

分片集群中的一个写入操作

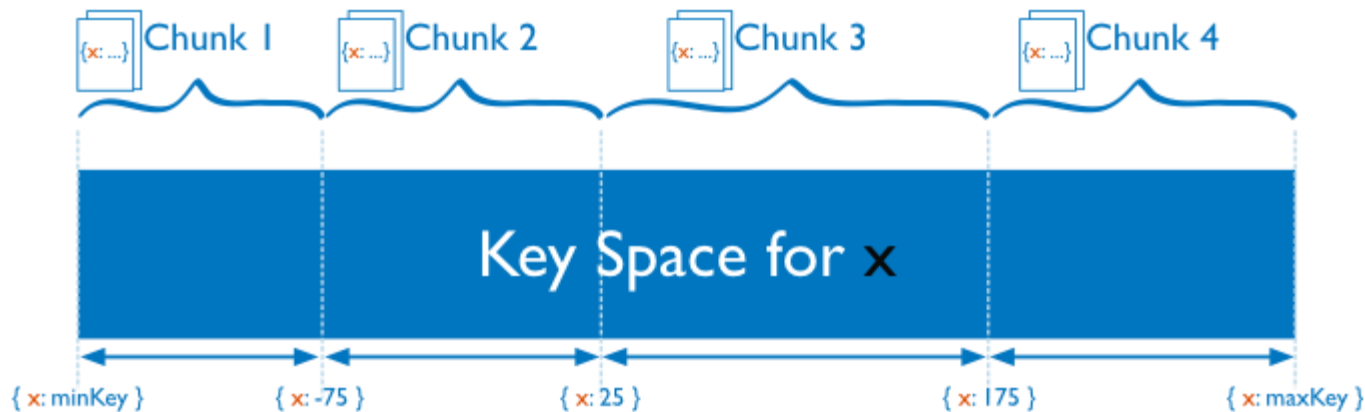
郭一然

(Based on Randolph's MDBW '16 talk)

