# [544] Kafka Reliability

Meenakshi Syamkumar

### Learning Objectives

- describe how leader and follower replicas work in Kafka (to record messages, handle failover, etc.)
- apply the definition of "committed" messages to reason about when messages (a) are acknowledged to producers and (b) can be read by consumers
- configure Kafka and write producer/consumer code to achieve "exactly once semantics"

### Outline: Kafka Reliability

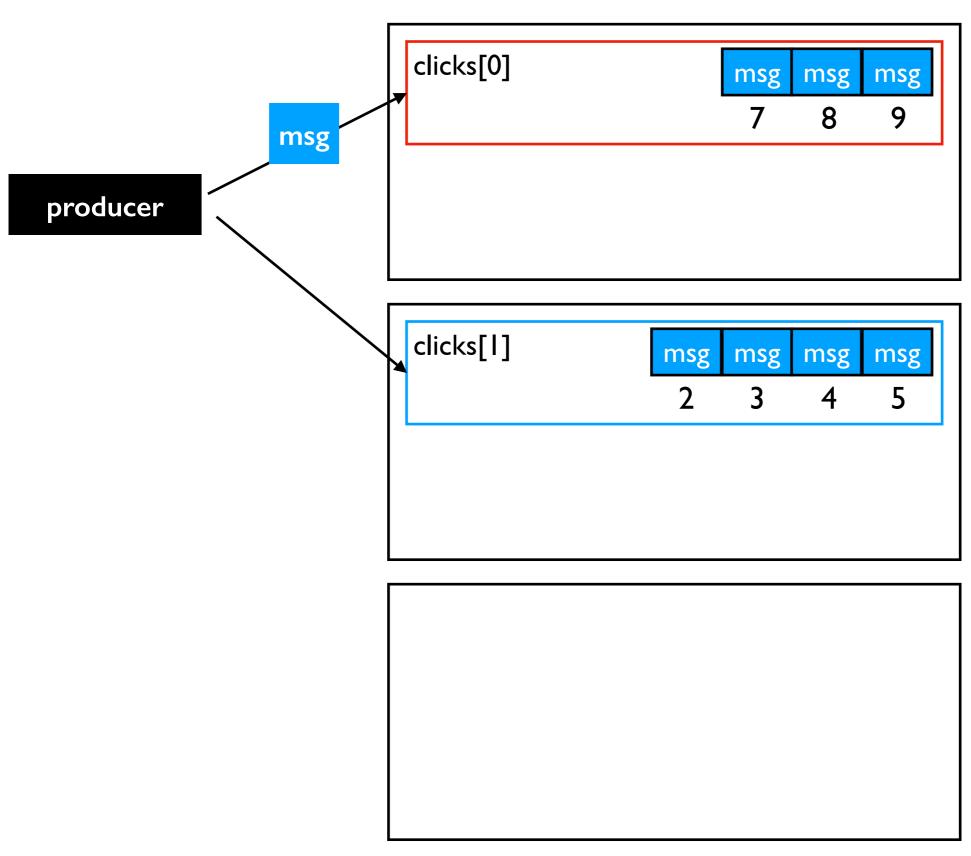
Kafka Replication

Fault Tolerance

**Exactly-Once Semantics** 

### Three brokers, 2 partitions, replication factor=1

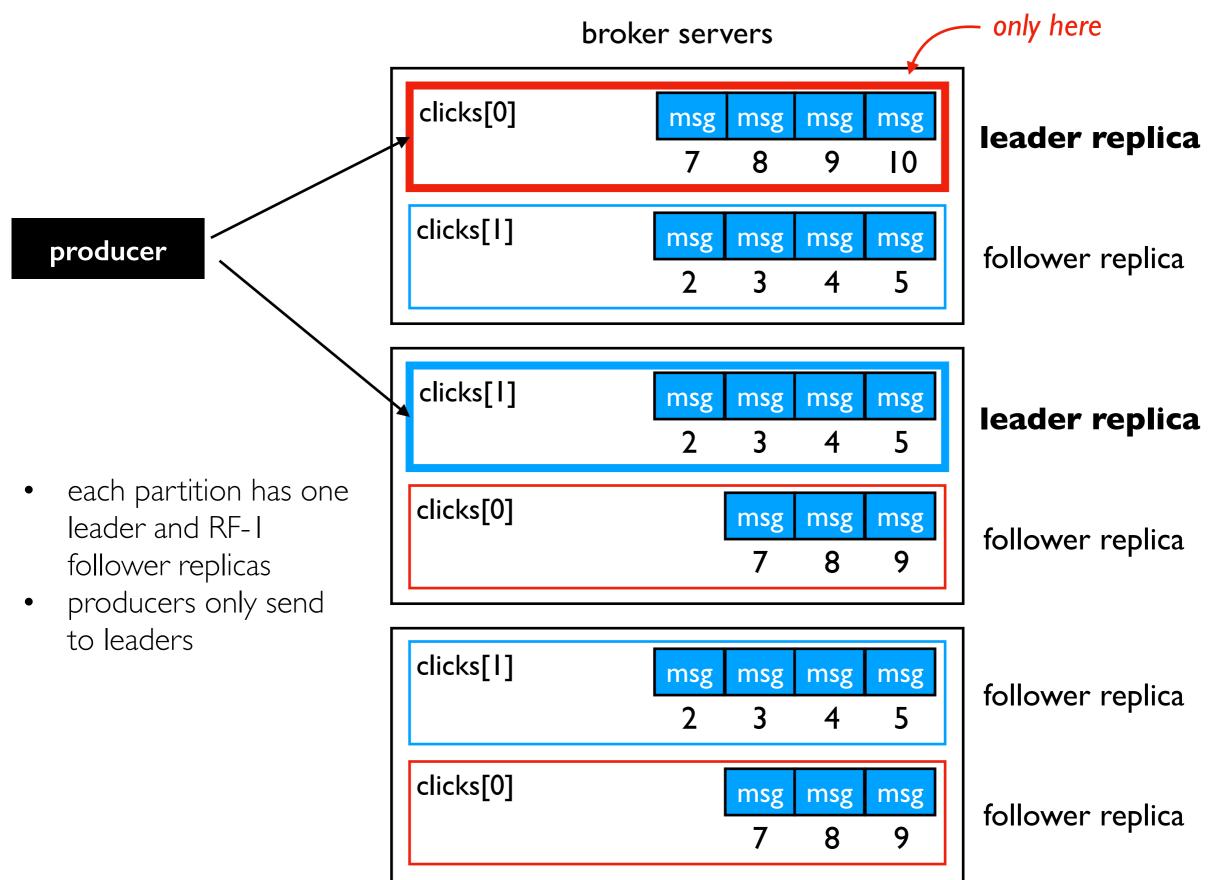
broker servers



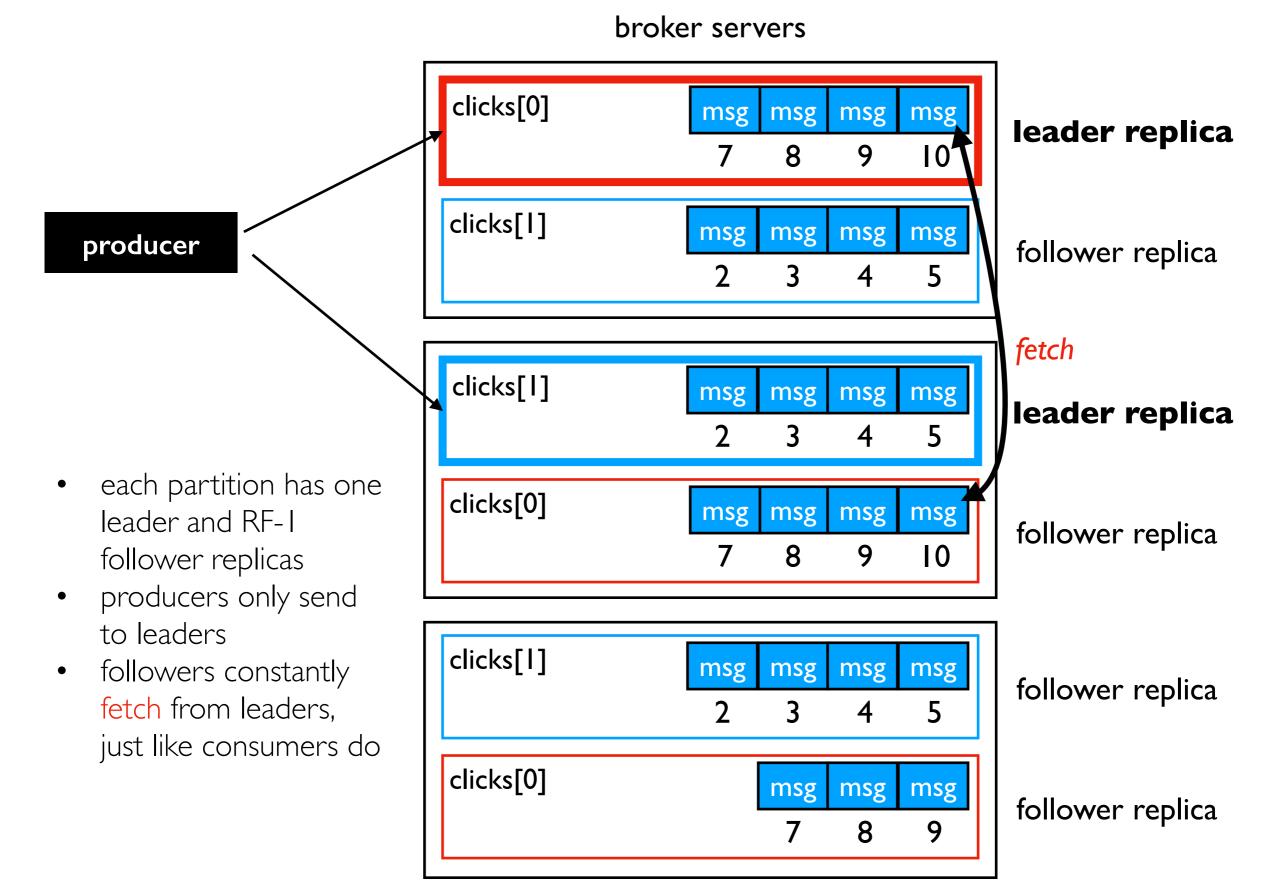
