

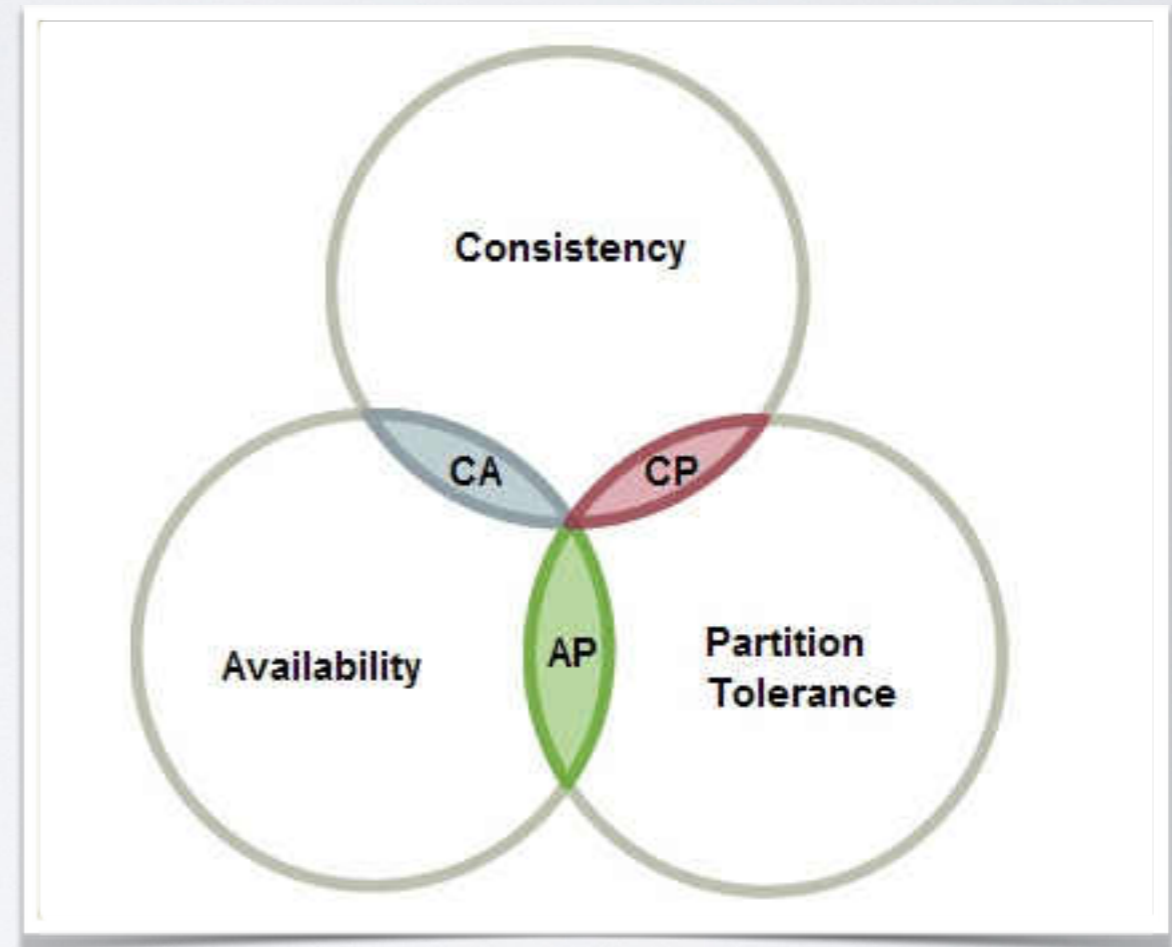


ZEPPELIN

High Available KV Storage Service

OVERVIEW

- CAP High Available
- Distributed KV storage



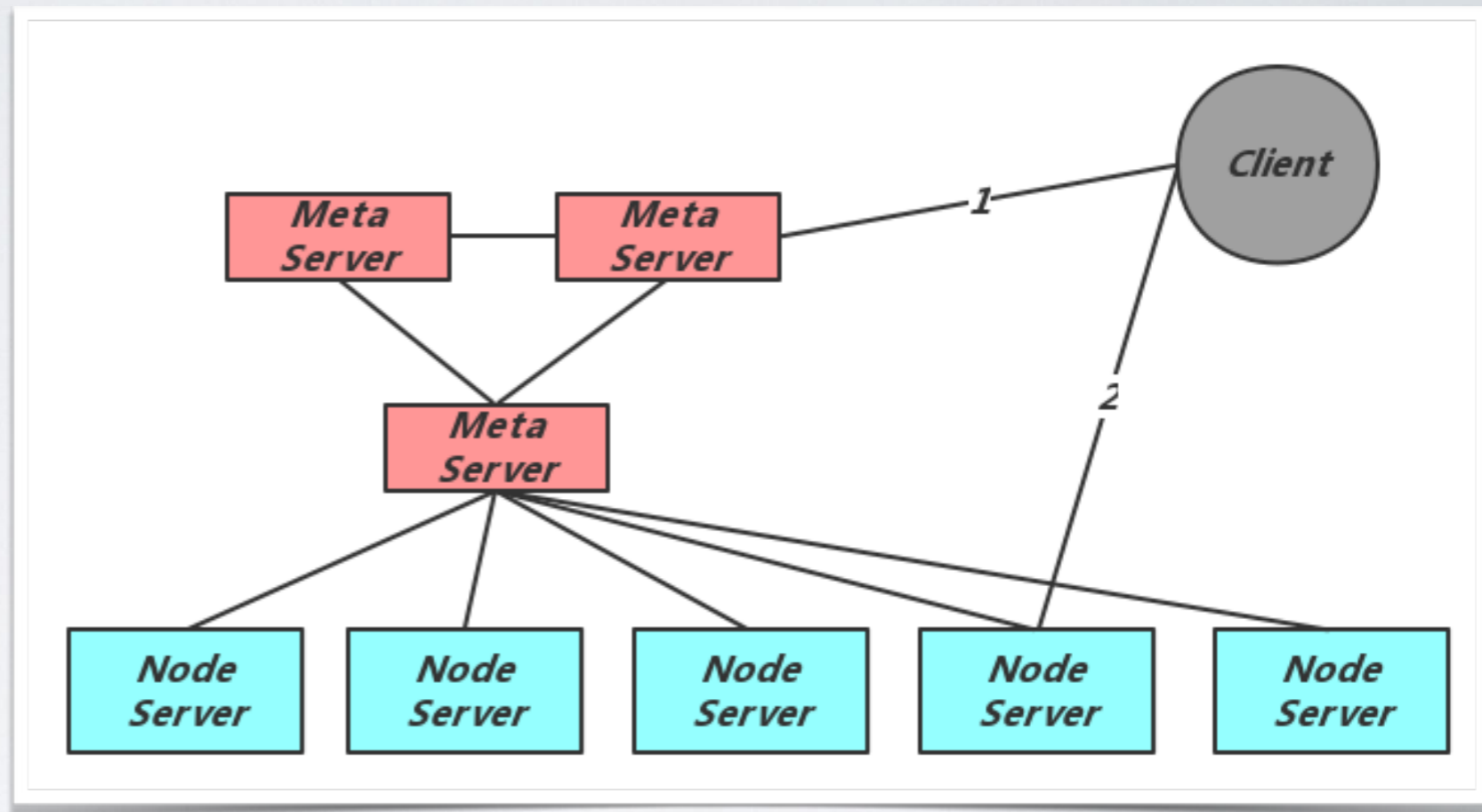
- Online Searching 600,000 QPS
 - 24 Physical Machine
 - 3 meta server 96 node server
 - 40+ tables
 - Highest Table Total query 500,000,000,000 times

OVERVIEW

- Interface Supported: SET, GET, DEL, MSET, MGET, INC
- TTL Supported
- Hashtag Supported

OVERVIEW

- Client pull meta Info
- Calculate partition
- Find node ip in meta info
- Send request to corresponding node server



