

KEYSPACE

2025 / 北京

Valkey 多可用区部署实践

刘家文 | 2025/12/13

腾讯云研发

01

Valkey 原生的架构

Valkey 标准版 / Valkey 集群版

02

Valkey 单可用区部署

单分片集群如何保证高可用

03

Valkey 多可用区部署

数据节点和仲裁节点分离

04

多可用区的关键技术

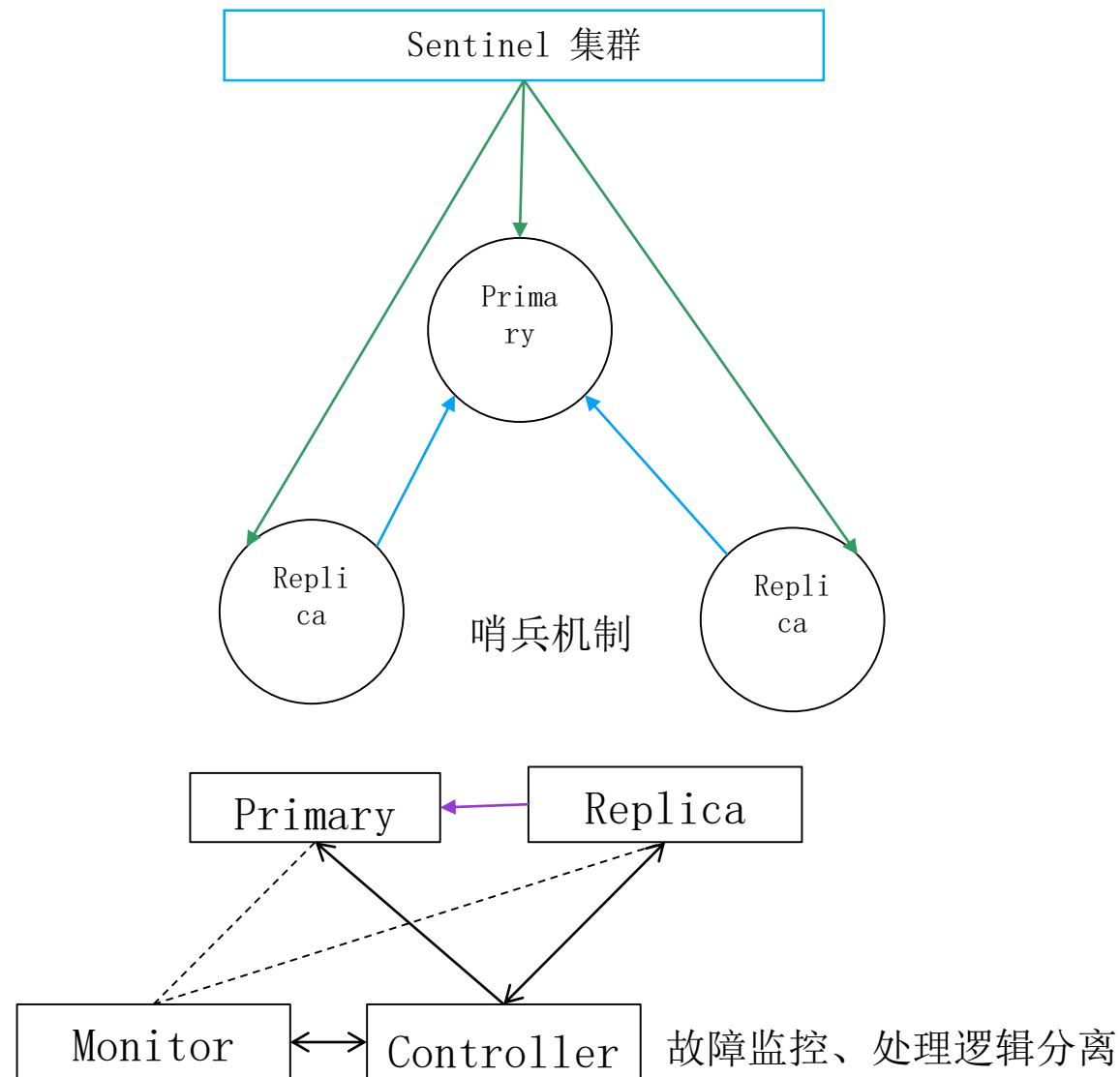
有序选举 / 本地优先

1. Valkey 原生的架构



Valkey 标准版

- 标准版：一个 Primary 和多个 Replica 组成
- Valkey 标准版高可用方案
 - 哨兵机制
 - 监控、切主、补从操作分离
- 标准版存在问题
 - 依赖外部仲裁机制
 - 高可用和切换速度受限于管理节点
 - 双写/丢数据的问题
 - Valkey 主节点下线有损
- Valkey 9.0 给哨兵引入了 Coordinated failover
 - Improve "SENTINEL FAILOVER" by using the "FAILOVER" command #1292



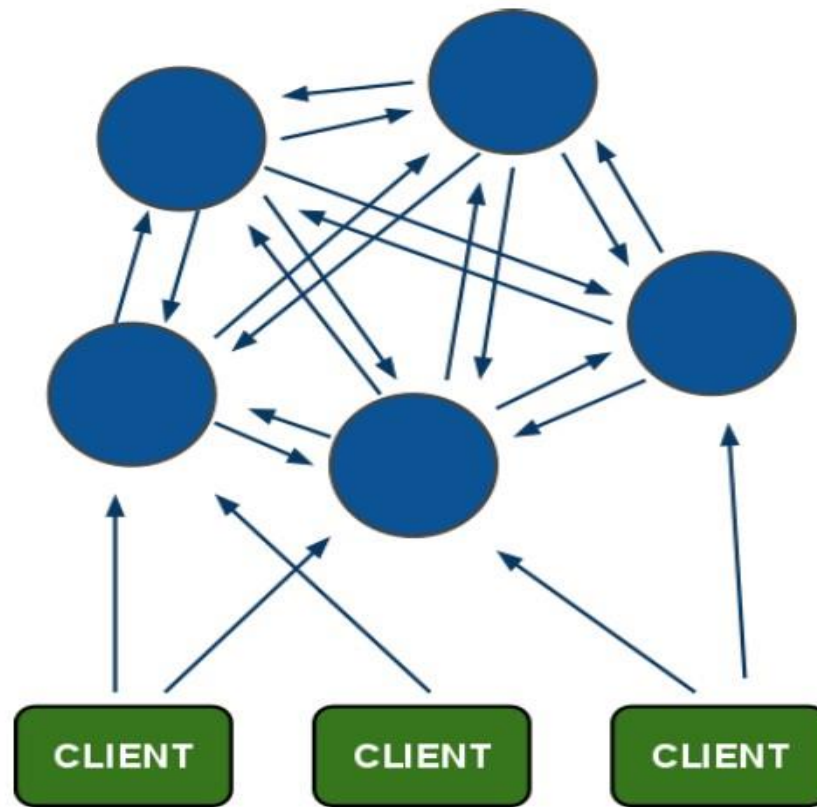
Valkey 集群版

■ 分片

- 节点独立
- 基于 Slot 的分片管理
- 支持平滑迁移
- 动态分片 (hot sharding)
- 数据均衡 (rebalance)

