

# 基于MINIO的对象存储方案在探探的实践

于乐

探探科技

[yule@tantanapp.com](mailto:yule@tantanapp.com)



探探 Gopher China 2019



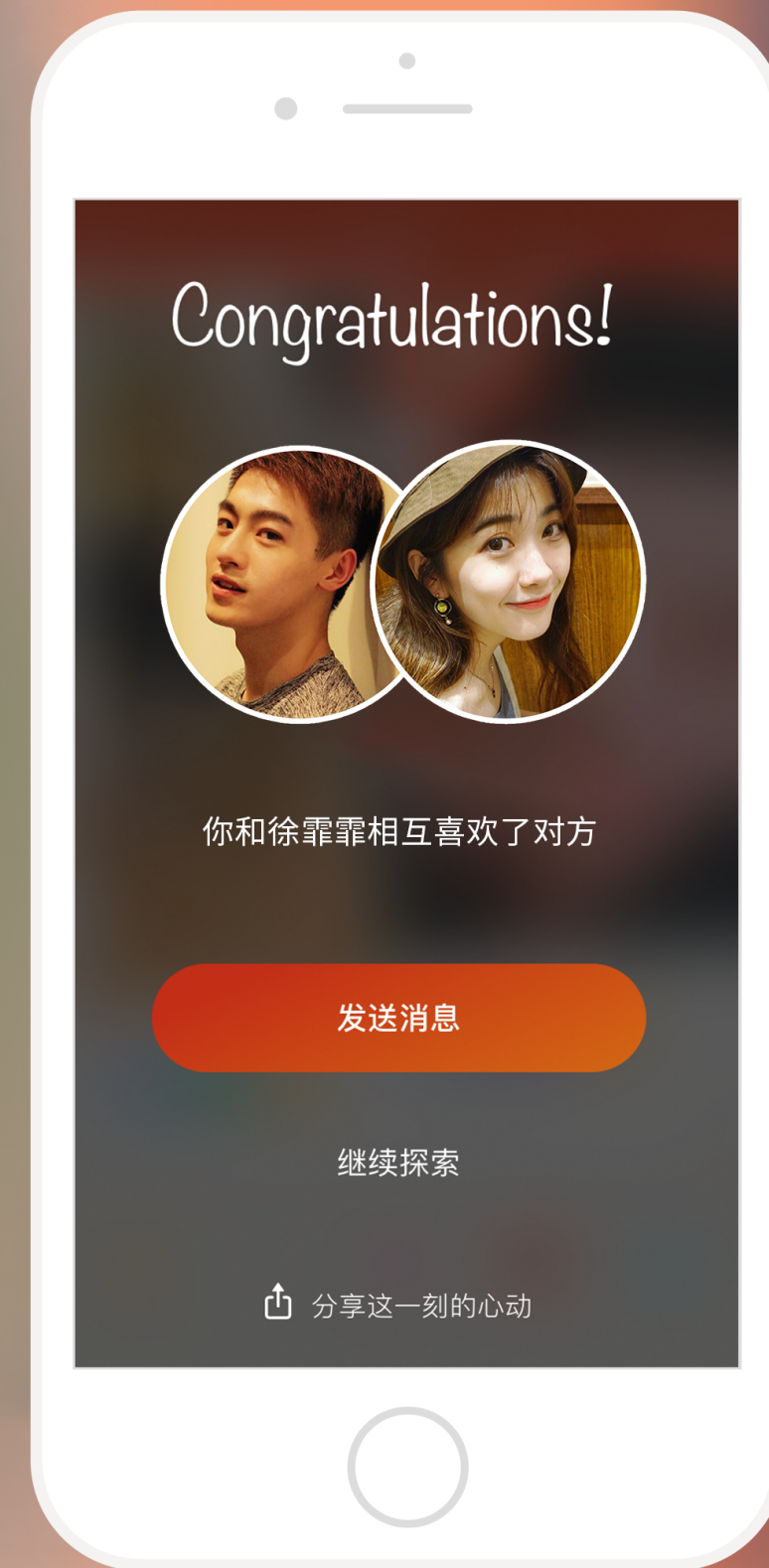
# 左滑无感 右滑喜欢

操作简单 超快认识人



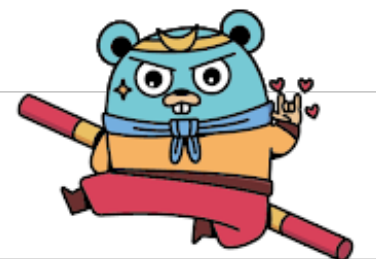
# 互赞开始聊天 零骚扰

全新的相遇方式 不再错过缘分

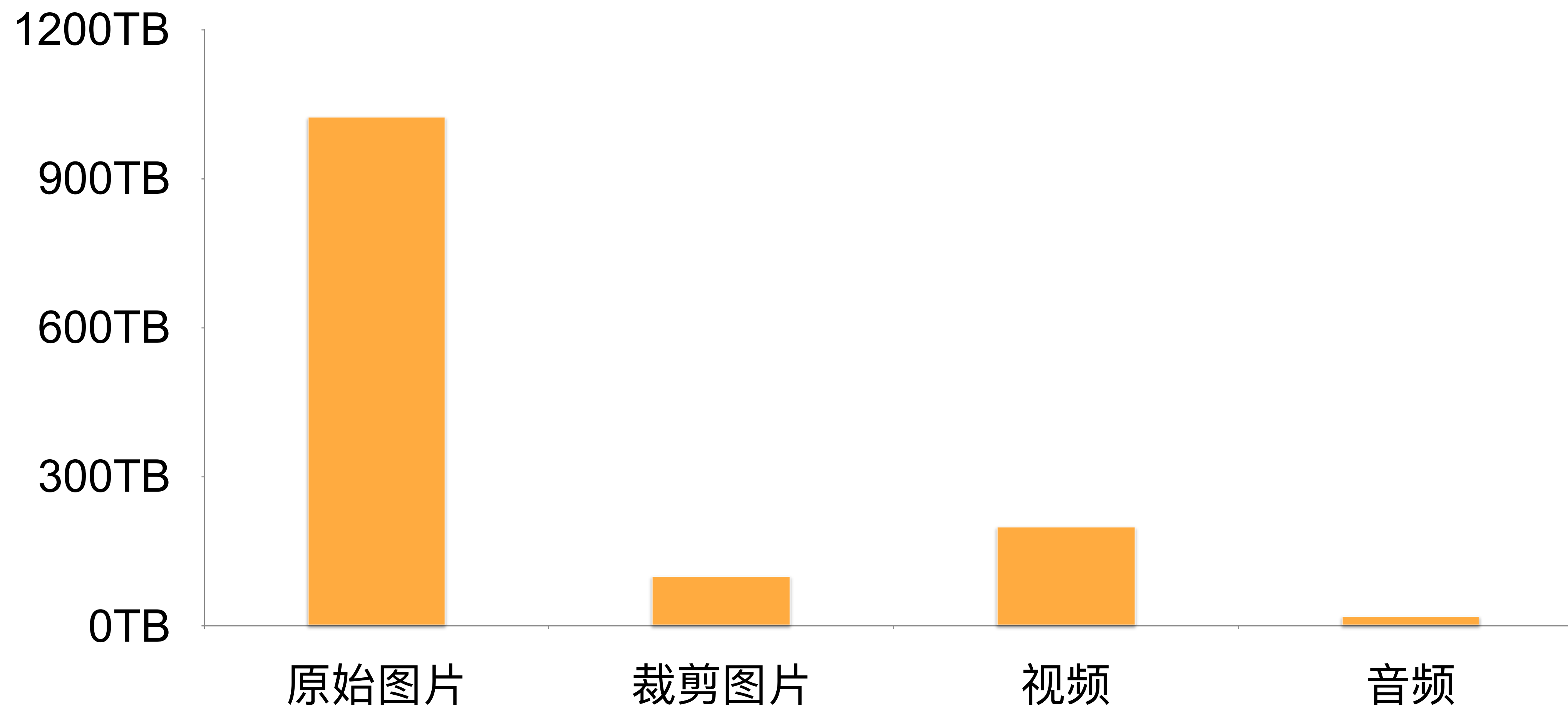


# 破冰利器 随心畅聊

还有好玩的真心话大冒险



# 存储现状



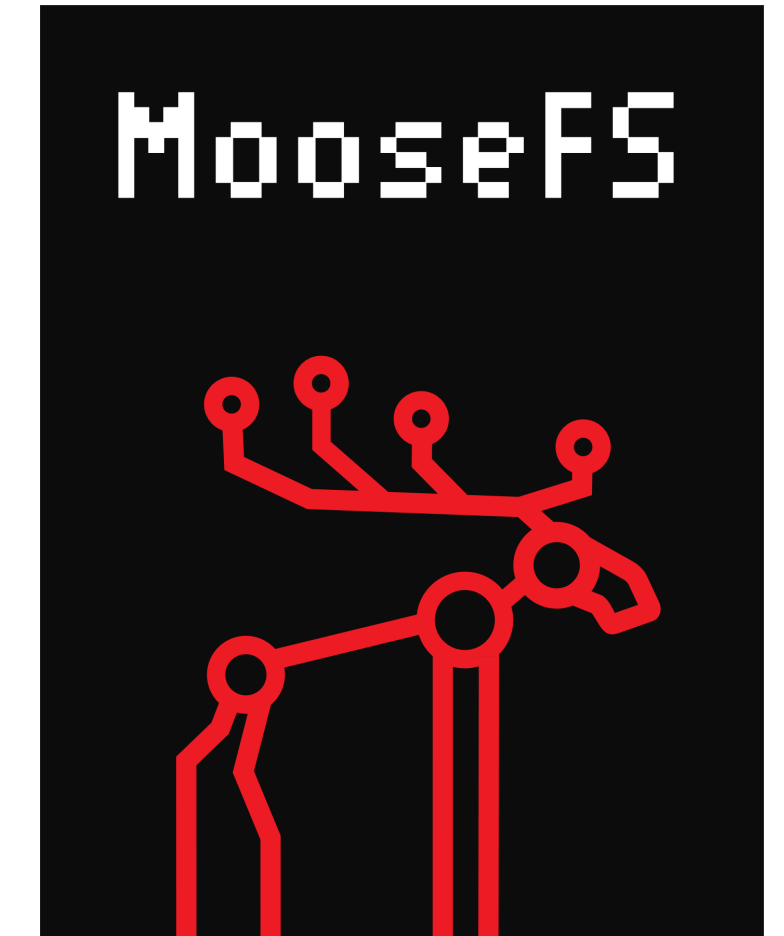


# 调研目前已有的开源方案

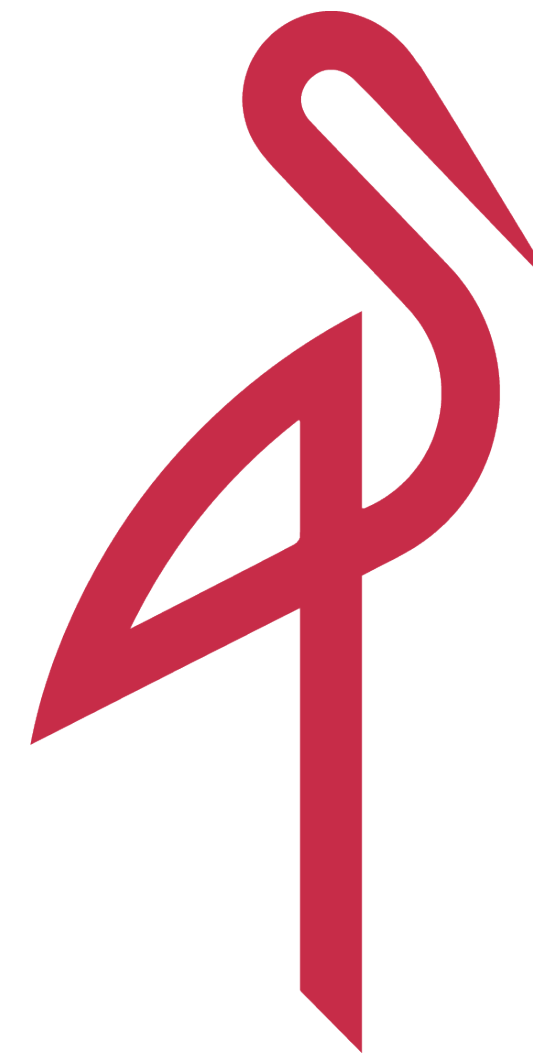


Ambry

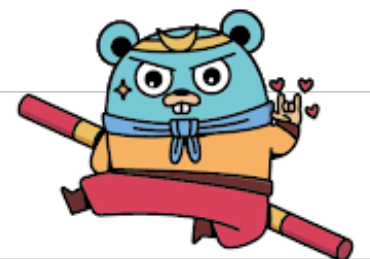
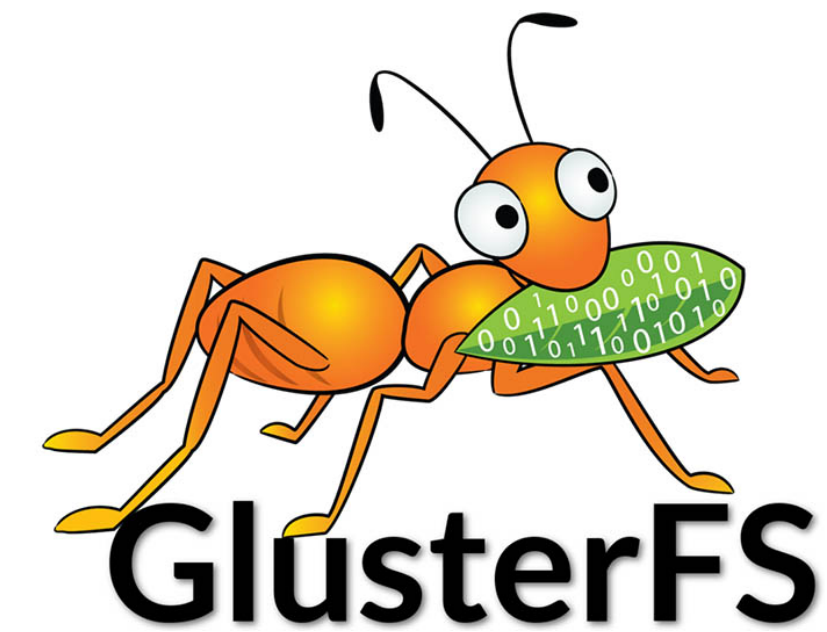
FastFS



MogileFS



TFS






Congratulations!

