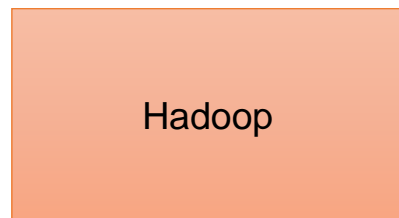
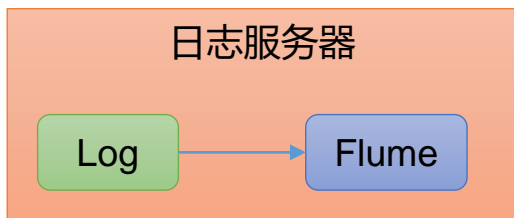




前端埋点记录用户**购买海购人参丸的行为数据**（浏览、点赞、收藏、评论等）。



日常：Flume采集速度，**小于100m/s。**

11.11 活动：Flume采集速度，**大于200m/s。**

Kafka传统定义：Kafka是一个**分布式**的基于**发布/订阅模式**的消息队列（Message Queue），主要应用于大数据实时处理领域。

发布/订阅：消息的发布者不会将消息直接发送给特定的订阅者，而是**将发布的消息分为不同的类别**，订阅者**只接收感兴趣的消息**。

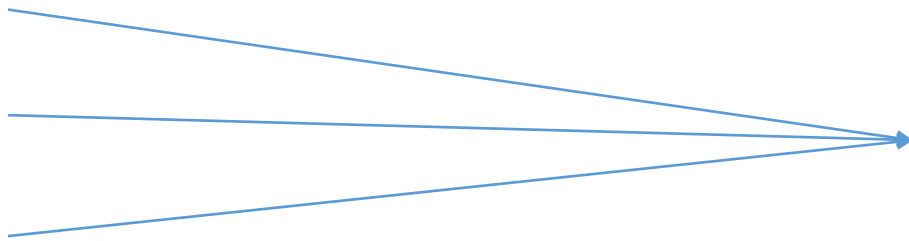
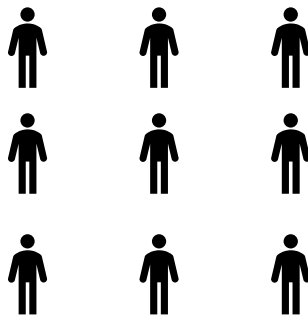
Kafka最新定义：Kafka是一个开源的**分布式事件流平台**（Event Streaming Platform），被数千家公司用于高性能**数据管道**、**流分析**、**数据集成**和**关键任务应用**。

Hadoop上传速度**100m/s左右。**



缓冲/消峰：有助于控制和优化数据流经过系统的速度，解决生产消息和消费消息的处理速度不一致的情况。

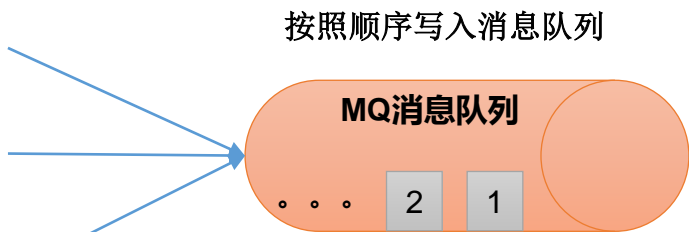
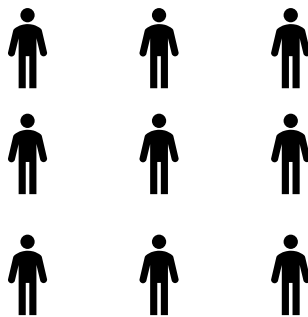
双十一参与用户：10亿人/s



处理能力：1千万人/s

秒杀系统
(秒杀海购人参丸)

双十一参与用户：10亿人/s

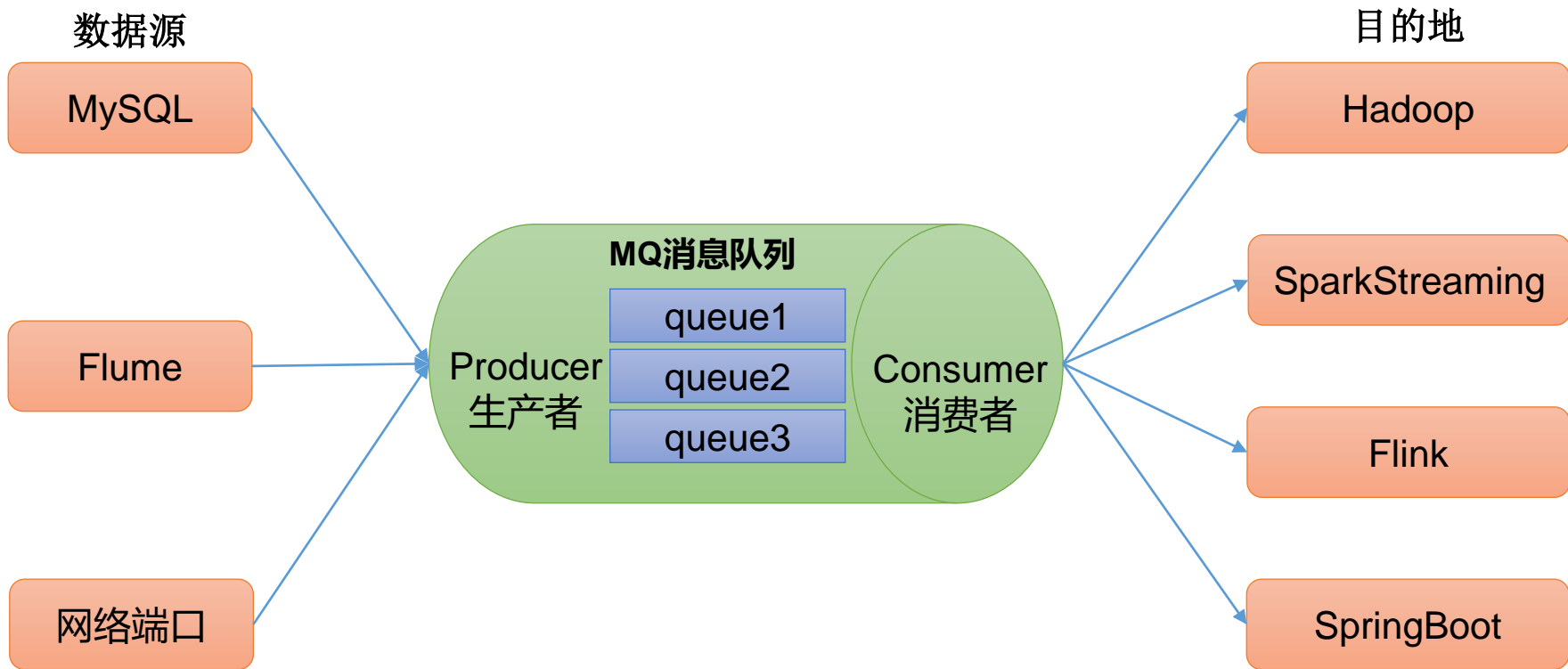


处理能力：1千万人/s

秒杀系统
(秒杀海购人参丸)



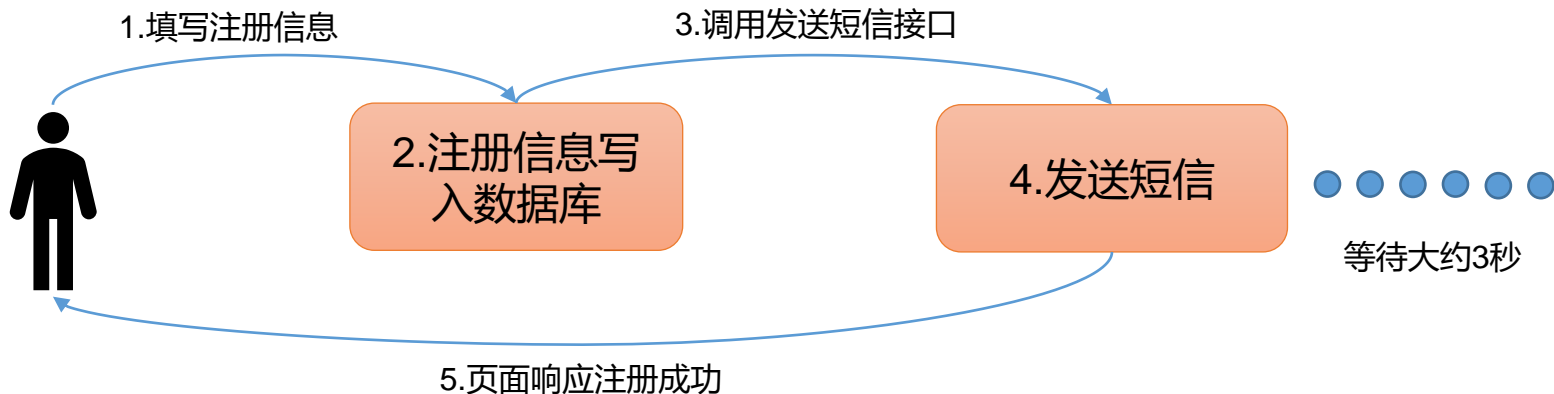
解耦： 允许你独立的扩展或修改两边的处理过程，只要确保它们遵守同样的接口约束。



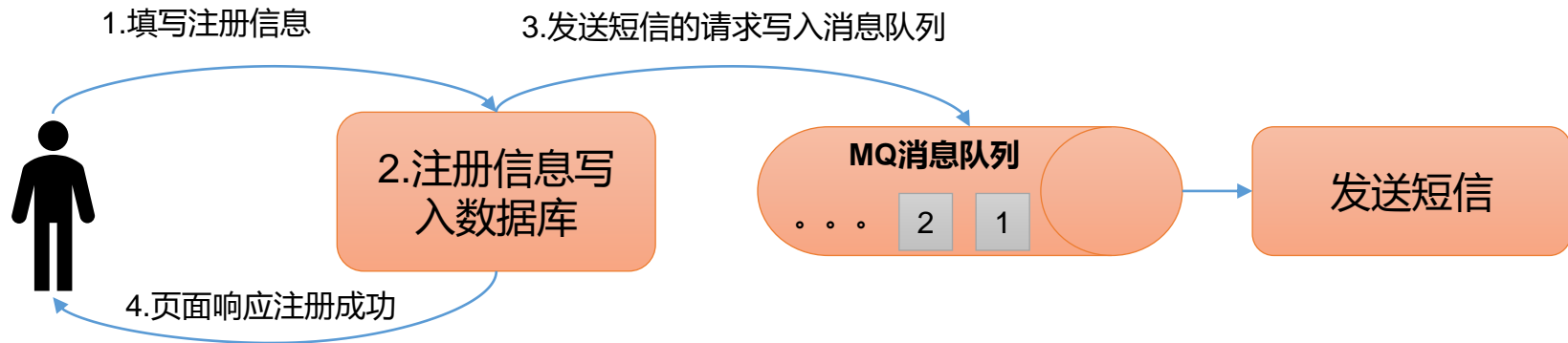


异步通信：允许用户把一个消息放入队列，但并不立即处理它，然后在需要的时候再去处理它们。

同步处理



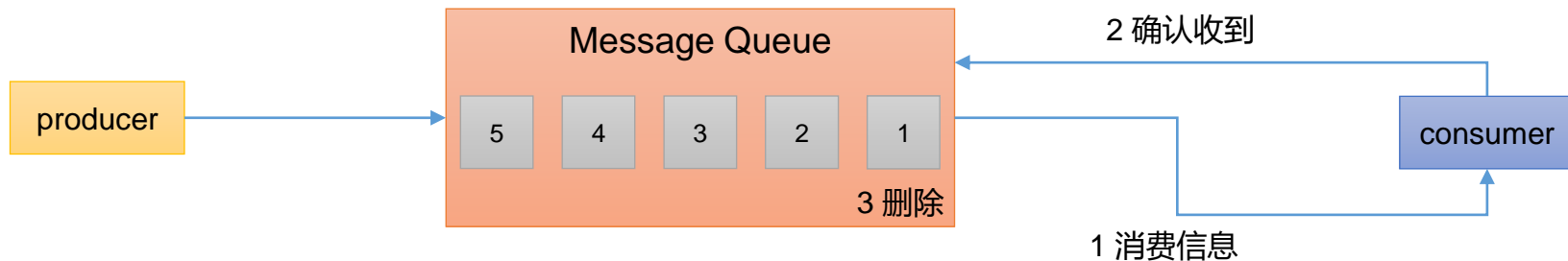
异步处理





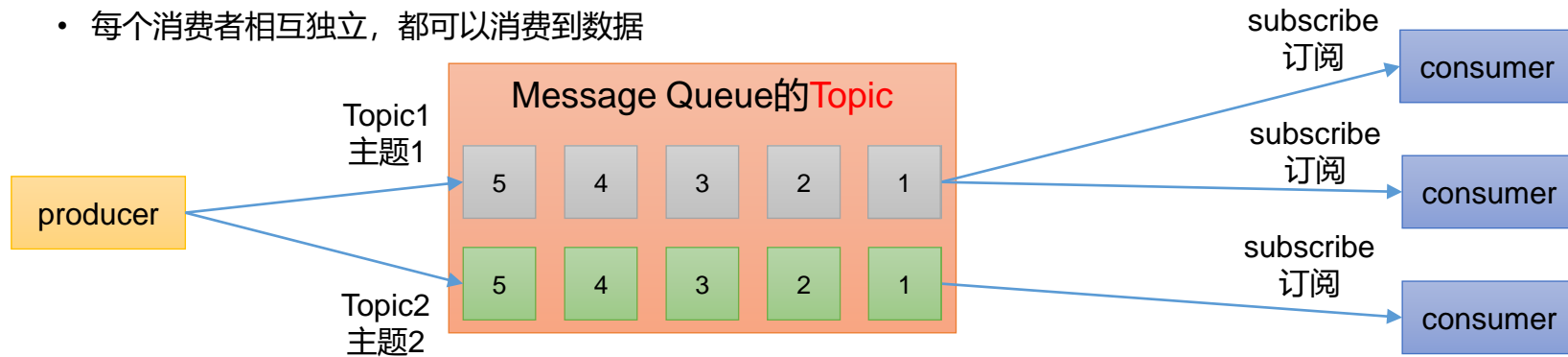
1) 点对点模式

- 消费者主动拉取数据，消息收到后清除消息



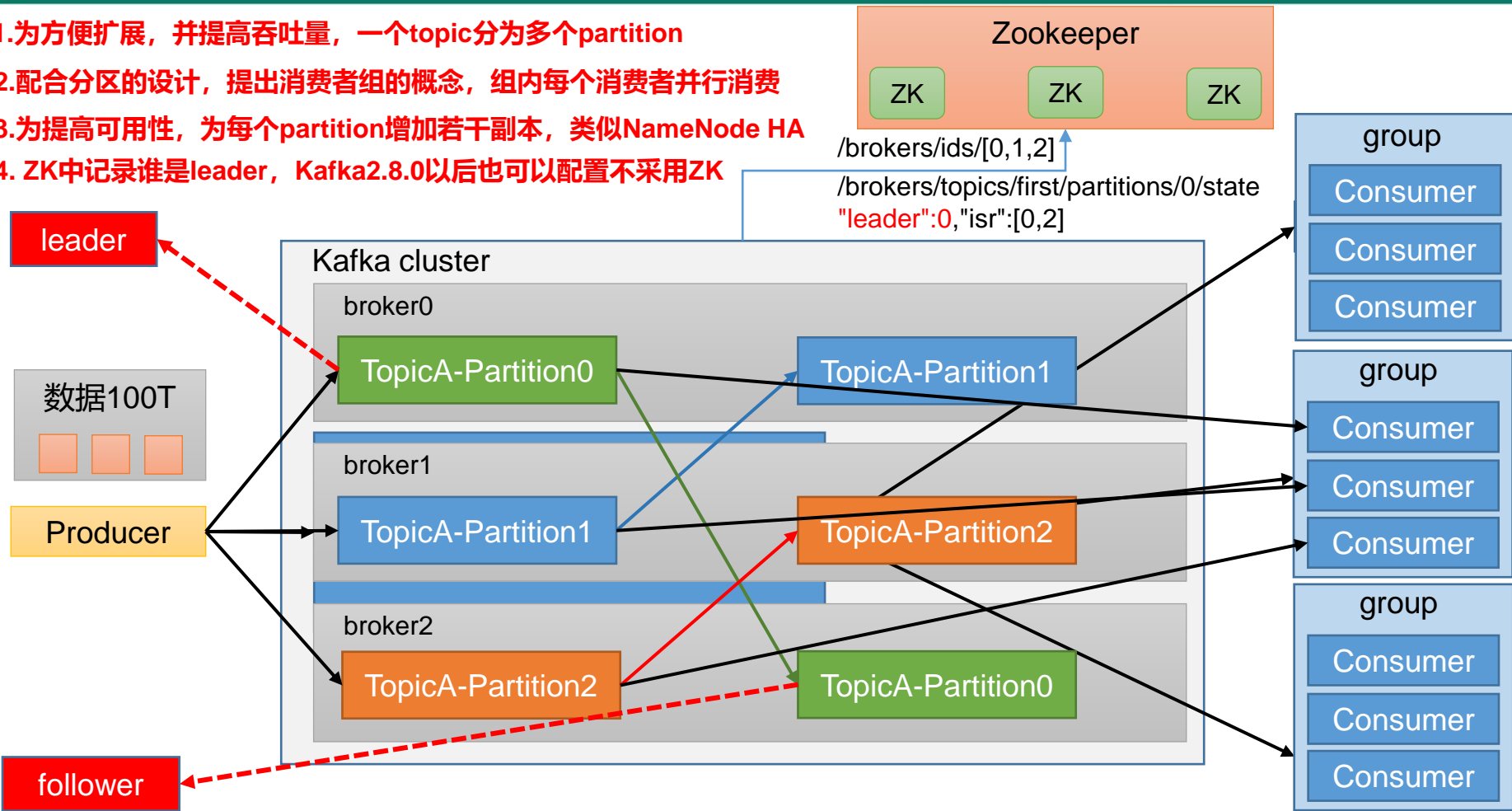
2) 发布/订阅模式

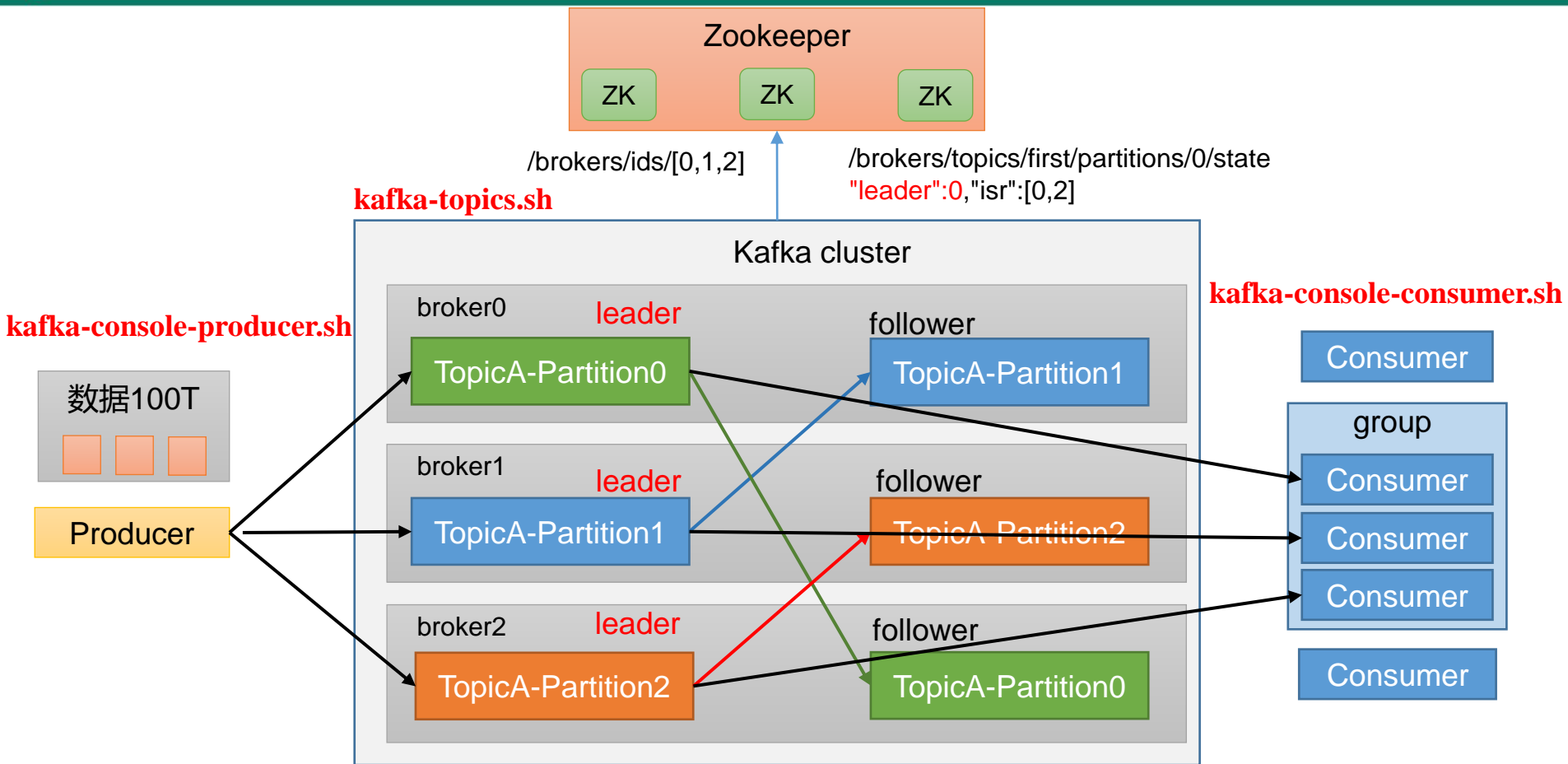
- 可以有多个topic主题 (浏览、点赞、收藏、评论等)
- 消费者消费数据之后，不删除数据
- 每个消费者相互独立，都可以消费到数据