### Three brokers, 2 partitions, replication factor=3

broker servers clicks[0] msg msg msg leader replica 8 9 msg clicks[1] msg msg msg msg producer follower replica 2 3 4 5 clicks[1] msg msg msg msg leader replica 3 each partition has one clicks[0] msg msg msg leader and RF-I follower replica 8 9 follower replicas clicks[1] msg msg msg msg follower replica 3 4 5 clicks[0] msg msg msg follower replica 8

### Three brokers, 2 partitions, replication factor=3

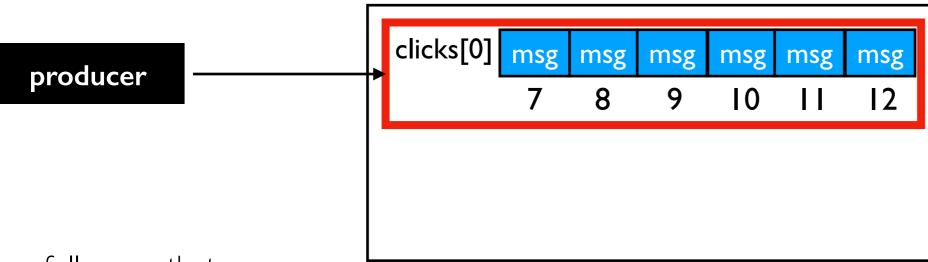


### Fetch Requests



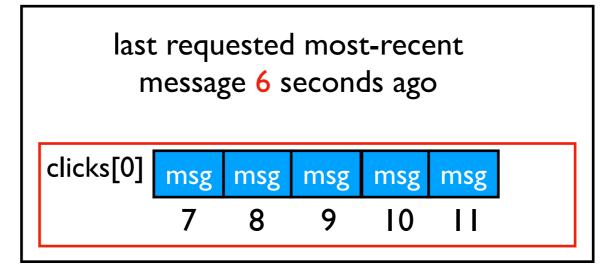
# Followers: In-Sync and/or Lagging Behind

broker servers



leader replica

- followers that are "keeping up" with leader messages are called "in-sync"
- definition is tunable, and depends on factors like how recently a follower got a batch with the most recent messages
- some flexibility: insync followers might be a little behind the leader



follower replica (in-sync)