■ Gossip

- 自动发现 (auto detect)
- 自动容灾 (auto failover)



Valkey 集群版

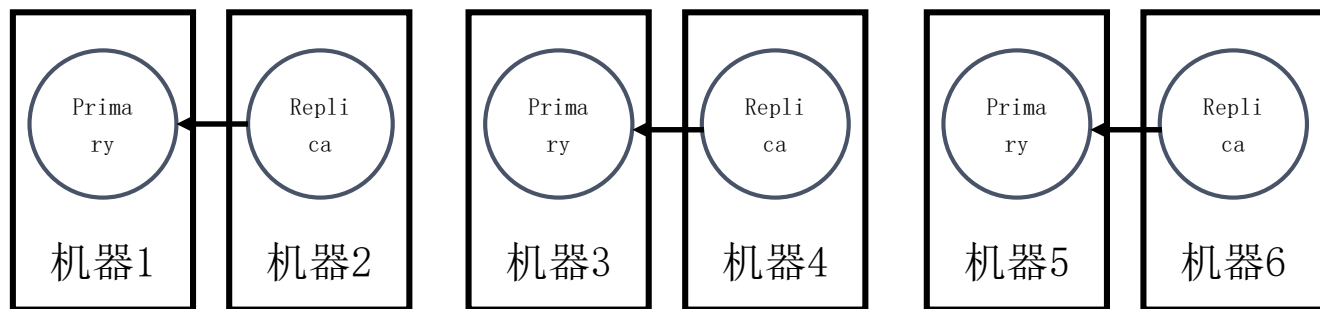
Valkey 集群版单可用区部署

■ Valkey 自动容灾:

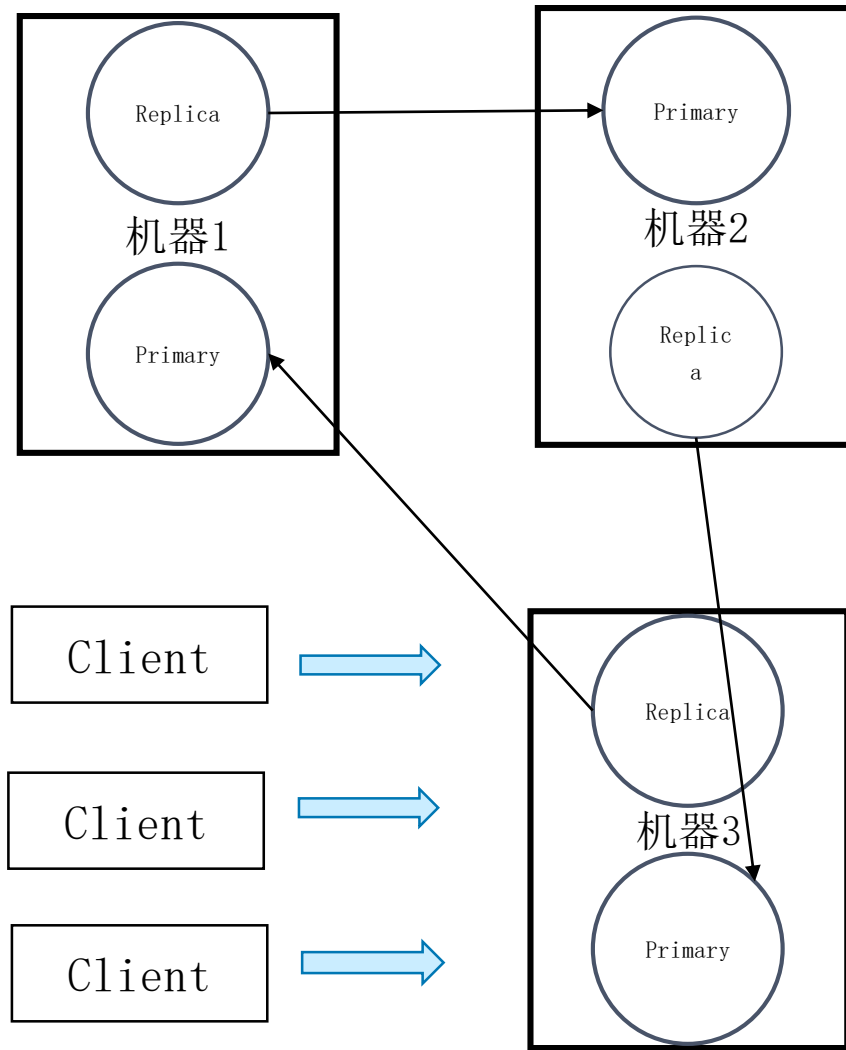
- 判死: 超过一半的主节点认定该节点故障
- 提主: 副本发起投票, 获得超过一半的主节点授权该节点。

■ 部署基本要求:

- 主从不能同机
- 单机**节点数**不能超过分片数的一半



Valkey 三分片集群正确部署



Valkey 三分片集群错误部署

2. Valkey 单可用区部署



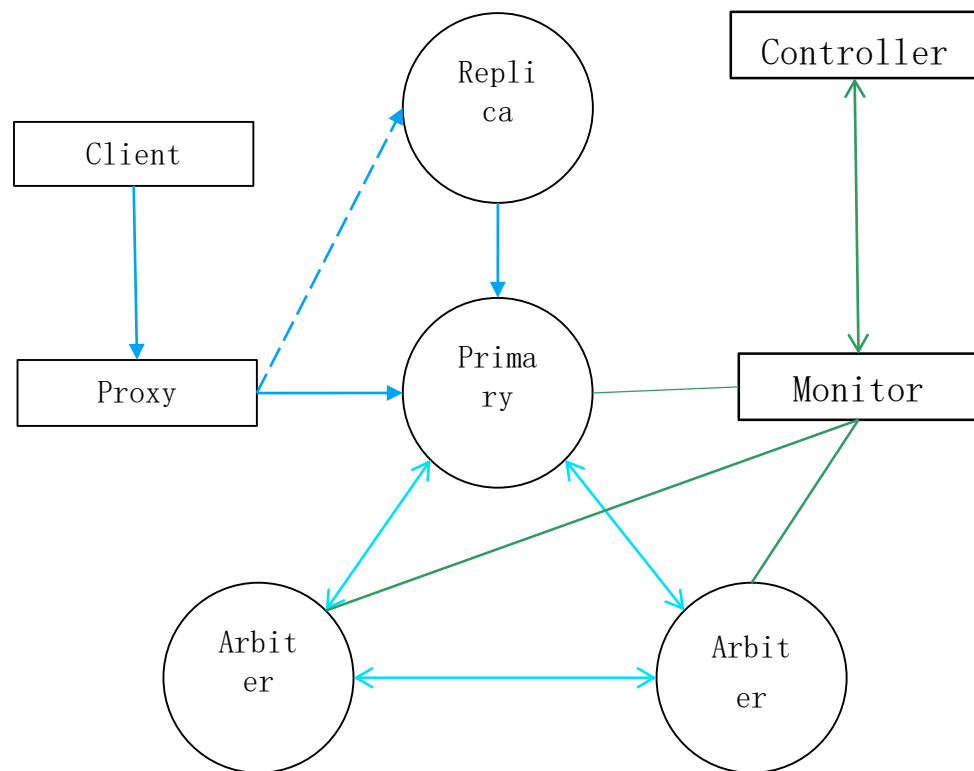
单分片集群如何保证高可用

■ 单分片集群（集群模式的标准版）

- 引入 voting empty primary (arbiter), 和主节点共同提供仲裁权
- 支持多 DB (Valkey 9.0 支持集群多 DB)
- 取消跨 slot 限制
- Smart proxy

■ 带来的优势

- 自动容灾
- 无损提主
- 兼容主从版
- 仲裁节点不受业务流量影响
- 无缝切换集群版，可水平扩缩容



单分片集群

3. Valkey 多可用区部署



原生 Valkey 多可用区部署缺陷

■ 部署基本要求：

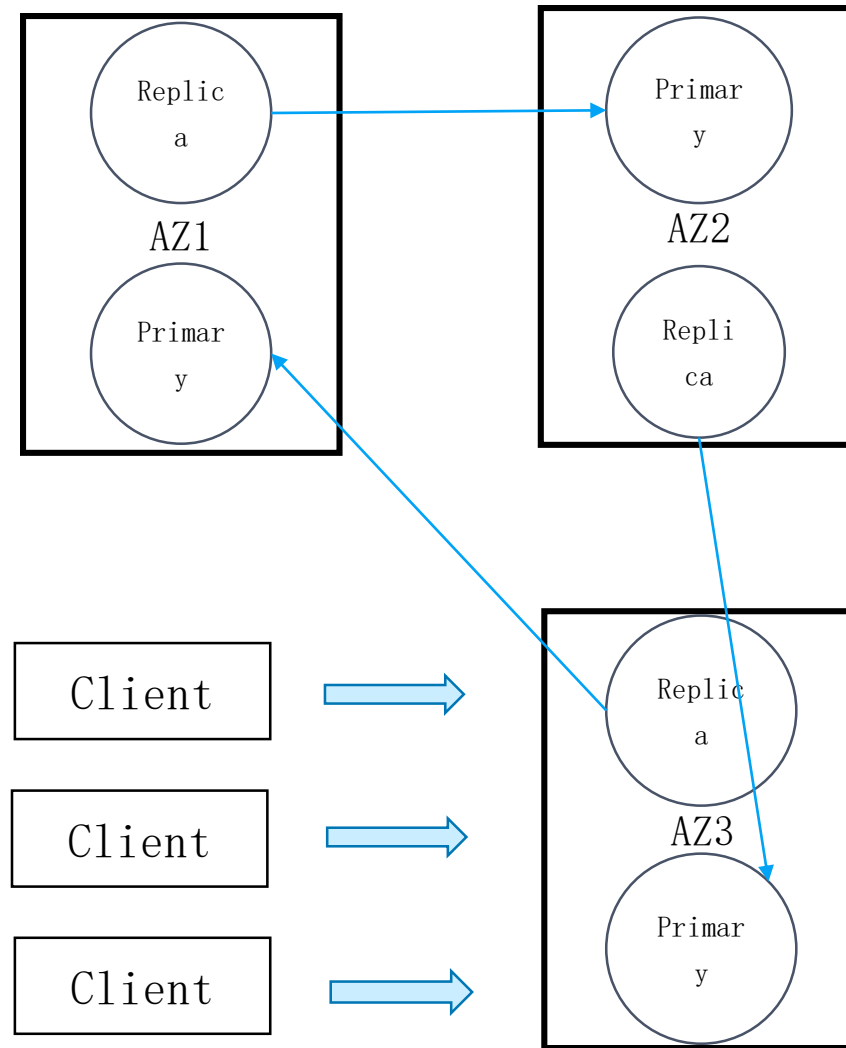
- 主节点和副本节点不能在同一可用区
- 一个可用区的**节点数**不能超过分片数的一半

■ 右边的部署问题：

- 无法保证高可用
- 跨可用区调用时延增加

■ 高可用问题

- 主节点不固定（故障转移随时可能发生）
- 仲裁权限与主节点绑定（只有主节点才有仲裁权）



Valkey 多可用区部署

社区数据节点和仲裁节点分离的讨论

[NEW] Non-voting primaries, voting empty primaries, voting replicas

Open



zuiderkwast opened on Jan 22

Member



The problem/use-case that the feature addresses

The cluster bus traffic overhead is huge for huge clusters because:

1. All nodes ping each other in a full mesh
2. Voting nodes = all primaries with slots

Description of the feature

Decouple right-to-vote from primary-serving-slots, so any node can be a voting or non-voting node.

This can allow non-voting primaries, voting replicas and voting empty nodes. Modifying these details may be easier than to rewrite the cluster protocol to Raft ([#384](#)).