NODE SERVER

- Data Distribution & Replication
- Thread Model
- Synchronization

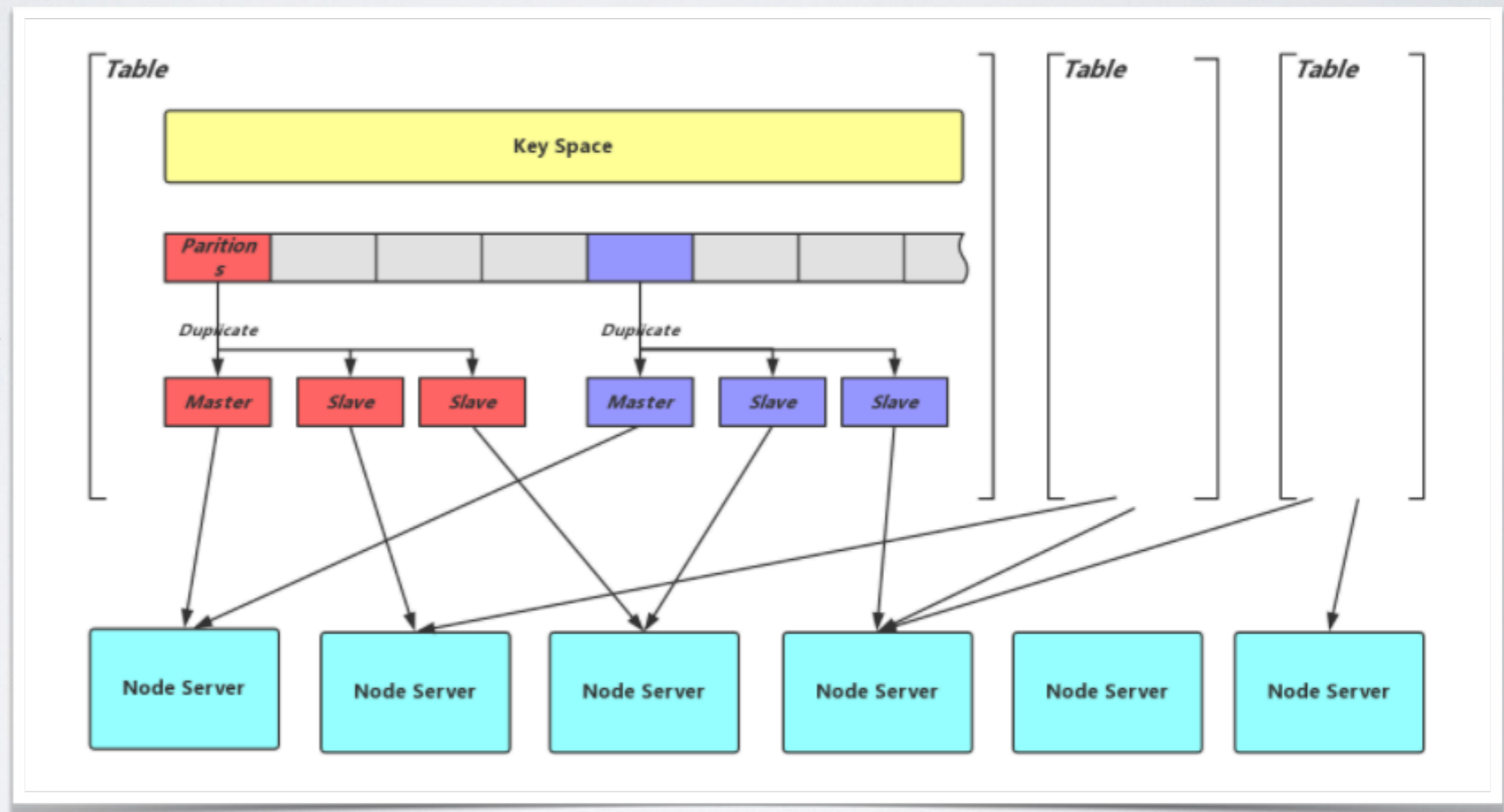
NODE SERVER

- Data Distribution & Replication

NODE SERVER

Data Distribution & Replication

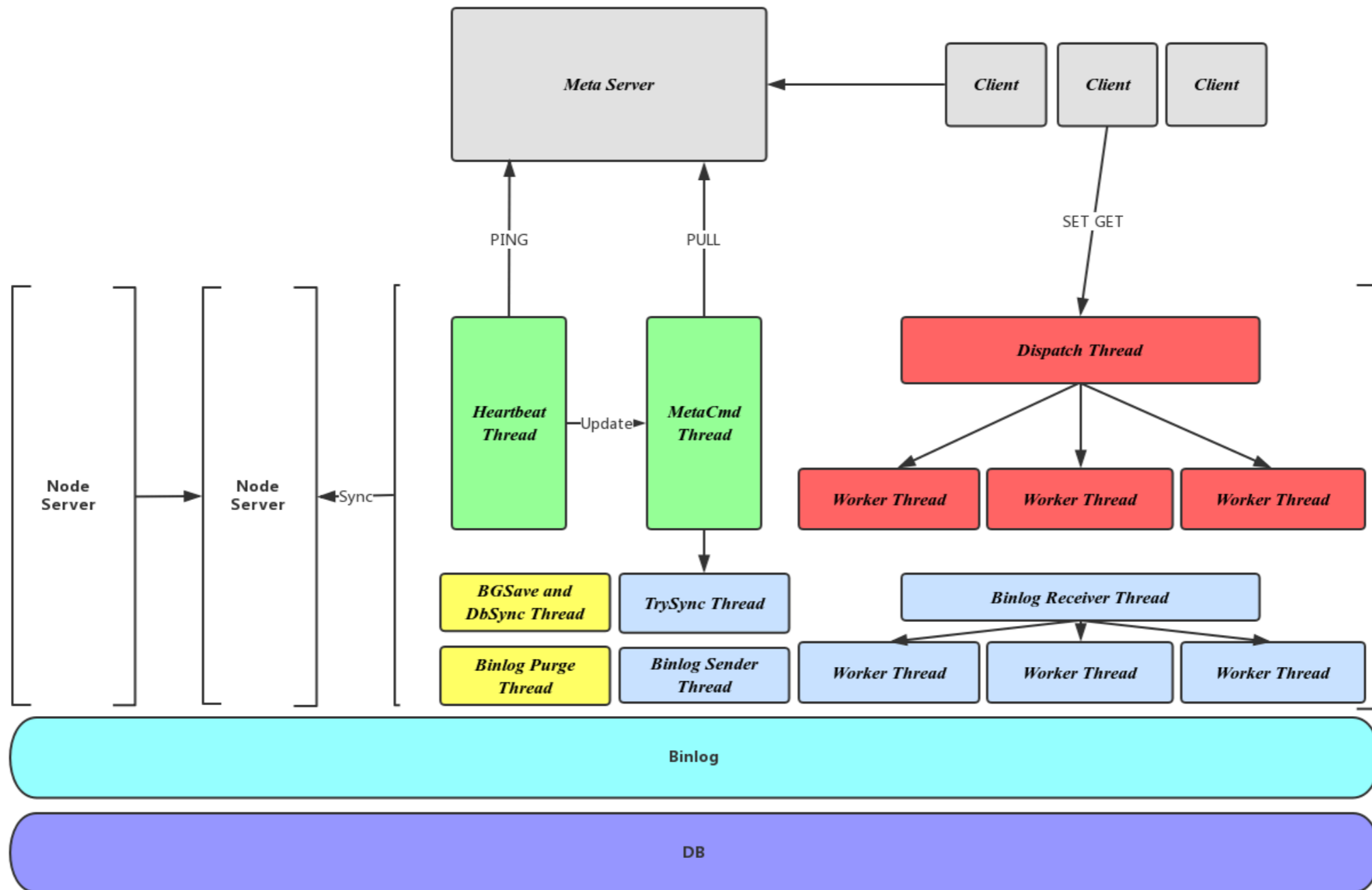
- Partition
- Master
- Slave



NODE SERVER

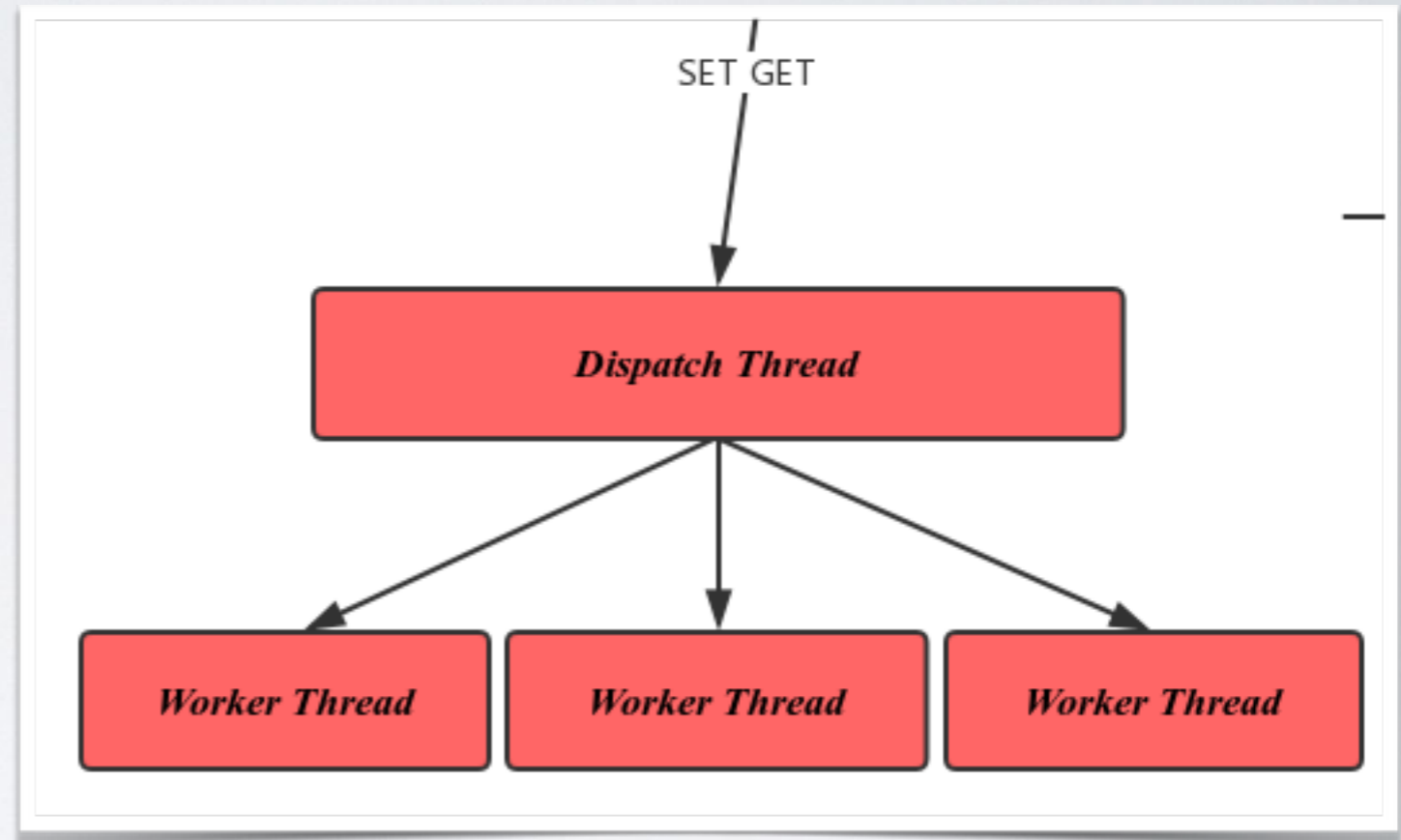
- Thread Model

NODE SERVER

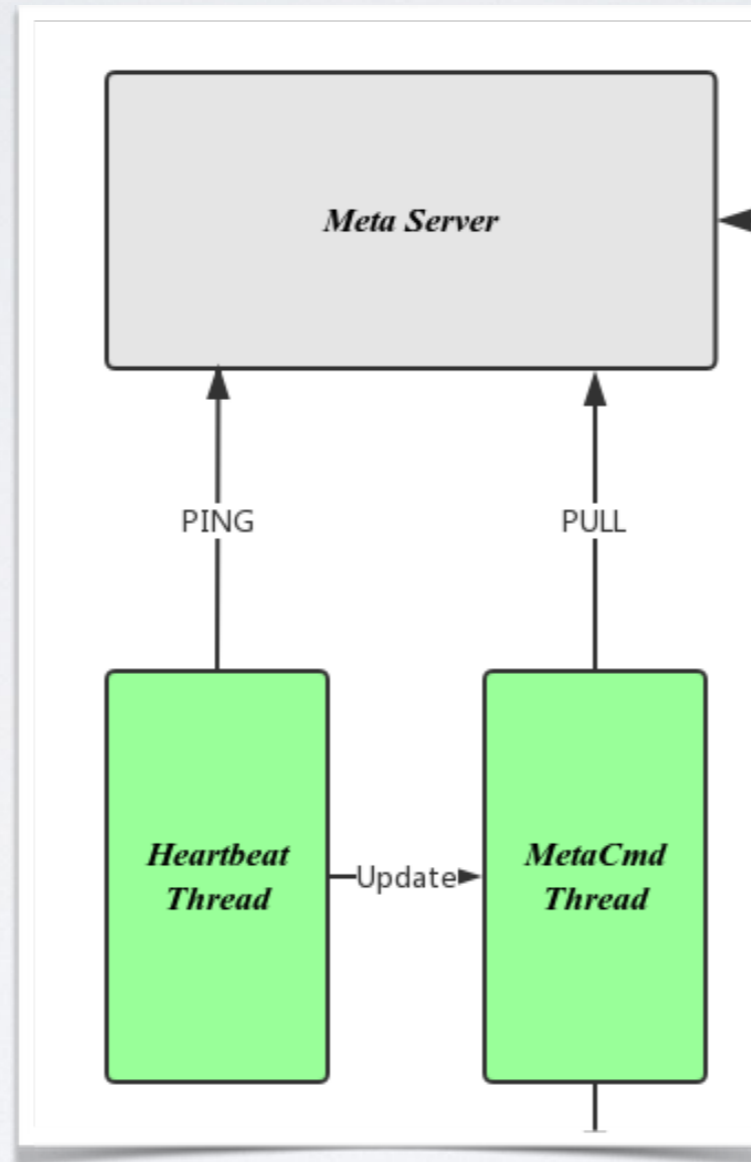


NODE SERVER

```
zp_data_partiton.cc
if (!cmd->is_suspend()) {
    // read lock
    pthread_rwlock_rdlock(&suspend_rw_);
}
if (cmd->is_write) {
    // lock this key
    mutex_record_.Lock(key);
}
cmd->Do(req, &res);
if (cmd->is_write) {
    if (res->code() == client::StatusCode::kOk) {
        logger_->Put(raw);
    }
    mutex_record_.Unlock(key);
}
if (!cmd->is_suspend()) {
    pthread_rwlock_unlock(&suspend_rw_);
}
```



