# Fault Tolerance and High Availability
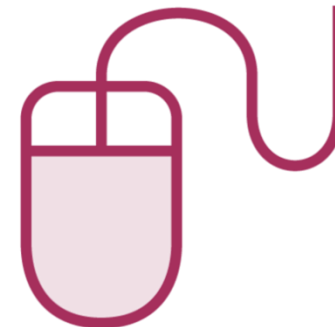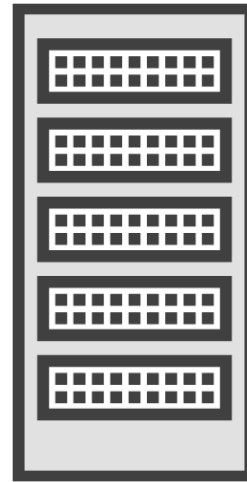
**Paweł Kordek**

DATA ENGINEER

@pawel_kordek   https://kordek.github.io

**Using a single broker means accepting:**

- Message loss
- Downtime

# Desired Properties

**Fault Tolerance**

**High Availability**

**Consistency**

# ZooKeeper

Source of truth about cluster members

Point of contact for new brokers

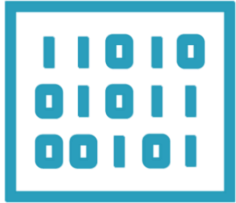Single ensemble can be shared with other applications

# Message Persistence

Producer

"acks" = "none"

topic-partition-0

Producer

ACK

topic-partition-0

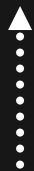"acks" = "1"

# Producer Delivery Guarantees

```java
Properties props = new Properties();

. . .

KafkaProducer<String, String> = new KafkaProducer<>(props);

ProducerRecord<String, String> r = new ProducerRecord("topic", "key", "value");

RecordMetadata rm = producer.send(r).get();
```
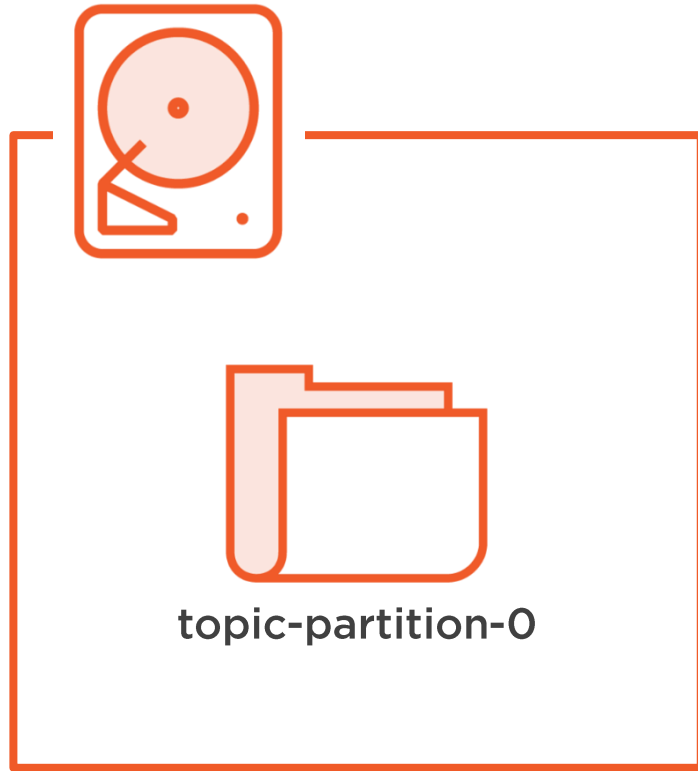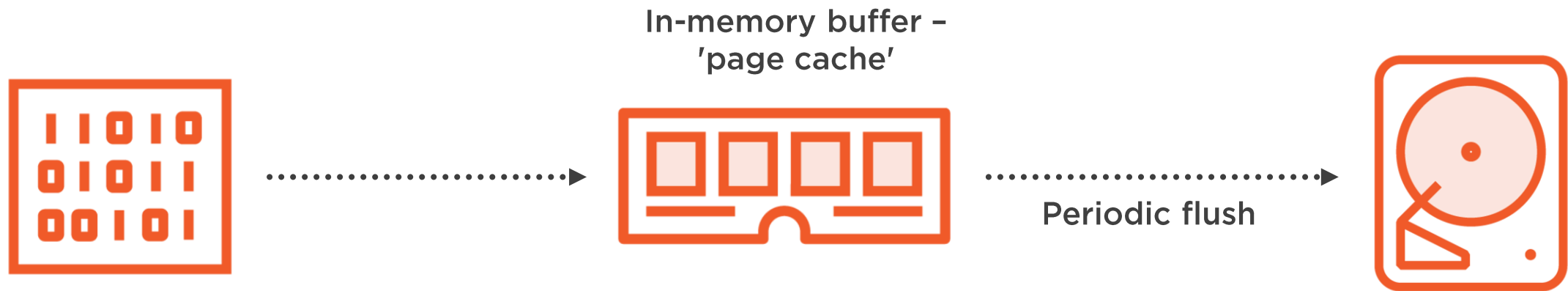
**When it returns depends on 'acks' value**