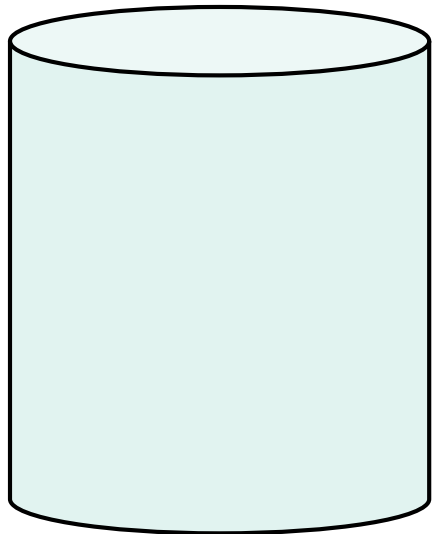
MONGODB
WORLD'16

FOR GIANT IDEAS

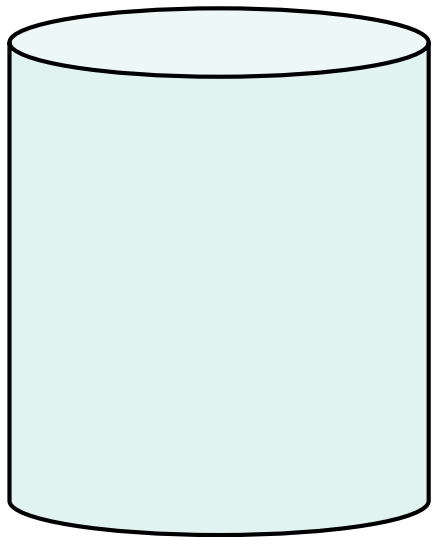
分片概念



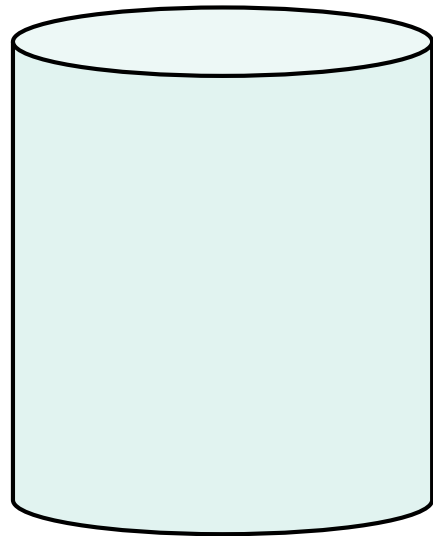
分片概念



分片0

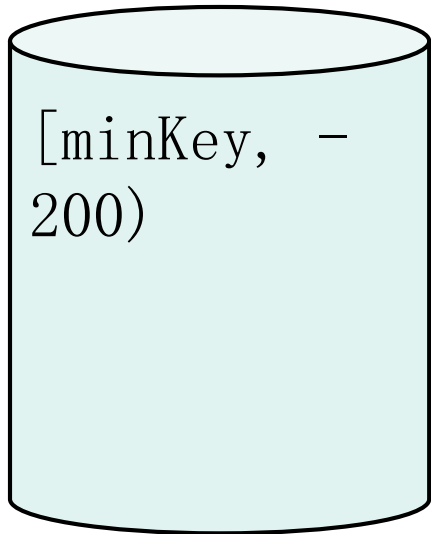


分片1

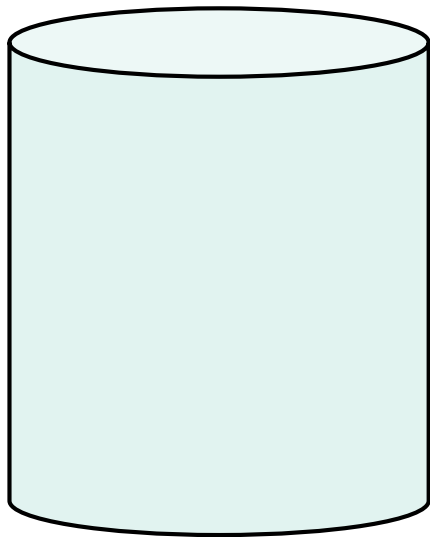


分片2

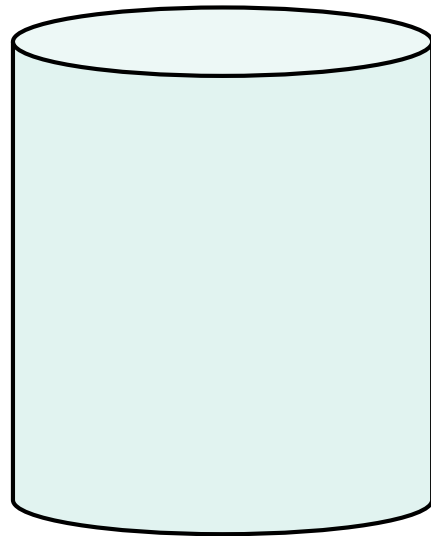
分片概念



分片0

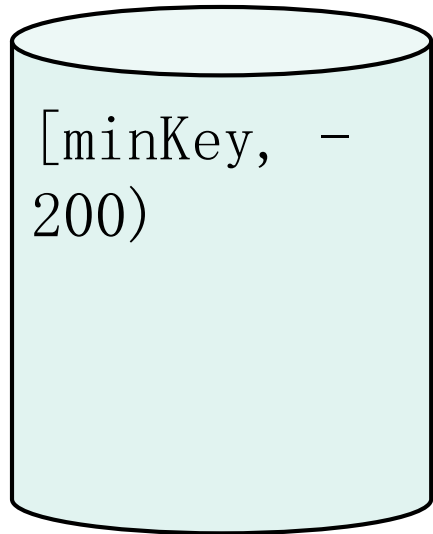


分片1

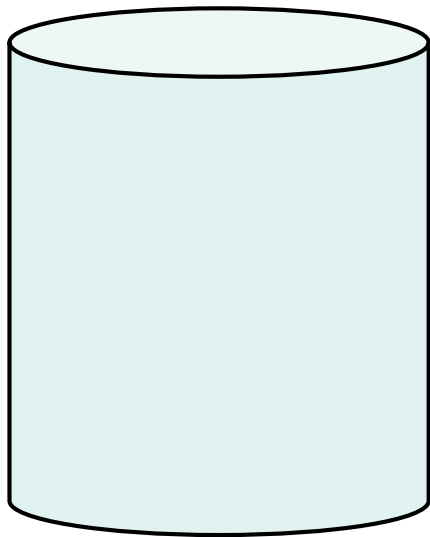


分片2

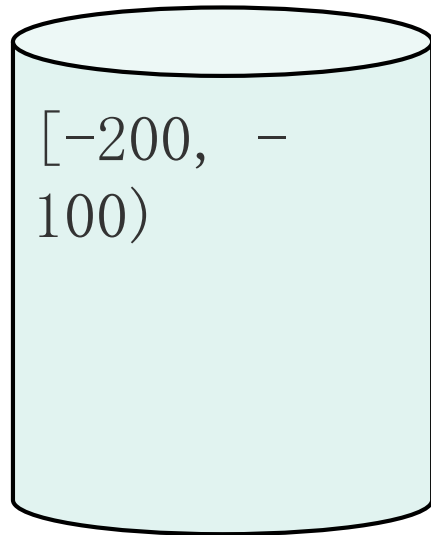
分片概念



分片0

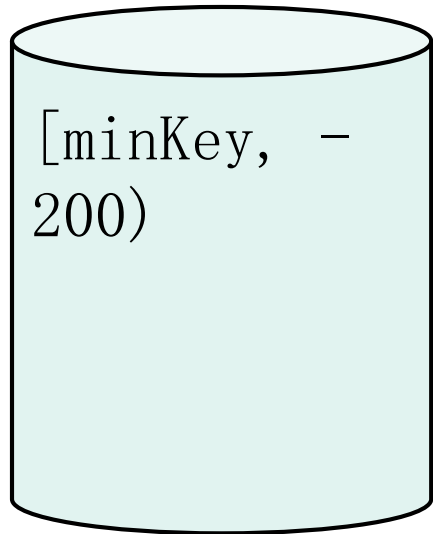


分片1

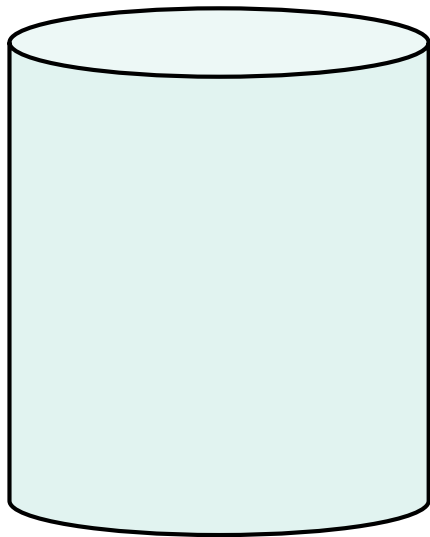


分片2

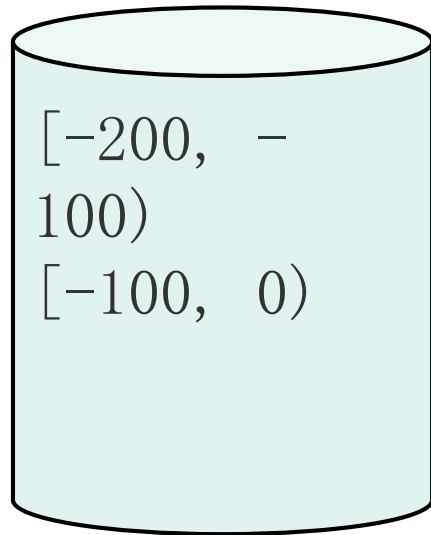
分片概念



分片0

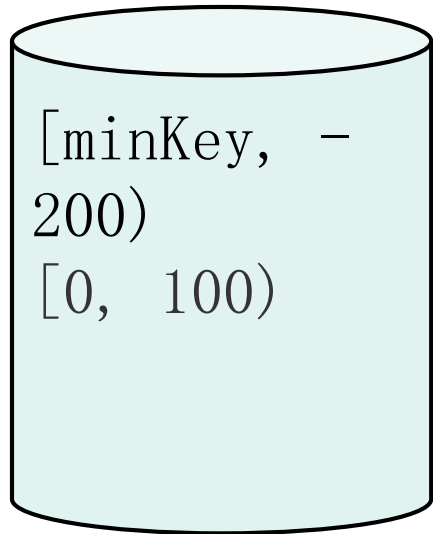


分片1

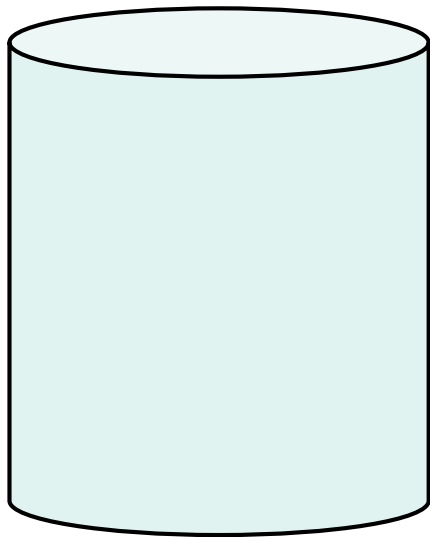


分片2

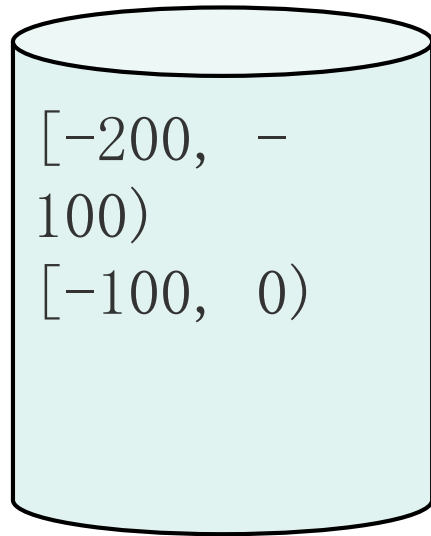
分片概念



分片0

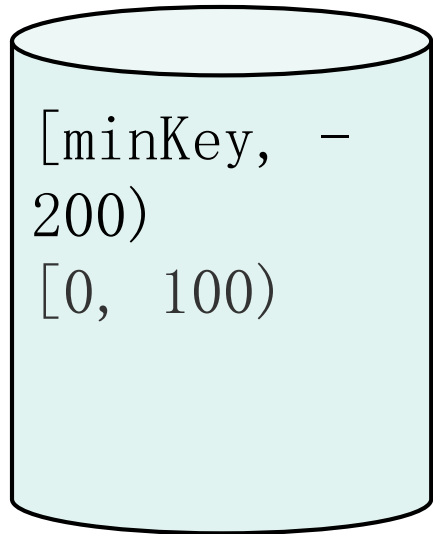


分片1

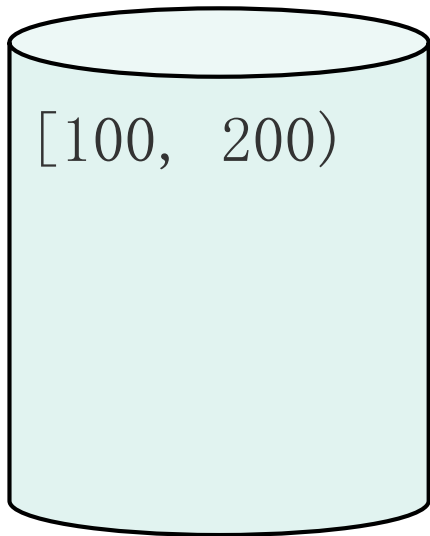


分片2

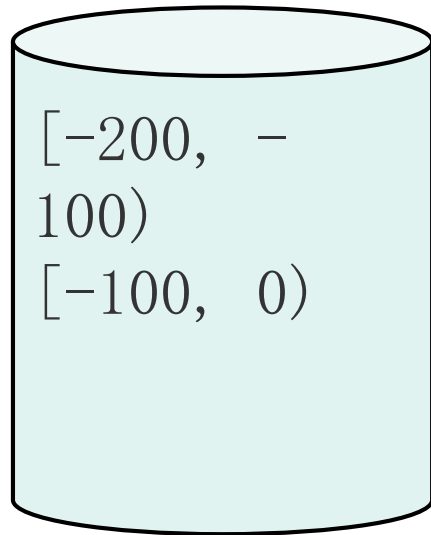
分片概念



分片0

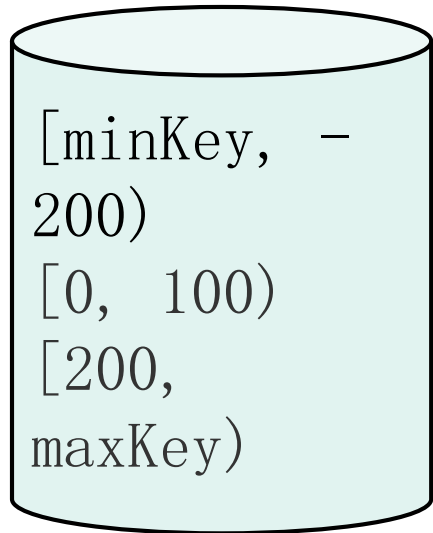


分片1

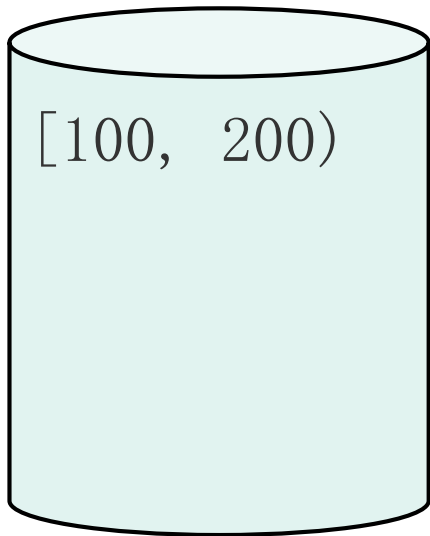


分片2

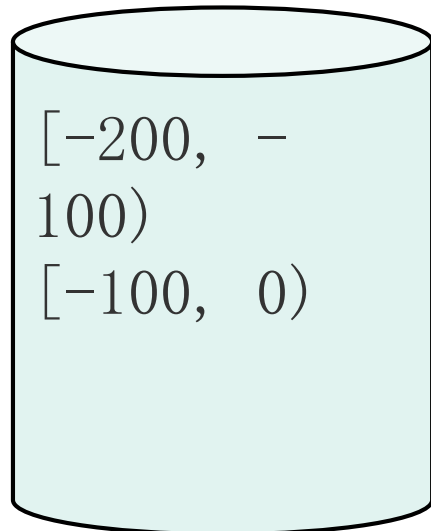
分片概念



分片0



分片1

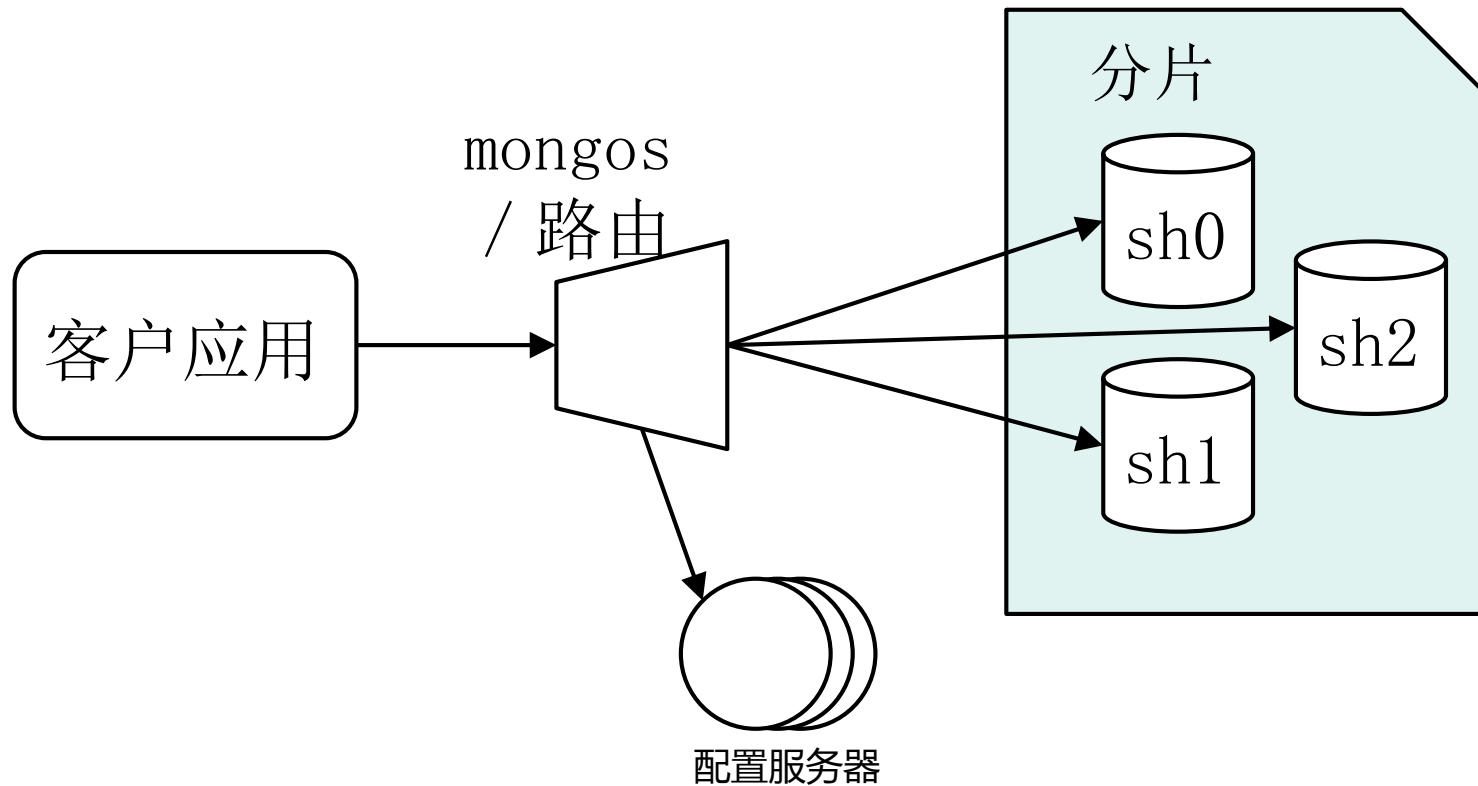


分片2

路由表

最小值	最大值	所在分片
MinKey	-200	分片0
-200	-100	分片2
-100	0	分片2
0	100	分片0
100	200	分片1
200	MaxKey	分片0

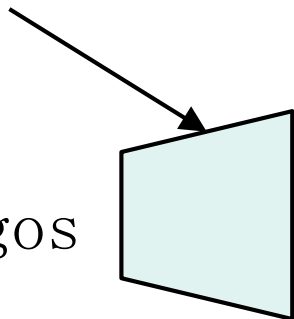
分片集群架构



单文档写入

```
db.rio.insert({ 金牌:  
8 })
```

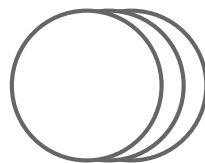
mongos



分片0

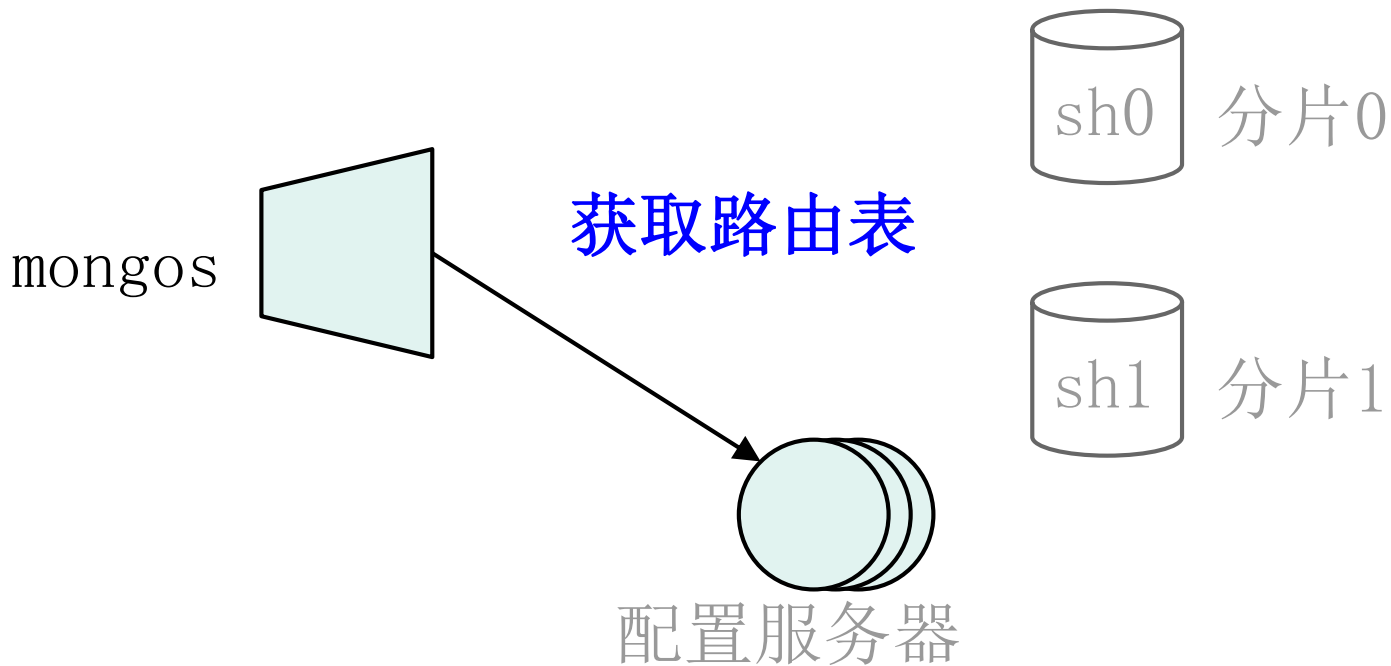


分片1



配置服务器

单文档写入



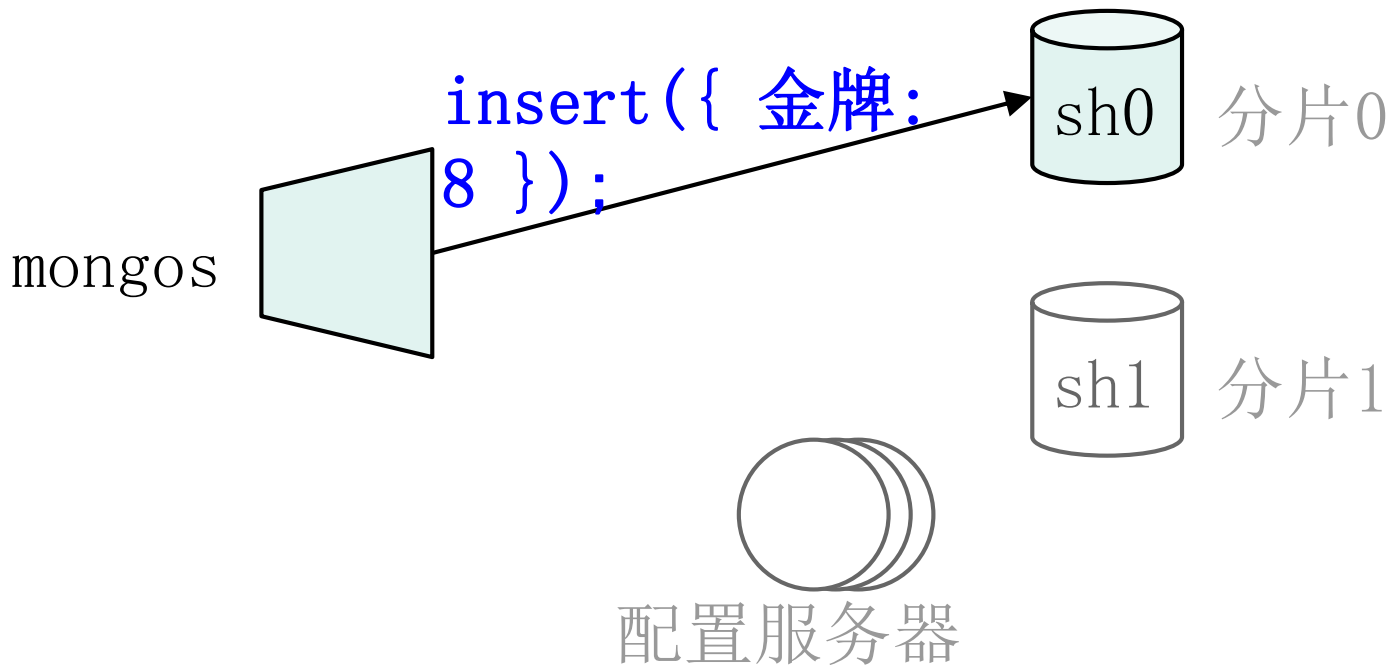
路由表

db.rio.insert({ 金

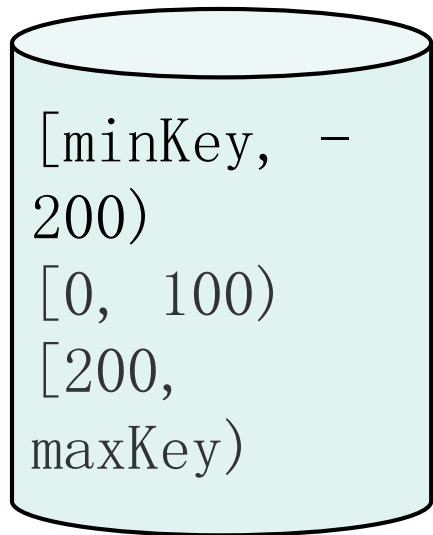
牌: 8
最小值 最大值

最小值	最大值	所在分片
MinKey	-200	分片0
-200	-100	分片2
-100	0	分片2
0	100	分片0
100	200	分片1
200	MaxKey	分片0

