## [544] Kafka Streaming

Meenakshi Syamkumar

### Learning Objectives

- describe the benefits of using streaming for ETL (extract transform load) work
- write code for Kafka consumers and producers in order to interact with topic data that stored by brokers
- scale out brokers and consumers by configuring topic partitions and consumer groups, respectively

### Outline: Kafka Streaming

#### Sending/Receiving Messages

- RPC (Remote Procedure Calls)
- Streaming

ETL (Extract Transform Load)

Kafka Design

#### Procedure Calls

```
counts = {
   "A": 123, ...
}

def increase(key, amt):
   counts[key] += amt
   return counts[key]

curr = increase("A", 5)
print(curr) # 128
```

what if we want many programs running on different computers to have access to this dict and the increase function?

### Remote Procedure Calls (RPCs)

client

curr = increase("A", 5)
print(curr) # 128

server

```
counts = {
   "A": 123, ...
}

def increase(key, amt):
   counts[key] += amt
   return counts[key]
```

client

move counts and increase to a server accessible to many client programs on different computers

• • •

### Remote Procedure Calls (RPCs)

#### client

```
def increase(key, amt):
    ...code to send

curr = increase("A", 5)
print(curr) # 128
```

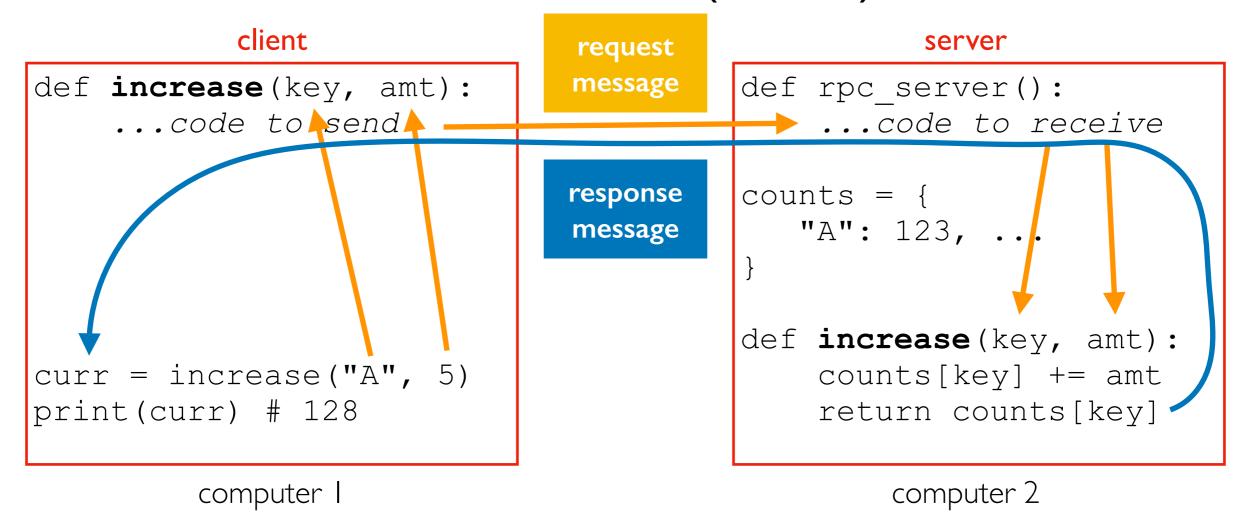
computer I

#### server

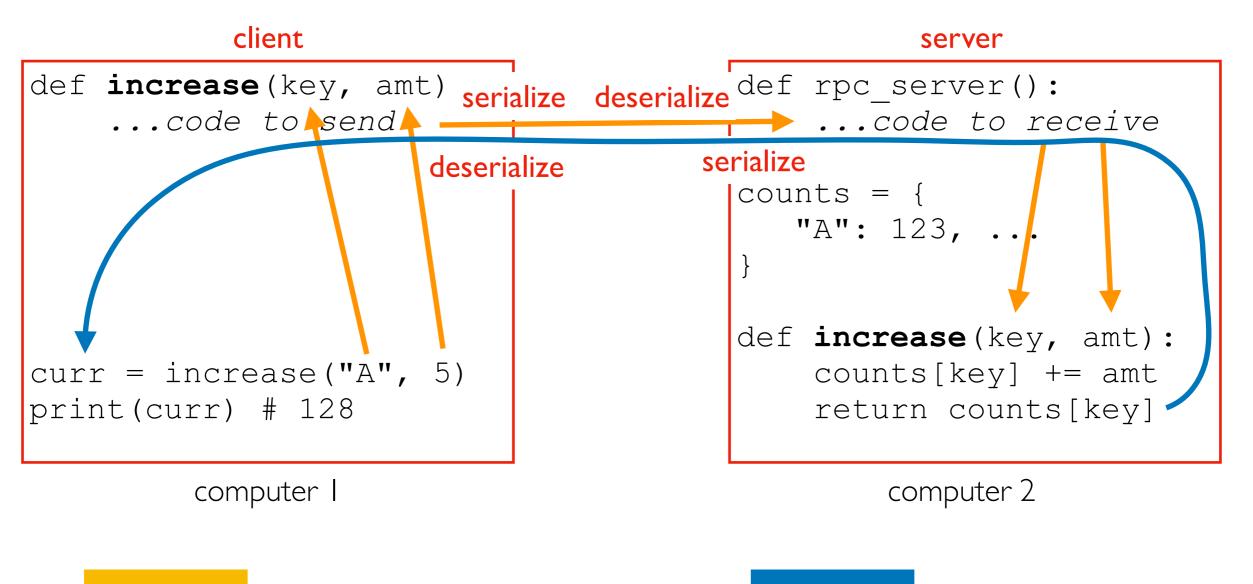
computer 2

need some extra functions to make calling a remote function feel the same as calling a regular one

### Remote Procedure Calls (RPCs)



#### Serialization/Deserialization

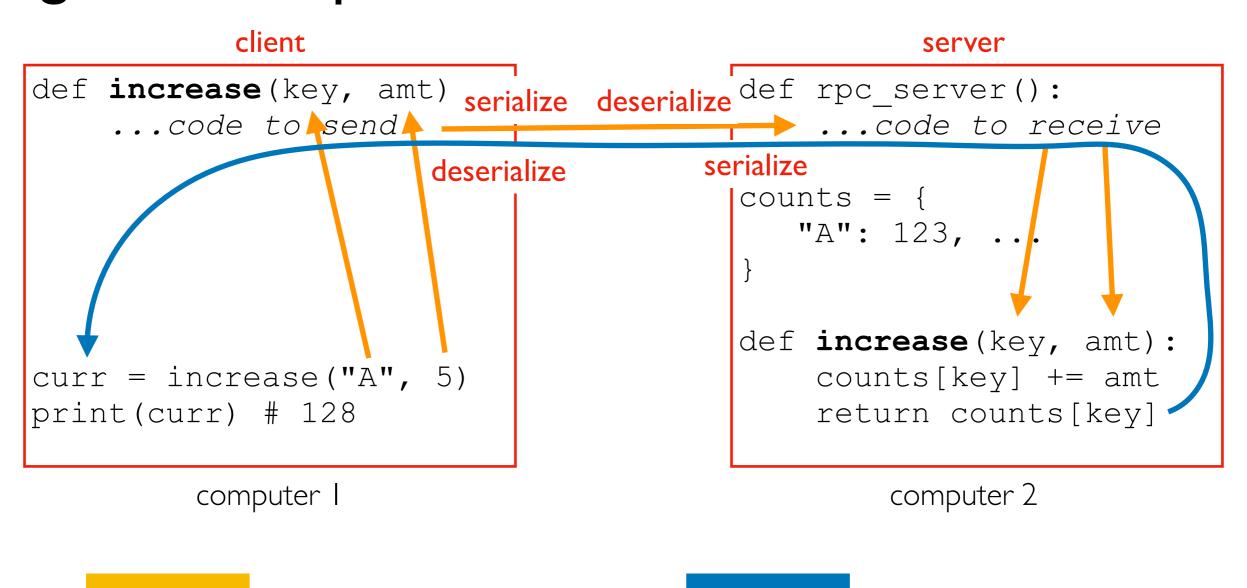


# request message args somehow encoded as bytes: b'{"key": "A" "amt": 5}'

#### response message

```
return val as bytes:
b'5'
```

#### gRPC uses protocol buffers for wire format



#### request message

```
protobuf (args to bytes)
b'1001000101011111'
(contains "A" and 5)
```

#### response message

```
protobuf (ret val to bytes)
b'01000000'
(contains 128)
```

#### Synchronous vs. Asynchronous Communication

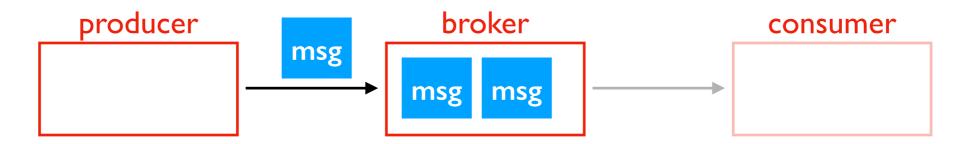
#### **Synchronous**

- both parties have to participate at the same time
- examples: phone call, RPC call



#### Asynchronous

- one party can send any time, the other can receive later
- examples: email, streaming