topic-partition-0

Only after a successful write, data is committed and made available to consumers.

# Writing to Disk

**In-memory buffer –
'page cache'**

**Periodic flush**

# Replication

# Creating a New Topic

```java
Admin admin = Admin.create(
        Map.of("bootstrap.servers", "localhost:9092")
);

NewTopic topic = new NewTopic("quote-feedback", 2, (short) 1)

admin.createTopics(List.of(topic)).get()
```
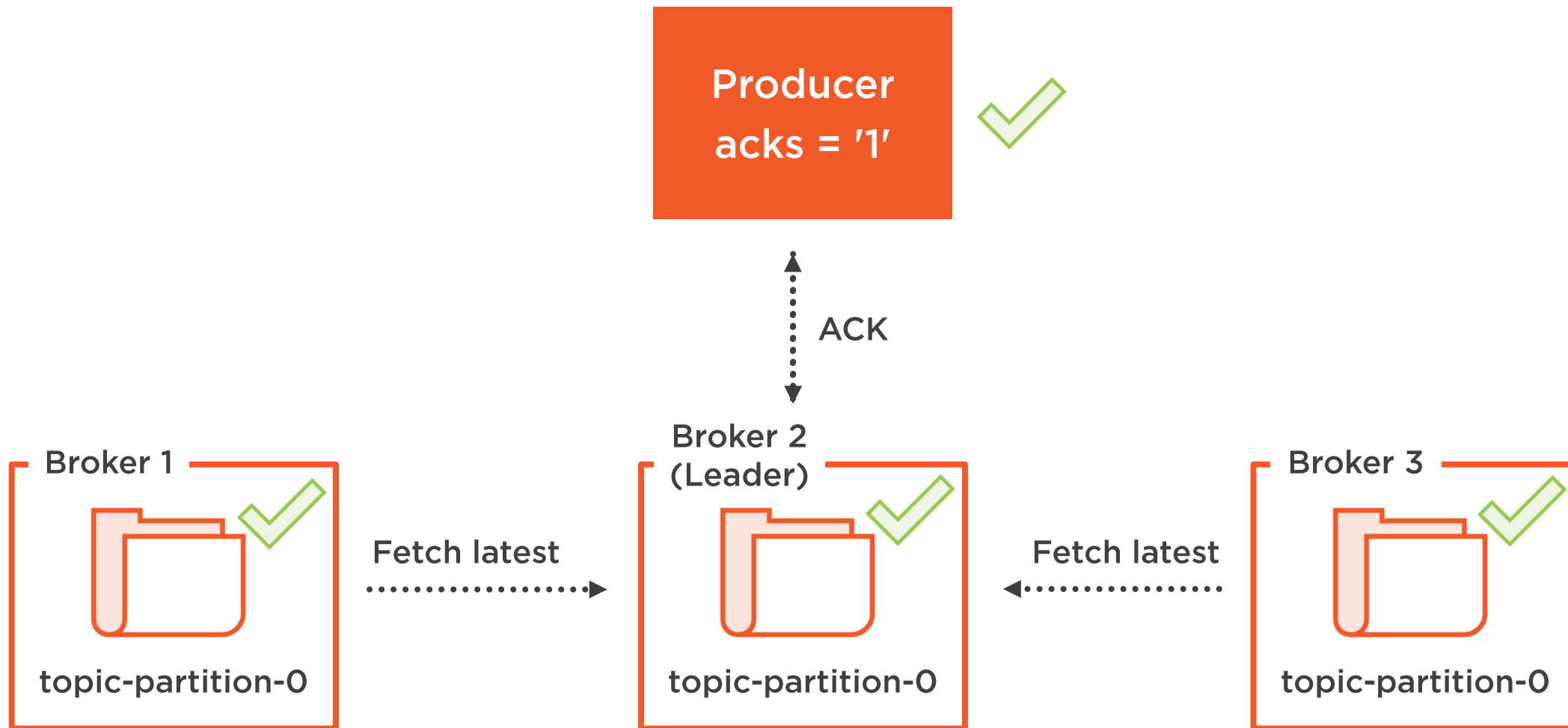
**Replication factor**