NODE SERVER

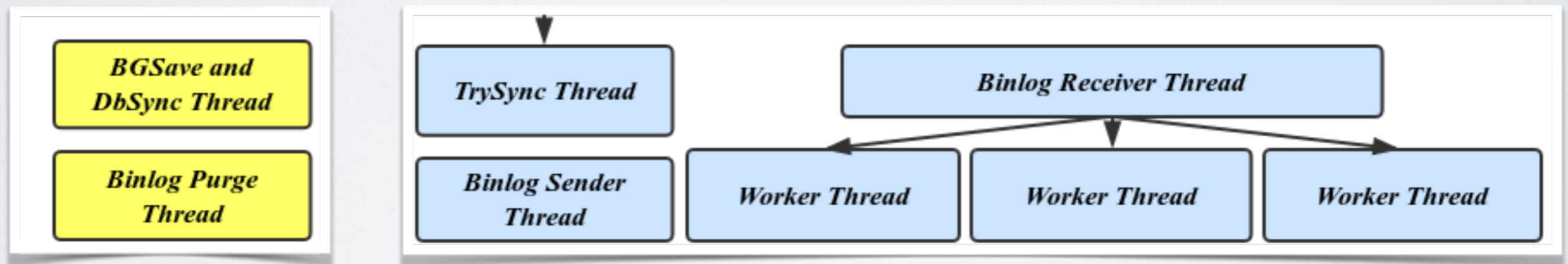


NODE SERVER

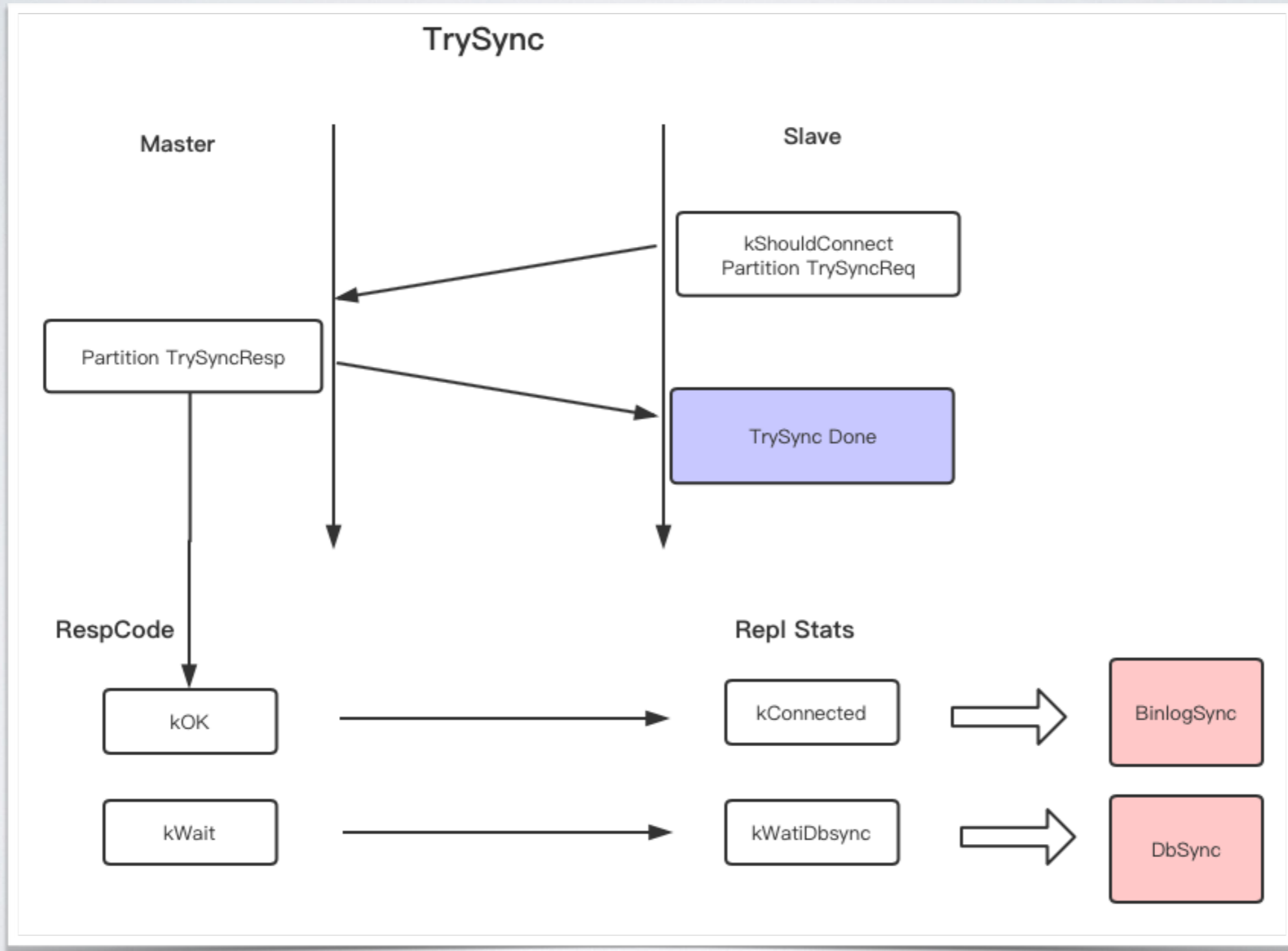
- Synchronization

NODE SERVER

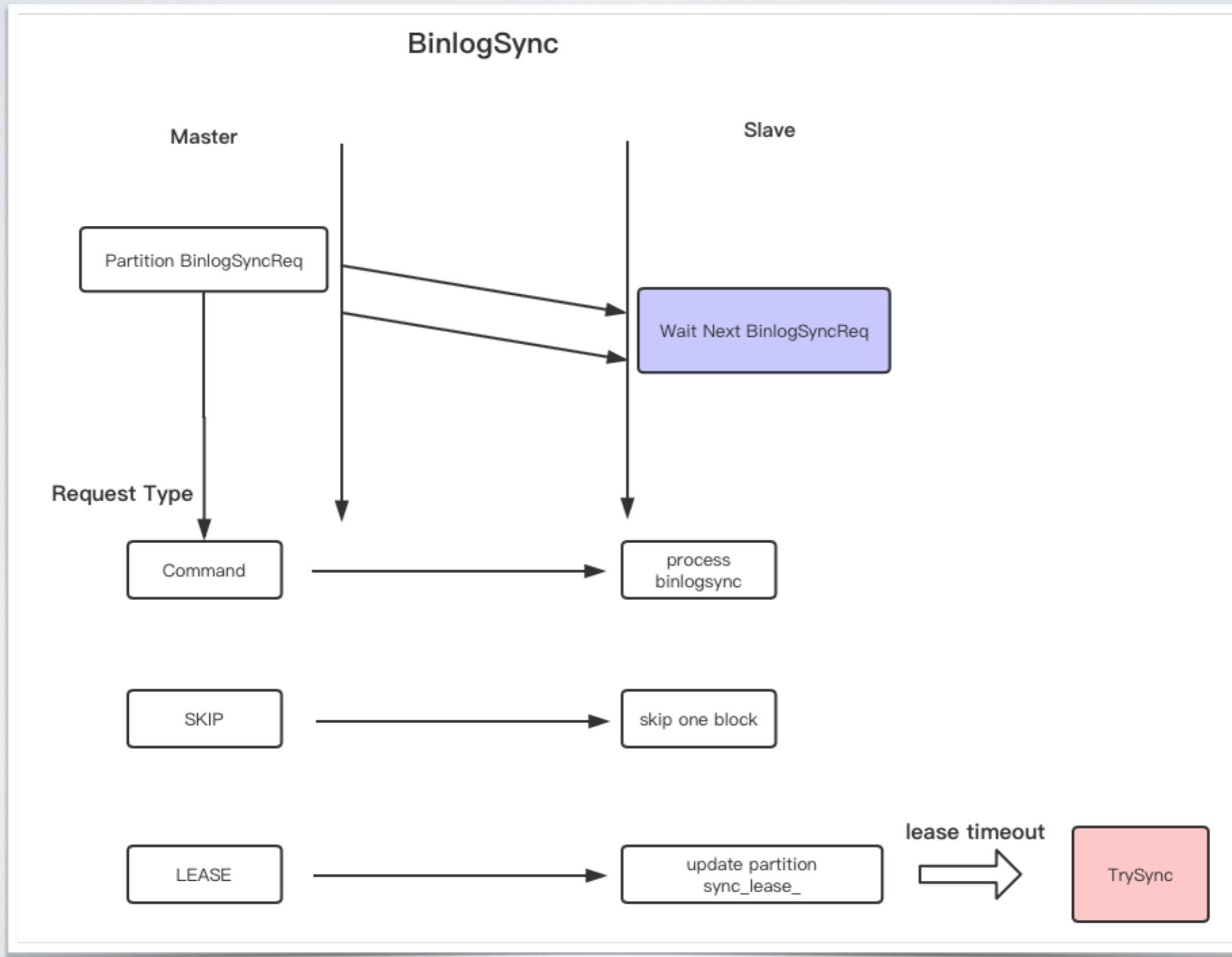
- Binlog
- DBSync & Binlog Sync