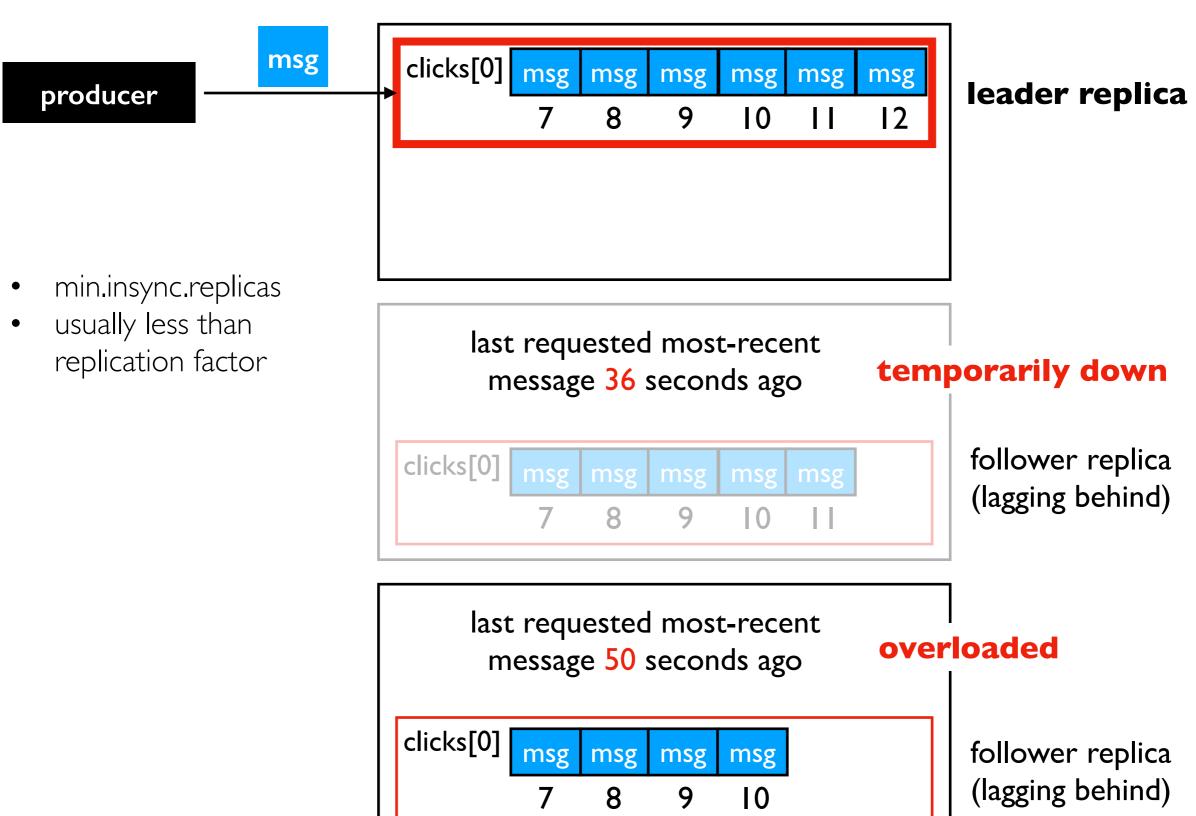
last requested most-recent message 20 seconds ago

clicks[0] msg msg 7 8

follower replica (lagging behind)

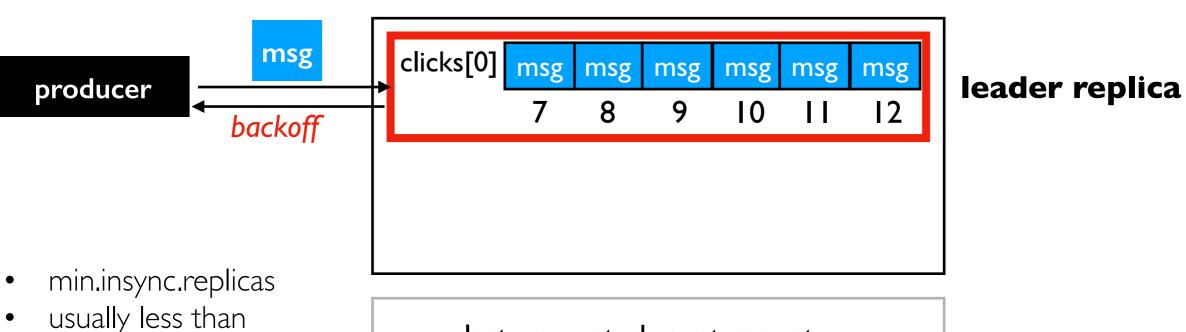
### Minimum In-Sync Replicas (Assume 2 Here)

broker servers



### Backoff: Not Enough Replicas Exception

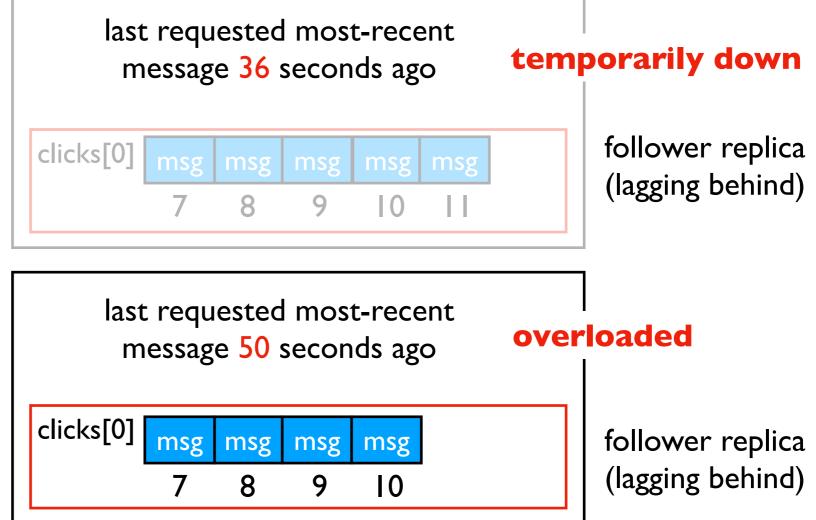
broker servers



 reject some messages until we have enough responsive in-sync replicas: NotEnough-ReplicasException

