

尚硅谷大数据技术之 Kafka

(作者：尚硅谷大数据研发部)

第 1 章 Kafka 概述

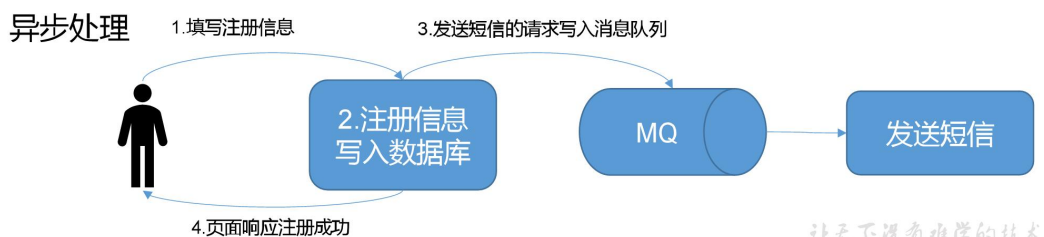
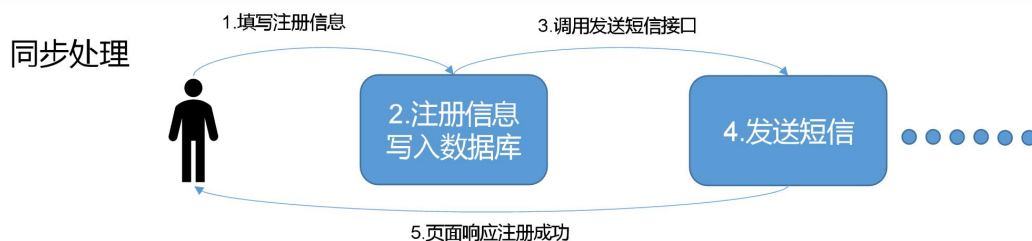
1.1 定义

Kafka 是一个分布式的基于发布/订阅模式的消息队列，主要应用于大数据实时处理领域。

1.2 消息队列（Message Queue）

1.2.1 传统消息队列的应用场景

MQ传统应用场景之异步处理

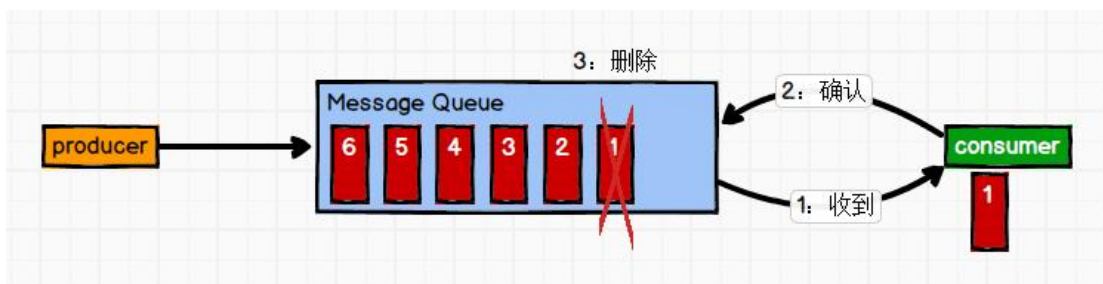


让天下没有难学的技术

1.2.2 消息队列的两种模式

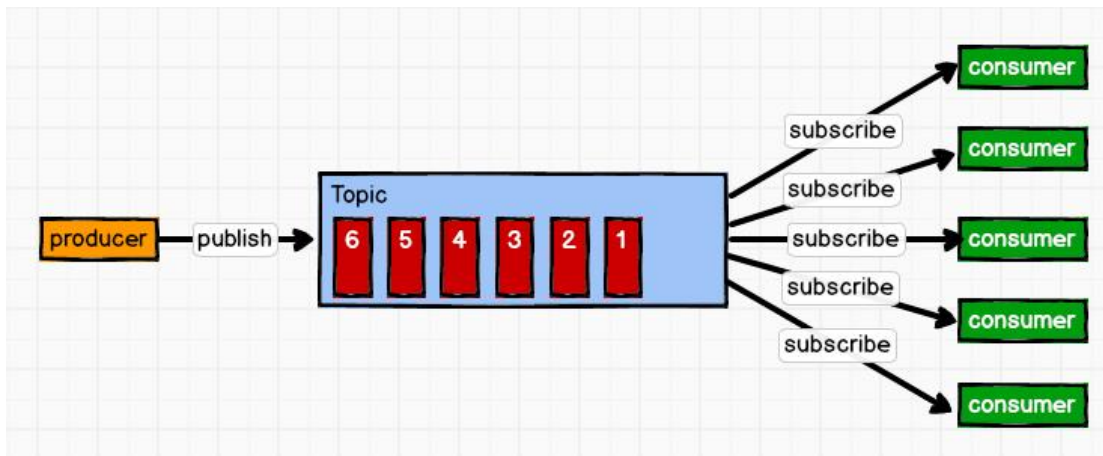
(1) 点对点模式（一对一，消费者主动拉取数据，消息收到后消息清除）

消息生产者生产消息发送到 Queue 中，然后消息消费者从 Queue 中取出并且消费消息。消息被消费以后，queue 中不再有存储，所以消息消费者不可能消费到已经被消费的消息。Queue 支持存在多个消费者，但是对一个消息而言，只会有一个消费者可以消费。



(2) 发布/订阅模式（一对多，消费者消费数据之后不会清除消息）

消息生产者（发布）将消息发布到 topic 中，同时有多个消息消费者（订阅）消费该消息。和点对点方式不同，发布到 topic 的消息会被所有订阅者消费。

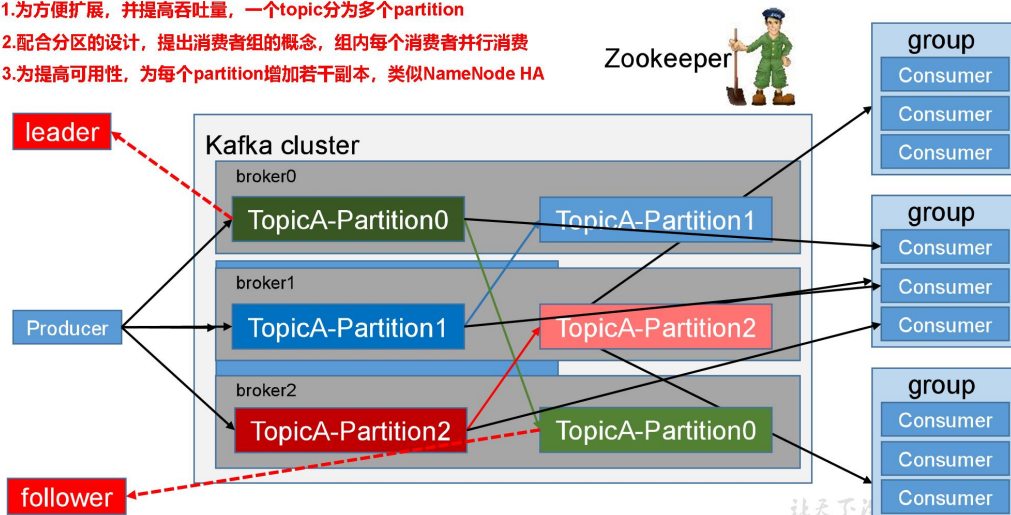


1.3 Kafka 基础架构



Kafka 基础架构

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



- 1) **Producer**：消息生产者，就是向 kafka broker 发消息的客户端；
- 2) **Consumer**：消息消费者，向 kafka broker 取消息的客户端；
- 3) **Consumer Group (CG)**：消费者组，由多个 consumer 组成。**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。**
- 4) **Broker**：一台 kafka 服务器就是一个 broker。一个集群由多个 broker 组成。一个 broker 可以容纳多个 topic。
- 5) **Topic**：可以理解为一个队列，**生产者和消费者面向的都是一个 topic；**
- 6) **Partition**：为了实现扩展性，一个非常大的 topic 可以分布到多个 broker（即服务器）上，