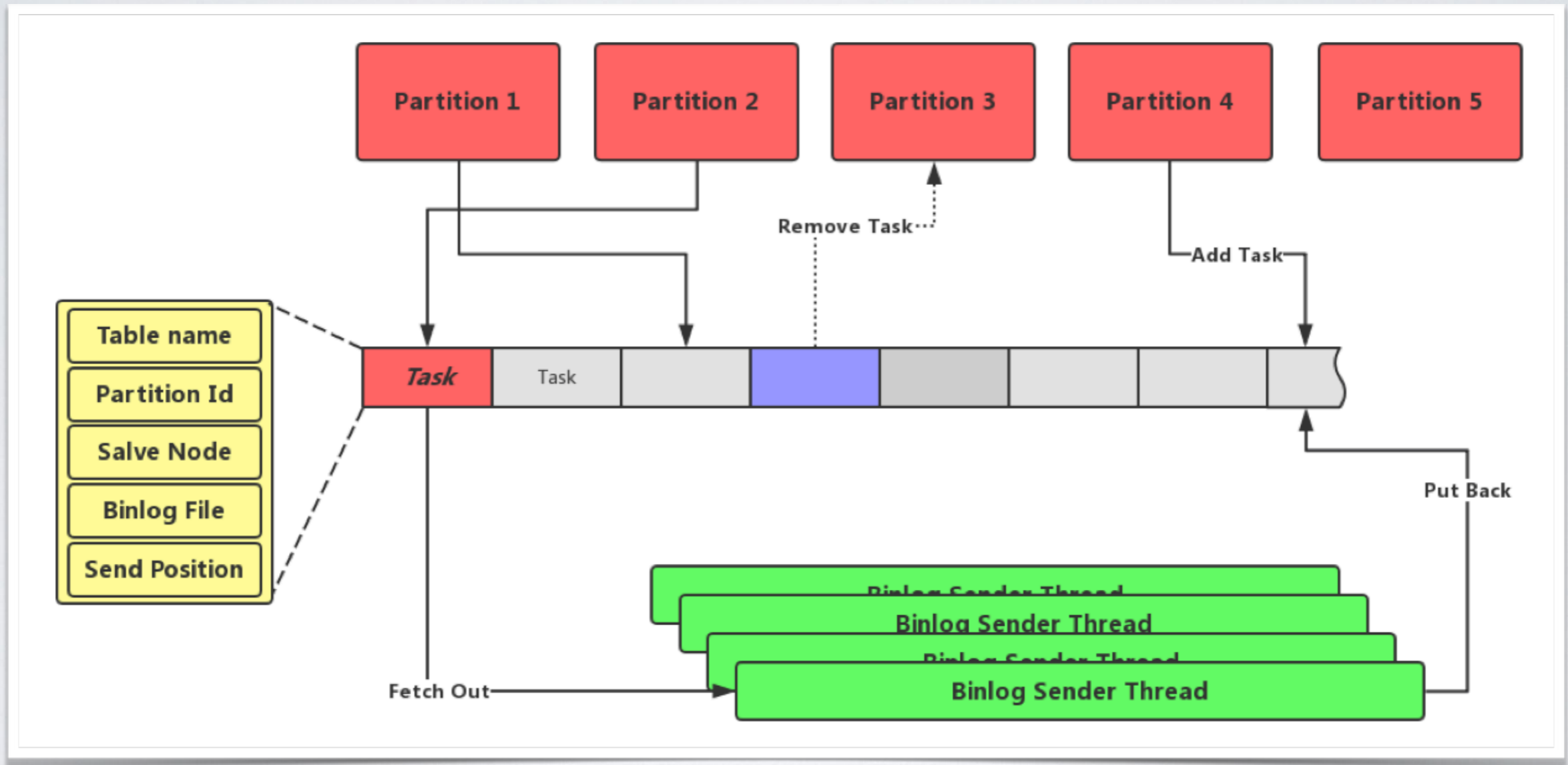
NODE SERVER



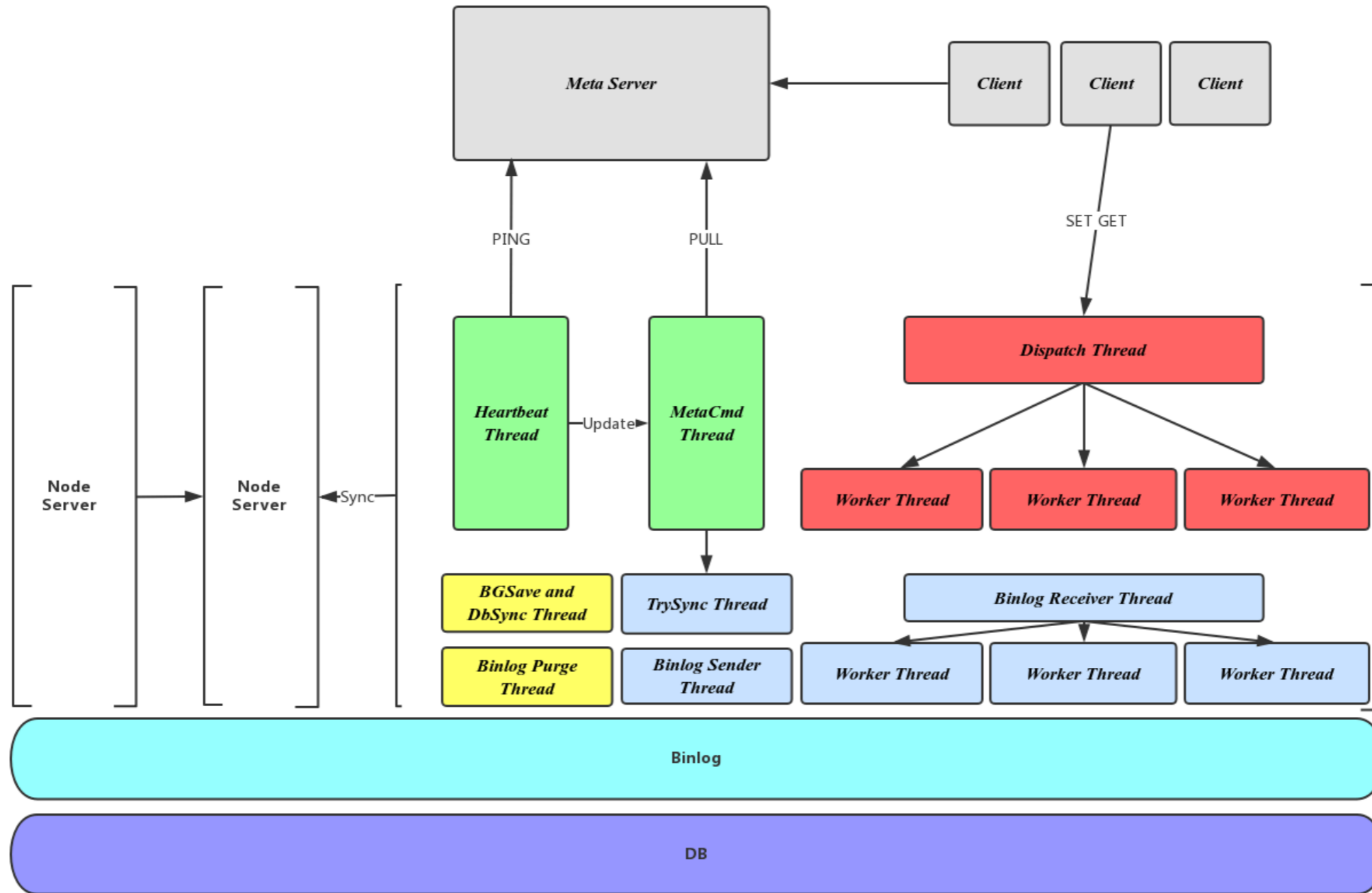
NODE SERVER



NODE SERVER



NODE SERVER



META SERVER

- Meta Info
- Thread Model
- Cluster Management(Migrate)