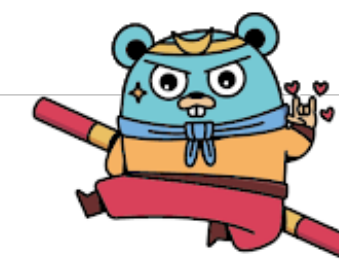


You and MINIO liked each other.

Let's Chat

Keep Swiping

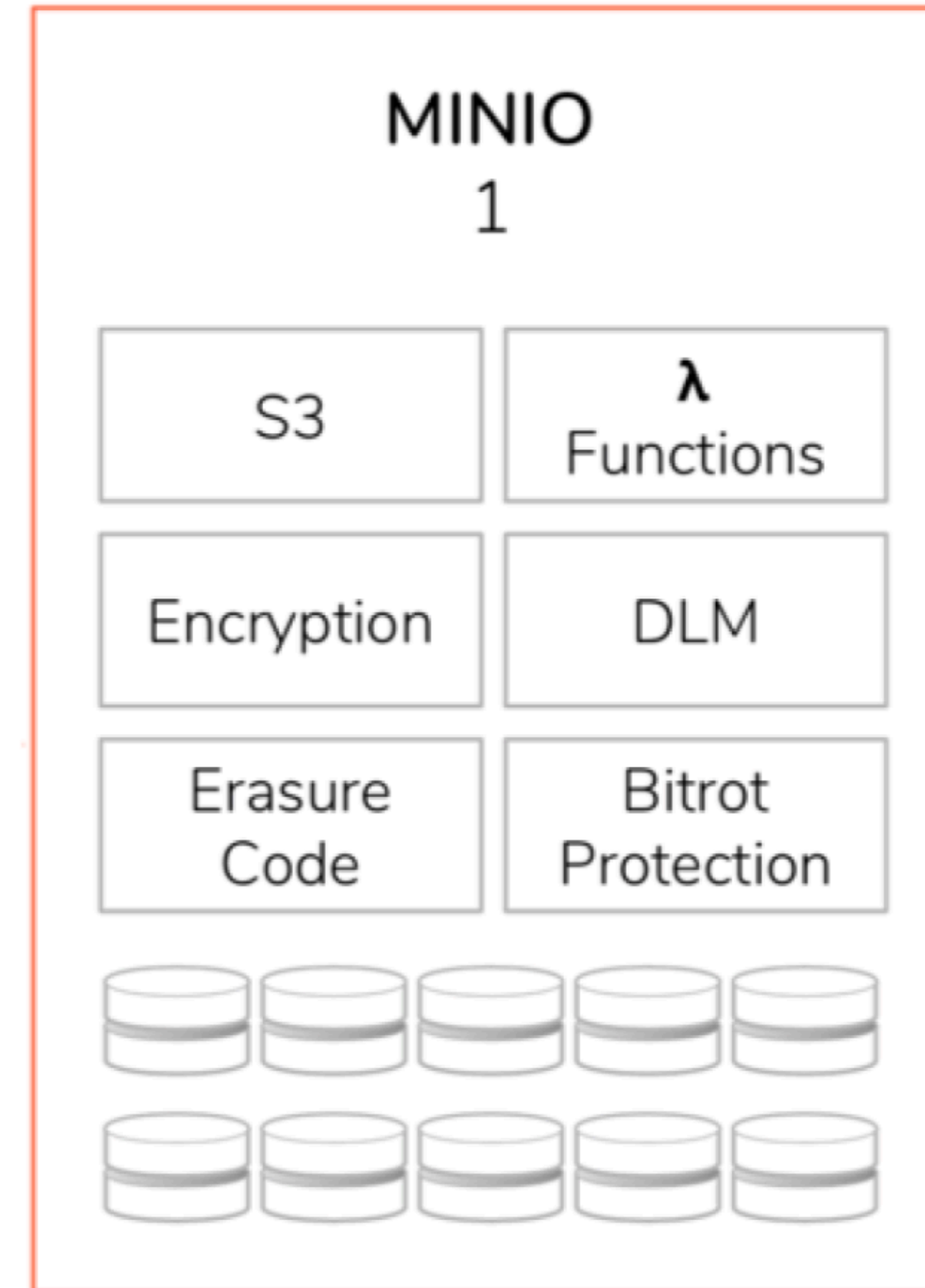
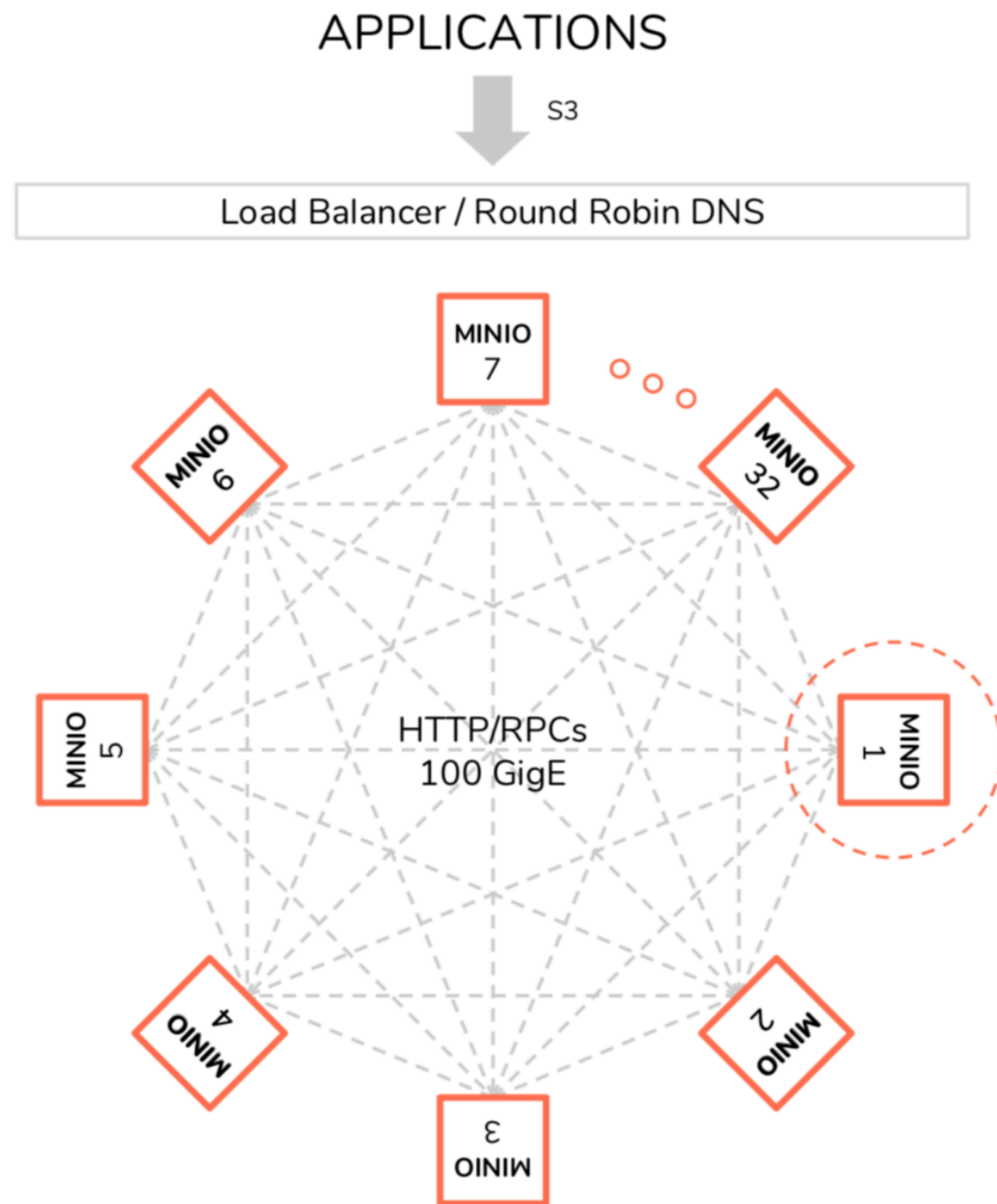
 Tell Your Friends



# 初探 MINIO



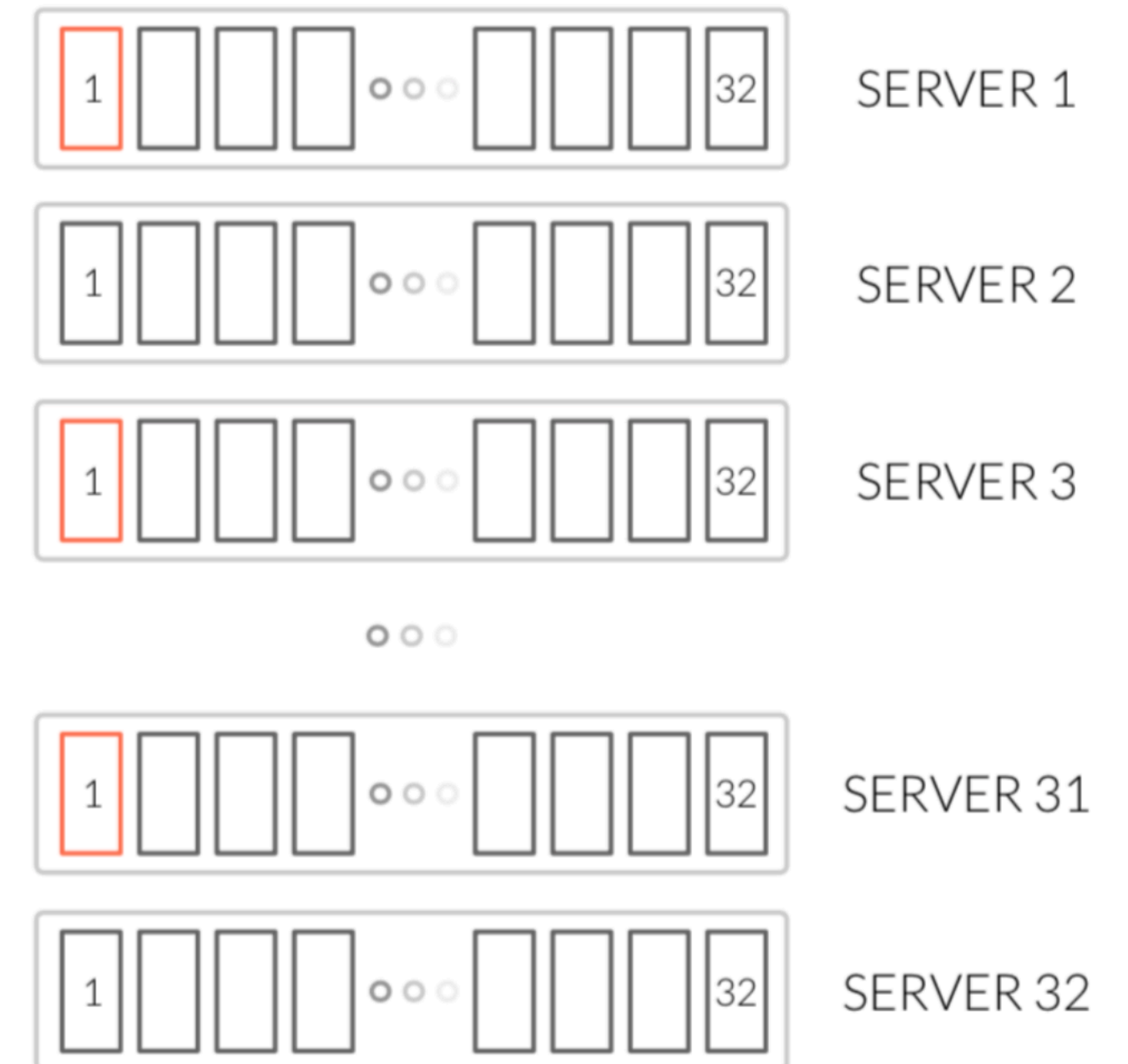
# MINIO 架构





# 数据分布以及可靠性

- Drive, 可以简单理解为一块硬盘
- Set , 一组 Drive 的集合
  - 一个对象存储在一个 Set 上
  - 一个集群划分为多个 Set
  - 一个 Set 包含的 Drive 数量是固定的
    - 默认由系统根据集群规模自动计算得出
    - MINIO\_ERASURE\_SET\_DRIVE\_COUNT
  - 一个 SET 中的 Drive 尽可能分布在不同节点上



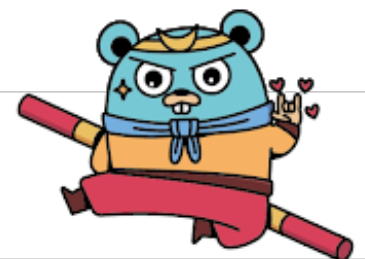
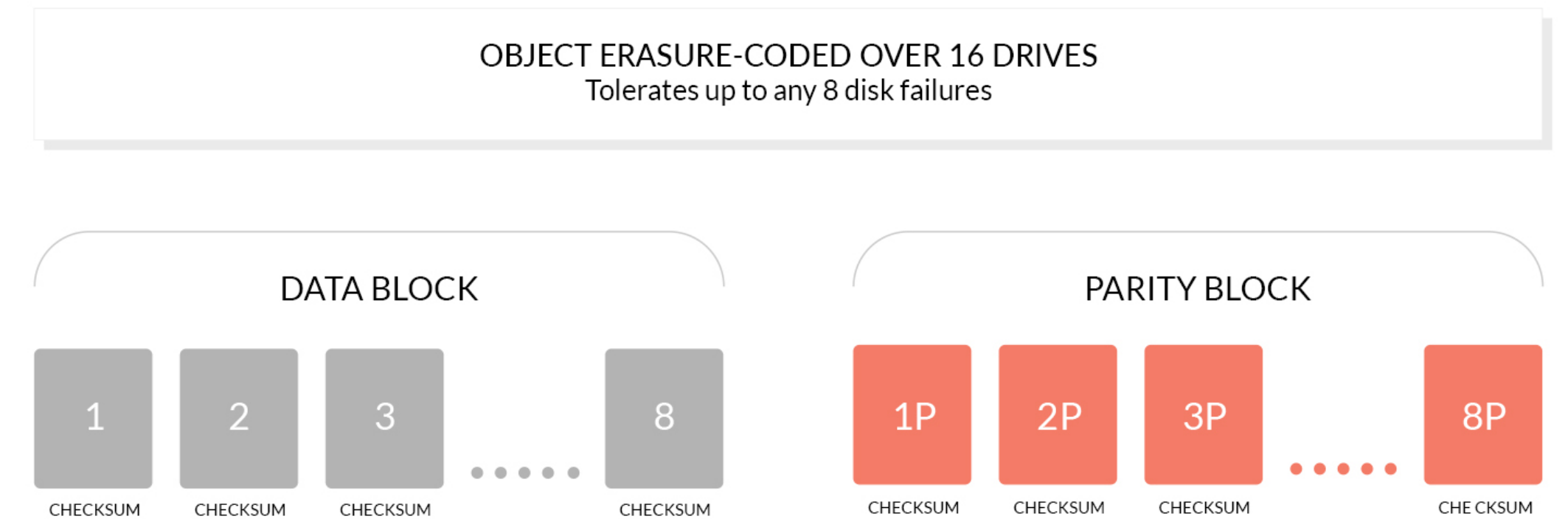
# 数据编码

Erasure Code, a mathematical algorithm to reconstruct missing or corrupted data.