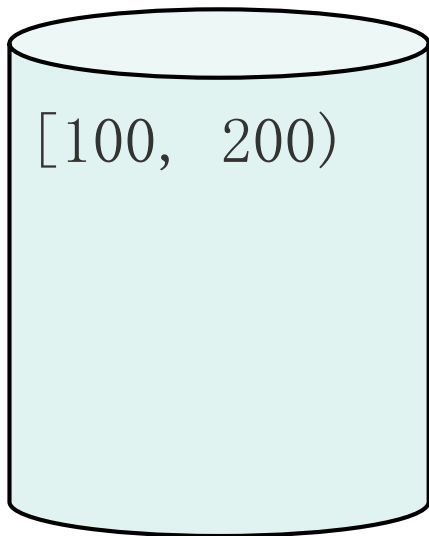
单文档写入



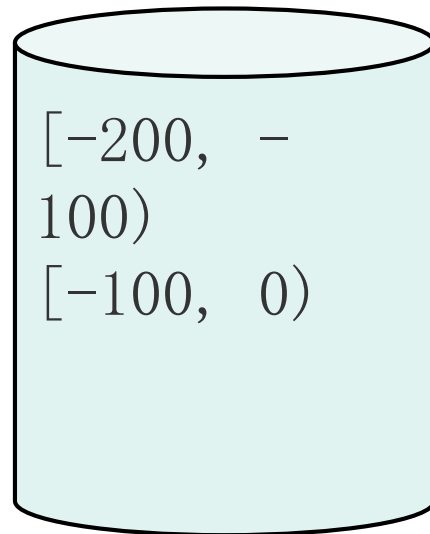
数据段迁徙



分片0

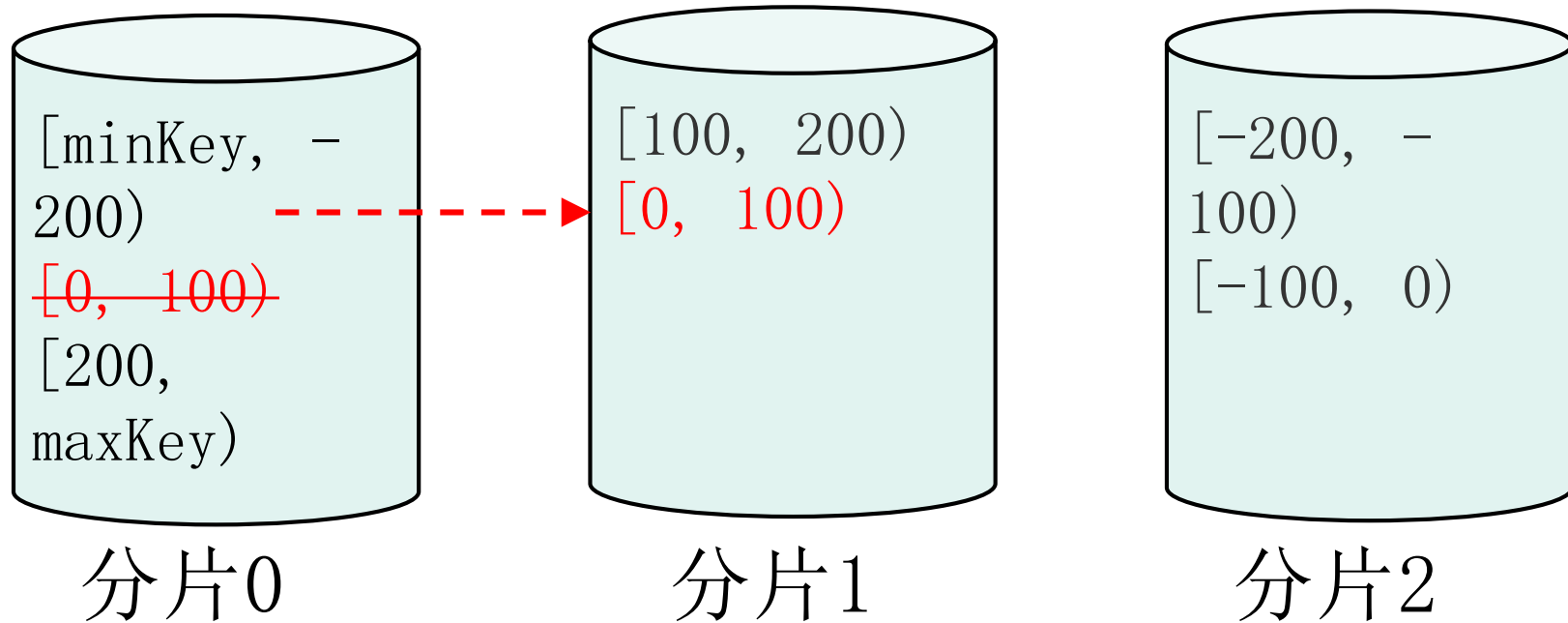


分片1



分片2

数据段迁徙



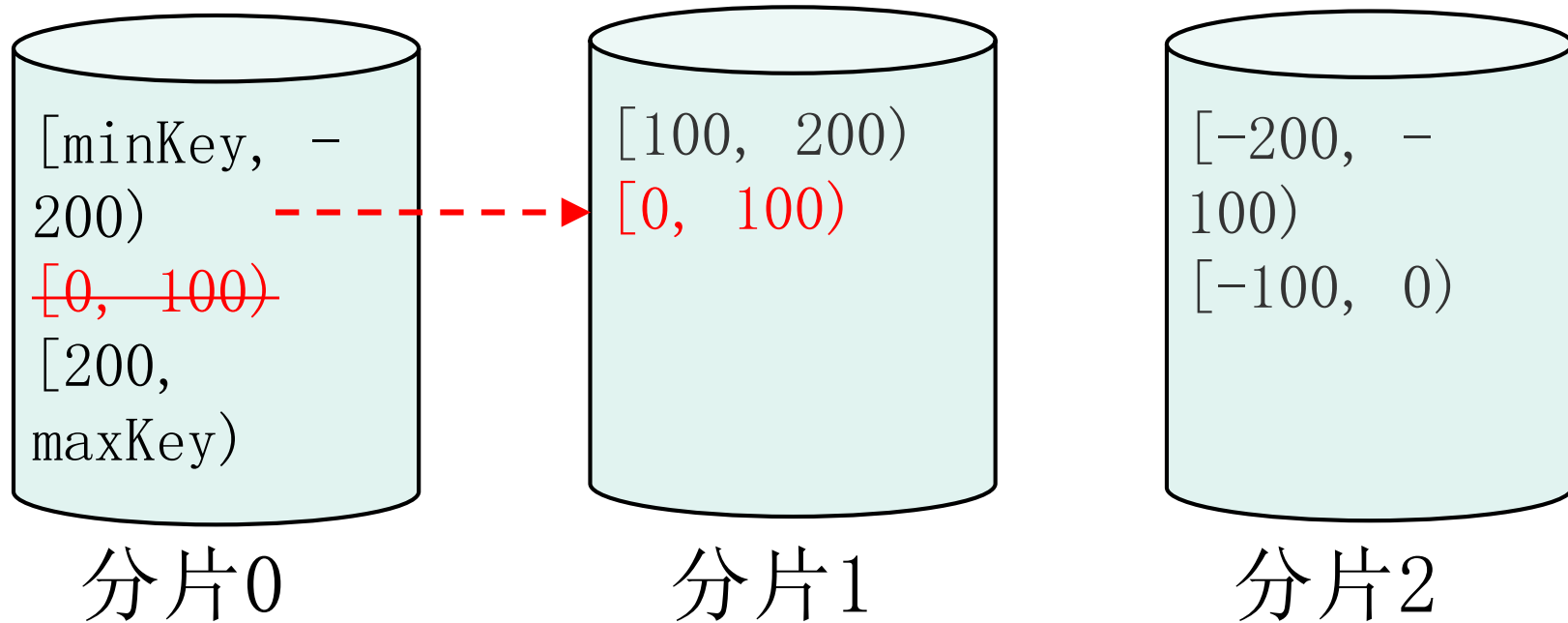
分片版本协议

版本	最小值	最大值	所在分片
1	MinKey	-200	分片0
2	-200	-100	分片2
3	-100	0	分片2
4	0	100	分片0
5	100	200	分片1
6	200	MaxKey	分片0

分片版本协议

$$\text{新版本} = \max(\text{路由表版本}) + 1$$

数据段迁徙



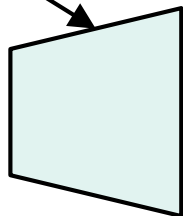
分片版本协议

版本	最小值	最大值	所在分片
1	MinKey	-200	分片0
2	-200	-100	分片2
3	-100	0	分片2
4 -> 7	0	100	分片0 -> 分片1
5	100	200	分片1
6	200	MaxKey	分片0

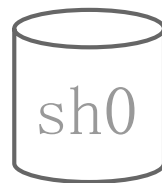
单文档写入

```
db.rio.insert({ 金牌: 8 });
```

mongos



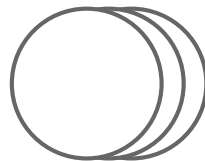
路由表v6



分片0
路由表v7



分片1
路由表 v7

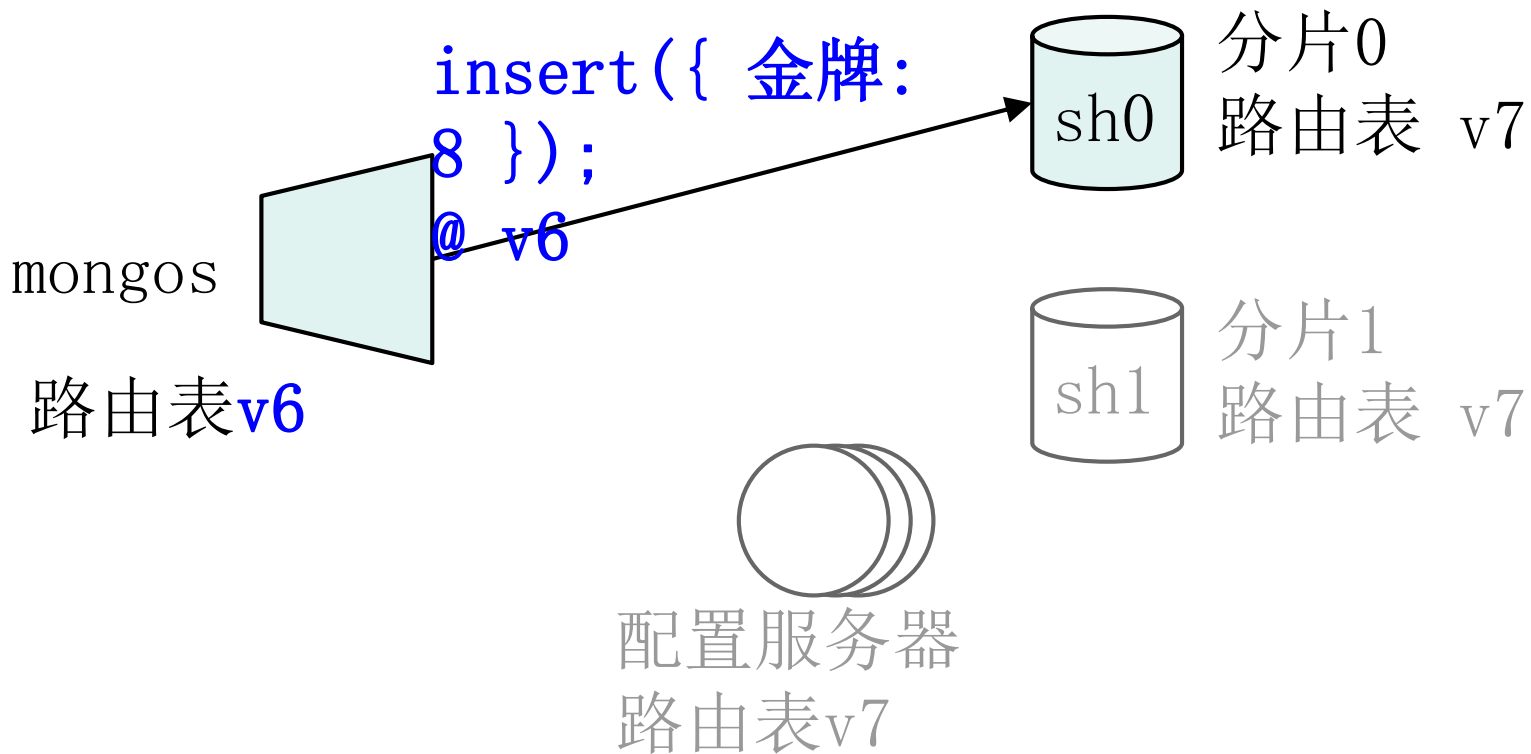


配置服务器路由
表v7

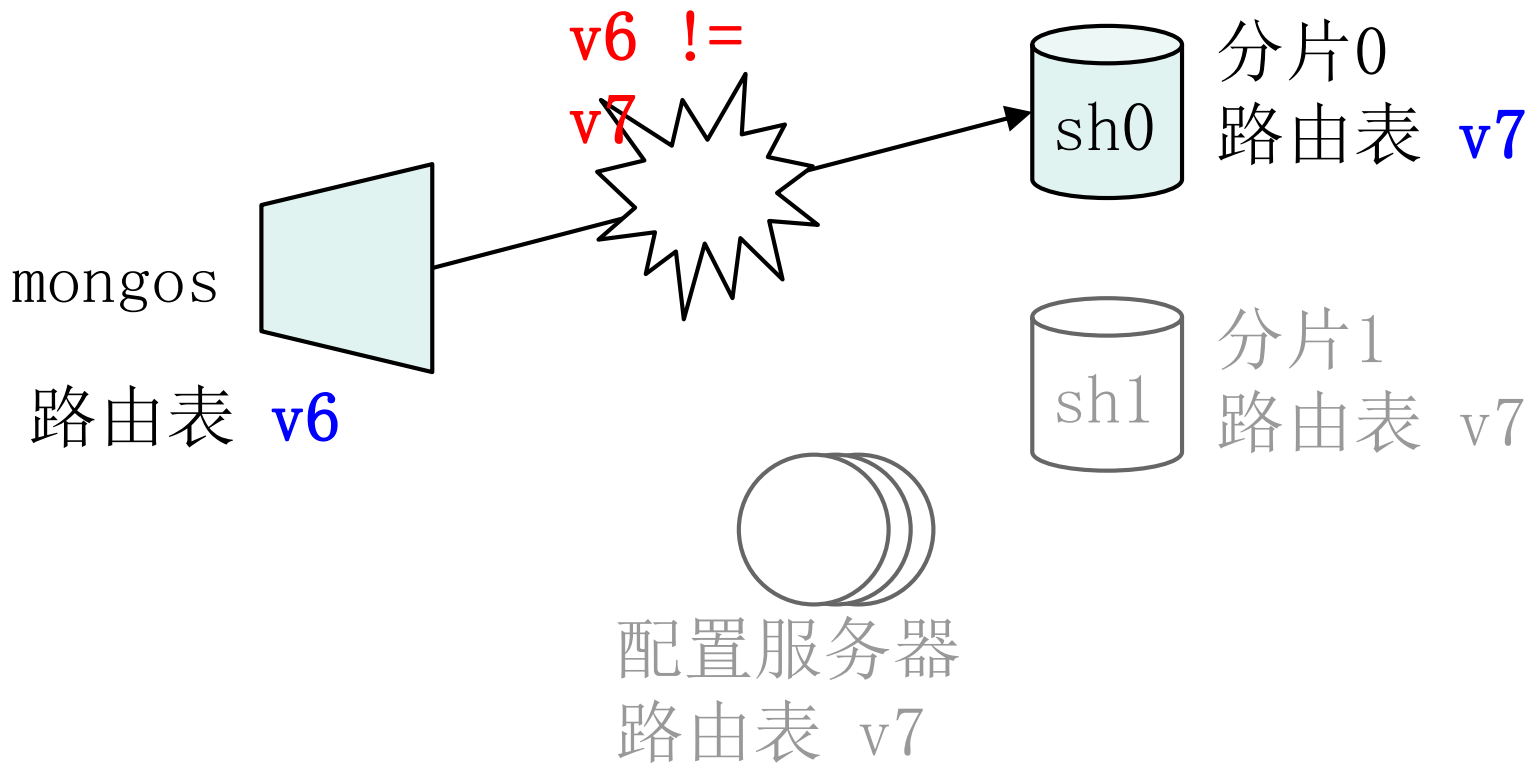
分片版本协议

版本	最小值	最大值	所在分片
1	MinKey	-200	分片0
2	-200	-100	分片2
3	-100	0	分片2
4	0	100	分片0
5	100	200	分片1
6	200	MaxKey	分片0