# Creating a New Topic

```java
Admin admin = Admin.create(
        Map.of("bootstrap.servers", "localhost:9092")
);

NewTopic topic = new NewTopic("quote-feedback", 2, (short) 3)

admin.createTopics(List.of(topic)).get()
```
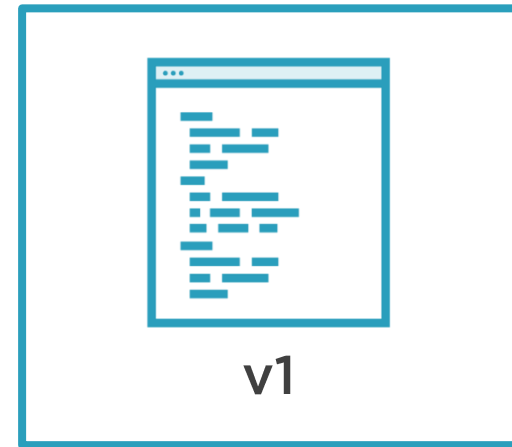
**Replication factor**

# Data Retention and Cluster Sizing

**Broker**

**Consumer**

v1

**Broker**

**Consumer**

Deployment

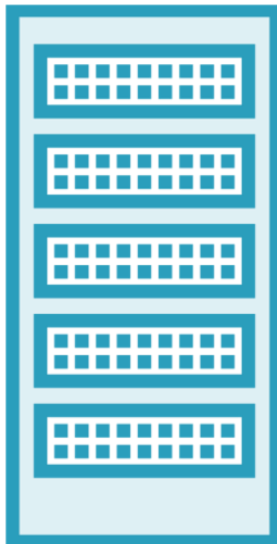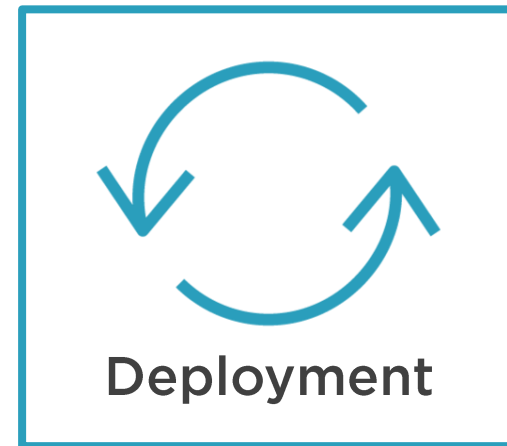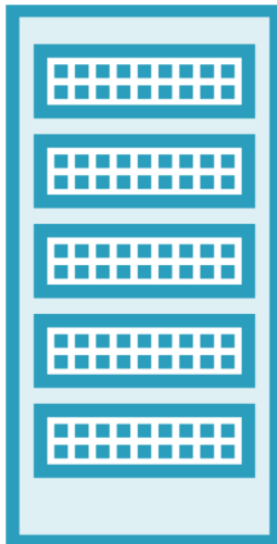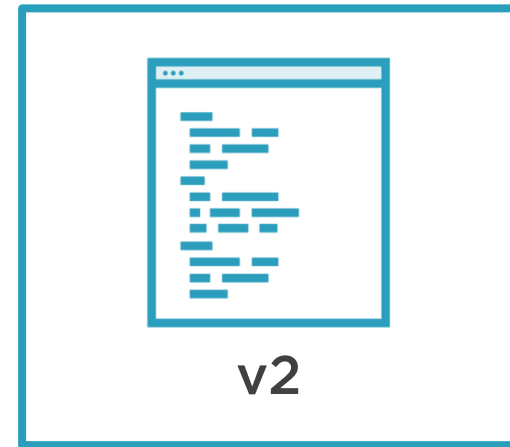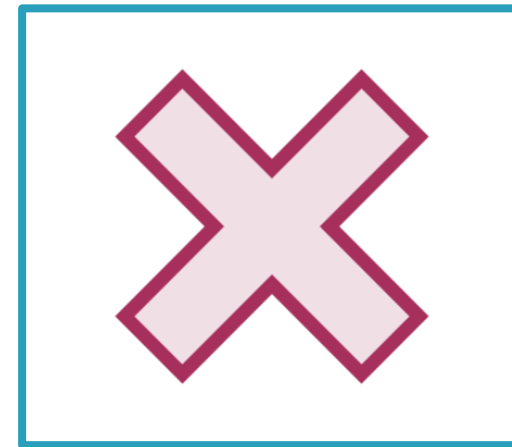Broker