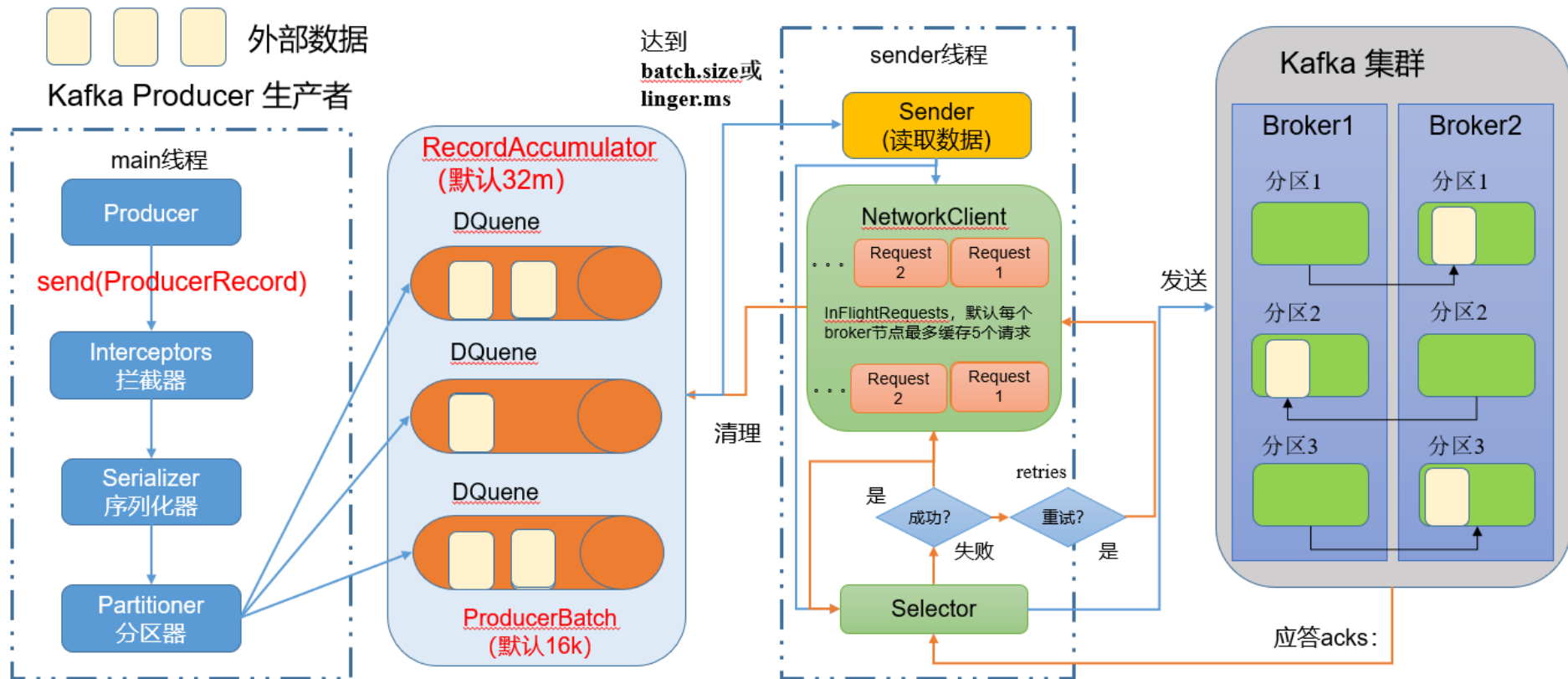




1. 为方便扩展, 并提高吞吐量, 一个topic分为多个partition
2. 配合分区的设计, 提出消费者组的概念, 组内每个消费者并行消费
3. 为提高可用性, 为每个partition增加若干副本, 类似NameNode HA
4. ZK中记录谁是leader, Kafka2.8.0以后也可以配置不采用ZK

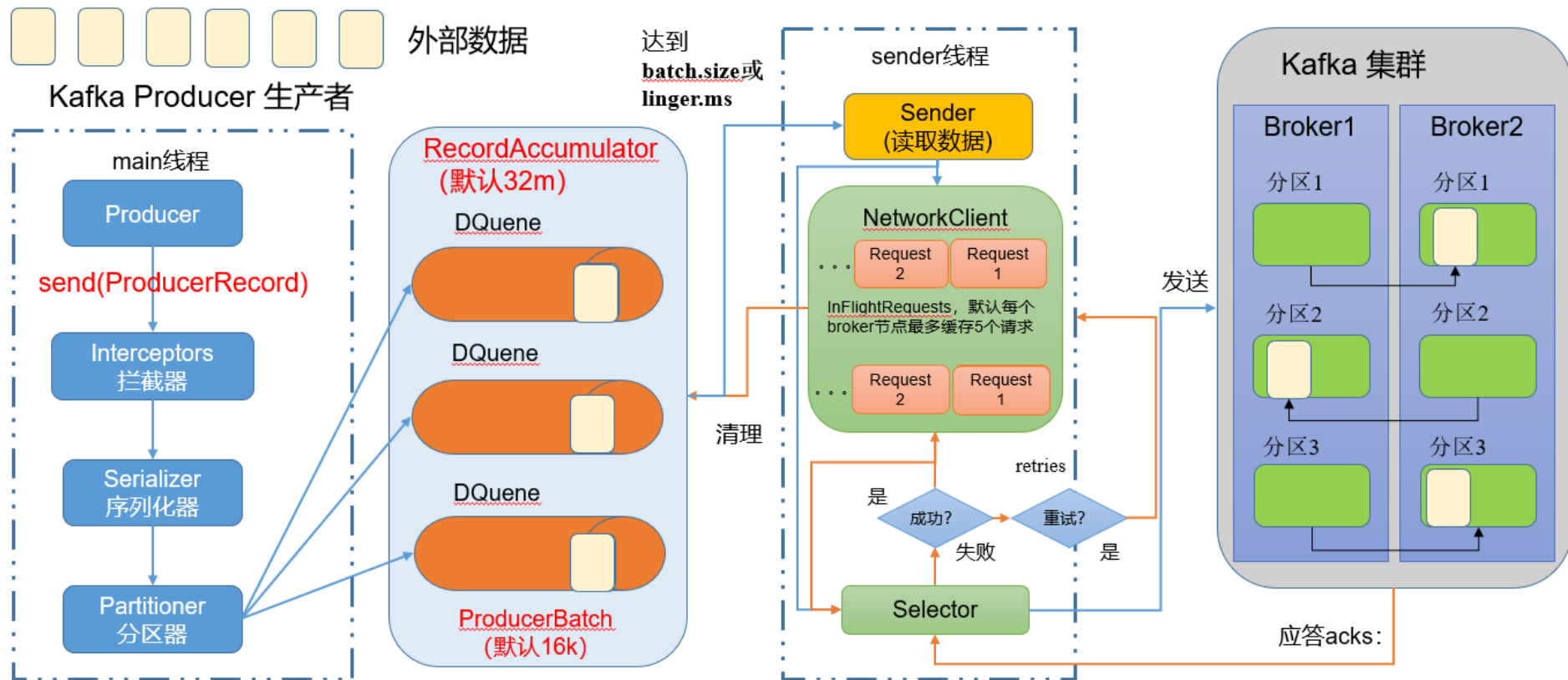






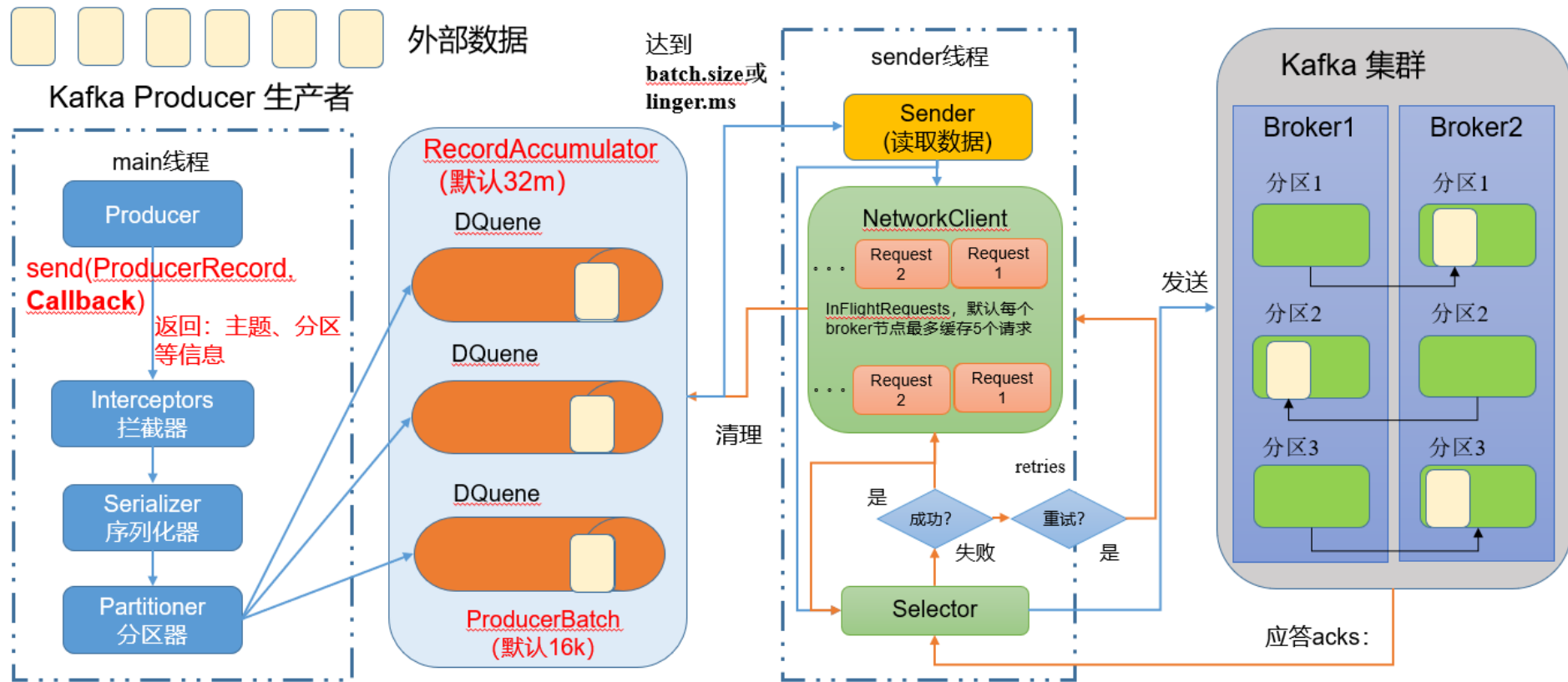
- **batch.size**: 只有数据积累到 **batch.size** 之后，sender 才会发送数据。默认 16k
- **linger.ms**: 如果数据迟迟未达到 **batch.size**，sender 等待 **linger.ms** 设置的时间到了之后就会发送数据。单位 **ms**，默认值是 0ms，表示没有延迟。

- **0**: 生产者发送过来的数据，不需要等数据落盘应答。
- **1**: 生产者发送过来的数据，Leader 收到数据后应答。
- **-1 (all)**: 生产者发送过来的数据，Leader 和 ISR 队列表中的所有节点对齐数据后应答。-1 和 all 等价。



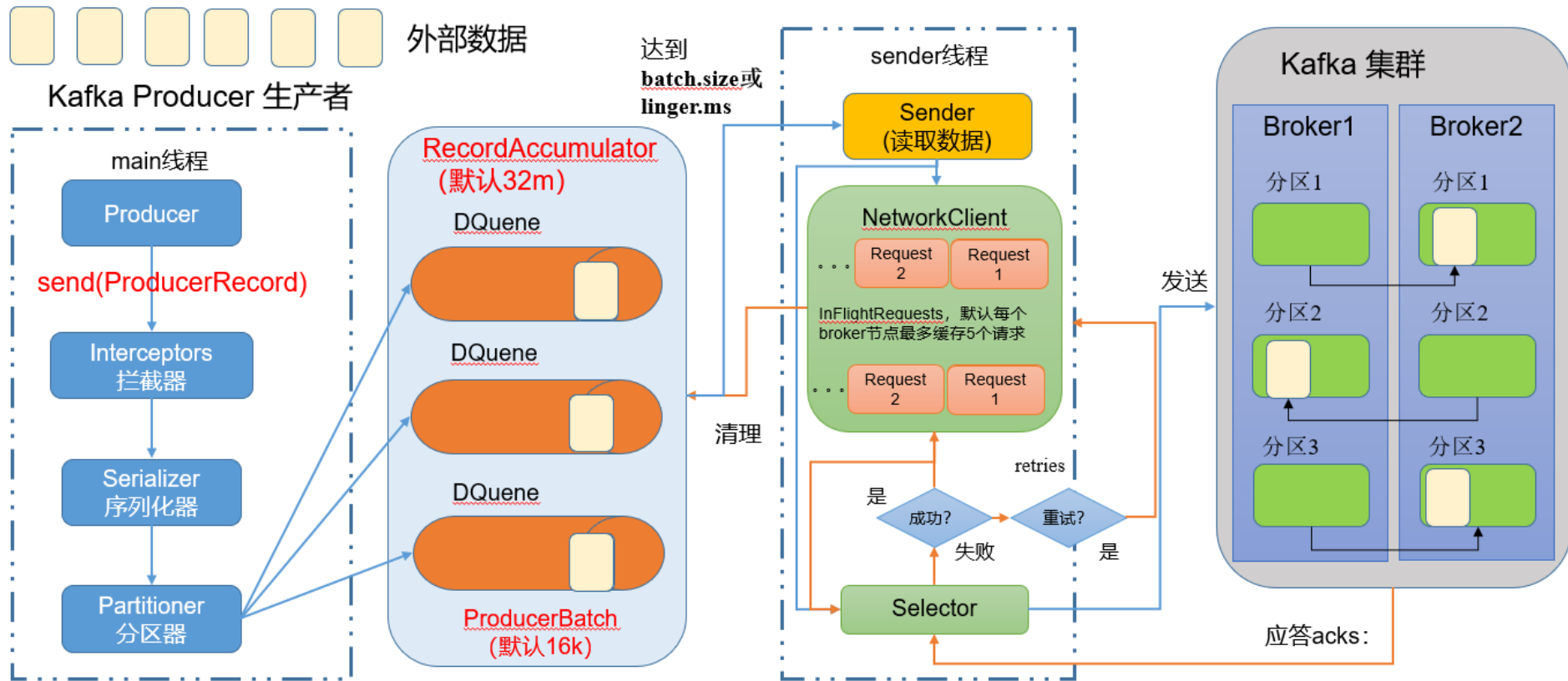
- **batch.size:** 只有数据积累到`batch.size`之后, sender才会发送数据。默认16k
- **linger.ms:** 如果数据迟迟未达到`batch.size`, sender等待`linger.ms`设置的时间到了之后就会发送数据。单位`ms`, 默认值是0ms, 表示没有延迟。

- **0:** 生产者发送过来的数据, 不需要等数据落盘应答。
- **1:** 生产者发送过来的数据, Leader收到数据后应答。
- **-1 (all):** 生产者发送过来的数据, Leader和ISR队里面的所有节点收齐数据后应答。-1和all等价。



- **batch.size**: 只有数据积累到batch.size之后, sender才会发送数据。默认16k
- **linger.ms**: 如果数据迟迟未达到batch.size, sender等待linger.ms设置的时间到了之后就会发送数据。单位ms, 默认值是0ms, 表示没有延迟。

- **0**: 生产者发送过来的数据, 不需要等数据落盘应答。
- **1**: 生产者发送过来的数据, Leader收到数据后应答。
- **-1 (all)**: 生产者发送过来的数据, Leader和ISR队里面的所有节点收齐数据后应答。-1和all等价。



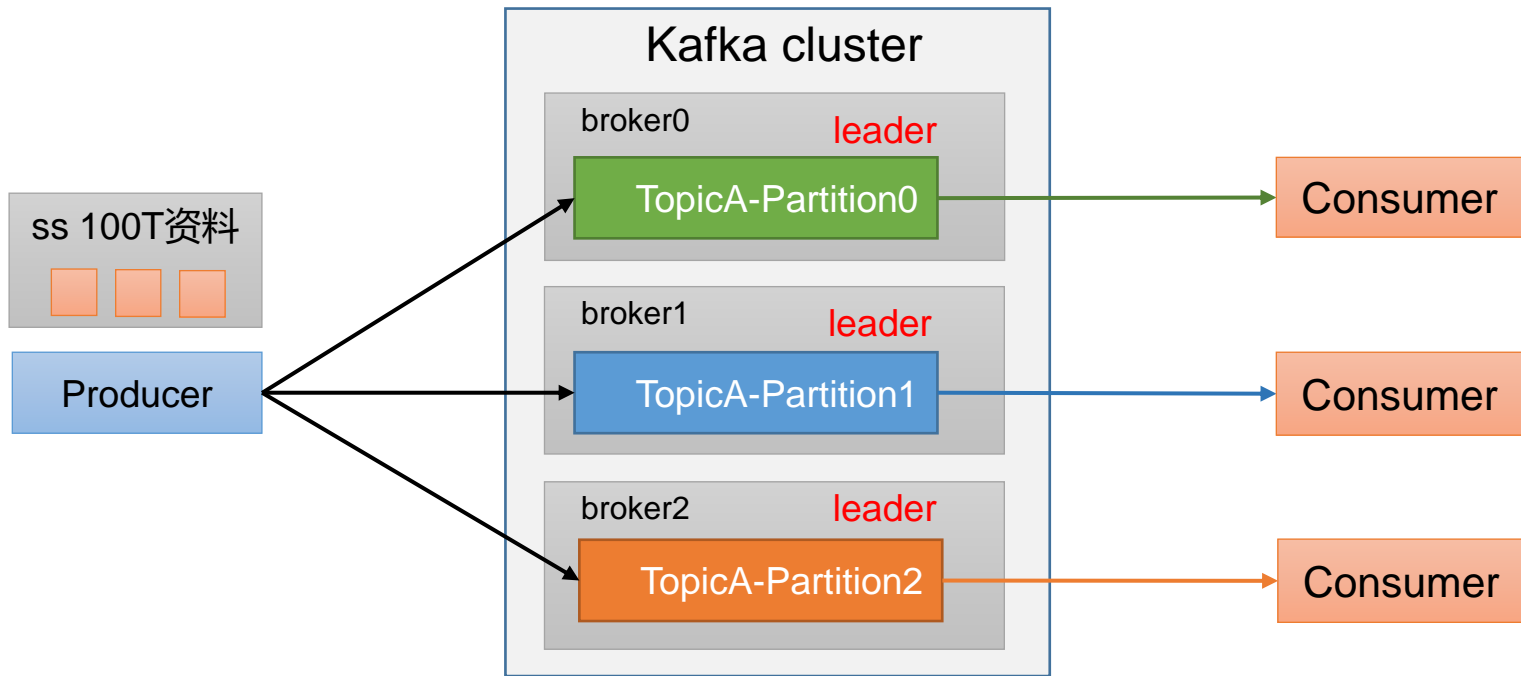
- **batch.size**: 只有数据积累到batch.size之后, sender才会发送数据。默认16k
- **linger.ms**: 如果数据迟迟未达到batch.size, sender等待linger.ms设置的时间到了之后就会发送数据。单位ms, 默认值是0ms, 表示没有延迟。

- **0**: 生产者发送过来的数据, 不需要等数据落盘应答。
- **1**: 生产者发送过来的数据, Leader收到数据后应答。
- **-1 (all)**: 生产者发送过来的数据, Leader和ISR队里面的所有节点对齐数据后应答。-1和all等价。



(1) **便于合理使用存储资源**，每个Partition在一个Broker上存储，可以把海量的数据按照分区切割成一块一块数据存储在多台Broker上。合理控制分区的任务，可以实现**负载均衡**的效果。

(2) **提高并行度**，生产者可以以分区为单位**发送数据**；消费者可以以分区为单位进行**消费数据**。



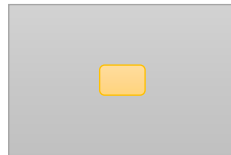


ss家仓库,默认32m

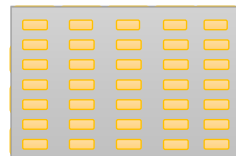
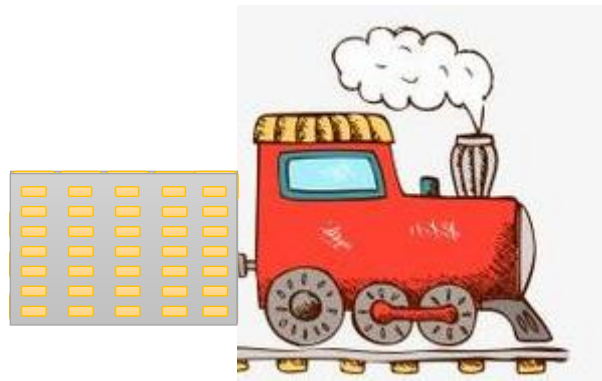


- batch.size: 批次大小, 默认16k
- linger.ms: 等待时间, 修改为5-100ms
- compression.type: 压缩snappy
- RecordAccumulator: 缓冲区大小, 修改为64m