By having only a few voting nodes (such as five) in a large cluster with hundreds of nodes, voting becomes faster.

Then we can also limit the ping-pong traffic between non-voting nodes. For example, non-voting nodes don't need to ping other nodes. They can just reply to pings from other nodes. This makes the voting nodes a team of leaders and the non-voters a team of followers.

<https://github.com/valkey-io/valkey/issues/1600>

数据节点和仲裁节点分离

- 只有主节点有仲裁权：

- 仲裁权会随着节点故障转移而漂移
- 大规模集群中投票效率低下

- 数据节点和仲裁节点分离：

- 仲裁节点不会受业务流量影响
- 仲裁权不会随着节点故障转移而漂移
- 大规模集群中只需要少量仲裁节点即可维持仲裁权
- 减少集群 Ping-Pong 通信开销，数据节点之间无需相互通信，只需要回复仲裁节点的 Ping

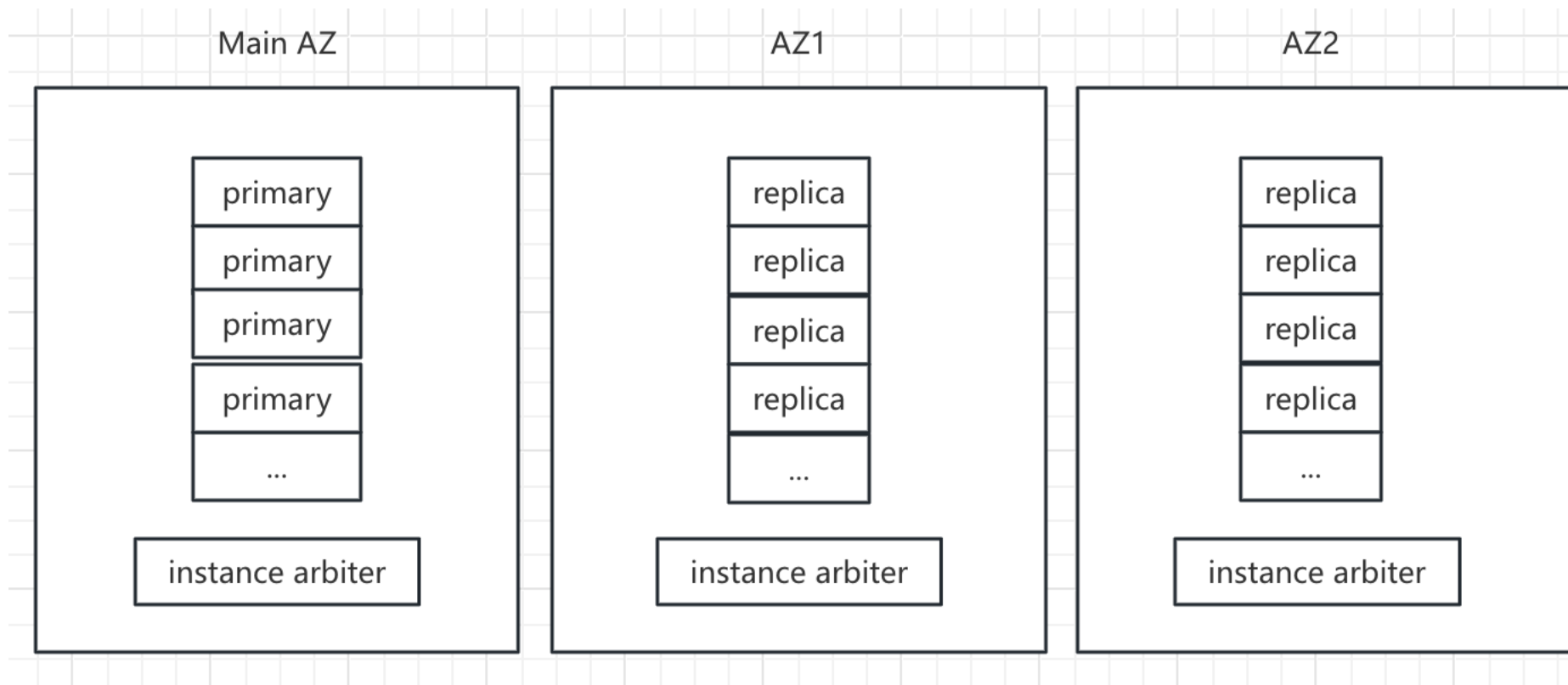
如何改造 Valkey 以实现多可用区部署

■ 如何解决高可用：

- 固定节点投票权
- 引入 instance arbiter
- 数据节点和仲裁节点分离

■ 多可用区的关键问题：

- 节点部署问题
- 跨可用区性能损耗
- 可用区故障提主



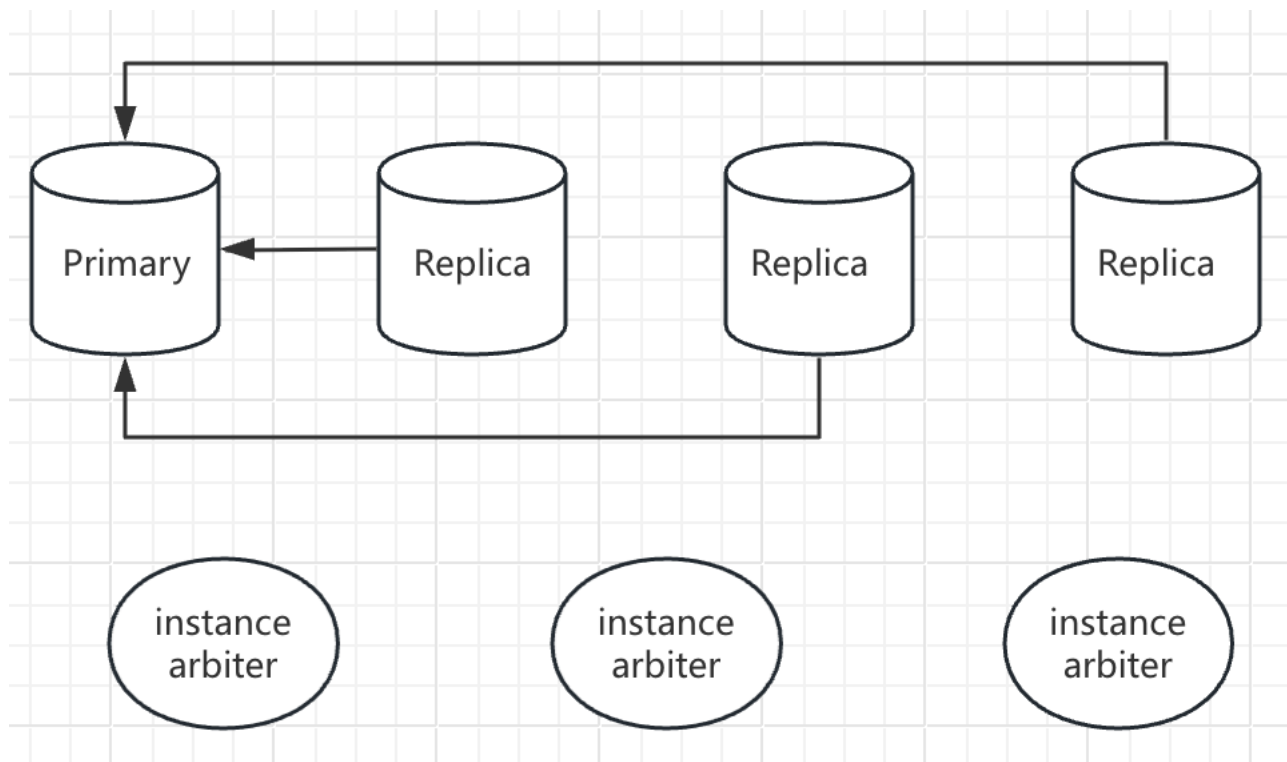
4. 多可用区的关键技术



多可用区节点部署

■ 部署准则：

- 主节点和副本不同区
- instance arbiter 节点位于不同可用区
- instance arbiter 至少需要部署 3 个节点



单分片集群多可用区部署

客户端就近接入

■ 数据访问：

- 当集群在跨多个可用区环境中运行时，出于时延和成本方面的考虑，最好将流量限制在客户端所在的可用区内

■ 关键点：

- 感知区域信息
- 新配置项：availability-zone
- 新 info field: availability_zone
- 新 hello response field:
availability_zone

<https://github.com/valkey-io/valkey/pull/700>

Upstream the availability zone info string from KeyDB #700

Merged

PingXie merged 1 commit into [valkey-io:unstable](#) from [JohnSully:availability-zone](#) on Jun 28, 2024

Conversation 39

Commits 1

Checks 0

Files changed 4



JohnSully commented on Jun 27, 2024 • edited by enjoy-binbin

Contributor

When Redis/Valkey/KeyDB is run in a cloud environment across multiple AZ's it is preferable to keep traffic local to an AZ both for cost reasons and for latency. This is typically done when you are enabling reads on replicas with the READONLY command.

For this change we are creating a setting that is echo'd back in the info command. We do not want to add the cloud SDKs as dependencies and this is the easiest way around that. It is fairly trivial to grab the AZ from the cloud and push that into your setting file.

Currently at Snapchat we have a custom client that after connecting reads this from the server and will preferentially use that server if the AZ string matches its internally configured AZ.

In the future it would be ideal if we used this information when performing failover or even exposed it in cluster nodes.

New configuration item: availability-zone

New info field: availability_zone

availability_zone was also added in HELLO response in 8.1, see [#1487](#) for more details.

```
> ./src/valkey-cli config set availability-zone beijing
OK
> ./src/valkey-cli info | grep availability_zone
availability_zone:beijing
> ./src/valkey-cli hello | grep availability_zone
availability_zone
```


