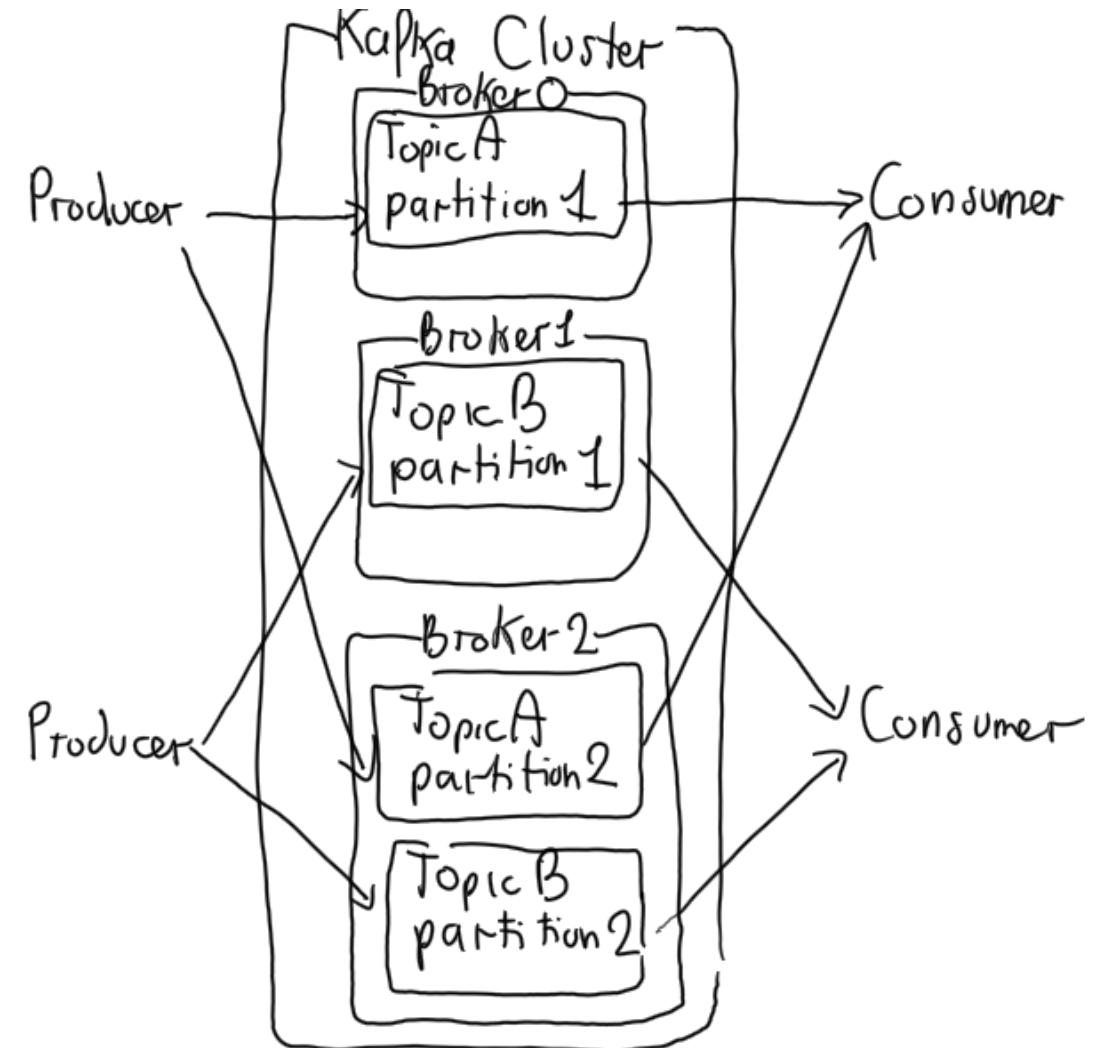


More on brokers

- Brokers receive and store messages when they are sent by the Producers
- A production Kafka cluster will have three or more Brokers
 - Each can handle hundreds of thousands, or millions, of messages per second
- Typically, a Broker manages multiple Partitions

Reconciling the two views

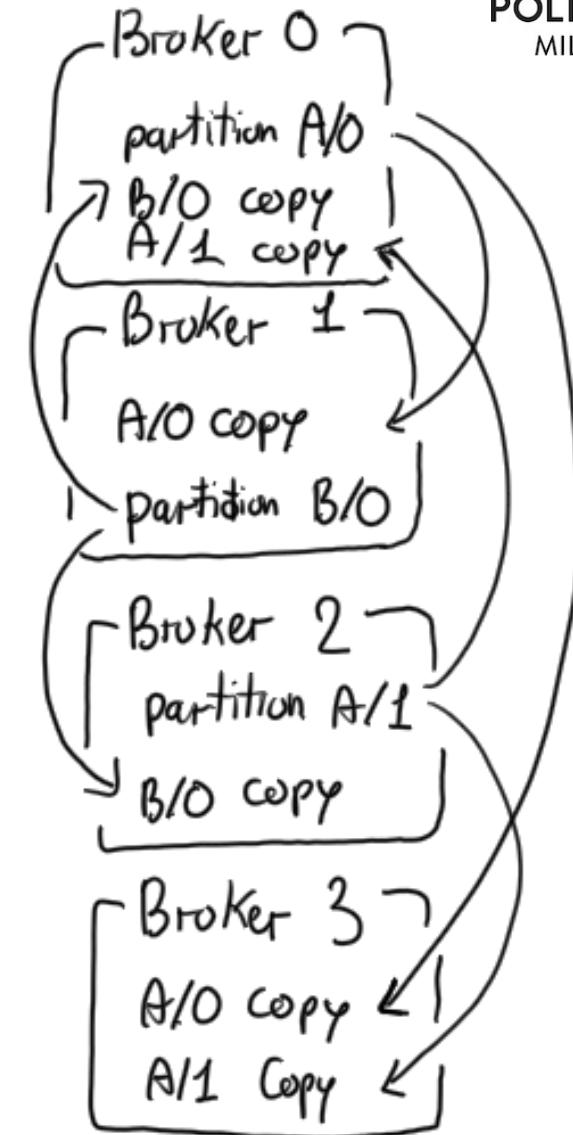
- Topics are **partitioned** across brokers
- **Producers shard messages** over the partitions of a certain topic
- Typically, hash partitioning, based on the message's key, determines the partition a message is assigned to
 - If the key is null, then round-robin partitioning is adopted