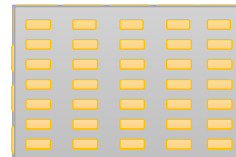
一次拉一个,
来了就走



broker



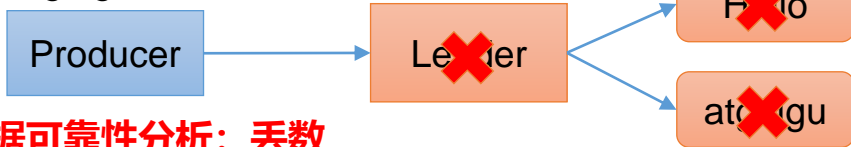
...



acks:

0: 生产者发送过来的数据，不需要等数据落盘应答

atguigu Hello

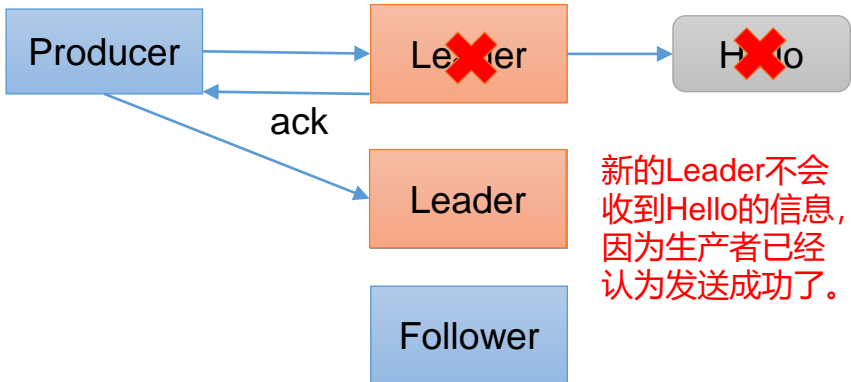


数据可靠性分析：丢数

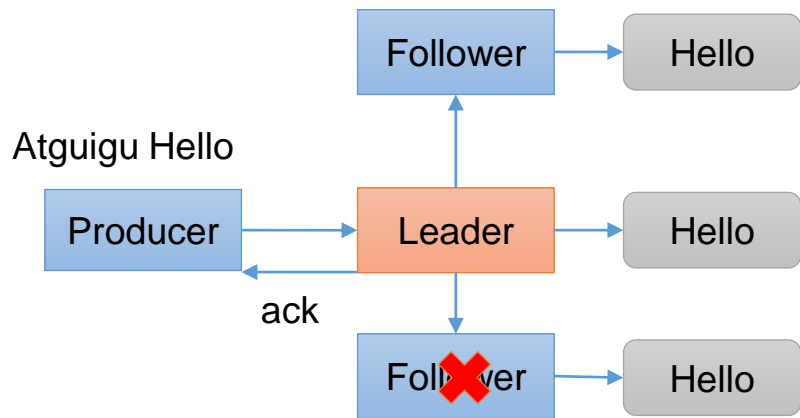
1: 生产者发送过来的数据，Leader收到数据后应答。

数据可靠性分析：丢数

Hello



-1 (all): 生产者发送过来的数据，Leader和ISR队列里面的所有节点收齐数据后应答。

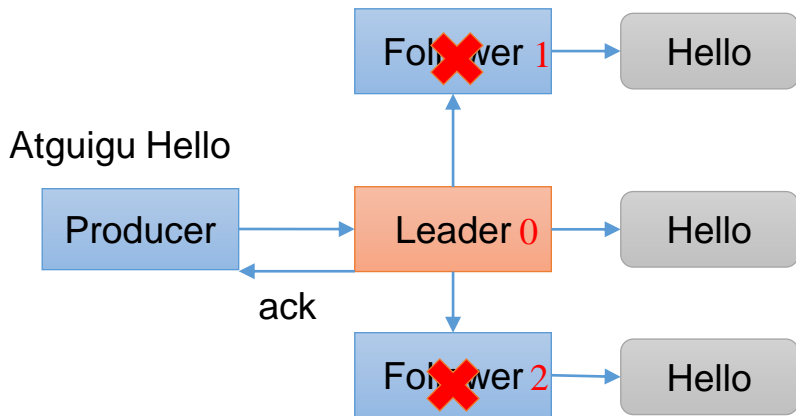


思考：Leader收到数据，所有Follower都开始同步数据，但有一个Follower，因为某种故障，迟迟不能与Leader进行同步，那这个问题怎么解决呢？



acks:

-1 (all)：生产者发送过来的数据，Leader和ISR队列里面的所有节点收齐数据后应答。



思考：Leader收到数据，所有Follower都开始同步数据，但有一个Follower，因为某种故障，迟迟不能与Leader进行同步，那这个问题怎么解决呢？

Leader维护了一个动态的in-sync replica set (**ISR**)，意为和Leader保持同步的Follower+Leader集合(leader: 0, isr:0,1,2)。

如果Follower长时间未向Leader发送通信请求或同步数据，则该Follower将被踢出ISR。该时间阈值由**replica.lag.time.max.ms**参数设定，默认30s。例如2超时，(leader:0, isr:0,1)。

这样就不用等长期联系不上或者已经故障的节点。

数据可靠性分析：

如果分区副本设置为1个，或者ISR里应答的最小副本数量（`min.insync.replicas` 默认为1）设置为1，和ack=1的效果是一样的，仍然有丢数的风险 (leader: 0, isr:0)。

- 数据完全可靠条件 = ACK级别设置为-1 + 分区副本大于等于2 + ISR里应答的最小副本数量大于等于2**



可靠性总结:

acks=0, 生产者发送过来数据就不管了, 可靠性差, 效率高;

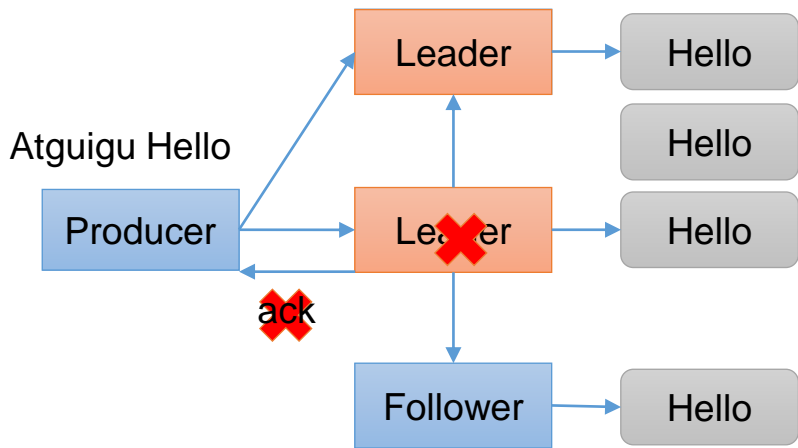
acks=1, 生产者发送过来数据Leader应答, 可靠性中等, 效率中等;

acks=-1, 生产者发送过来数据Leader和ISR队列里面所有Follower应答, 可靠性高, 效率低;

在生产环境中, **acks=0**很少使用; **acks=1**, 一般用于传输普通日志, 允许丢个别数据; **acks=-1**, 一般用于传输和钱相关的数据, 对可靠性要求比较高的场景。

数据重复分析:

acks: -1 (all): 生产者发送过来的数据, Leader和ISR队列里面的所有节点收齐数据后应答。



接收了两份Hello数据, 导致数据重复

具体如何解决数据重复? 下回分解。

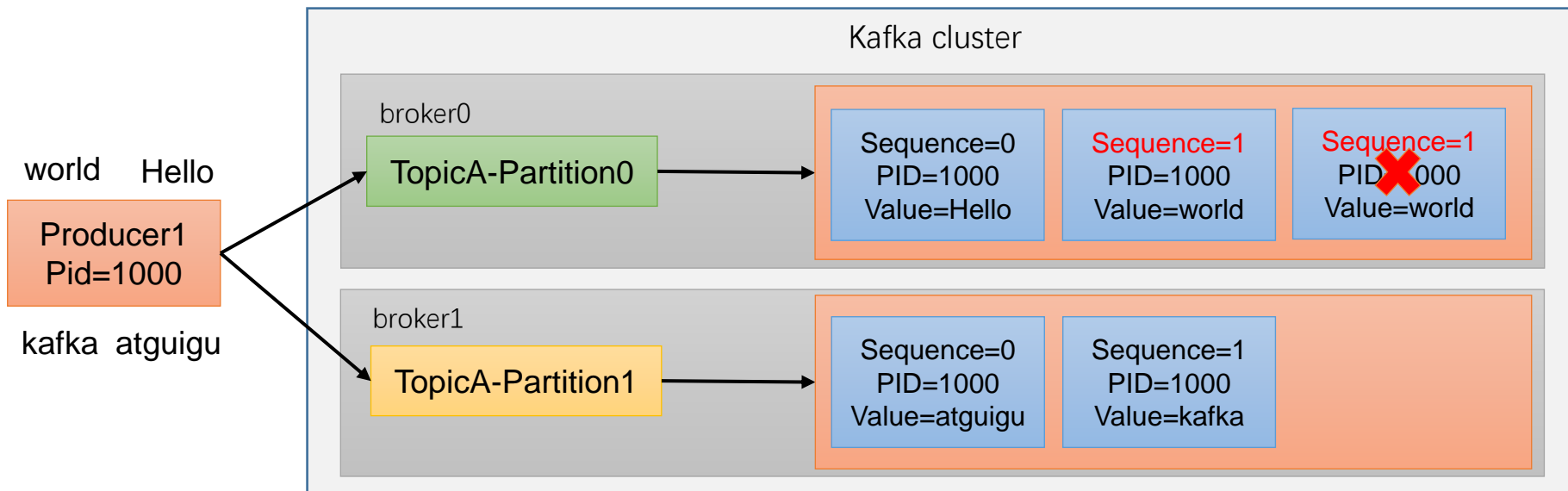


幂等性就是指Producer不论向Broker发送多少次重复数据，Broker端都只会持久化一条，保证了不重复。

精确一次（Exactly Once）= 幂等性 + 至少一次（ $\text{ack}=-1 + \text{分区副本数} \geq 2 + \text{ISR最小副本数量} \geq 2$ ）。

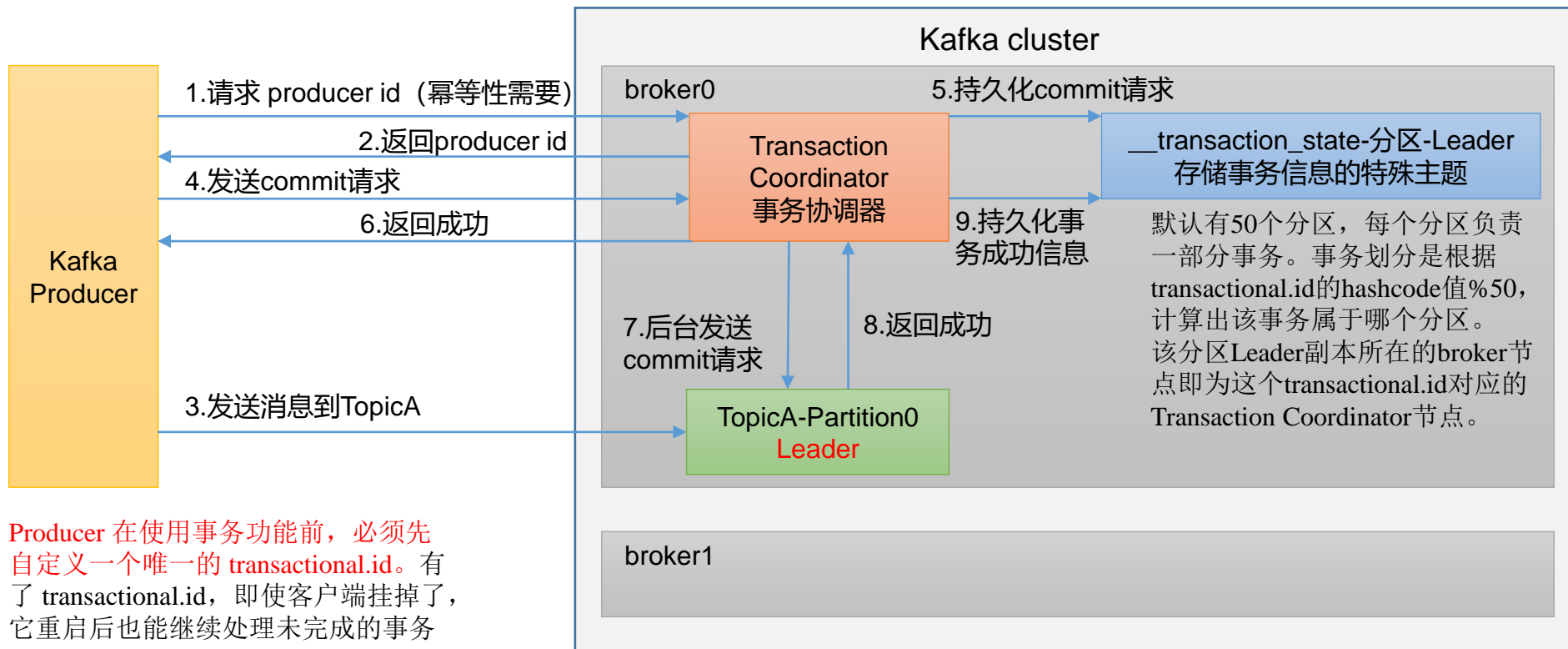
重复数据的判断标准：具有<PID, Partition, SeqNumber>相同主键的消息提交时，Broker只会持久化一条。其中PID是Kafka每次重启都会分配一个新的；Partition表示分区号；Sequence Number是单调自增的。

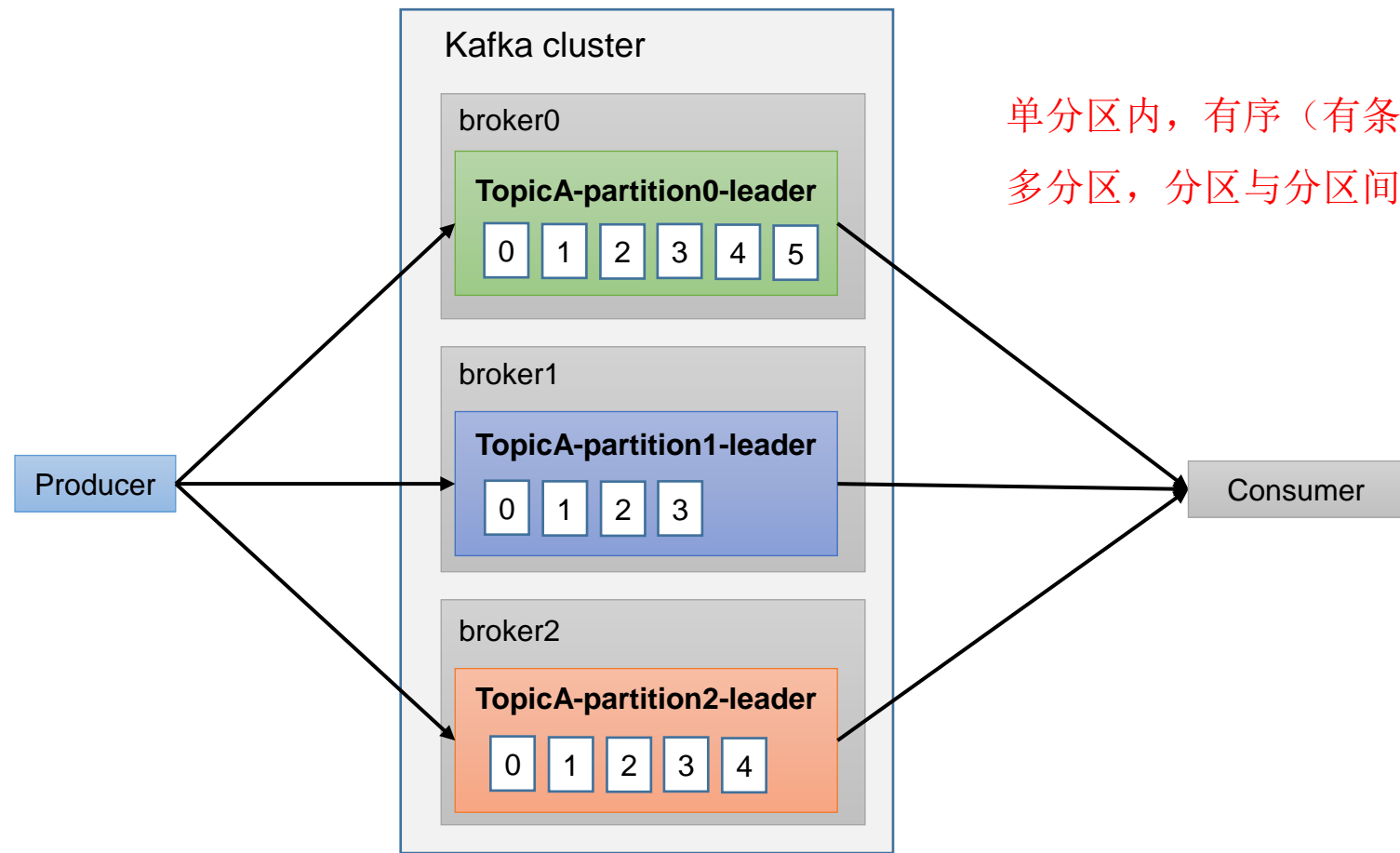
所以幂等性只能保证的是在单分区单会话内不重复。





说明：开启事务，必须开启幂等性。





单分区内，有序（有条件的，详见下节）；
多分区，分区与分区间无序；



1) kafka在1.x版本之前保证数据单分区有序，条件如下：

`max.in.flight.requests.per.connection=1`（不需要考虑是否开启幂等性）。

2) kafka在1.x及以后版本保证数据单分区有序，条件如下：

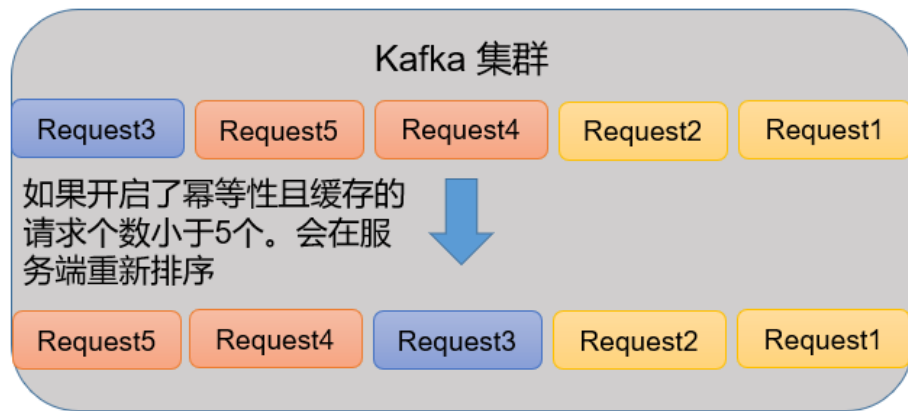
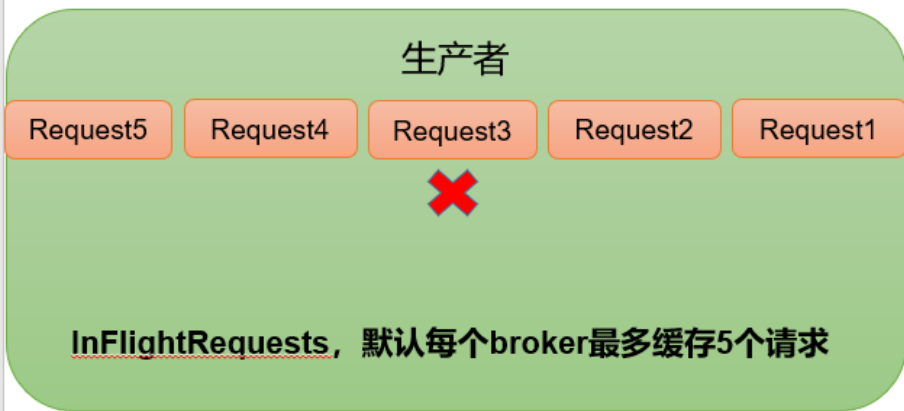
(1) 未开启幂等性

`max.in.flight.requests.per.connection`需要设置为1。

(2) 开启幂等性

`max.in.flight.requests.per.connection`需要设置小于等于5。

原因说明：因为在kafka1.x以后，启用幂等后，kafka服务端会缓存producer发来的最近5个request的元数据，故无论如何，都可以保证最近5个request的数据都是有序的。