replication factor

- bigger min: stronger durability
- smaller min: better write availability



### Outline: Kafka Reliability

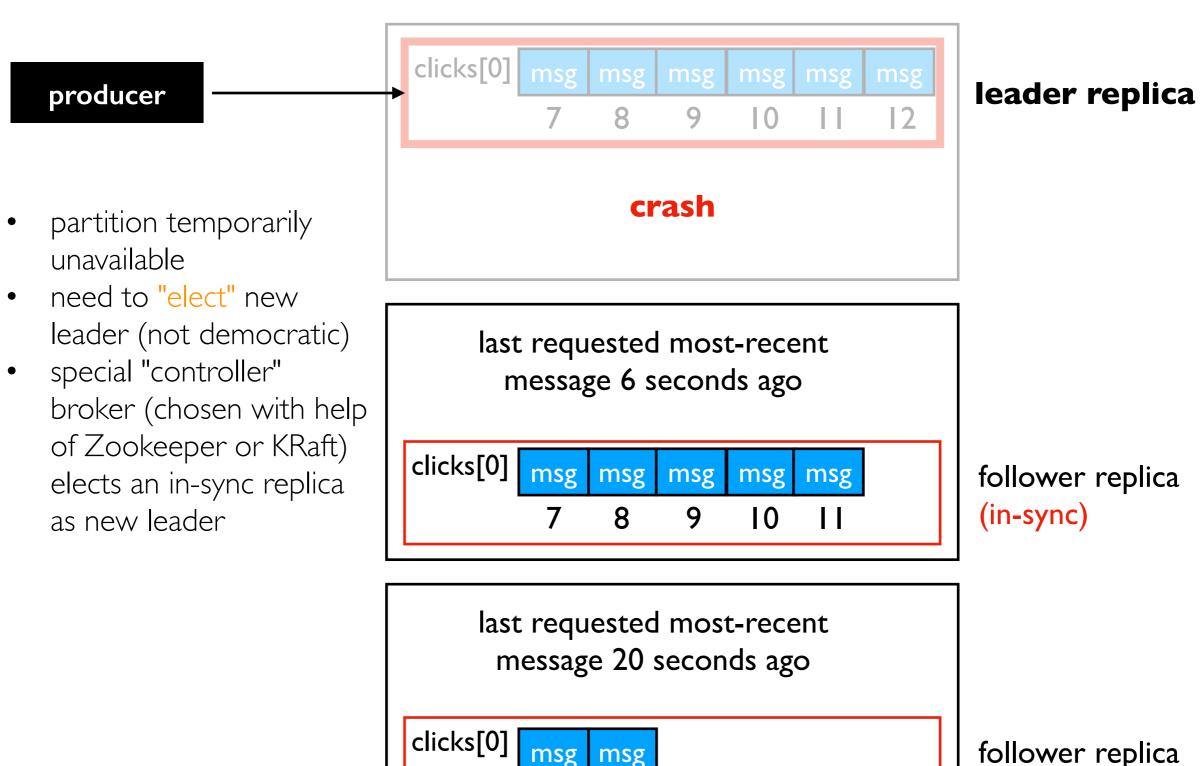
Kafka Replication

Fault Tolerance

**Exactly-Once Semantics** 

### What if the leader fails? Elect a new one!

broker servers

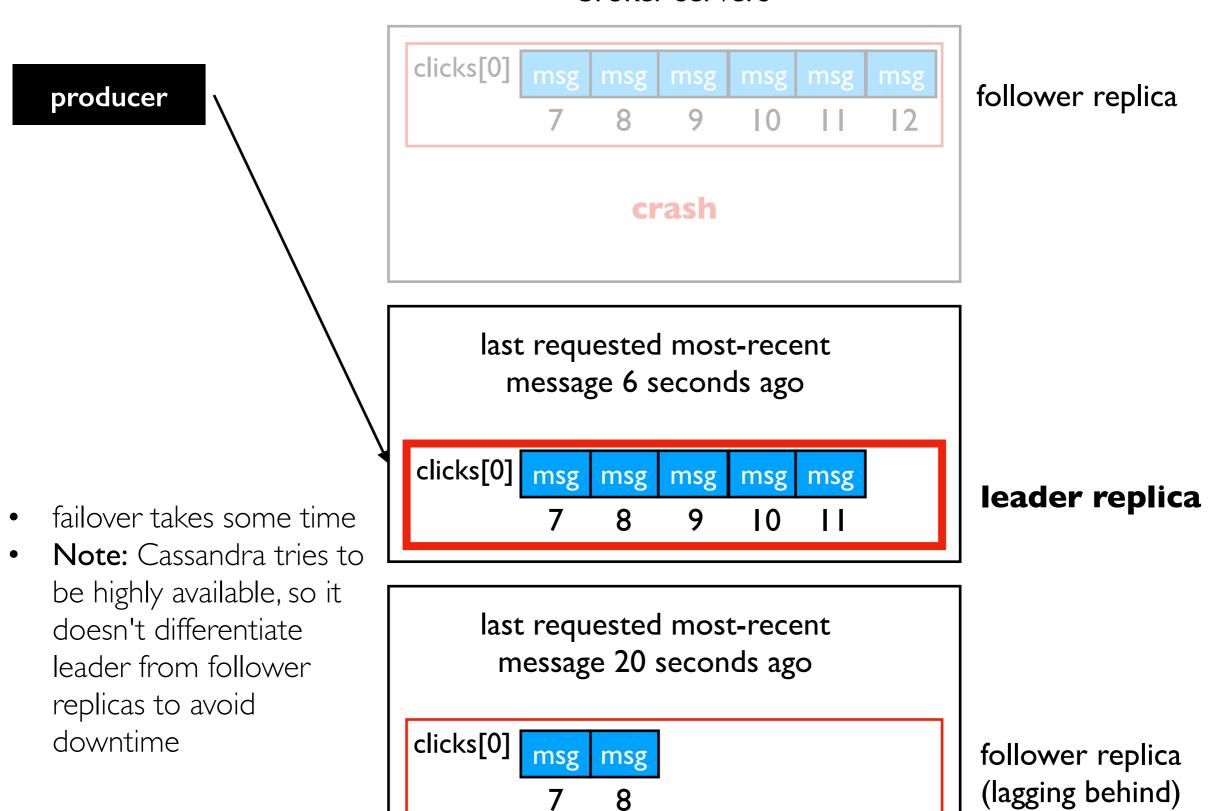


8

follower replica (lagging behind)