Fault Tolerance via a Replicated Log

- Kafka maintains replicas of each partition on other Brokers in the cluster
 - The number of replicas is configurable
- One Broker is the leader for that Partition
 - All writes and reads go to and from the leader
 - Other Brokers are followers
- Replication provides fault tolerance in case a Broker goes down



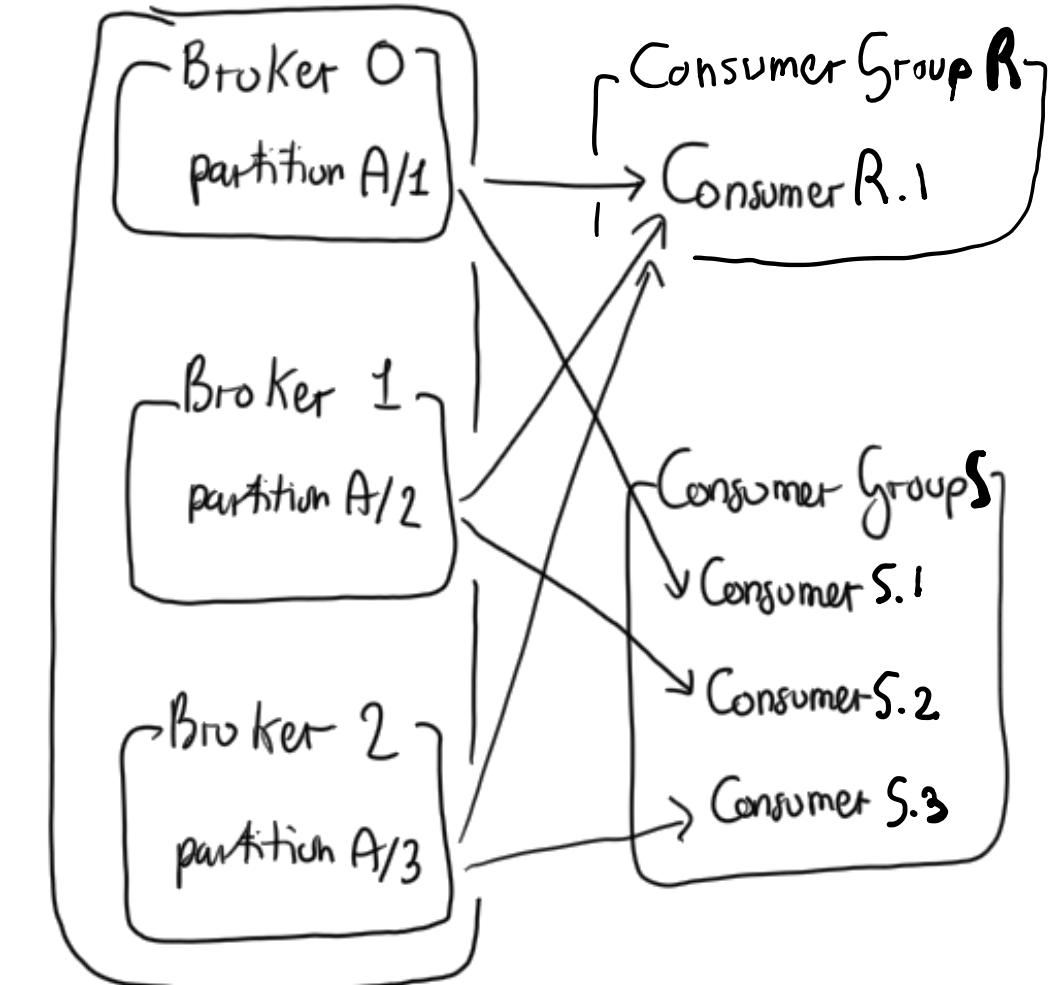
Producers, fault tolerance and durability

- Producers can control durability by requiring the leader a number of acknowledgments before considering the request complete.
 - **acks=0**
Producer will not wait for any acknowledgment from the broker
 - **acks=1**
Producer will wait until the leader has written the record to its local log
 - **acks=all**
Producer will wait until all insync replicas have acknowledged receipt of the record

15 : 50

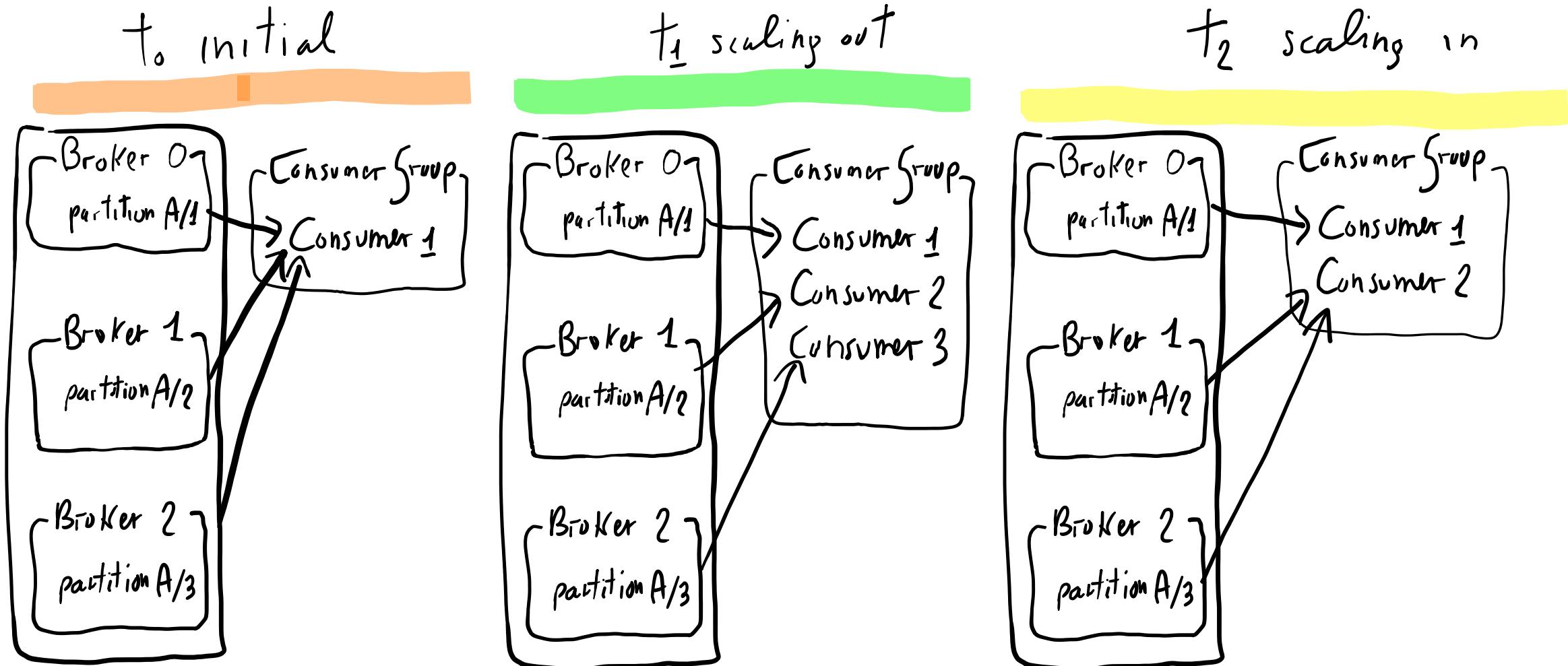
Topic partitioning invites distributed consumption