META SERVER

- Meta Info

META SERVER

- Meta Info
- Cluster Status
- Cluster Topology

```
class ZPMetaInfoStore {
    // -2 for uninitialed
    // -1 for initialed but no table
    // Otherwise non-negative integer and monotone increasing
    // epoch records topology info changes
    std::atomic<int> epoch_;
    // table => ZPMeta::Table set
    std::unordered_map<std::string, ZPMeta::Table> table_info_;
    // node => alive time + offset set, 0 means already down node
    // only valid for leader
    std::unordered_map<std::string, NodeInfo> node_infos_;
}

message Partitions {
    required int32 id = 1;
    required PState state = 2;
    required Node master = 3;
    repeated Node slaves = 4;
}

ZPMeta::Table
message Table {
    required string name = 1;
    repeated Partitions partitions = 2;
}

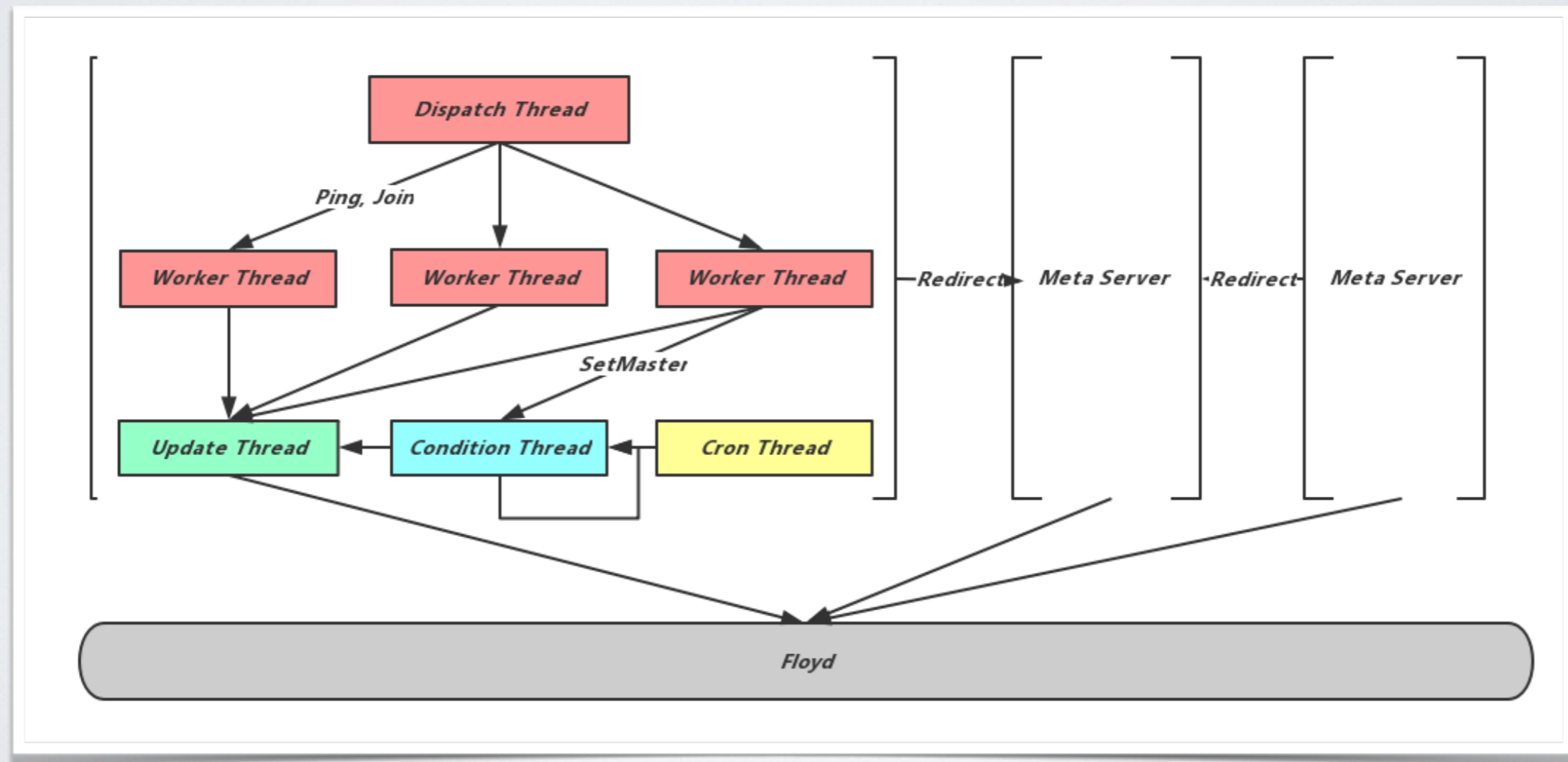
struct NodeInfo {
    uint64_t last_active_time;
    // table_partition -> offset
    std::map<std::string, NodeOffset> offsets;
}
```

META SERVER

- Thread Model

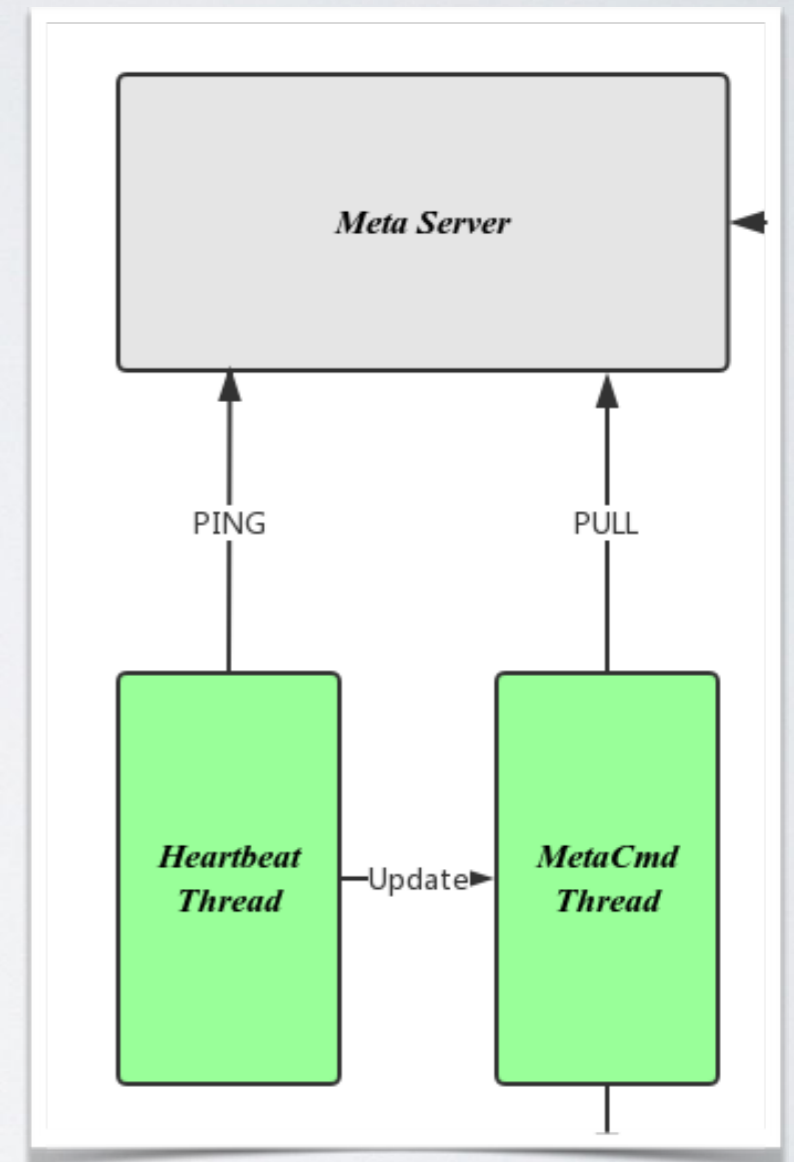
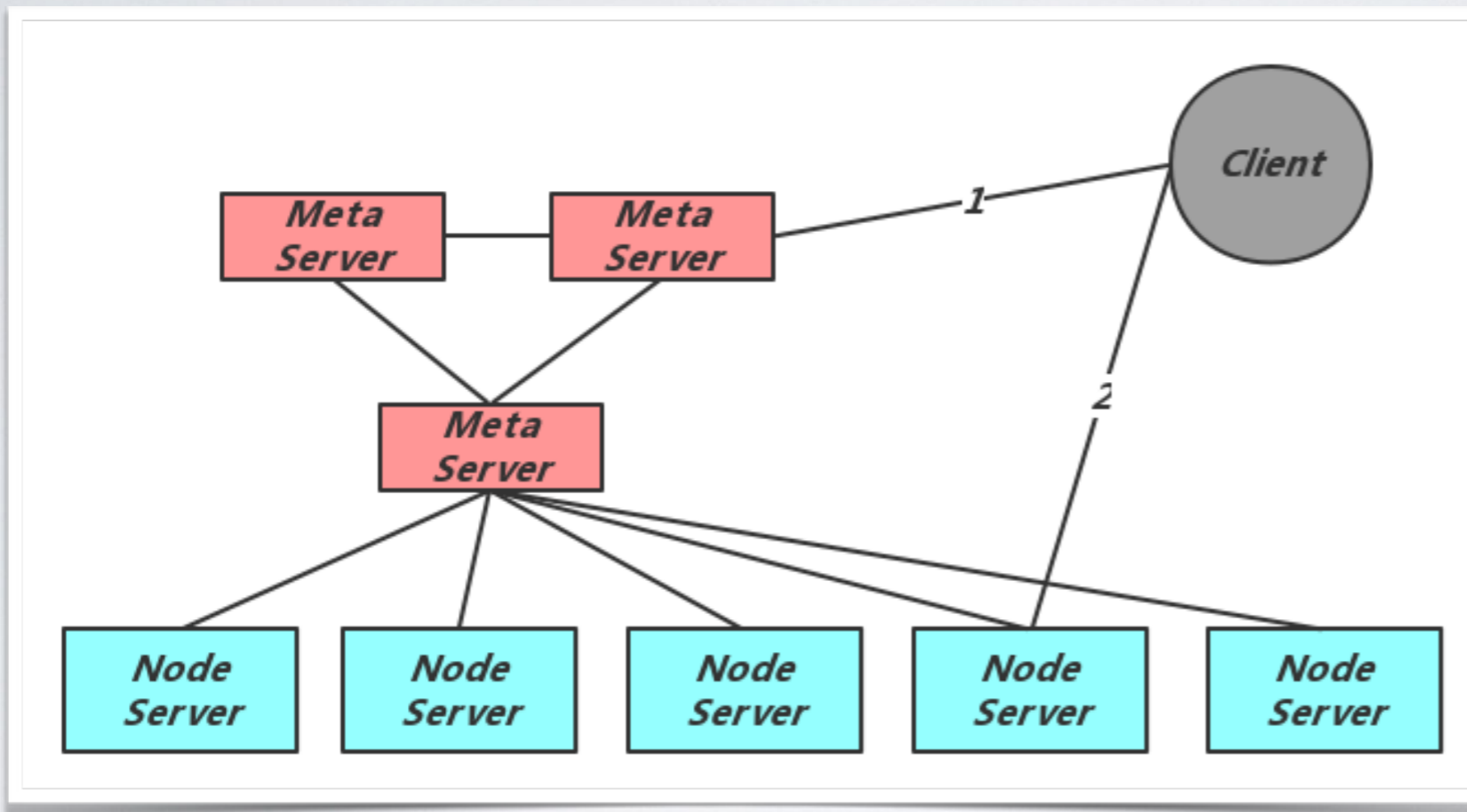
META SERVER