多可用区故障

■ 多可用区和单可用区故障差异：

- 节点故障数量
- 部分节点故障切主导致的性能波动

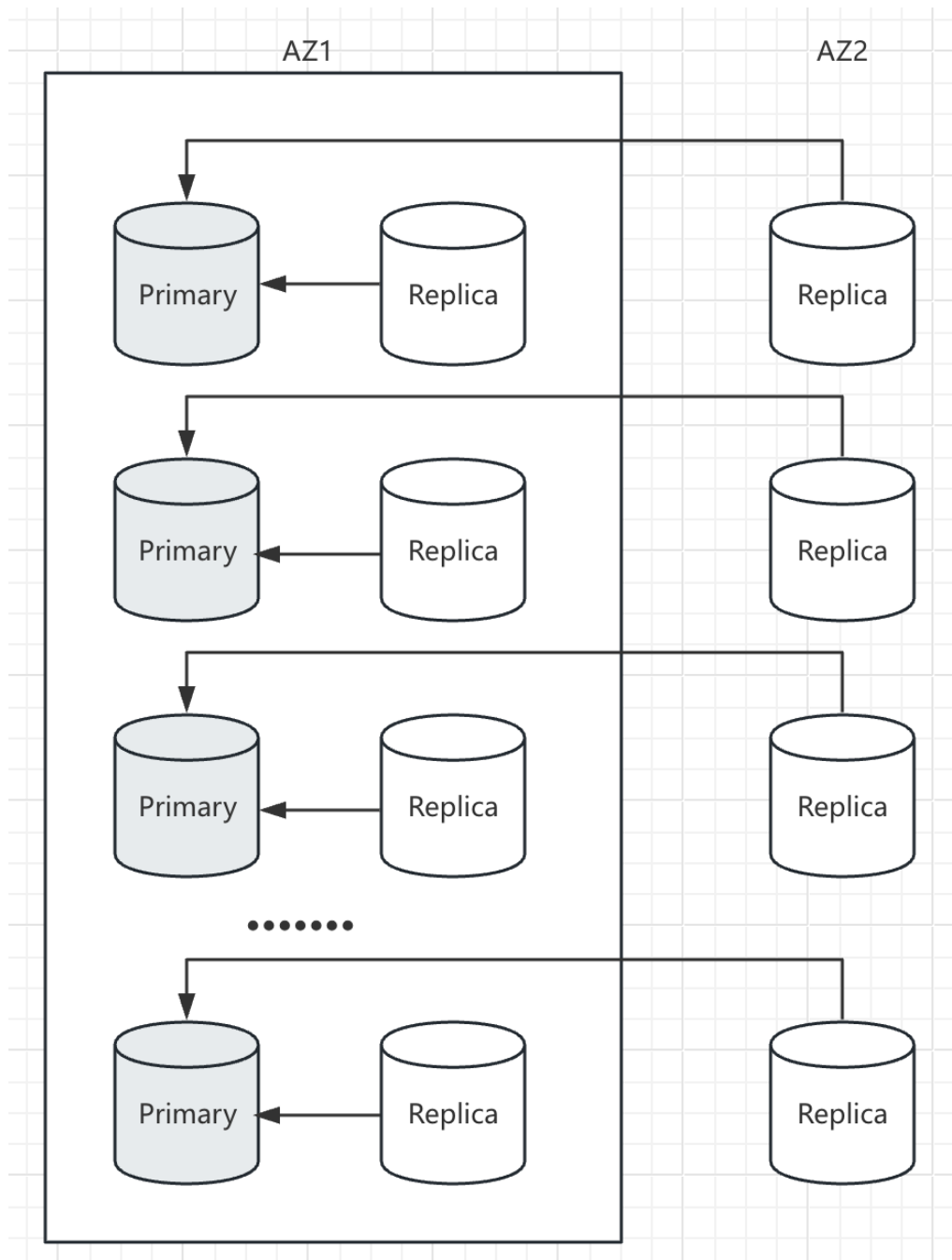
■ 多节点故障投票测试：

- Valkey 8.1 之前：128 分片，63 个主节点故障，在默认配置项下，99% 的情况集群无法自动恢复

- Do election in order based on failed primary rank to avoid voting conflicts #1018

■ 投票机制的关键点：

- Gossip 投票机制



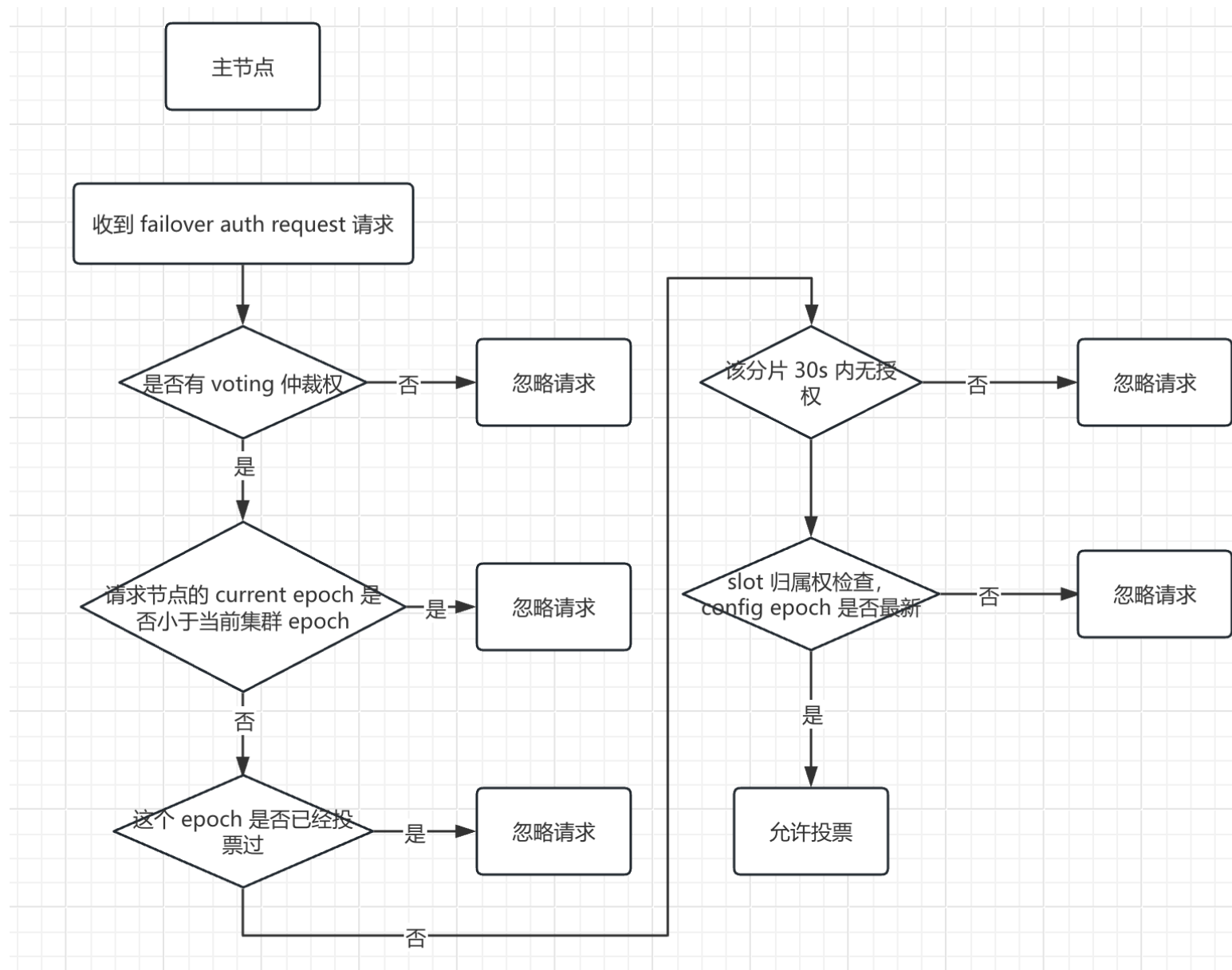
选主机制-主节点

■ 备注:

- Current epoch: 集群的 epoch
- Config epoch: 每个节点的 epoch

■ 允许投票的要求:

- 节点的集群 epoch 是最新的
- 这个 epoch 没有授权过。
- 同一分片在 $2 * \text{cluster_node_timeout}$ 内（默认 30s）内没有投票过。该限制在 8.1 里被移除。
- 这个节点 slot 归属权是该节点的。