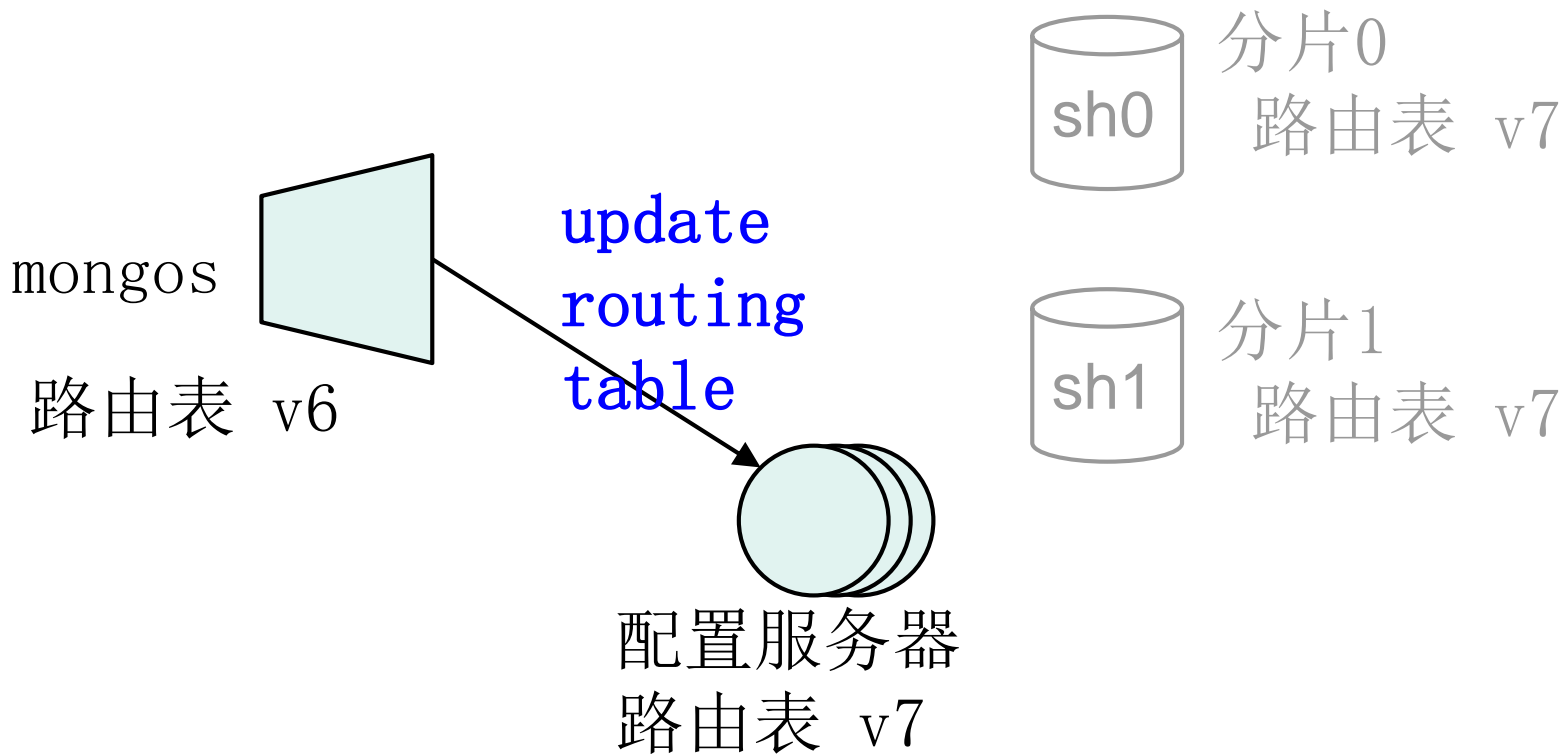
单文档写入



单文档写入



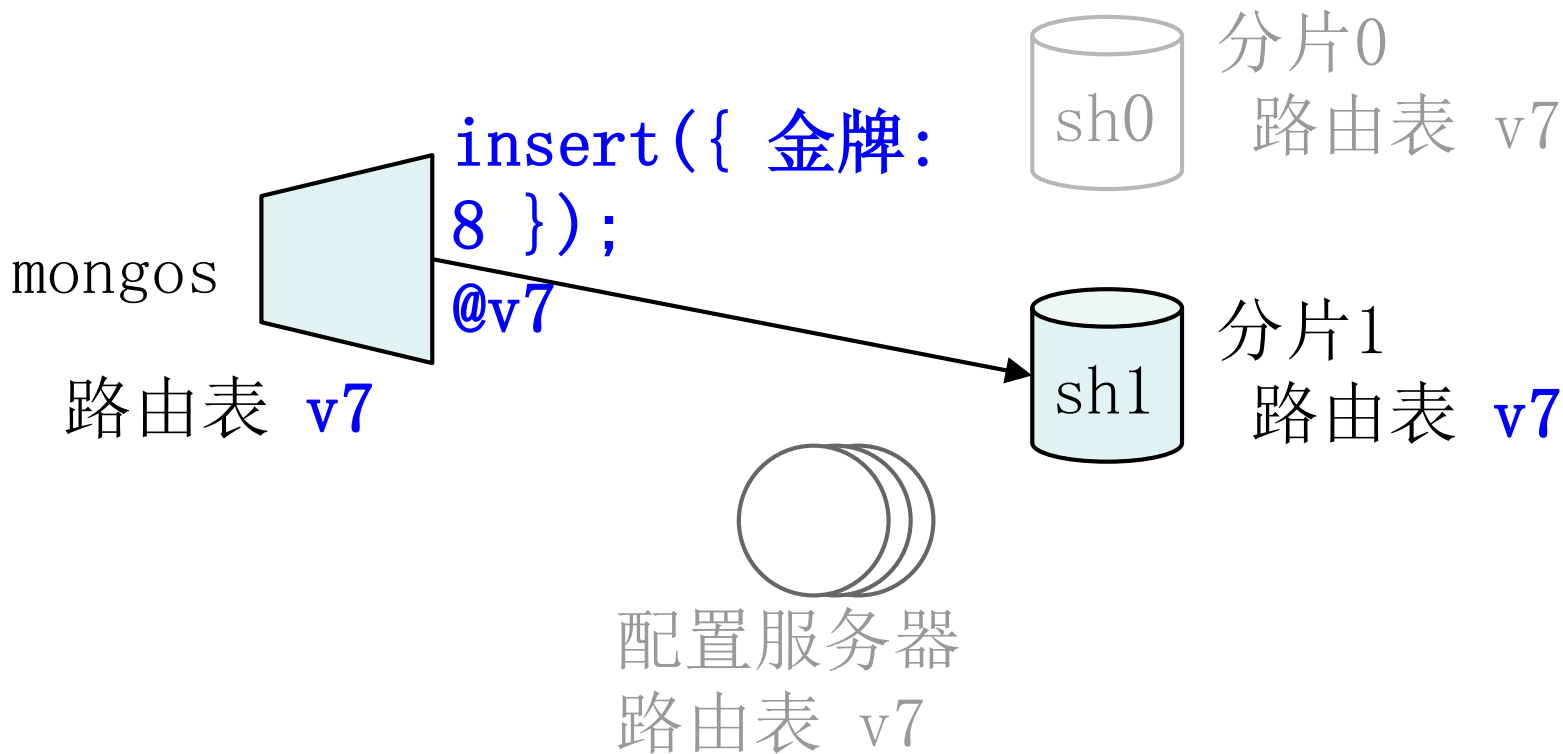
单文档写入



分片版本协议

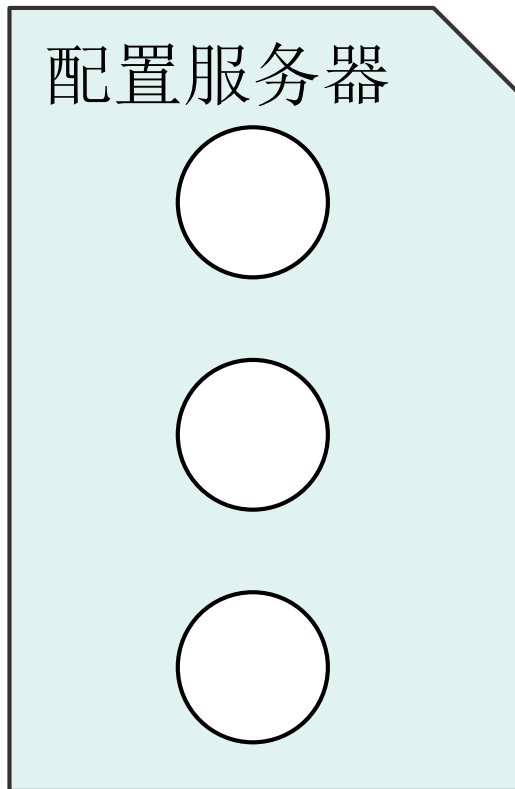
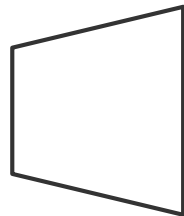
版本	最小值	最大值	所在分片
1	MinKey	-200	分片0
2	-200	-100	分片2
3	-100	0	分片2
7	0	100	分片1
5	100	200	分片1
6	200	MaxKey	分片0