- 将一个对象编码成若干个数据块和跟校验块
- Reed-Solomon
- 低冗余，高可靠

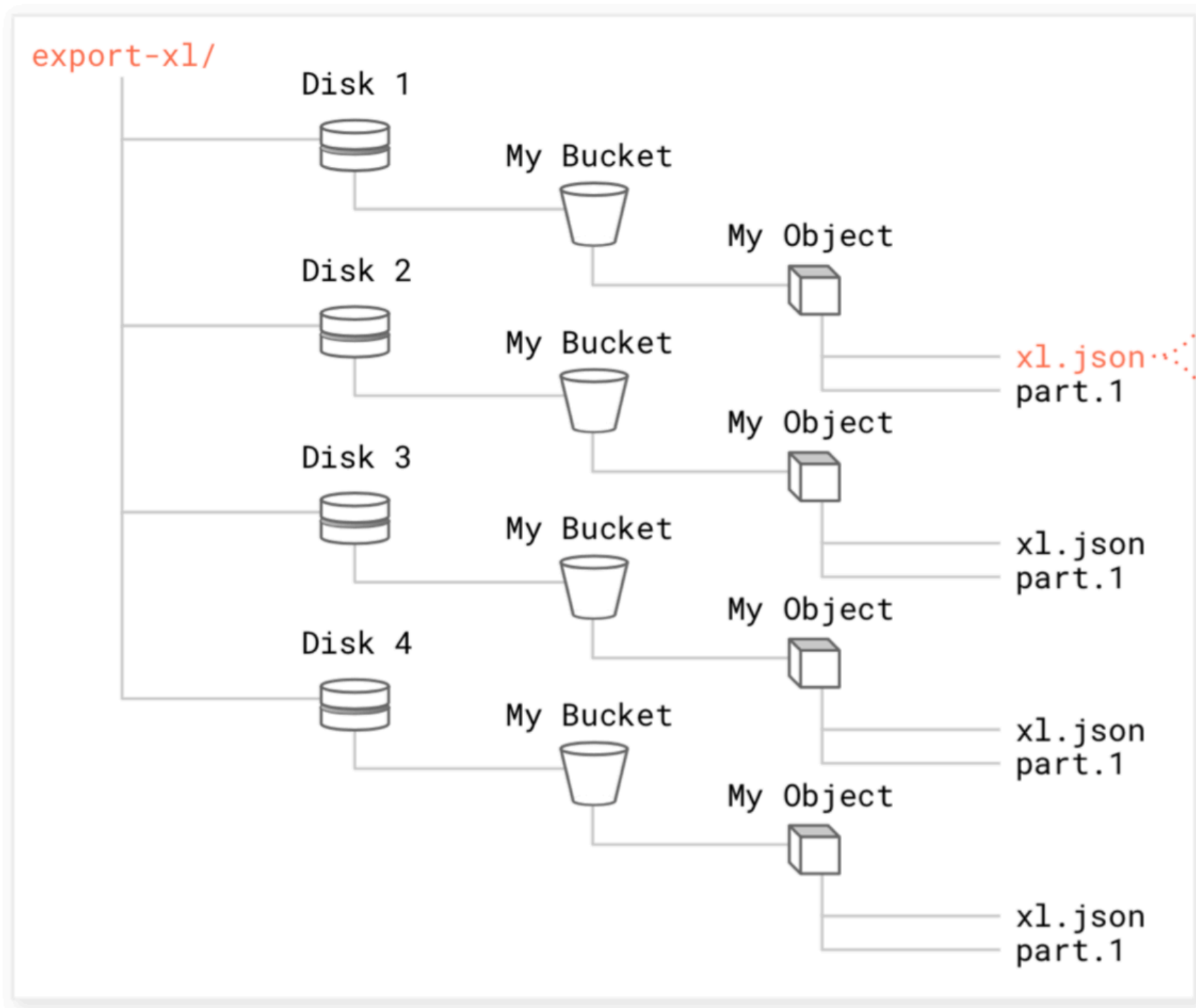
## Bit Rot Protection

- HighwayHash

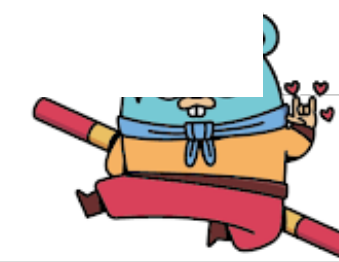


# MINIO

## Object Erasure Code

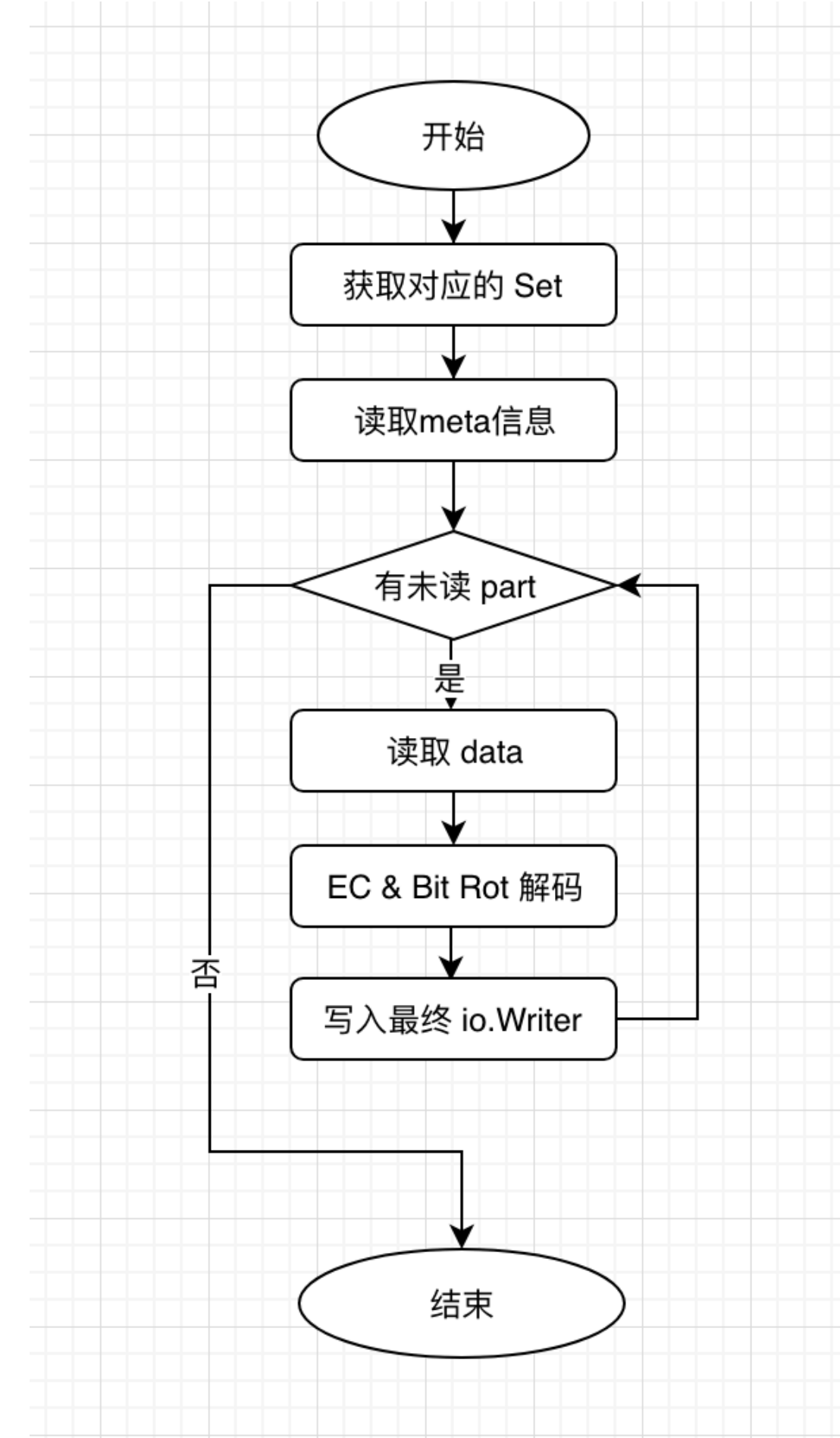
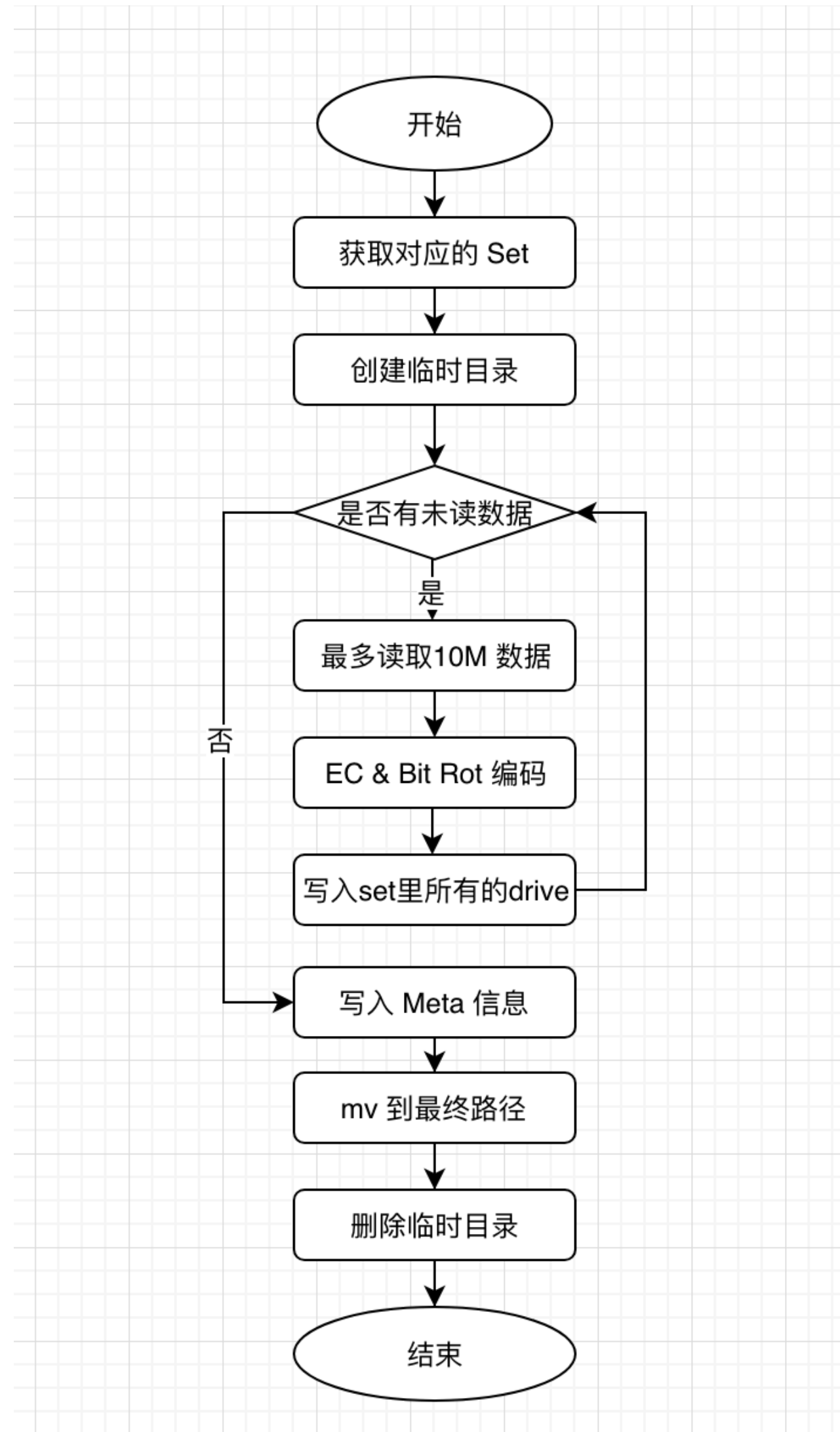


```
{  
  "version": "1.0.1",  
  "format": "xl",  
  "stat": {  
    "size": 2286,  
    "modTime": "2017-12-02T00:24:20.975968336Z",  
  },  
  "erasure": {  
    "algorithm": "klauspost/reedsolomon/vandermonde",  
    "data": 2,  
    "parity": 2,  
    "blockSize": 10485760,  
    "index": 2,  
    "distribution": "[  
      2,  
      3,  
      4,  
      1  
    ]",  
    "checksum": [  
      {  
        "name": "part.1",  
        "algorithm": "blake2b",  
        "hash": "c24fa0451fd85a3a482c...b672b7f"  
      }  
    ],  
    "minio": {  
      "release": "DEVELOPMENT.GOGET"  
    },  
    "meta": {  
      "content-type": "application/octet-stream",  
      "etag": "c1d217c52d44c9eab00e81496b2b91b6"  
    },  
    "parts": [  
      {  
        "number": 1,  
        "name": "part.1",  
        "etag": "",  
        "size": 2286  
      }  
    ]  
  }  
}
```