Consumer

v2

Broker

Consumer

# Data Retention

quote-feedback-x

**log == partition replica**
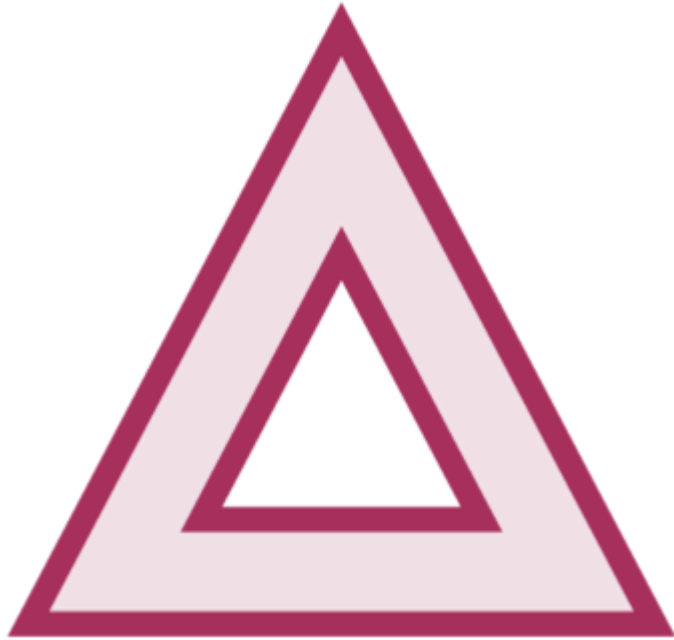