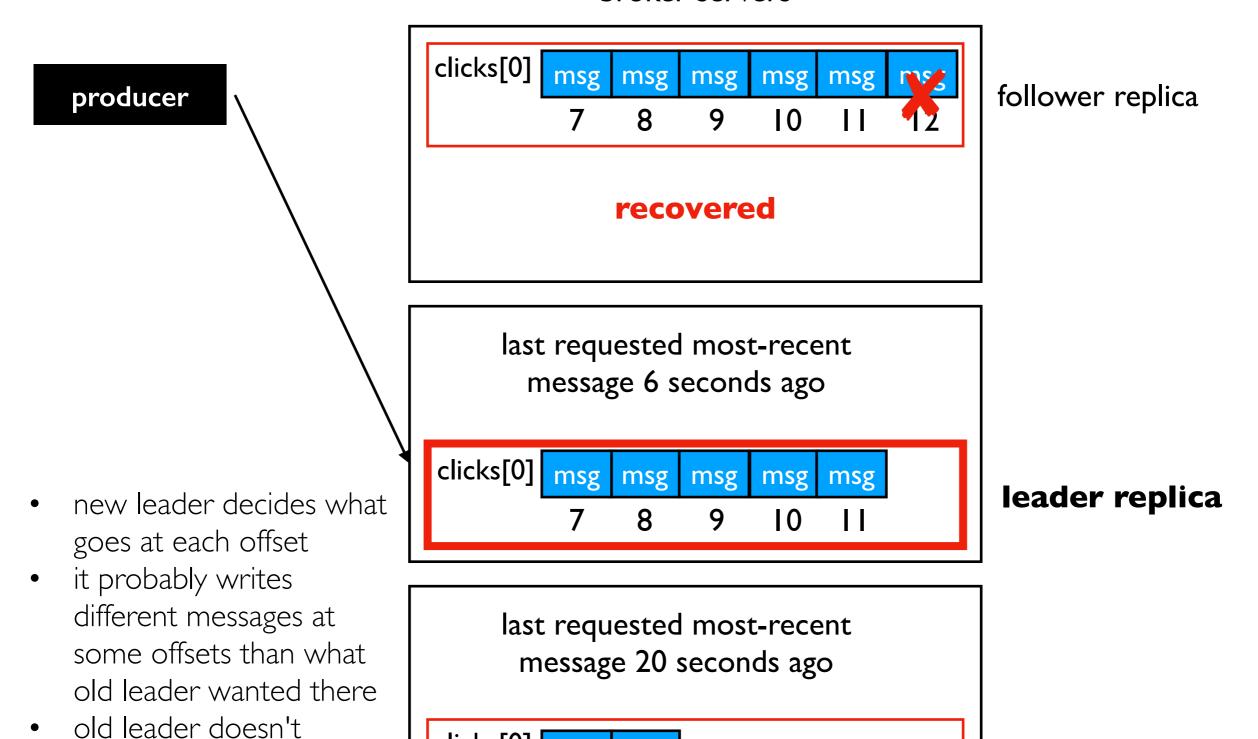
### Kafka Replica Failover

broker servers



### Some Messages Seen by Old Leader Lost

broker servers



clicks[0]

immediately get its job

back upon recovery

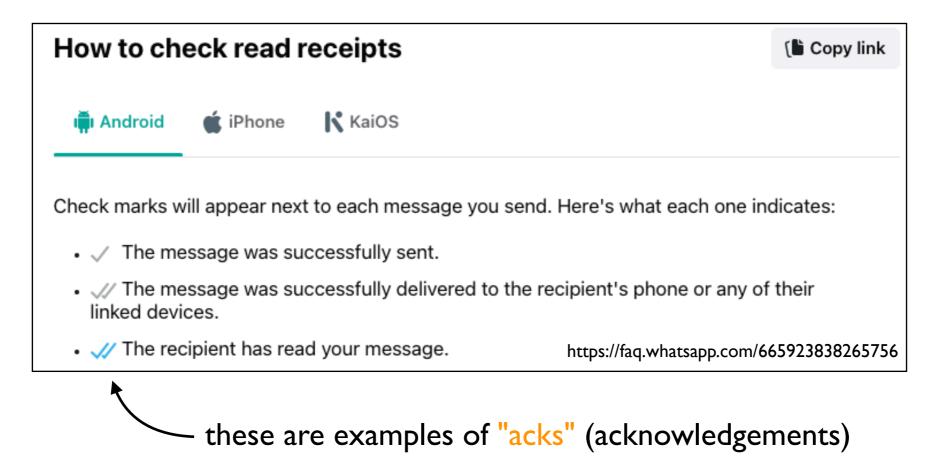
msg

msg

8

follower replica (lagging behind)

### Review "Committed": WhatsApp Acks Example



In distributed storage systems/databases, an ack means our data is committed.

"Committed" means our data is "safe", even if bad things happen. The definition varies system to system, based on what bad things are considered. For example:

- a node could hang until rebooted; a node's disk could permanently fail
- a rack could lose power; a data center could be destroyed

In Kafka's leader/follower replica design, what are some "bad things" we might worry about?

### Kafka: Committed Messages

Messages are "committed" when written to ALL in-sync replicas.

Depending on how many are in-sync, the strength of the guarantee can vary, but min.insync.replicas lets us specify a worst case.

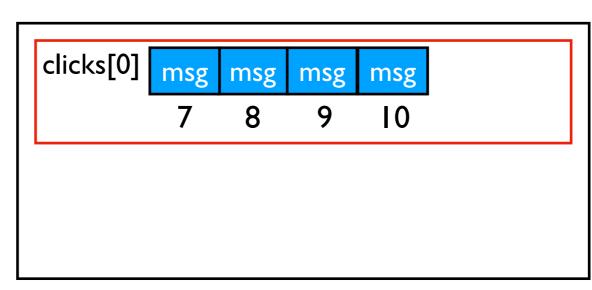
If number of concurrent broker failures < min.insync.replicas, then our committed data is safe, even if the leader fails (because we can elect another in-sync replica, and all in-sync replicas have all committed data).

### Committed Messages

broker servers

What is committed?

- assume RF=3 and minimum in-sync=2
- is message 8 committed?
- message 10?
- message | 1?



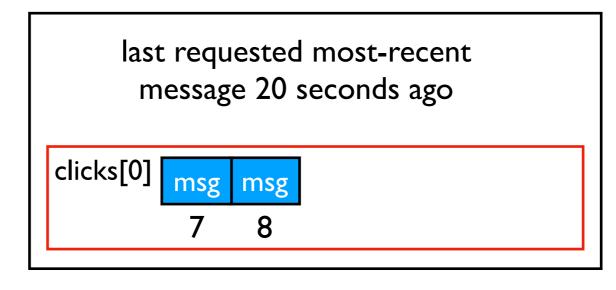
follower replica (in-sync)

last requested most-recent message 6 seconds ago

clicks[0] msg msg msg msg msg
7 8 9 10 11

leader replica

**TopHat** 



follower replica (lagging behind)

### Working with Committed Data

How can we avoid "anomalies" (unexpected system behavior) by taking advantage of committed data?

### Example 1:Write Anomaly

#### Scenario:

- producer writes a message
- produce receives an ACK (acknowledgement) from the broker
- consumers never see the message

**Cause**: maybe the leader sent an ACK back, then crashed, before replicating the message to the followers.

How to avoid it? Use strong acks.