### Outline: Kafka Streaming

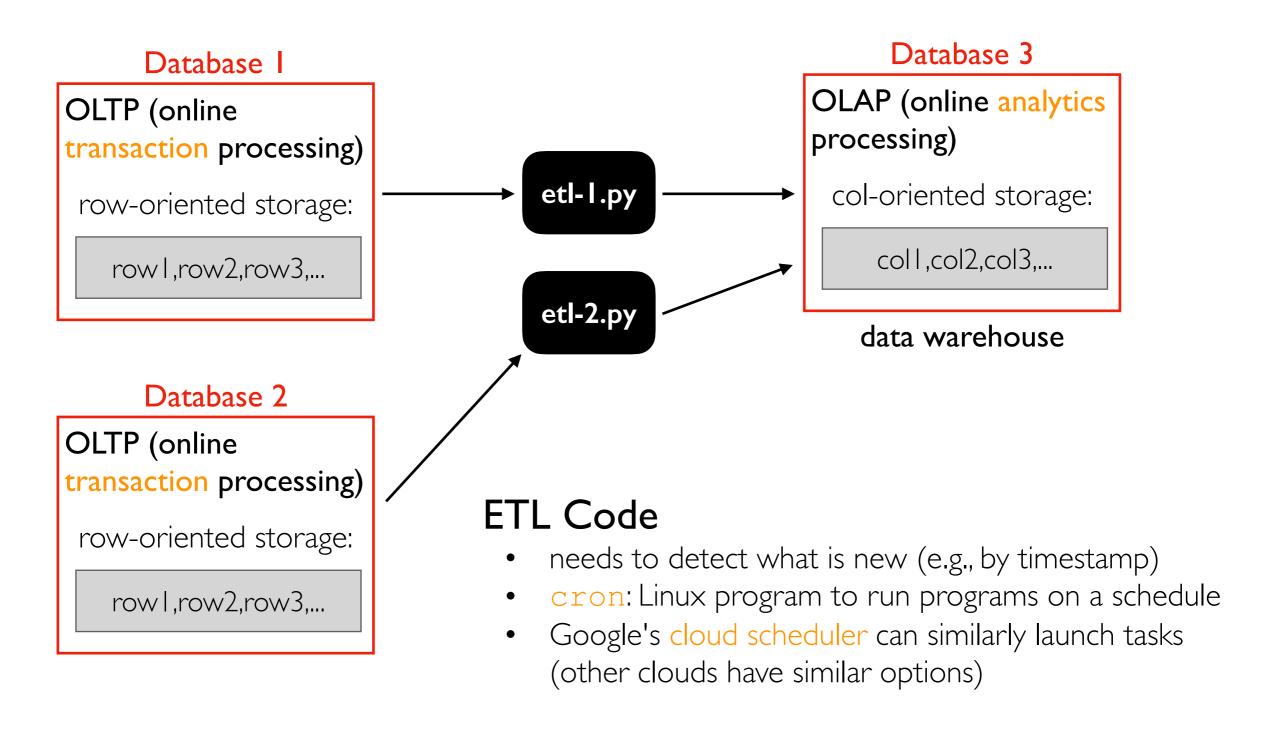
Sending/Receiving Messages

#### ETL (Extract Transform Load)

- Batch
- Streaming

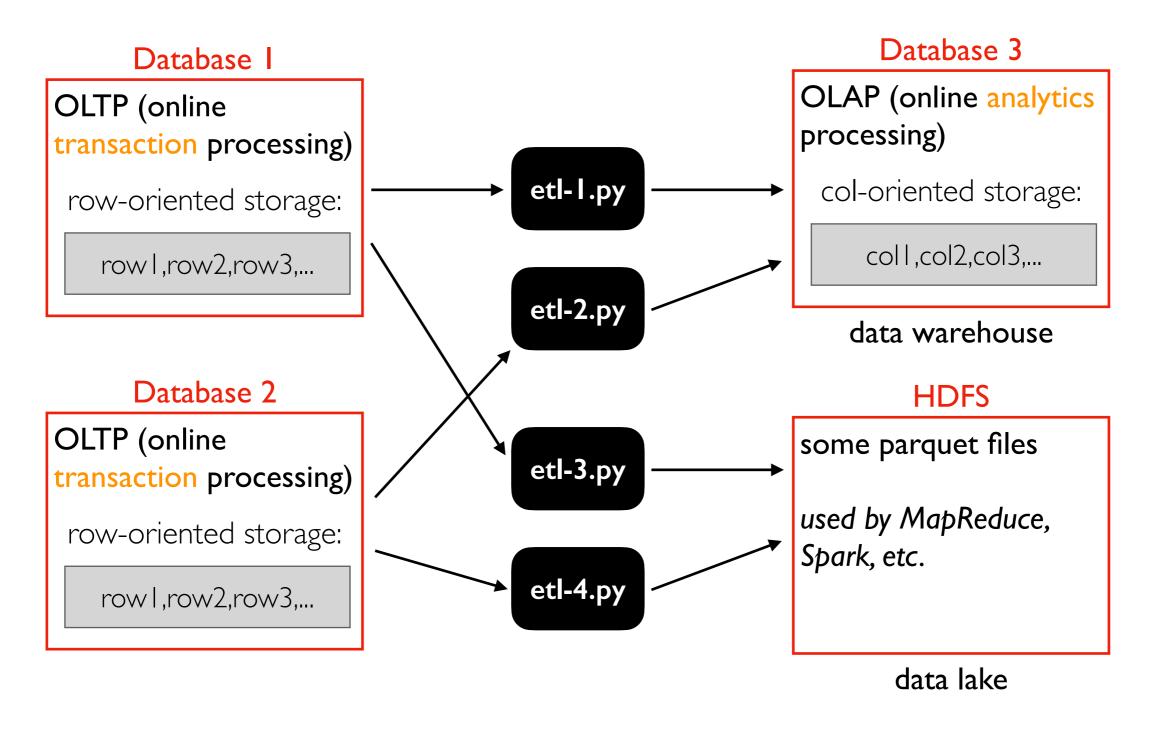
Kafka Design

### Extract Transform Load (ETL)



issue I: data freshness

### Extract Transform Load (ETL)

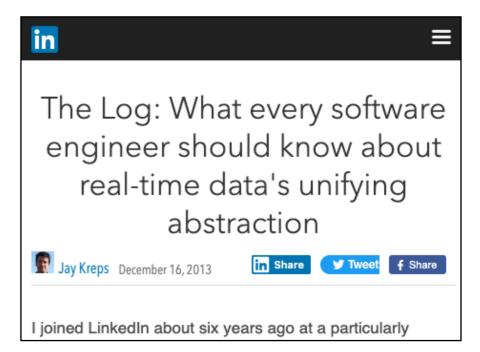


if we have **X** OLTP databases and **Y** derivative stores, how many ETL programs must we write?

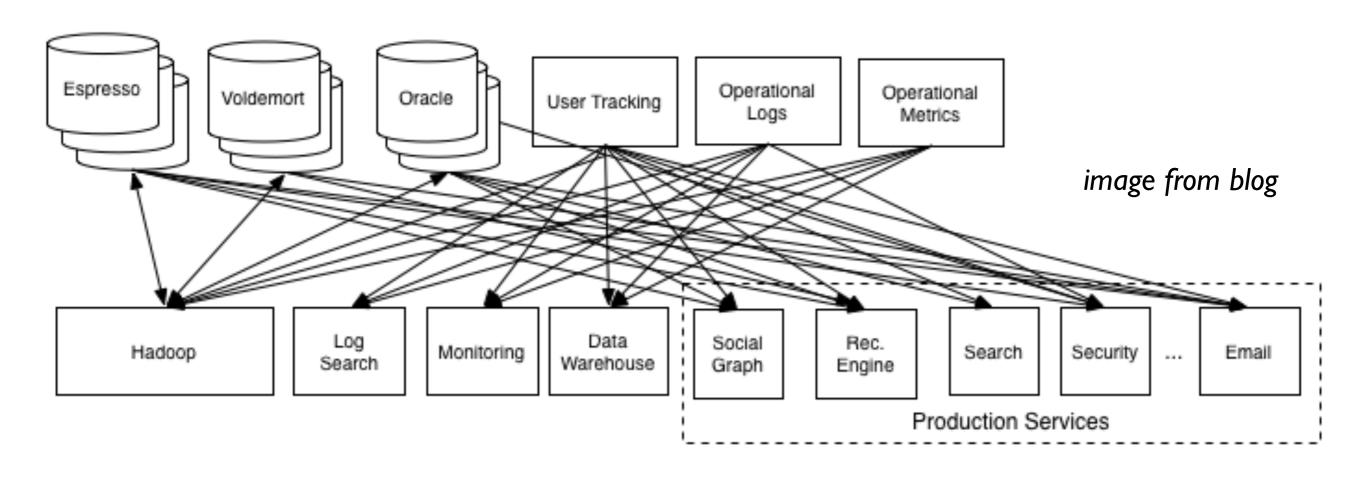
#### Too much ETL...

Don't want data transfer between every pair of DB/services

- Jay Krepps helped build Kafka at LinkedIn
- Later co-founded Confluent (Kafka-based company)
- Partners with cloud providers to provide Kafka as a service



https://engineering.linkedin.com/distributed-systems/log-what-every-software-engineer-should-know-about-real-time-datas-unifying

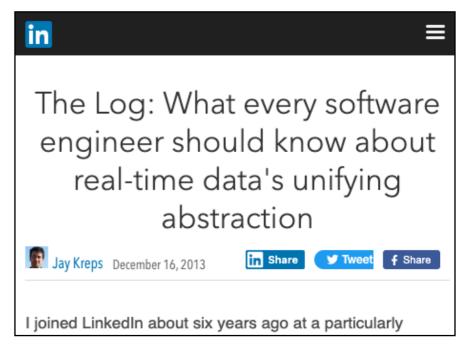


### Unified Log

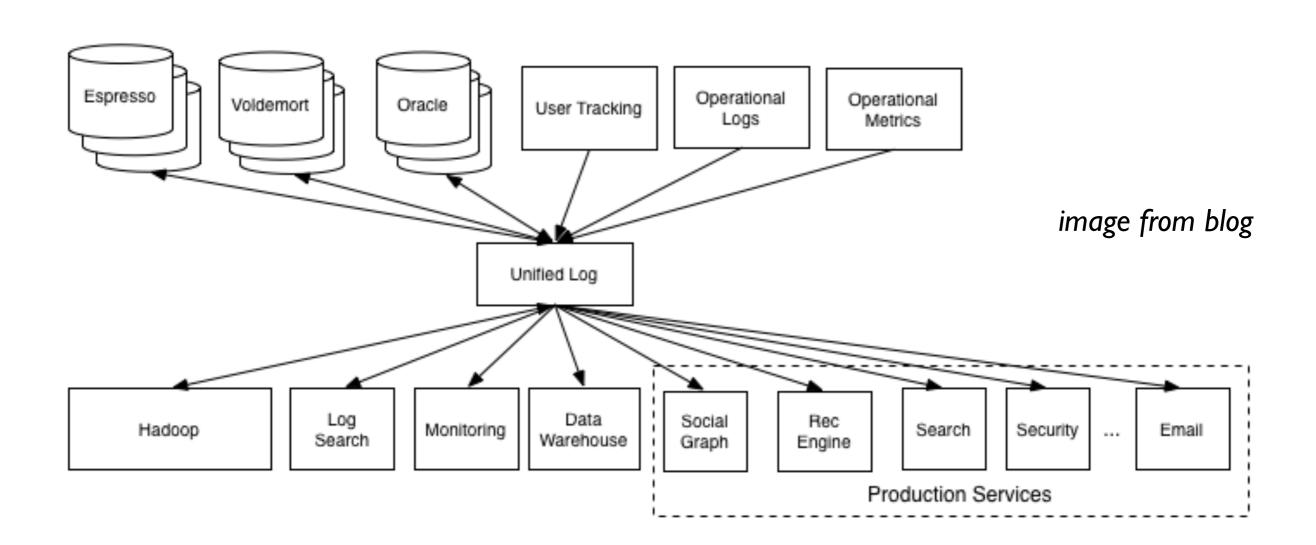
Centralize changes in a distributed logging service

- Many writers (called producers)
- Many readers (called consumers)

Data is constantly flowing, so ETL can be done in realtime (instead of batch jobs with cron)



https://engineering.linkedin.com/distributed-systems/log-what-every-software-engineer-should-know-about-real-time-datas-unifying



### Outline: Kafka Streaming

Sending/Receiving Messages

ETL (Extract Transform Load)