一个 **topic** 可以分为多个 **partition**，每个 **partition** 是一个有序的队列；

7) **Replica**: 副本，为保证集群中的某个节点发生故障时，该节点上的 **partition** 数据不丢失，且 **kafka** 仍然能够继续工作，**kafka** 提供了副本机制，一个 **topic** 的每个分区都有若干个副本，一个 **leader** 和若干个 **follower**。

8) **leader**: 每个分区多个副本的“主”，生产者发送数据的对象，以及消费者消费数据的对象都是 **leader**。

9) **follower**: 每个分区多个副本中的“从”，实时从 **leader** 中同步数据，保持和 **leader** 数据的同步。**leader** 发生故障时，某个 **follower** 会成为新的 **follower**。

第 2 章 Kafka 快速入门

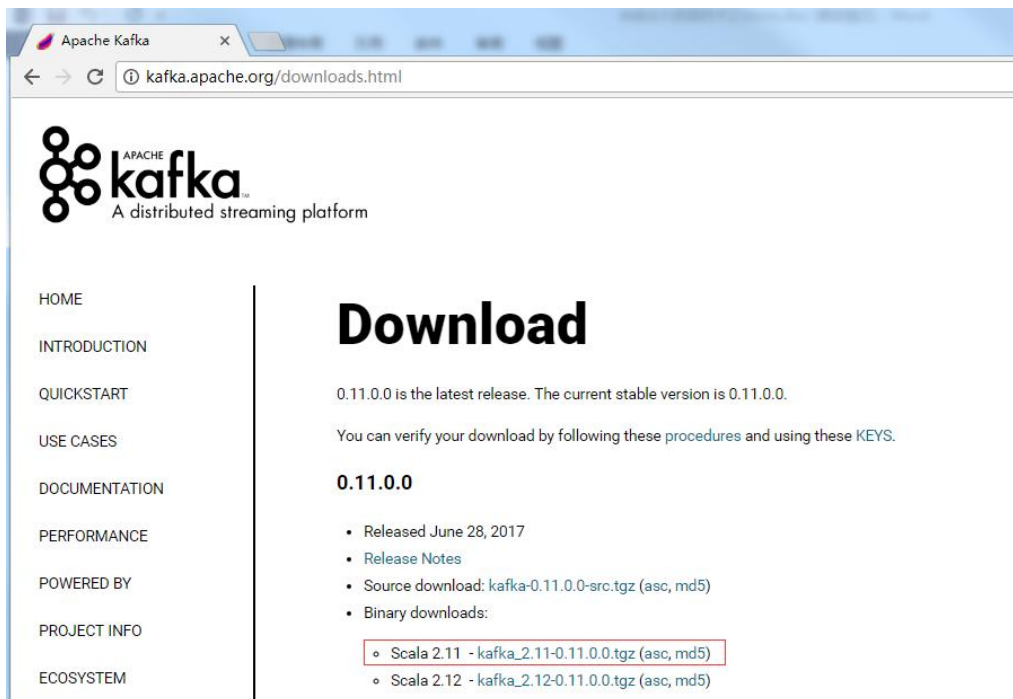
2.1 安装部署

2.1.1 集群规划

hadoop102	hadoop103	hadoop104
zk	zk	zk
kafka	kafka	kafka

2.1.2 jar 包下载

<http://kafka.apache.org/downloads.html>



2.1.3 集群部署

1) 解压安装包

```
[atguigu@hadoop102 software]$ tar -zxvf kafka_2.11-0.11.0.0.tgz -C /opt/module/
```

2) 修改解压后的文件名称

```
[atguigu@hadoop102 module]$ mv kafka_2.11-0.11.0.0/ kafka
```

3) 在/opt/module/kafka 目录下创建 logs 文件夹

```
[atguigu@hadoop102 kafka]$ mkdir logs
```

4) 修改配置文件

```
[atguigu@hadoop102 kafka]$ cd config/  
[atguigu@hadoop102 config]$ vi server.properties
```

输入以下内容:

```
#broker 的全局唯一编号, 不能重复  
broker.id=0  
#删除 topic 功能使能  
delete.topic.enable=true  
#处理网络请求的线程数量  
num.network.threads=3  
#用来处理磁盘 IO 的线程数量  
num.io.threads=8  
#发送套接字的缓冲区大小  
socket.send.buffer.bytes=102400  
#接收套接字的缓冲区大小  
socket.receive.buffer.bytes=102400  
#请求套接字的缓冲区大小  
socket.request.max.bytes=104857600  
#kafka 运行日志存放的路径  
log.dirs=/opt/module/kafka/logs  
#topic 在当前 broker 上的分区个数  
num.partitions=1  
#用来恢复和清理 data 下数据的线程数量  
num.recovery.threads.per.data.dir=1  
#segment 文件保留的最长时间, 超时将被删除  
log.retention.hours=168  
#配置连接 zookeeper 集群地址  
zookeeper.connect=hadoop102:2181,hadoop103:2181,hadoop104:2181
```

5) 配置环境变量

```
[atguigu@hadoop102 module]$ sudo vi /etc/profile  
  
#KAFKA_HOME  
export KAFKA_HOME=/opt/module/kafka  
export PATH=$PATH:$KAFKA_HOME/bin  
  
[atguigu@hadoop102 module]$ source /etc/profile
```

6) 分发安装包

```
[atguigu@hadoop102 module]$ xsync kafka/
```

注意: 分发之后记得配置其他机器的环境变量

7) 分别在 hadoop103 和 hadoop104 上修改配置文件/opt/module/kafka/config/server.properties 中的 **broker.id=1**、**broker.id=2**

注: broker.id 不得重复

8) 启动集群

依次在 hadoop102、hadoop103、hadoop104 节点上启动 kafka

```
[atguigu@hadoop102 kafka]$ bin/kafka-server-start.sh -daemon config/server.properties
[atguigu@hadoop103 kafka]$ bin/kafka-server-start.sh -daemon config/server.properties
[atguigu@hadoop104 kafka]$ bin/kafka-server-start.sh -daemon config/server.properties
```

9) 关闭集群

```
[atguigu@hadoop102 kafka]$ bin/kafka-server-stop.sh stop
[atguigu@hadoop103 kafka]$ bin/kafka-server-stop.sh stop
[atguigu@hadoop104 kafka]$ bin/kafka-server-stop.sh stop
```

10) kafka 群起脚本

```
for i in `cat /opt/module/hadoop-2.7.2/etc/hadoop/slaves`
do
echo "===== $i ====="
ssh $i 'source /etc/profile&&/opt/module/kafka_2.11-0.11.0.2/bin/kafka-server-start.sh -daemon /opt/module/kafka_2.11-0.11.0.2/config/server.properties'
echo $?
done
```

2.2 Kafka 命令行操作

1) 查看当前服务器中的所有 topic

```
[atguigu@hadoop102 kafka]$ bin/kafka-topics.sh --zookeeper hadoop102:2181 --list
```

2) 创建 topic

```
[atguigu@hadoop102 kafka]$ bin/kafka-topics.sh --zookeeper hadoop102:2181 \
--create --replication-factor 3 --partitions 1 --topic first
```