单文档写入

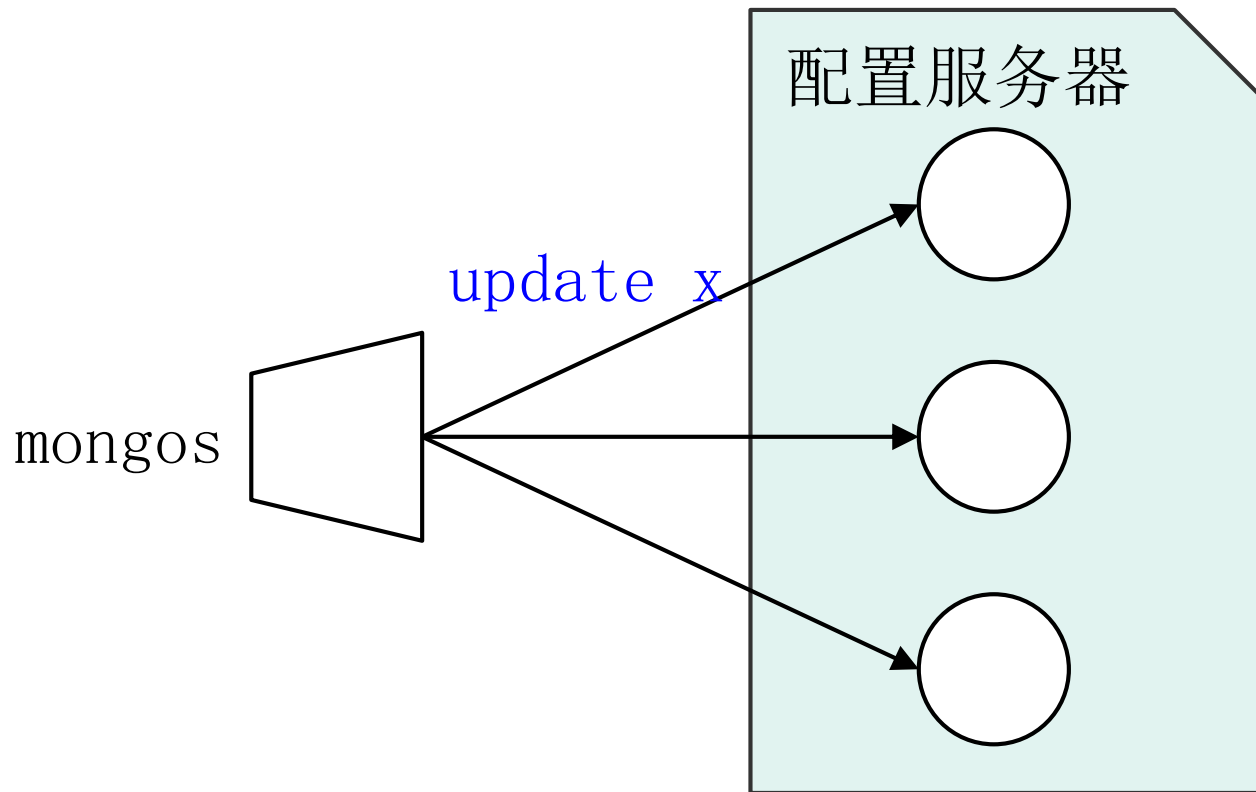


镜像配置服务器

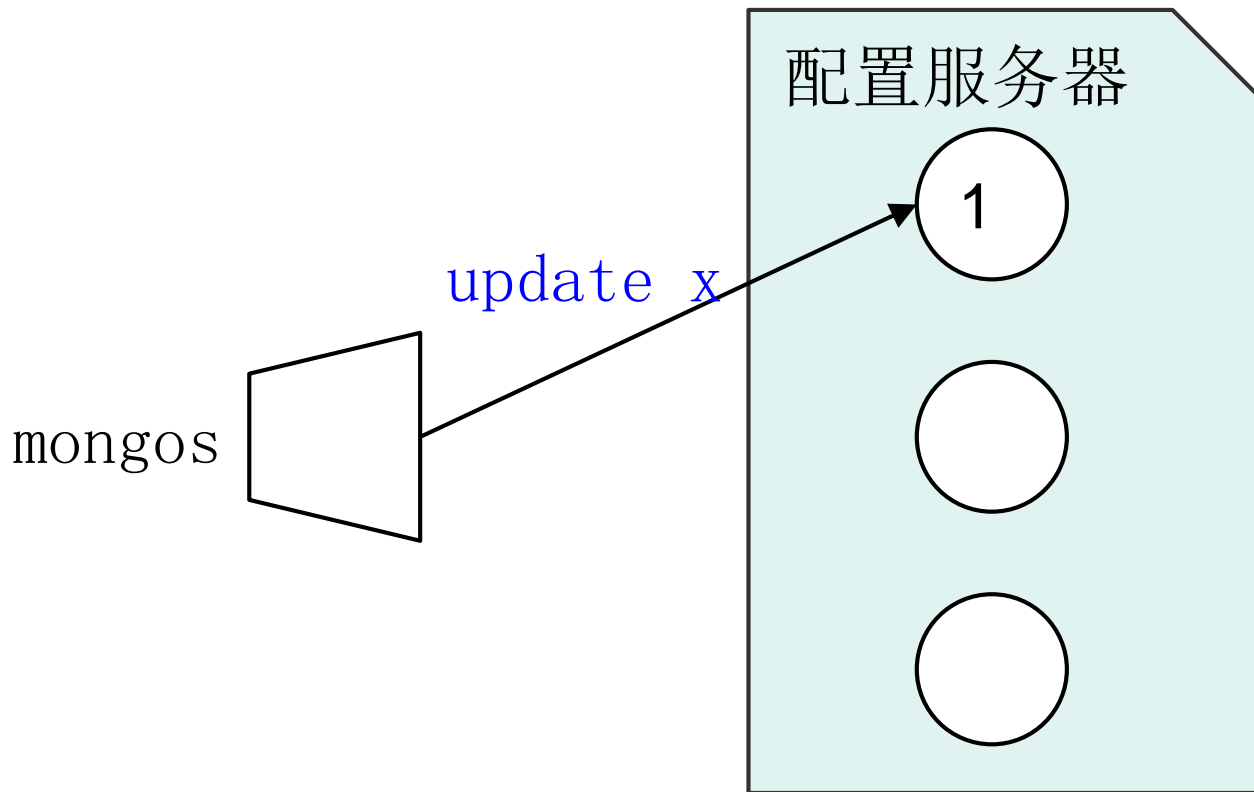
mongos



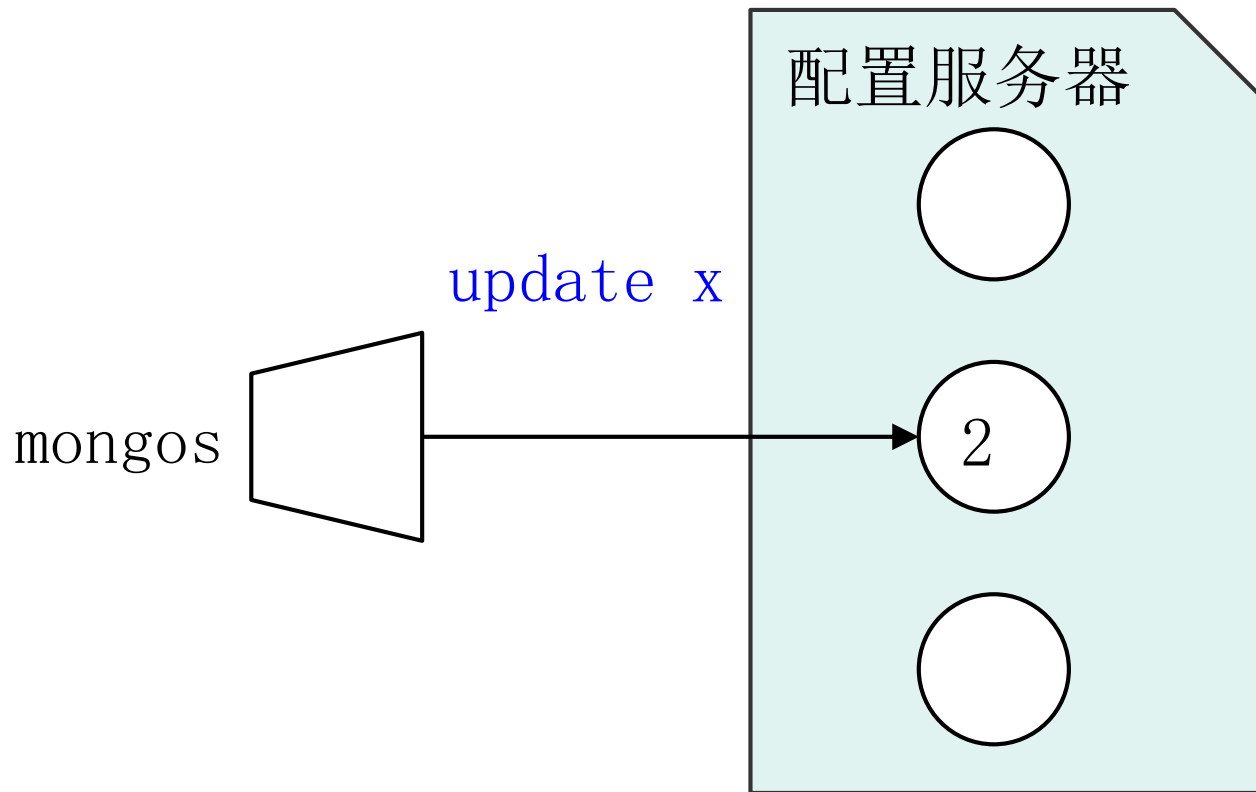
镜像配置服务器



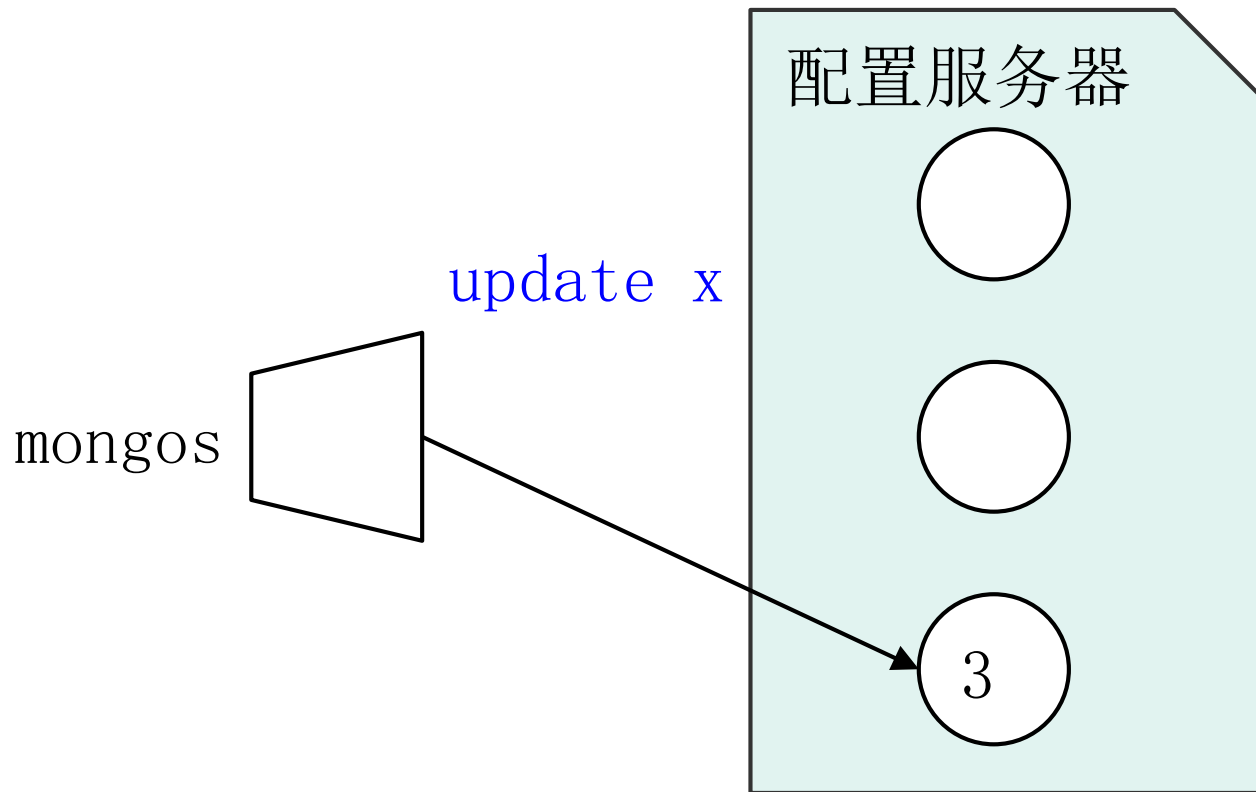
镜像配置服务器



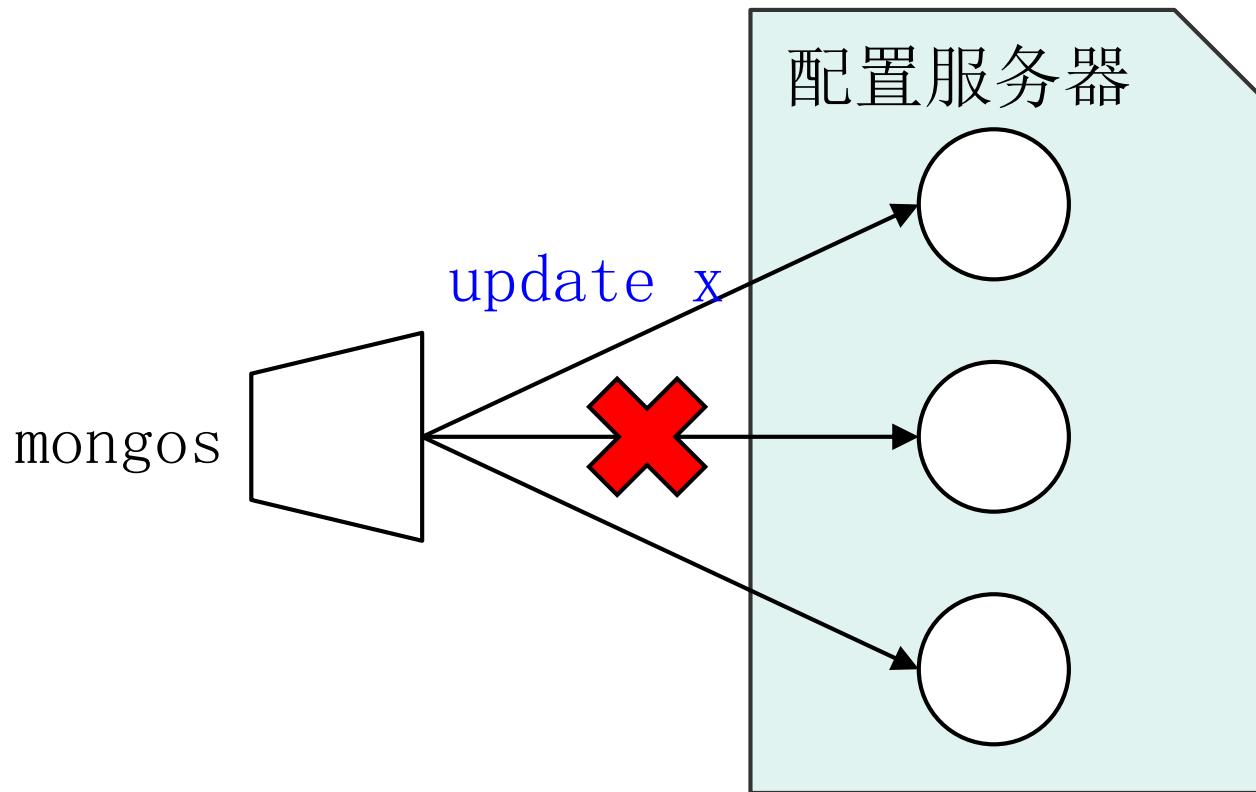
镜像配置服务器



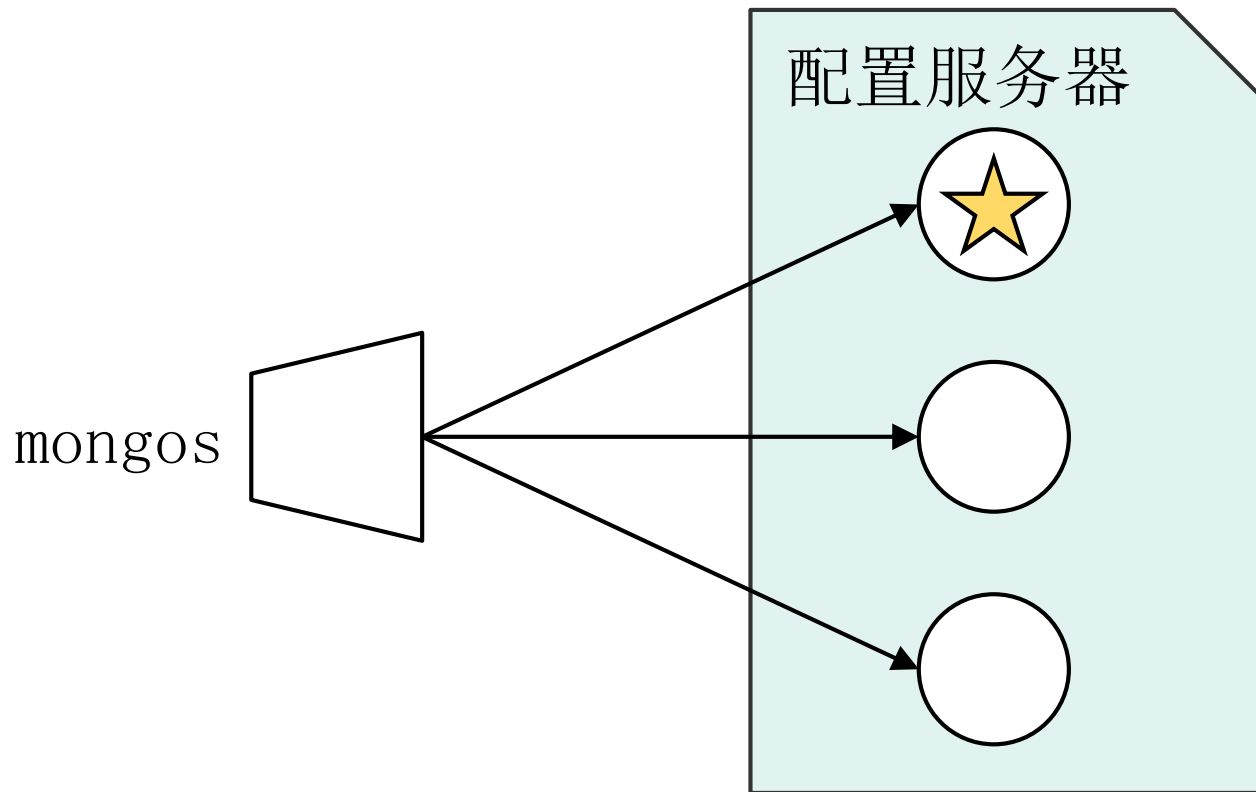
镜像配置服务器



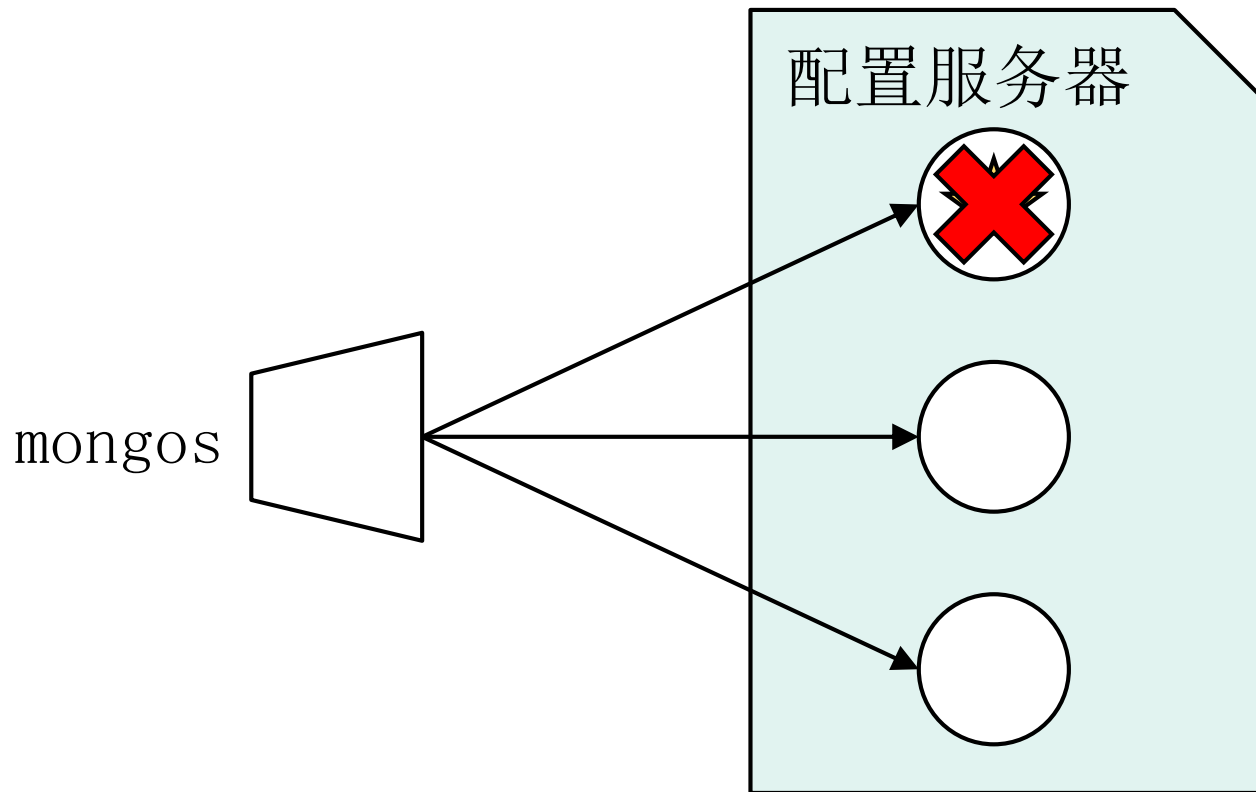
镜像配置服务器



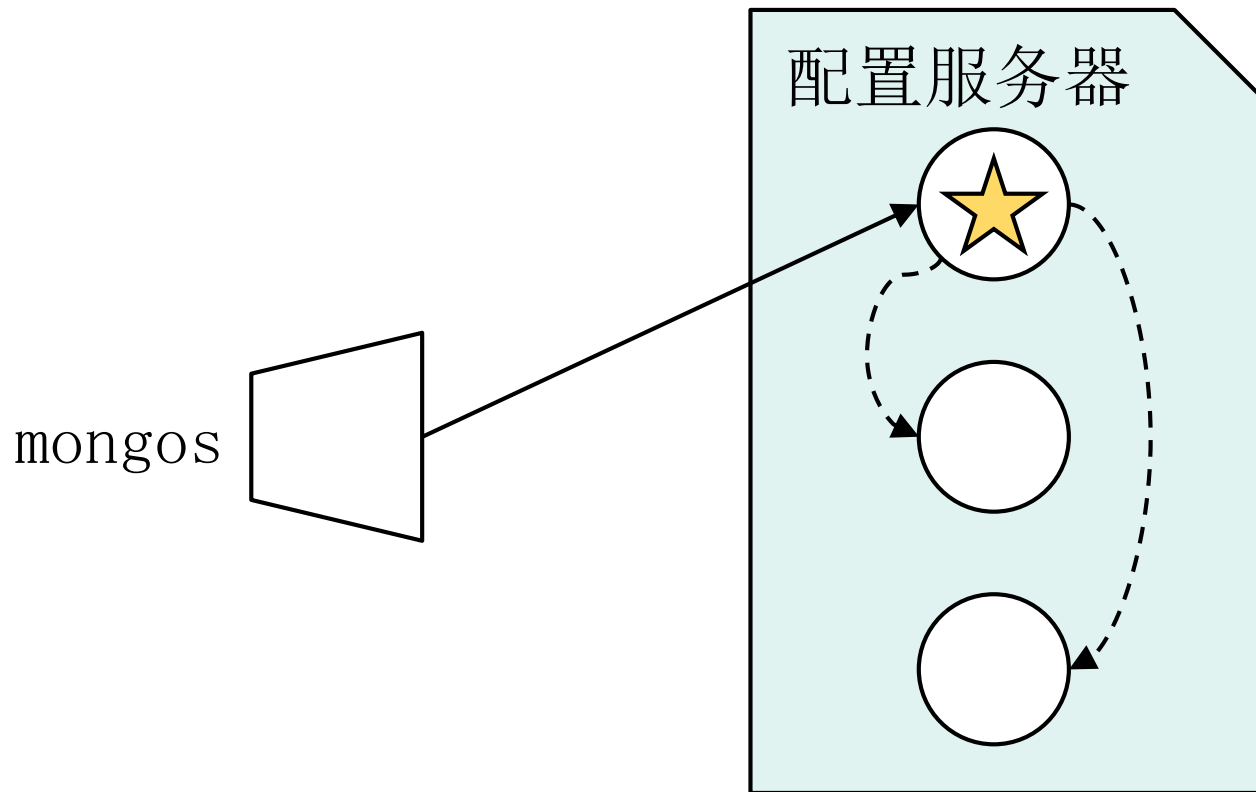
镜像配置服务器



镜像配置服务器



复制集配置服务器



镜像 VS 复制集配置服务器

镜像

复制集

失去数据一致性

镜像 VS 复制集

镜像

失去数据一致性

复制集

单一数据来源

镜像 VS 复制集

镜像

失去数据一致性

失去成员后变成只读

复制集

单一数据来源

镜像 VS 复制集

镜像

失去数据一致性

失去成员后变成只读

复制集

单一数据来源

竞选新主节点

镜像 VS 复制集

镜像

失去数据一致性

失去成员后变成只读

蛋疼的分布式锁

复制集

单一数据来源

竞选新主节点

镜像 VS 复制集

镜像

失去数据一致性

失去成员后变成只读

蛋疼的分布式锁

复制集

单一数据来源

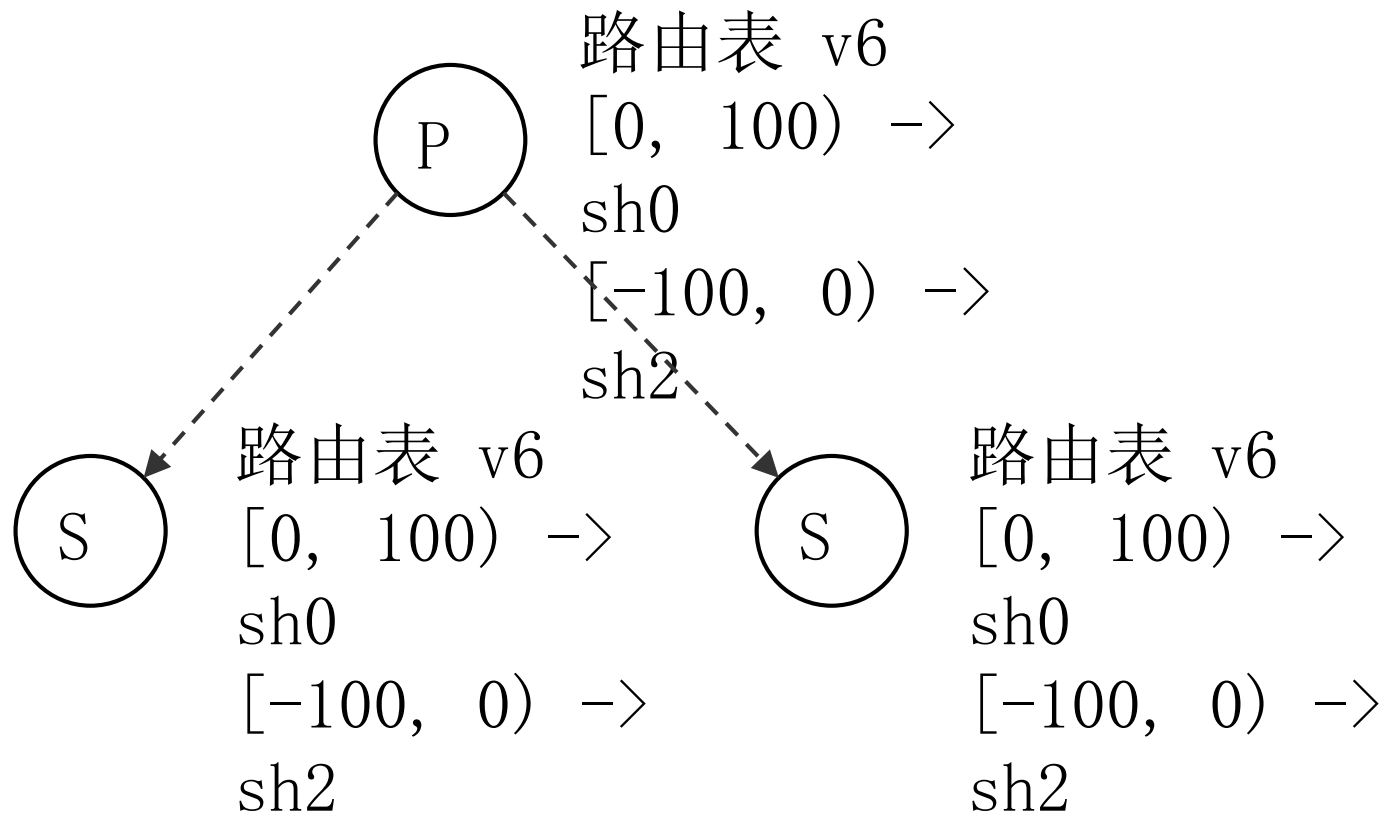
竞选新主节点

只需锁住主节点

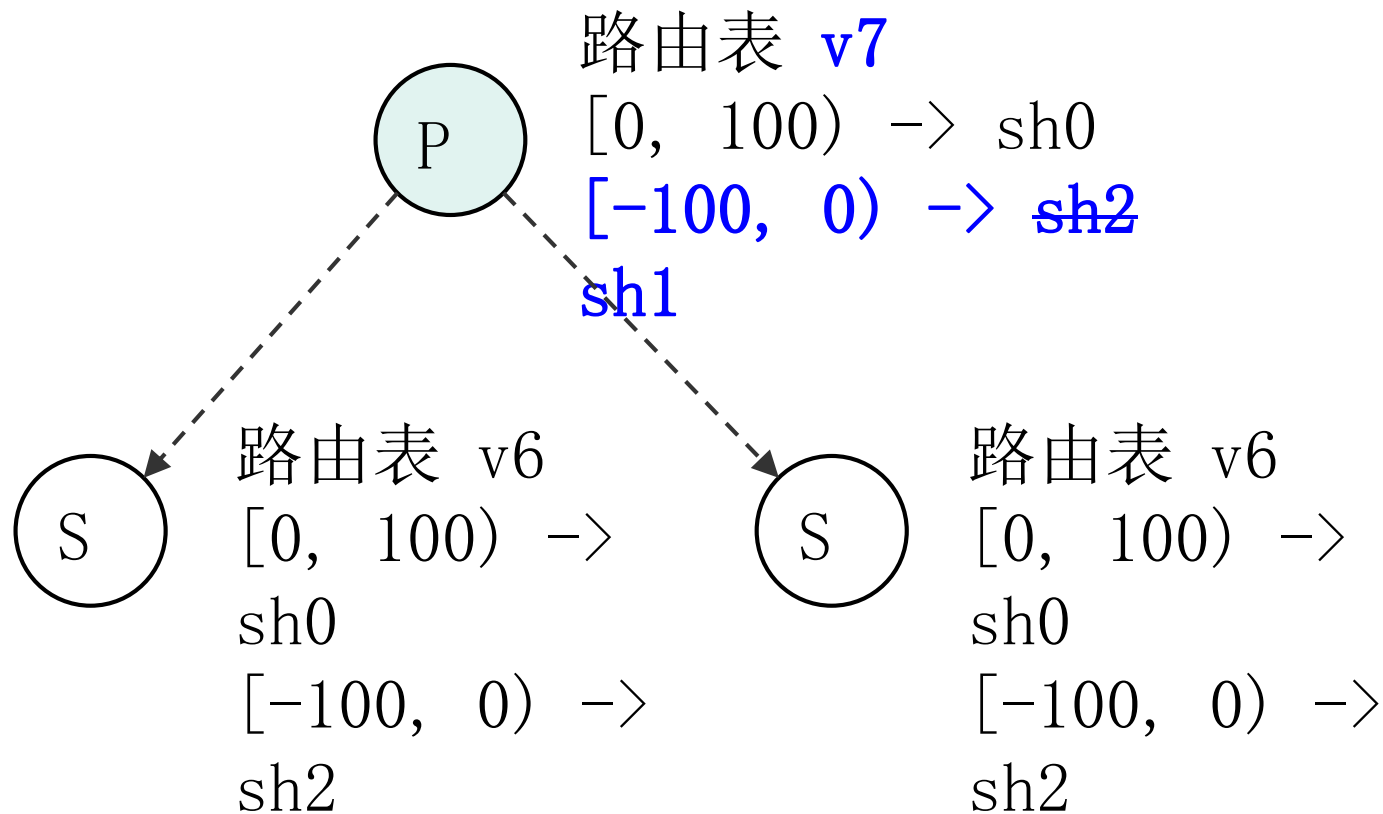
复制集配置服务器面临的挑战

- 已读取的数据可能回滚
- 从节点上的数据会过于陈旧

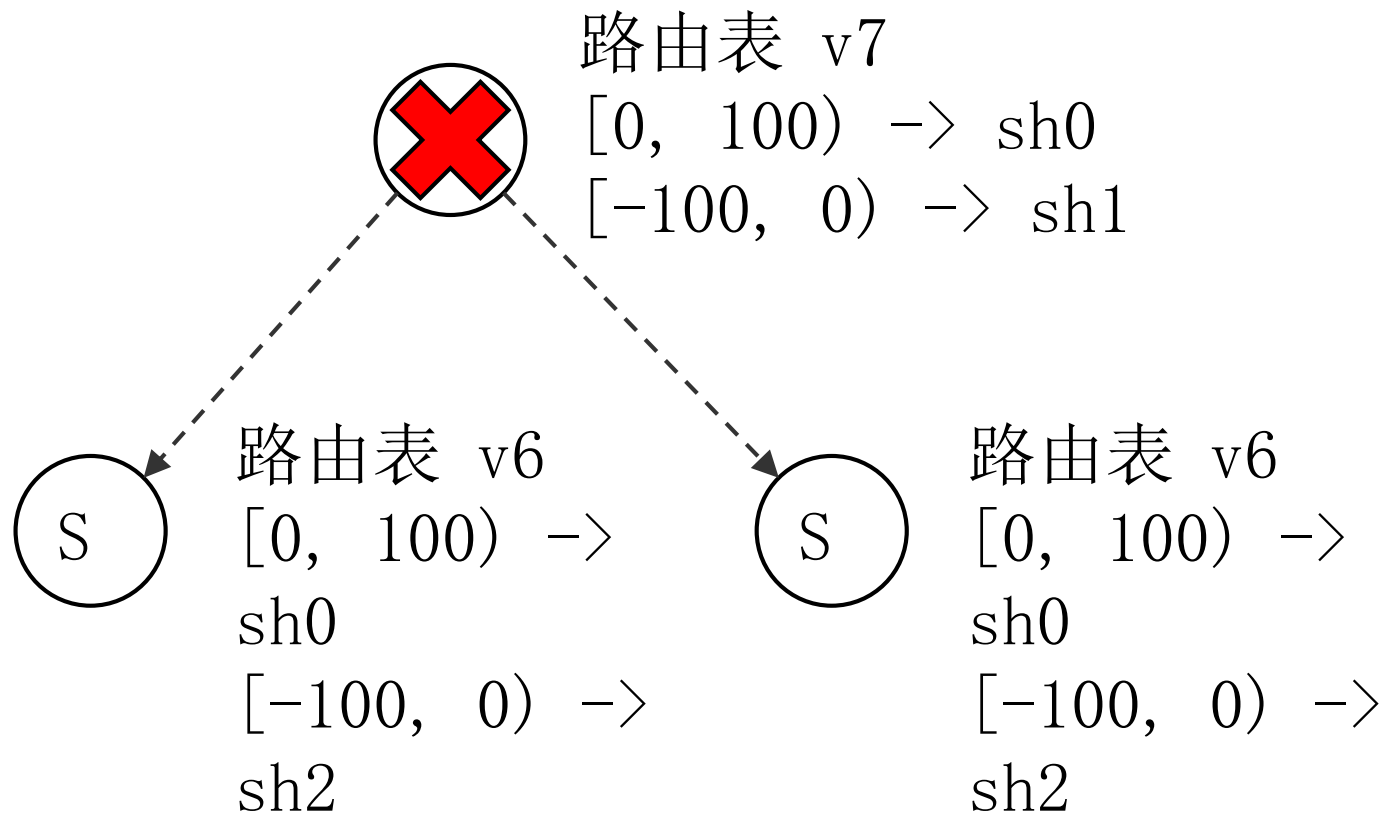
回滚案例



回滚案例



回滚案例



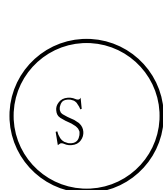
回滚案例



路由表 v7

$[0, 100) \rightarrow \text{sh0}$