#### Kafka Design

- Topics
- Producers, Consumers, Brokers
- Scalability with Partitioning

### **Topics**

Kafka topics (managed by servers called brokers)

weather msg msg

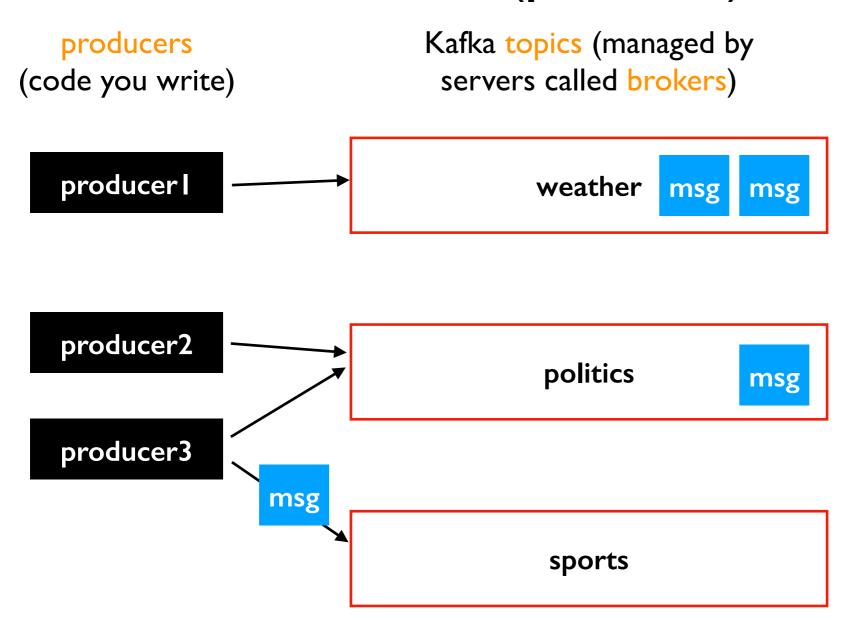
politics msg

sports

admin = **KafkaAdminClient**(...)
admin.create\_topics([NewTopic("sports", ...)])

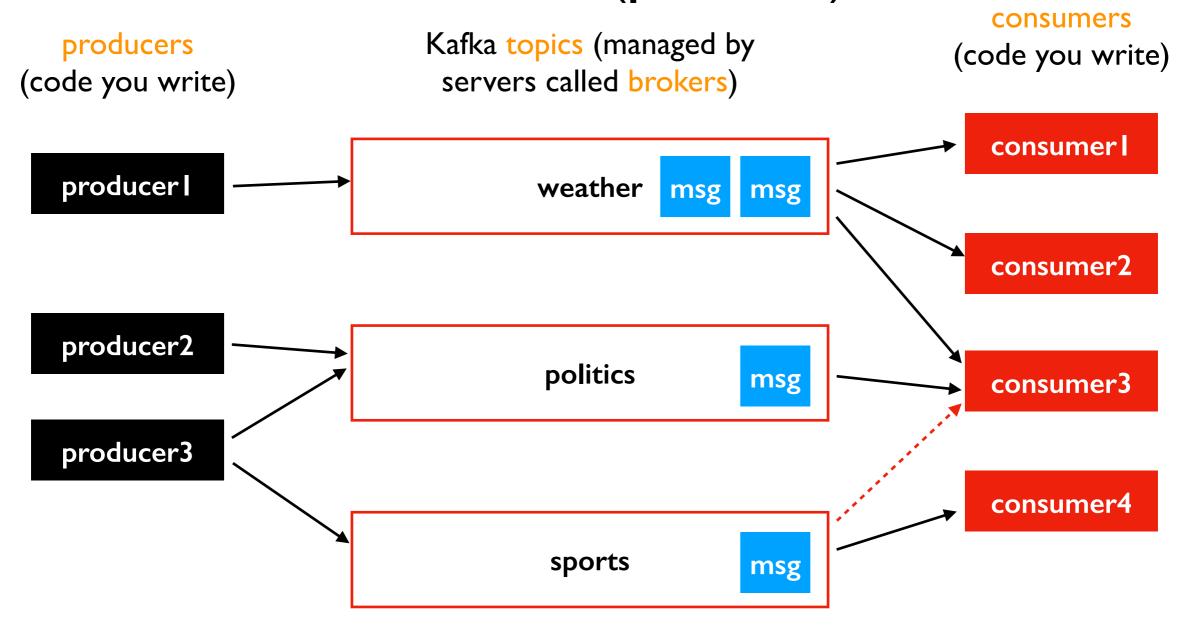
pip install kafka-python

### Producers Publish (pub/sub)



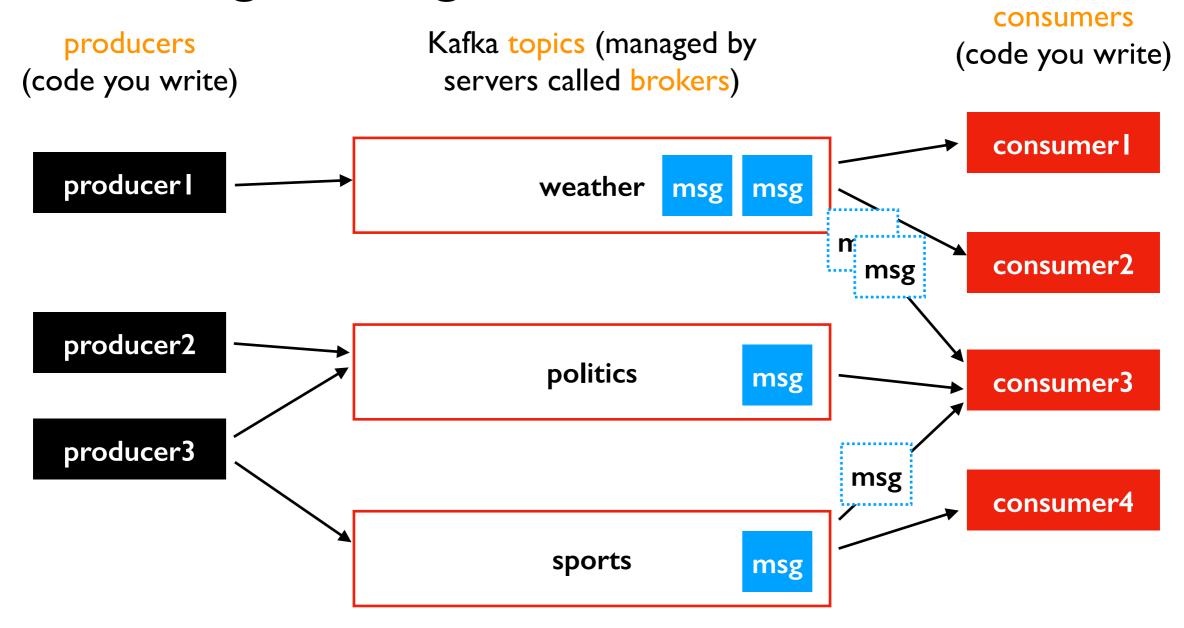
producer3 = KafkaProducer(...)
producer3.send("sports", ...)

### Consumers Subscribe (pub/sub)



consumer3 = KafkaConsumer(...)
consumer3.subscribe(["sports"])

### Receiving Messages



#### poll() loop

- generally runs forever
- poll (ideally) returns some messages the consumer hasn't seen before, from any subscribed topic
- leaves messages intact on brokers (for other consumers), unlike many prior streaming systems

```
consumer3 = KafkaConsumer(...)
while True:
  batch = consumer3.poll(????)
for topic, messages in batch.items():
  for msg in messages:
```

### What's in a Message?

#### Message parts

- key (optional): some bytes
- value (required): some bytes
- other stuff...

```
producer.send("topic", value=????)
OR
producer.send("topic", value=????, key=????)
```

Common usage: the value is usually some kind of structure with many values. The key is used for partitioning and is usually one of the entries in the value structure.

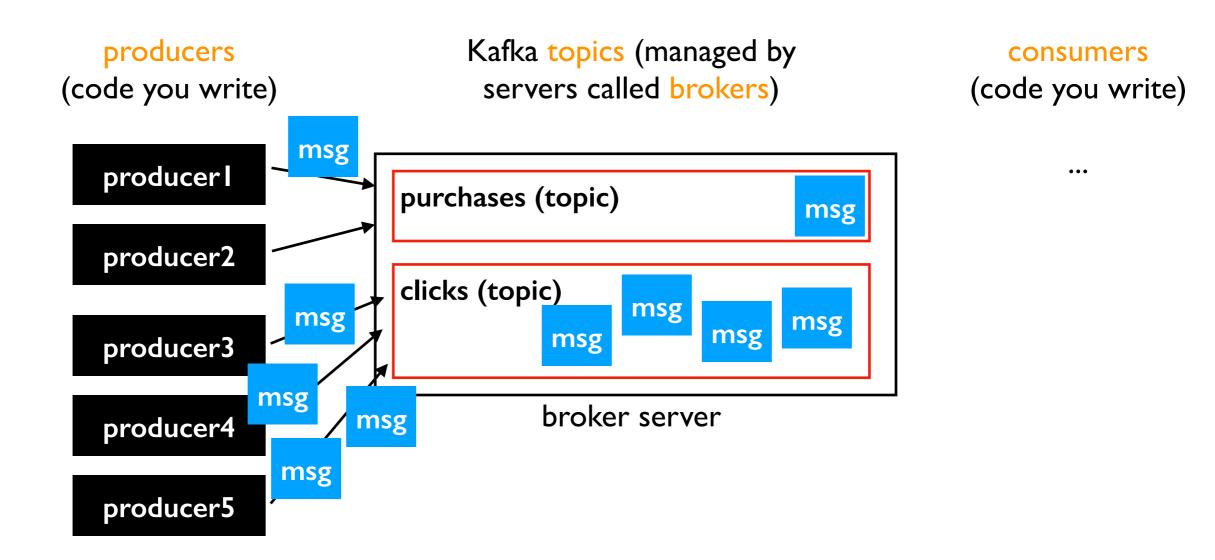
#### Python dict => bytes:

```
d = {...}
value = bytes(json.dumps(d), "utf-8")
```

#### Protobuf => bytes:

```
msg = mymod_pb2.MyMessage(...)
value = msg.SerializeToString() # actually bytes, not str
```

### Scaling the Brokers



problem: some topics might have too many messages for one machine (or set of machines with replicas) to keep up