- Different Consumers can read data from the same Topic
 - By default, each Consumer will receive all the messages in the Topic
- Multiple Consumers can be *logically* combined into a **Consumer Group**
 - Consumer Groups provide scaling capabilities
 - Each Consumer is assigned a subset of Partitions for consumption



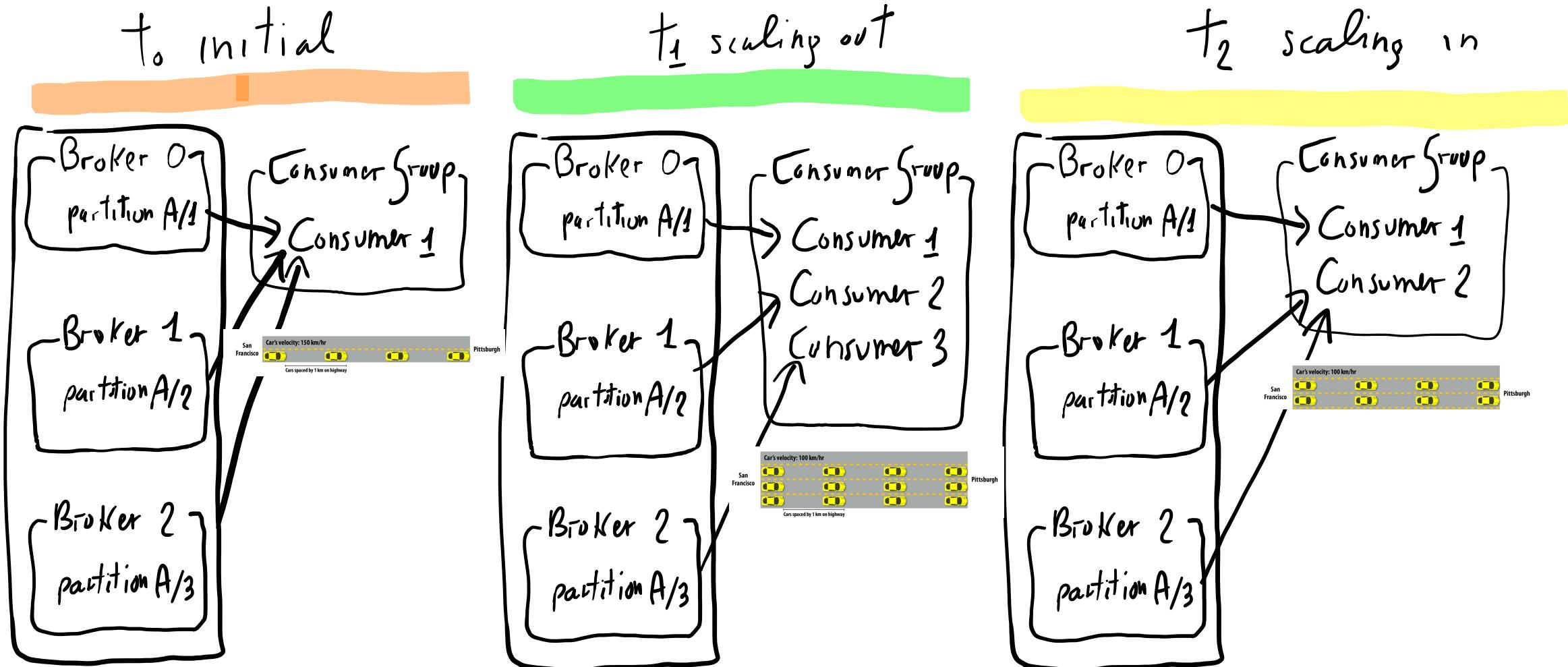


Consumer Group and scalability



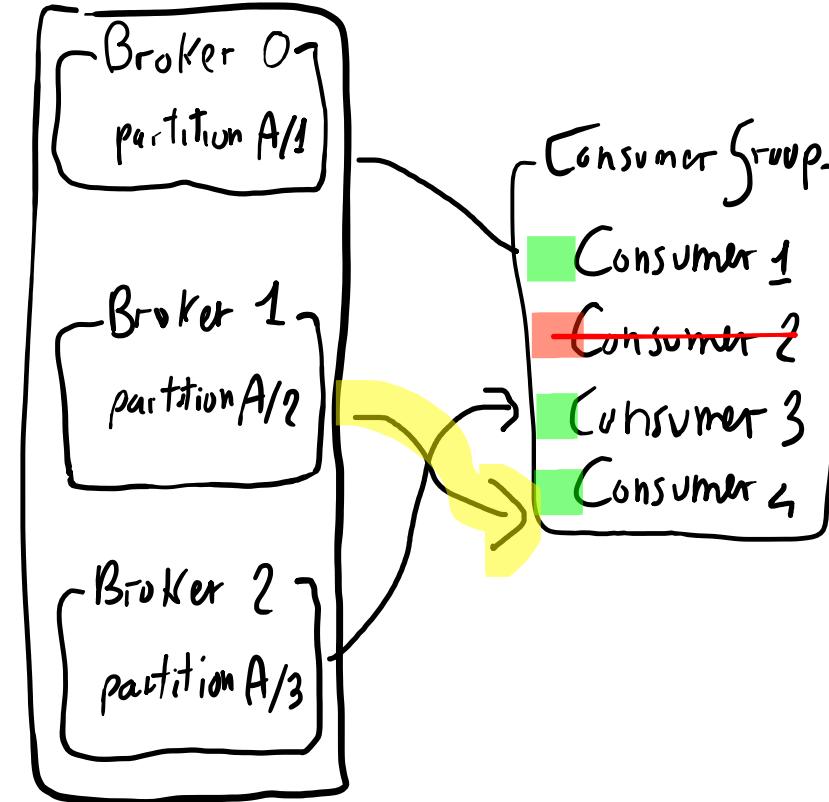