选项说明:

--topic 定义 topic 名

--replication-factor 定义副本数

--partitions 定义分区数

3) 删除 topic

```
[atguigu@hadoop102 kafka]$ bin/kafka-topics.sh --zookeeper hadoop102:2181 \
--delete --topic first
```

需要 server.properties 中设置 delete.topic.enable=true 否则只是标记删除。

4) 发送消息

```
[atguigu@hadoop102 kafka]$ bin/kafka-console-producer.sh \
--broker-list hadoop102:9092 --topic first
>hello world
>atguigu atguigu
```

5) 消费消息

```
[atguigu@hadoop103 kafka]$ bin/kafka-console-consumer.sh \
--bootstrap-server hadoop102:9092 --from-beginning --topic first

[atguigu@hadoop103 kafka]$ bin/kafka-console-consumer.sh \
--bootstrap-server hadoop102:9092 --from-beginning --topic first
```

--from-beginning: 会把主题中以往所有的数据都读取出来。

6) 查看某个 Topic 的详情

```
[atguigu@hadoop102 kafka]$ bin/kafka-topics.sh --zookeeper
hadoop102:2181 \
--describe --topic first
```

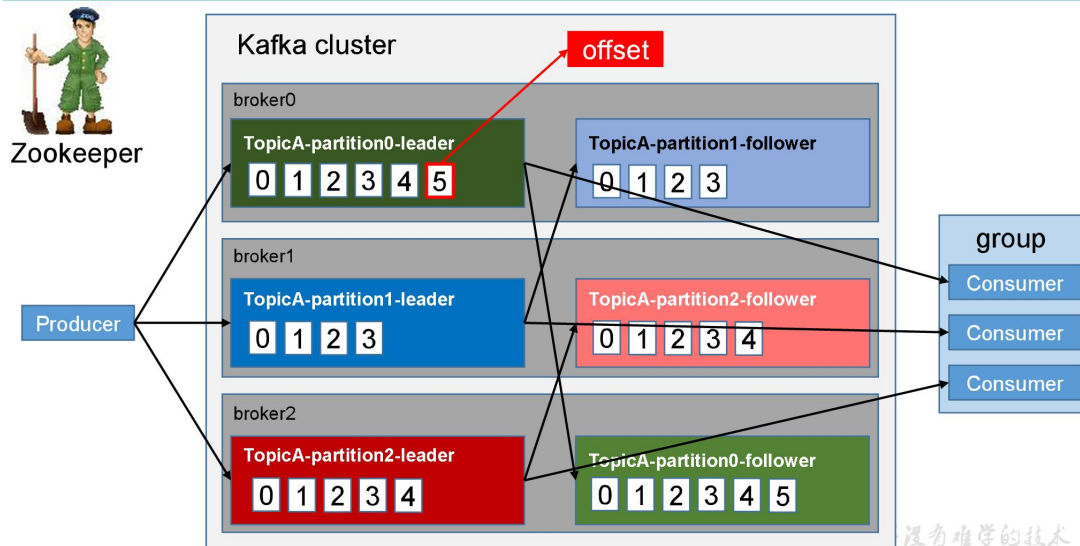
7) 修改分区数

```
[atguigu@hadoop102 kafka]$ bin/kafka-topics.sh --zookeeper
hadoop102:2181 --alter --topic first --partitions 6
```

第 3 章 Kafka 架构深入

3.1 Kafka 工作流程及文件存储机制

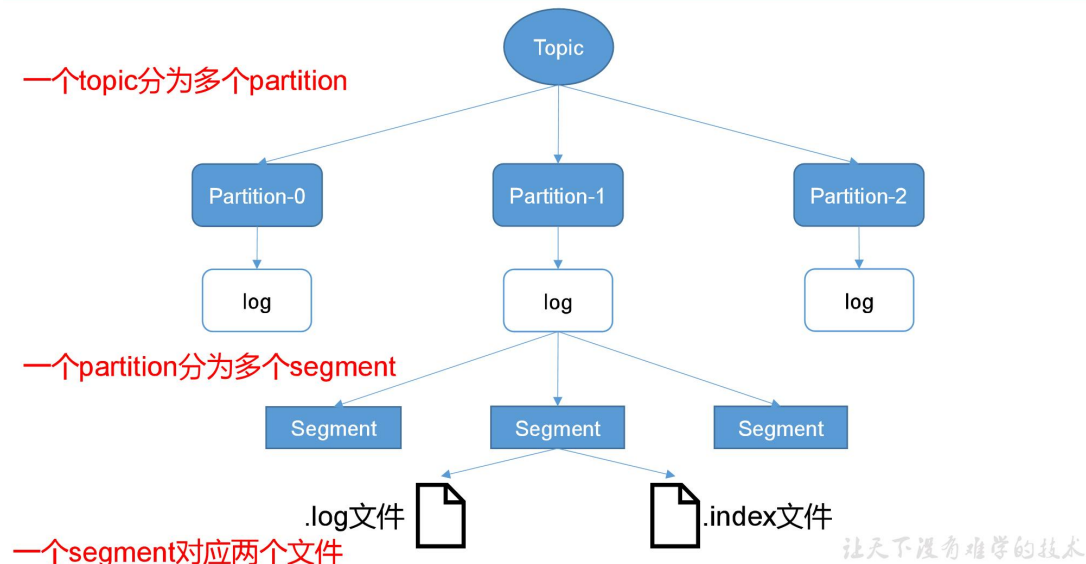
Kafka 工作流程



Kafka 中消息是以 **topic** 进行分类的，生产者生产消息，消费者消费消息，都是面向 topic 的。

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，且每条数据都有自己的 offset。消费者组中的每个消费者，都会实时记录自己消费到了哪个 offset，以便出错恢复时，从上次的位置继续消费。



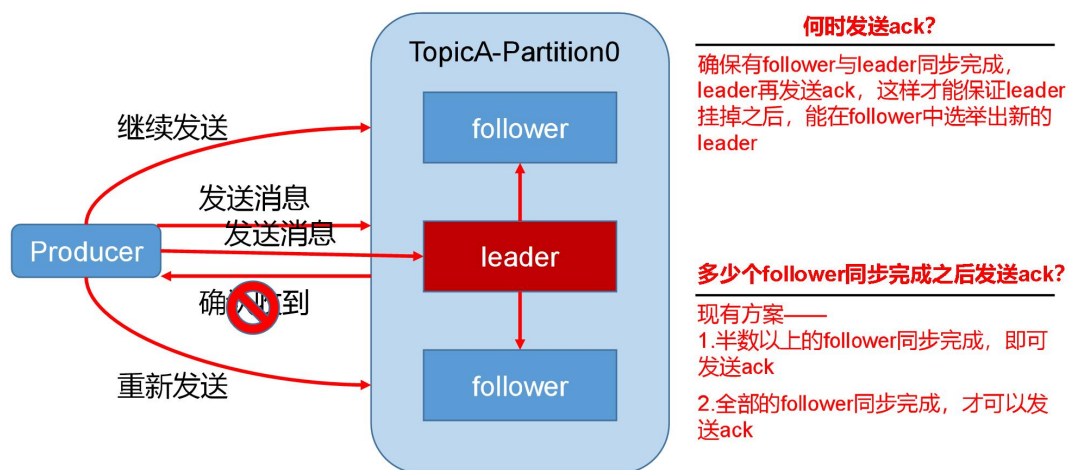
由于生产者生产的消息会不断追加到 log 文件末尾，为防止 log 文件过大导致数据定位效率低下，Kafka 采取了分片和索引机制，将每个 partition 分为多个 segment。每个 segment 对应两个文件——“.index”文件和“.log”文件。这些文件位于一个文件夹下，该文件夹的命名规则为：topic 名称+分区序号。例如，first 这个 topic 有三个分区，则其对应的文件夹为 first-0,first-1,first-2。

```
00000000000000000000000000000000.index
00000000000000000000000000000000.log
00000000000000000000000000000000.index
00000000000000000000000000000000.log
00000000000000000000000000000000.index
00000000000000000000000000000000.log
```

index 和 log 文件以当前 segment 的第一条消息的 offset 命名。下图为 index 文件和 log 文件的结构示意图。

3.2.3 数据可靠性保证

为保证 producer 发送的数据，能可靠的发送到指定的 topic，topic 的每个 partition 收到 producer 发送的数据后，都需要向 producer 发送 ack（acknowledgement 确认收到），如果 producer 收到 ack，就会进行下一轮的发送，否则重新发送数据。