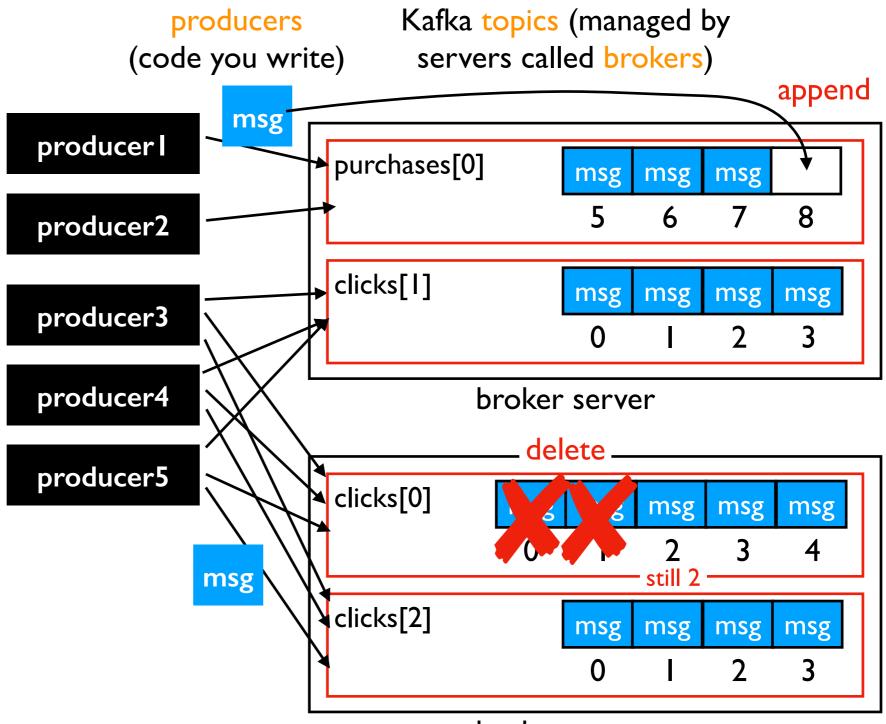
#### **Partitions**

producers Kafka topics (managed by (code you write) servers called brokers) msg producer l purchases[0] msg msg msg 5 producer2 clicks[1] msg msg msg msg producer3 0 3 producer4 broker server producer5 clicks[0] msg msg msg msg msg 0 4 msg clicks[2] msg msg msg msg 0 3 broker server

Topics can be created with N partitions

- each partition is like an array of messages
- partitions are assigned to brokers
- each producer using a stream works with all partitions

# Changing Partitions



Changes

- append right
- delete left (depends on "retention" policy)
- delete does NOT change indexes

### Selecting **Partitions**

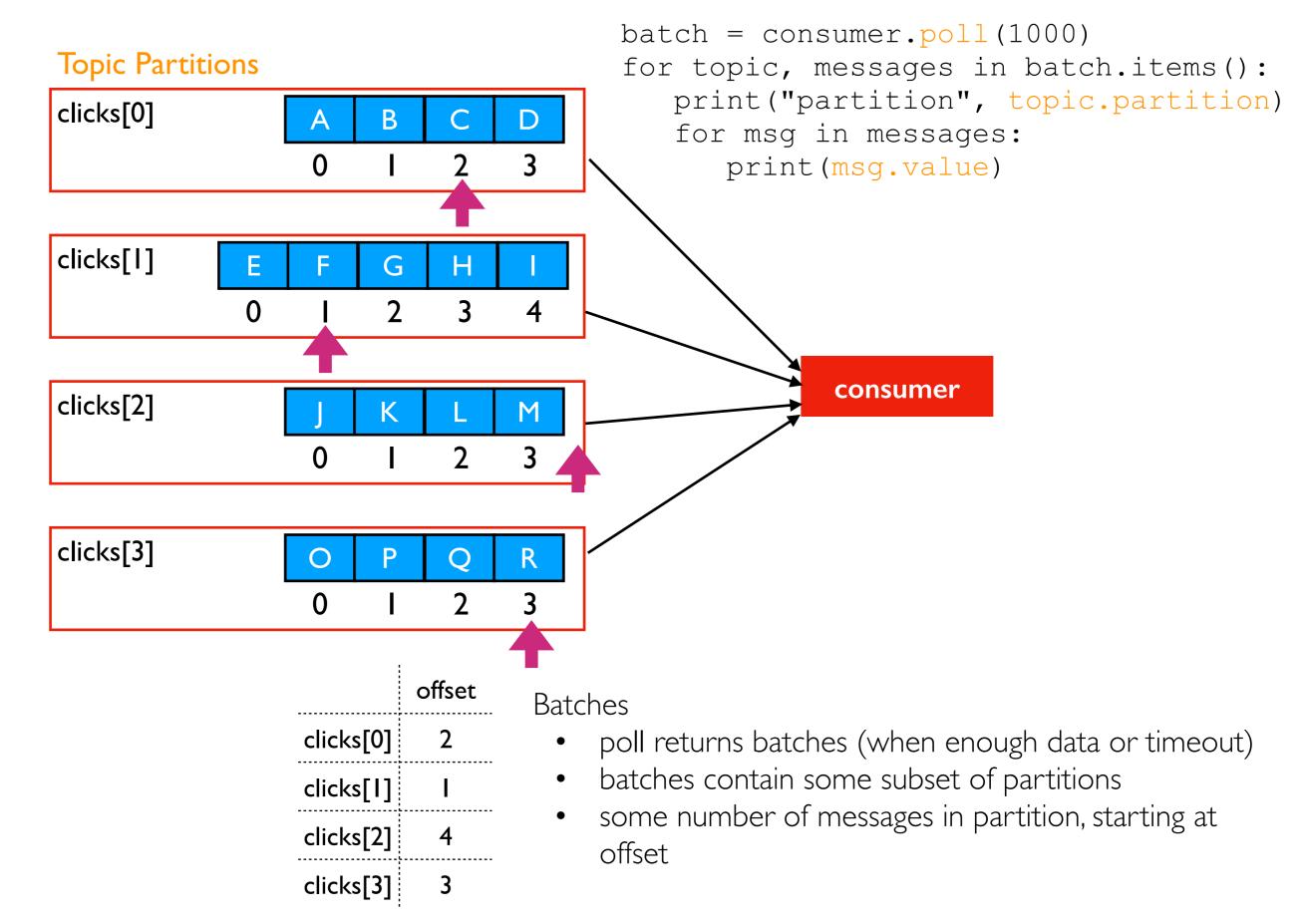
producers Kafka topics (managed by (code you write) servers called brokers) msg producer l purchases[0] msg msg msg 5 producer2 clicks[1] msg msg msg msg producer3 0 3 producer4 broker server producer5 clicks[0] msg msg msg msg msg 0 msg case I: message only has value clicks[2] producer rotates between partitions msg msg msg msg called "round robin" policy 3 0

broker server

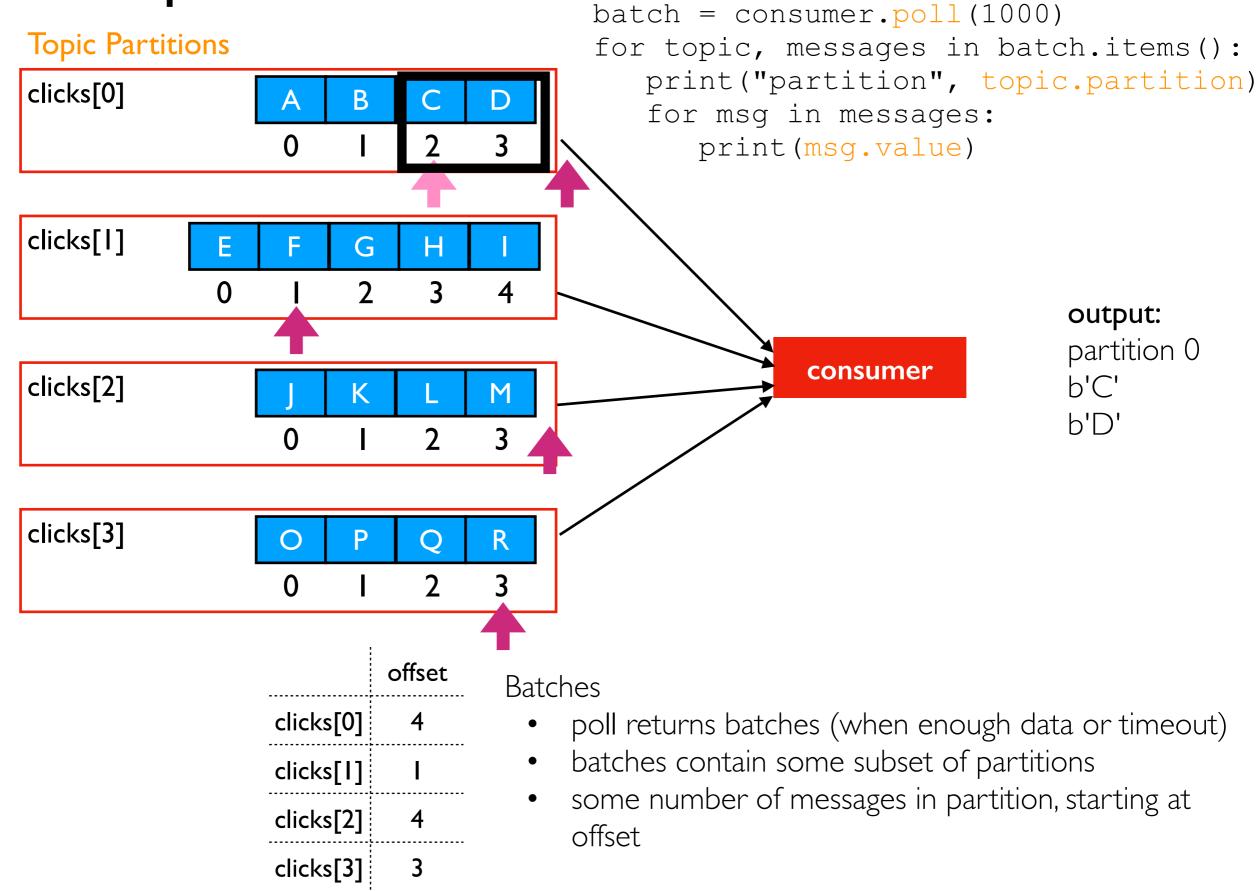
case 2: message has key and value

- calculate partition, for example: hash(key) % partition count
- same keys will go to the same partition
- can plug in alternative partitioning schemes