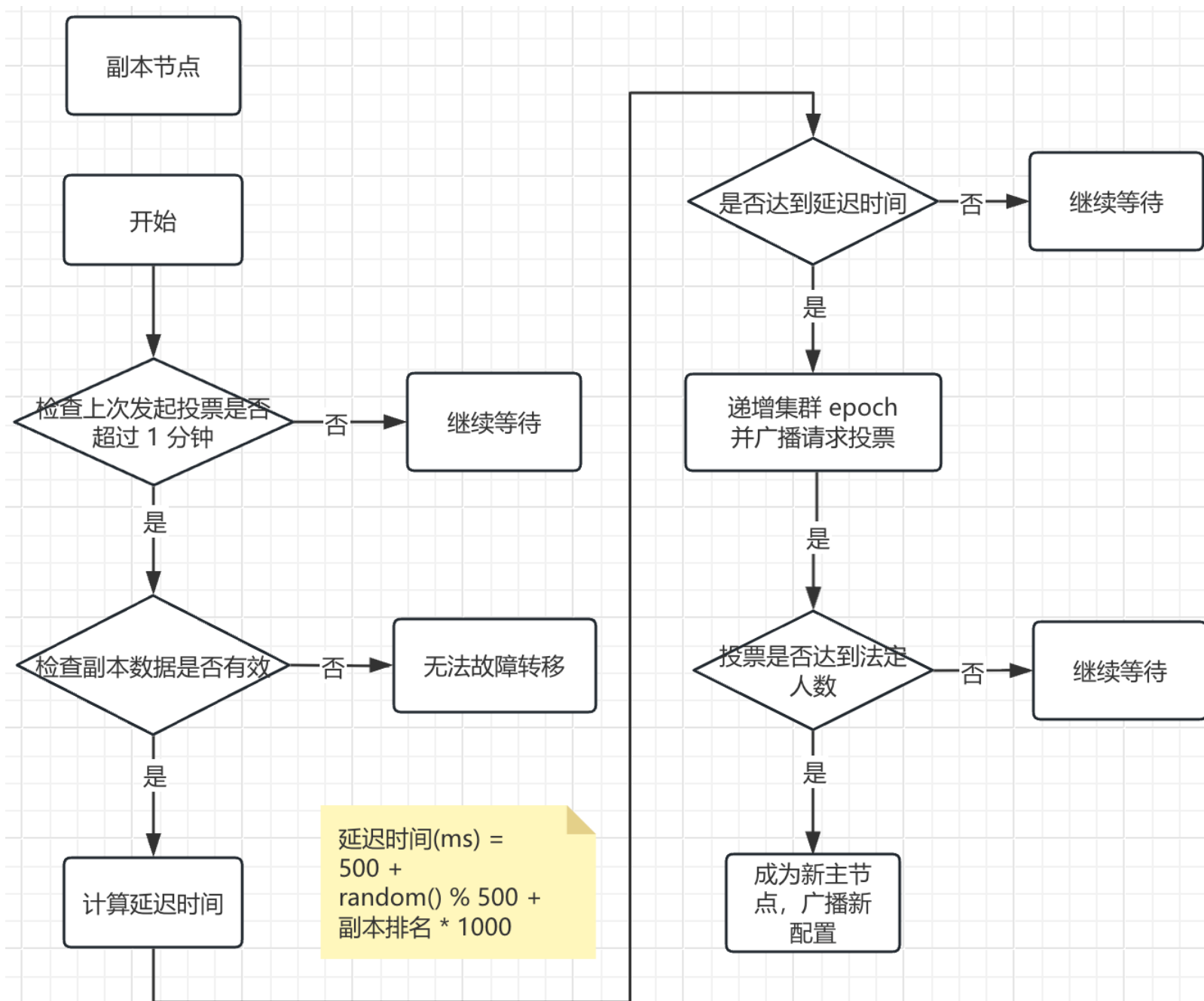
Automatic failover vote is not limited by two times the node timeout #1356

Manual failover vote is not limited by two times the node timeout #1305

选主机制-副本节点

■ 发起选举投票的要求:

- a) 一分钟只能重试一次投票。最大重试160s，从节点最多可以投票3次
- b) 当超过一半的主节点投票，则提主成功，否则故障转移超时
- c) 适当的延迟时间以避免选票瓜分