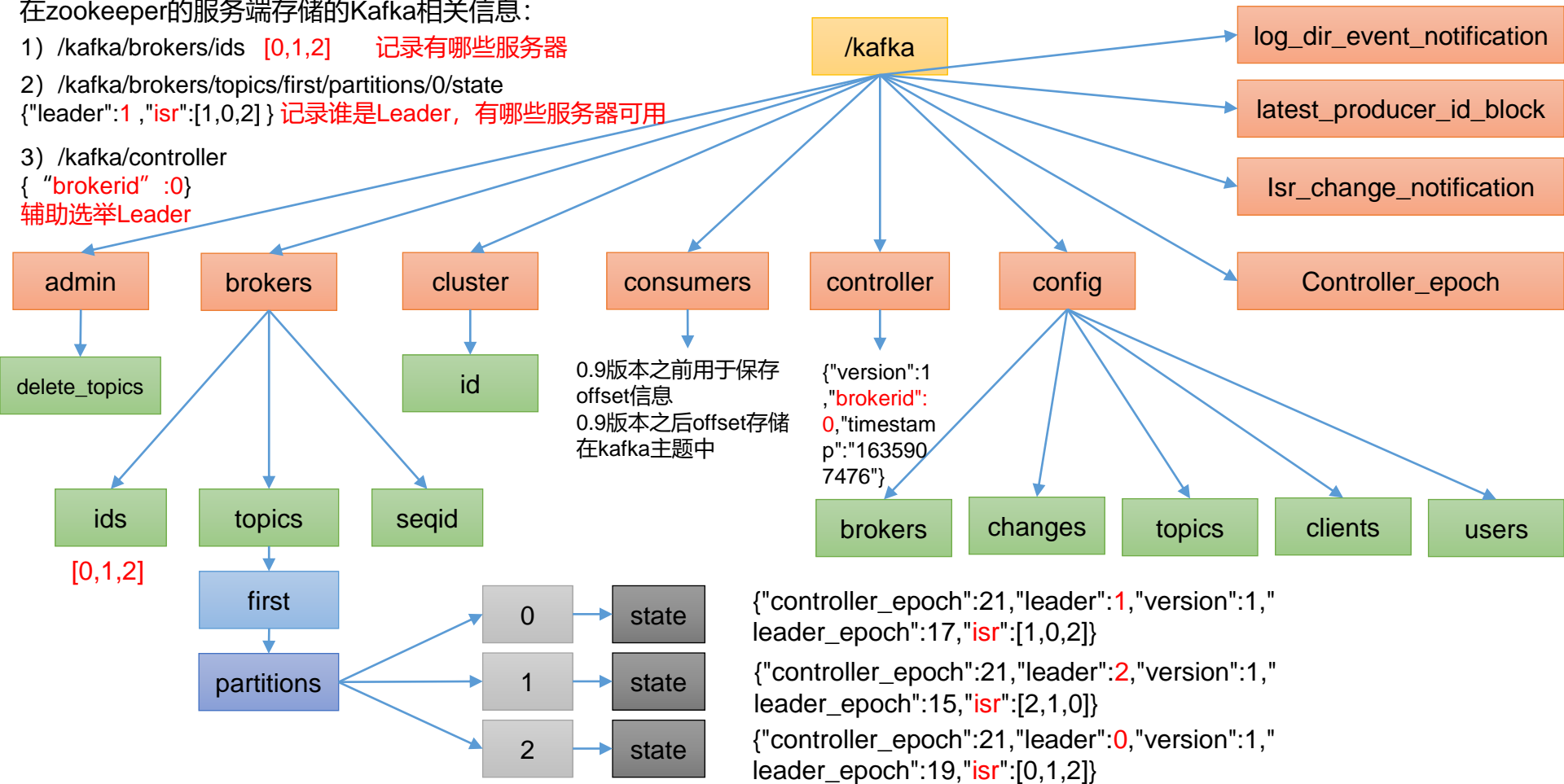


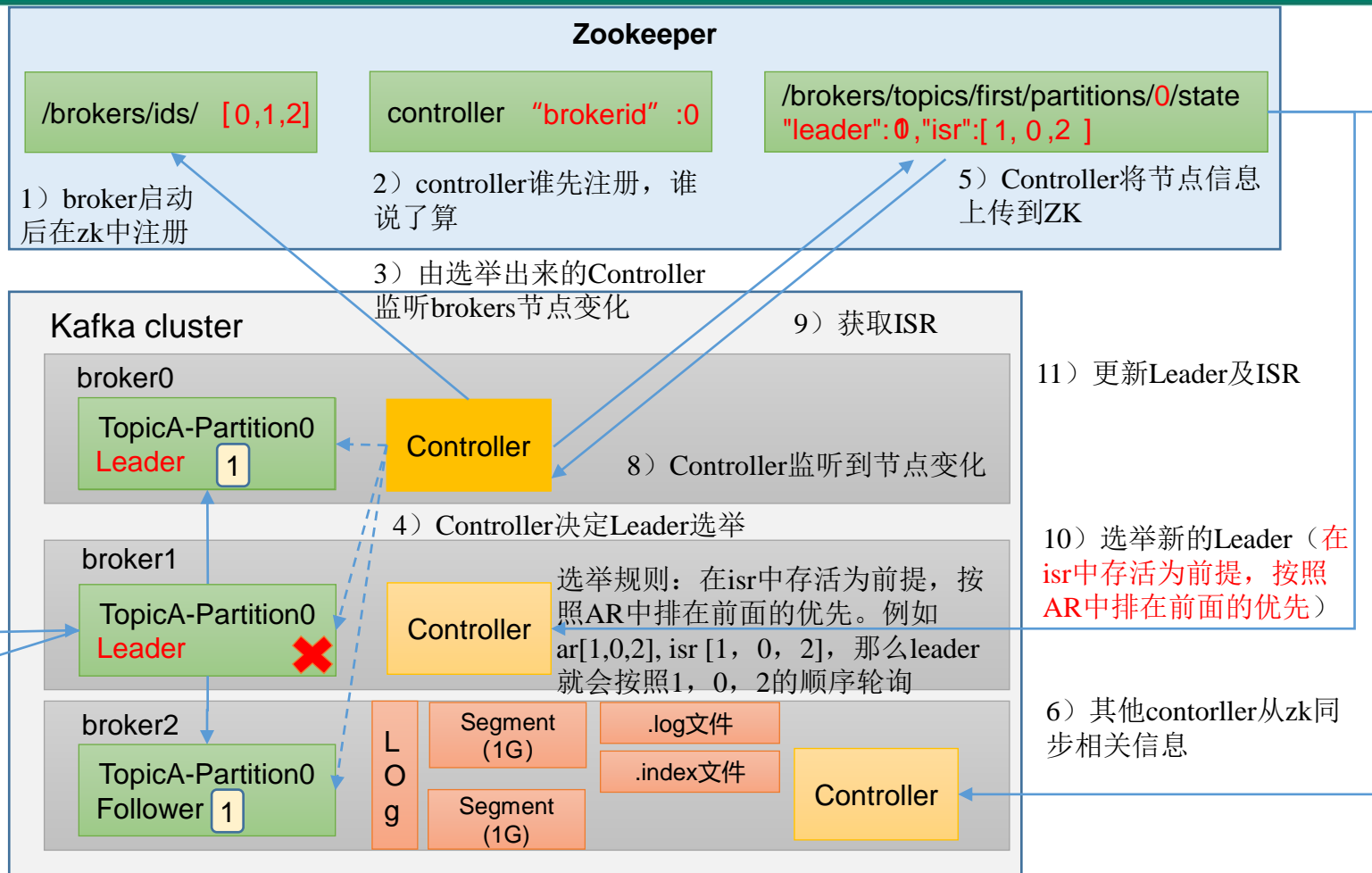
在zookeeper的服务端存储的Kafka相关信息:

1) /kafka/brokers/ids [0,1,2] 记录有哪些服务器

2) /kafka/brokers/topics/first/partitions/0/state
{"leader":1,"isr":[1,0,2]} 记录谁是Leader, 有哪些服务器可用

3) /kafka/controller
{ "brokerid" :0}
辅助选举Leader





AR: Kafka分区中的所有副本统称

7) 假设Broker1中Leader挂了

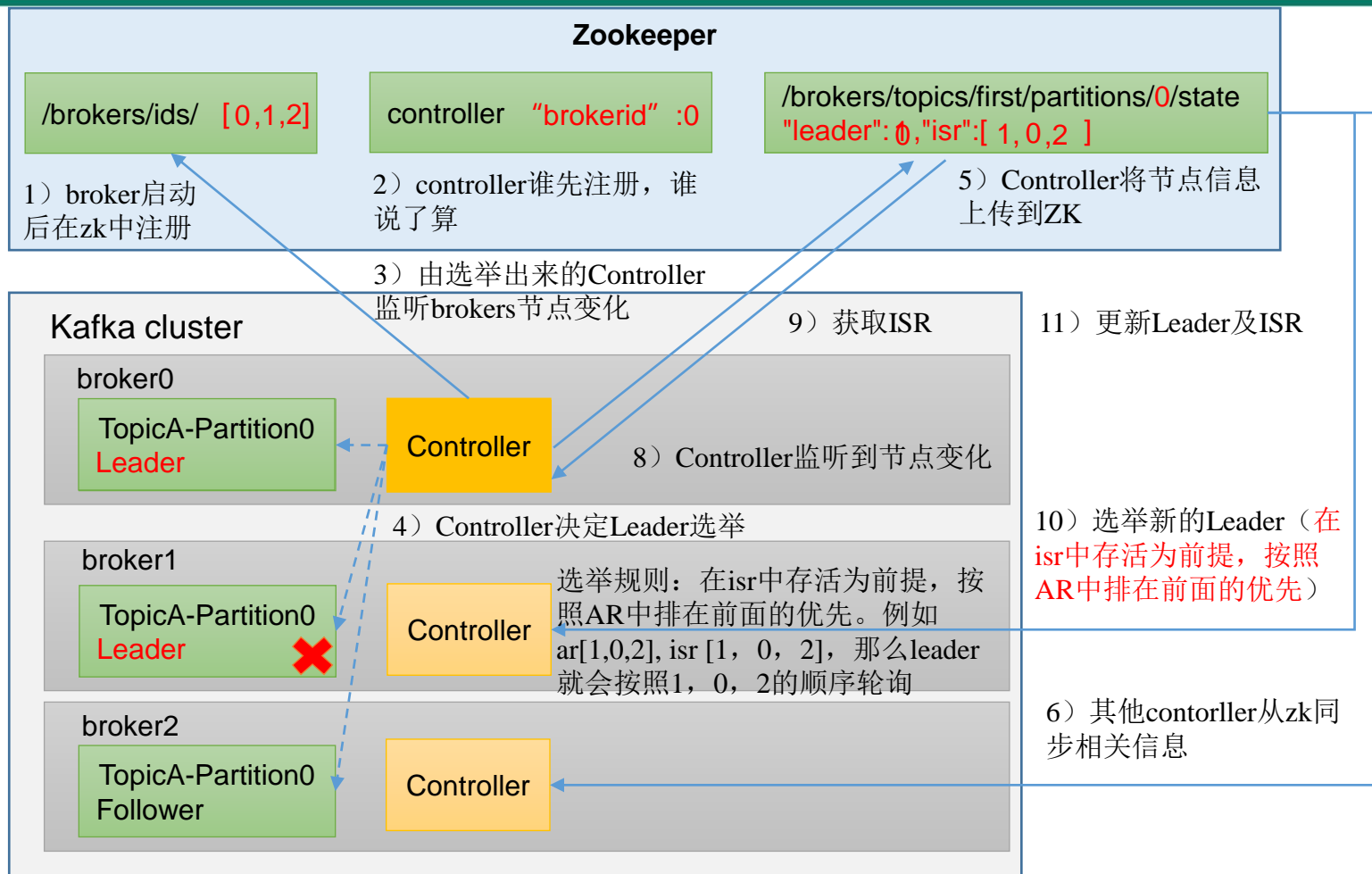
发送消息

应答消息

11) 更新Leader及ISR

10) 选举新的Leader (在isr中存活为前提, 按照AR中排在前面的优先)

6) 其他controller从zk同步相关信息



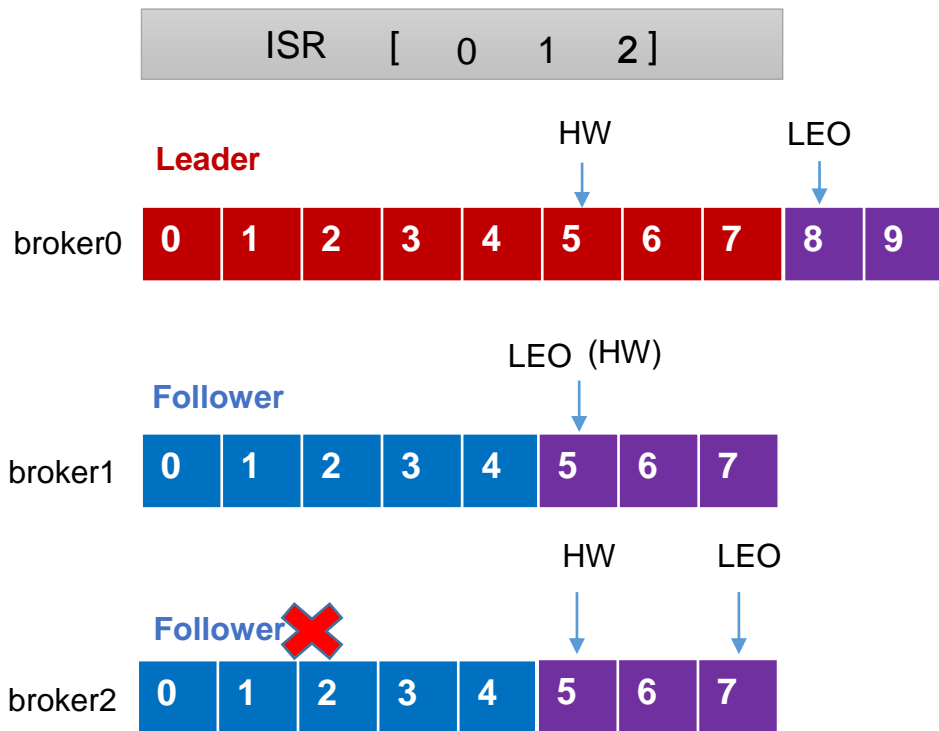
AR: Kafka分区中的所有副本统称

7) 假设Broker1中Leader挂了



LEO (Log End Offset) : 每个副本的最后一个offset, LEO其实就是最新的offset + 1。

HW (High Watermark) : 所有副本中最小的LEO。



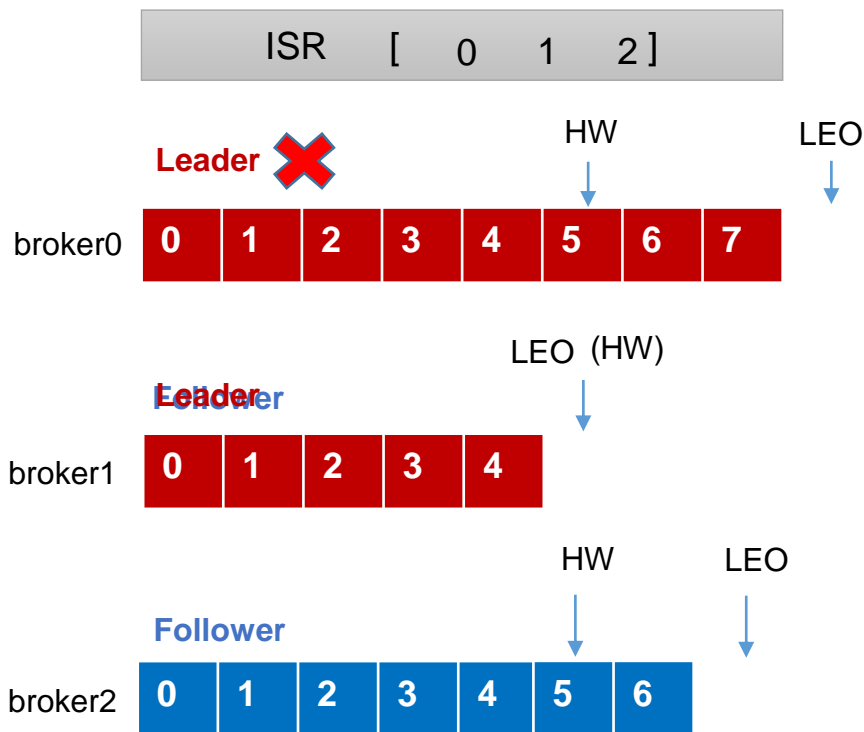
1) Follower故障

- (1) Follower发生故障后会被临时踢出ISR
- (2) 这个期间Leader和Follower继续接收数据
- (3) 待该Follower恢复后, Follower会读取本地磁盘记录的上次的HW, 并将log文件高于HW的部分截取掉, 从HW开始向Leader进行同步。
- (4) 等该**Follower**的**LEO**大于等于该**Partition**的**HW**, 即Follower追上Leader之后, 就可以重新加入ISR了。



LEO (Log End Offset) : 每个副本的最后一个offset, LEO其实就是最新的offset + 1

HW (High Watermark) : 所有副本中最小的LEO

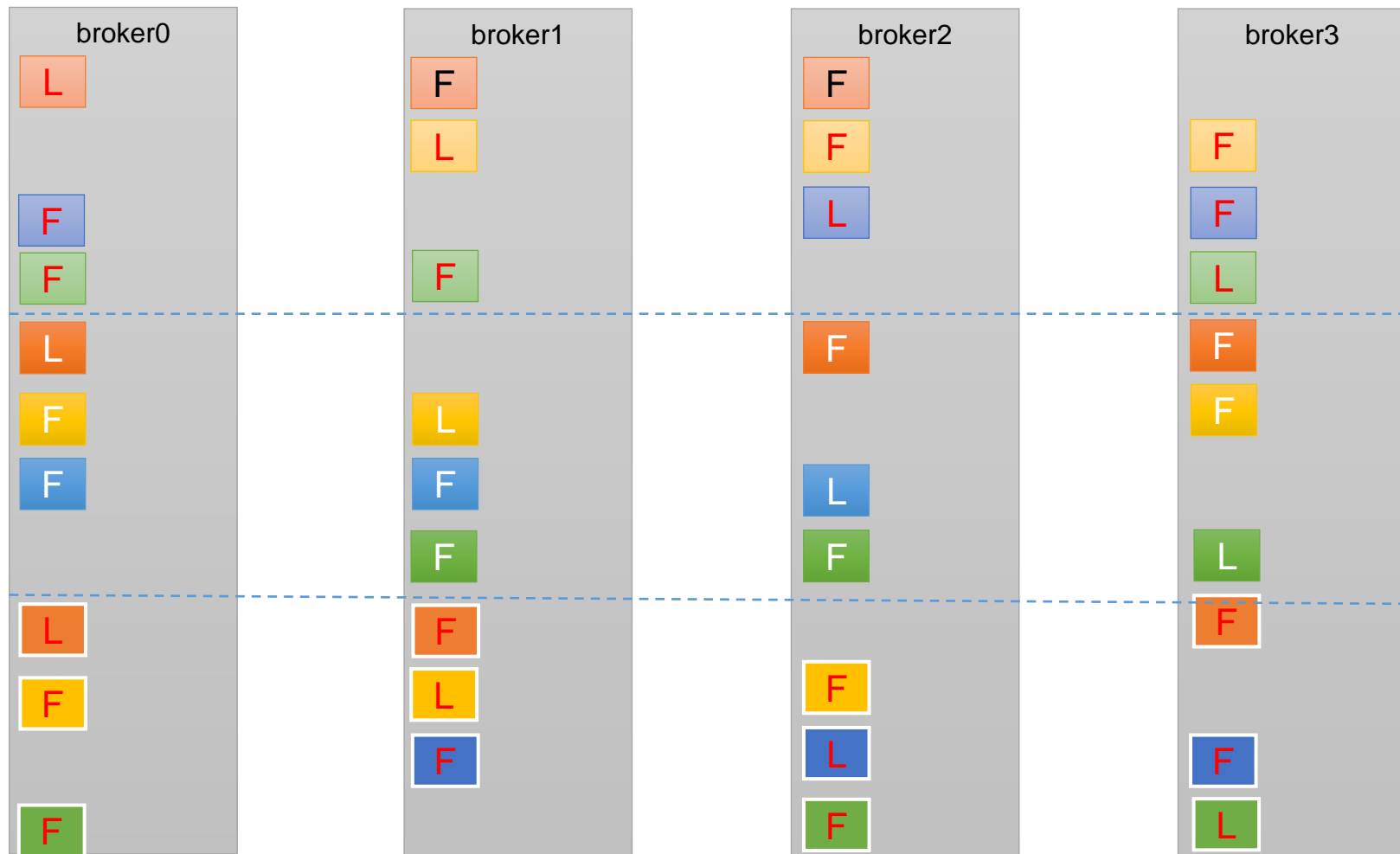


1) Leader故障

(1) Leader发生故障之后, 会从ISR中选出一个新的Leader

(2) 为保证多个副本之间的数据一致性, 其余的Follower会先将各自的log文件高于HW的部分截掉, 然后从新的Leader同步数据。

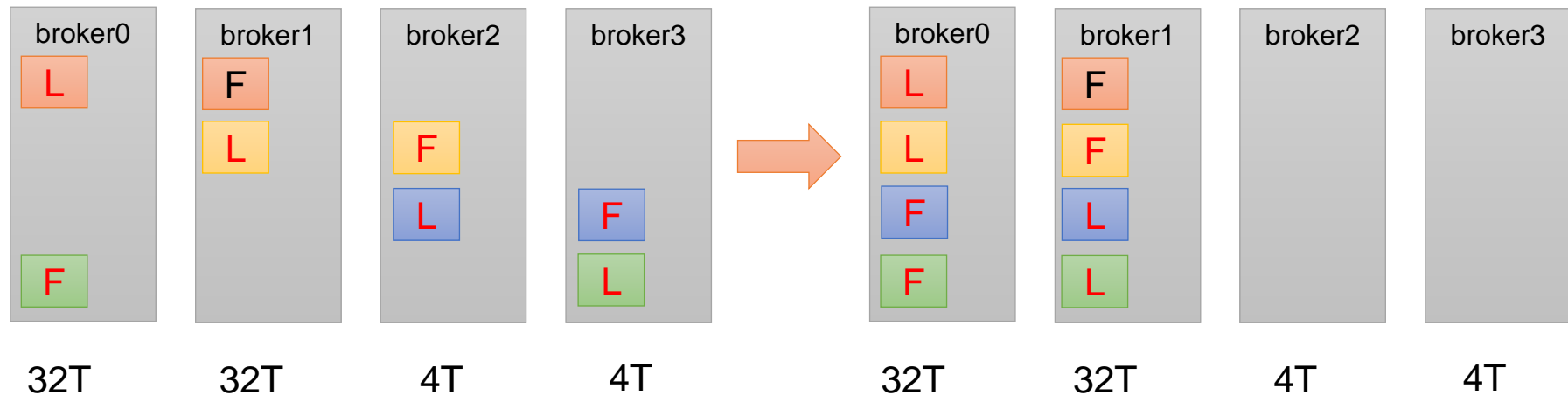
注意: 这只能保证副本之间的数据一致性, 并不能保证数据不丢失或者不重复。





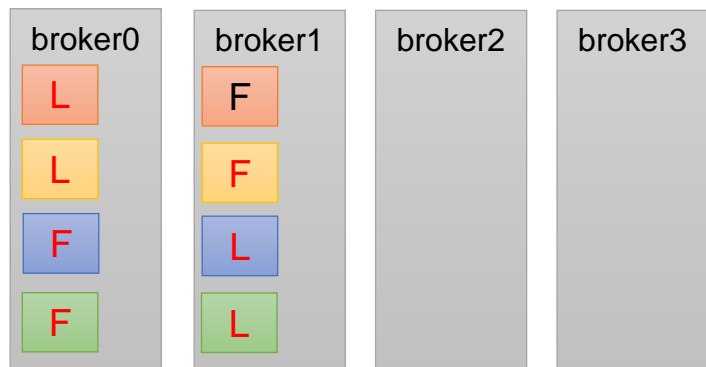
在生产环境中，每台服务器的配置和性能不一致，但是Kafka只会根据自己的代码规则创建对应的分区副本，就会导致个别服务器存储压力较大。所有需要手动调整分区副本的存储。