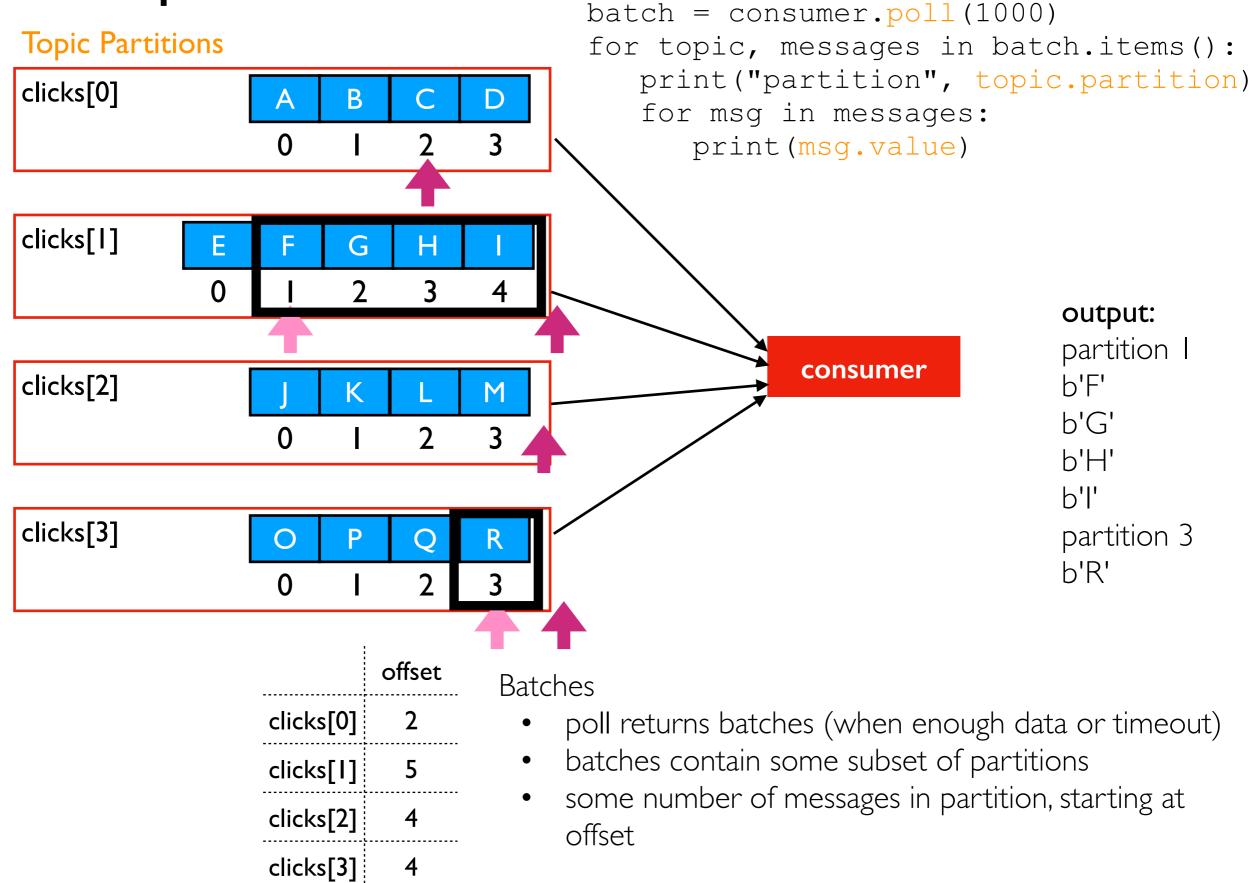
#### Consumers: Read Offsets



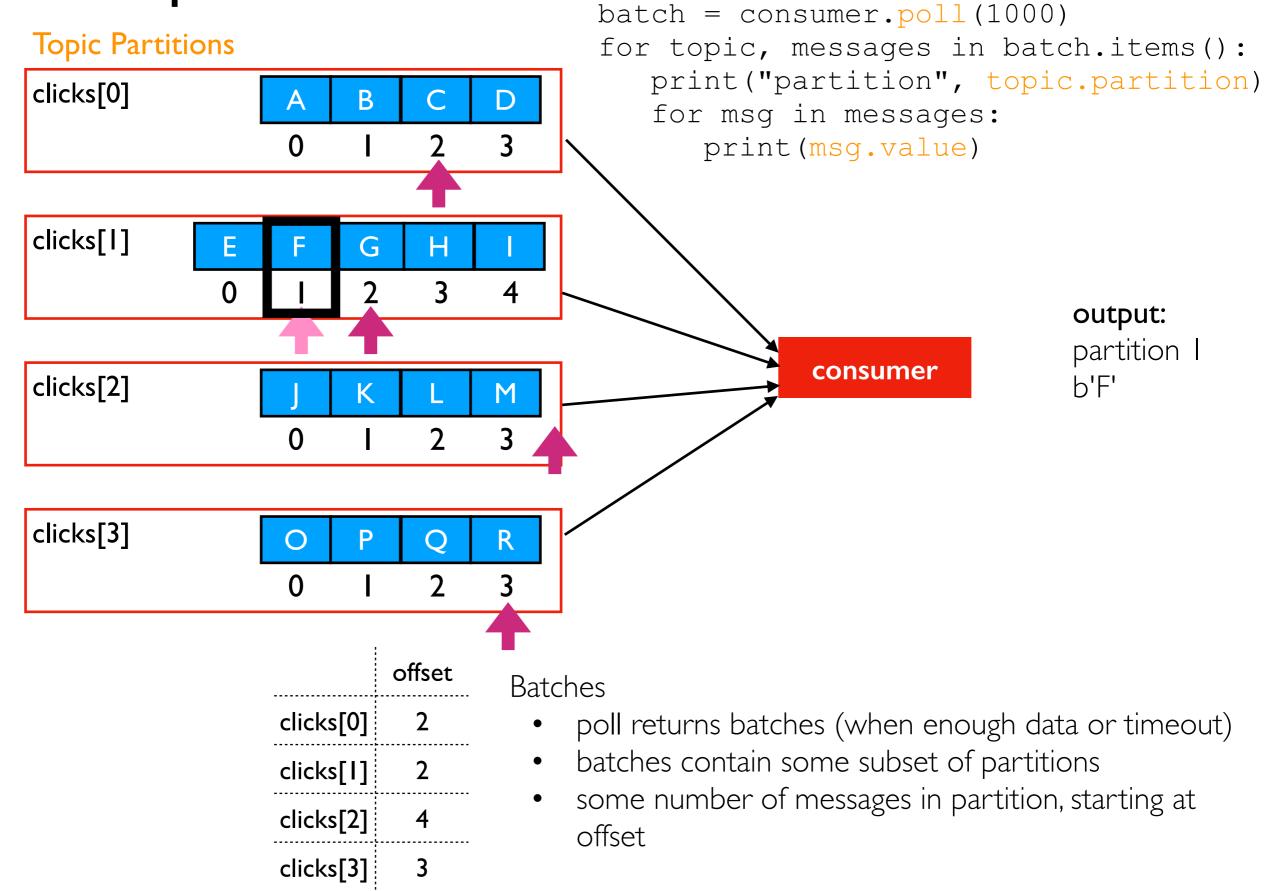
### Example I



#### Example 2



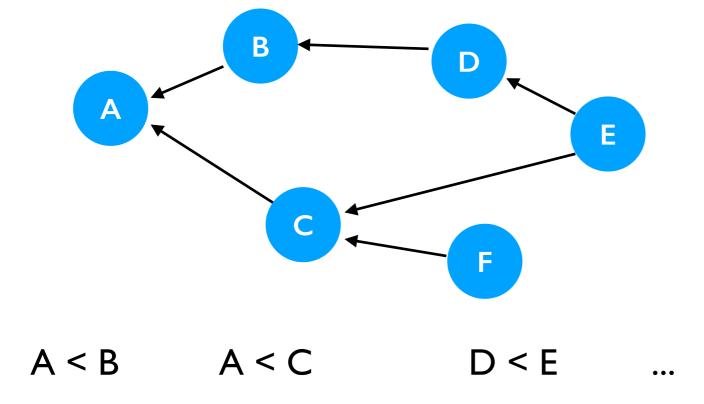
#### Example 3



### Partially vs. Totally Ordered

Some things are totally ordered, like integers. Either x < y or y >= x.

Other things are partially ordered, like git commits. Sometimes you can compare, sometimes you can't!



Can't compare B and C Can't compare D and F

• • •

### Ordering Kafka Messages

Kafka Messages are partially ordered. Messages are consumed from a partition in the order they were written to that partition (no guarantees across topics or across partitions).

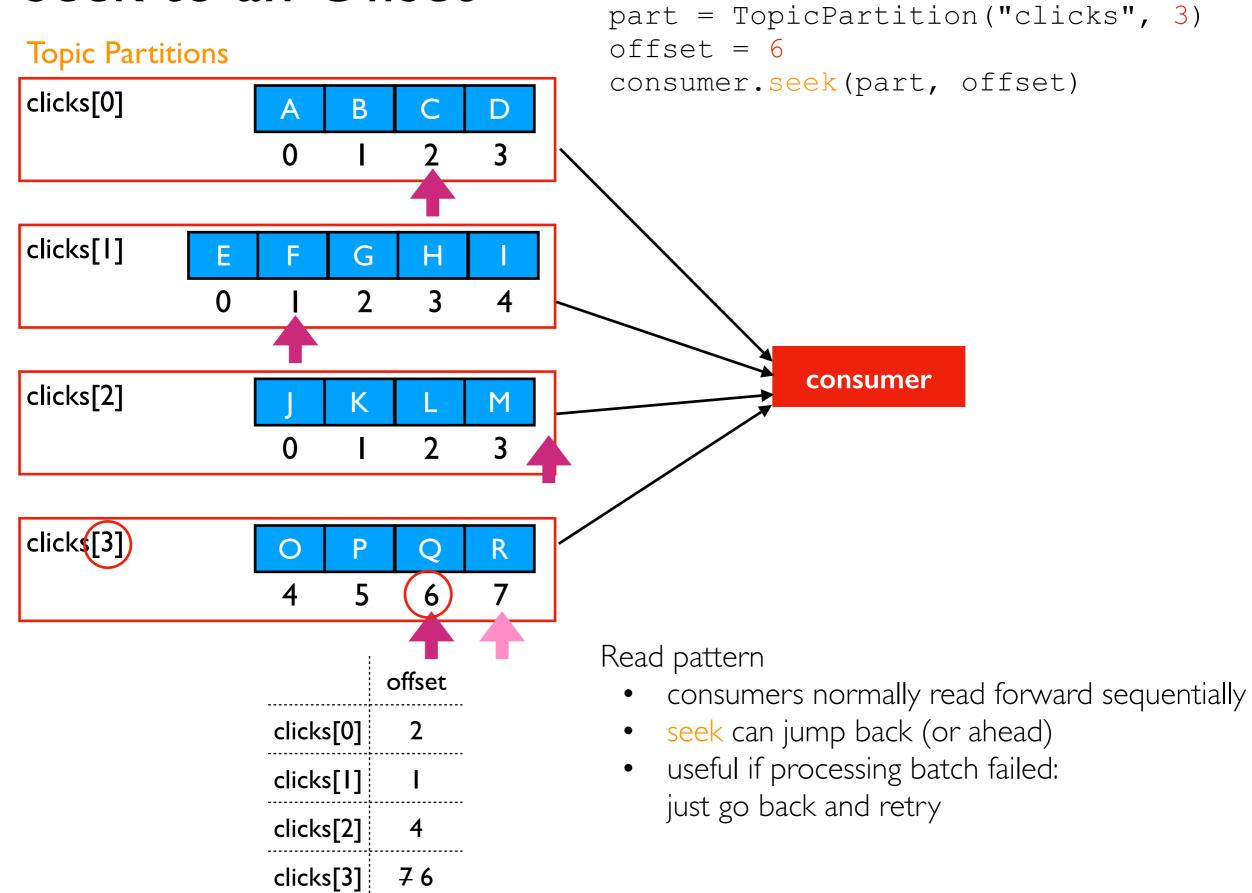
#### If A and B share the same topic and key, and B was produced after A, then:

- we say B "happened after" A
- A and B will be in the same partition (assuming partition count is constant)
- each consumer group of the topic will consume A before B

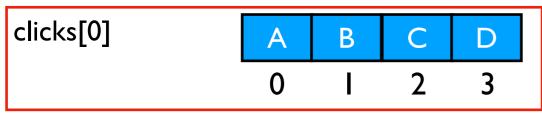
Choose your key carefully! Try to create enough partitions initially and never change it.