让天下没有难学的技术

1) 副本数据同步策略

方案	优点	缺点
半数以上完成同步，就发送 ack	延迟低	选举新的 leader 时，容忍 n 台节点的故障，需要 $2n+1$ 个副本
全部完成同步，才发送 ack	选举新的 leader 时，容忍 n 台节点的故障，需要 $n+1$ 个副本	延迟高

Kafka 选择了第二种方案，原因如下：

- 同样为了容忍 n 台节点的故障，第一种方案需要 $2n+1$ 个副本，而第二种方案只需要 $n+1$ 个副本，而 Kafka 的每个分区都有大量的数据，第一种方案会造成大量数据的冗余。
- 虽然第二种方案的网络延迟会比较高，但网络延迟对 Kafka 的影响较小。

2) ISR

采用第二种方案之后，设想以下情景：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 leader 就要一直等下去，直到它完成同步，才能发送 ack。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 集合。当 ISR 中的 follower 完成数据的同步之后，leader 就会给 follower 发送 ack。如果 follower 长时间未向 leader 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。Leader 发生故障之后，就会从 ISR 中选举新的 leader。

3) ack 应答机制

对于某些不太重要的数据，对数据的可靠性要求不是很高，能够容忍数据的少量丢失，所以没必要等 ISR 中的 follower 全部接收成功。

所以 Kafka 为用户提供了三种可靠性级别，用户根据对可靠性和延迟的要求进行权衡，选择以下的配置。

acks 参数配置：

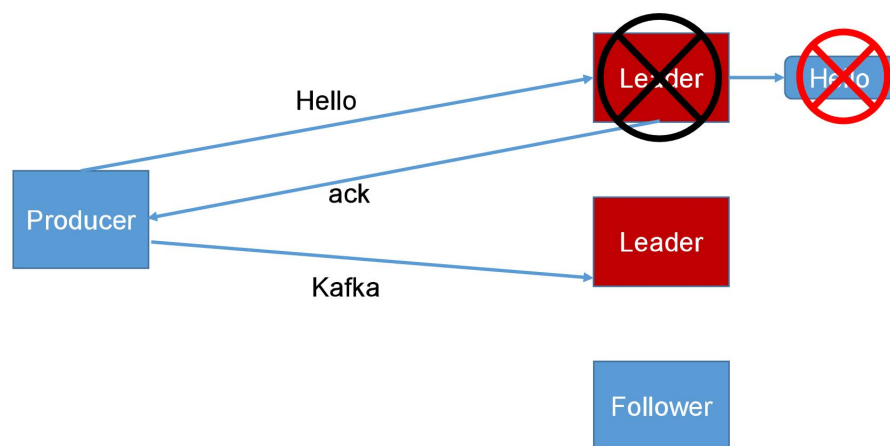
acks:

0: producer 不等待 broker 的 ack，这一操作提供了一个最低的延迟，broker 一接收到还没有写入磁盘就已经返回，当 broker 故障时有可能 **丢失数据**；

1: producer 等待 broker 的 ack，partition 的 leader 落盘成功后返回 ack，如果在 follower 同步成功之前 leader 故障，那么将会 **丢失数据**；

 **acks = 1 数据丢失案例**

 尚硅谷

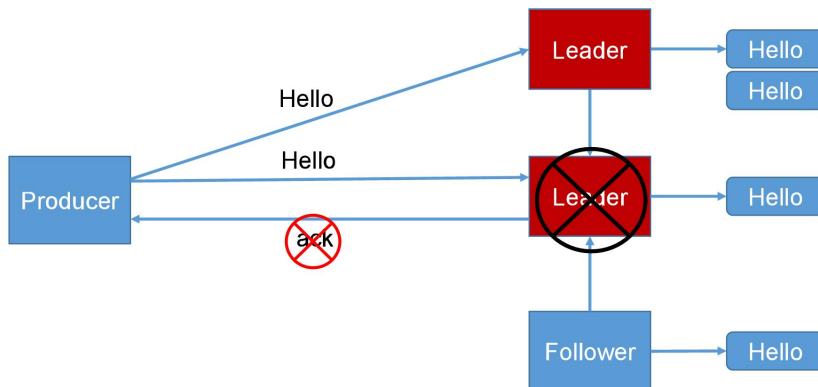


让天下没有难学的技术

-1 (all) : producer 等待 broker 的 ack，partition 的 leader 和 follower 全部落盘成功后才

返回 ack。但是如果在 follower 同步完成后, broker 发送 ack 之前, leader 发生故障, 那么会造成数据重复。

acks = -1 数据重复案例



让天下没有难学的技术

4) 故障处理细节

Log文件中的HW和LEO

LEO: 每个副本的最后一个offset

HW: 所有副本中最小的LEO

HW(High Watermark)

LEO(Log End Offset)

leader



LEO(HW)

follower



HW之前的数据才对Consumer可见

HW

LEO

follower



让天下没有难学的技术

(1) follower 故障

follower 发生故障后会被临时踢出 ISR, 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。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。

(2) leader 故障

leader 发生故障之后，会从 ISR 中选出一个新的 leader，之后，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。

3.2.4 Exactly Once 语义

对于某些比较重要的消息，我们需要保证 exactly once 语义，即**保证每条消息被发送且仅被发送一次**。

在 0.11 版本之后，Kafka 引入了幂等性机制（idempotent），配合 `acks = -1` 时的 at least once 语义，实现了 producer 到 broker 的 exactly once 语义。

idempotent + at least once = exactly once

使用时，只需将 `enable.idempotence` 属性设置为 `true`，kafka 自动将 `acks` 属性设为 -1。

3.3 Kafka 消费者

3.3.1 消费方式

consumer 采用 pull（拉）模式从 broker 中读取数据。

push（推）模式很难适应消费速率不同的消费者，因为消息发送速率是由 broker 决定的。

它的目标是尽可能以最快速度传递消息，但是这样很容易造成 consumer 来不及处理消息，典型的表现就是拒绝服务以及网络拥塞。而 pull 模式则可以根据 consumer 的消费能力以适当的速率消费消息。

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。针对这一点，Kafka 的消费者在消费数据时会传入一个时长参数 `timeout`，如果当前没有数据可供消费，consumer 会等待一段时间之后再返回，这段时长即为 `timeout`。