需求：创建一个新的topic，4个分区，两个副本，名称为three。将该topic的所有副本都存储到broker0和broker1两台服务器上。





正常情况下，Kafka本身会自动把Leader Partition均匀分散在各个机器上，来保证每台机器的读写吞吐量都是均匀的。但是如果某些broker宕机，会导致Leader Partition过于集中在其他少部分几台broker上，这会导致少数几台broker的读写请求压力过高，其他宕机的broker重启之后都是follower partition，读写请求很低，造成集群负载不均衡。



- `auto.leader.rebalance.enable`, 默认是`true`。
自动Leader Partition 平衡
- `leader.imbalance.per.broker.percentage`,
默认是`10%`。每个broker允许的不平衡的leader的比率。如果每个broker超过了这个值，控制器会触发leader的平衡。
- `leader.imbalance.check.interval.seconds`,
默认值`300秒`。检查leader负载是否平衡的间隔时间。

下面拿一个主题举例说明，假设集群只有一个主题如下图所示：

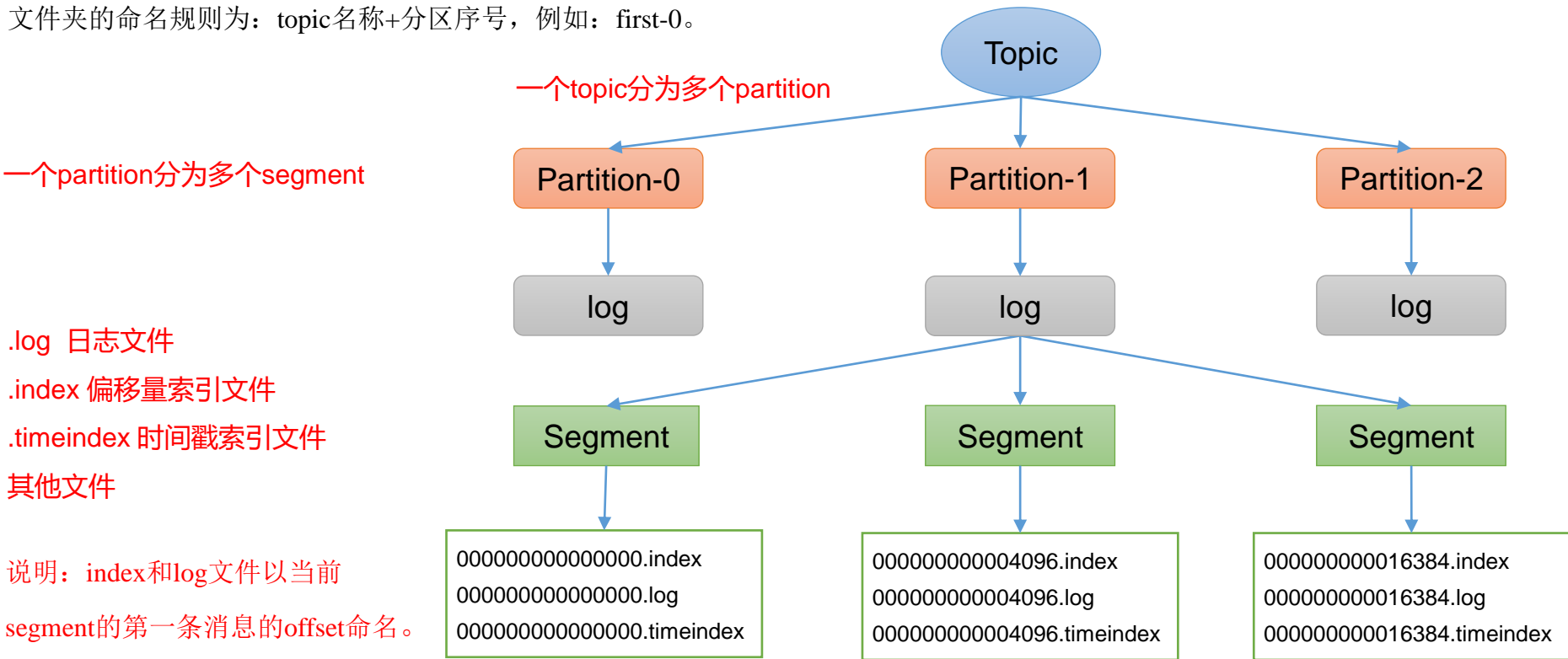
```
Topic: atguigu1 Partition: 0    Leader: 0    Replicas: 3,0,2,1    Isr: 3,0,2,1
Topic: atguigu1 Partition: 1    Leader: 1    Replicas: 1,2,3,0    Isr: 1,2,3,0
Topic: atguigu1 Partition: 2    Leader: 2    Replicas: 0,3,1,2    Isr: 0,3,1,2
Topic: atguigu1 Partition: 3    Leader: 3    Replicas: 2,1,0,3    Isr: 2,1,0,3
```

针对broker0节点，分区2的AR优先副本是0节点，但是0节点却不是Leader节点，所以不平衡数加1，AR副本总数是4
所以broker0节点不平衡率为 $1/4 > 10\%$ ，需要再平衡。

broker2和broker3节点和broker0不平衡率一样，需要再平衡。
Broker1的不平衡数为0，不需要再平衡



Topic是逻辑上的概念，而partition是物理上的概念，**每个partition对应于一个log文件**，该log文件中存储的就是Producer生产的数据。**Producer**生产的数据会被不断**追加**到该log文件末端，为防止log文件过大导致数据定位效率低下，Kafka采取了**分片**和**索引**机制，将**每个partition分为多个segment**。**每个segment包括**：“.index”文件、“.log”文件和.timeindex等文件。这些文件位于一个文件夹下，该文件夹的命名规则为：topic名称+分区序号，例如：first-0。





Log文件和Index文件详解

Segment-0 [offset:0-521]

00000000000000000000.index

00000000000000000000.log

Segment-1 [offset:522-1004]

00000000000000000522.index

00000000000000000522.log

绝对Offset	相对Offset	Position
587	65	6410
639	117	13795
691	169	21060
743	221	28367

RecordBatch[baseOffset~lastOffset]		
baseOffset	lastOffset	position
522	522	0
523	523	200
524	536	819
537	562	3092
563	587	6410
588	613	10090
614	639	13795
640	665	17481
666	691	21060
692	717	24781
728	743	28367

如何在log文件中定位到
offset=600的Record?

- 1.根据目标offset定位Segment文件
- 2.找到小于等于目标offset的最大offset对应的索引项
- 3.定位到log文件
- 4.向下遍历找到目标Record

注意:

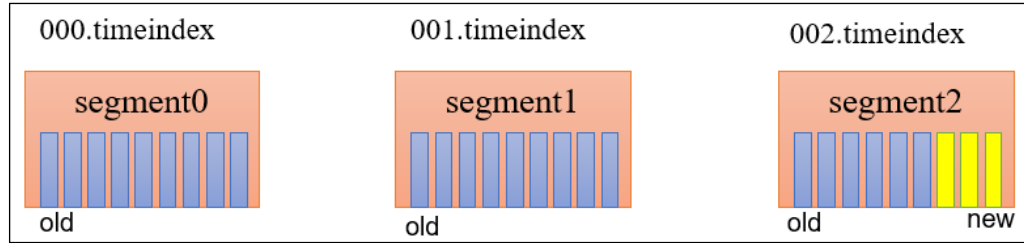
1.index为稀疏索引, 大约每往log文件写入4kb数据, 会往index文件写入一条索引。
参数log.index.interval.bytes默认4kb。

2.Index文件中保存的offset为相对offset, 这样能确保offset的值所占空间不会过大,
因此能将offset的值控制在固定大小

Segment-2 [offset:1005-]

00000000000000001005.index

00000000000000001005.log





compact日志压缩：对于相同key的不同value值，只保留最后一个版本。

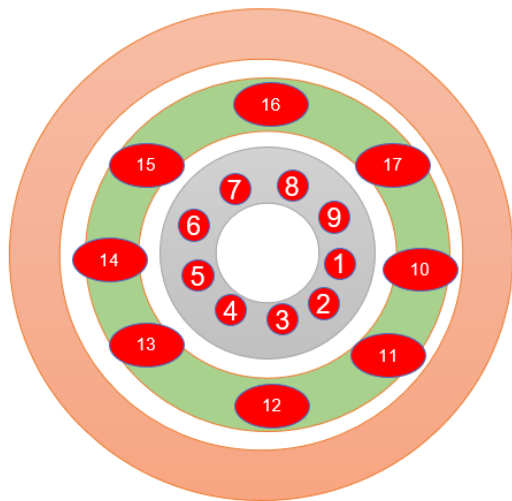
- log.cleanup.policy = compact 所有数据启用压缩策略

压缩之前的数据	Offset	0	1	2	3	4	5	6	7	8
	key	K1	K2	K1	K1	K3	K4	K5	K5	K2
	value	V1	V2	V3	V4	V5	V6	V7	V8	V9

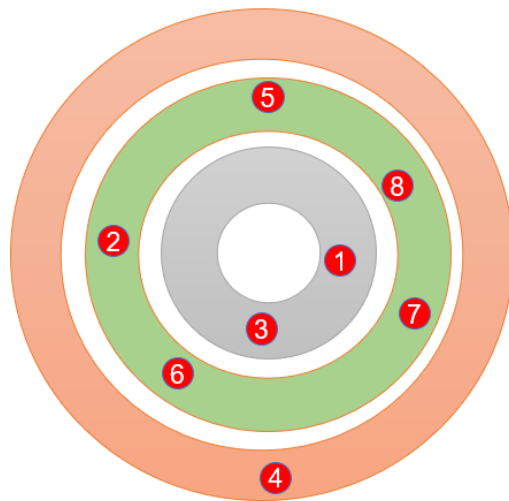
压缩之后的数据	Offset	3	4	5	7	8
	keys	k1	K3	K4	K5	K2
	values	V4	V5	V6	V8	V9

压缩后的offset可能是不连续的，比如上图中没有6，当从这些offset消费消息时，将会拿到比这个offset大的offset对应的消息，实际上会拿到offset为7的消息，并从这个位置开始消费。

这种策略只适合特殊场景，比如消息的key是用户ID，value是用户的资料，通过这种压缩策略，整个消息集里就保存了所有用户最新的资料。



顺序写磁盘能到600M/s



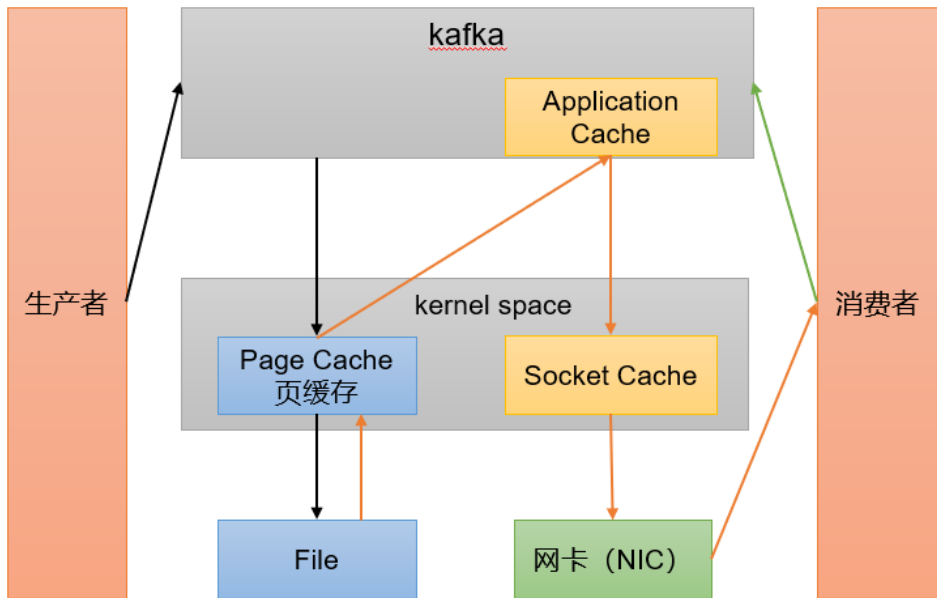
随机写磁盘能到100K/s



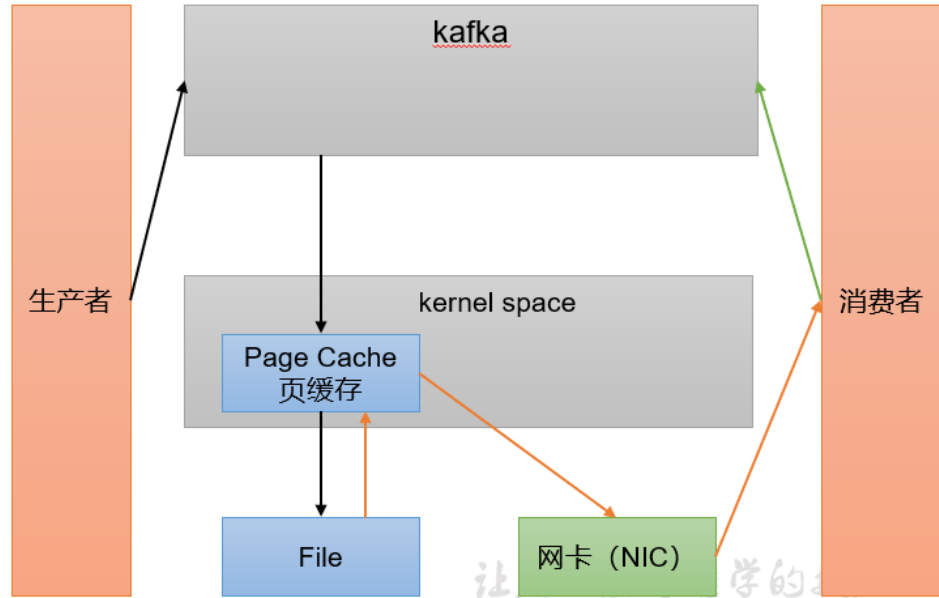
零拷贝：Kafka的数据加工处理操作交由Kafka生产者和Kafka消费者处理。**Kafka Broker**应用层不关心存储的数据，所以就不用走应用层，传输效率高。

PageCache页缓存：Kafka重度依赖底层操作系统提供的PageCache功能。当上层有写操作时，操作系统只是将数据写入PageCache。当读操作发生时，先从PageCache中查找，如果找不到，再去磁盘中读取。实际上PageCache是把尽可能多的空闲内存都当做了磁盘缓存来使用。

非零拷贝工作流程（假设）



零拷贝工作流程（实际）





➤ pull（拉）模式：

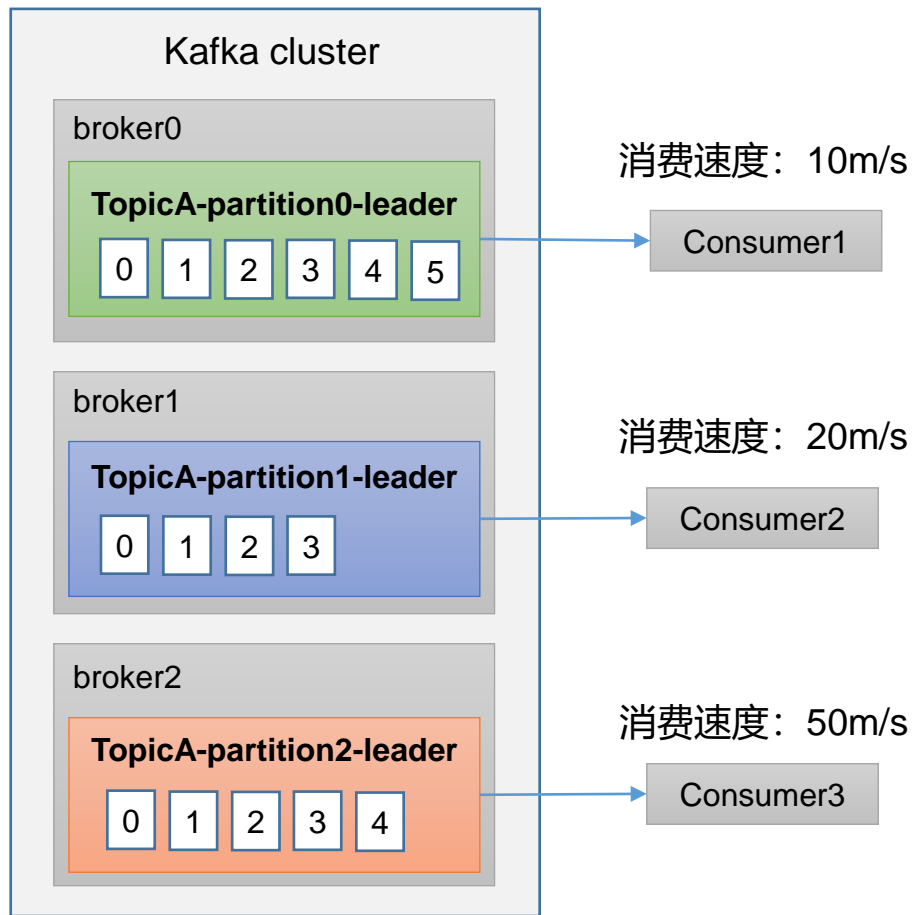
consumer采用从broker中主动拉取数据。

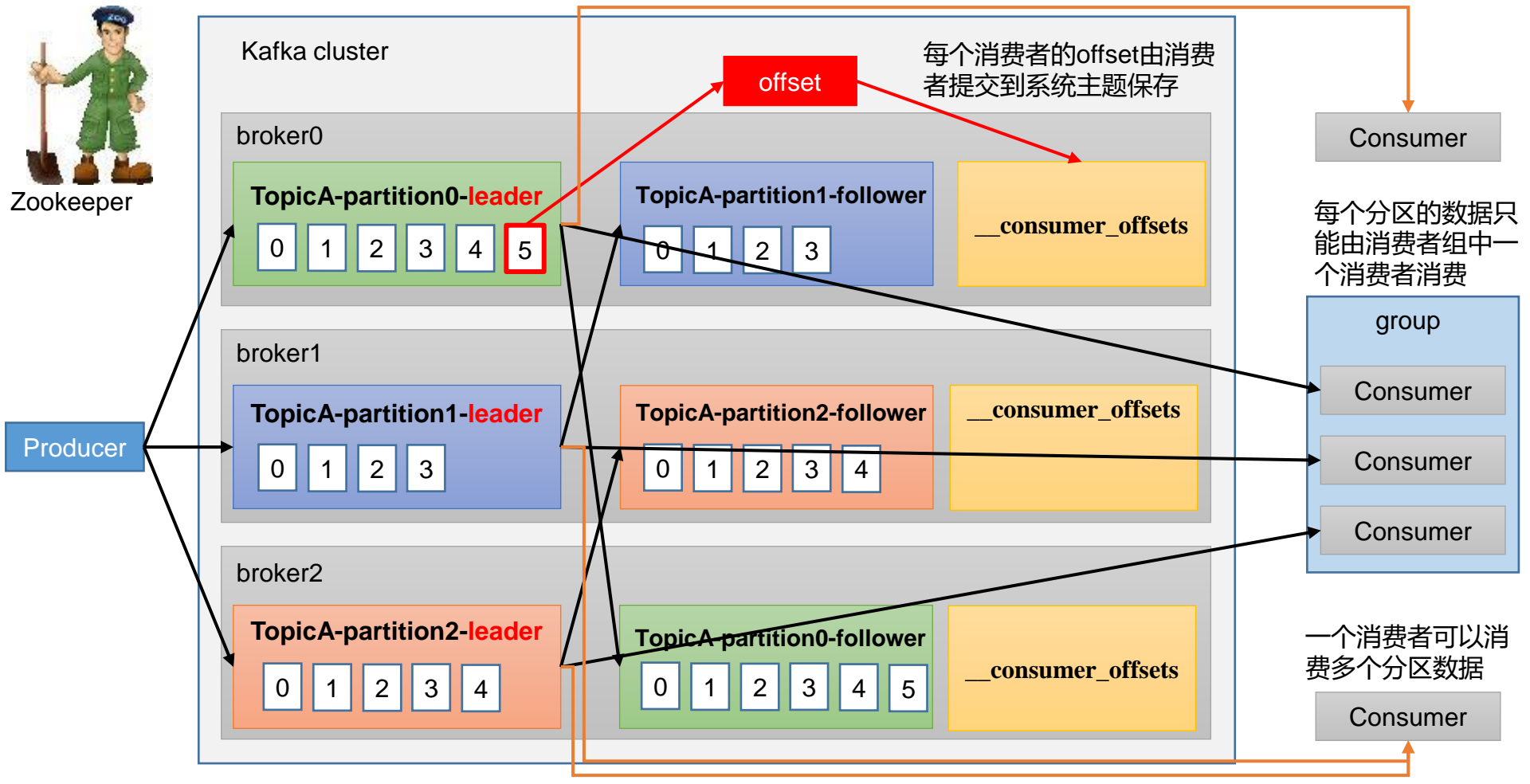
Kafka采用这种方式。

➤ push（推）模式：

Kafka没有采用这种方式，因为由broker决定消息发送速率，很难适应所有消费者的消费速率。例如推送的速度是50m/s，Consumer1、Consumer2就来不及处理消息。