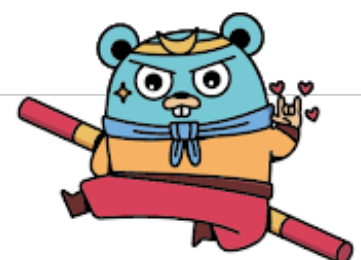
# 上传下载流程



# IO 实现细节

```
var onlineDisks []StorageAPI
// ...
writers := make([]io.Writer, len(onlineDisks))
for i, disk := range onlineDisks {
    writers[i] = newBitrotWriter(
        disk, minioMetaTmpBucket, tempErasureObj,
        erasure.ShardFileSize(curPartSize),
        DefaultBitrotAlgorithm,
        erasure.ShardSize())
}
// ...
n, erasureErr := erasure.Encode(
    ctx, curPartReader, writers, buffer, erasure.dataBlocks+1)
// ...
```

```
type Writer interface {
    Write(p []byte) (n int, err error)
}
type Reader interface {
    Read(p []byte) (n int, err error)
}
```



# 巧用 channel

```
readTriggerCh := make(chan bool, len(p.readers))
for i := 0; i < p.dataBlocks; i++ {
    readTriggerCh <- true
}
for readTrigger := range readTriggerCh {
    canDecode := p.canDecode(newBuf)
    if canDecode {
        break
    }
    if !readTrigger {
        continue
    }
    wg.Add(1)
    go func(i int) {
        defer wg.Done()
        disk := p.readers[i]
        if disk == nil {
            readTriggerCh <- true
            return
        }
        _, err := disk.ReadAt(p.buf[i], p.offset)
        if err != nil {
            readTriggerCh <- true
            return
        }
        // ...
        readTriggerCh <- false
    }(readerIndex)
    readerIndex++
}
wg.Wait()
```





# 如何部署

```
# 启动一个12个节点的 minio 集群，每个节点有12个drive  
minio server http://10.3.1.{1...12}/data{1...12}
```

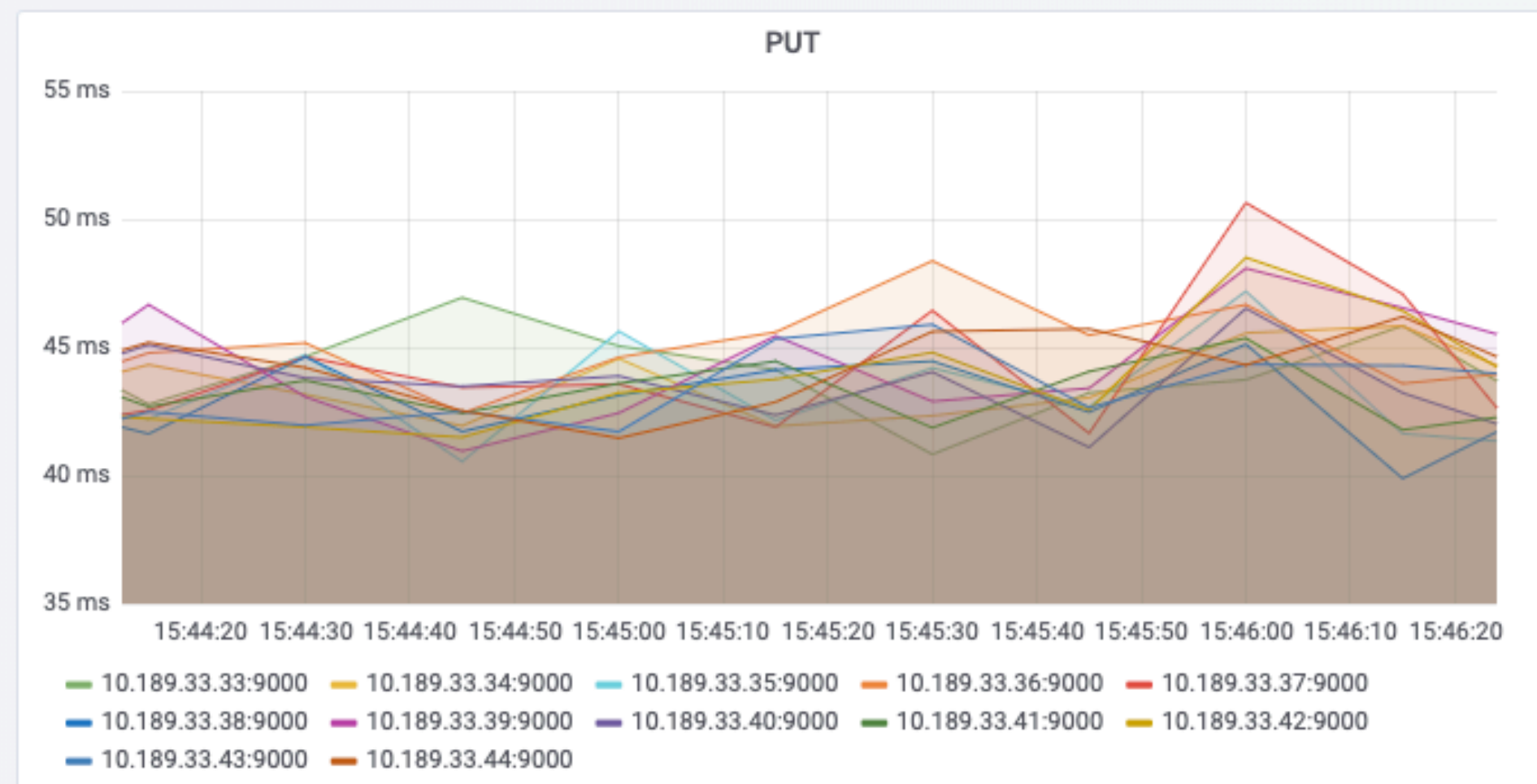
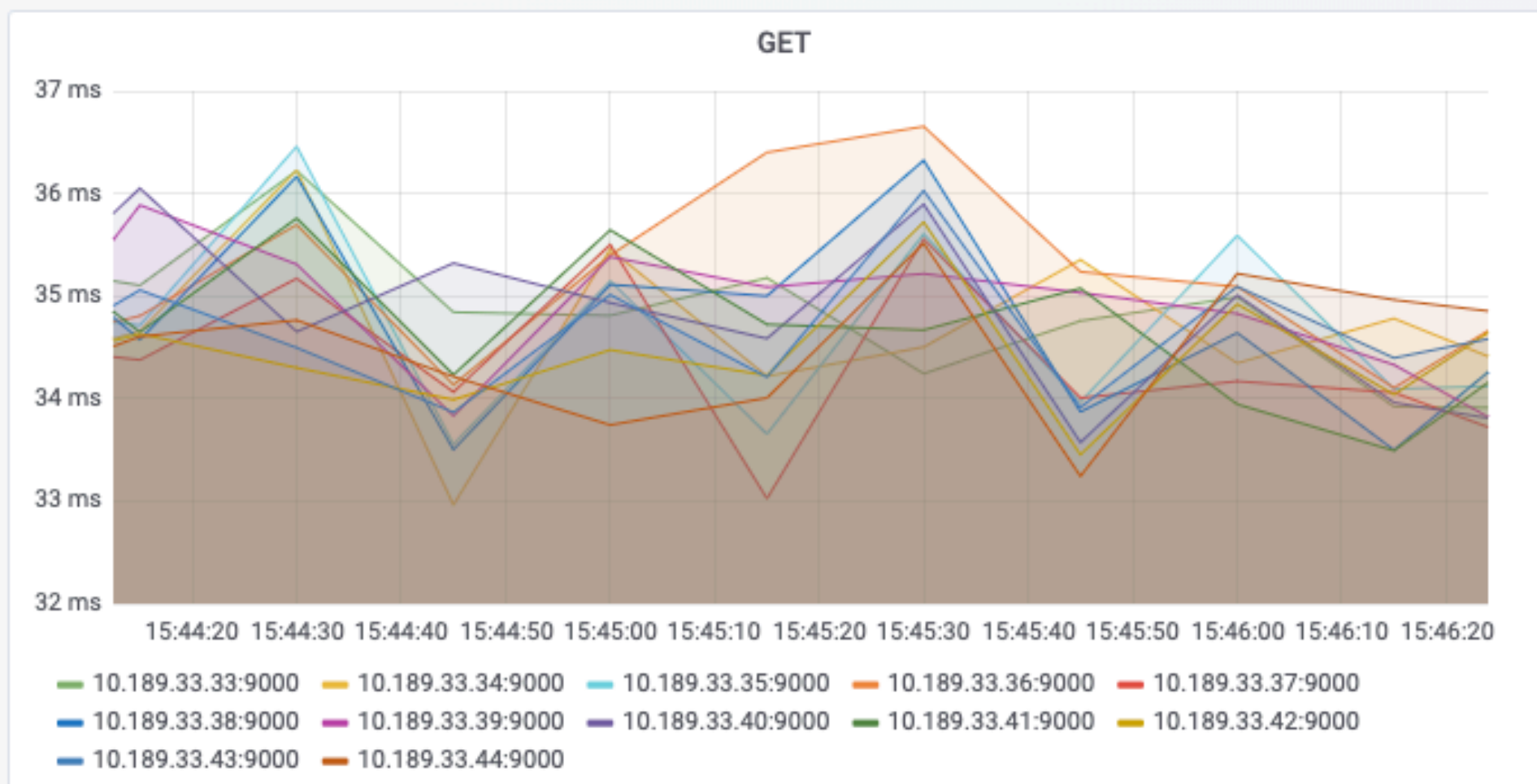