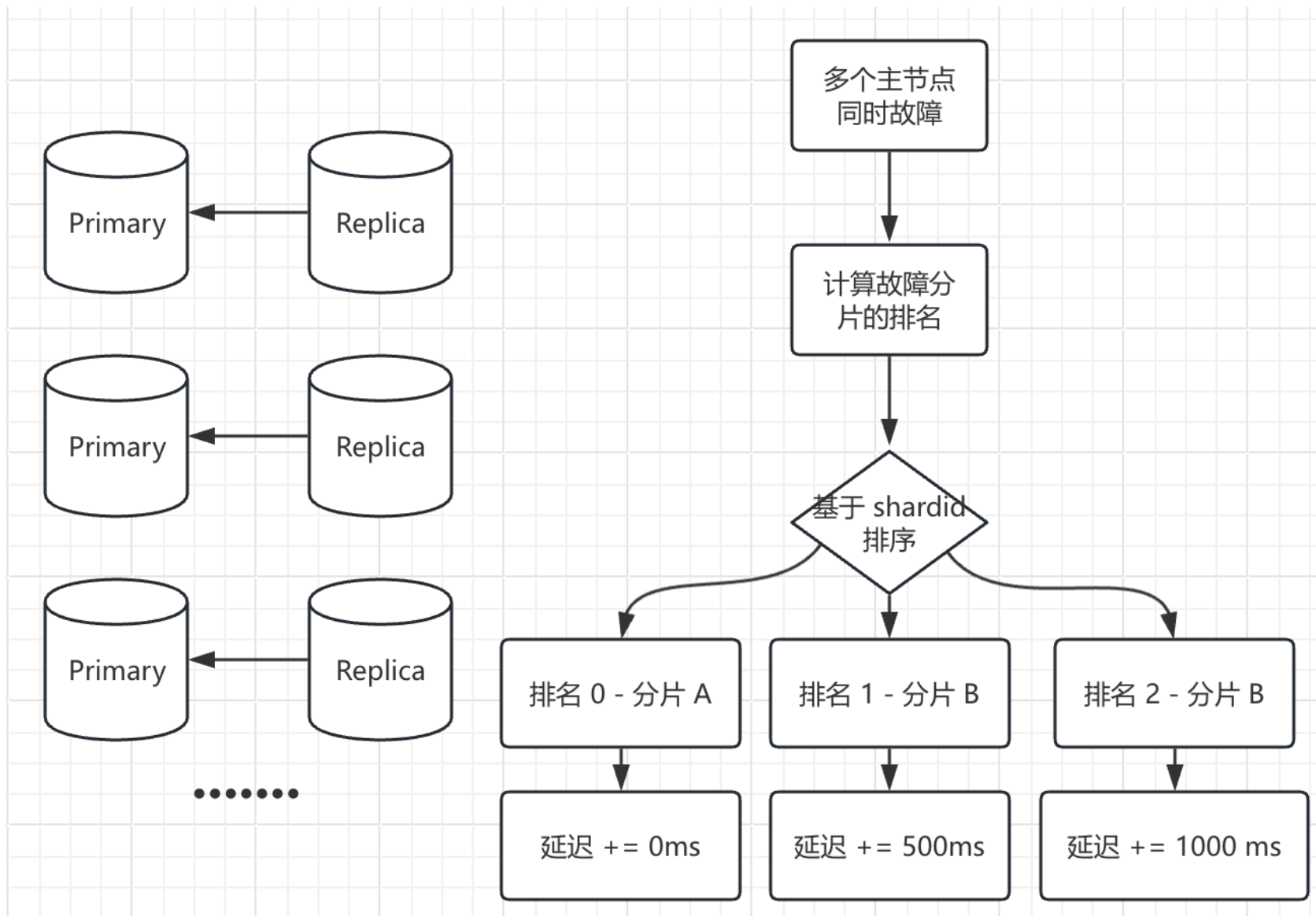
多分片有序投票

■ 多分片投票:

- 随机 500ms 以内发起投票
- 一个节点失败，集群无法可用

■ 按节点有序来选主

- 外部控制
- 仲裁节顺序授权:
 - 需要打破一个 epoch 投一个从的原则
- 从节点顺序投票:
 - 排序投票 (shard id)
 - 延时投票



Do election in order based on failed primary rank to avoid voting conflicts #1018

同可用区副本优先选举

■ 单分片的投票时机

- 随机延长 500ms 以内的时间
- 副本排名 * 1s

■ Replication Offset 相同存在的问题:

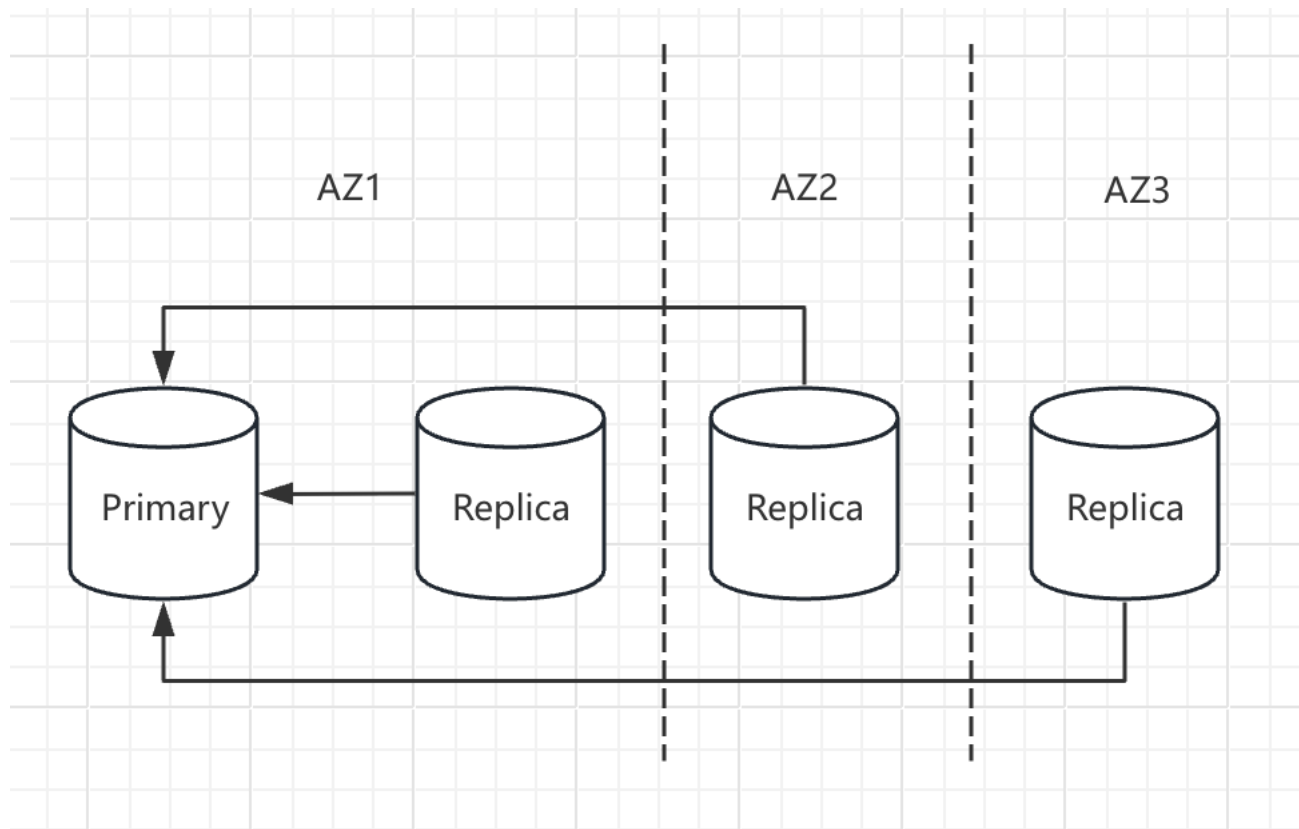
- 副本节点可能漂移到其它可用区

■ 解决办法:

- 同可用区副本优先
- 引入同可用区排名
- 可以借助 Gossip 扩展字段

■ 新排名方式为:

- (数据 offset, 是否同可用区)



Add cluster-replica-priority to allow better ranking in auto failover #2204

Valkey 新版本跟故障转移有关的优化

- 避免无法发起选举 / 避免无法完成选举
 - Optimize failover time when the new primary node is down again #782
 - Fix replica not able to initiate election in time when epoch fails #1009
 - Make manual failover reset the on-going election to promote failover #1274
 - Broadcast a PONG to all node in cluster when role changed #1295
 - Manual failover vote is not limited by two times the node timeout #1305
 - Automatic failover vote is not limited by two times the node timeout #1356
 - Fix replica can't finish failover when config epoch is outdated #2178
- 尽可能快的发起选举
 - Do the failover immediately if the replica is the best ranked replica #2227(未合入)
 - Make cluster failover delay relative to node timeout #2449
- 尽可能有序的发起选举
 - Replicas with the same offset queue up for election #762
 - Do election in order based on failed primary rank to avoid voting conflicts #1018
- 保障故障转移过程中数据安全性
 - Fix data loss when replica do a failover with a old history repl offset #885



THANKS