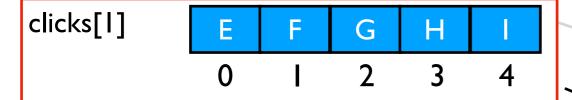
No keys specified => no guarantee about what order messages are consumed.

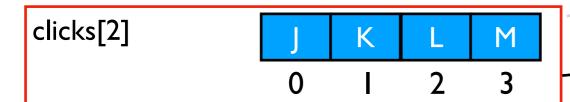
#### Seek to an Offset

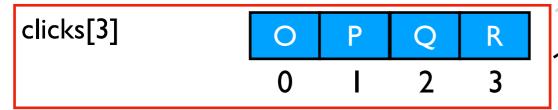


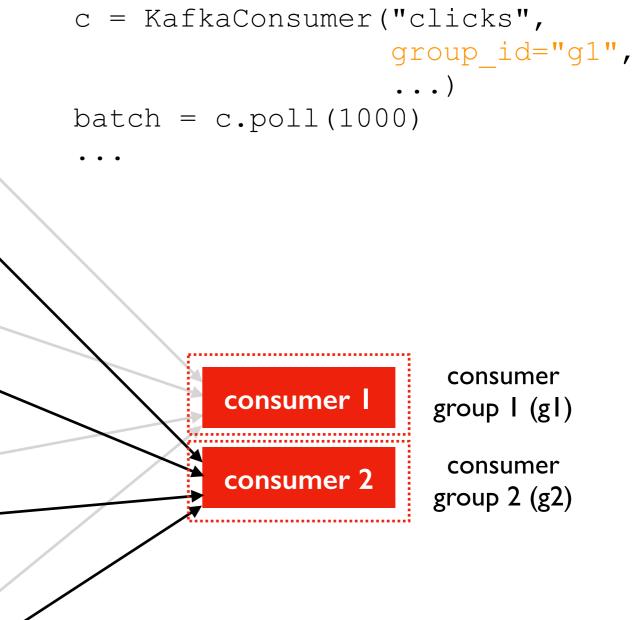
#### **Topic Partitions**









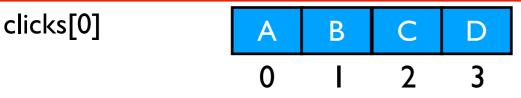


	g1 offsets	g2 offsets
clicks[0]	2	3
clicks[1]	l	2
clicks[2]	4	4
clicks[3]	3	3

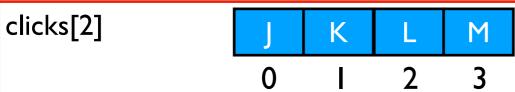
#### Groups

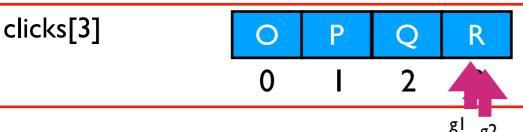
- different applications might operate independently
- they should ALL get a chance to consume messages
- need offsets for each topic/partition/consumer group combination

#### **Topic Partitions**







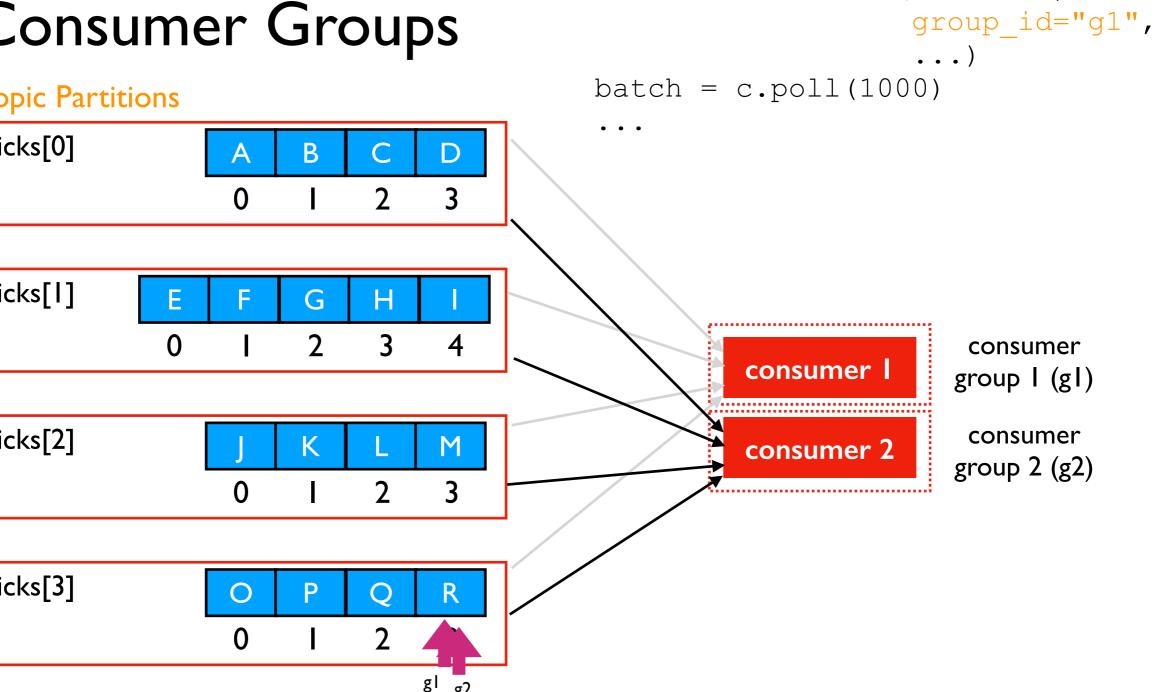


Grou	ıns

different applications might operate independently

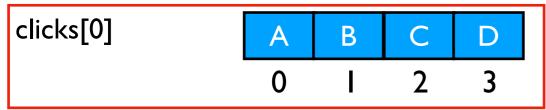
c = KafkaConsumer("clicks",

- they should ALL get a chance to consume messages
- need offsets for each topic/partition/consumer group combination

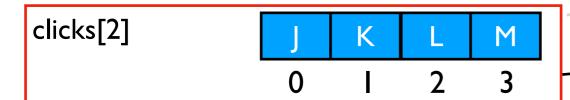


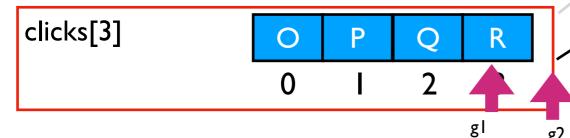
	gronsets	gz onsets
clicks[0]	2	3
clicks[1]	I	2
clicks[2]	4	4
clicks[3]	3	3

#### **Topic Partitions**









	g2 offsets
	3
- ;	

clicks[0]	2	3
clicks[1]	I	2
clicks[2]	4	4

clicks[3]

gl offsets

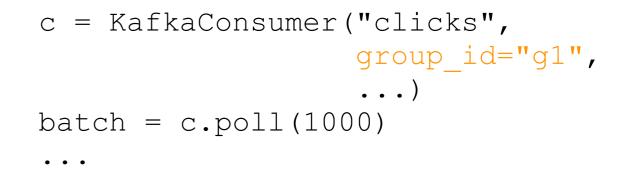


- different applications might operate independently
- they should ALL get a chance to consume messages
- need offsets for each topic/partition/consumer group combination

consumer

consumer 2

R



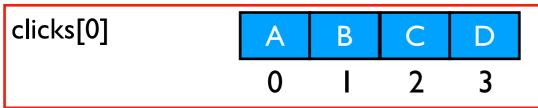
consumer

group I (gI)

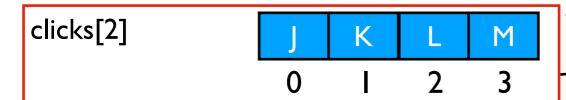
consumer

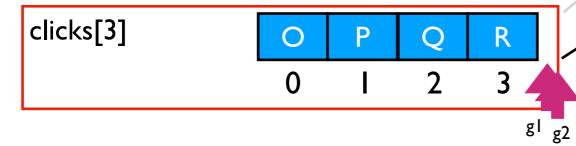
group 2 (g2)