- Thread Model



META SERVER

- How does meta Info apply to node server?



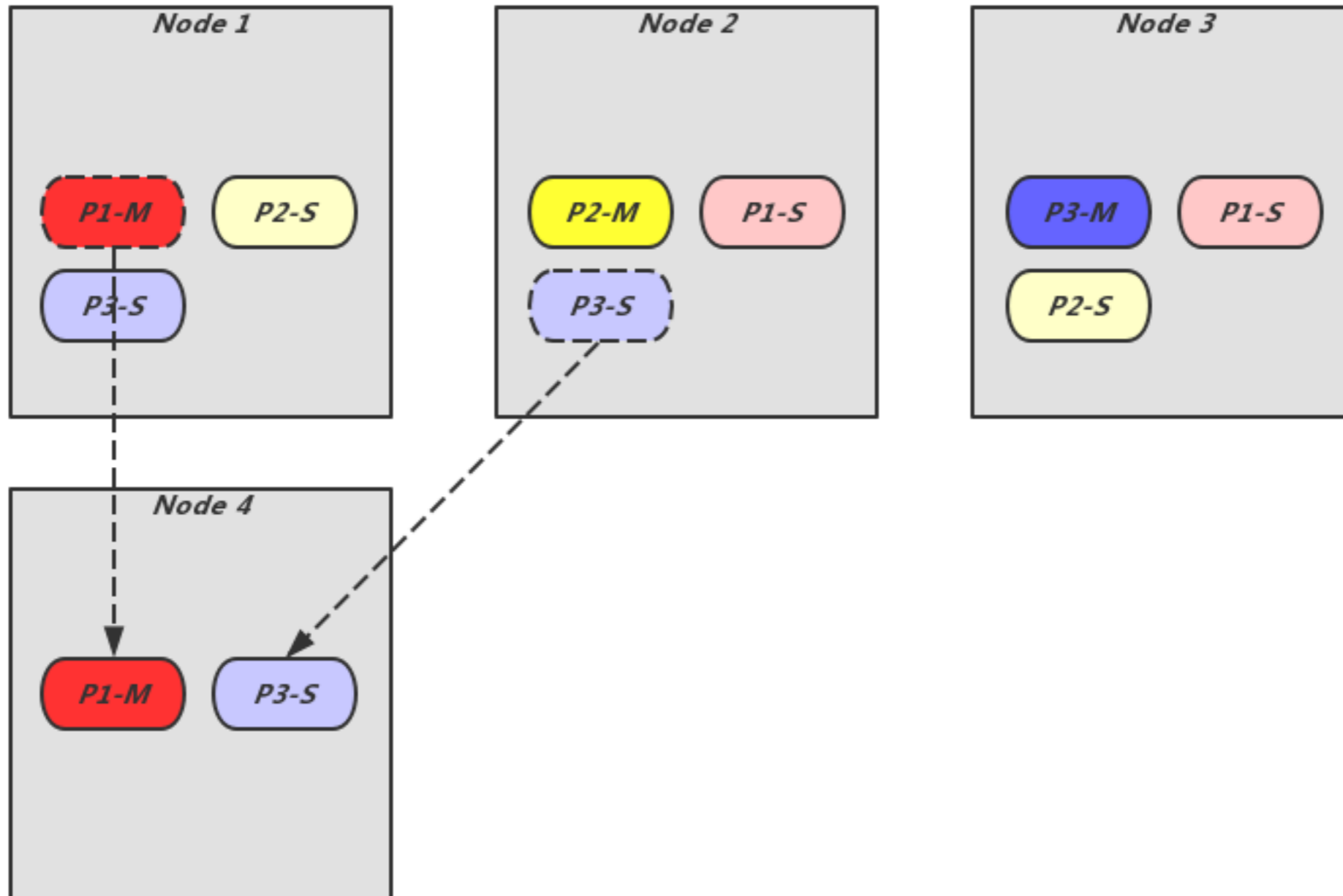
- **1. Client connect to leader Meta Server and modify meta info.**
- **2. As raft exist, follower and leader could reach a consensus about meta info**
- **3. Node server do ping routine. Find epoch differ from Meta Server. Pull newest Meta Info and apply this modification locally.**

META SERVER

-

Migrate

META SERVER



META SERVER

- Original Meta

```
Table_test:  
Partition4  
Status: Active  
Master 1.1.1.1:9221  
Slave: 2.2.2.2:9221
```

- Send “migreate table_test 1.1.1.1:9221 4 3.3.3.3:9221” to meta server

update_thread

```
Table_test:  
Partition4  
Status: SlowDown  
Master 1.1.1.1:9221  
Slave: 2.2.2.2:9221, 3.3.3.3:9221
```

condition_thread

task:
condition: Partition 4 BinlogSync Same file
actiion: Forward Partition 4 sutck to
update_thread

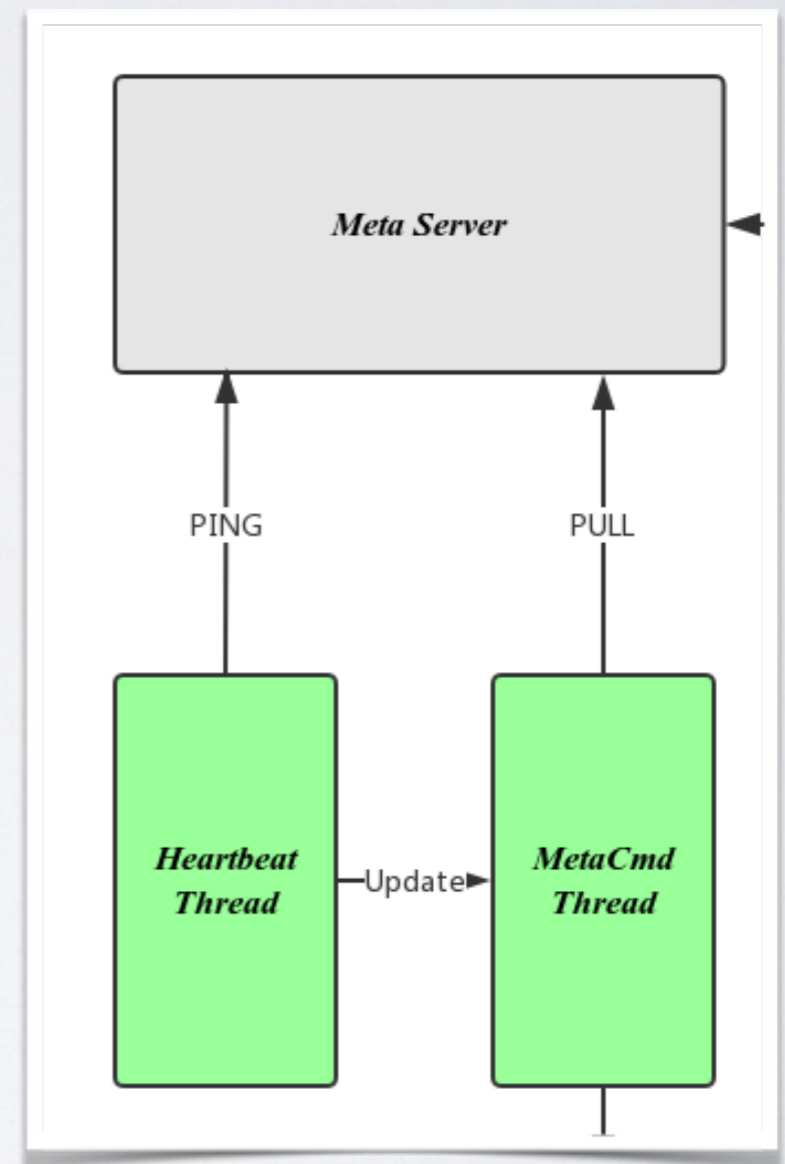
task:
condition: Partition 4 BinlogSync Offset same
actiion: Forward Partition 4 handover & active
to update_thread