$[-100, 0) \rightarrow \text{sh1}$



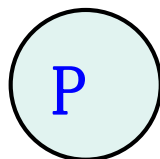
路由表 v6

$[0, 100) \rightarrow$

sh0

$[-100, 0) \rightarrow$

sh2



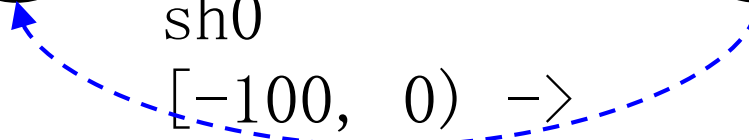
路由表 v6

$[0, 100) \rightarrow$

sh0

$[-100, 0) \rightarrow$

sh2



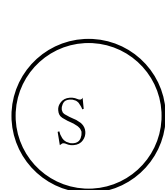
回滚案例



路由表 v7

$[0, 100) \rightarrow \text{sh0}$

$[-100, 0) \rightarrow \text{sh1}$



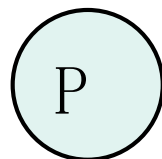
路由表 v6

$[0, 100) \rightarrow$

sh0

$[-100, 0) \rightarrow$

sh2



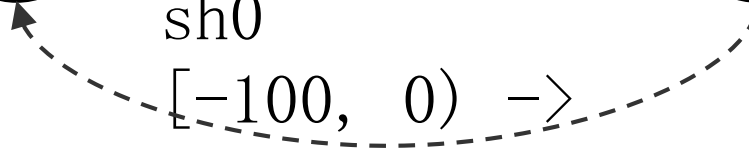
路由表 v7

$[0, 100) \rightarrow \text{sh0}$

sh1

$[-100, 0) \rightarrow$

sh2



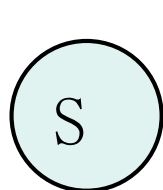
回滚案例



路由表 v7

$[0, 100) \rightarrow \text{sh0}$

$[-100, 0) \rightarrow \text{sh1}$



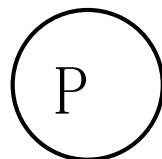
路由表 v7

$[0, 100) \rightarrow$

sh1

$[-100, 0) \rightarrow$

sh2

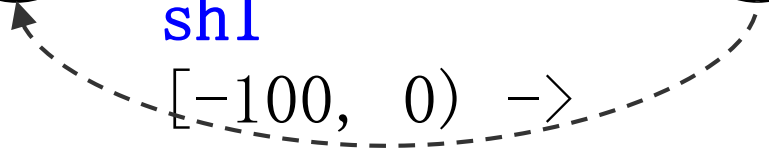


路由表 v7

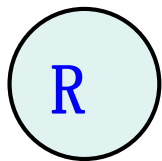
$[0, 100) \rightarrow \text{sh1}$

$[-100, 0) \rightarrow$