pull模式不足之处是，如果Kafka没有数据，消费者可能会陷入循环中，一直返回空数据。

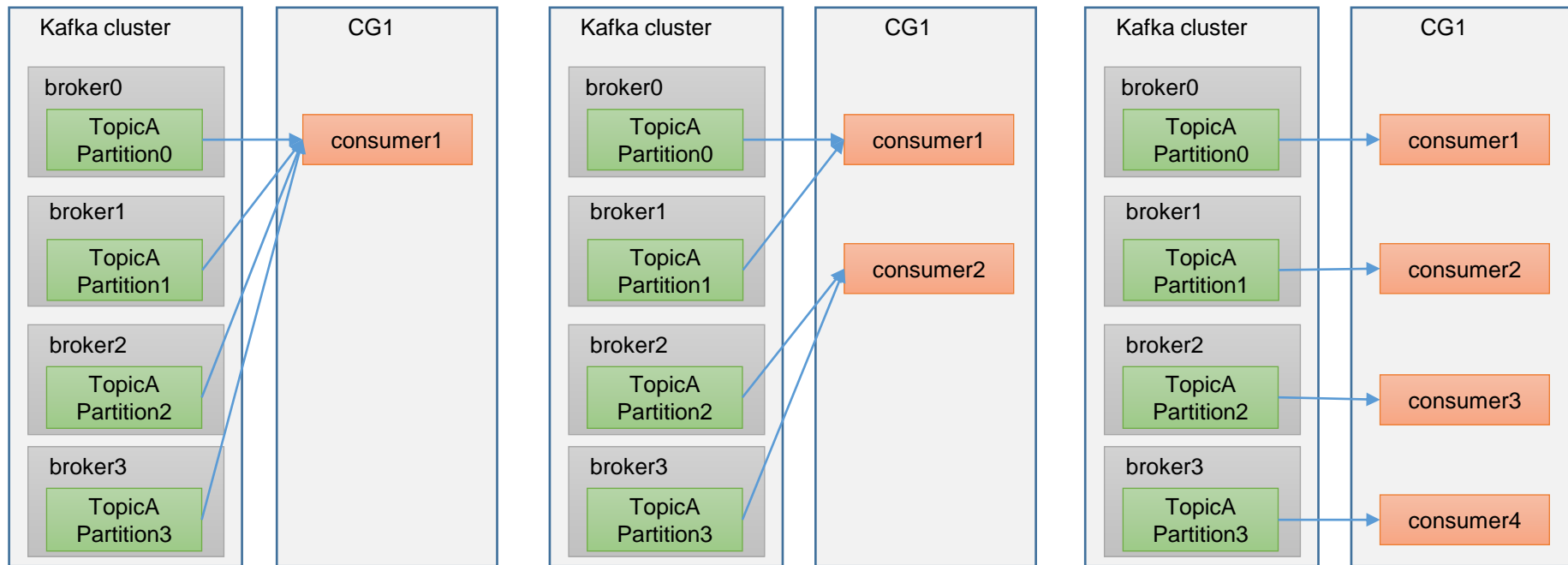


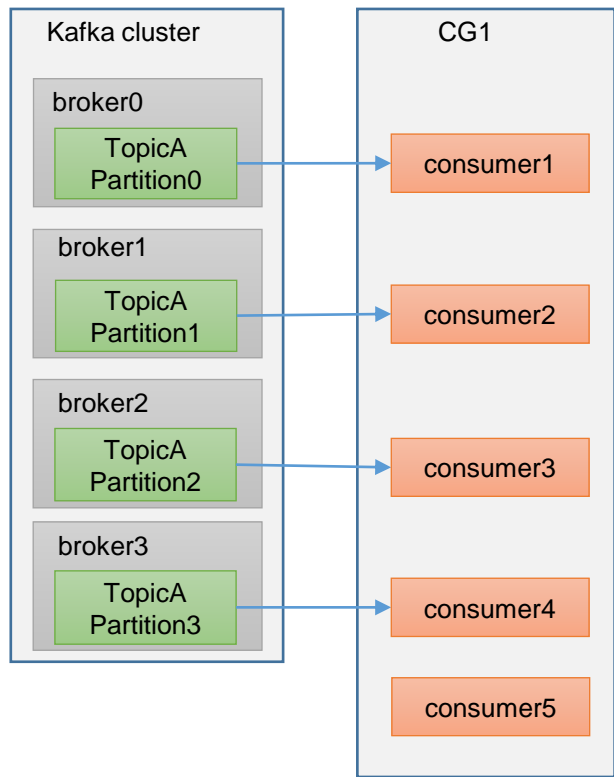




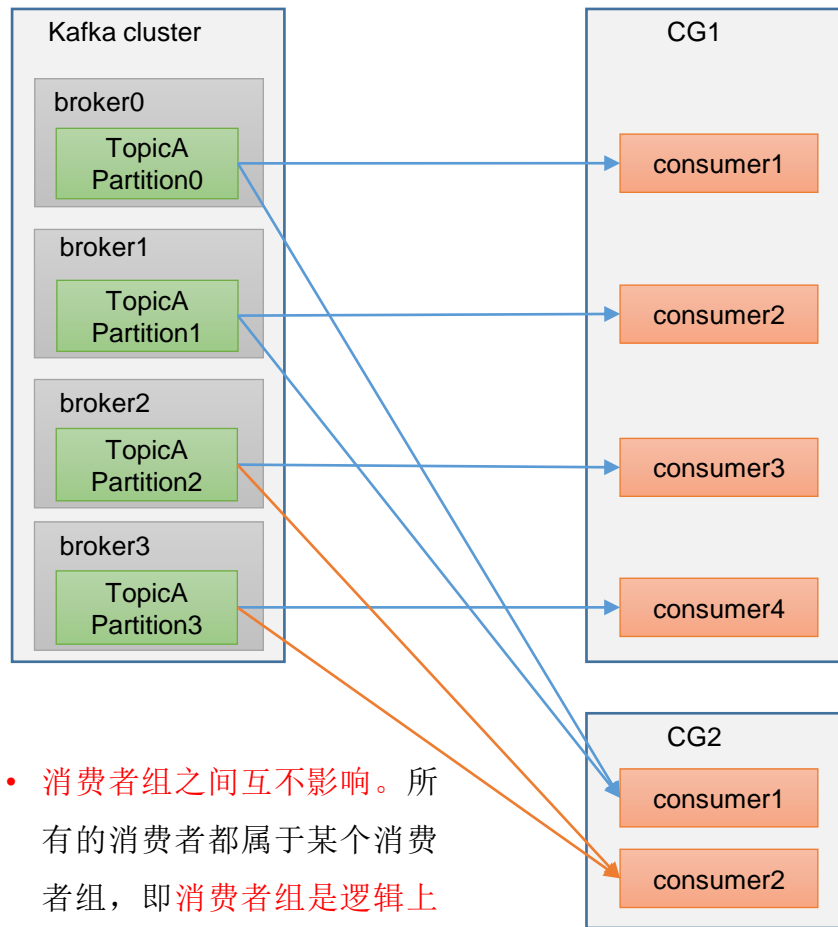
Consumer Group (CG)：消费者组，由多个consumer组成。形成一个消费者组的条件，是所有消费者的groupid相同。

- 消费者组内每个消费者负责消费不同分区的数据，一个分区只能由一个组内消费者消费。
- 消费者组之间互不影响。所有的消费者都属于某个消费者组，即消费者组是逻辑上的一个订阅者。





- 如果向消费组中添加更多的消费者，超过主题分区数量，则有一部分消费者就会闲置，不会接收任何消息。



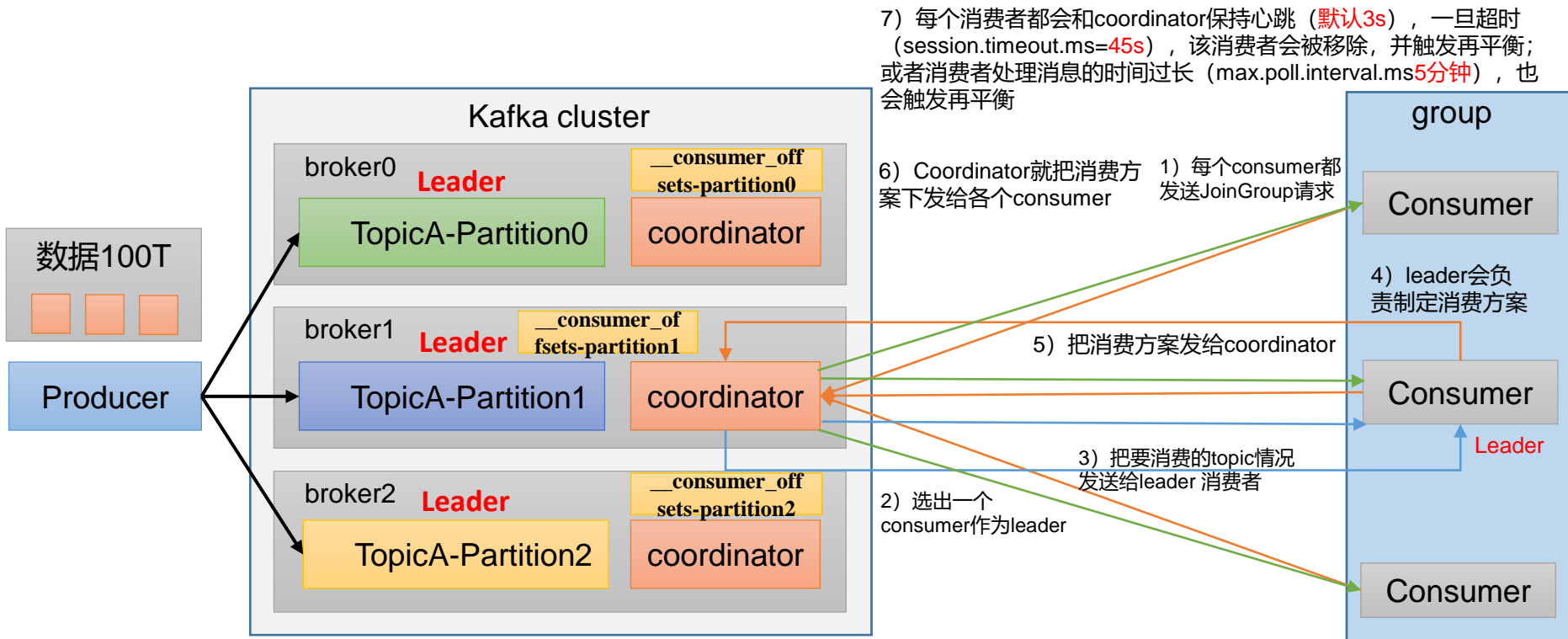
- 消费者组之间互不影响。所有的消费者都属于某个消费者组，即消费者组是逻辑上的一个订阅者。

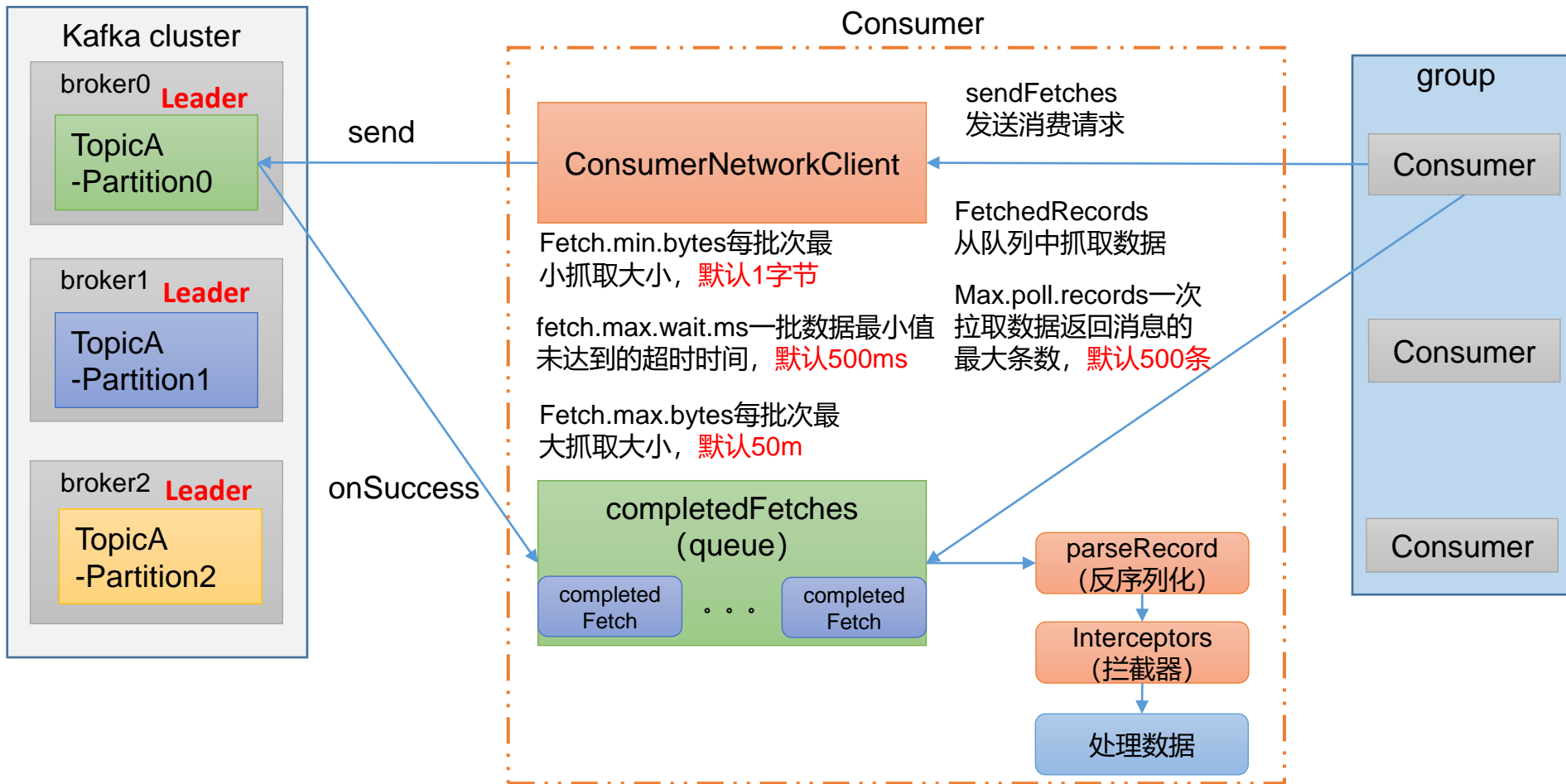


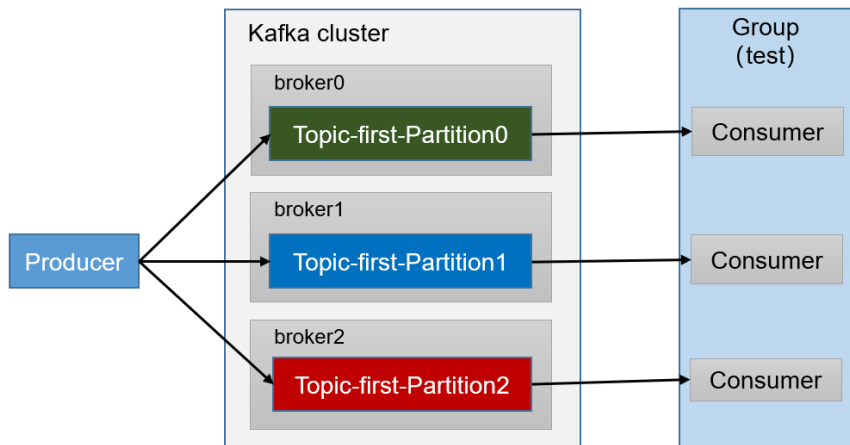
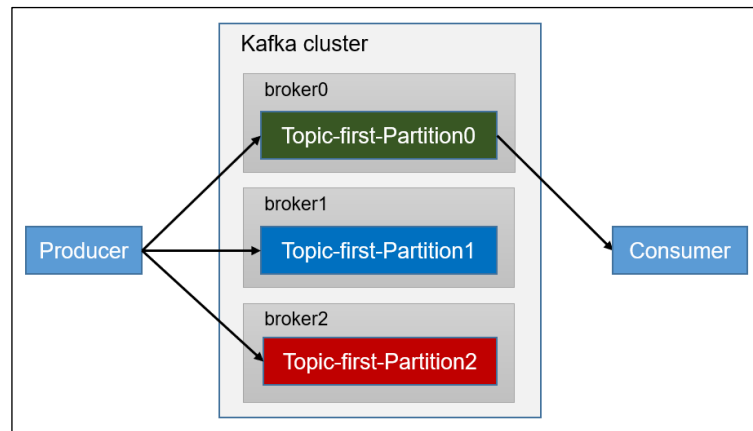
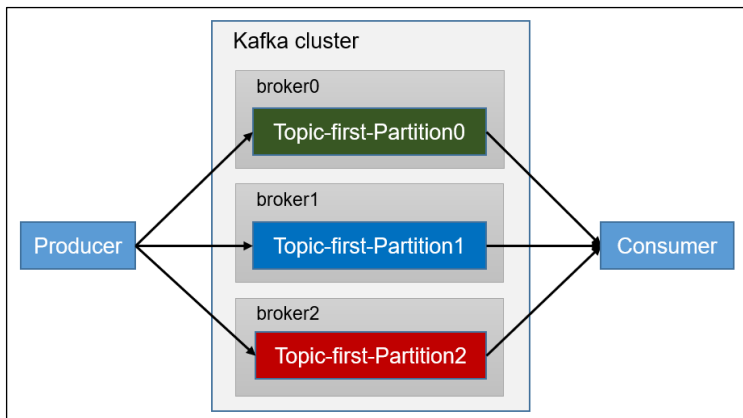
1、coordinator：辅助实现消费者组的初始化和分区的分配。

coordinator节点选择 = $\text{groupid的hashCode值} \% 50$ ($__\text{consumer_offsets}$ 的分区数量)

例如：groupid的hashCode值 = 1, $1 \% 50 = 1$, 那么 $__\text{consumer_offsets}$ 主题的1号分区, 在哪个broker上, 就选择这个节点的coordinator作为这个消费者组的老大。消费者组下的所有的消费者提交offset的时候就往这个分区去提交offset。







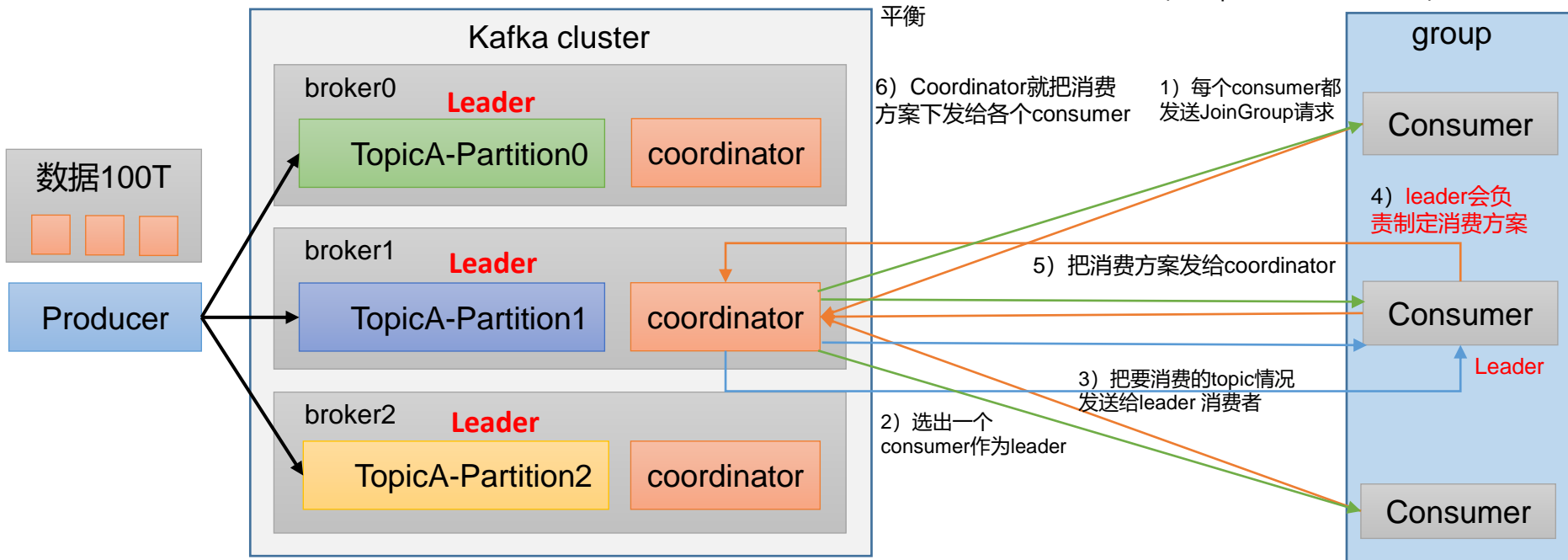


1、一个consumer group中有多个consumer组成，一个 topic有多个partition组成，现在的问题是，到底由哪个consumer来消费哪个partition的数据。