# 性能压测

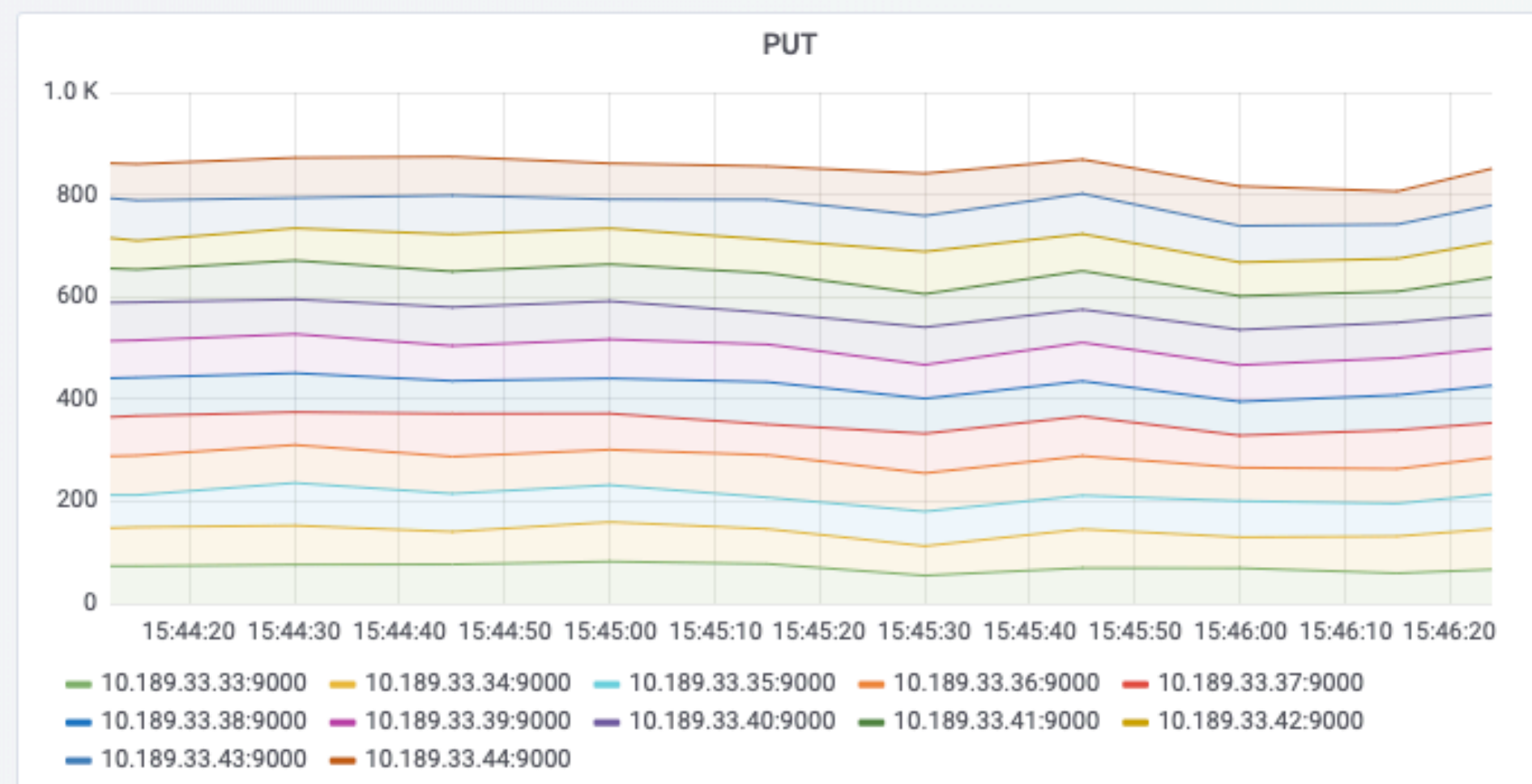
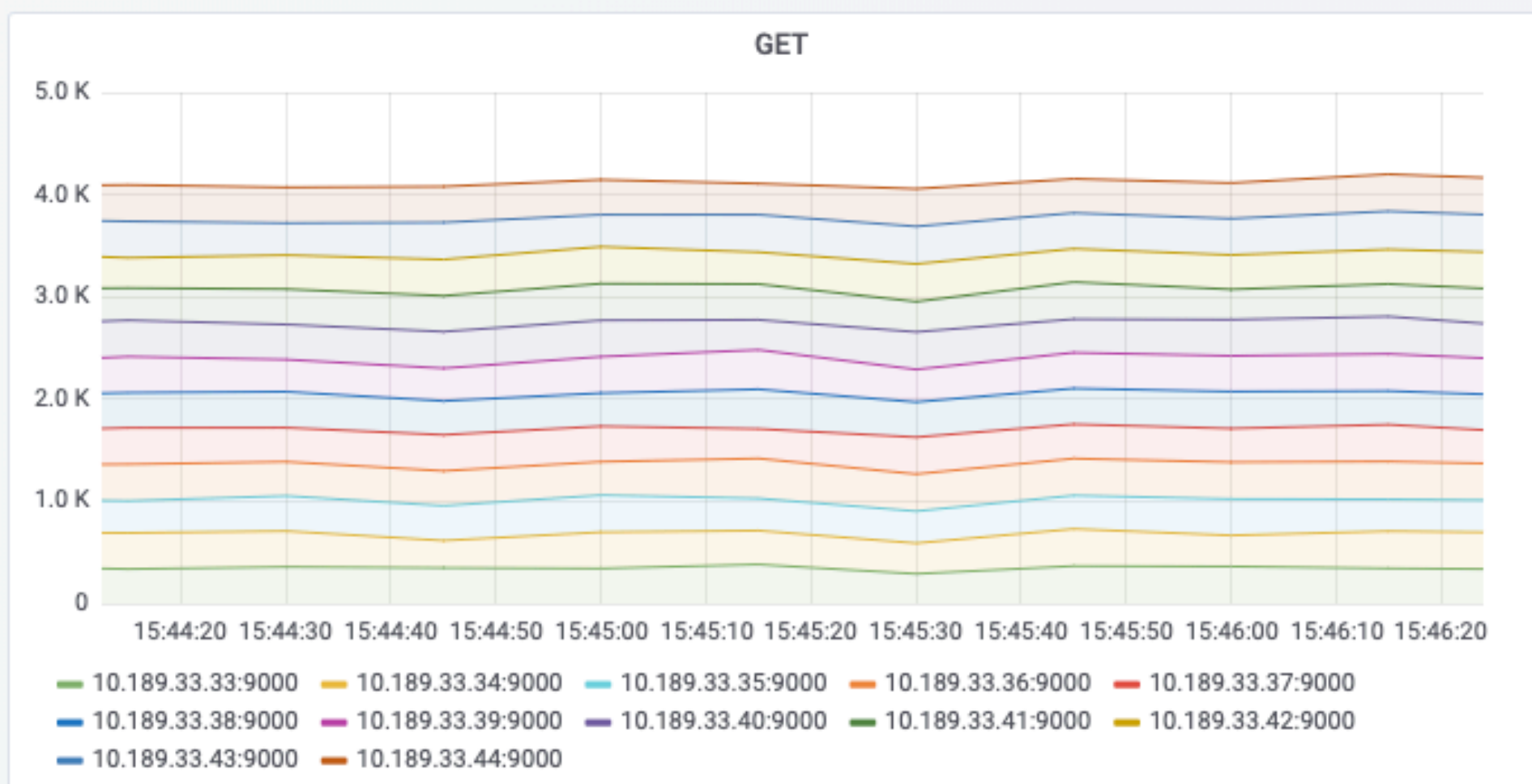
- CosBench
- 参考 Minio [官方压测方案](#)
- 不同大小，不同并发，不同读写比例组合测试



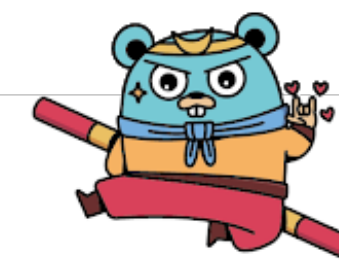
HTTP Request Average Duration



HTTP Requests Count per second



12 个节点， 每个节点 12 块盘， 混合读写 64KB





# 多集群扩容



# MINIO Limitation

- 集群搭建后之后不允许扩容
- 大集群最大节点数为 32
- Dsync 分布式锁



“

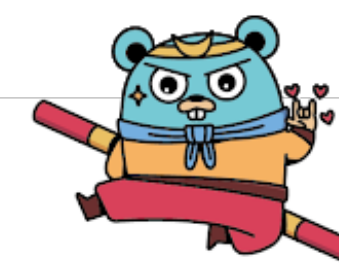
*We can certainly do single large namespace, but then taking one cluster down can potentially be a large downtime and a very complex deployment. This is exactly what we did with GlusterFS for 10yrs. This why most storage solutions never scale beyond few PBs.*

*Minio's fundamental design is to avoid such practical issues in production and make managing clusters easier. Failures are common in real world scenarios and when it's a gigantic pool, it is very hard to know what is going on.*

”

*You can avoid these details in your application, if you are clever about it on how you lookup URLs.*

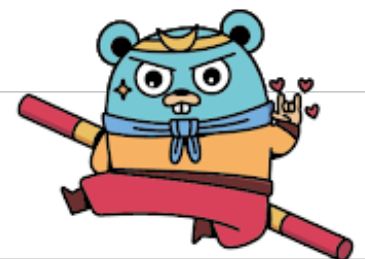
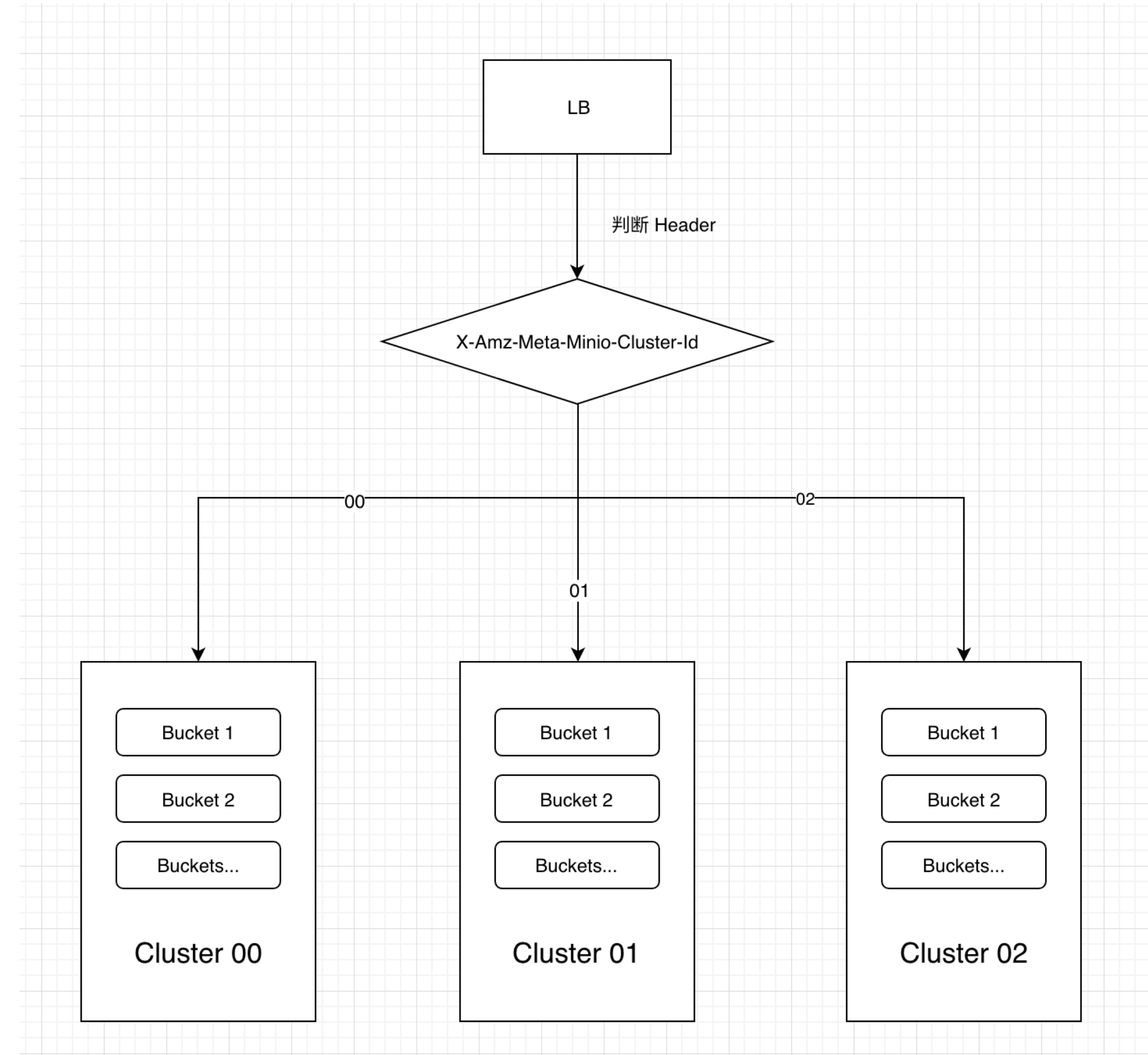
— Harshavardhana





# 扩容方案

- 可同时向多个集群写入
- 集群容量到达阈值，从写集群中剔除，变成只读
- 读取的时候根据编码到文件名上的集群信息，从而路由到正确的集群

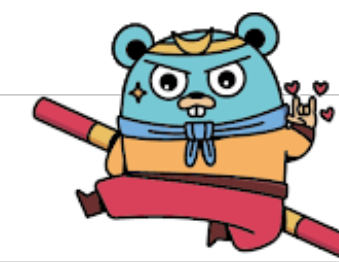


# 万事俱备，只欠东风

- 筹划数据迁移
- 准备接入服务
- 采购机器

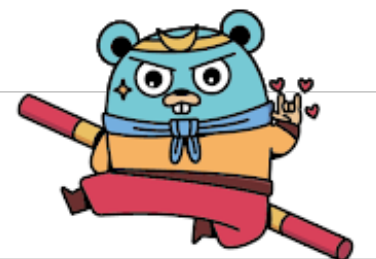


生活并不是一帆风顺的



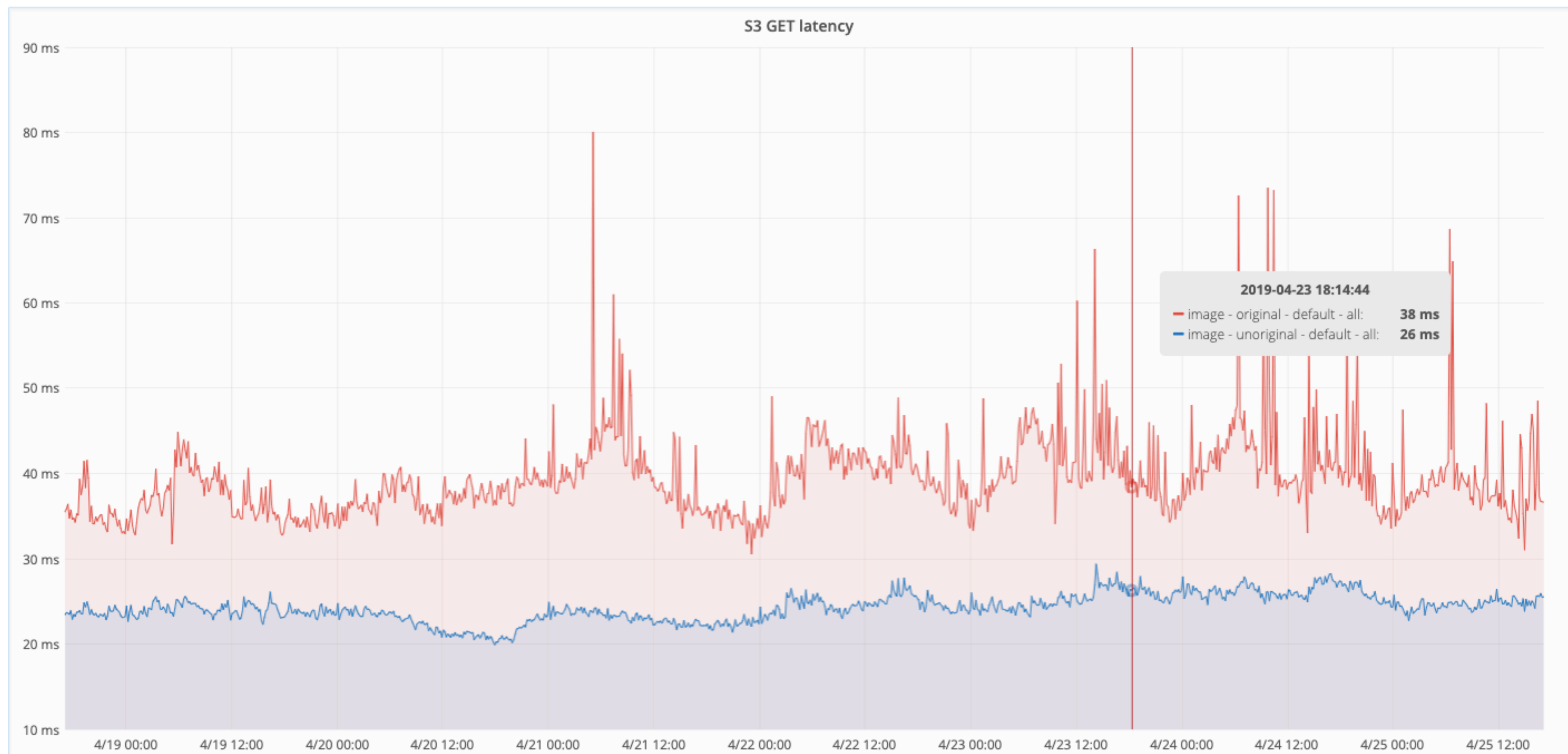
# 实际流量测试

- 100ms
- 200ms
- 500ms
- 报警!

