
fpchen note

发布 *V0.1.0*

方平

2021 年 10 月 08 日

1	常用命令脚本	1
1.1	Linux 命令	1
1.2	git 操作	3
1.3	SSH 远程操作	4
2	BlockChain	7
2.1	Consensus	7
2.2	ETH	60
2.3	BTC	71
2.4	Crypto	93
2.5	Wallet	98
2.6	Cosmos	109
2.7	Rust	114
3	算法与数据结构	121
3.1	链表	121
3.2	Tree	123
4	数据存储与消费	125
4.1	Mysql	125
4.2	Redis	135
4.3	Pulsar	147
5	Golang 开发	151
5.1	基础知识	151
5.2	开发常用	157
5.3	同步原语与锁	160
5.4	RPC	162
5.5	Protobuf3	163
6	服务部署运行	177
6.1	Docker	177
6.2	Aplo	187
6.3	Eureka	208
6.4	Log	222
6.5	开发踩坑	224

7	网络通信	225
7.1	网络协议	226
7.2	IO 通信	230
7.3	通信框架	234
7.4	Nginx	236
8	设计模式	241
8.1	类图关系	241
8.2	PUML 图形	242
8.3	单例模式	242

1.1 Linux 命令

1.1.1 进程常用命令

```
# 后台不停止运行test脚本，并将脚本日志写到test.log
nohup ./test.sh > test.log &

# 搜索当前运行的进程
ps aux | grep start

# 查看端口占用
lsof -i:端口号

#列出所有这个进程关联的文件句柄
lsof -p $pid

# 杀死进程
kill -9 ppid

# 进程监控工具
top
```

1.1.2 文件常用命令

```
# 查看日志中最近出现的字符str
tail -f test.log | grep str

# 使用通配符*(0或者任意多个)。表示在/etc目录下查找文件名中含有字符串'srm'的文件
find /etc -name '*srm*'

# 修改文件目录权限
chmod -R 755 /upload
# -rwxr-xr-x (755) 只有所有者才有读, 写, 执行的权限, 组群和其他人只有读和执行的权限

# 更改文件目录所属用户
chown -R $USER /tmp/

# 查看系统用户:
cat /etc/passwd

# 文件比较异常
cat origin_text_file1 | sort > file1.out
cat origin_text_file2 | sort > file2.out
vimdiff file2.out file2.out
```

1.1.3 网络相关命令

```
# 查询、配置网络卡与 IP 网域等相关参数
ifconfig

# 域名解析命令,进行域名与IP地址解析
nslookup [主机名或者IP]
> server

# 查看连接的网络服务
netstat -an | grep ESTABLISHED

# 查看路由表(网关)
netstat -rn

cat /etc/hosts:      域名到 IP 地址的映射。
cat /etc/networks:  网络名称到 IP 地址的映射。
cat /etc/resolv.conf: DNS域名地址
cat /etc/protocols: 协议名称到协议编号的映射。
cat /etc/services:  TCP/UDP 服务名称到端口号的映射。
```

1.2 git 操作

1.2.1 本地分支回退版本

```
git reset --hard 22f8aae 。22f8aae 为某次提交的提交号。
```

```
git reset --hard HEAD~3 （回退3次提交）
```

```
--
```

↪hard: 本地的源码和本地未提交的源码都会回退到某个版本, 包括commit内容, 和git自己对代码的索引都会回退
↪any local changes will be lost.

```
--
```

↪soft: 保留源码, 只能回退到commit信息到某个版本, 不涉及到index的回退, 如果还需要提交, 直接commit即可

↪--mixed: 会保留源码, 只是将git commit和index信息回退到某个版本。

1.2.2 远程分支回退版本

方法一, 先 git reset 回滚到本地, 然后再强制 push 到远程。

```
git push -u origin master -f          origin: 远程仓库名  master: 分支名称  -
```

↪f: force, 意为强制、强行

1.2.3 删除提交分支

```
# 删除当前提交分支
```

```
git remote remove origin
```

```
# 添加新的分支
```

```
git remote add origin https://gitee.com/kingCould/HelloWord.git
```

1.2.4 分支文件冲突

```
# 保存未修改文件, 添加备注, 方便查找
```

```
git stash save "save message"
```

```
# 应用并删除第二个存储
```

```
git stash pop stash@{1}
```

```
# 应用第二个存储, 但不会把存储从存储列表中删除
```

```
git stash apply stash@{1}
```

```
# 从列表中删除第二个存储
```

```
git stash drop stash@{1}
```

```
# 先切换到主分支 (或即将被覆盖的分支)
```

```
$ git checkout master
```

```
# 将 test 分支覆盖到当前分支
```

```
$ git merge test
```

(下页继续)

