

Kafka Replication Design

The purpose of adding replication in Kafka is for stronger durability and higher availability. We want to guarantee that any successfully published message will not be lost and can be consumed, even when there are server failures. Such failures can be caused by machine error, program error, or more commonly, software upgrades. We have the following high-level goals:

1. Configurable durability guarantees: For example, an application with critical data can choose stronger durability, with increased write latency, and another application generating a large volume of soft-state data can choose weaker durability but better write response time.
2. Automated replica management: We want to simplify the assignment of replicas to broker servers and be able to grow the cluster incrementally.

There are mainly two problems that we need to solve here: (1) How to assign replicas of a partition to broker servers evenly? (2) For a given partition, how to propagate every message to all replicas?

1. Replica placements

1.1. Initial placement

We first use an administrative api to create the initial set of brokers:

```
create cluster with brokers broker-0, broker-1, broker2
```

We then use another administrative api to create a new topic:

```
create topic topicX with 100 partitions
```

After that, the following information will be registered in zookeeper: (1) a list of brokers; (2) a list of topics and for each topic, a list of partitions.

For better load balancing, we want to over partition a topic. Typically, there will be many more partitions than servers. For each topic, we want to divide the partitions evenly among all the brokers. We sort the list of brokers and the list of partitions. If there are n brokers, we assign the i th partition to the $(i \bmod n)$ th broker. The first replica of this partition will reside on this assigned broker and is referred to as the preferred replica of this partition. We want to place the other replicas in such a way that if a broker is down, its load is spread evenly to all surviving brokers, instead of to a single one. In order to achieve that, suppose there are m partitions assigned to a broker i . The j th replica of partition k will be assigned to broker $(i + j + k) \bmod n$. The following figure illustrates the replica assignments for partitions p0 to p14 on brokers broker-0 to broker-4. In this example, if broker-0 goes down, partitions p0, p1, and p2 can be served from all remaining 4 brokers. We store the information about the replica assignment for each partition in Zookeeper.

broker-0	broker-1	broker-2	broker-3	broker-4	
p0	p3	p6	p9	p12	(1st replica)
p1	p4	p7	p10	p13	(1st replica)
p2	p5	p8	p11	p14	(1st replica)
p12	p0	p3	p6	p9	(2nd replica of row 1, right-shift 1)
p9	p12	p0	p3	p6	(3rd replica of row 1, right-shift 2)
p10	p13	p1	p4	p7	(2nd replica of row 2, right-shift 2)
p7	p10	p13	p1	p4	(3rd replica of row 2, right-shift 3)
p8	p11	p14	p2	p5	(2nd replica of row 3, right-shift 3)
p5	p8	p11	p14	p2	(3rd replica of row 3, right-shift 4)

1.2. Incrementally add brokers online

We'd like to be able to incrementally grow the set of brokers using an administrative command like the following.

```
alter cluster add brokers broker-3, broker-4
```

When a new broker is added, we will automatically move some partitions from existing brokers to the new one. Our goal is to minimize the amount of data movement while maintaining a balanced load on each broker. We use a standalone coordinator process to do the rebalance and the algorithm is given below.

1. For each new broker b {
2. suppose there are m partitions and the cluster is growing to n brokers
3. randomly select m/n partitions to move to b
4. For each partition p to be moved to b {
5. initialize the new replicas corresponding to the move
6. let the new replicas catch up from the current leader of the partition
7. }
8. wait until all new replicas have fully caught up
9. de-register existing replicas of p
10. register new replicas of p and trigger the leader election
11. physically delete the old replicas
12. }

2. Data replication

We'd like to allow a client to choose either asynchronous or synchronous replication. In the former case, a message to be published is acknowledged as soon as it reaches 1 replica. In the latter case, we will make our best effort to make sure that a message is only acknowledged after it reaches multiple replicas. When a client tries to publish a message to a partition of a topic, we

need to propagate the message to all replicas. We have to decide: (1) how to propagate a message; (2) how many replicas receive the message before we acknowledge to the client; (3) what happens when a replica goes down; (4) what happens when a failed replica comes back again. We introduce existing replication strategies in Section 2.1. We then describe our synchronous and asynchronous replication in Section 2.2 and Section 2.3, respectively.

2.1 Related works

There are two common strategies for keeping replicas in sync, primary-backup replication and quorum-based replication. In both cases, one replica is designated as the leader and the rest of the replicas are called followers. All write requests go through the leader and the leader propagates the writes to the follower.

In primary-backup replication, the leader waits until the write completes on every replica in the group before acknowledging the client. If one of the replicas is down, the leader drops it from the current group and continues to write to the remaining replicas. A failed replica is allowed to rejoin the group if it comes back and catches up with the leader. With f replicas, primary-backup replication can tolerate $f-1$ failures.

In the quorum-based approach, the leader waits until a write completes on a majority of the replicas. The size of the replica group doesn't change even when some replicas are down. If there are $2f+1$ replicas, quorum-based replication can tolerate f replica failures. If the leader fails, it needs at least $f+1$ replicas to elect a new leader.