#### Consumer initialization:

- KafkaProducer(..., acks=0)
   don't wait for leader to send back ACK
- KafkaProducer(..., acks=1)

  ACK after leader writes to its own log
- KafkaProducer(..., acks="all")
   ACK after data is committed (slowest but strongest)

If you don't get an ACK that data is committed, usually best to retry in a loop (Kafka can be configured to do this for you).

### Example 2: Read Anomaly

#### Scenario:

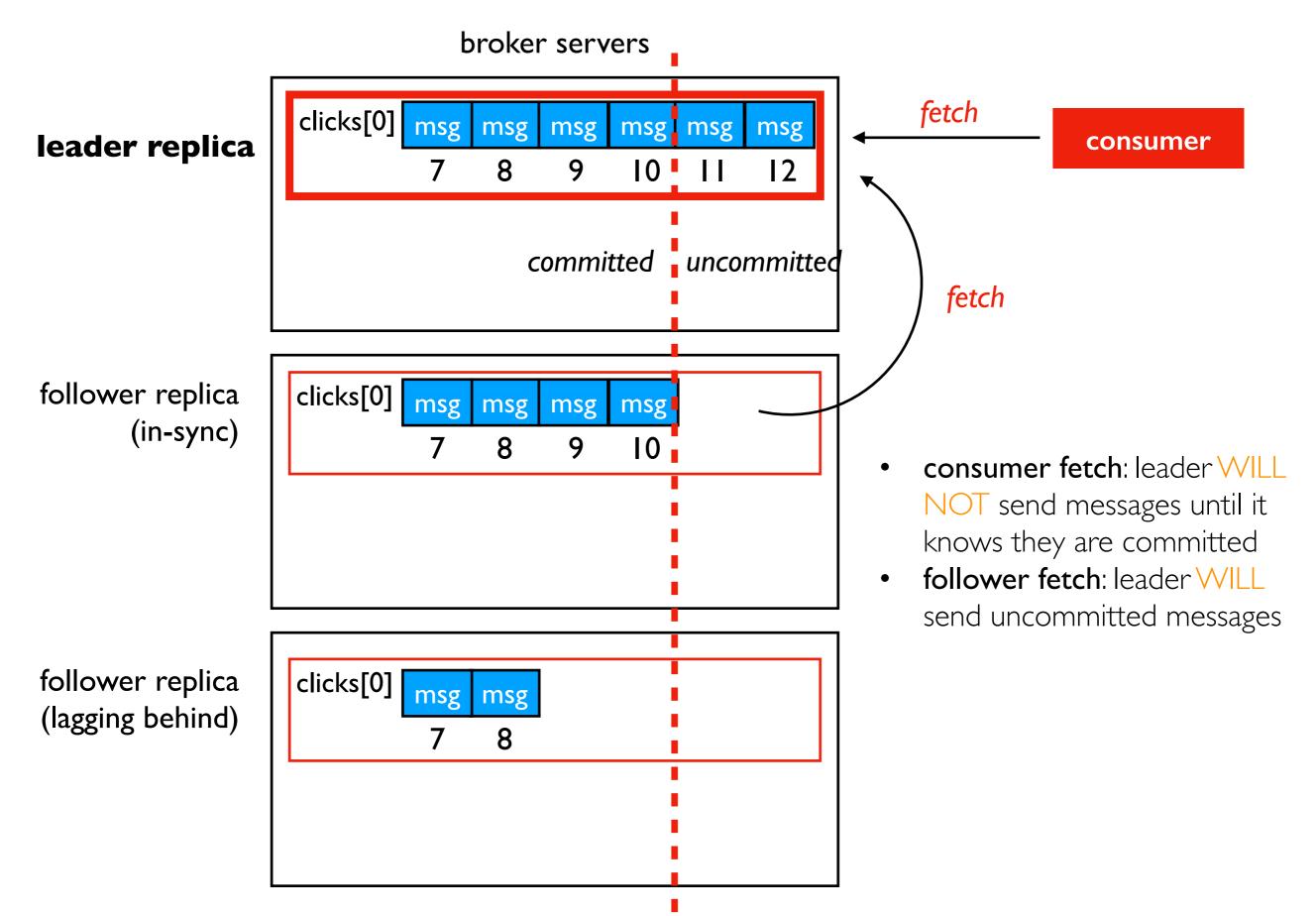
- a consumer reads a message
- there is an attempt to read the message again later (same consumer, or other)
- message is gone, or changed

**Cause**: maybe the message was consumed from the leader before it was replicated to the followers; then the leader crashed and the new leader doesn't have that message for future consumption.

How to avoid it? Never read un-committed data.

The leader does this automatically.

### Fetch Behavior: Consumer vs. Follower



### Outline: Kafka Reliability

Kafka Replication

Fault Tolerance

**Exactly-Once Semantics** 

# Semantics (Meaning)

#### Dictionary

Definitions from Oxford Languages · Learn more



#### se·man·tics

noun

noun: semantics; noun: logical semantics; noun: lexical semantics

the branch of linguistics and logic concerned with meaning. There are a number of branches and

#### Programming Example:

- Runtime bug: the program crashed, there was clearly a problem
- Semantic bug: you need to understand the meaning of the results to say whether or not the program behaved correctly

#### In Systems:

- what does it mean when we get we get an ACK, or a write returns?
- the meaning depends on how we configured things...

### At-most-once semantics

```
producer = KafkaProducer(..., acks=1)
producer.send("my-topic", b"some-value")
```

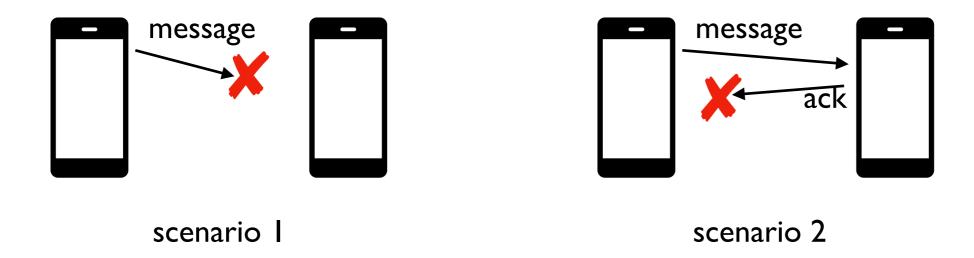
With acks as 0 or 1 and no retry, a successful write means the data was recorded at most once (ideally once, but if the leader crashes at a bad time, maybe zero times).