2、Kafka有四种主流的分区分配策略：Range、RoundRobin、Sticky、CooperativeSticky。

可以通过配置参数`partition.assignment.strategy`，修改分区的分配策略。默认策略是Range + CooperativeSticky。Kafka可以同时使用多个分区分配策略。

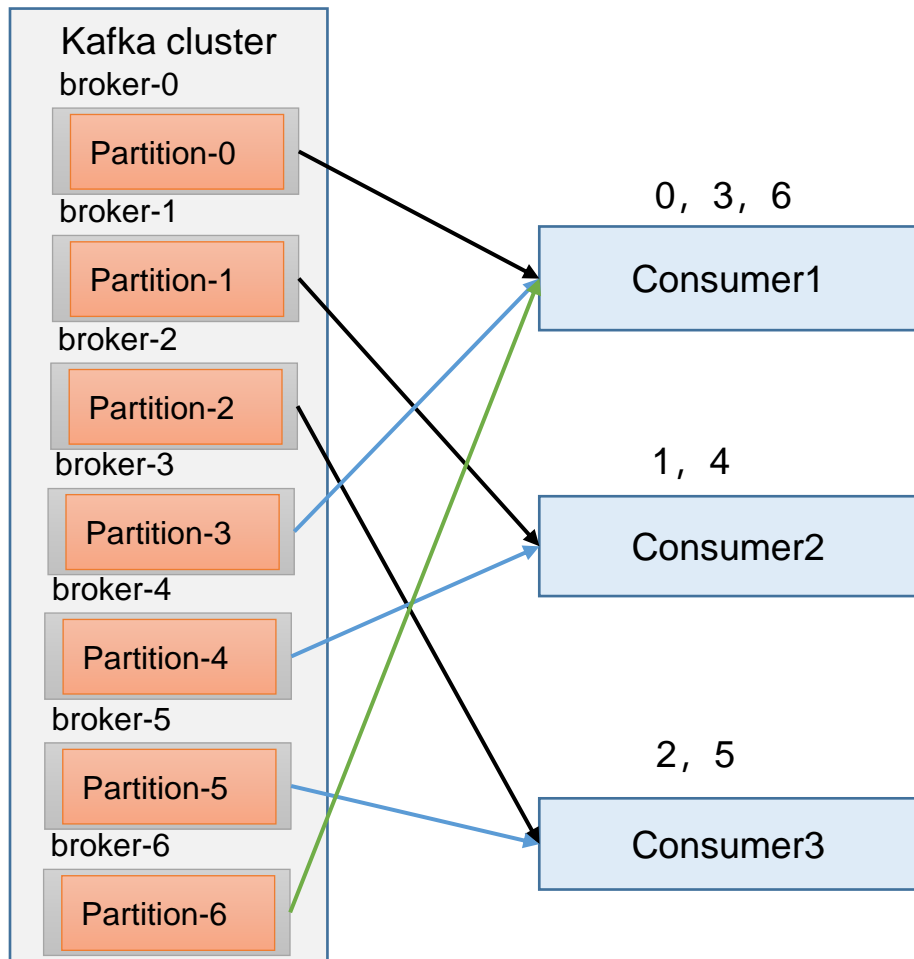
7) 每个消费者都会和coordinator保持心跳（默认3s），一旦超时（`session.timeout.ms=45s`），该消费者会被移除，并触发再平衡；或者消费者处理消息的过长（`max.poll.interval.ms5分钟`），也会触发再平衡





RoundRobin 针对集群中**所有Topic**而言。

RoundRobin 轮询分区策略，是把**所有的 partition 和所有的 consumer 都列出来**，然后**按照 hashCode 进行排序**，最后通过**轮询算法**来分配 partition 给到各个消费者。

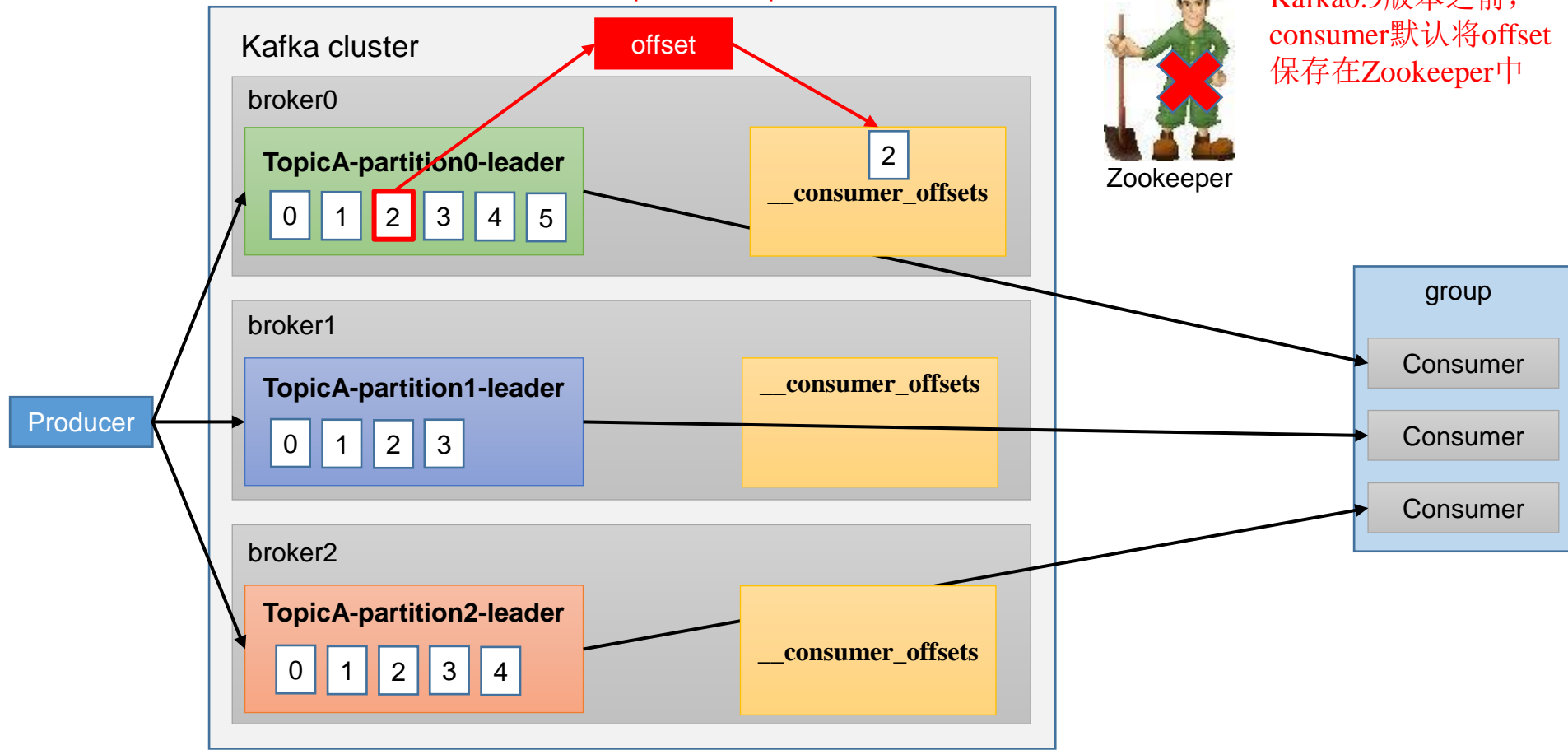




从0.9版本开始, consumer默认将offset保存在Kafka
一个内置的topic中, 该topic为__consumer_offsets



Kafka0.9版本之前,
consumer默认将offset
保存在Zookeeper中

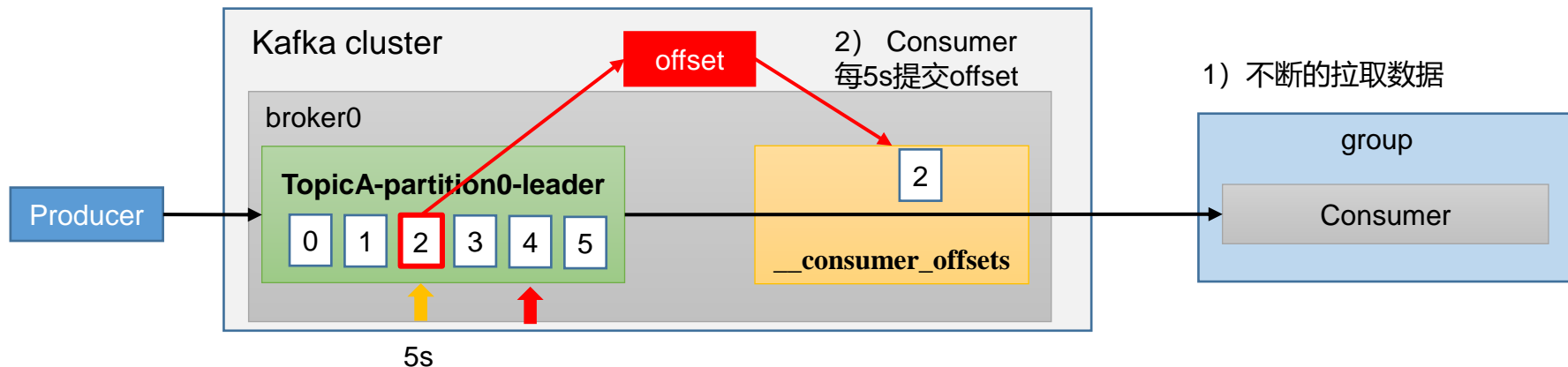




为了使我们能够专注于自己的业务逻辑，Kafka提供了自动提交offset的功能。

自动提交offset的相关参数：

- **enable.auto.commit**：是否开启自动提交offset功能，默认是true
- **auto.commit.interval.ms**：自动提交offset的时间间隔，默认是5s

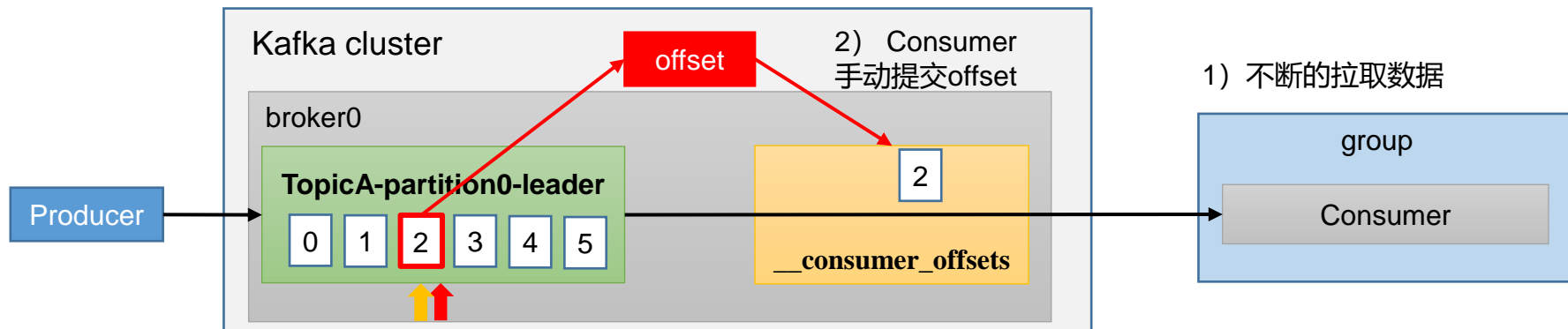


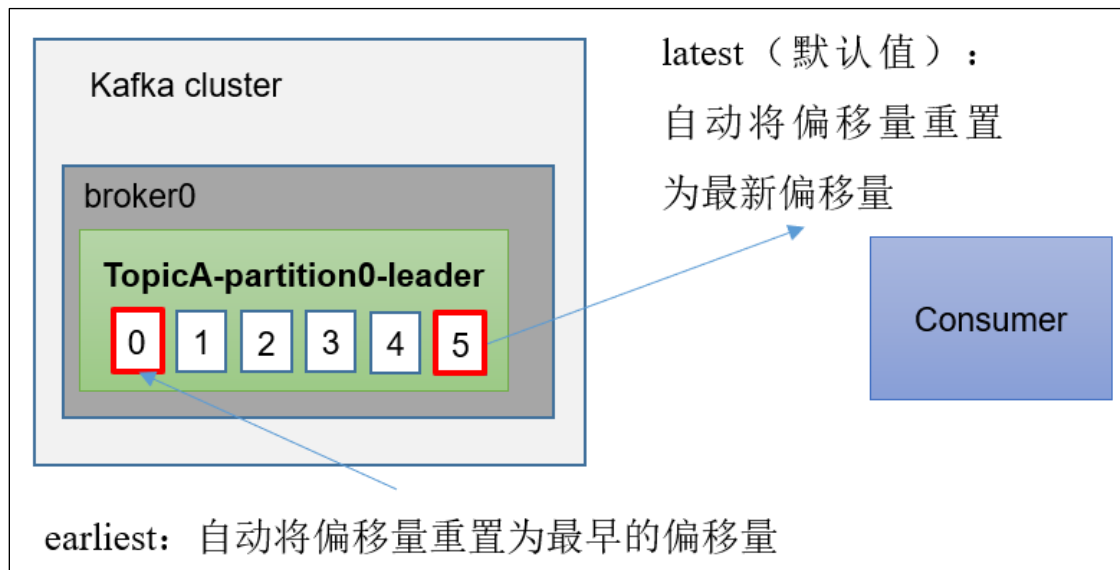


虽然自动提交offset十分简单便利，但由于其是基于时间提交的，开发人员难以把握offset提交的时机。因此Kafka还提供了手动提交offset的API。

手动提交offset的方法有两种：分别是`commitSync`（同步提交）和`commitAsync`（异步提交）。两者的相同点是，都会将本次提交的一批数据最高的偏移量提交；不同点是，同步提交阻塞当前线程，一直到提交成功，并且会自动失败重试（由不可控因素导致，也会出现提交失败）；而异步提交则没有失败重试机制，故有可能提交失败。

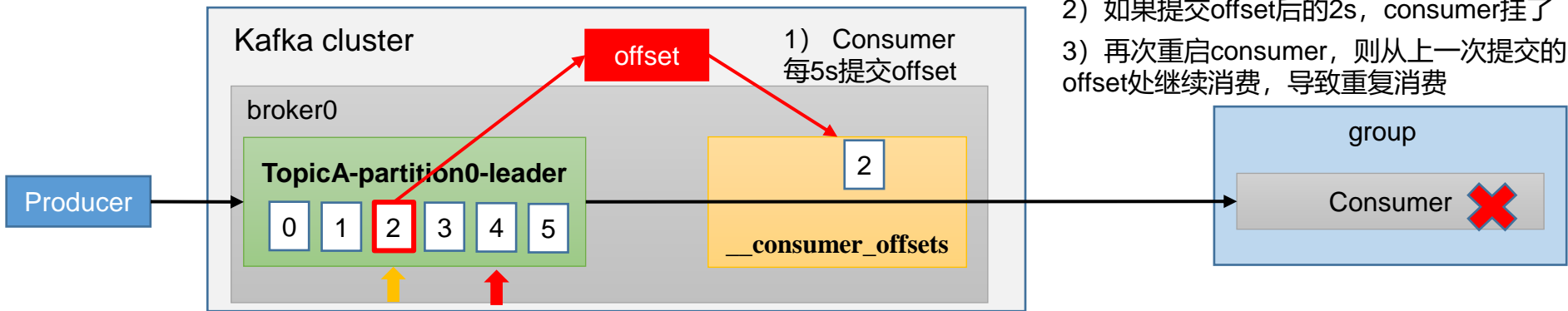
- `commitSync`（同步提交）：必须等待offset提交完毕，再去消费下一批数据。
- `commitAsync`（异步提交）：发送完提交offset请求后，就开始消费下一批数据了。



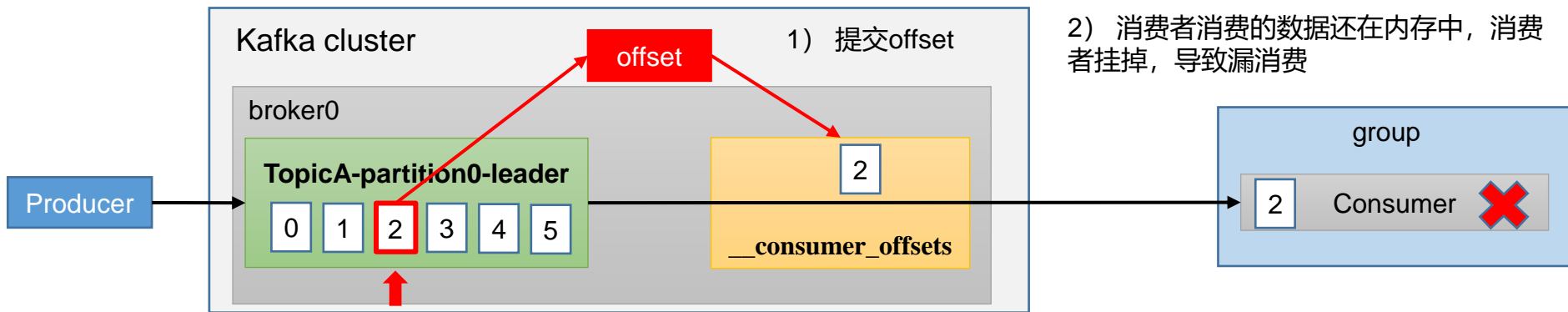




(1) 场景1: **重复消费**。自动提交offset引起。

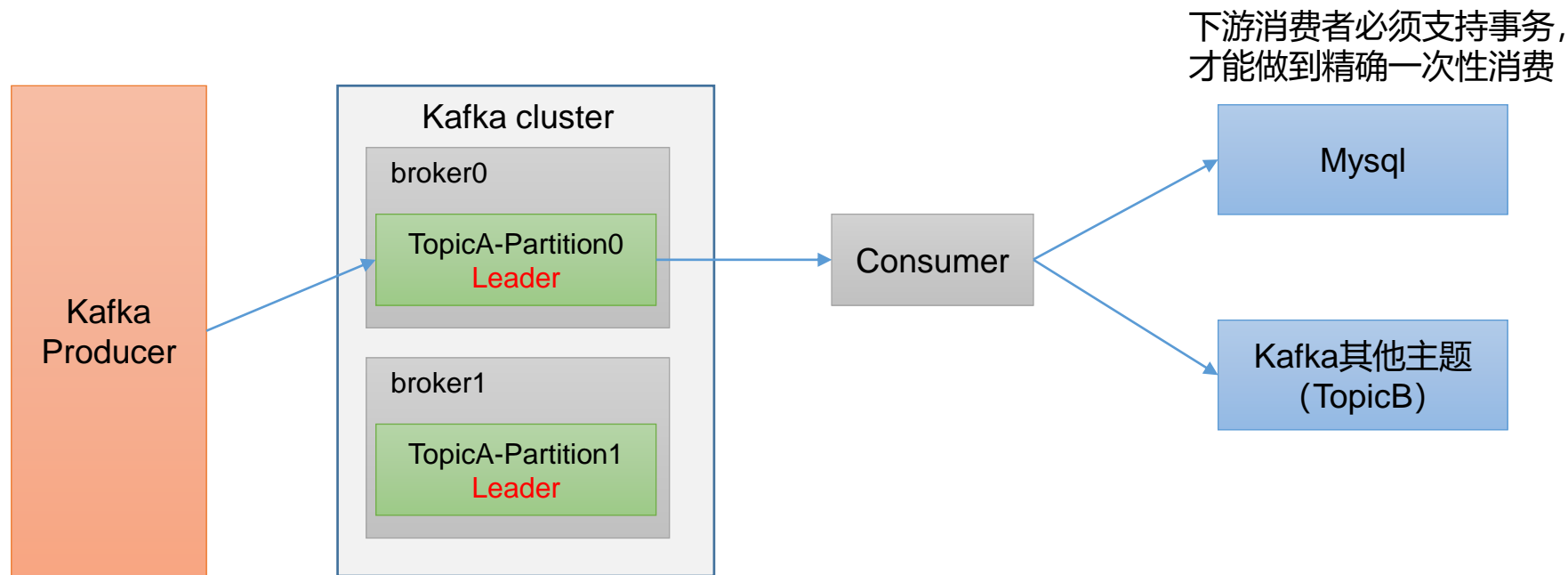


(2) 场景1: **漏消费**。设置offset为手动提交, 当offset被提交时, 数据还在内存中未落盘, 此时刚好消费者线程被kill掉, 那么offset已经提交, 但是数据未处理, 导致这部分内存中的数据丢失。