There are tradeoffs between the 2 approaches:

1. The quorum-based approach has better write latency than the primary-backup one. A delay (e.g., long GC) in any replica increases the write latency in the latter, but not the former.
2. Given the same number of replicas, the primary-backup approach tolerates more concurrent failures.
3. A replication factor of 2 works well with the primary-backup approach. In quorum-based replication, both replicas have to be up for the system to be available.

We choose the primary-backup replication in Kafka since it tolerates more failures and works well with 2 replicas. A hiccup can happen when a replica is down or becomes slow. However, those are relatively rare events and the hiccup time can be reduced by tuning various timeout parameters.

2.2 Synchronous replication

Our synchronous replication follows the typical primary-backup approach. Each partition has n replicas and can tolerate $n-1$ replica failures. One of the replicas is elected as the leader and the rest of the replicas are followers. The leader maintains a set of in-sync replicas (ISR): the set of replicas that have fully caught up with the leader. For each partition, we store in Zookeeper the current leader and the current ISR.

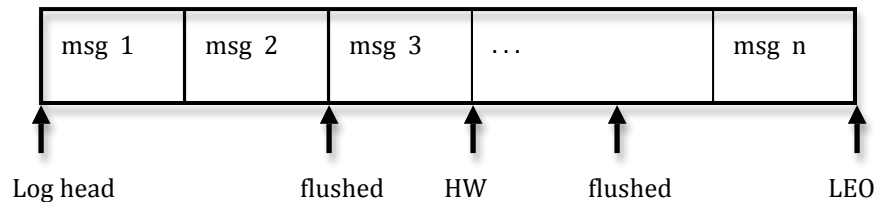


Figure 1

Each replica stores messages in a local log and maintains a few important offset positions in the log (depicted in Figure 1). The log end offset (LEO) represents the tail of the log. The high watermark (HW) is the offset of the last committed message. Each log is periodically synced to disks. Data before the flushed offset is guaranteed to be persisted on disks. As we will see, the flush offset can be before or after HW.

Writes:

To publish a message to a partition, the client first finds the leader of the partition from Zookeeper and sends the message to the leader. The leader writes the message to its local log. Each follower constantly pulls new messages from the leader using a single socket channel. That way, the follower receives all messages in the same order as written in the leader. The follower writes each received message to its own log and sends an acknowledgment back to the leader. Once the leader receives the acknowledgment from all replicas in ISR, the message is committed. The leader advances the HW and sends an acknowledgment to the client. For better performance, each follower sends an acknowledgment after the message is written to memory. So, for each committed message, we guarantee that the message is stored in multiple replicas in memory. However, there is no guarantee that any replica has persisted the commit message to disks though. Given that correlated failures are relatively rare, this approach gives us a good balance between response time and durability. In the future, we may consider adding options that provide even stronger guarantees. The leader also periodically broadcasts the HW to all followers. The broadcasting can be piggybacked on the return value of the fetch requests from the followers. From time to time, each replica checkpoints its HW to its disk.

Reads:

For simplicity, reads are always served from the leader. Only messages up to the HW are exposed to the reader.

Failure scenario:**a. Follower failure:**

After a configured timeout period, the leader will drop the failed follower from its ISR and writes will continue on the remaining replicas in ISR. If the failed follower comes back, it first truncates its log to the last checkpointed HW. It then starts to catch up all messages after its HW from the leader. When the follower fully catches up, the leader will add it back to the current ISR.

b. Leader failure:

When this happens, we need to perform the following steps to elect a new leader.

1. Each surviving replica in ISR registers itself in Zookeeper.
2. The replica that registers first becomes the new leader. The new leader chooses its LEO as the new HW.
3. Each replica registers a listener in Zookeeper so that it will be informed of any leader change. Everytime a replica is notified about a new leader:
 - a. If the replica is not the new leader (it must be a follower), it truncates its log to its HW and then starts to catch up from the new leader.
4. The leader waits until all surviving replicas in ISR have caught up or a configured time has passed. The leader writes the current ISR to Zookeeper and opens itself up for both reads and writes.

(Note, during the initial startup when ISR is empty, any replica can become the leader.)

Load balancing:

If a broker goes down and then comes back, no replica on this broker is a leader even after they have fully caught up. Since typically the leader does a bit of more work, the load on this broker could be less than others. We can periodically trigger a leader election and try to force the leader to the preferred replica. To do that, each leader does the following periodically.

```
if (leader is not on the preferred replica && the preferred replica is in ISR) {  
    de-register itself as the leader to trigger a leader election  
    delay its participation in the new leader election a bit  
    (this gives the preferred replica a better chance of becoming the new leader)  
}  
}
```

2.3 Asynchronous replication

To support asynchronous replication, the leader can acknowledge the client as soon it finishes writing the message to its local log. The only caveat is that during the catchup phase, the follower may have to truncate the data before its HW. Since the replication is asynchronous, there is no guarantee that a commit message can survive any broker failure.

3. Open issues

3.1. If the brokers are in multiple racks, how to guarantee that at least one replica goes to a different rack?

3.2. How to avoid LEO/LMT directly embedding a rich client?

A simple solution is to reject the produce request if the broker has no leader replicas and let the client retry.