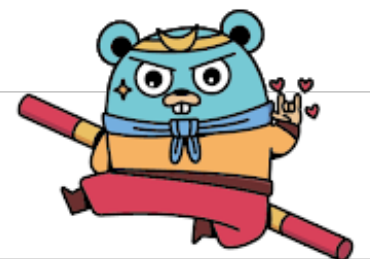
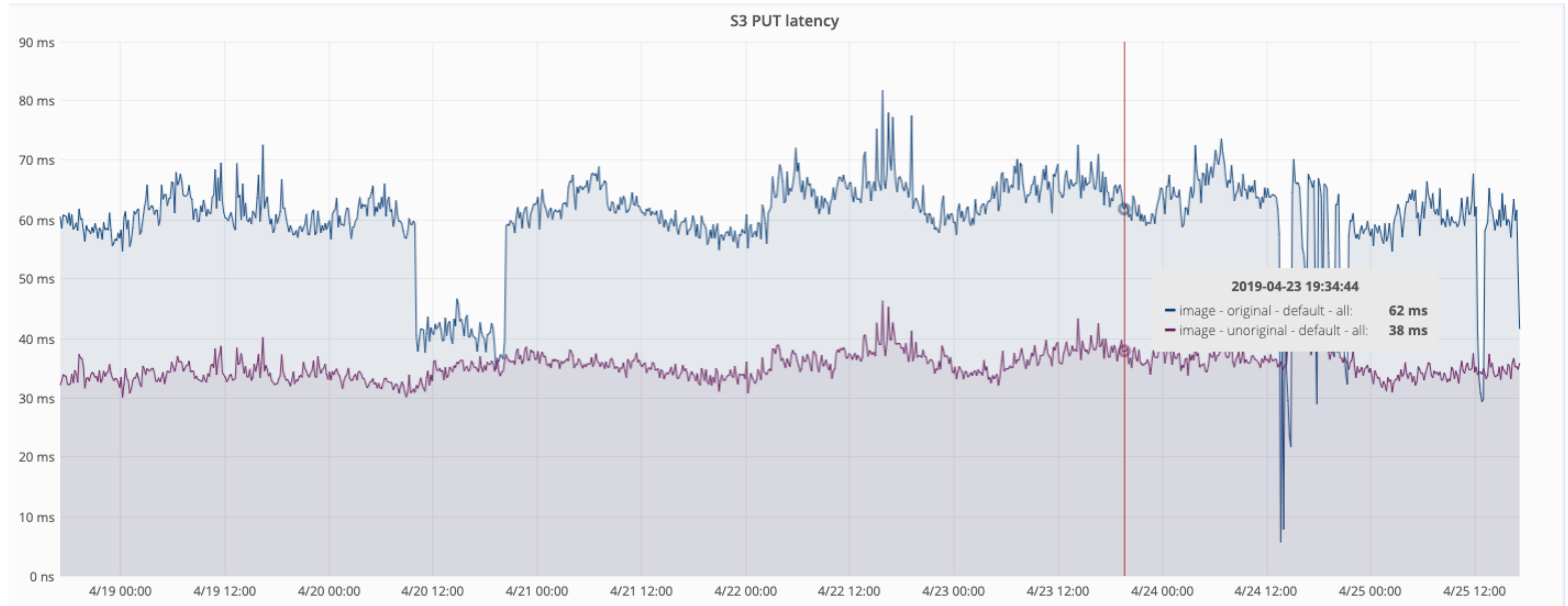




# S3 下载响应时间



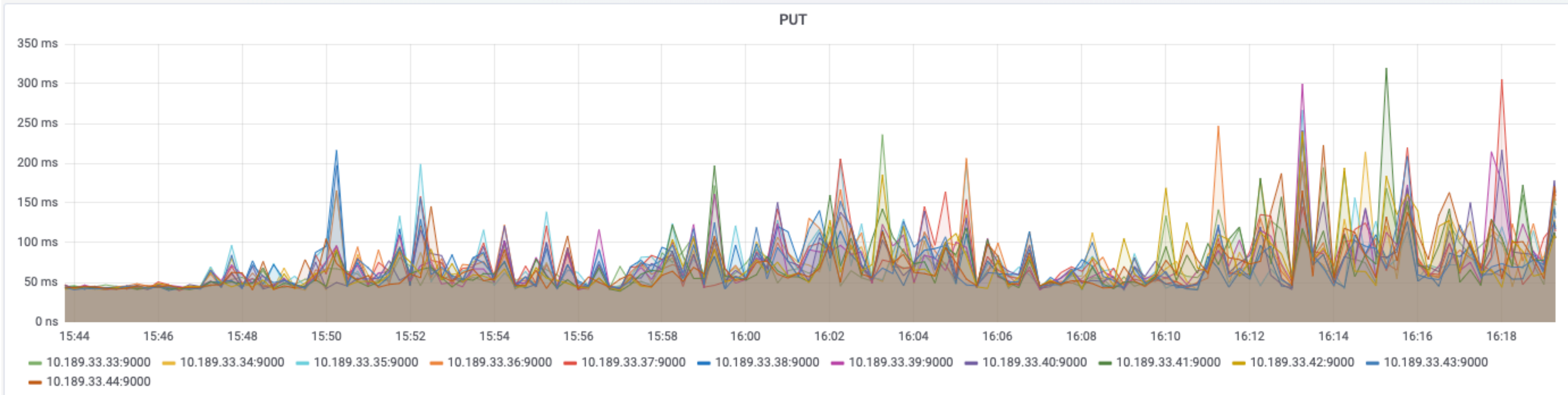
# S3 上传响应时间



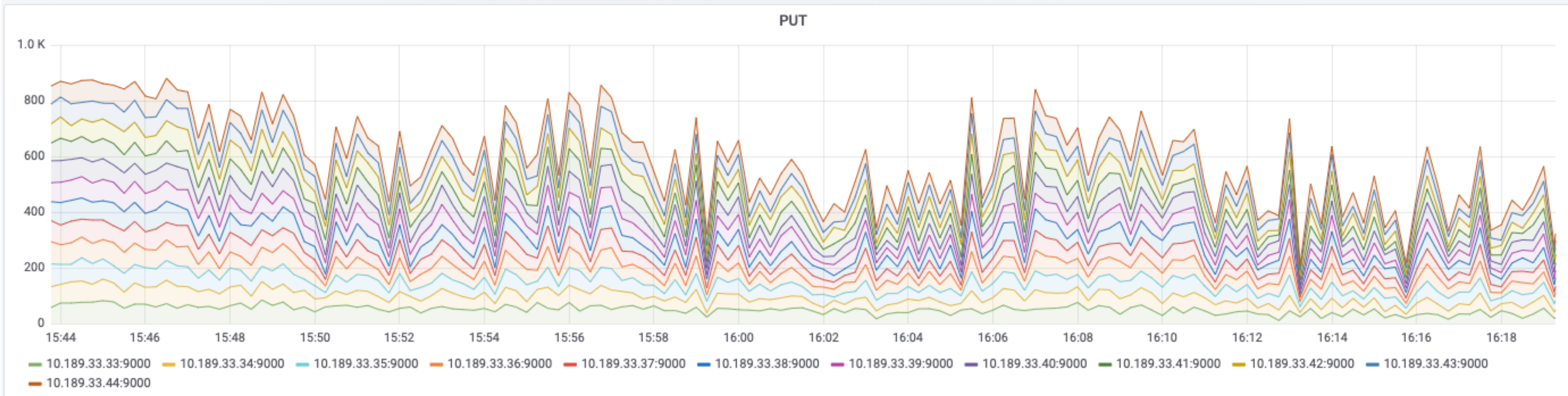


# Minio 上传性能

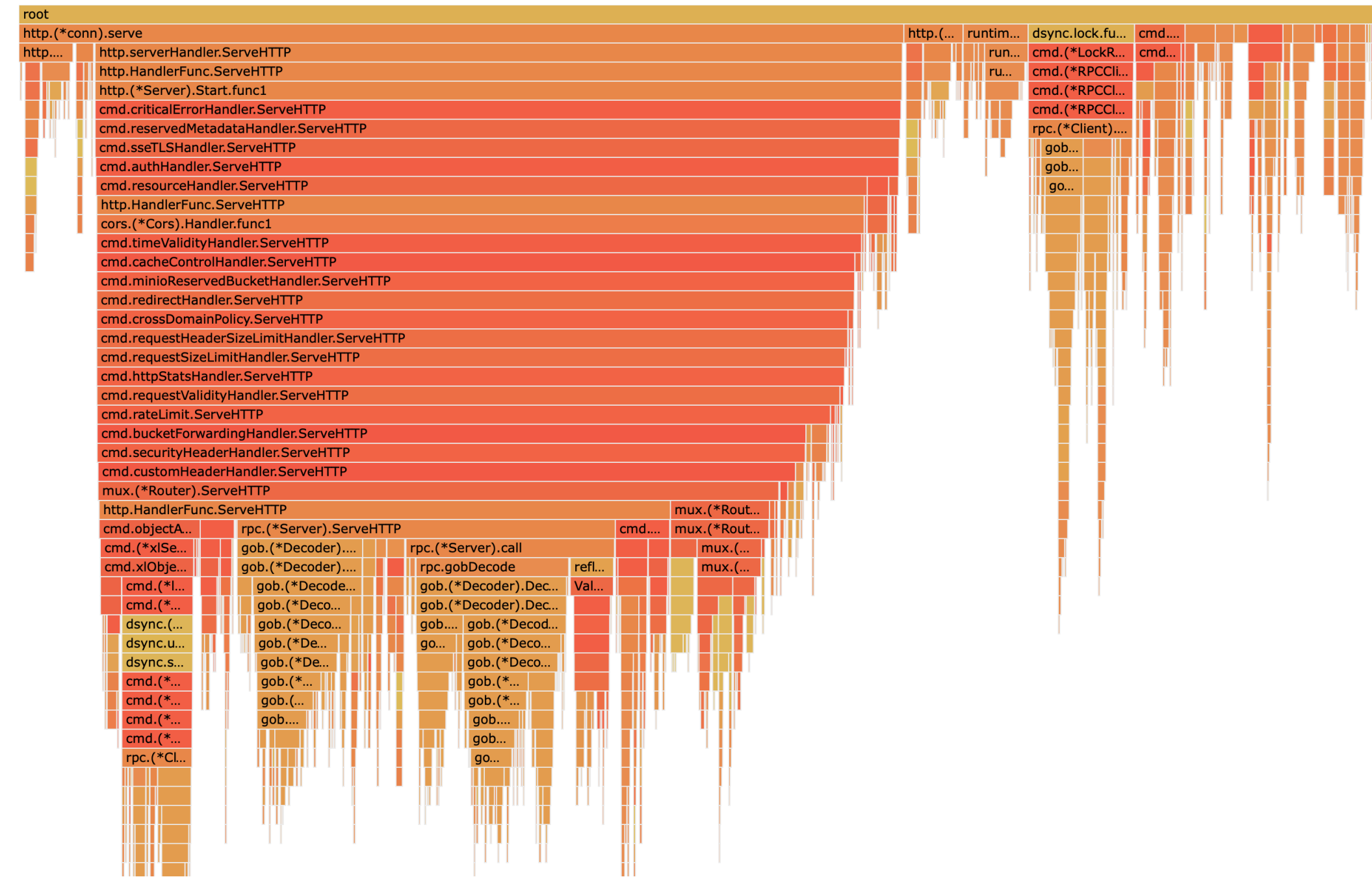
## HTTP Request Average Duration



## HTTP Requests Count per second



# MINIO Profile

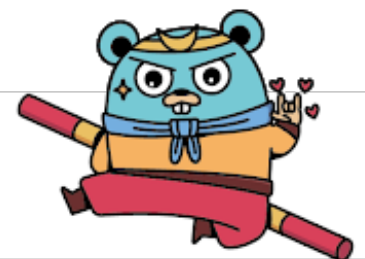


```
# 开始采集整个集群指定 profile 类型的数据, 类型可选 'cpu', 'mem', 'block', 'mutex', 'trace'  
mc admin profile start -type cpu  
# 结束采集, 自动下载采集结果文件  
mc admin profile stop  
# 使用 google pprof 工具进行可视化展示  
pprof -http=0.0.0.0:2222 ./profiling-10.189.33.60\ :9000.pprof
```