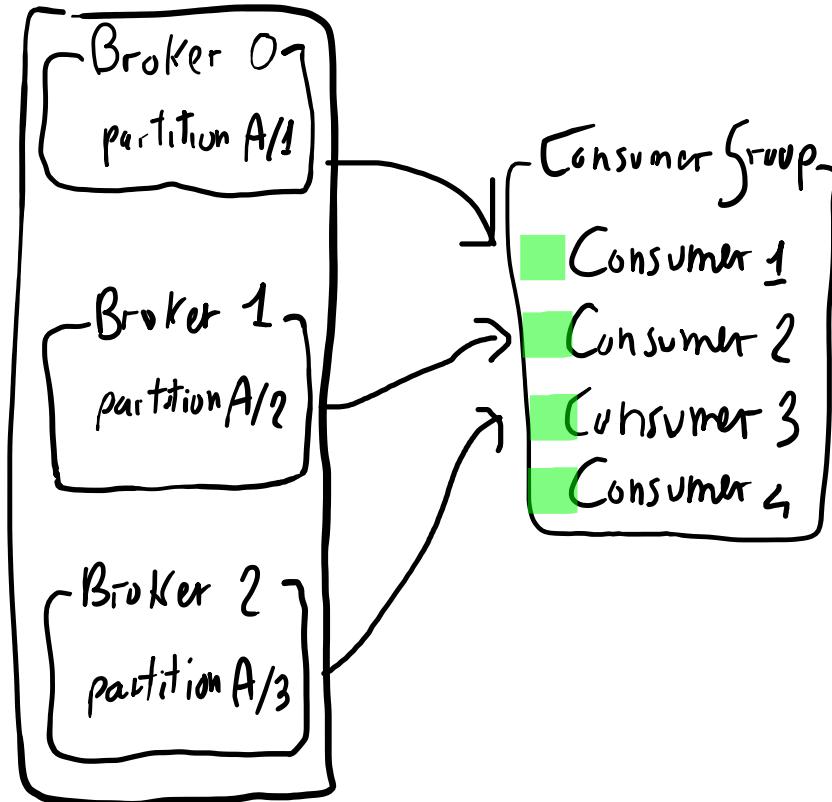


Once again, notice the similarity



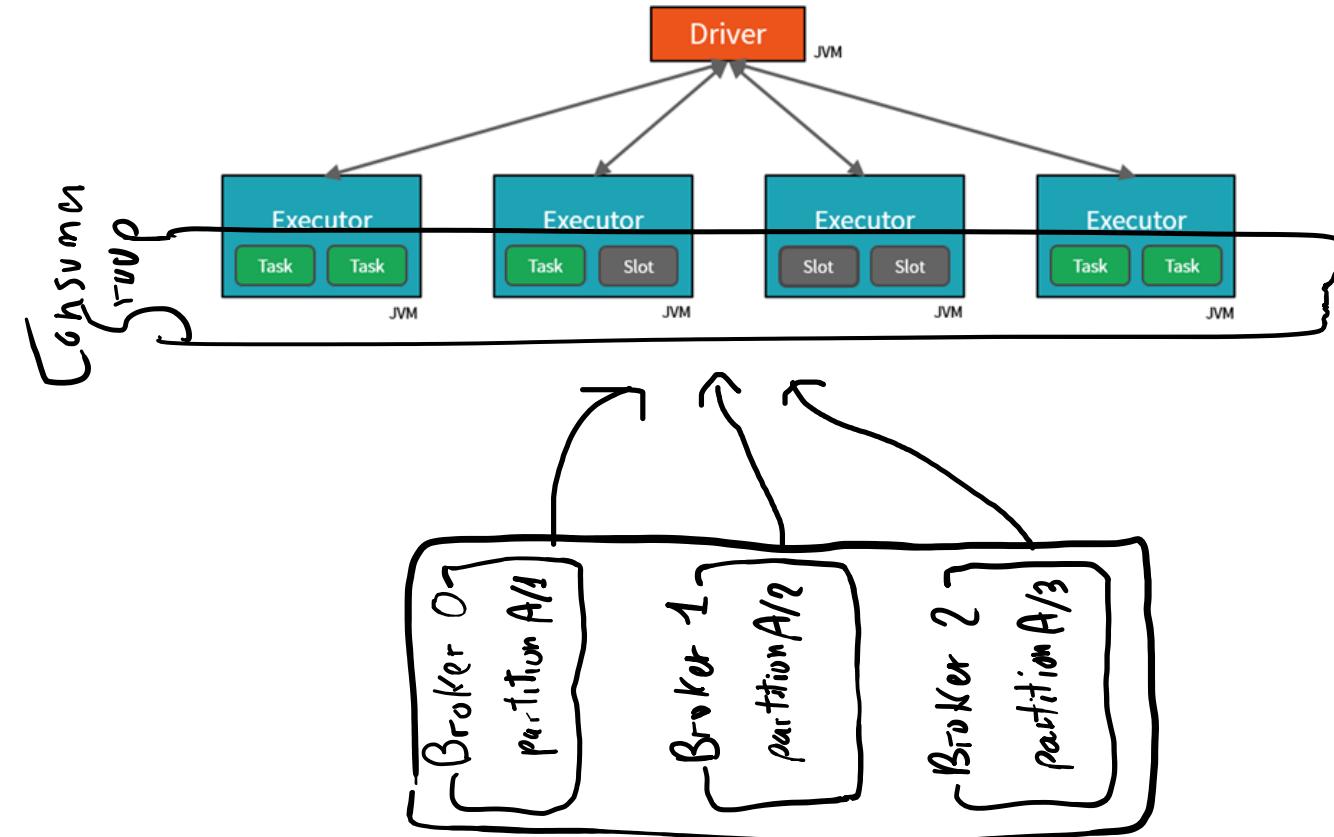


Consumer Group and fault tolerance





Spark, the perfect consumer group



Quiz

- Provide the correct relationship — 1:1, 1:N, N:1, or N:N —
 - Broker to Partition — ?
 - Key to Partition — ?
 - Producer to Topic — ?
 - Consumer Group to Topic — ?
 - Consumer to Consumer Group — ?
 - Consumer (in a Consumer Group) to Partition — ?