META SERVER

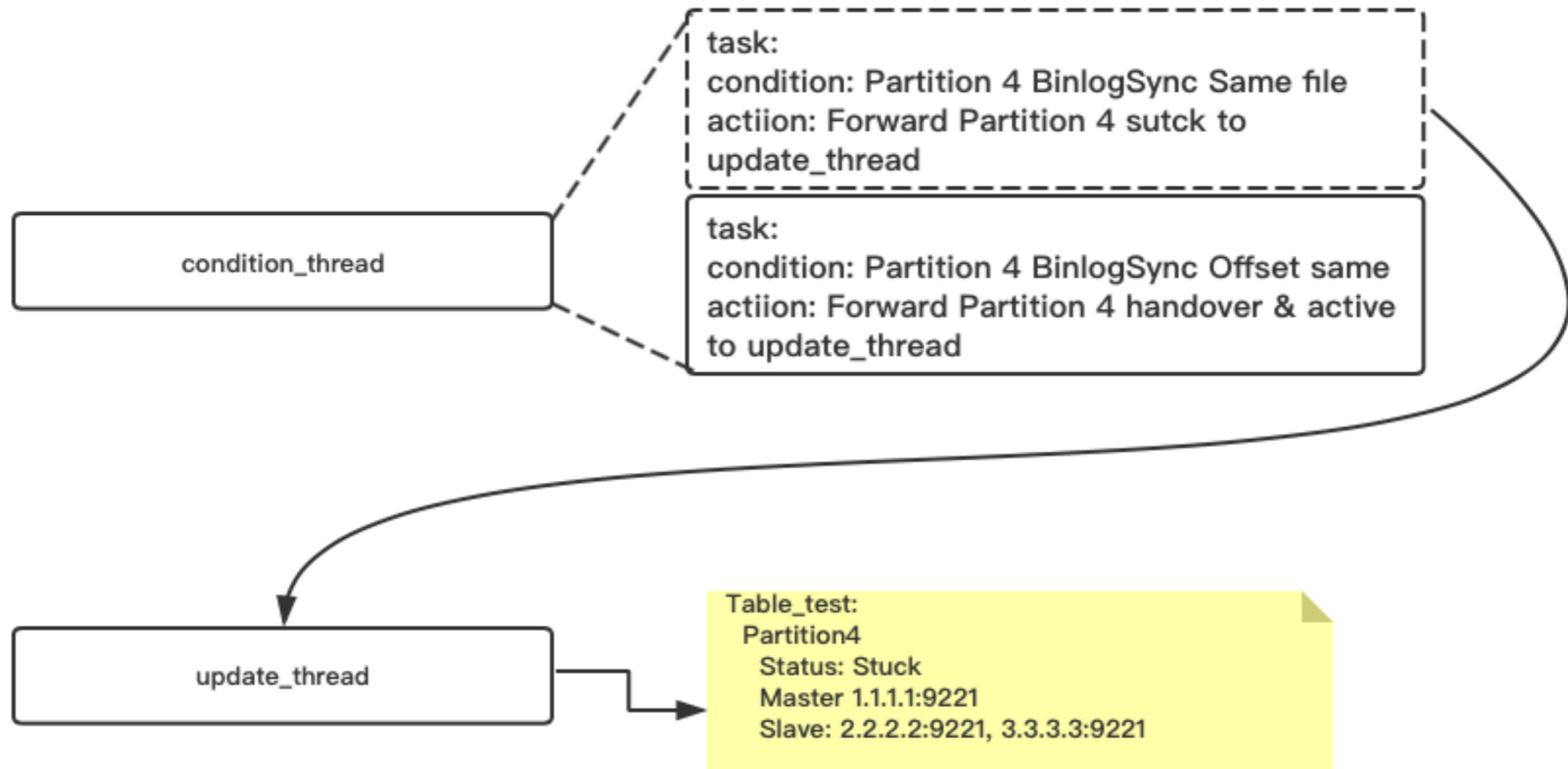
- node server apply meta info changes locally

```
Table_test:  
Partition4  
Status: SlowDown  
Master 1.1.1.1:9221  
Slave: 2.2.2.2:9221, 3.3.3.3:9221
```

- continue ping...

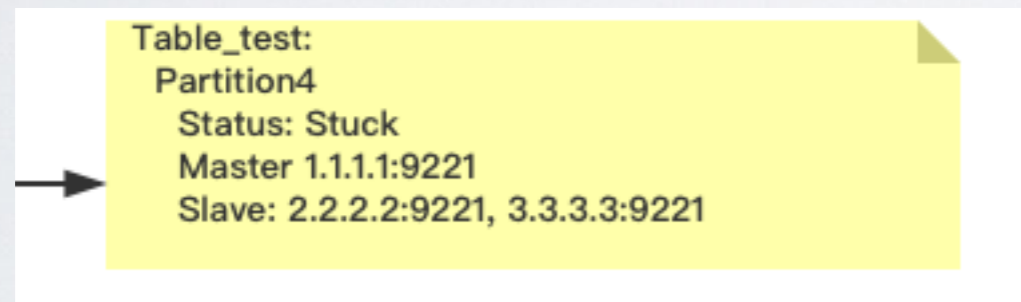


META SERVER



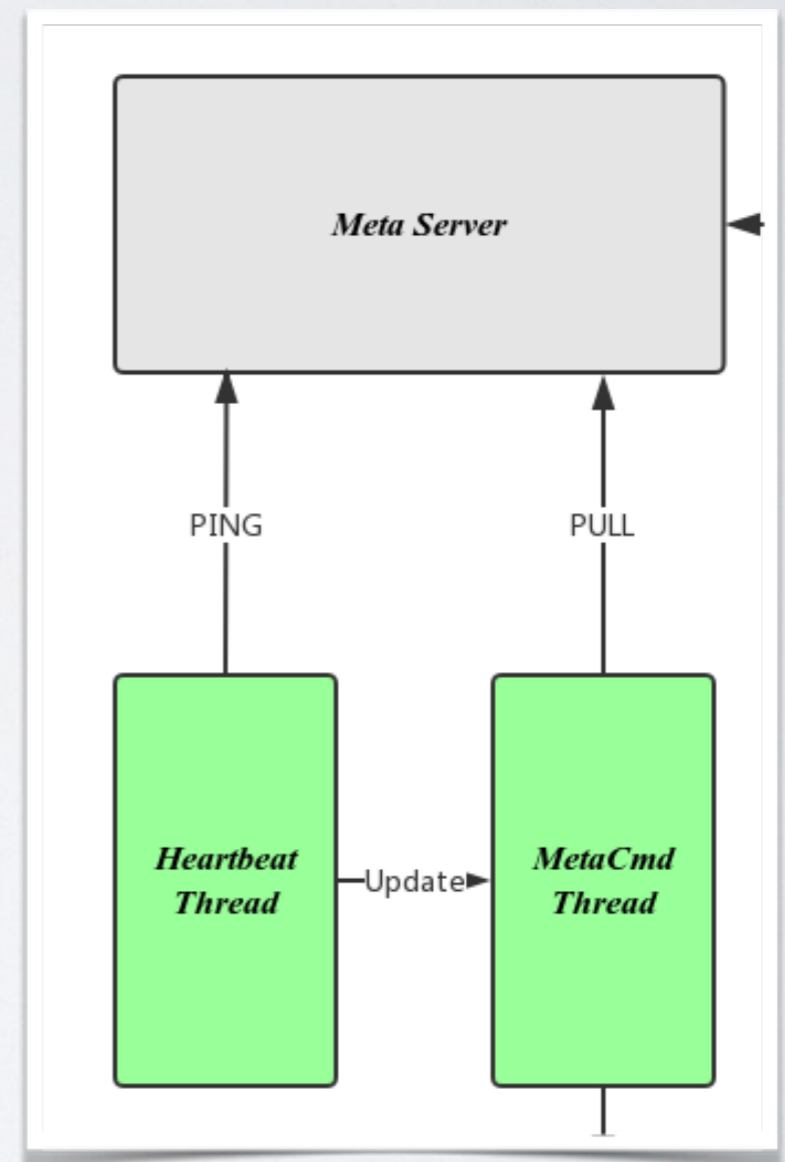
META SERVER

- node server apply meta info changes locally

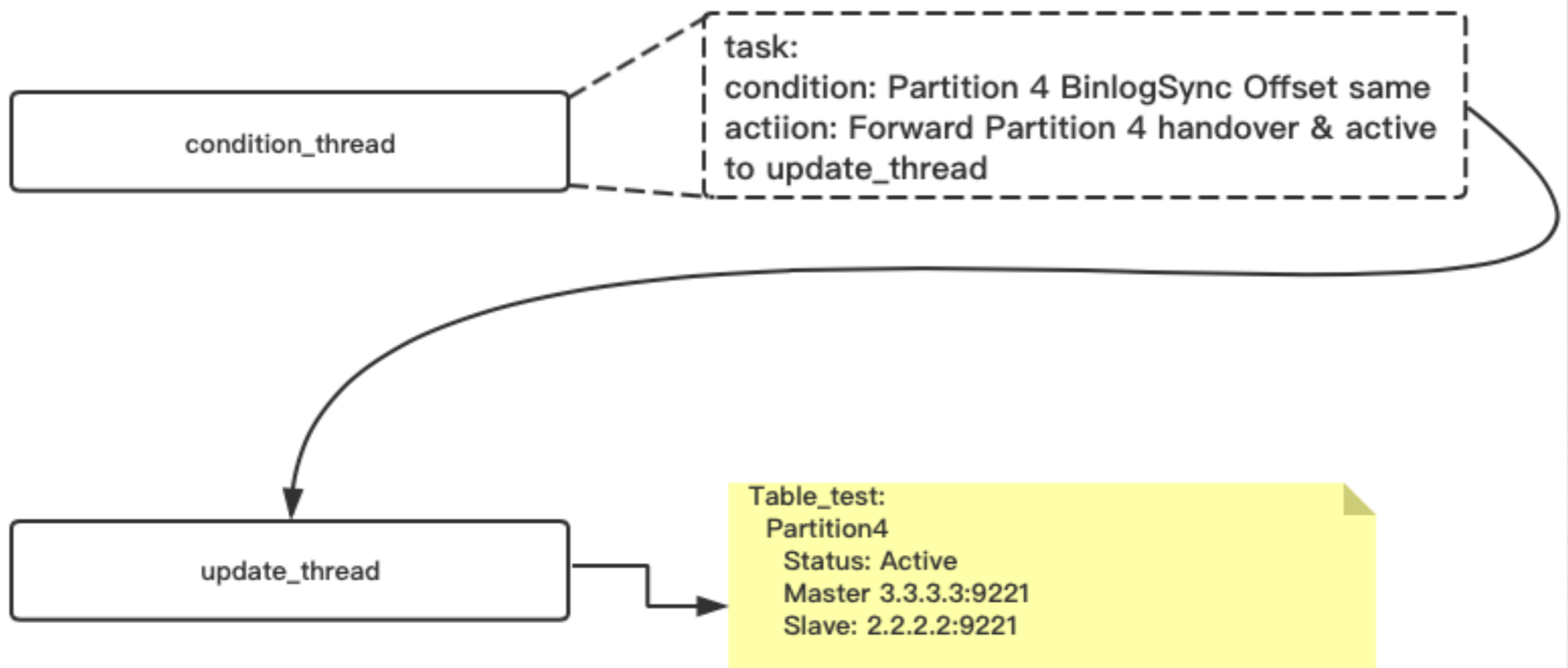


```
Table_test:  
Partition4  
Status: Stuck  
Master 1.1.1.1:9221  
Slave: 2.2.2.2:9221, 3.3.3.3:9221
```