(续上页)

```
# 删除之前的分支 (已经没用了)
$ git branch -d <branch-name>
```

1.2.5 git lfs 突破大文件

```
git init

git lfs install

git lfs track "*.pdf"

git lfs track "*.epub"

cat .gitattributes

git add .gitattributes

git commit -m "track *.pdf *epub files using Git LFS"

git add .

git commit -m "first add books"

git lfs ls-files
```

1.3 SSH 远程操作

1.3.1 远程连接

连接远程服务器

```
ssh root@192.168.163.129
```

iterm 添加连接文件

- 可以在 ~/.ssh/ 下面写一个脚本, ~/.ssh/server-test

```
#!/usr/bin/expect -f
set user root
set host 192.168.1.110
set password 123456
set timeout -1

spawn ssh $user@$host
expect "*assword:*"
send "$password\r"
interact
expect eof
```

- iTerm2 的 Profiles 新建一个 profile, Command 里填入 expect ~/.ssh/server-test

1.3.2 文件传输

1、上传本地文件到服务器

```
scp /path/filename username@servername:/path/
```

例如 `scp /var/www/test.php root@192.168.0.101:/var/www/` 把本机 `/var/www/` 目录下的 `test.php` 文件上传到 `192.168.0.101` 这台服务器上的 `/var/www/` 目录中

2、从服务器上下载文件

下载文件我们经常使用 `wget`，但是如果没有 `http` 服务，如何从服务器上下载文件呢？

```
scp username@servername:/path/filename /var/www/local_dir (本地目录)
```

例如 `scp root@192.168.0.101:/var/www/test.txt /var/www/local_dir` 把 `192.168.0.101` 上的 `/var/www/test.txt` 的文件下载到 `/var/www/local_dir` (本地目录)

3、从服务器下载整个目录

```
scp -r username@servername:/var/www/remote_dir/ (远程目录) /var/www/local_
↳dir (本地目录)
```

例如：`scp -r root@192.168.0.101:/var/www/test /var/www/`

4、上传目录到服务器

```
scp -r local_dir username@servername:remote_dir
```

例如：`scp -r test root@192.168.0.101:/var/www/` 把当前目录下的 `test` 目录上传到服务器的 `/var/www/` 目录

1.3.3 Bash 脚本连接

- 执行多条命令

```
#!/bin/bash
ssh user@remoteNode << eeoff
cd /home
touch abcdefg.txt
exit
eeoff
echo done!
```

```
#!/bin/bash
#变量定义
ip_array=("192.168.1.1" "192.168.1.2" "192.168.1.3")
user="test1"
remote_cmd="/home/test/1.sh"
```

(下页继续)

(续上页)

```
#本地通过ssh执行远程服务器的脚本
for ip in ${ip_array[*]}
do
    if [ $ip = "192.168.1.1" ]; then
        port="7777"
    else
        port="22"
    fi
    ssh -t -p $port $user@$ip "remote_cmd"
done
```

2.1 Consensus

2.1.1 ACID 原则与多阶段算法

ACID 原则

具体来说，ACID 原则描述了分布式数据库需要满足的一致性需求，同时允许付出可用性的代价。

- **Atomicity**(原子性): 每次事务是原子的，事务包含的所有操作要么全部成功，要么全部不执行。一旦有操作失败，则需要回退状态到执行事务之前；
- **Consistency** (一致性): 数据库的状态在事务执行前后的状态是一致的和完整的，无中间状态。即只能处于成功事务提交后的状态；
- **Isolation** (隔离性): 各种事务可以并发执行，但彼此之间互相不影响。按照标准 SQL 规范，从弱到强可以分为未授权读取、授权读取、可重复读取和串行化四种隔离等级；
- **Durability** (持久性): 状态的改变是持久的，不会失效。一旦某个事务提交，则它造成的状态变更就是永久性的。

两阶段提交算法 (Two-phase Commit, 2PC)

其基本思想十分简单，既然在分布式场景下，直接提交事务可能出现各种故障和冲突，那么可将其分解为预提交和正式提交两个阶段，规避冲突的风险。

- **预提交**: 协调者 (Coordinator) 发起提交某个事务的申请，各参与执行者 (Participant) 需要尝试进行提交并反馈是否能完成；
- **正式提交**: 协调者如果得到所有执行者的成功答复，则发出正式提交请求。如果成功完成，则算法执行成功。

在此过程中任何步骤出现问题 (例如预提交阶段有执行者回复预计无法完成提交)，则需要回退。

两阶段提交算法因为其简单容易实现的优点，在关系型数据库等系统中被广泛应用。当然，其缺点也很明显。整个过程需要同步阻塞导致性能一般较差；同时存在单点问题，较坏情况下可能一直无法完成提交；另外可能产生数据不一致的情况（例如协调者和执行者在第二个阶段出现故障）。