### Using strong ACKs and retry

```
producer = KafkaProducer(..., acks="all", retries=10)
producer.send("my-topic", b"some-value")
```

Keep retrying until success (within reason -- for example, 10 times)

**Problem:** there are two reasons we might not get an ACK:

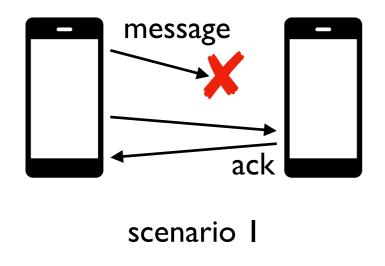


### Using strong ACKs and retry

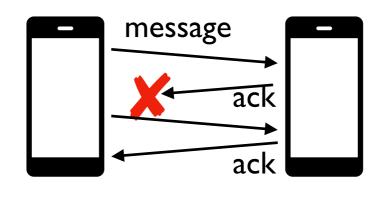
```
producer = KafkaProducer(..., acks="all", retries=10)
producer.send("my-topic", b"some-value")
```

Keep retrying until success (within reason -- for example, 10 times)

**Problem:** there are two reasons we might not get an ACK:







scenario 2

message written twice

A strong ACK with retry provides at-least-once semantics because we're guaranteed I+ messages upon success

### Are duplicate messages OK?

#### Yes, if they're idempotent.

"An operation is called idempotent when the effect of performing the operation multiple times is equivalent to the effect of performing the operation a single time" ~ Operating Systems: Three Easy Pieces, by Arpaci-Dusseau

```
x = 0
y = 0

def set_x(value):
    global x
    x = value

def inc_y(value):
    global y
    y += value
```

```
# if we just do once, is it the same?
set_x(123)
set_x(123)

# if we just do once, is it the same?
inc_y(3)
inc_y(3)
inc_y(3)
```

### TopHat

### Suppressing Duplicates

With some cleverness, we can make ANYTHING idempotent.

```
\Delta = 0
completed ops = set()
def inc y (value, operation id):
    global y
    if not operation id in completed ops:
        y += value
        completed ops.add(operation id)
inc y(3, 1251253)
inc y(3, 1251253) # no effect
inc y(3, 1251253) # no effect
inc y(3, 9876)
               # no effect
inc y(3, 9876)
inc y(1, 5454)
```

### Exactly-Once Semantics: Producer Side

Upon a successful write, the message will be considered exactly once (duplicates will be suppressed by brokers or consumers).

#### Producer settings:

- acks="all"
- retry=N
- enable.idempotence=True

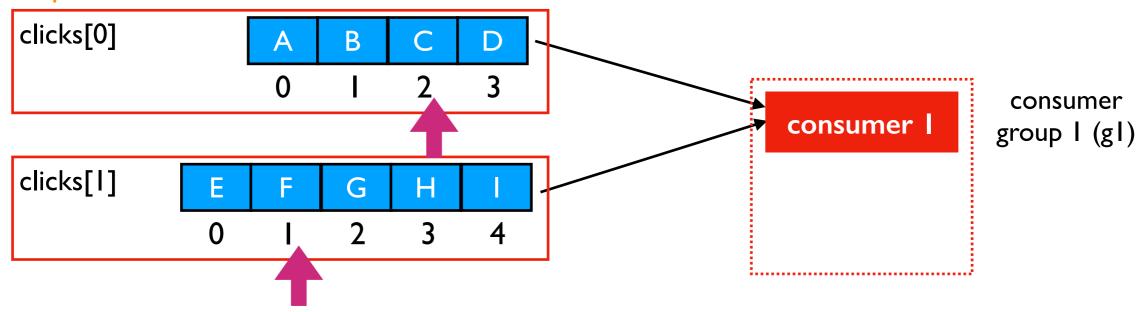
With idempontence enabled, producers automatically generate unique operation IDs and brokers suppress duplicates (this has an extra cost).

You can use enable.idempotence in Java, but the kafka-python package doesn't support it.:

- need to handle it yourself
- often, messages have a unique ID anyway, so consumers can ignore dups
- Example: weather stations that emit one record per day -- if a consumer sees a date for a station it has seen before, ignore it

