

kafka的基本概念



昂迪梵德

关注

4 人赞同了该文章

kafka的基本概念

kafka是什么？

kafka是一个多分区、多副本且基于zookeeper协调的分布式消息系统。也是一个分布式流式处理平台，它以高吞吐、可持久化、可水平扩展、支持流数据处理等多种特性而被广泛使用。

kafka扮演的三大角色

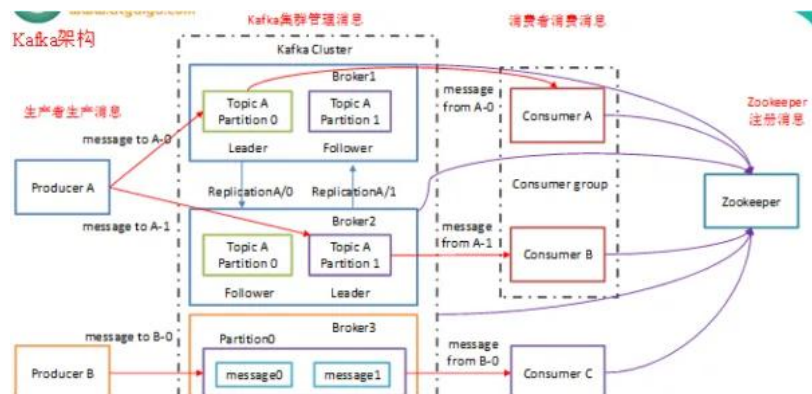
消息系统：kafka具备系统解耦、冗余存储、流量削峰、缓冲、异步通信、扩展性、可恢复性等功能。与此同时，Kafka 还提供了大多数消息系统难以实现的消息顺序性保障及回溯消费的功能。

存储系统：Kafka 把消息持久化到磁盘，相比于其他基于内存存储的系统而言，有效地降低了数据丢失的风险。也正是得益于Kafka 的消息持久化功能和多副本机制，我们可以把Kafka作为长期的数据存储系统来使用，只需要把对应的数据保留策略设置为“永久”或启用主题的日志压缩功能即可。

流式处理平台：Kafka 不仅为每个流行的流式处理框架提供了可靠的数据来源，还提供了一个完整的流式处理类库，比如窗口、连接、变换和聚合等各类操作。

kafka的基本架构

一个典型的 Kafka 体系架构包括若干 **Producer**、若干 **Broker**、若干 **Consumer**，以及一个 **ZooKeeper** 集群，如下图所示。其中ZooKeeper是Kafka用来负责集群元数据的管理、控制器的选举等操作的。Producer将消息发送到Broker，Broker负责将收到的消息存储到磁盘中，而Consumer负责从Broker订阅并消费消息。



(1) Producer：生产者，也就是发送消息的一方。生产者负责创建消息，然后将其投递到Kafka中。

(2) Consumer：消费者，也就是接收消息的一方。消费者连接到Kafka上并接收消息，进而进行相应的业务逻辑处理。

(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。

(4) Broker：服务代理节点。对于Kafka而言，Broker可以简单地看作一个独立的Kafka服务节点，容纳多个 topic。



(5) Topic: Kafka中的消息以主题为单位进行归类,生产者负责将消息发送到特定的主题(发送到Kafka集群中的每一条消息都要指定一个主题),而消费者负责订阅主题并进行消费

(6) Partition: 主题是一个逻辑上的概念,它还可以细分为多个分区,一个分区只属于单个主题,很多时候也会把分区称为主题分区(Topic-Partition)。同一主题下的不同分区包含的消息是不同的,分区在存储层面可以看作一个可追加的日志(Log)文件,消息在被追加到分区日志文件的时候都会分配一个特定的偏移量(offset)。offset是消息在分区中的唯一标识,Kafka通过它来保证消息在分区内的顺序性,不过offset并不跨越分区,也就是说,Kafka保证的是分区有序而不是主题有序。

Kafka中的分区可以分布在不同的服务器(broker)上,也就是说,一个主题可以横跨多个broker,以此来提供比单个broker更强大的性能。

每一条消息被发送到broker之前,会根据分区规则选择存储到哪个具体的分区。如果分区规则设定得合理,所有的消息都可以均匀地分配到不同的分区中。如果一个主题只对应一个文件,那么这个文件所在的机器I/O将会成为这个主题的性能瓶颈,而分区解决了这个问题。(7) Replica: Kafka为分区引入了多副本(Replica)机制,通过增加副本数量可以提升容灾能力。同一分区的不同副本中保存的是相同的消息(在同一时刻,副本之间并非完全一样),副本之间是“一主多从”的关系,其中leader副本负责处理读写请求,follower副本只负责与leader副本的消息同步。副本处于不同的broker中,当leader副本出现故障时,从follower副本中重新选举新的leader副本对外提供服务。Kafka通过多副本机制实现了故障的自动转移,当Kafka集群中某个broker失效时仍然能保证服务可用。

Kafka消费端也具备一定的容灾能力。Consumer使用拉(Pull)模式从服务端拉取消息,并且保存消费的具体位置,当消费者宕机后恢复上线时可以根据之前保存的消费位置重新拉取需要的消息进行消费,这样就不会造成消息丢失。

AR、ISR、OSR、HW、LEO

分区中的所有副本统称为**AR**(Assigned Replicas)。所有与leader副本保持一定程度同步的副本(包括leader副本在内)组成**ISR**(In-Sync Replicas),ISR集合是AR集合中的一个子集。消息会先发送到leader副本,然后follower副本才能从leader副本中拉取消息进行同步,同步期间内follower副本相对于leader副本而言会有一定程度的滞后。

与leader副本同步滞后过多的副本(不包括leader副本)组成**OSR**(Out-of-Sync Replicas),由此可见, $AR = ISR + OSR$ 。在正常情况下,所有的follower副本都应该与leader副本保持一定程度的同步,即 $AR = ISR$, OSR集合为空。

leader副本负责维护和跟踪ISR集合中所有follower副本的滞后状态,当follower副本落后太多或失效时,leader副本会把它从ISR集合中剔除。如果OSR集合中有follower副本“追上”了leader副本,那么leader副本会把它从OSR集合转移至ISR集合。默认情况下,当leader副本发生故障时,只有在ISR集合中的副本才有资格被选举为新的leader,而在OSR集合中的副本则没有任何机会(不过这个原则也可以通过修改相应的参数配置来改变)。

ISR与HW和LEO也有紧密的关系。**HW**是High Watermark的缩写,俗称高水位,它标识了一个特定的消息偏移量(offset),消费者只能拉取到这个offset之前的消息。

LEO是Log End Offset的缩写,它标识当前日志文件中下一条待写入消息的offset,分区ISR集合中的每个副本都会维护自身的LEO,而ISR集合中最小的LEO即为分区的HW,对消费者而言只能消费HW之前的消息。

kafka的复制机制

Kafka的复制机制既不是完全的同步复制,也不是单纯的异步复制。事实上,同步复制要求所有能工作的follower副本都复制完,这条消息才会被确认为已成功提交,这种复制方式极大地影响了性能。而在异步复制方式下,follower副本异步地从leader副本中复制数据,数据只要被leader副本写入就被认为已经成功提交。在这种情况下,如果follower副本都还没有复制完而落后于leader副本,突然leader副本宕机,则会造成数据丢失。Kafka使用的这种ISR的方式则有效地权衡了数据可靠性和性能之间的关系