sh2



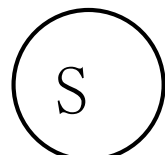
回滚案例



路由表 v7

$[0, 100) \rightarrow \text{sh0}$

$[-100, 0) \rightarrow \text{sh1}$



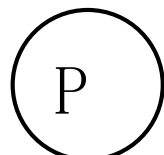
路由表 v7

$[0, 100) \rightarrow$

sh1

$[-100, 0) \rightarrow$

sh2

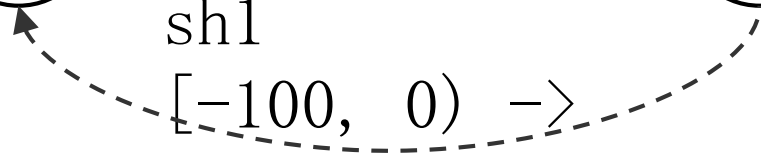


路由表 v7

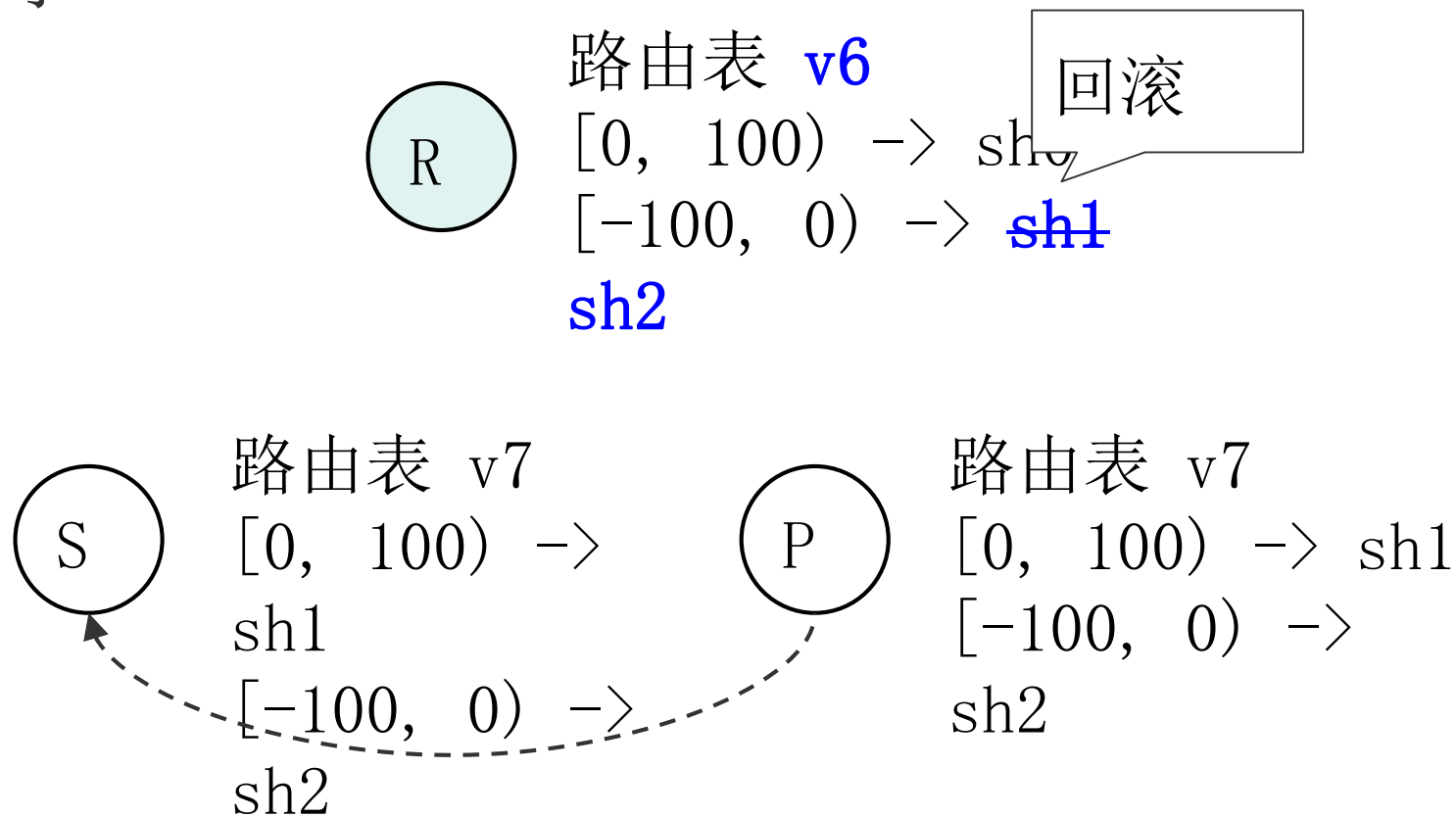
$[0, 100) \rightarrow \text{sh1}$

$[-100, 0) \rightarrow$

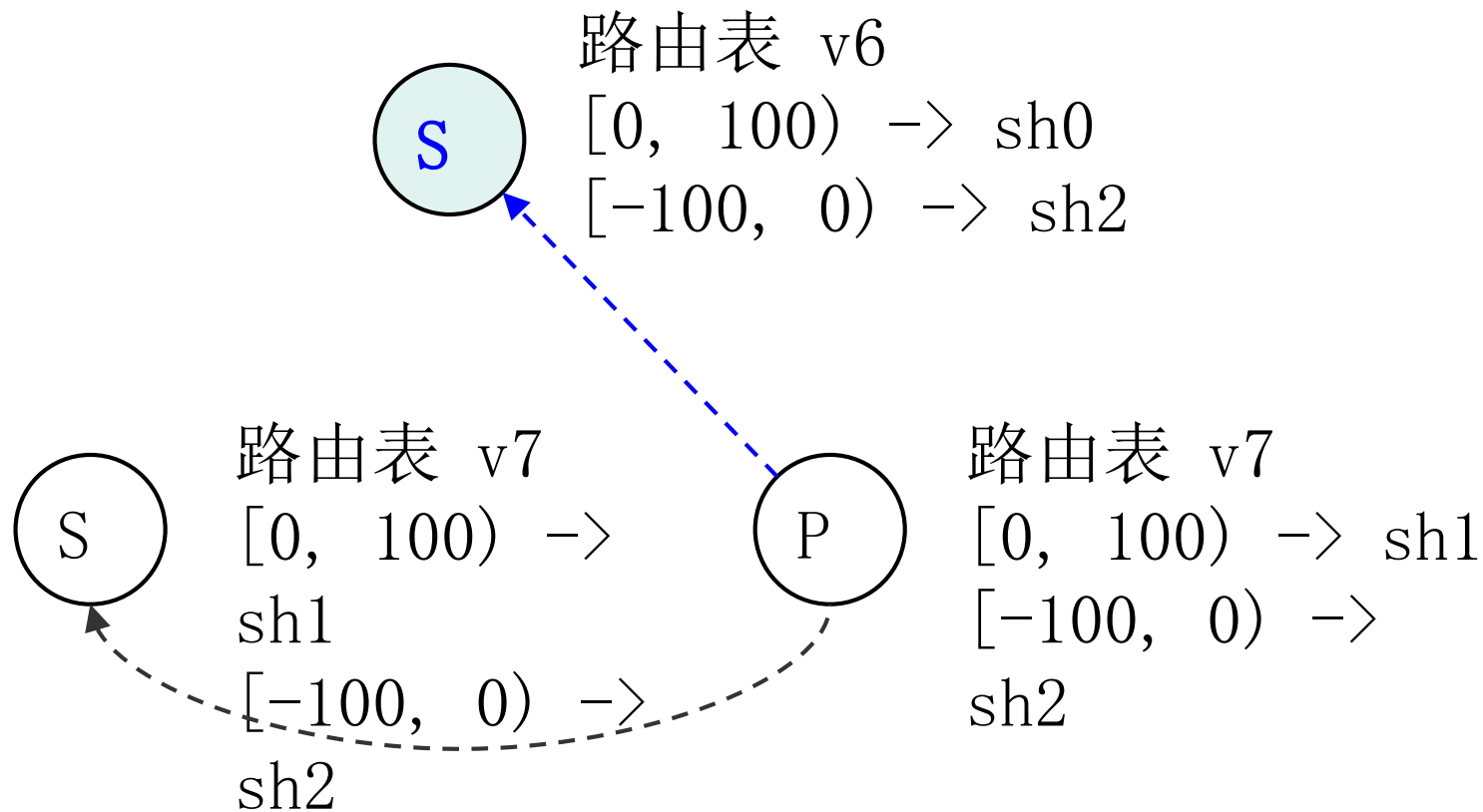
sh2



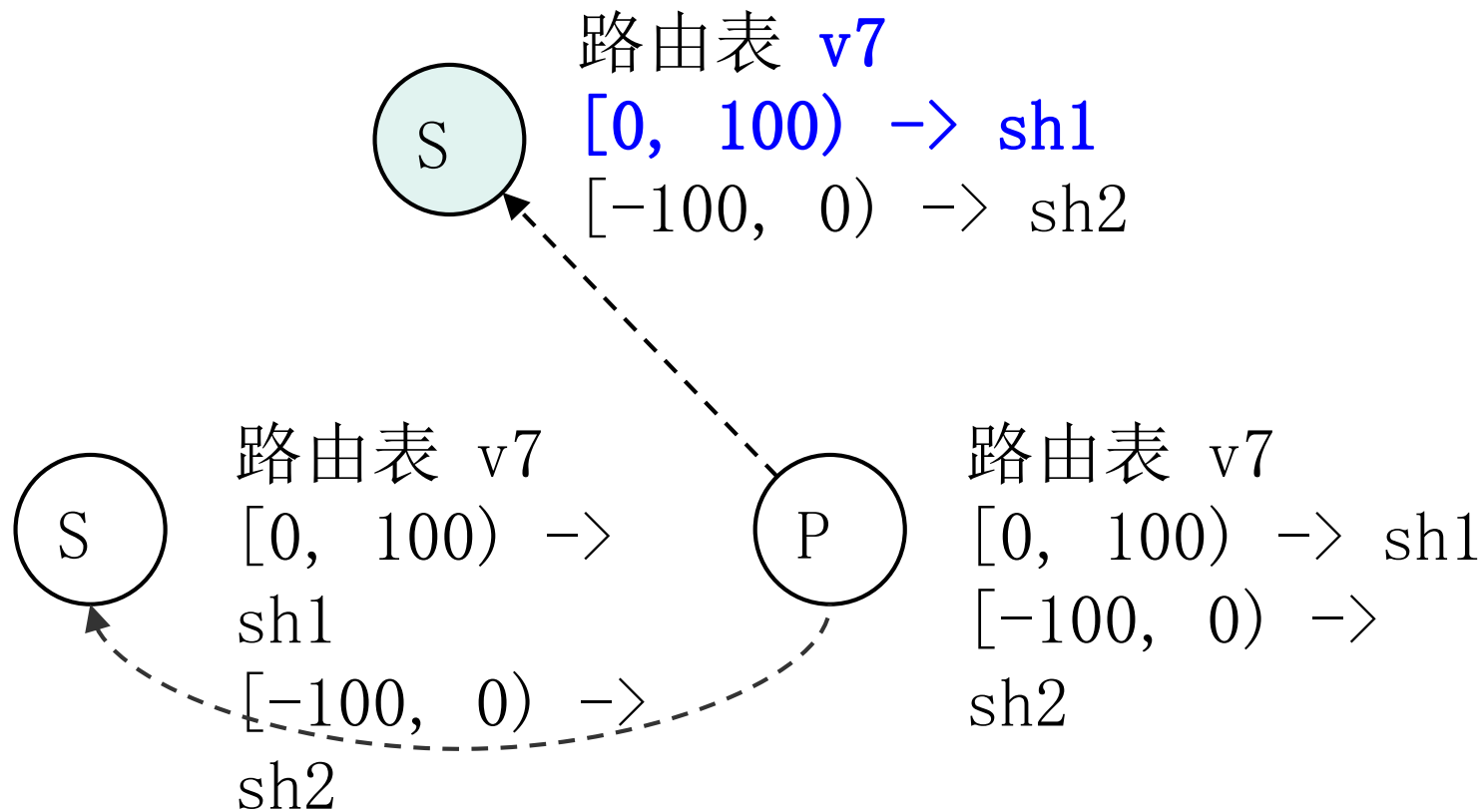
回滚案例



回滚案例



回滚案例



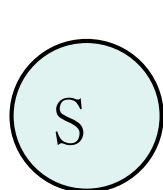
回滚案例



路由表 v7

$[0, 100) \rightarrow \text{sh0}$

$[-100, 0) \rightarrow \text{sh1}$



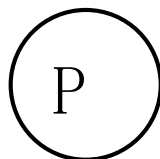
路由表 v7

$[0, 100) \rightarrow$

sh1

$[-100, 0) \rightarrow$

sh2

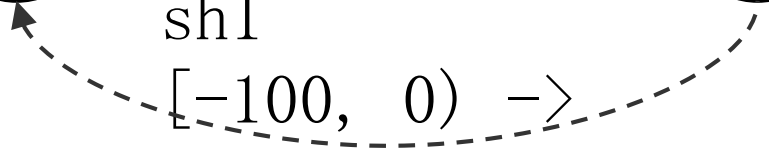


路由表 v7

$[0, 100) \rightarrow \text{sh1}$

$[-100, 0) \rightarrow$

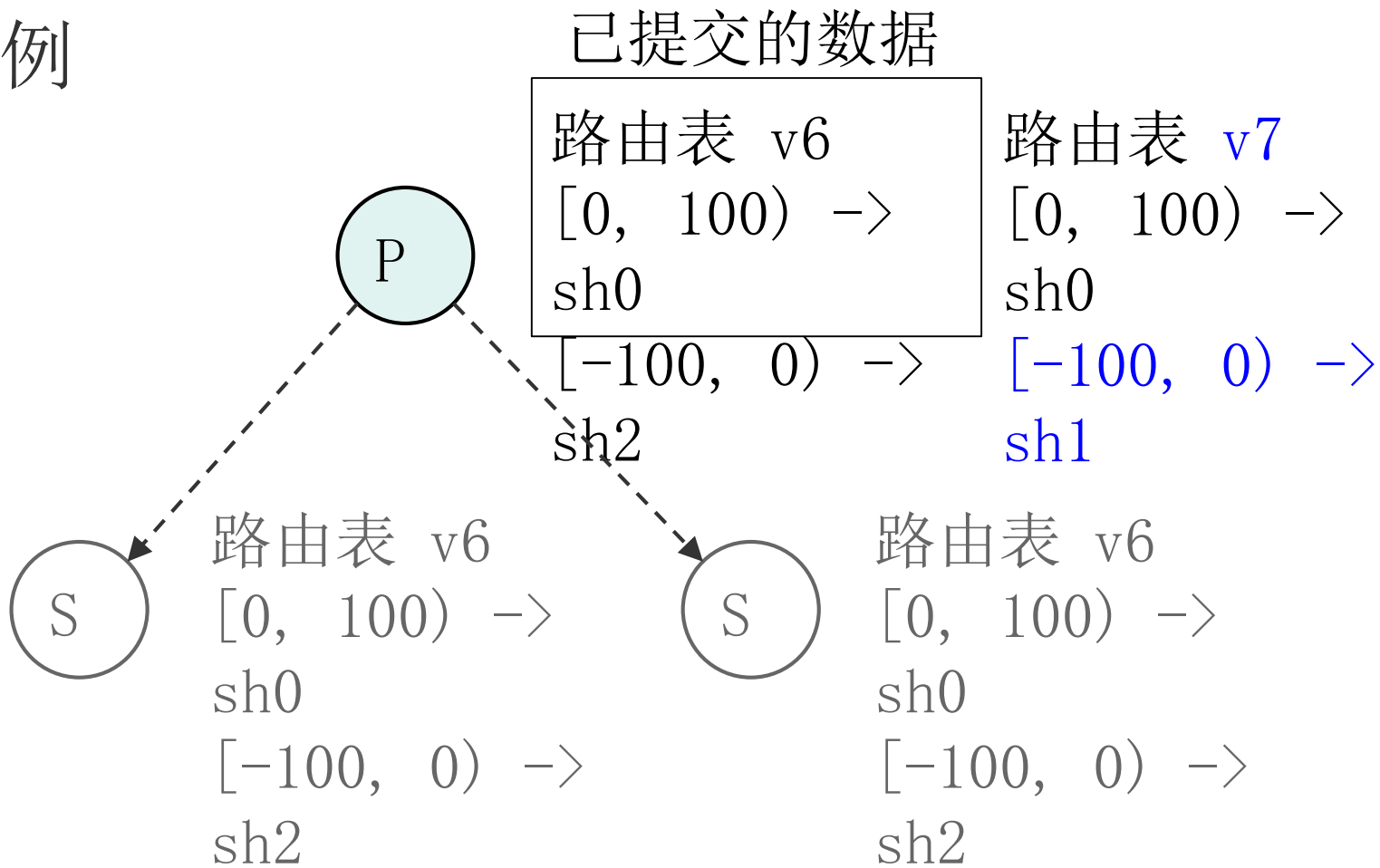
sh2



安全读级别：多数（新功能）

```
db.runCommand({
  find: 'foo',
  filter: { 金牌: 1 },
  readConcern: { level: 'majority' }
});
```