# 杀手锏



# 全面了解IO



# 全面了解 IO

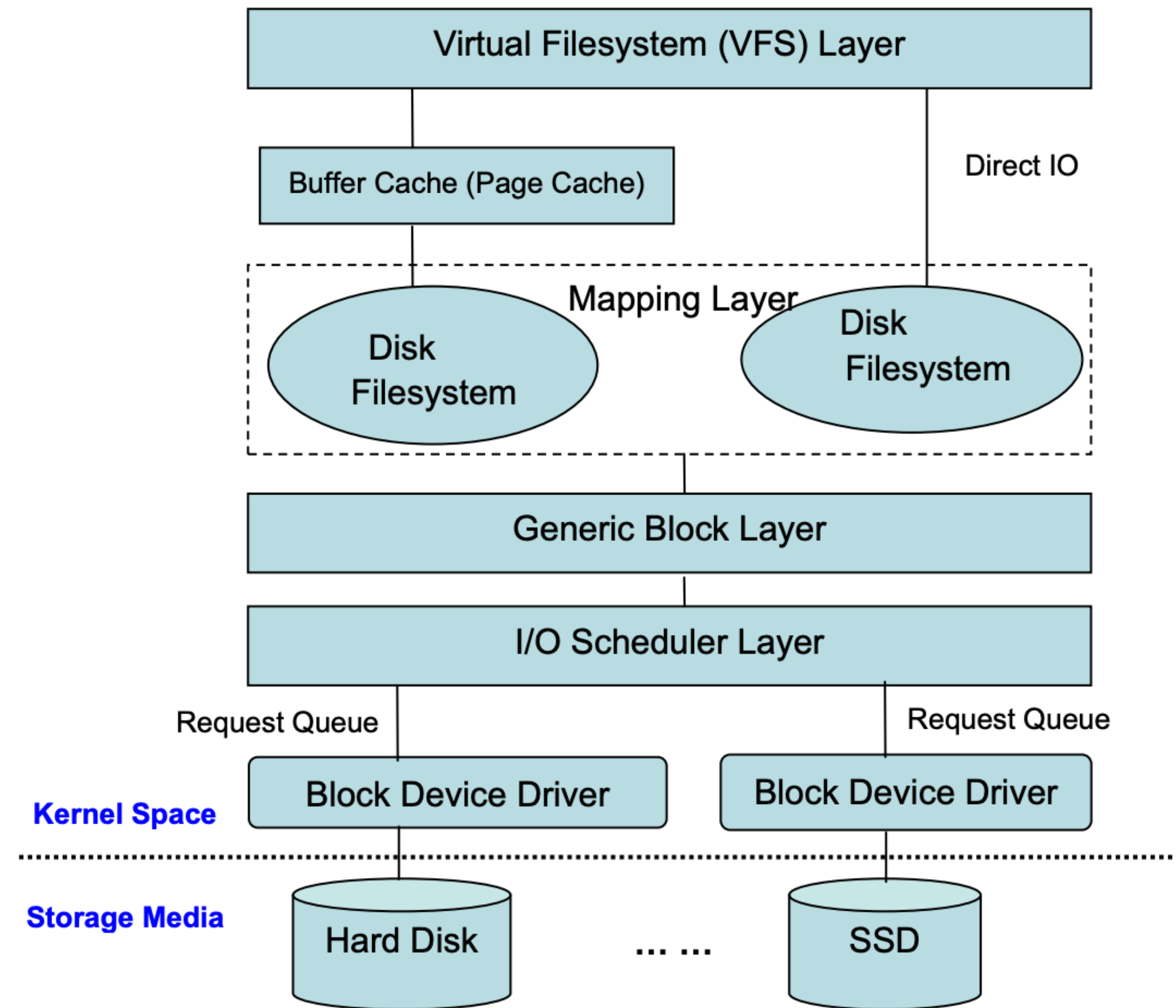
- IO 流程
- Buffered IO vs Direct IO
- Page Cache
- IO Scheduler
- 磁盘性能
- 参数调优



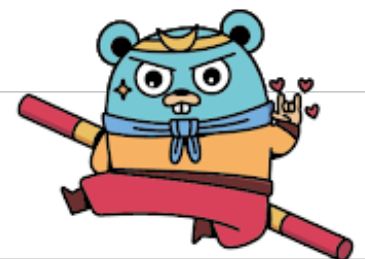




# Block IO 流程

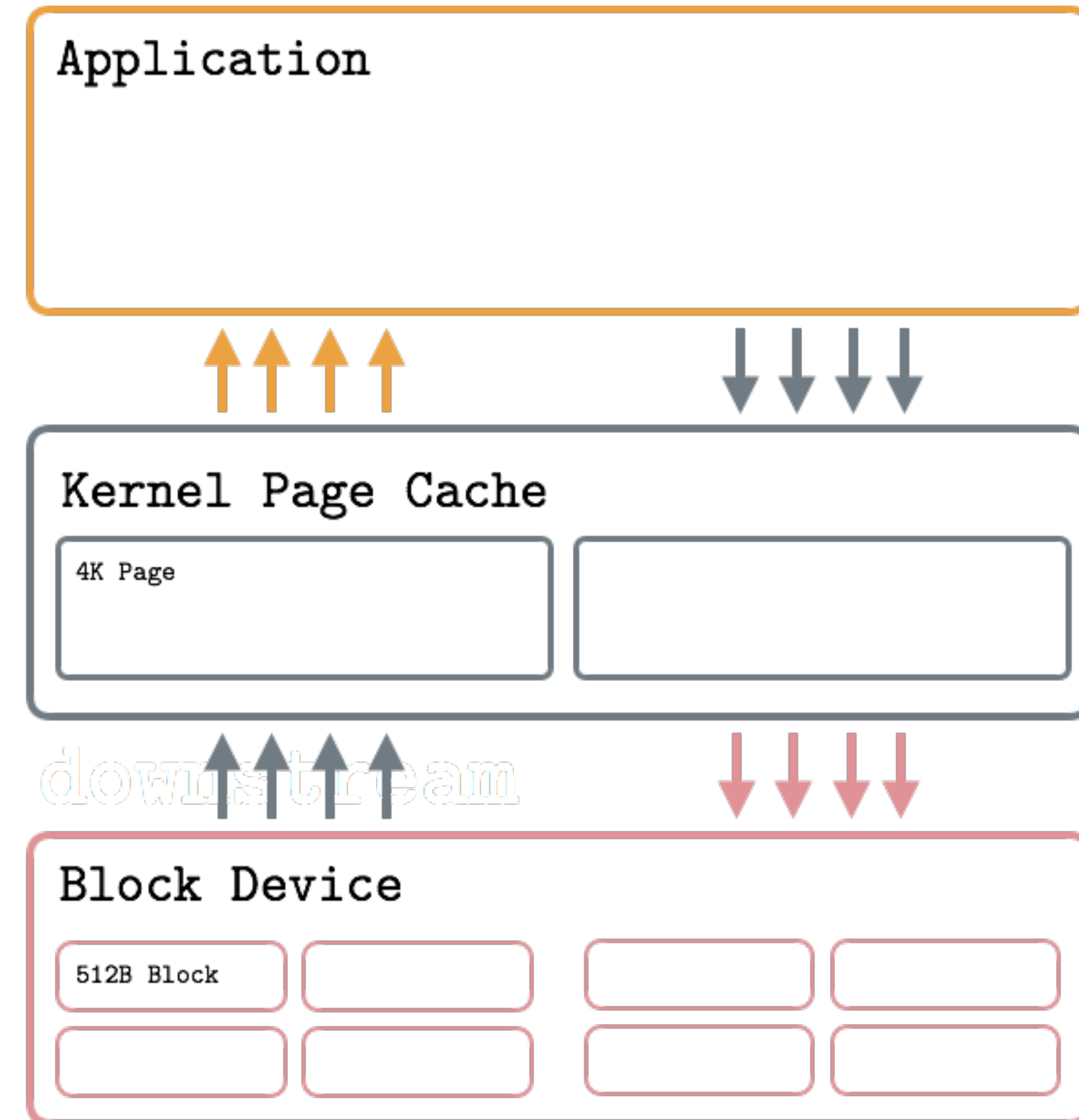


<http://www.ilinuxkernel.com/files/Linux.Generic.Block.Layer.pdf>



# Buffered IO vs Direct IO

- Buffered IO
  - read (LRU)
  - flush or writeback
- Direct IO
  - O\_DIRECT
  - Block Alignment
  - WAL(write ahead log)

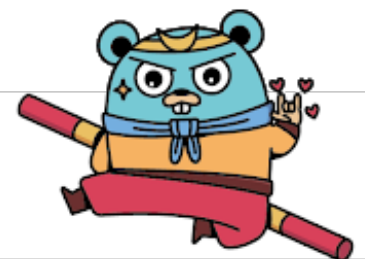


<https://medium.com/databass/on-disk-io-part-1-flavours-of-io-8e1ace1de017>



# Page Cache

- **Temporal Locality**, recently accessed pages will be accessed again at some point in nearest future.
- **Spatial Locality**, the elements physically located nearby have a good chance of being located close to each other.



# IO Scheduler

- IO Queue
- elevator: Deadline, CFQ, Noop
- how to

# 查看 sda 的 scheduler

```
cat /sys/block/sda/queue/scheduler
```

# 设置 sdb 的 scheduler 为 noop 算法

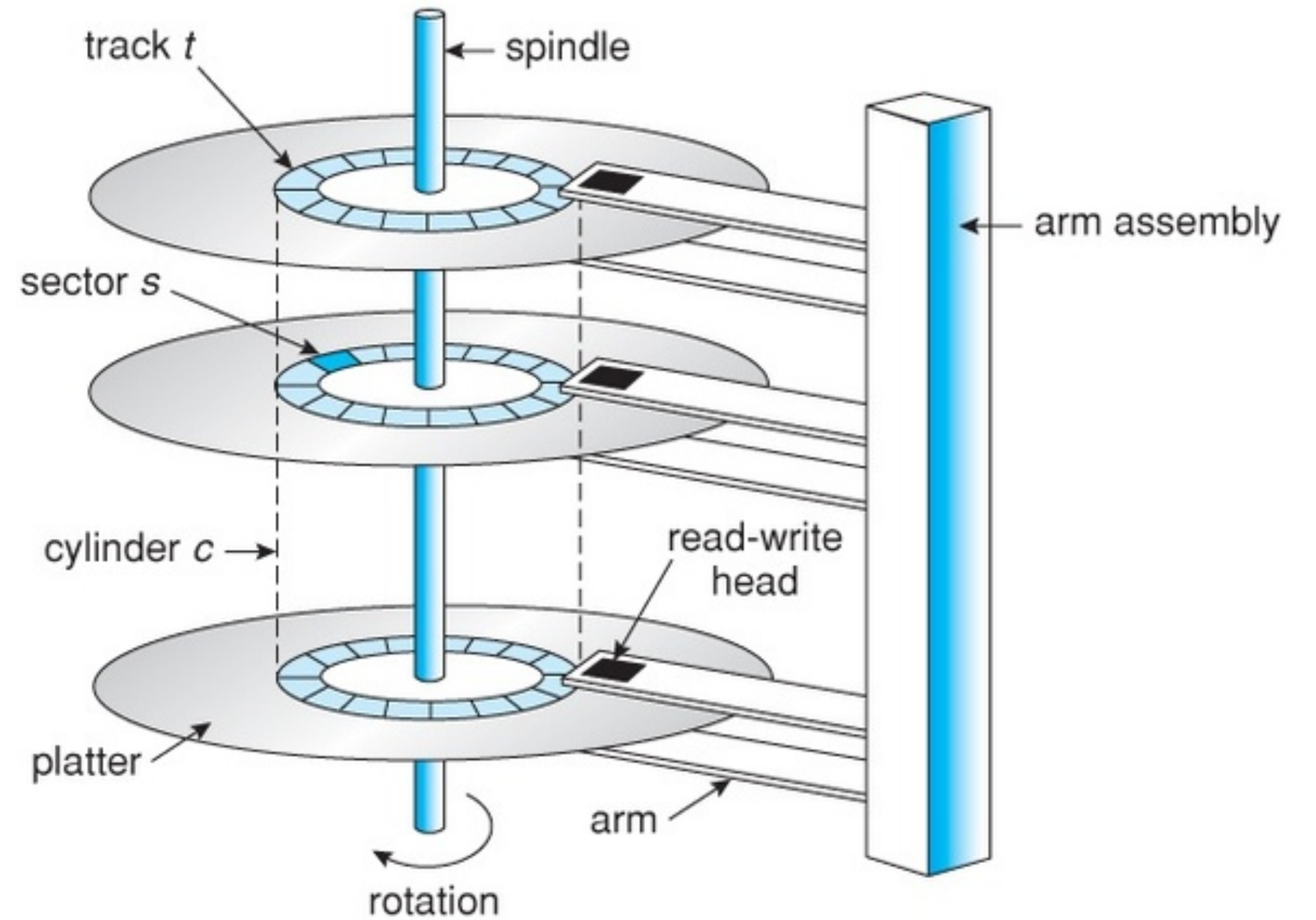
```
echo noop > /sys/block/sdb/queue/scheduler
```





# 磁盘

- 顺序读写快
- 随机读写慢



<https://qph.fs.quoracdn.net/main-qimg-f0bf80ac2007d1f39e8851a5e25b870c>



# 回顾两个问题

- 之前的压测结果中，为什么刚开始的一段时间响应时间特别低？
- 为什么之后响应时间变高？并且波动幅度很大



# 参数调优

- Virtual Memory
  - `/proc/sys/vm/dirty_ratio`
  - `/proc/sys/vm/dirty_background_ratio`
  - `/proc/sys/vm/dirty_writeback_centiseecs`
  - `/proc/sys/vm/dirty_expire_centiseecs`
- IO
  - `cat /sys/block/sda/queue/scheduler`
  - `/sys/block/sda/queue/read_ahead_kb`
  - `nr_requests`
  - `max_sectors_kb`