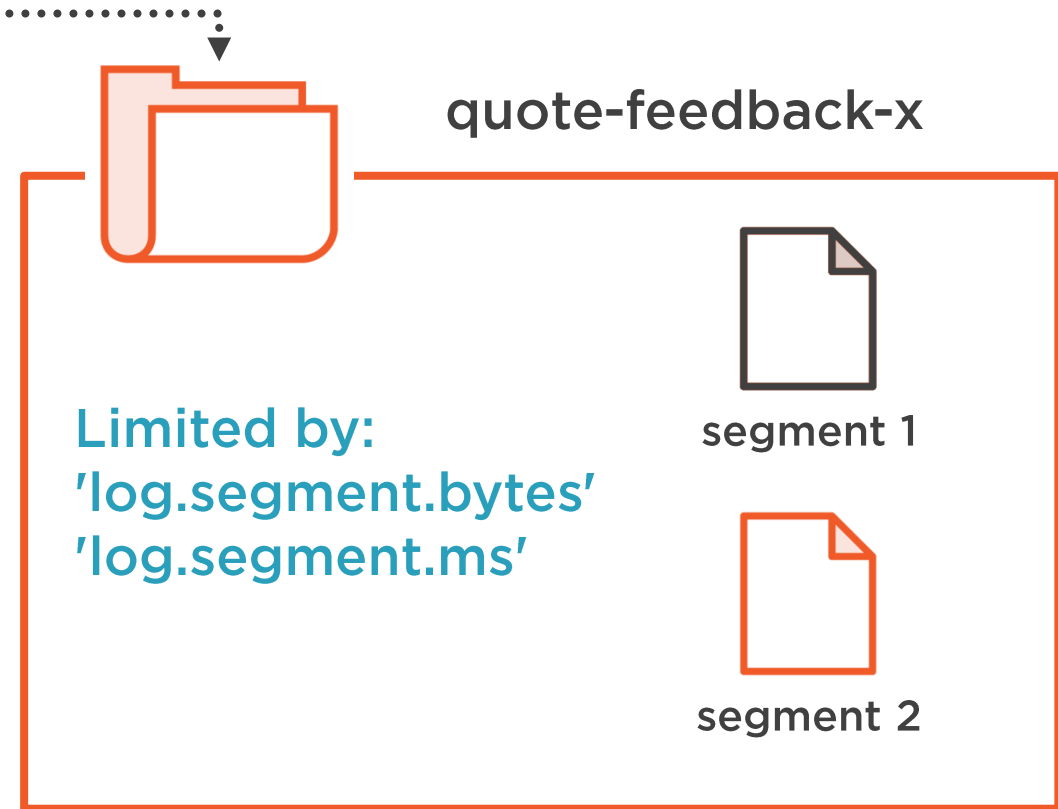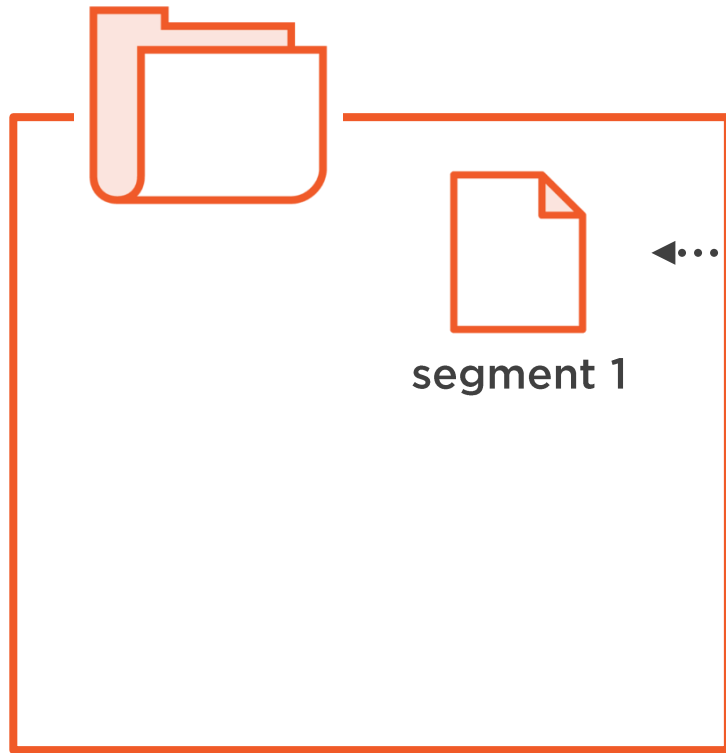Think about possible log size when choosing the number of partitions for a topic