A Physical View of topic partition

- Each Partition is stored on the Broker's disk as one or more log files
- Each message in the log is identified by its **offset** number

Commit log

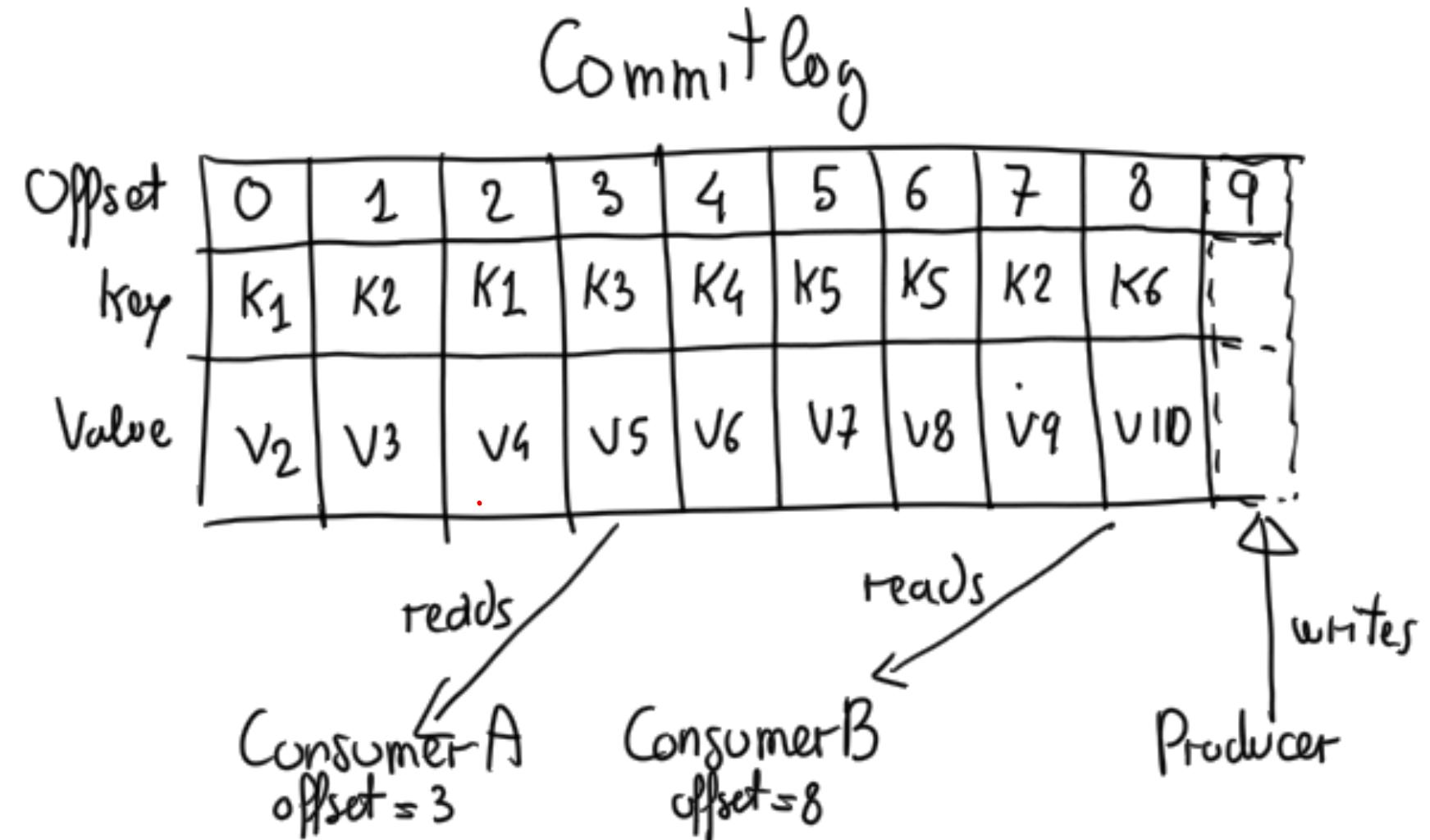
Offset	0	1	2	3	4	5	6	7	8
key	k_1	k_2	k_1	k_3	k_4	k_5	k_5	k_2	k_6
Value	v_2	v_3	v_4	v_5	v_6	v_7	v_8	v_9	v_{10}