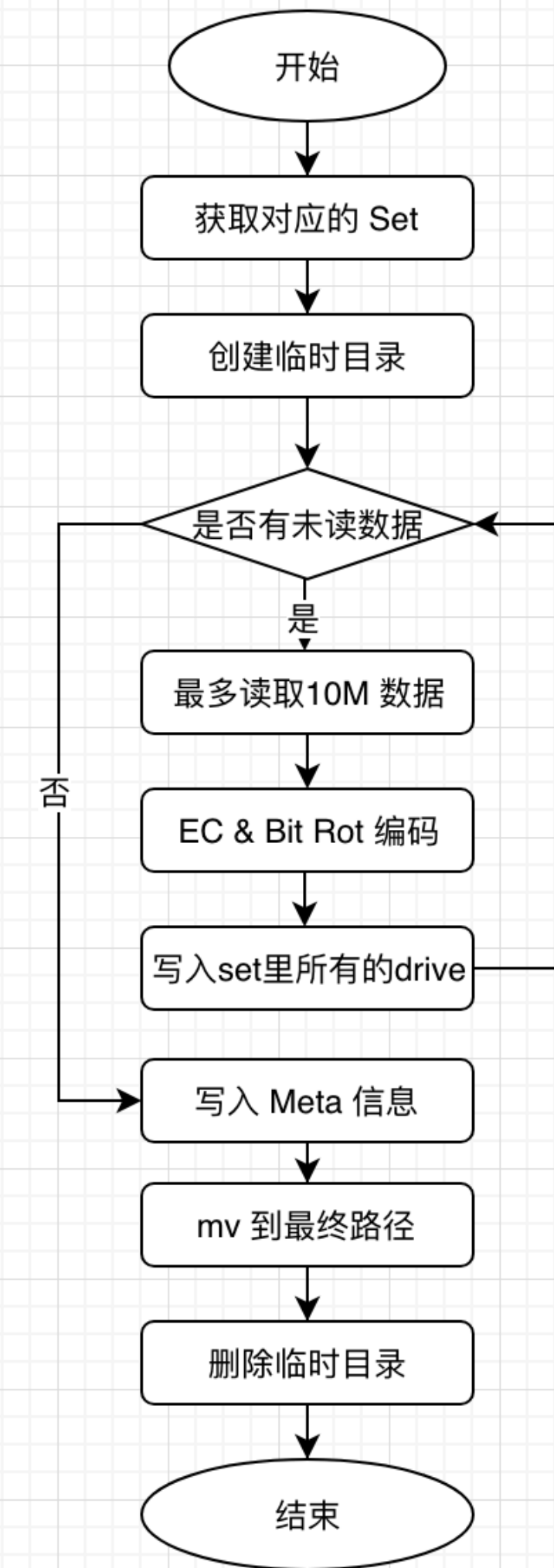


# 深度优化 Minio



# 再回顾 MINIO 写流程

- 随机写
- 写次数放大





# 优化方案 — 小文件合并大文件

- fallocate 预分配空间
- 顺序写入
- 小文件的索引 (volume id, offset, size) 存 leveledb



# FIO 测试效果

# 测试随机写性能

```
fio -name sequentialwrite -rw=write -bs=10k -runtime=100 -iodepth 128 \  
-ioengine libaio -direct=1 --size=200G --randrepeat=1 \  
--directory=/data10/fio --nrfiles=1 --openfiles=1 -group_reporting
```

# 测试顺序写性能

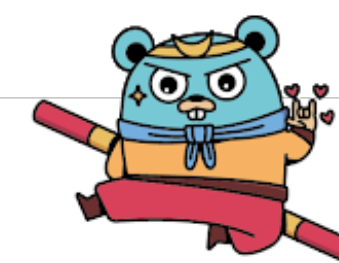
```
fio -name randomwrite -rw=randwrite -bs=10k -runtime=100 -iodepth 128 \  
-ioengine libaio -direct=1 --size=200G --randrepeat=1 \  
--directory=/data10/fio --nrfiles=1 --openfiles=1 --group_reporting
```



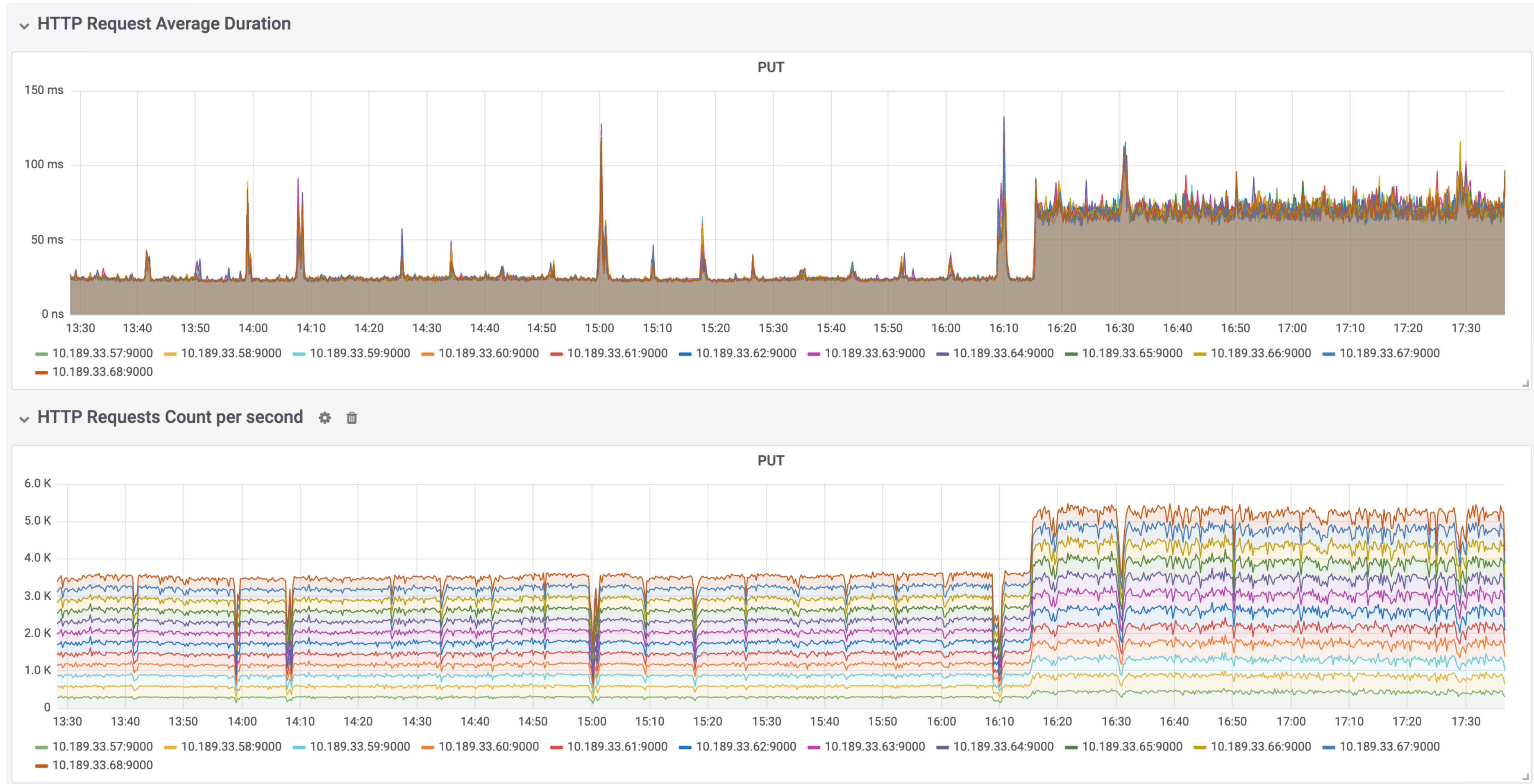
# 再谈 Golang IO

- 不同 IO 介质之间数据互传
  - req.Body
  - \*os.File
  - io.Copy

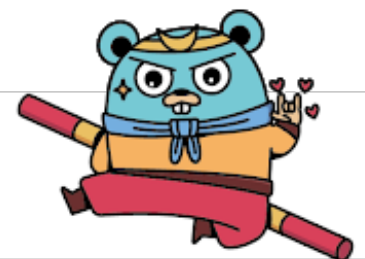
```
func (v *Volume) WriteAll(key string, size int64, r io.Reader) error {  
    // ...  
    v.lock.Lock()  
    defer v.lock.Unlock()  
    _, err := io.Copy(v.file, r)  
    if err != nil {  
        return err  
    }  
    // ...  
    return nil  
}
```



# 性能压测



12 个节点，每个节点 12 块盘，写 64KB





# 测试读性能

## 分析业务场景

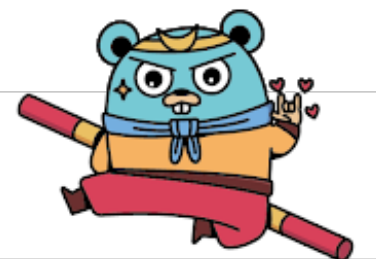
- 20% 数据为刚写入，大部分在 Page Cache 中
- 80% 数据随机分布在硬盘上

## 制定测试方案

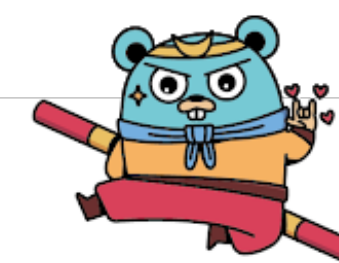
- 利用 MINIO 消息机制，将所上传的对象的列表写入 kafka
- 写入一段时间之后，再根据 kafka 中的文件列表进行读压测



你猜性能如何？

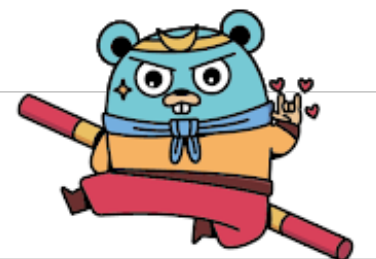


性能高的离谱!



# 事实本该如此

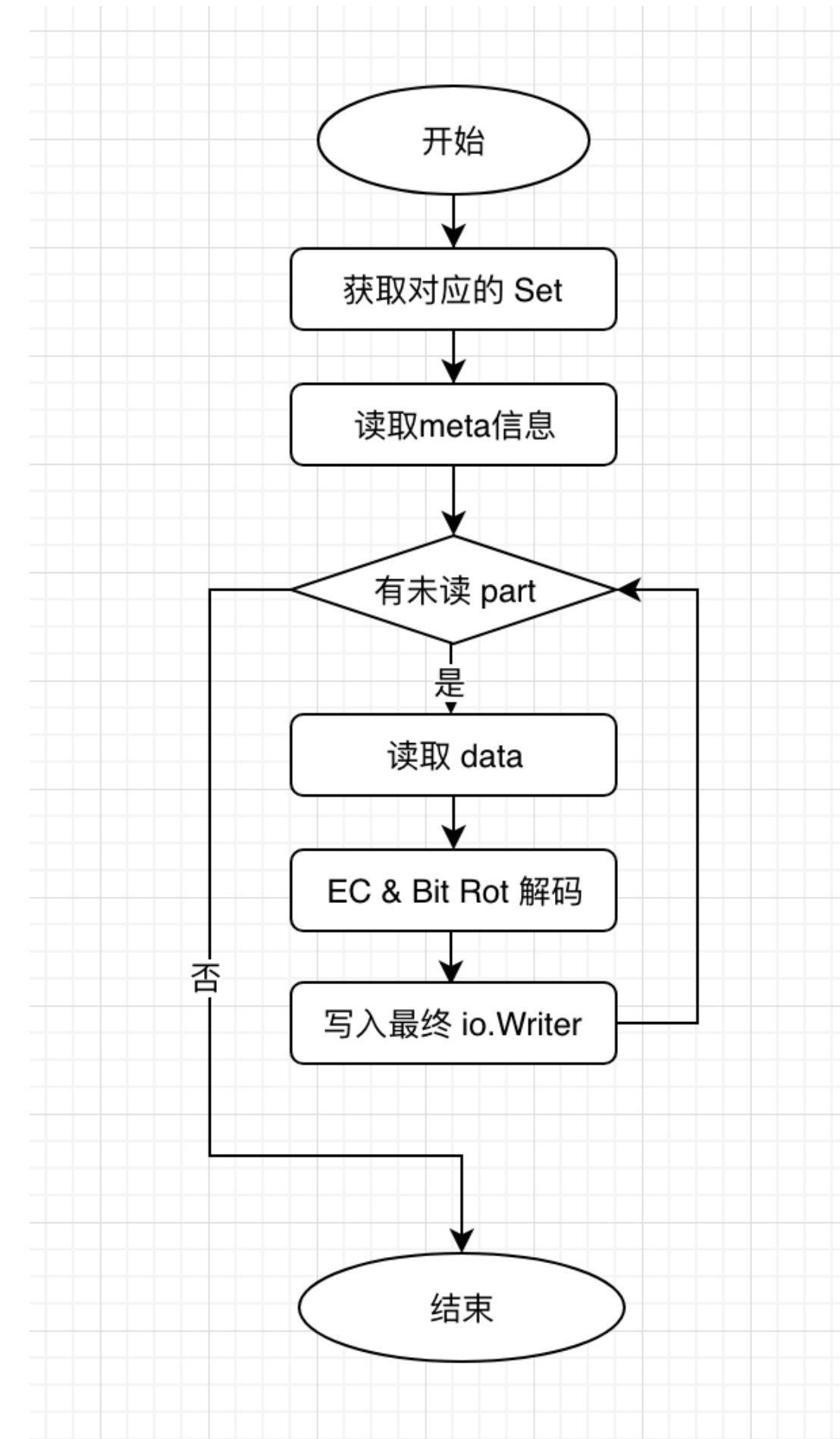
- 因为数据顺序写入，其文件列表发到kafka的时候也是顺序的，所以命中了顺序读的策略
- 加入随机因子之后，读性能急剧下降，符合预期





# 再回顾 MINIO 读流程

- 从 N 个节点上读取 meta 信息
- 至少从 N/2 个节点读取 data
- 每一次读取的时间包括两个部分
  - 从 levelDB 读取索引
  - 从大文件中读取对应的数据
- 并发读的最终时间取决于最长的一个



# 优化方案

- 将 Meta 信息压缩后直接存入 levelDB, 减少一组 IO
- 将 levelDB 存在性能更高的 SSD
- 降低单个 Set 的 Drive 数量



# 压测结果



12 个节点，每个节点 12 块盘，混合读写 64KB





# 路漫漫其修远兮

- 如何解决尖峰问题
- 长时间压测，确保集群容量很大以后仍然可以保持不错的性能
- levelDB 是否适合这个场景？





# 总结



# 总结

- 再回顾 MINIO
- 磨刀不误砍柴工



谢谢