3.3.3 分区分配策略

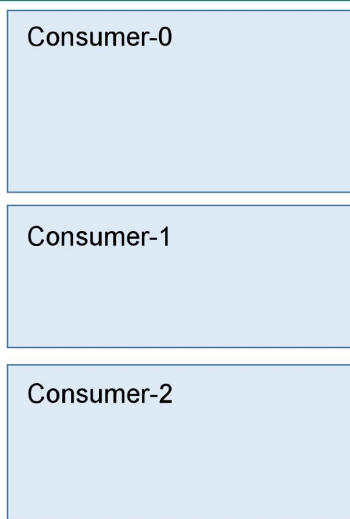
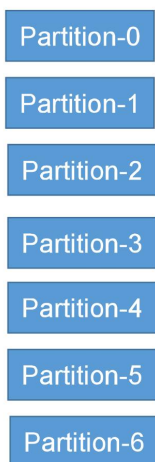
一个 consumer group 中有多个 consumer，一个 topic 有多个 partition，所以必然会涉及到 partition 的分配问题，即确定那个 partition 由哪个 consumer 来消费。

Kafka 有两种分配策略，一是 roundrobin，一是 range。

1) roundrobin



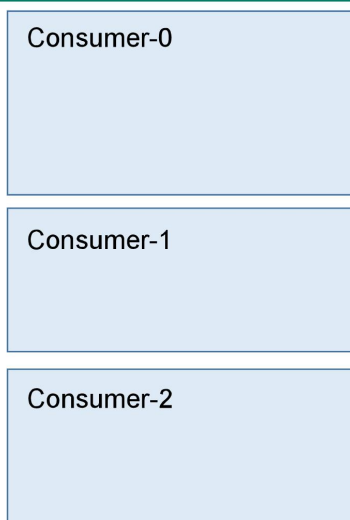
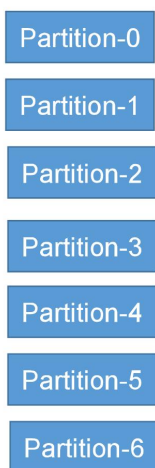
分区分配策略之roundrobin



2) range



分区分配策略之range



3.3.4 offset 的维护

由于 consumer 在消费过程中可能会出现断电宕机等故障，consumer 恢复后，需要从故障前的位置的继续消费，所以 consumer 需要实时记录自己消费到了哪个 offset，以便故障恢复后继续消费。

Kafka 0.9 版本之前，consumer 默认将 offset 保存在 Zookeeper 中，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`。

3.4 Kafka 高效读写数据

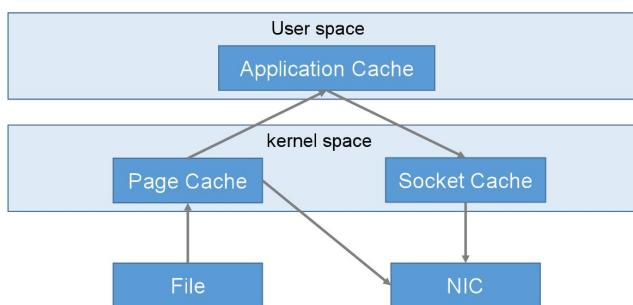
1) 顺序写磁盘

Kafka 的 producer 生产数据，要写入到 log 文件中，写的过程是一直追加到文件末端，为顺序写。官网有数据表明，同样的磁盘，顺序写能到 600M/s，而随机写只有 100k/s。这与磁盘的机械机构有关，顺序写之所以快，是因为其省去了大量磁头寻址的时间。

2) 零复制技术



零拷贝



让天下没有难学的技术

3.5 Zookeeper 在 Kafka 中的作用

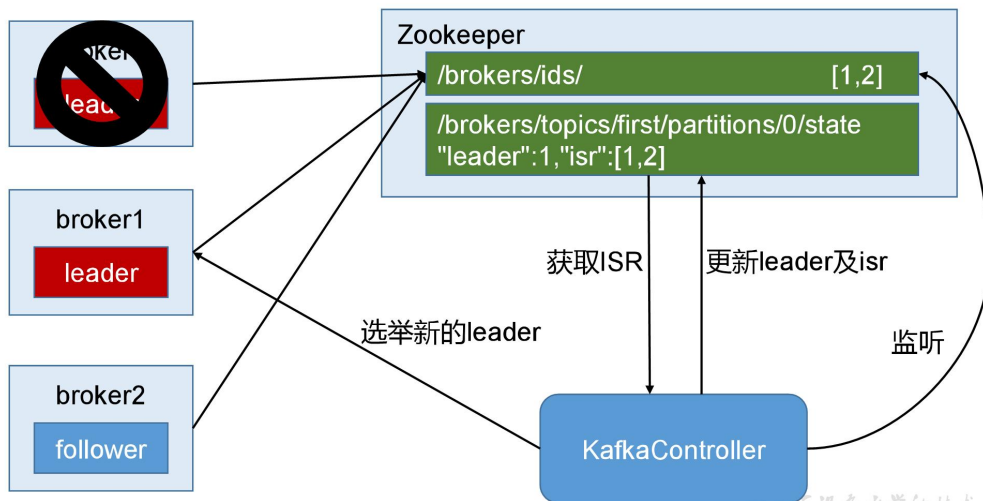
Kafka 集群中有一个 broker 会被选举为 Controller，负责管理集群 broker 的上下线，所有 topic 的分区副本分配和 leader 选举等工作。

Controller 的管理工作都是依赖于 Zookeeper 的。

以下为 partition 的 leader 选举过程：



Leader选举流程



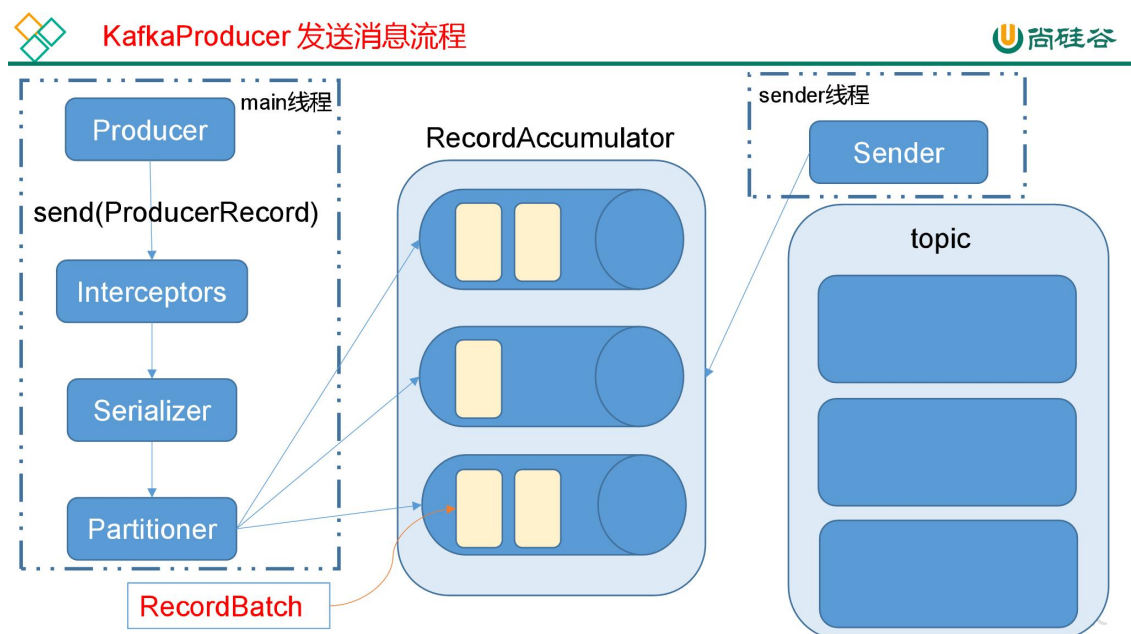
让天下没有难学的技术

第 4 章 Kafka API

4.1 Producer API

4.1.1 消息发送流程

Kafka 的 Producer 发送消息采用的是**异步发送**的方式。在消息发送的过程中，涉及到了**两个线程——main 线程和 Sender 线程**，以及一个**线程共享变量——RecordAccumulator**。main 线程将消息发送给 RecordAccumulator，Sender 线程不断从 RecordAccumulator 中拉取消息发送到 Kafka broker。