17:05



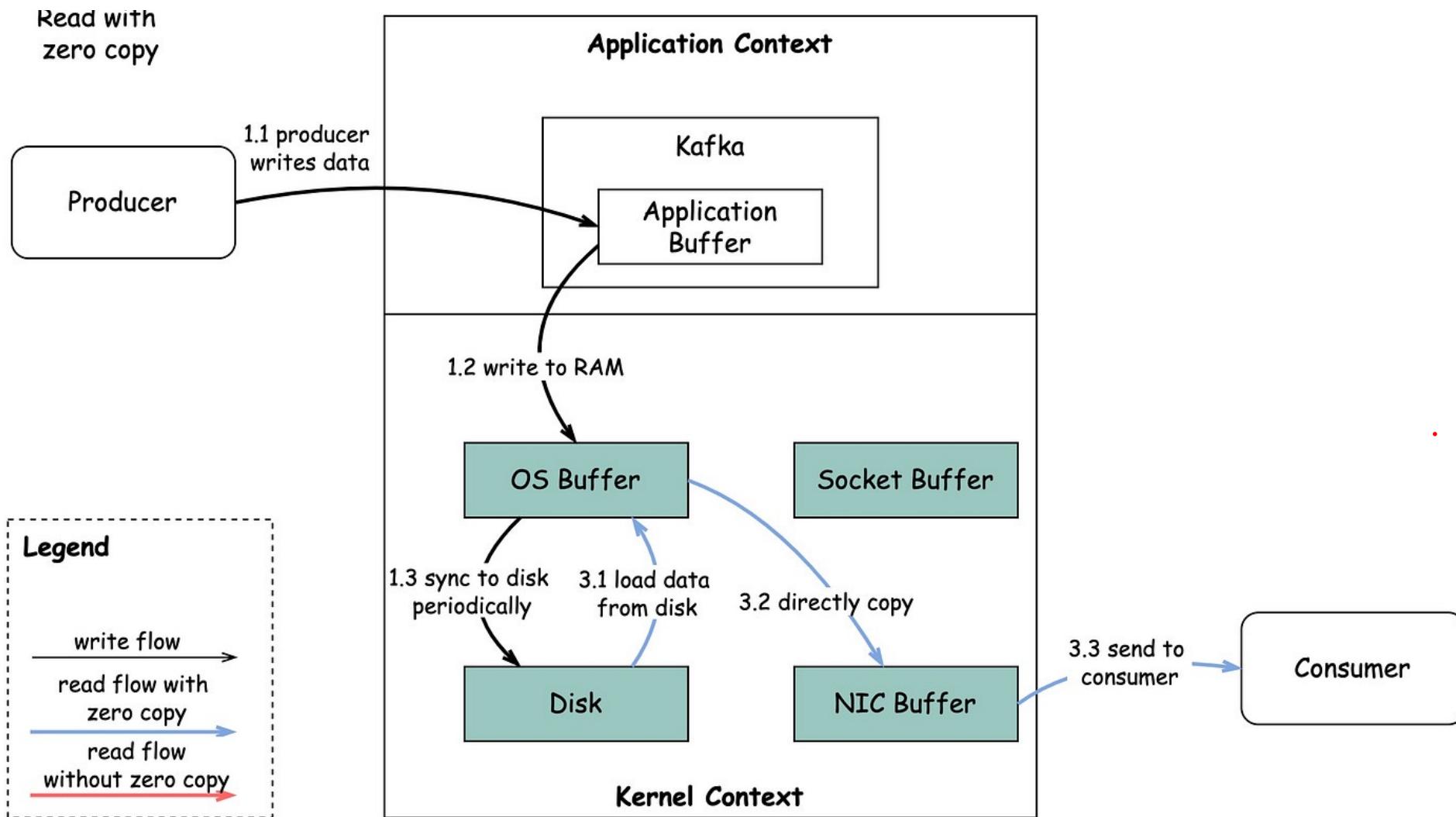
A Physical View of topic partition (cont.)

- Messages are always appended
- Consumers can consume from different offset
- Brokers are single thread to guarantee consistency





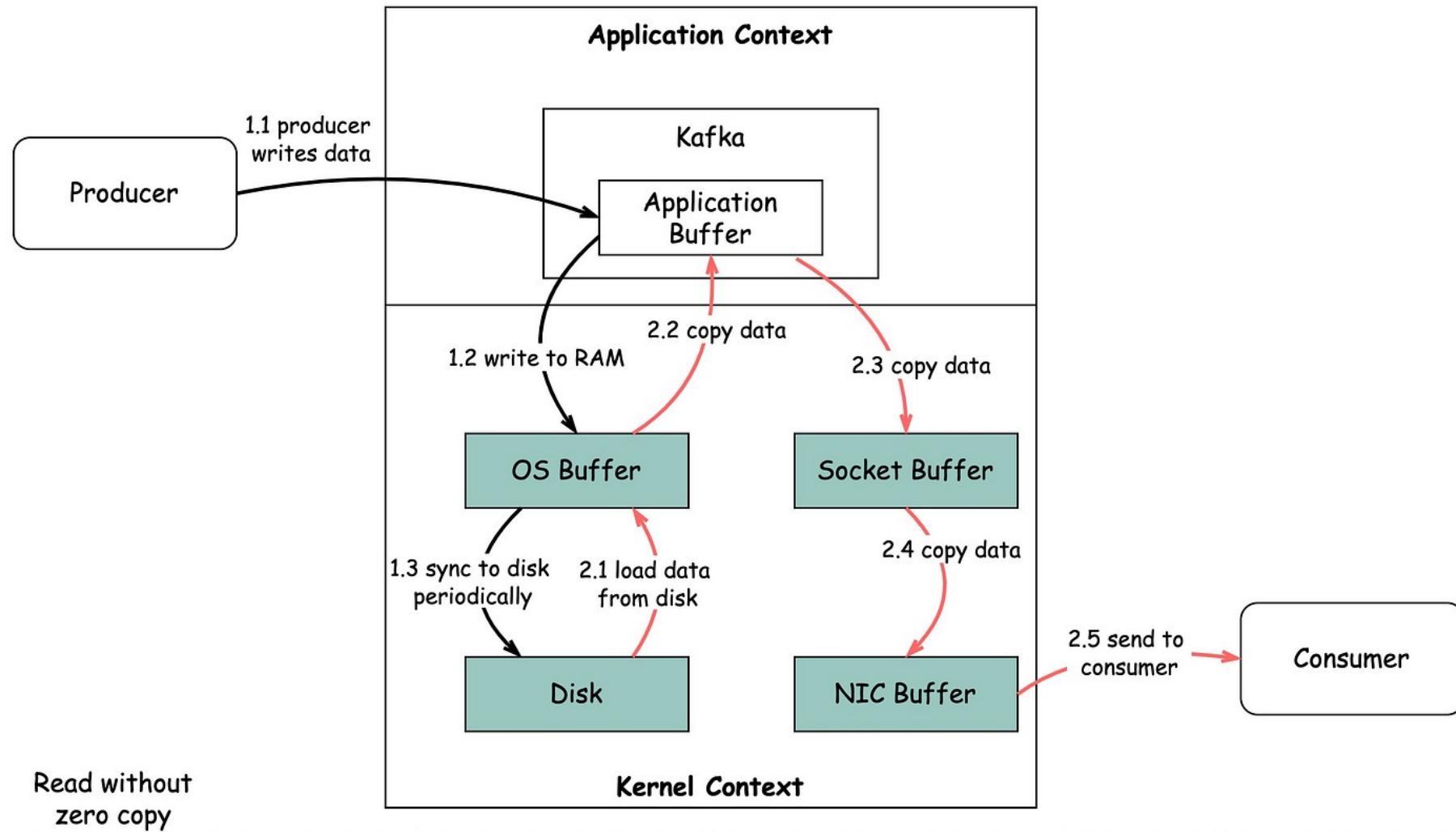
Zero-copy!