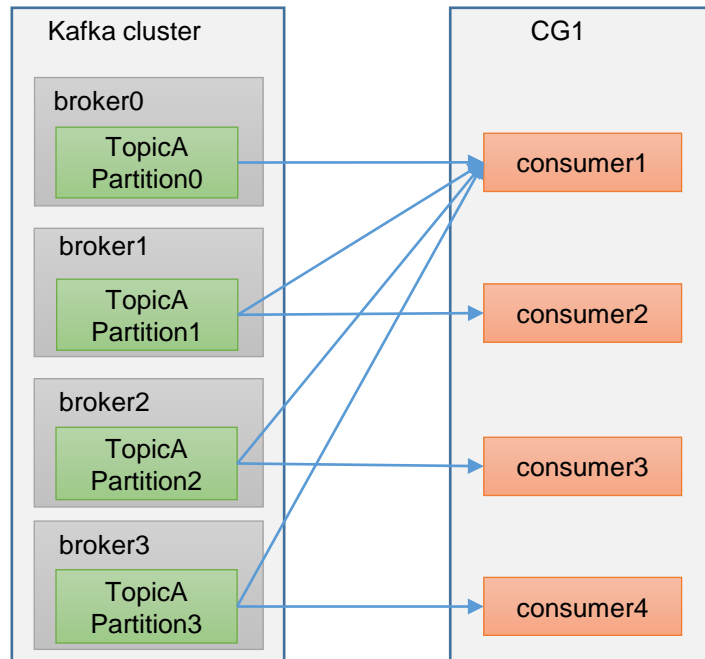




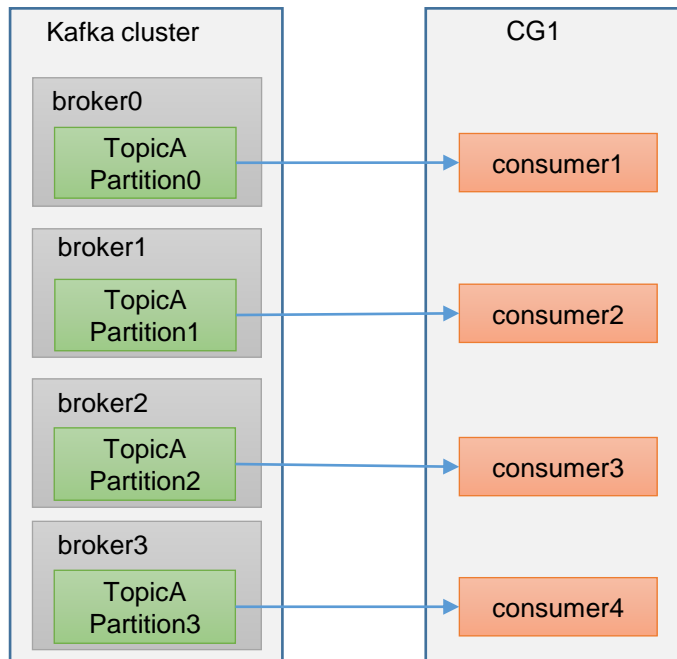
如果想完成Consumer端的精准一次性消费，那么需要Kafka消费端将消费过程和提交offset过程做原子绑定。此时我们需要将Kafka的offset保存到支持事务的自定义介质（比如MySQL）。这部分知识会在后续项目部分涉及。



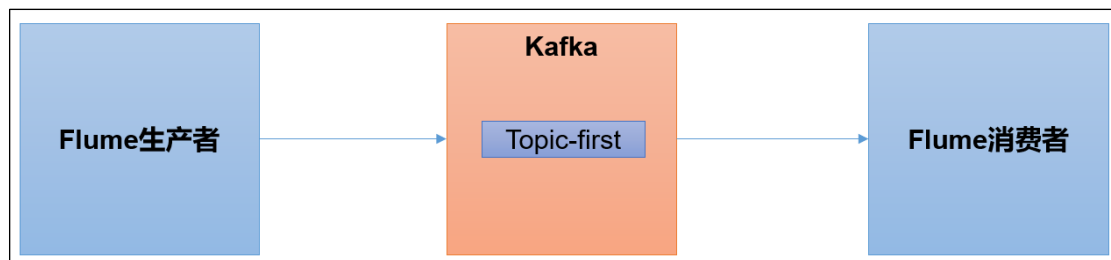
1) 如果是Kafka消费能力不足，则可以考虑增加Topic的分区数，并且同时提升消费组的消费者数量，消费者数 = 分区数。（两者缺一不可）

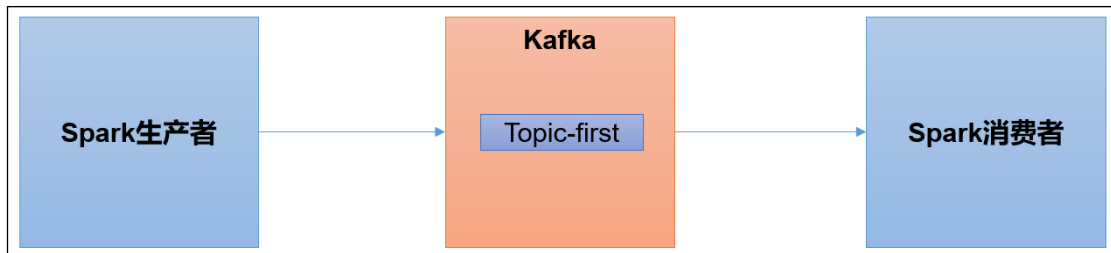
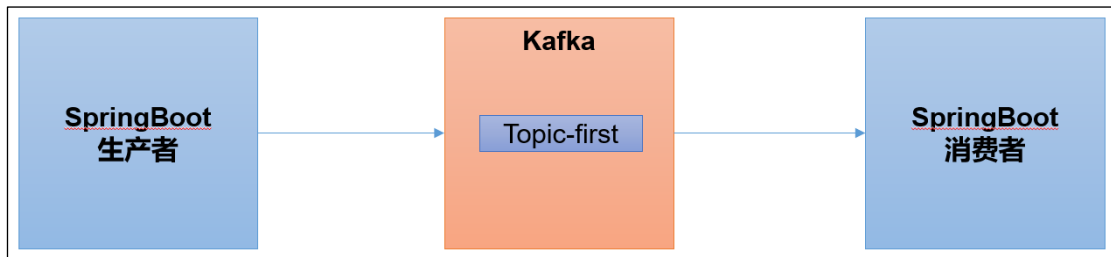
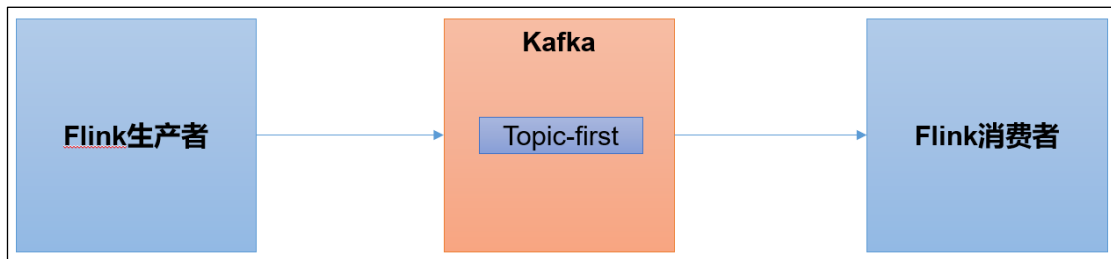


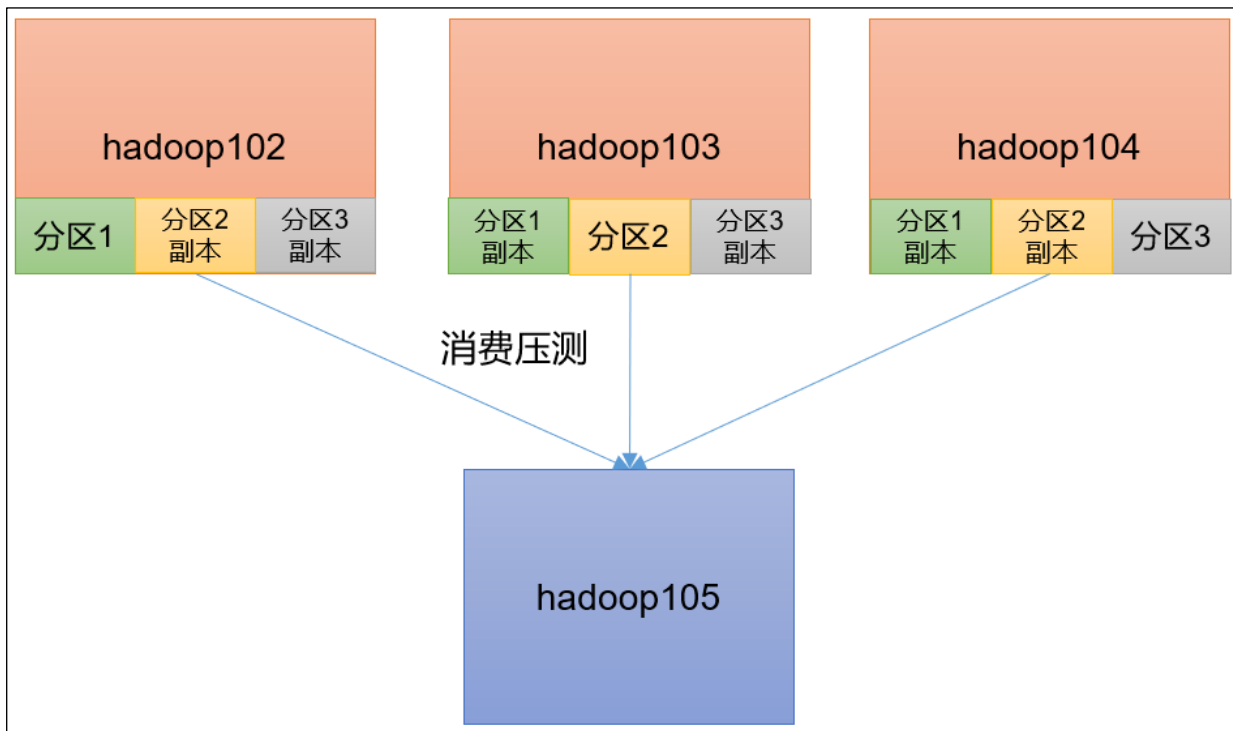
2) 如果是下游的数据处理不及时：提高每批次拉取的数量。批次拉取数据过少（拉取数据/处理时间 < 生产速度），使处理的数据小于生产的数据，也会造成数据积压。



从一次最多拉取500条，调整为一次最多拉取1000条









用户自定义生产者
CustomProducer.java

main

创建Kafka生产者对象
new KafkaProducer

连续点击三次this构造器

获取事务id
transactionalId

获取客户端id
clientId

监控相关配置
new JmxReporter()

分区器配置
this.partition器

序列化配置
keySerializer
valueSerializer

拦截器配置
interceptorList

单条信息的最大值,
默认1m
maxRequestSize

缓存大小, 默认32m
totalMemorySize

创建缓存队列
new RecordAccumulator()

连接kafka集群
BOOTSTRAP_SERVERS
_CONFIG

从Kafka集群获取元数据
this.metadata

创建sender线程
new Sender

启动发送线程
this.ioThread.start()

批次大小, 默认16k
BATCH_SIZE_CONFIG

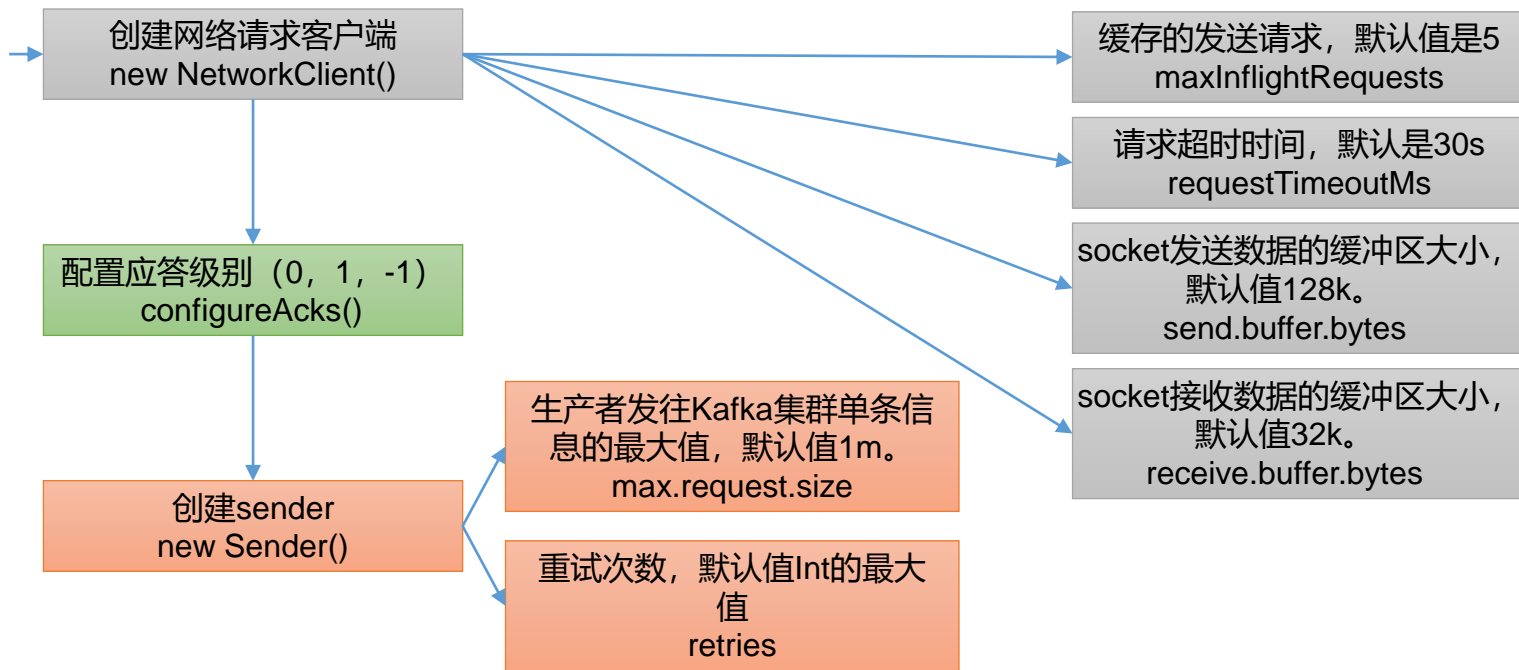
是否压缩, 默认none
compressionType

linger.ms, 默认值0
lingerMs()

重试间隔时间, 默认值
100ms
retryBackoffMs



1、初始化sender线程 `this.sender = new Sender();`



2、启动sender线程

```
this.ioThread = new  
KafkaThread(ioThreadName,  
this.sender, true);  
this.ioThread.start();
```



```
@Override  
public void run() {  
    while (running) {  
        // sender线程从缓冲区准备拉取数据, 刚启动拉不到数据  
        runOnce();  
    }  
}
```

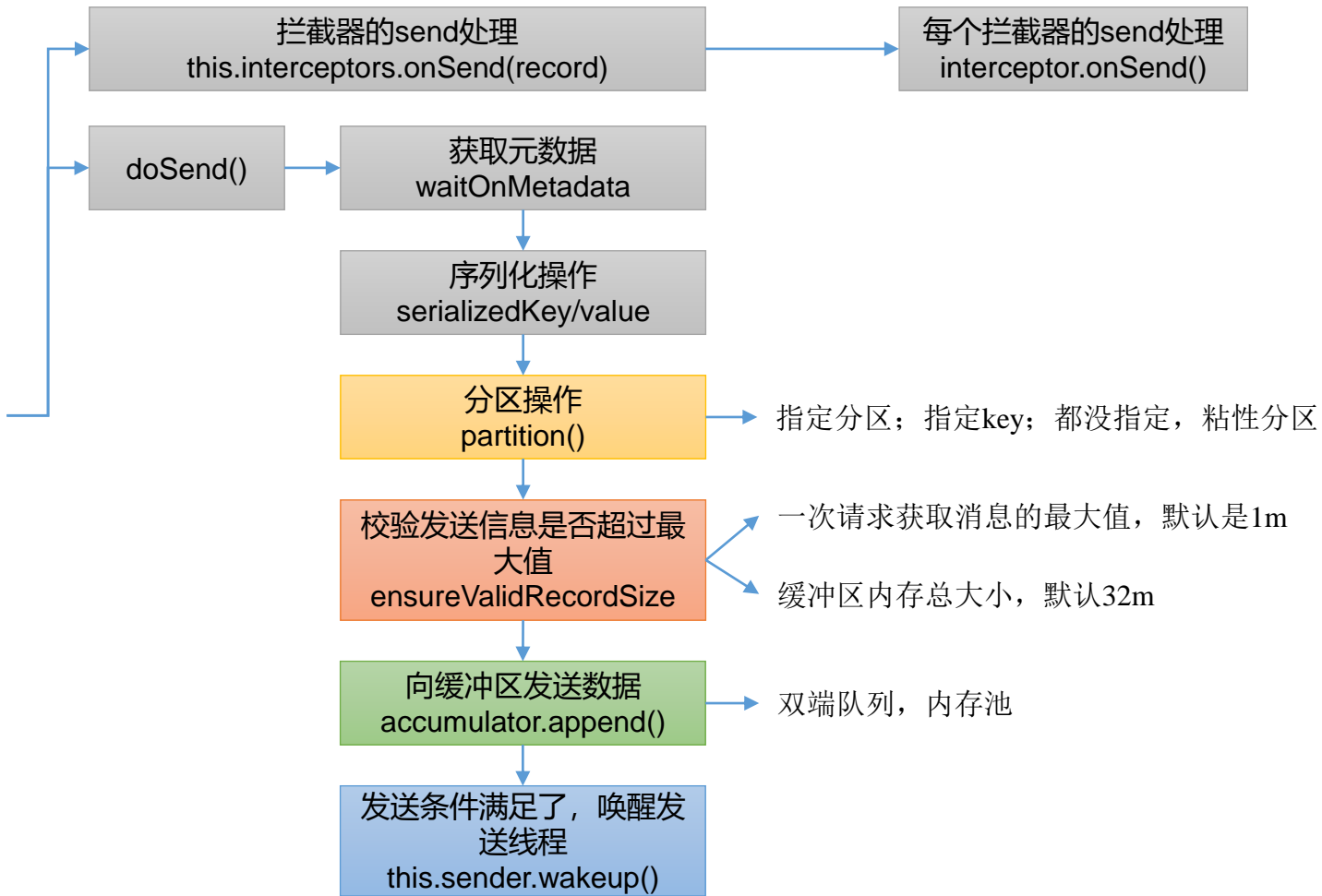


用户自定义生产者
CustomProducer.java

main

创建Kafka生产者对象
new KafkaProducer

发送数据
kafkaProducer.send()





唤醒发送请求
`this.sender.wakeup();`

```
@Override  
public void run() {  
    while (running) {  
        runOnce();  
    }  
}
```

将准备好的数据发
送到服务器端
`sendProducerData()`

等待发送响应
`client.poll()`

获取元数据
`metadata.fetch()`

检查32m缓存是否准备好
`this.accumulator.ready()`

lingerMs时间到

发往同一个broker节点的数
据，打包为一个请求批次
`this.accumulator.drain()`

`batches.put(node.id(), ready);`

发送请求
`sendProduceRequests()`

`client.newClientRequest`

`client.send()` → **doSend()**

`this.inFlightRequests.add(inFlightRequest);`

`selector.send()`

`this.selector.poll()`



用户自定义消费者
CustomConsumer.java

main

创建Kafka消费者对象
new KafkaConsumer

连续点击三次this构造器

获取消费者组
this.groupId

获取客户端id
this.clientId

拦截器配置
interceptorList

key和value反序列化
keyDeserializer/valueDes
erializer

offset从什么位置开始消费
offsetResetStrategy

获取元数据
this.metadata

连接Kafka集群
BOOTSTRAP_SERVERS
_CONFIG

心跳时间,默认3s
heartbeatIntervalMs

创建网络客户端
new NetworkClient()

创建一个消费者客户端
new
ConsumerNetworkClient()

获取消费者分区分配策略
this.assignors

创建消费者协调器
new
ConsumerCoordinator()

抓取数据配置
new Fetcher<>()



用户自定义消费者
CustomConsumer.java



main



创建Kafka消费者对象
new KafkaConsumer



订阅主题

```
ArrayList<String> topics = new ArrayList<>();  
topics.add("first");  
kafkaConsumer.subscribe(topics);
```

判断是否需要更改订阅主题，如果需要更改主题，则更新元数据信息
`this.subscriptions.subscribe(new HashSet<>(topics), listener)`

注册负载均衡监听（例如消费者组中，其他消费者退出触发再平衡）
`registerRebalanceListener(listener)`

修改订阅主题信息
`changeSubscription(topics)`

如果订阅的主题和以前订阅的一致，就不需要修改订阅信息。如果不一致，就需要修改

如果订阅的和以前不一致，需要更新元数据信息
`metadata.requestUpdateForNewTopics()`



用户自定义消费者
CustomConsumer.java

main

1 创建Kafka消费者对象
new KafkaConsumer

2 订阅主题

```
ArrayList<String> topics = new  
ArrayList<>();  
topics.add("first");  
kafkaConsumer.subscribe(topics);
```

3 消费数据

```
while (true){
```

```
    ConsumerRecords<String, String> consumerRecords = kafkaConsumer.poll(Duration.ofSeconds(1));  
    for (ConsumerRecord<String, String> consumerRecord : consumerRecords) {  
        System.out.println(consumerRecord);  
    }  
}
```

消费者or消费者组初始化

updateAssignmentMetadataIfNeeded()

协调器交互

coordinator.poll()

获取最新元数据

3s发送一次心跳
pollHeartbeat()

判断是否需要加入消费者组()

拉取数据

pollForFetches()

发送请求并抓取数据

fetcher.sendFetches()

拦截器处理消息

interceptors.onConsume()

初始化抓取参数

FetchRequest.Builder

发送拉取数据请求

client.send()

监听服务器端返回的数据

future.addListener()