回滚案例



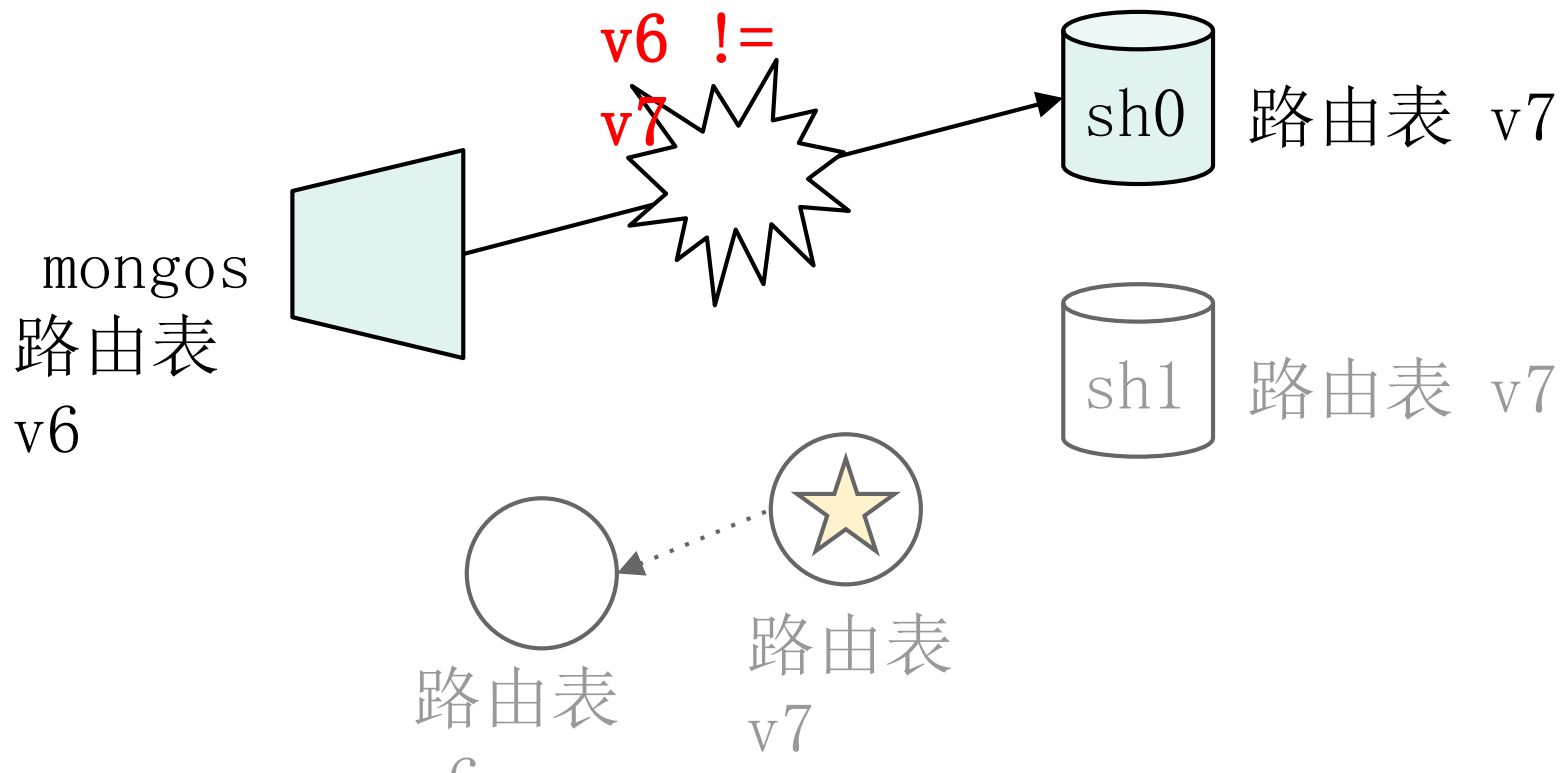
复制集配置服务器面临的挑战

- ~~已读取的数据可能回滚~~

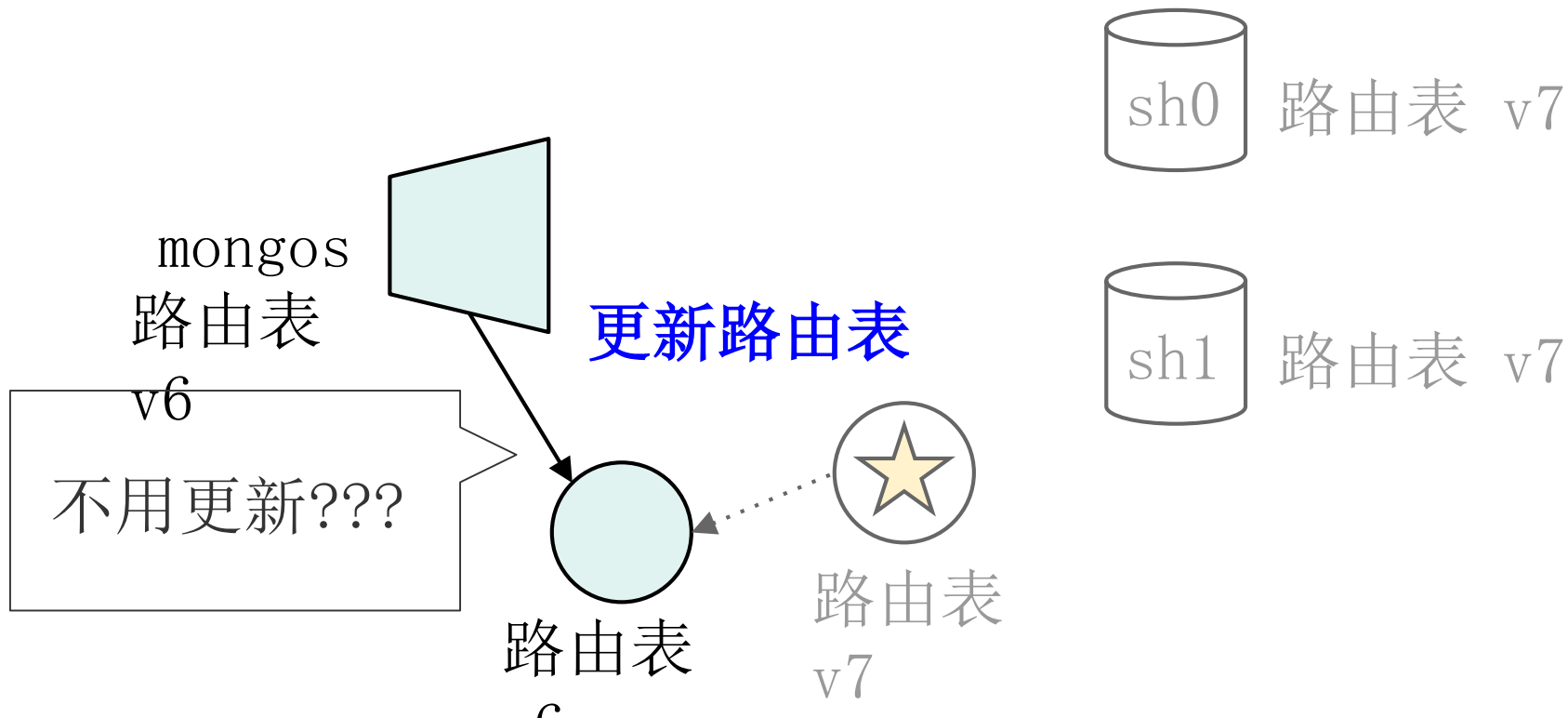
安全读级别：多数

- 从节点上的数据会过于陈旧

读取陈旧数据案例



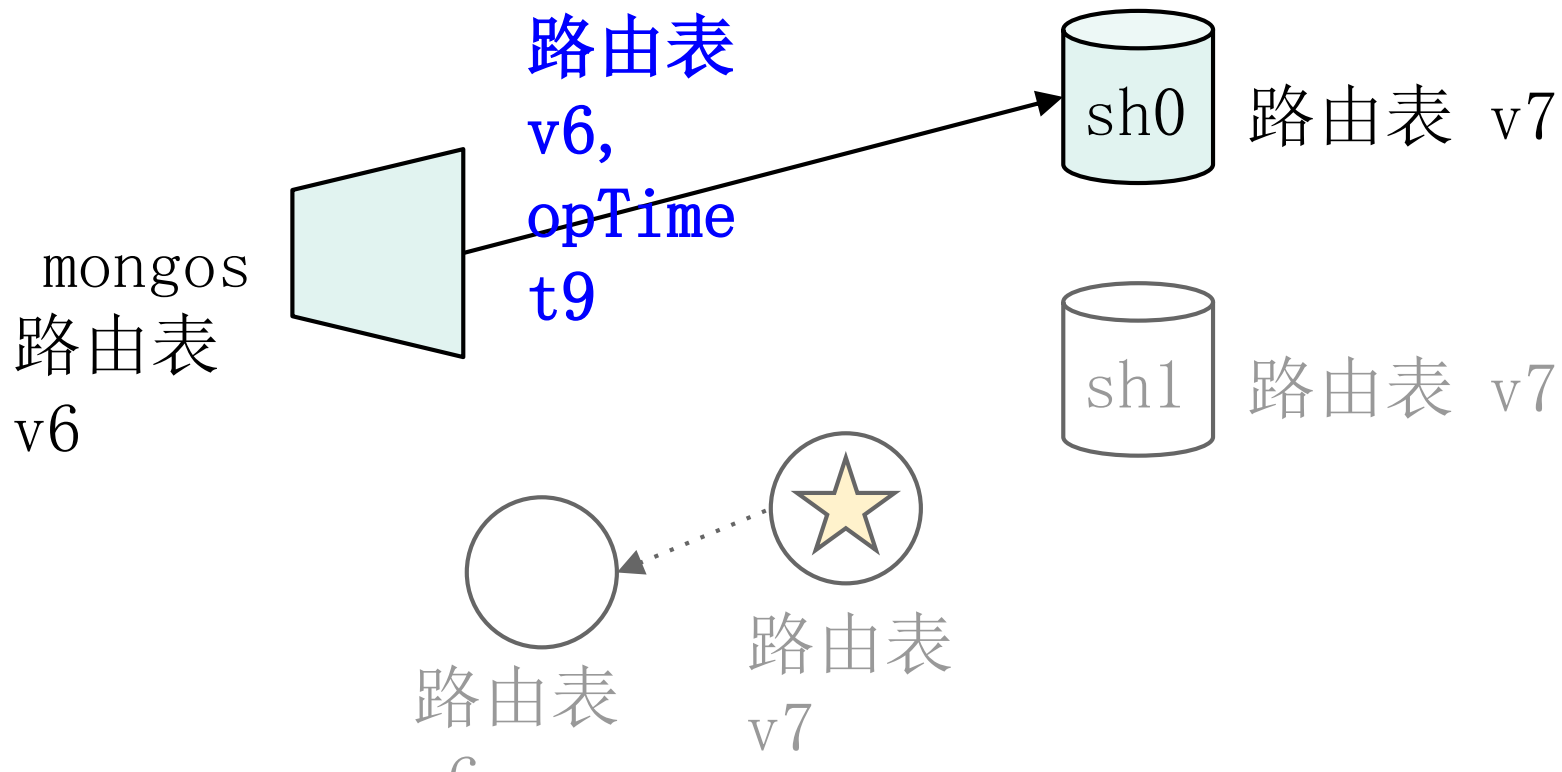
读取陈旧数据案例



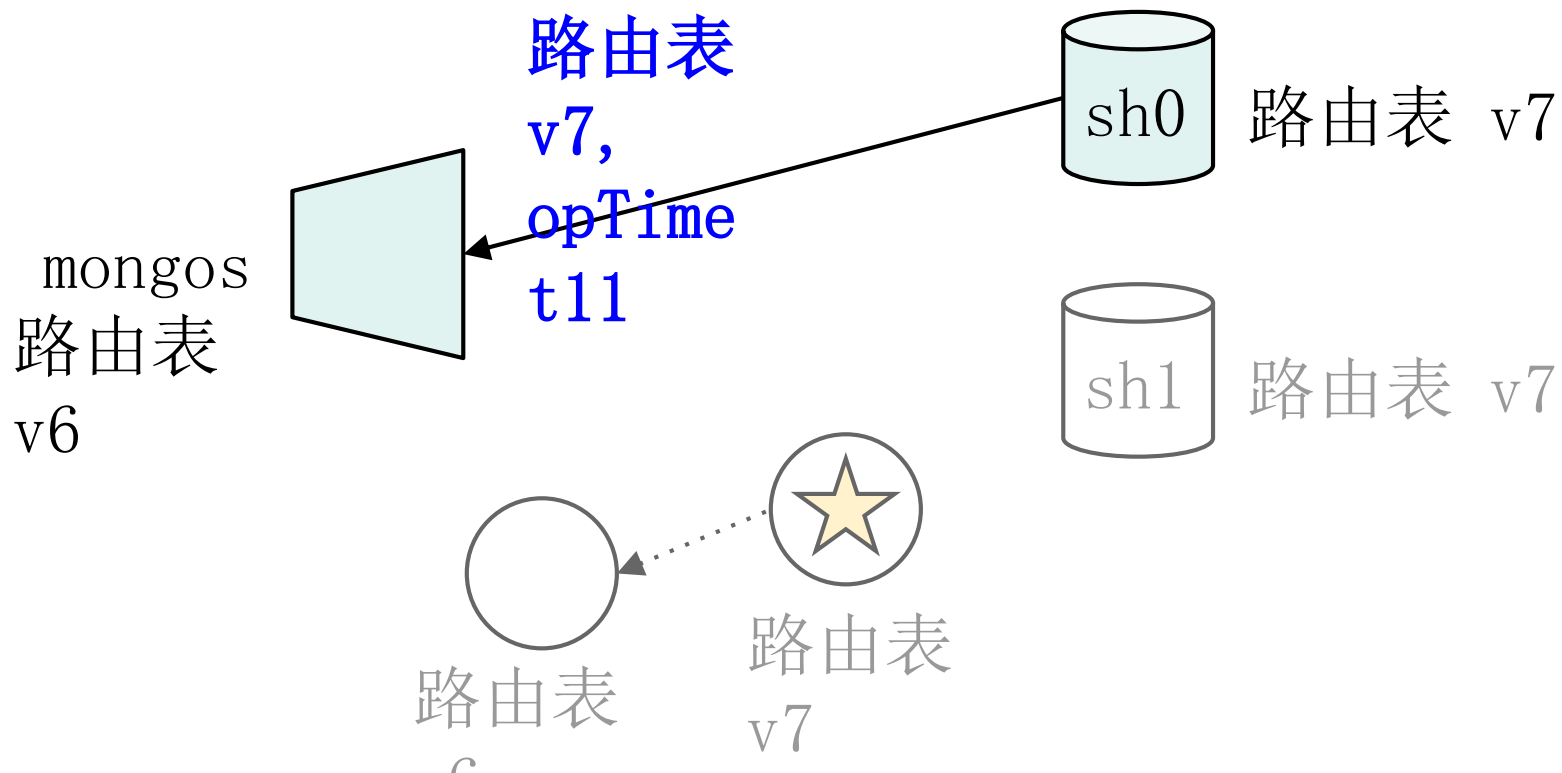
新功能（仅限内部使用）

```
db.runCommand({
  find: 'foo',
  filter: { 金牌: 1 },
  readConcern: {
    level: 'majority',
    afterOpTime: <opTime>
  })
```

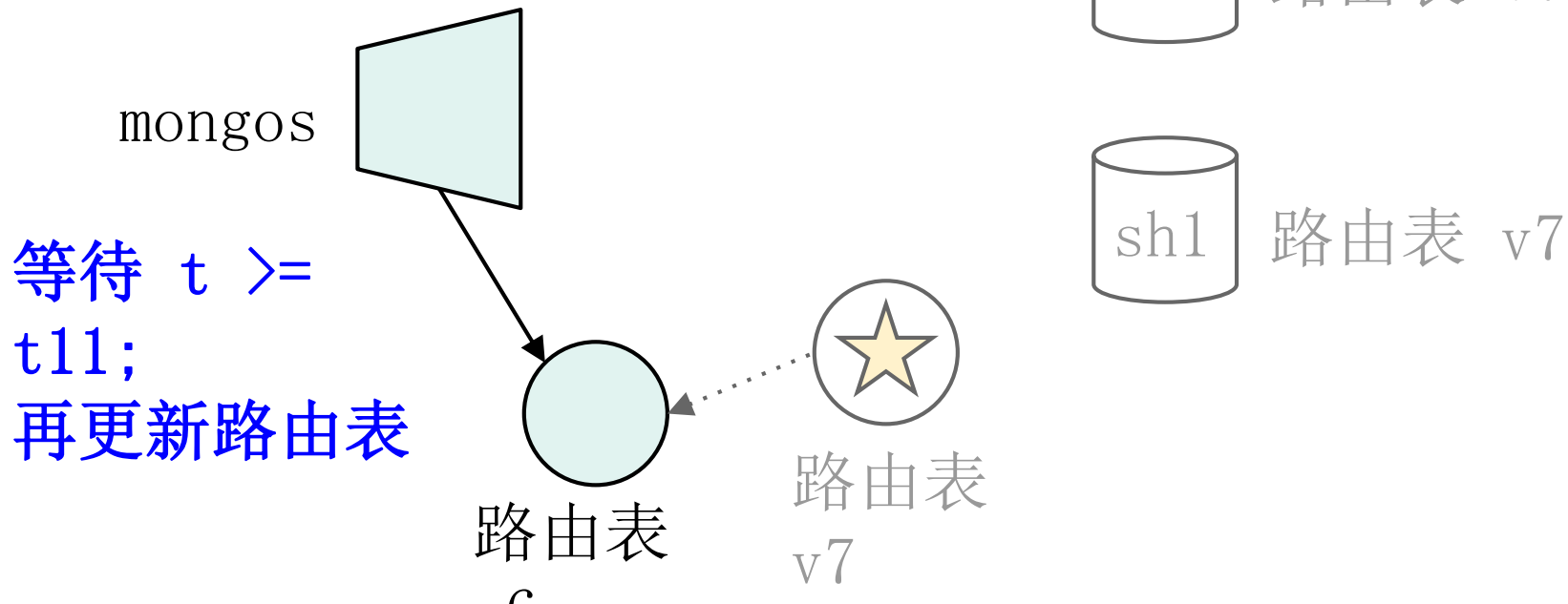
读取晚于OpTime



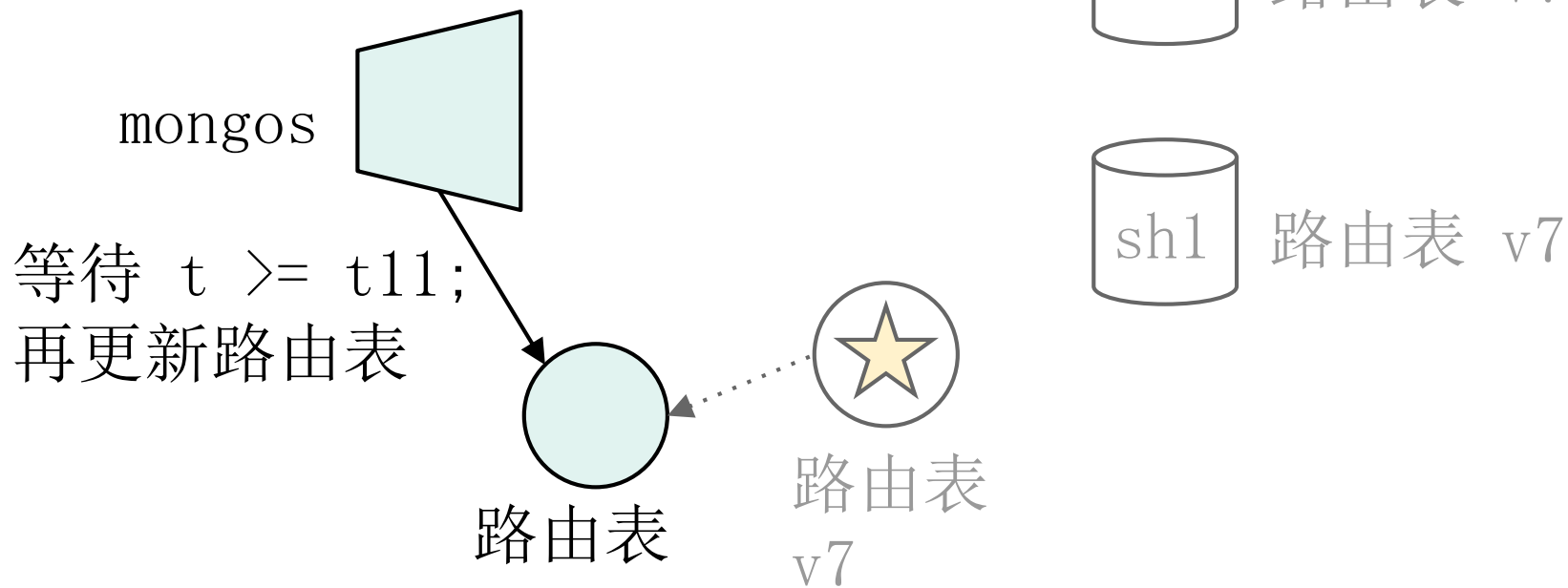
读取晚于OpTime



读取晚于OpTime



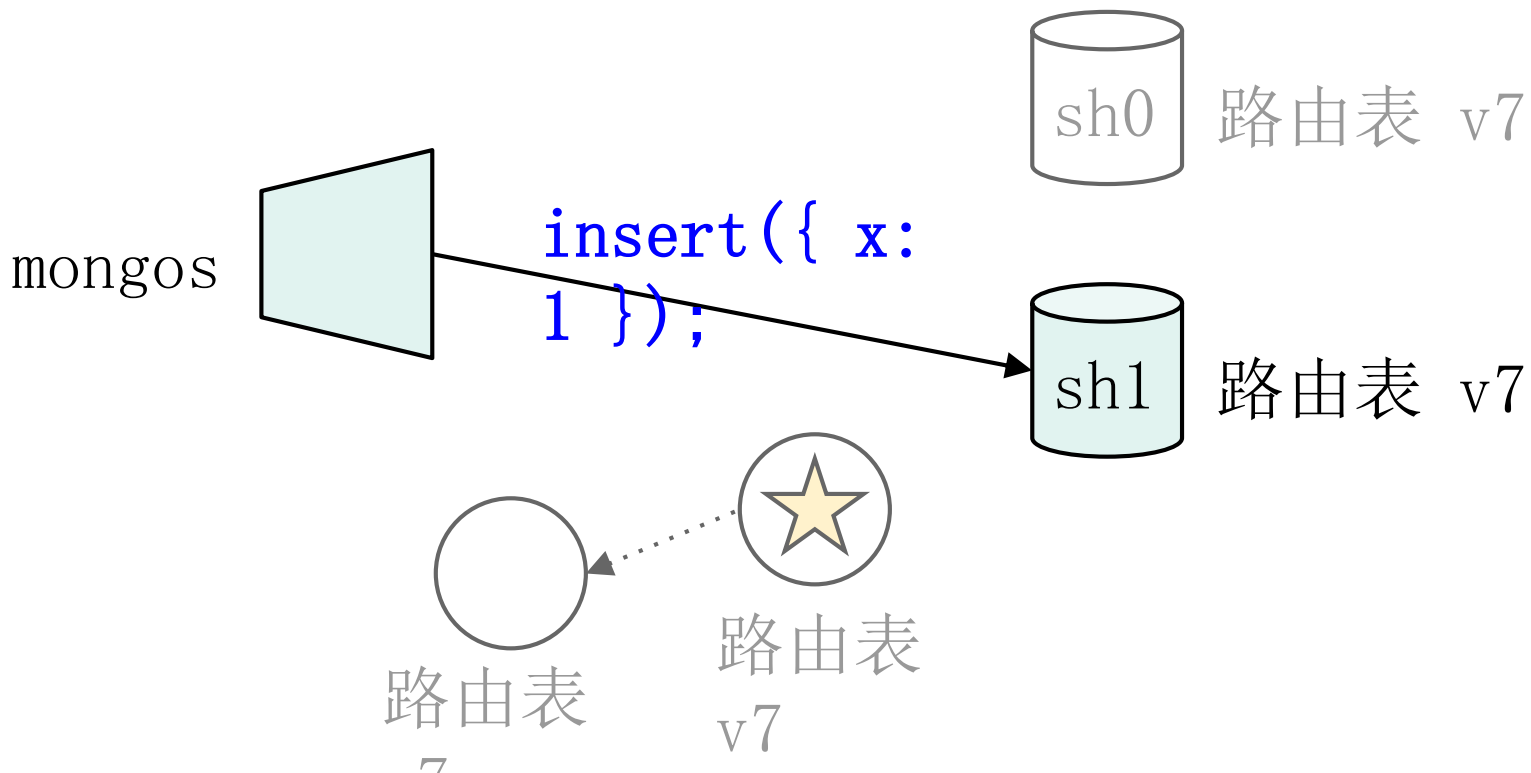
读取晚于OpTime



分片版本协议

版本	最小值	最大值	所在分片
1	MinKey	-200	分片0
2	-200	-100	分片2
3	-100	0	分片2
7	0	100	分片1
5	100	200	分片1
6	200	MaxKey	分片0

读取晚于OpTime



复制集配置服务器面临的挑战

- ~~已读取的数据可能回滚~~

安全读级别：多数

- 从节点上的数据会过于陈旧

读取晚于optime

MongoDB 3.2 中新功能

- 配置服务器可以使用复制集
 - V3.4中将强制使用复制集配置服务器
- 新功能：读取安全级别

展望未来

- 开放读取晚于optime

举手问题环节