#### **Topic Partitions**

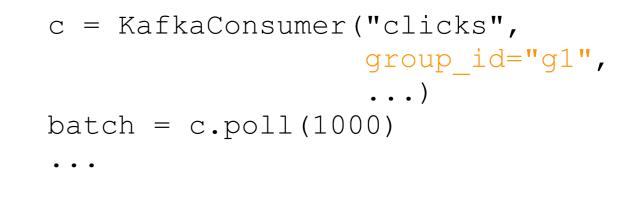








	g1 offsets	g2 offsets
clicks[0]	2	3
clicks[1]	I	2
clicks[2]	4	4
clicks[3]	4	4

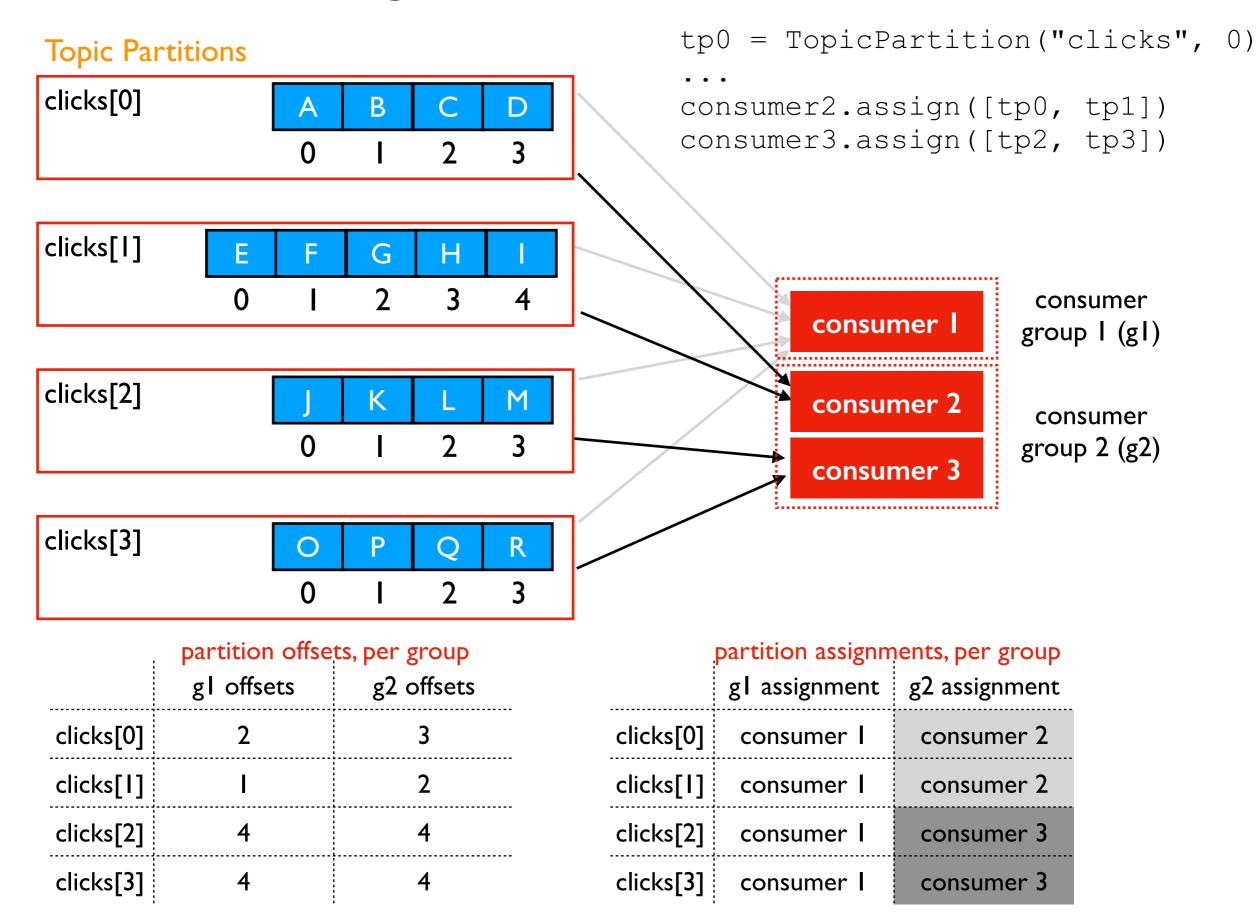




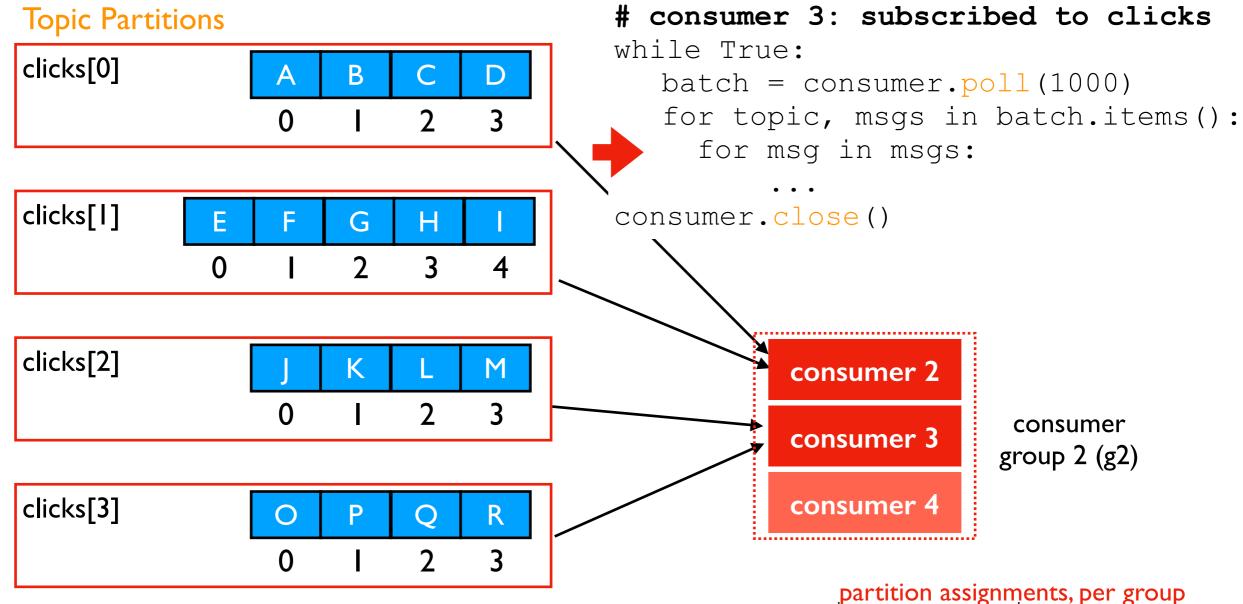
#### Groups

- different applications might operate independently
- they should ALL get a chance to consume messages
- need offsets for each topic/partition/consumer group combination

### Partition Assignment: Manual



### Partition Assignment: Automatic

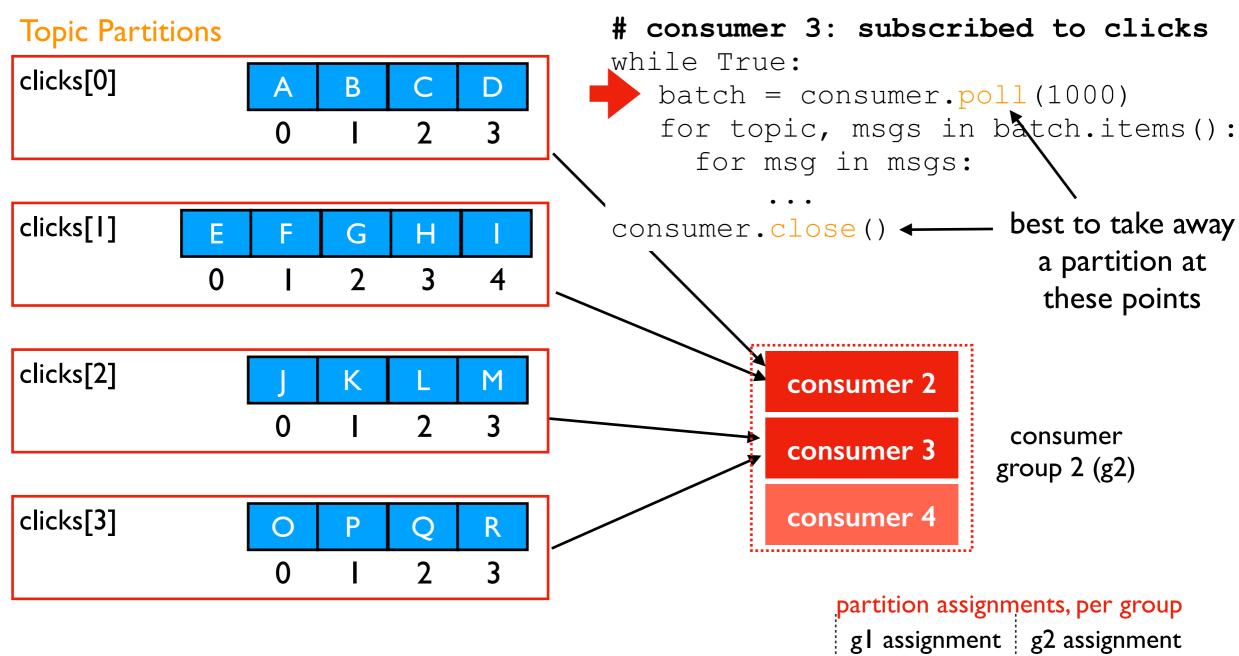


#### Assignment and re-assignment

- by default, consumers are automatically assigned partitions when they start polling
- **challenge:** Kafka shouldn't re-assign a partition in the ... middle of a batch (might double process messages)

	partition assignments, per group	
	g I assignment	g2 assignment
clicks[0]	consumer I	consumer 2
clicks[1]	consumer l	consumer 2
clicks[2]	consumer l	consumer 3
clicks[3]	consumer I	consumer 3

### Partition Assignment: Automatic

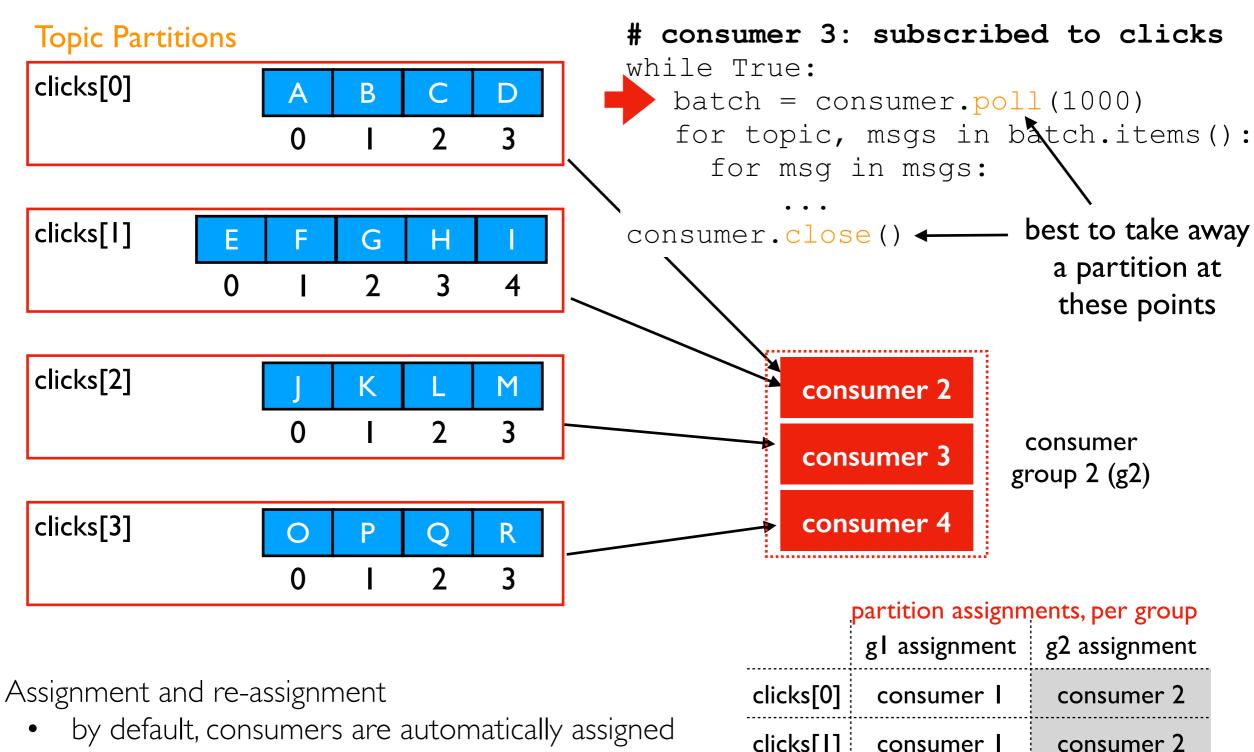


Assignment and re-assignment

- by default, consumers are automatically assigned partitions when they start polling
- **challenge:** Kafka shouldn't re-assign a partition in the ... middle of a batch (might double process messages)

	pai didon assignin	icitis, per group
	g lassignment	g2 assignment
clicks[0]	consumer l	consumer 2
clicks[1]	consumer l	consumer 2
clicks[2]	consumer l	consumer 3
clicks[3]	consumer I	consumer 3