# Data Retention

**Per-log limits:**
**'log.retention.bytes'**
**'log.retention.ms'**

**quote-feedback-x**

Limited by:
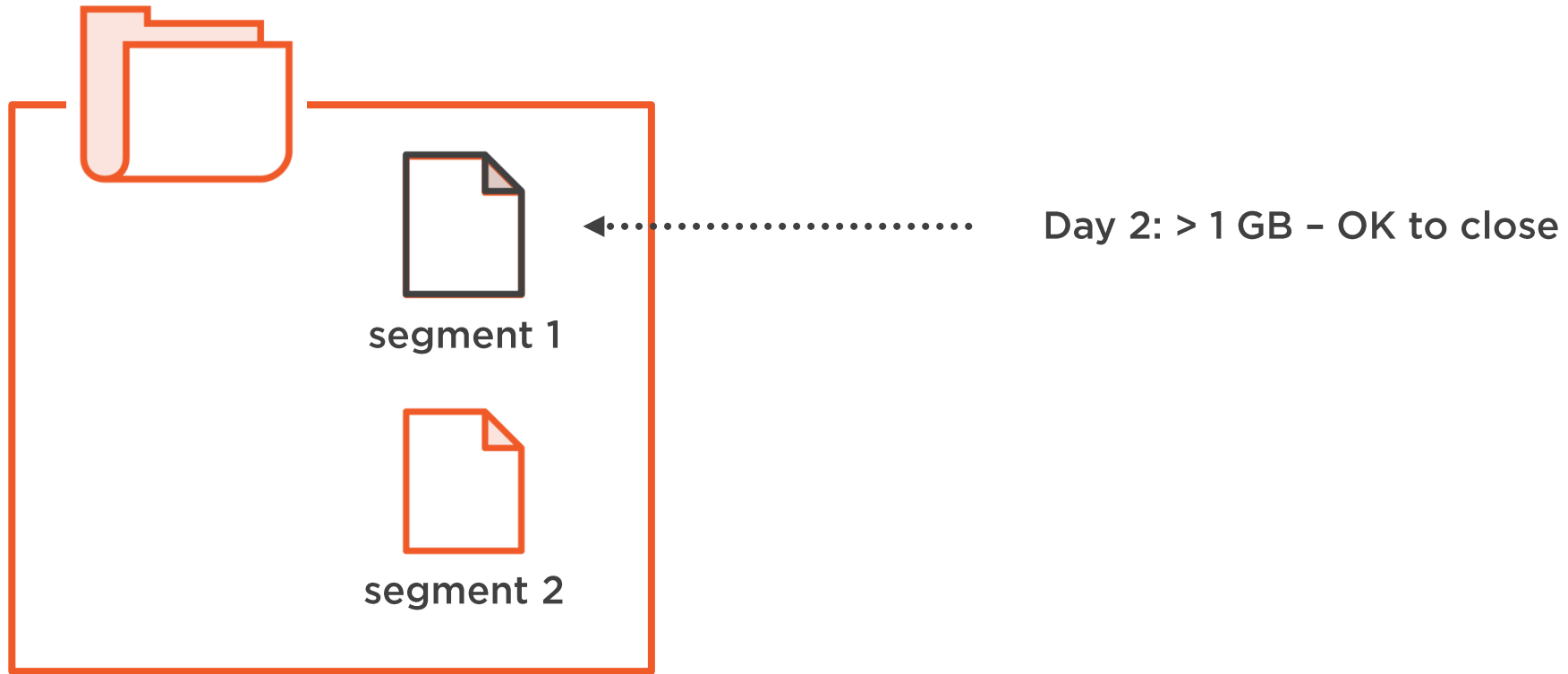'log.segment.bytes'
'log.segment.ms'

segment 1

segment 2

segment 1

Day 1: ~500MB – cannot close

log.segment.bytes: 1GB
log.segment.ms: 604800000 (1 week)
log.retention.bytes: 2GB
log.retention.ms: 86400000 (1 day)

Day 2: > 1 GB – OK to close

segment 1

segment 2

log.segment.bytes: 1GB
log.segment.ms: 604800000 (1 week)
log.retention.bytes: 2GB
log.retention.ms: 86400000 (1 day)

Day 2: Remove based on log.retention.ms

segment 1

segment 2

log.segment.bytes: 1GB
log.segment.ms: 604800000 (1 week)
log.retention.bytes: 2GB
log.retention.ms: 86400000 (1 day)