[src: <https://blog.bytebytogo.com/p/why-is-kafka-fast>]



Without zero-copy ...





Log Retention

- Duration default: messages will be retained for seven days
- Duration is configurable per Broker by setting
 - a time period
 - a size limit (number of messages or size in GB)
- The topic configuration can override a Broker's retention policy
- When cleaning up a log
 - the default policy is to delete
 - An **alternative** policy is compact



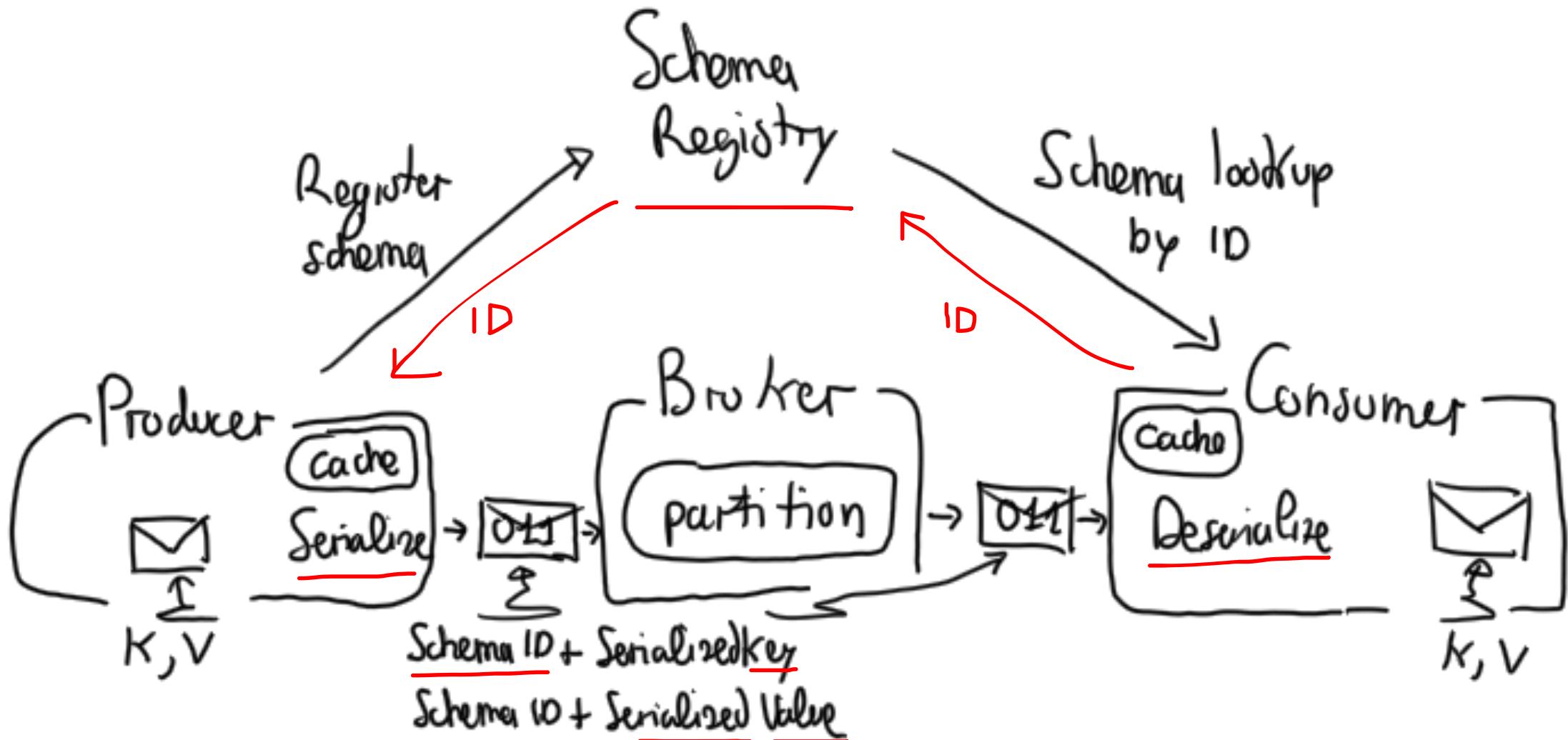
Log compaction

A compacted log retains at least the last known message value for each key within the Partition

Before compaction									After compaction						
Offset	0	1	2	3	4	5	6	7	8	0	1	2	3	4	5
key	<u>K₁</u>	K ₂	<u>K₁</u>	K ₃	K ₄	<u>K₅</u>	<u>K₅</u>	K ₂	K ₆	K ₁	K ₃	K ₄	K ₅	K ₂	K ₆
Value	V₂	V₃	V₄	V ₅	V ₆	V₇	V ₈	V₉	V ₁₀	V ₄	V ₅	V ₆	V ₈	V ₉	V ₁₀

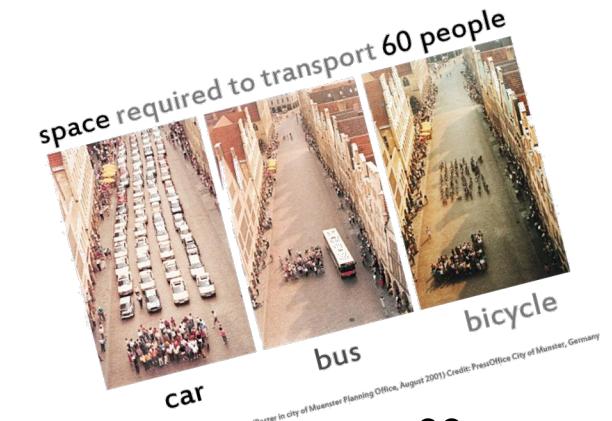


Avro and Schema Registry at work



Avro

- Avro is an Apache open source project
- Provides data serialization into a binary format
 - so stores data efficiently
- Data is defined with a self-describing schema allowing for
 - code generation for serializers and de-serializers in multiple languages
 - type checking at write time