相关参数：

batch.size: 只有数据积累到 batch.size 之后，sender 才会发送数据。

linger.ms: 如果数据迟迟未达到 batch.size，sender 等待 linger.time 之后就会发送数据。

4.1.1 异步发送 API

1) 导入依赖

```
<dependency>
  <groupId>org.apache.kafka</groupId>
  <artifactId>kafka-clients</artifactId>
  <version>0.11.0.0</version>
</dependency>
```

2) 编写代码

需要用到的类：

KafkaProducer: 需要创建一个生产者对象，用来发送数据

ProducerConfig: 获取所需的一系列配置参数

ProducerRecord: 每条数据都要封装成一个 ProducerRecord 对象

1.不带回调函数的 API

```
package com.atguigu.kafka;

import org.apache.kafka.clients.producer.*;

import java.util.Properties;
import java.util.concurrent.ExecutionException;

public class CustomProducer {

    public static void main(String[] args) throws ExecutionException,
        InterruptedException {
        Properties props = new Properties();
        props.put("bootstrap.servers", "hadoop102:9092");//kafka 集群,
        broker-list
        props.put("acks", "all");
        props.put("retries", 1);//重试次数
        props.put("batch.size", 16384);//批次大小
        props.put("linger.ms", 1);//等待时间
        props.put("buffer.memory", 33554432);//RecordAccumulator 缓冲
        区大小
        props.put("key.serializer",
            "org.apache.kafka.common.serialization.StringSerializer");
        props.put("value.serializer",
            "org.apache.kafka.common.serialization.StringSerializer");

        Producer<String, String> producer = new
        KafkaProducer<>(props);
        for (int i = 0; i < 100; i++) {
            producer.send(new ProducerRecord<String, String>("first",
                Integer.toString(i), Integer.toString(i)));
        }
        producer.close();
    }
}
```

2.带回调函数的 API

回调函数会在 producer 收到 ack 时调用，为异步调用，该方法有两个参数，分别是 RecordMetadata 和 Exception，如果 Exception 为 null，说明消息发送成功，如果 Exception 不为 null，说明消息发送失败。

注意：消息发送失败会自动重试，不需要我们在回调函数中手动重试。

```
package com.atguigu.kafka;

import org.apache.kafka.clients.producer.*;

import java.util.Properties;
```

```
import java.util.concurrent.ExecutionException;

public class CustomProducer {

    public static void main(String[] args) throws ExecutionException,
    InterruptedException {
        Properties props = new Properties();
        props.put("bootstrap.servers", "hadoop102:9092");//kafka 集群,
broker-list
        props.put("acks", "all");
        props.put("retries", 1);//重试次数
        props.put("batch.size", 16384);//批次大小
        props.put("linger.ms", 1);//等待时间
        props.put("buffer.memory", 33554432);//RecordAccumulator 缓冲
区大小
        props.put("key.serializer",
"org.apache.kafka.common.serialization.StringSerializer");
        props.put("value.serializer",
"org.apache.kafka.common.serialization.StringSerializer");

        Producer<String, String> producer = new
KafkaProducer<>(props);
        for (int i = 0; i < 100; i++) {
            producer.send(new ProducerRecord<String, String>("first",
Integer.toString(i), Integer.toString(i)), new Callback() {

                //回调函数，该方法会在 Producer 收到 ack 时调用，为异步调用
                @Override
                public void onCompletion(RecordMetadata metadata,
Exception exception) {
                    if (exception == null) {
                        System.out.println("success->" +
metadata.offset());
                    } else {
                        exception.printStackTrace();
                    }
                }
            });
        }
        producer.close();
    }
}
```

4.1.2 同步发送 API

同步发送的意思就是，一条消息发送之后，会阻塞当前线程，直至返回 `ack`。

由于 `send` 方法返回的是一个 `Future` 对象，根据 `Future` 对象的特点，我们也可以实现同步发送的效果，只需在调用 `Future` 对象的 `get` 方法即可。

```
package com.atguigu.kafka;

import org.apache.kafka.clients.producer.KafkaProducer;
import org.apache.kafka.clients.producer.Producer;
import org.apache.kafka.clients.producer.ProducerRecord;

import java.util.Properties;
```

```
import java.util.concurrent.ExecutionException;

public class CustomProducer {

    public static void main(String[] args) throws ExecutionException,
    InterruptedException {
        Properties props = new Properties();
        props.put("bootstrap.servers", "hadoop102:9092");//kafka 集群,
broker-list
        props.put("acks", "all");
        props.put("retries", 1);//重试次数
        props.put("batch.size", 16384);//批次大小
        props.put("linger.ms", 1);//等待时间
        props.put("buffer.memory", 33554432);//RecordAccumulator 缓冲
区大小
        props.put("key.serializer",
"org.apache.kafka.common.serialization.StringSerializer");
        props.put("value.serializer",
"org.apache.kafka.common.serialization.StringSerializer");

        Producer<String, String> producer = new
KafkaProducer<>(props);
        for (int i = 0; i < 100; i++) {
            producer.send(new ProducerRecord<String, String>("first",
Integer.toString(i), Integer.toString(i))).get();
        }
        producer.close();
    }
}
```

4.2 Consumer API

Consumer 消费数据时的可靠性是很容易保证的，因为数据在 Kafka 中是持久化的，故不用担心数据丢失问题。

由于 consumer 在消费过程中可能会出现断电宕机等故障，consumer 恢复后，需要从故障前的位置的继续消费，所以 consumer 需要实时记录自己消费到了哪个 offset，以便故障恢复后继续消费。

所以 offset 的维护是 Consumer 消费数据是必须考虑的问题。

4.2.1 手动提交 offset

1) 导入依赖

```
<dependency>
    <groupId>org.apache.kafka</groupId>
    <artifactId>kafka-clients</artifactId>
    <version>0.11.0.0</version>
</dependency>
```

2) 编写代码

需要用到的类:

KafkaConsumer: 需要创建一个消费者对象, 用来消费数据

ConsumerConfig: 获取所需的一系列配置参数

ConsumerRecord: 每条数据都要封装成一个 ConsumerRecord 对象

```
package com.atguigu.kafka;

import org.apache.kafka.clients.consumer.ConsumerRecord;
import org.apache.kafka.clients.consumer.ConsumerRecords;
import org.apache.kafka.clients.consumer.KafkaConsumer;

import java.util.Arrays;
import java.util.Properties;

public class CustomConsumer {

    public static void main(String[] args) {
        Properties props = new Properties();
        props.put("bootstrap.servers", "hadoop102:9092");
        props.put("group.id", "test");//消费者组, 只要 group.id 相同, 就属于同一个消费者组
        props.put("enable.auto.commit", "false");//自动提交 offset

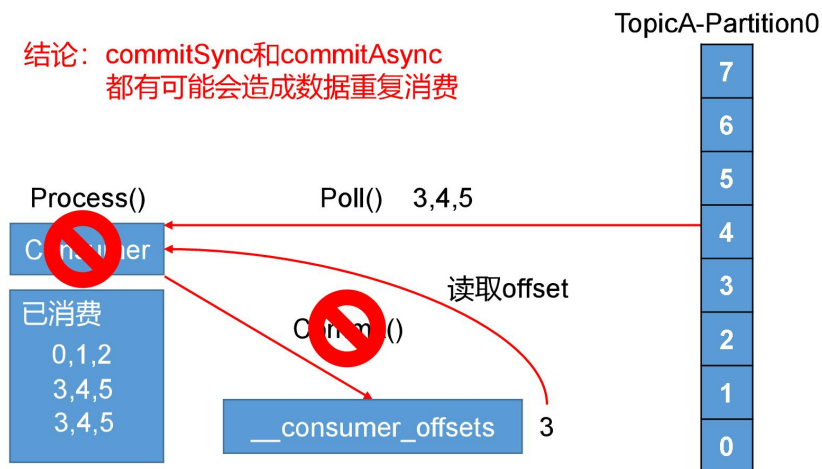
        props.put("key.deserializer",
"org.apache.kafka.common.serialization.StringDeserializer");
        props.put("value.deserializer",
"org.apache.kafka.common.serialization.StringDeserializer");
        KafkaConsumer<String, String> consumer = new
KafkaConsumer<>(props);
        consumer.subscribe(Arrays.asList("first"));
        while (true) {
            ConsumerRecords<String, String> records =
consumer.poll(100);
            for (ConsumerRecord<String, String> record : records) {
                System.out.printf("offset = %d, key = %s, value = %s\n",
record.offset(), record.key(), record.value());
            }
            consumer.commitSync();
        }
    }
}
```

3) 代码分析:

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 poll 的一批数据最高的偏移量提交**; 不同点是, `commitSync` 会失败重试, 一直到提交成功 (如果由于不可恢复原因导致, 也会提交失败); 而 `commitAsync` 则没有失败重试机制, 故有可能提交失败。

4) 数据重复消费问题

结论: `commitSync`和`commitAsync`都有可能会造成数据重复消费



让天下没有难学的技术

4.2.2 自动提交 offset

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

auto.commit.interval.ms: 自动提交 offset 的时间间隔

以下为自动提交 offset 的代码：

```
package com.atguigu.kafka;

import org.apache.kafka.clients.consumer.ConsumerRecord;
import org.apache.kafka.clients.consumer.ConsumerRecords;
import org.apache.kafka.clients.consumer.KafkaConsumer;

import java.util.Arrays;
import java.util.Properties;

public class CustomConsumer {

    public static void main(String[] args) {
        Properties props = new Properties();
        props.put("bootstrap.servers", "hadoop102:9092");
        props.put("group.id", "test");
        props.put("enable.auto.commit", "true");
        props.put("auto.commit.interval.ms", "1000");
        props.put("key.deserializer",
"org.apache.kafka.common.serialization.StringDeserializer");
        props.put("value.deserializer",
"org.apache.kafka.common.serialization.StringDeserializer");
        KafkaConsumer<String, String> consumer = new
KafkaConsumer<>(props);
        consumer.subscribe(Arrays.asList("first"));
        while (true) {
            ConsumerRecords<String, String> records =
```



```
consumer.poll(100);
    for (ConsumerRecord<String, String> record : records)
        System.out.printf("offset = %d, key = %s, value = %s%n",
record.offset(), record.key(), record.value());
    }
}
```

4.3 自定义 Interceptor

4.3.1 拦截器原理

Producer 拦截器(interceptor)是在 Kafka 0.10 版本被引入的, 主要用于实现 clients 端的定制化控制逻辑。

对于 producer 而言, interceptor 使得用户在消息发送前以及 producer 回调逻辑前有机会对消息做一些定制化需求, 比如修改消息等。同时, producer 允许用户指定多个 interceptor 按序作用于同一条消息从而形成一个拦截链(interceptor chain)。Interceptpor 的实现接口是 `org.apache.kafka.clients.producer.ProducerInterceptor`, 其定义的方法包括:

(1) `configure(configs)`

获取配置信息和初始化数据时调用。

(2) `onSend(ProducerRecord):`

该方法封装进 `KafkaProducer.send` 方法中, 即它运行在用户主线程中。Producer 确保在消息被序列化以及计算分区前调用该方法。用户可以在该方法中对消息做任何操作, 但最好保证不要修改消息所属的 topic 和分区, 否则会影响目标分区的计算。

(3) `onAcknowledgement(RecordMetadata, Exception):`

该方法会在消息从 `RecordAccumulator` 成功发送到 `Kafka Broker` 之后, 或者在发送过程中失败时调用。并且通常都是在 producer 回调逻辑触发之前。onAcknowledgement 运行在 producer 的 IO 线程中, 因此不要在该方法中放入很重的逻辑, 否则会拖慢 producer 的消息发送效率。

(4) `close:`

关闭 interceptor, 主要用于执行一些资源清理工作

如前所述, interceptor 可能被运行在多个线程中, 因此在具体实现时用户需要自行确保线程安全。另外倘若指定了多个 interceptor, 则 producer 将按照指定顺序调用它们, 并仅仅是捕获每个 interceptor 可能抛出的异常记录到错误日志中而非在向上传递。这在使用过程中要特别留意。

4.3.2 拦截器案例

1) 需求:

实现一个简单的双 interceptor 组成的拦截链。第一个 interceptor 会在消息发送前将时间戳信息加到消息 value 的最前部；第二个 interceptor 会在消息发送后更新成功发送消息数或失败发送消息数。



发送的数据	TimeInterceptor	CounterInterceptor	InterceptorProducer
	1) 实现ProducerInterceptor	1) 返回record	1) 构建拦截器链
	2) 获取record数据, 并在value前增加时间戳	2) 统计发送成功是失败次数 3) 关闭producer时, 打印统计次数 success:10 error:0	2) 发送数据
message0	1502102979120,message0	1502102979120,message0	
message1	1502102979242,message1	1502102979242,message1	
...	
message9	1502102979242,message9	1502102979242,message9	
message10	1502102979242,message10	1502102979242,message10	

让天下没有难学的技术

2) 案例实操

(1) 增加时间戳拦截器

```
package com.atguigu.kafka.interceptor;
import java.util.Map;
import org.apache.kafka.clients.producer.ProducerInterceptor;
import org.apache.kafka.clients.producer.ProducerRecord;
import org.apache.kafka.clients.producer.RecordMetadata;

public class TimeInterceptor implements
    ProducerInterceptor<String, String> {

    @Override
    public void configure(Map<String, ?> configs) {

    }

    @Override
    public ProducerRecord<String, String>
onSend(ProducerRecord<String, String> record) {
        // 创建一个新的 record, 把时间戳写入消息体的最前部
        return new ProducerRecord(record.topic(),
record.partition(), record.timestamp(), record.key(),
            System.currentTimeMillis() + "",
record.value().toString());
    }
}
```

```
@Override
    public void onAcknowledgement(RecordMetadata metadata,
Exception exception) {

    }

    @Override
    public void close() {

    }
}
```

(2) 统计发送消息成功和发送失败消息数，并在 **producer** 关闭时打印这两个计数器

```
package com.atguigu.kafka.interceptor;
import java.util.Map;
import org.apache.kafka.clients.producer.ProducerInterceptor;
import org.apache.kafka.clients.producer.ProducerRecord;
import org.apache.kafka.clients.producer.RecordMetadata;

public class CounterInterceptor implements
ProducerInterceptor<String, String>{
    private int errorCounter = 0;
    private int successCounter = 0;

    @Override
    public void configure(Map<String, ?> configs) {

    }

    @Override
    public ProducerRecord<String, String>
onSend(ProducerRecord<String, String> record) {
        return record;
    }

    @Override
    public void onAcknowledgement(RecordMetadata metadata,
Exception exception) {
        // 统计成功和失败的次数
        if (exception == null) {
            successCounter++;
        } else {
            errorCounter++;
        }
    }

    @Override
    public void close() {
        // 保存结果
        System.out.println("Successful sent: " + successCounter);
        System.out.println("Failed sent: " + errorCounter);
    }
}
```