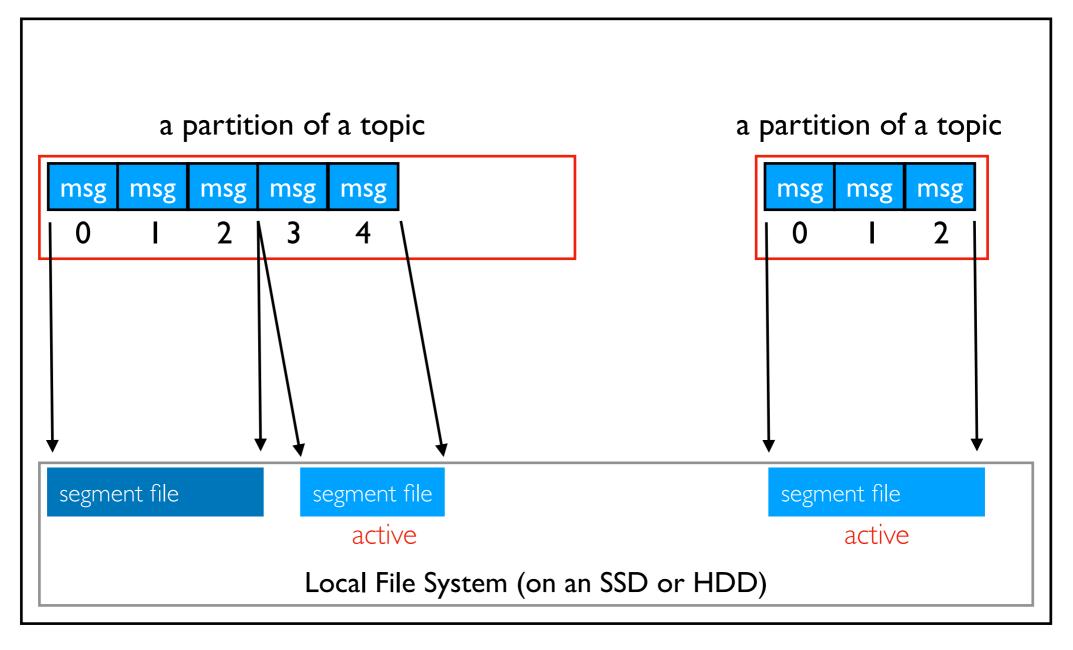
### Partition Assignment: Automatic



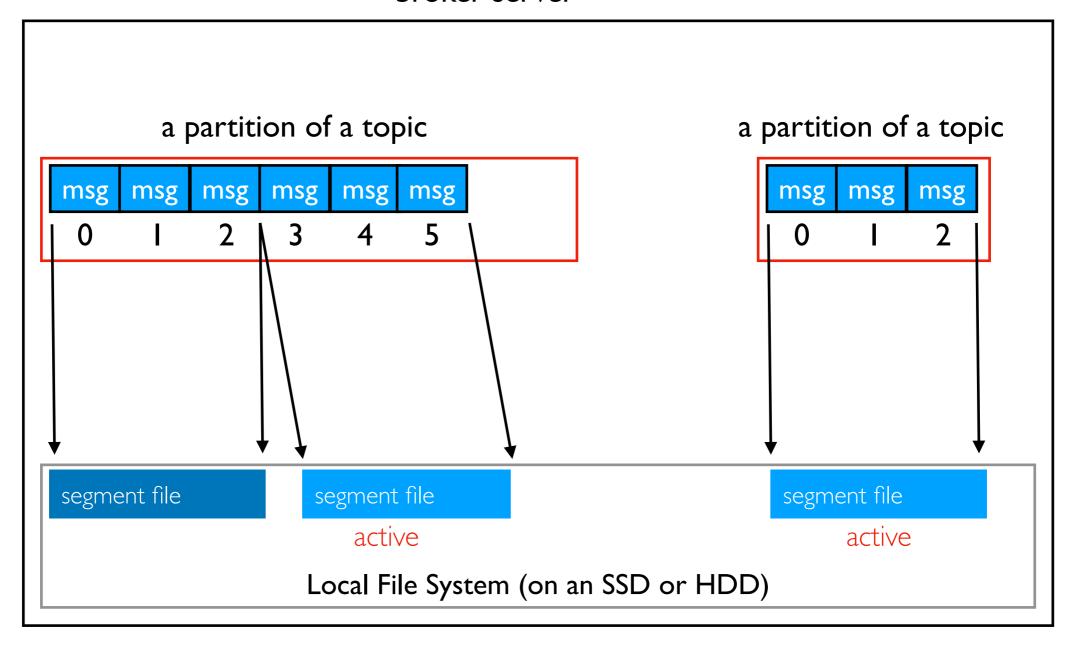
partitions when they start polling

challenge: Kafka shouldn't re-assign a partition in the middle of a batch (might double process messages)

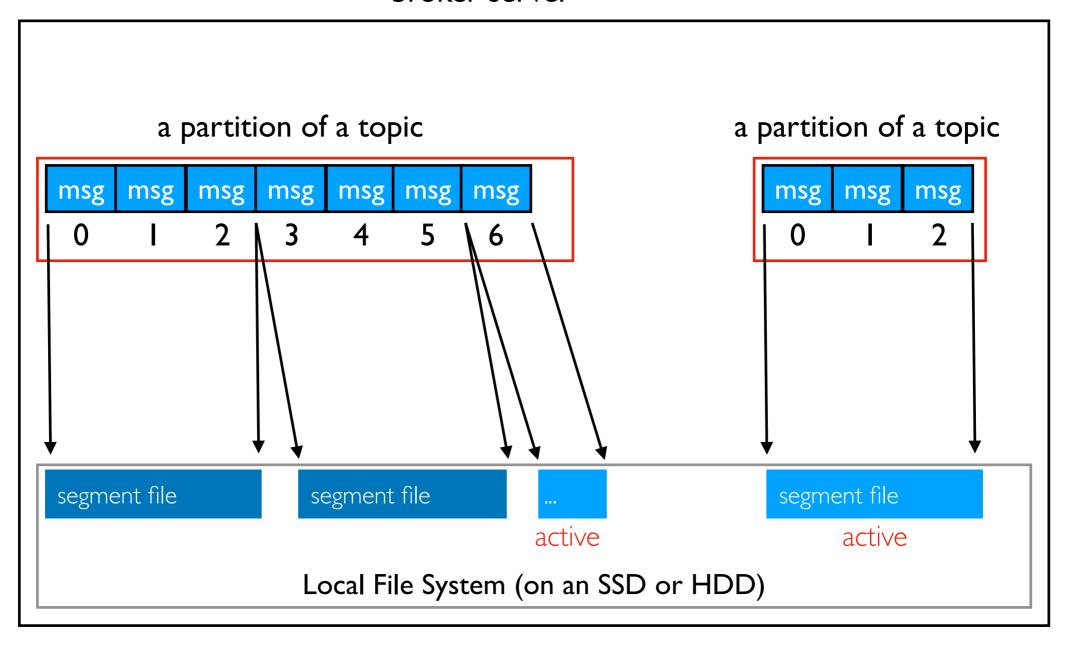
	g lassignment	g2 assignment
clicks[0]	consumer I	consumer 2
clicks[1]	consumer I	consumer 2
clicks[2]	consumer I	consumer 3
clicks[3]	consumer I	consumer 4



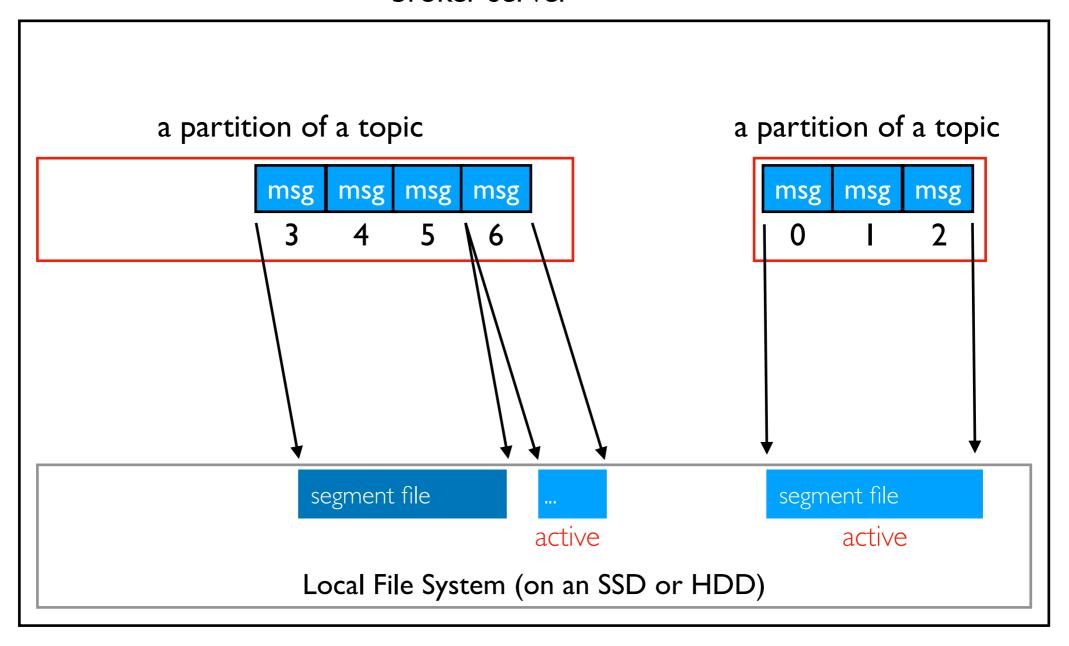
- partitions are divided into consecutive regions and saved in segment files
- all new data is sequentially written to the end of an active segment



- partitions are divided into consecutive regions and saved in segment files
- all new data is sequentially written to the end of an active segment



- rollover: current segment is finalized (no more changes)
- new segment is created and becomes active



- deletion: old segment is deleted
- always starts from smallest offset
- active segment is NEVER deleted

### Log Policy

Rollover and retention policies are configurable in Kafka.

#### Rollover

- setting I: max segment age (log.roll.hours=7 day by default)
- setting 2: max segment size (log.segment.bytes=IGB by default)
- rollover happens when segment gets too big or too old (whichever happens first)

#### Retention/Deletion

- setting I: log age cutoff (log.retention.hours=7 days by default)
- setting 2: log size cutoff (log.retention.bytes=disabled by default)
- deletion happens on oldest segment when log is too big or has records too old
- note: age cutoff applies to newest messages in a segment, so there will probably be some older ones in the same segment past the cutoff. Not useful for legal compliance with data retention laws.

### TopHat