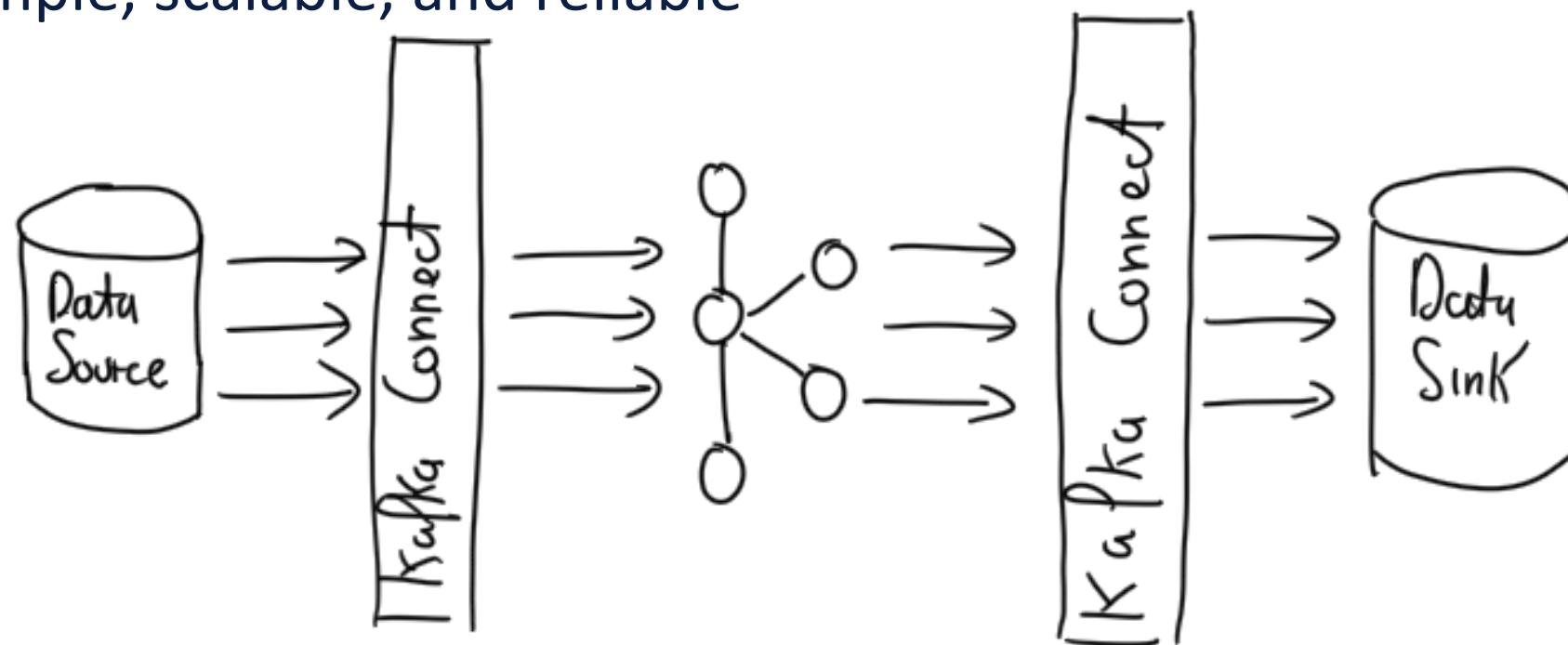
Schema Registry

- Sending the Avro schema with each message would be inefficient
- The Schema Registry
 - is a Confluent component
 - provides centralized management of schemas
 - At registration time each schema is given an ID
- The mapping schema/ID is stored in a special Kafka topic



Kafka Connect

- It is an open source framework for streaming data between Apache Kafka and other data systems
- It is simple, scalable, and reliable





Off-The-Shelf Connectors

<https://www.confluent.io/product/connectors/>

<p>Verified Standard</p>  <p>StreamSets Data Collector</p> <p>StreamSets</p> <p>Read More</p>	<p>Verified Gold</p>  <p>yugabyteDB</p> <p>Kafka Connect Yugabyte</p> <p>Yugabyte, Inc.</p> <p>Read More</p>	<p>Verified Standard</p>  <p>hazelcast JET</p> <p>Hazelcast Jet Kafka Connector</p> <p>Hazelcast</p> <p>Read More</p>	<p>Verified Standard</p>  <p>ORACLE GOLDEN GATE</p> <p>Oracle GoldenGate</p> <p>Oracle</p> <p>Read More</p>	<p>Confluent Supported</p>  <p>IBM MQ</p> <p>Kafka Connect IBM MQ Sink</p> <p>Confluent, Inc.</p> <p>Read More</p>	<p>Confluent Supported</p>  <p>Kafka Connect MQTT</p> <p>Confluent, Inc.</p> <p>Read More</p>	<p>Confluent Supported</p>  <p>Kafka Connect Google Cloud Functions Sink Connector</p> <p>Confluent, Inc.</p> <p>Read More</p>	<p>Confluent Supported</p>  <p>Kafka Connect S3</p> <p>Confluent, Inc.</p> <p>Read More</p>
<p>Verified Standard</p>  <p>druid</p> <p>Druid Kafka indexing service</p> <p>Imply</p> <p>Read More</p>	<p>Verified Standard</p>  <p>HVR</p> <p>HVR Change Data Capture</p> <p>HVR</p> <p>Read More</p>	<p>Verified Standard</p>  <p>B.O.S. Software tcVISION</p> <p>B.O.S. Software</p> <p>Read More</p>	<p>Verified Standard</p>  <p>SQData CDC Connector</p> <p>SQData Corporation</p> <p>Read More</p>	<p>Confluent Supported</p>  <p>Kafka Connect Pivotal Gemfire</p> <p>Confluent, Inc.</p> <p>Read More</p>	<p>Confluent Supported</p>  <p>JMS</p> <p>Kafka Connect JMS Source</p> <p>Confluent, Inc.</p> <p>Read More</p>	<p>Confluent Supported</p>  <p>TIBCO</p> <p>Kafka Connect TIBCO Sink</p> <p>Confluent, Inc.</p> <p>Read More</p>	<p>Confluent Supported</p>  <p>RabbitMQ</p> <p>Kafka Connect RabbitMQ</p> <p>Confluent, Inc.</p> <p>Read More</p>



POLITECNICO
MILANO 1863

Kafka

Scaling continuous ingestion

Emanuele Della Valle

emanuele.dellavalle@polimi.it