(3) **producer** 主程序

```
package com.atguigu.kafka.interceptor;
import java.util.ArrayList;
import java.util.List;
```

```
import java.util.Properties;
import org.apache.kafka.clients.producer.KafkaProducer;
import org.apache.kafka.clients.producer.Producer;
import org.apache.kafka.clients.producer.ProducerConfig;
import org.apache.kafka.clients.producer.ProducerRecord;

public class InterceptorProducer {

    public static void main(String[] args) throws Exception {
        // 1 设置配置信息
        Properties props = new Properties();
        props.put("bootstrap.servers", "hadoop102:9092");
        props.put("acks", "all");
        props.put("retries", 0);
        props.put("batch.size", 16384);
        props.put("linger.ms", 1);
        props.put("buffer.memory", 33554432);
        props.put("key.serializer",
"org.apache.kafka.common.serialization.StringSerializer");
        props.put("value.serializer",
"org.apache.kafka.common.serialization.StringSerializer");

        // 2 构建拦截链
        List<String> interceptors = new ArrayList<>();

        interceptors.add("com.atguigu.kafka.interceptor.TimeInterce
ptor");
        interceptors.add("com.atguigu.kafka.interceptor.CounterInte
rceptor");
        props.put(ProducerConfig.INTERCEPTOR_CLASSES_CONFIG,
interceptors);

        String topic = "first";
        Producer<String, String> producer = new
KafkaProducer<>(props);

        // 3 发送消息
        for (int i = 0; i < 10; i++) {

            ProducerRecord<String, String> record = new
ProducerRecord<>(topic, "message" + i);
            producer.send(record);
        }

        // 4 一定要关闭 producer，这样才会调用 interceptor 的 close 方法
        producer.close();
    }
}
```

3) 测试

(1) 在 kafka 上启动消费者，然后运行客户端 java 程序。

```
[atguigu@hadoop102 kafka]$ bin/kafka-console-consumer.sh \
--bootstrap-server hadoop102:9092 --from-beginning --topic first

1501904047034,message0
1501904047225,message1
```

```
1501904047230,message2
1501904047234,message3
1501904047236,message4
1501904047240,message5
1501904047243,message6
1501904047246,message7
1501904047249,message8
1501904047252,message9
```

第 5 章 Flume 对接 Kafka

1) 配置 flume(flume-kafka.conf)

```
# define
a1.sources = r1
a1.sinks = k1
a1.channels = c1

# source
a1.sources.r1.type = exec
a1.sources.r1.command = tail -F -c +0 /opt/module/datas/flume.log
a1.sources.r1.shell = /bin/bash -c

# sink
a1.sinks.k1.type = org.apache.flume.sink.kafka.KafkaSink
a1.sinks.k1.kafka.bootstrap.servers =
hadoop102:9092,hadoop103:9092,hadoop104:9092
a1.sinks.k1.kafka.topic = first
a1.sinks.k1.kafka.flumeBatchSize = 20
a1.sinks.k1.kafka.producer.acks = 1
a1.sinks.k1.kafka.producer.linger.ms = 1

# channel
a1.channels.c1.type = memory
a1.channels.c1.capacity = 1000
a1.channels.c1.transactionCapacity = 100

# bind
a1.sources.r1.channels = c1
a1.sinks.k1.channel = c1
```

2) 启动 kafkaIDEA 消费者

3) 进入 flume 根目录下，启动 flume

```
$ bin/flume-ng agent -c conf/ -n a1 -f jobs/flume-kafka.conf
```

4) 向 /opt/module/datas/flume.log 里追加数据，查看 kafka 消费者消费情况

```
$ echo hello >> /opt/module/datas/flume.log
```

第 6 章 Kafka 监控

6.1 Kafka Monitor

1.上传 jar 包 KafkaOffsetMonitor-assembly-0.4.6.jar 到集群

2.在/opt/module/下创建 kafka-offset-console 文件夹

3.将上传的 jar 包放入刚创建的目录下

4.在/opt/module/kafka-offset-console 目录下创建启动脚本 start.sh，内容如下：

```
#!/bin/bash
java -cp KafkaOffsetMonitor-assembly-0.4.6-SNAPSHOT.jar \
com.quantifind.kafka.offsetapp.OffsetGetterWeb \
--offsetStorage kafka \
--kafkaBrokers hadoop102:9092,hadoop103:9092,hadoop104:9092 \
--kafkaSecurityProtocol PLAINTEXT \
--zk hadoop102:2181,hadoop103:2181,hadoop104:2181 \
--port 8086 \
--refresh 10.seconds \
--retain 2.days \
--dbName offsetapp_kafka &
```

5.在/opt/module/kafka-offset-console 目录下创建 mobile-logs 文件夹

```
mkdir /opt/module/kafka-offset-console/mobile-logs
```

6.启动 KafkaMonitor

```
./start.sh
```

7.登录页面 hadoop102:8086 端口查看详情

6.2 Kafka Manager

1.上传压缩包 kafka-manager-1.3.3.15.zip 到集群

2.解压到/opt/module

3.修改配置文件 conf/application.conf

```
kafka-manager.zkhosts="kafka-manager-zookeeper:2181"
```

修改为：

```
kafka-manager.zkhosts="hadoop102:2181,hadoop103:2181,hadoop104:2181"
```

4.启动 kafka-manager

```
bin/kafka-manager
```

5.登录 hadoop102:9000 页面查看详细信息

第 7 章 Kafka 面试题

7.1 面试问题

1.Kafka 中的 ISR、AR 又代表什么？

2.Kafka 中的 HW、LEO 等分别代表什么？

3.Kafka 中是怎么体现消息顺序性的？

4.Kafka 中的分区器、序列化器、拦截器是否了解？它们之间的处理顺序是什么？

5.Kafka 生产者客户端的整体结构是什么样子的？使用了几个线程来处理？分别是什么？

6.“消费组中的消费者个数如果超过 topic 的分区，那么就会有消费者消费不到数据”这句

话是否正确？

7.消费者提交消费位移时提交的是当前消费到的最新消息的 offset 还是 offset+1？

8.有哪些情形会造成重复消费？

9.那些情景会造成消息漏消费？

10.当你使用 kafka-topics.sh 创建（删除）了一个 topic 之后，Kafka 背后会执行什么逻辑？

1) 会在 zookeeper 中的 /brokers/topics 节点下创建一个新的 topic 节点，如：
/brokers/topics/first

2) 触发 Controller 的监听程序

3) kafka Controller 负责 topic 的创建工作，并更新 metadata cache

11.topic 的分区数可不可以增加？如果可以怎么增加？如果不可以，那又是为什么？

12.topic 的分区数可不可以减少？如果可以怎么减少？如果不可以，那又是为什么？

13.Kafka 有内部的 topic 吗？如果有是什么？有什么所用？

14.Kafka 分区分配的概念？

15.简述 Kafka 的日志目录结构？

16.如果我指定了一个 offset，Kafka Controller 怎么查找到对应的消息？

17.聊一聊 Kafka Controller 的作用？

18.Kafka 中有那些地方需要选举？这些地方的选举策略又有哪些？

19.失效副本是指什么？有那些应对措施？

20.Kafka 的那些设计让它有如此高的性能？

7.2 参考答案



Kafka相关面试题
及答案.docx