### Exactly-Once Semantics: Consumer Side

#### Topic Partitions

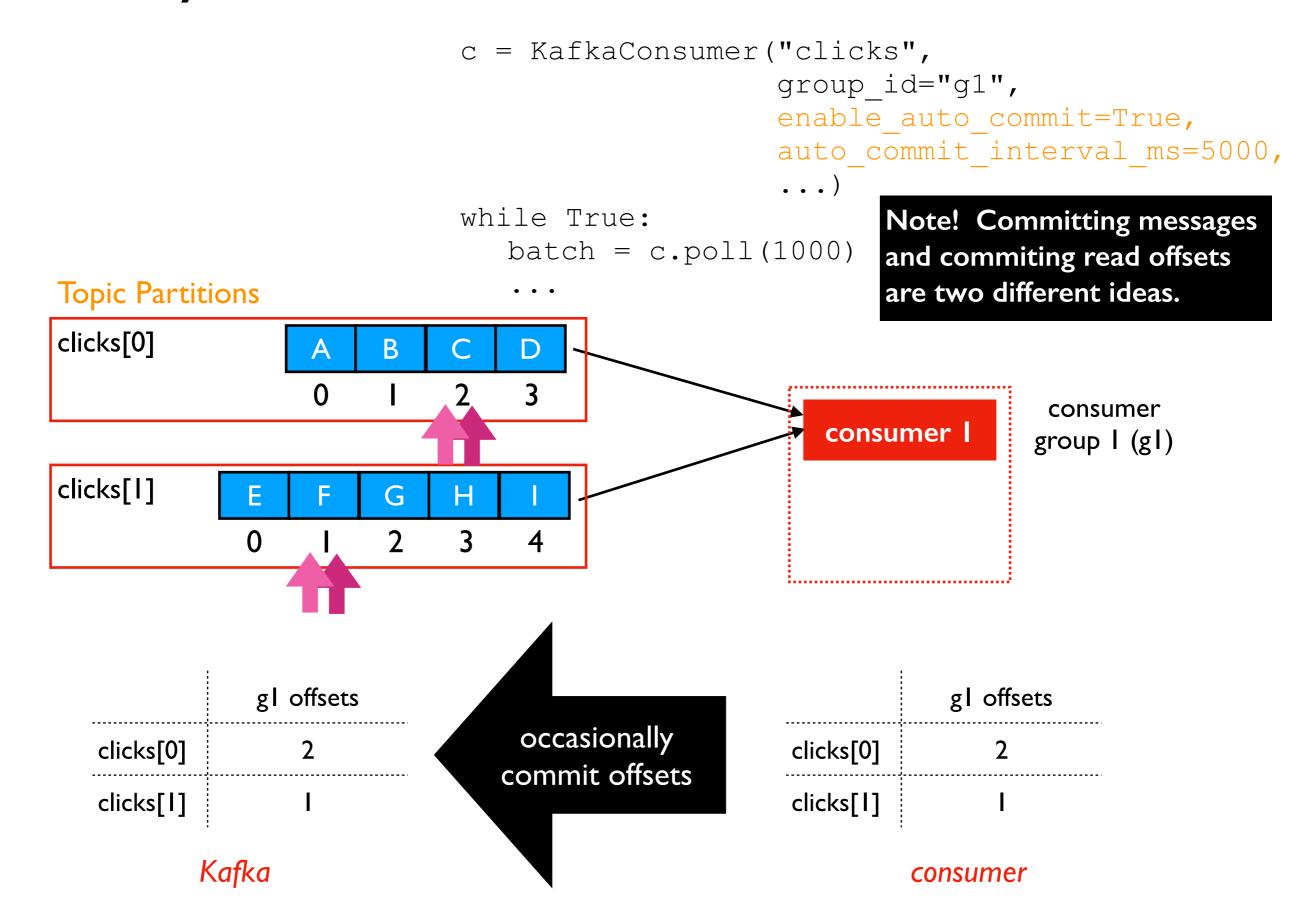


# Suppose consumer dies and is replaced by another in the same group

- don't want replacement to miss any messages
- don't want replacement to repeat any processing

	g1 offsets
clicks[0]	2
clicks[1]	I

### Exactly-Once Semantics: Consumer Side



### Exactly-Once Semantics: Consumer Side

Kafka

```
c = KafkaConsumer("clicks",
                                                     group id="g1",
                                                     enable auto commit=True,
                                                     auto commit interval ms=5000,
                                                     . . . )
                             while True:
                                 batch = c.poll(1000)
Topic Partitions
clicks[0]
                  0
                                                                         consumer
                                                         consumer
                                                                        group I (gI)
clicks[1]
                        G
                        2
             0
                                   4
              gl offsets
                                                                  gl offsets
                        If we crash at a bad time, the
                        offsets the next consumer
                                                      clicks[0]
  clicks[0]
                        gets from Kafka will only be
  clicks[1]
                                                      clicks[1]
                        approximately correct.
```

consumer

### Approach I: Manually Commit Offsets

```
c = KafkaConsumer("clicks",
                                                  group id="g1",
                                                  enable_auto_commit=False,
                                                  ...)
                            while True:
                               batch = c.poll(1000)
                               c.commit() # manually commit read offsets
Topic Partitions
clicks[0]
                 0
                                                                     consumer
                                                     consumer I
                                                                    group I (gl)
clicks[1]
                      G
                       2
             0
                                4
             gl offsets
                                                              gl offsets
  clicks[0]
                                                   clicks[0]
```

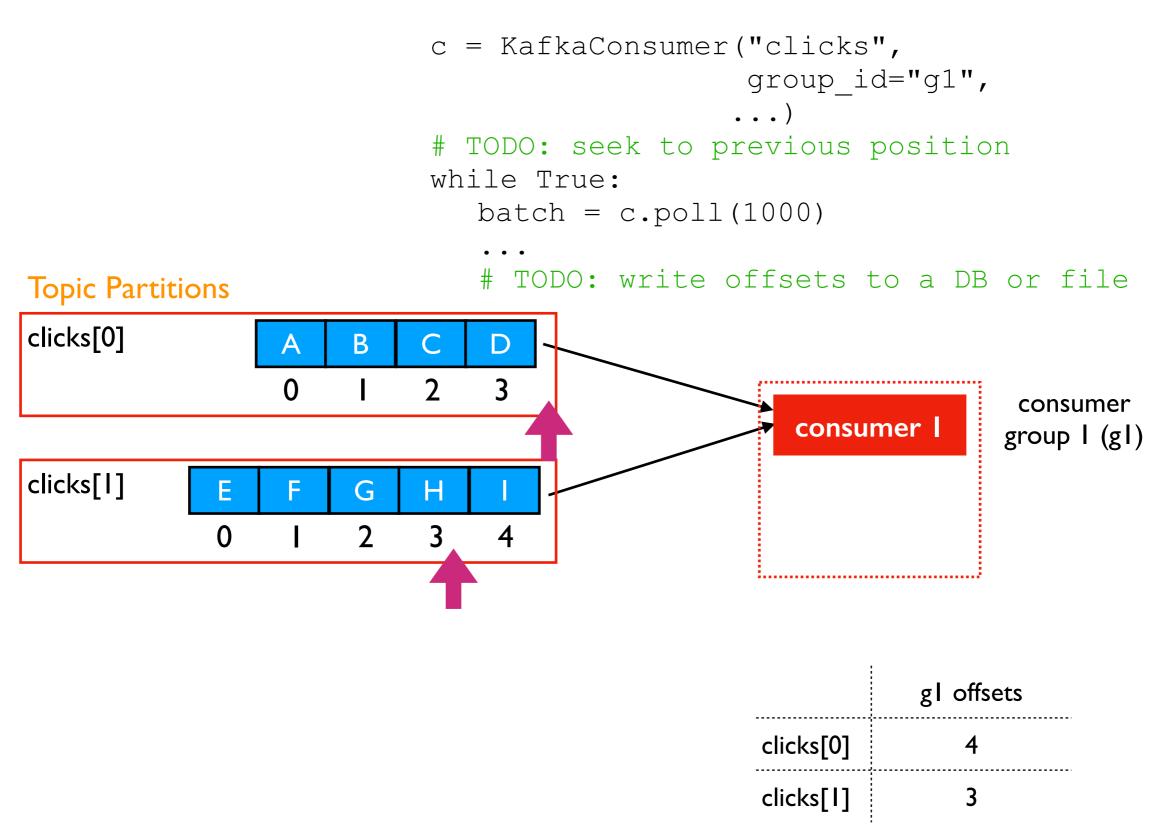
Kafka

clicks[1]

consumer

clicks[1]

# Approach 2: Externally Save Commits



consumer

### Conclusion

Every part of the system has a part to play in reliability and exactly-once semantics.

#### **Producer:**

- requesting strong acks
- retry
- idempotence

#### Broker:

- replicating data to followers
- failing over to new leader
- sending acks
- helping producer suppress duplicates
- keeping uncommitted data hidden from consumers

#### Consumer:

- carefully handling read offsets
- sometimes suppressing duplicates (if not handled by producers+brokers)