- continue ping...

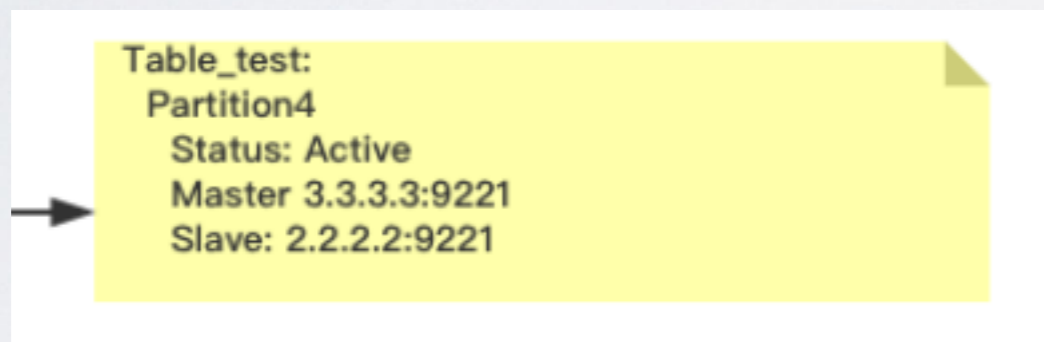


META SERVER

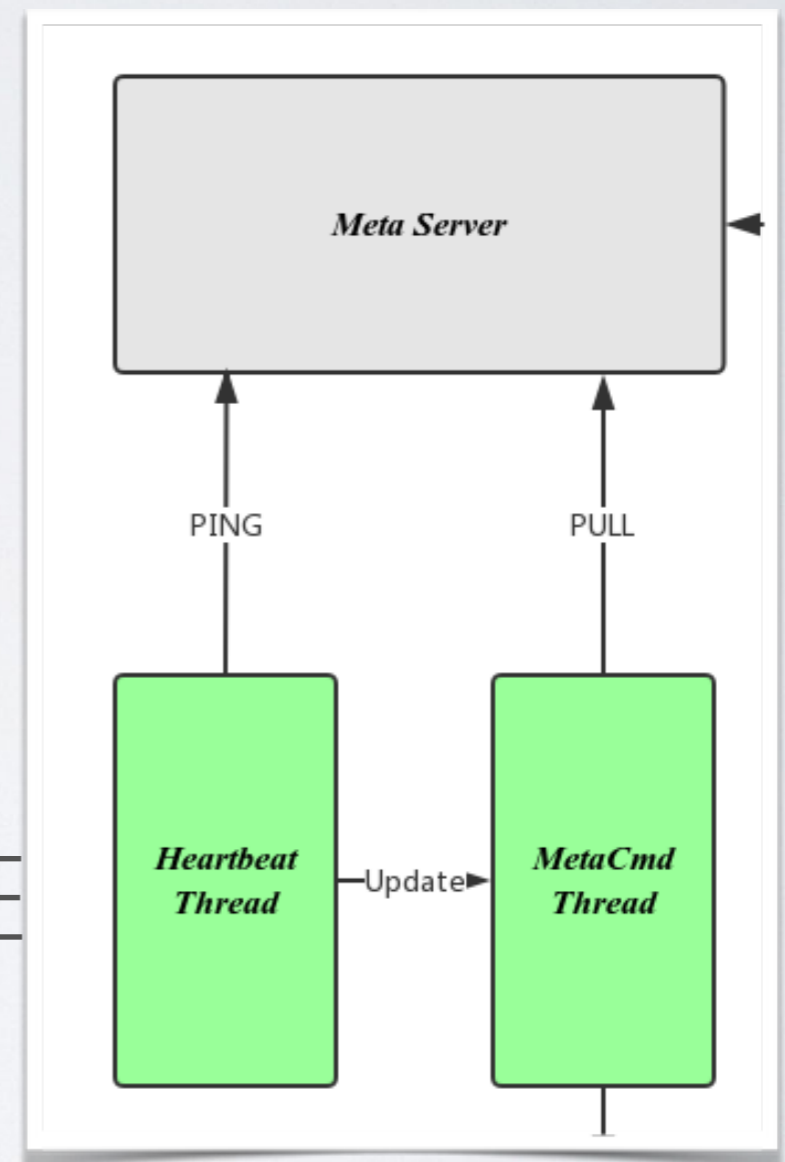


META SERVER

- node server apply meta info changes locally

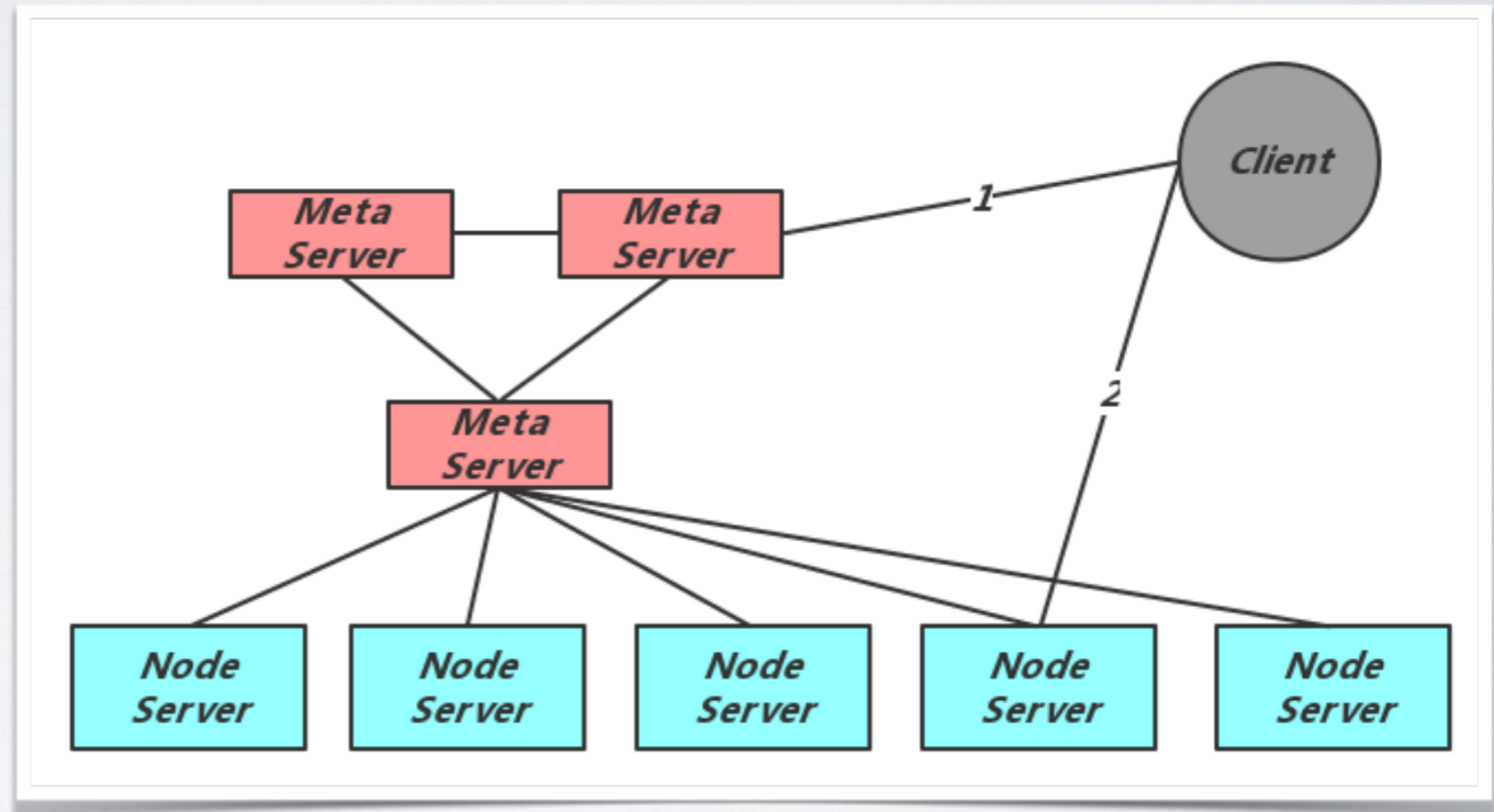


- New Partition4 Master is ACTIVE



ZEPPELIN OVERVIEW

- Pros & Cons
- Future.....



WHAT ABOUT DISTRIBUTED PIKA

- Synchronization Evolution Done

- Next.....





• THANKS