segment 1

segment 2

Defaults
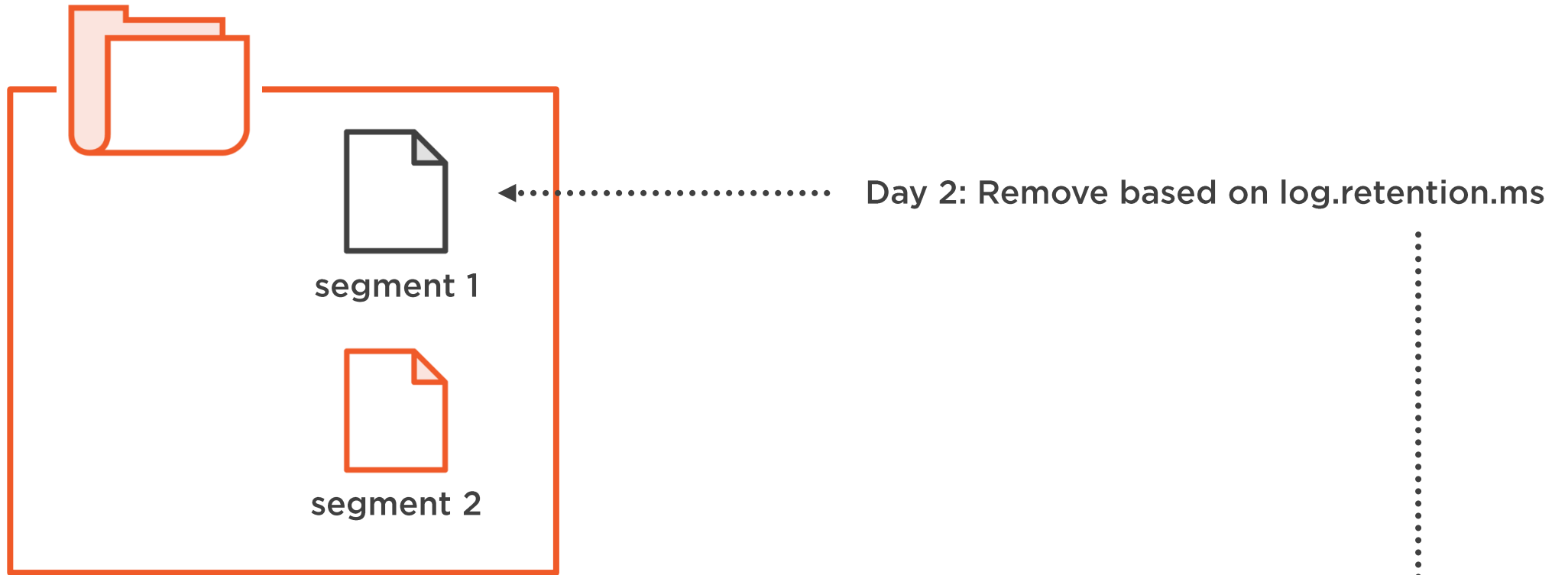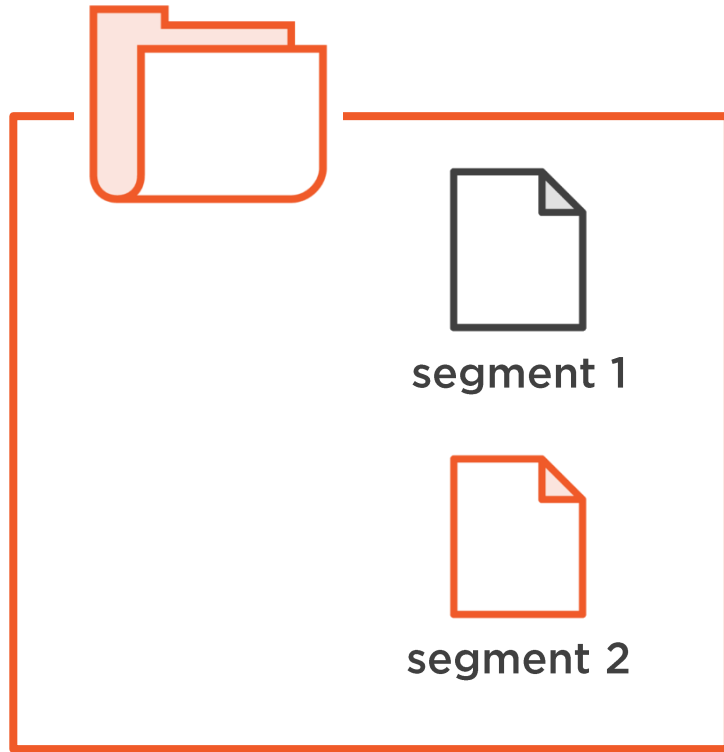
log.segment.bytes: 1GB
log.segment.ms: 604800000 (1 week)
log.retention.bytes: 2GB
log.retention.ms: 86400000 (1 day)

# Cluster Sizing

**Usual bottlenecks:**
- Disk size/throughput
- Network throughput

Always estimate the load
before deploying a cluster

# Summary

Why multiple brokers are needed

Reliability guarantees

Practical example

Basics of data retention and cluster sizing