三阶段提交算法 (Three-phase Commit, 3PC)。

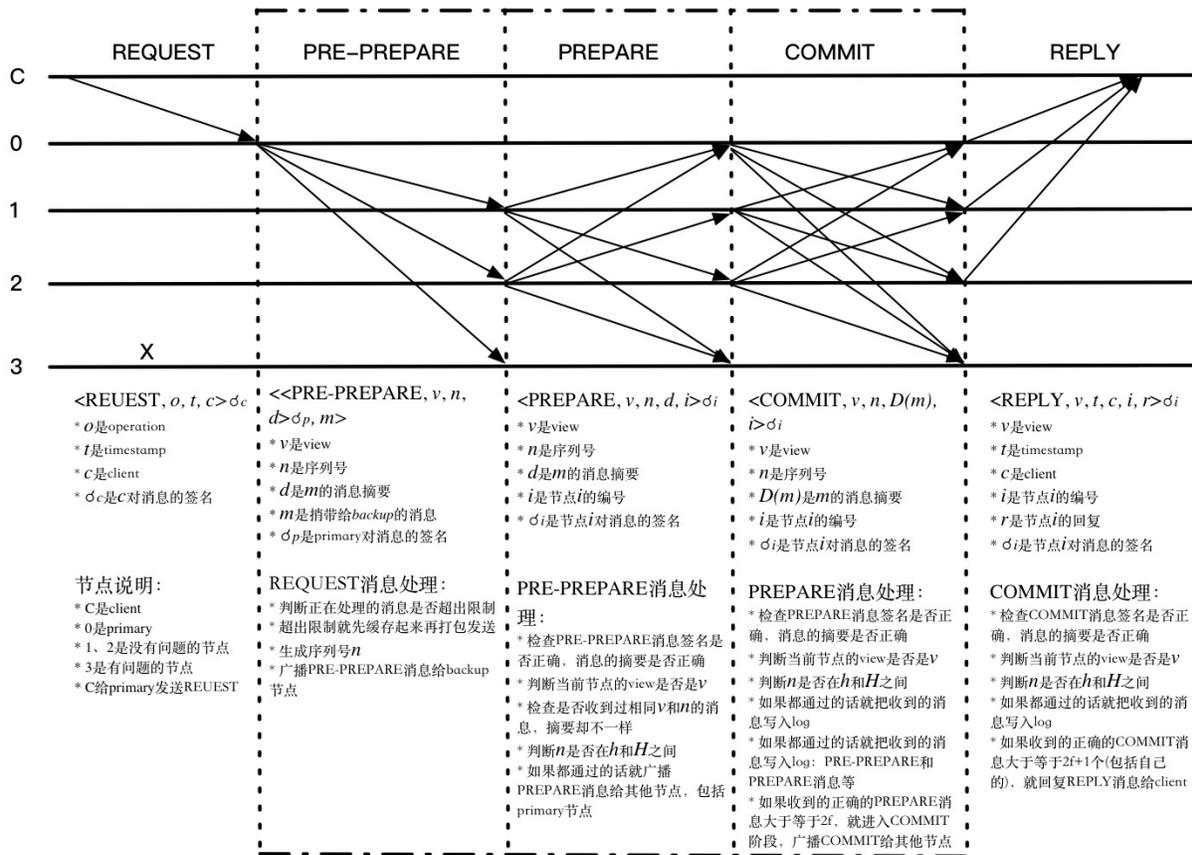
三阶段提交针对两阶段提交算法第一阶段中可能阻塞部分执行者的情况进行了优化。具体来说，将预提交阶段进一步拆成两个步骤：尝试预提交和预提交。

完整过程如下：

- 尝试预提交：协调者询问执行者是否能进行某个事务的提交。执行者需要返回答复，但无需执行提交。这就避免出现部分执行者被无效阻塞住的情况。
- 预提交：协调者检查收集到的答复，如果全部为真，则发起提交事务请求。各参与执行者 (Participant) 需要尝试进行提交并反馈是否能完成；
- 正式提交：协调者如果得到所有执行者的成功答复，则发出正式提交请求。如果成功完成，则算法执行成功。

其实，无论两阶段还是三阶段提交，都只是一定程度上缓解了提交冲突的问题，并无法一定保证系统的一致性。首个有效的算法是后来提出的 Paxos 算法。

3阶段协议



2.1.2 FLP and CAP

分布式系统的核心问题：一致性问题和共识问题引起不一致的因素：

- 节点间网络通信的不可靠, 消息延迟, 消息乱序, 内容错误.
- 节点处理时间无法保证: 结果可能错误, 或者节点宕机.

FLP 不可能原理

FLP 不可能原理：在网络可靠，但允许节点失效（即便只有一个）的最小化异步模型系统中，不存在一个可以解决一致性问题的确定性共识算法。（节点失效不可避免，所以不能共识）

FLP 不可能原理告诉大家不必浪费时间去追求完美的共识方案，而要根据实际情况设计可行的工程方案。

那么，退一步讲，在付出一些代价的情况下，共识能做到多好？

回答这一问题的是另一个很出名的原理：CAP 原理。

CAP 原理

CAP 原理：分布式计算系统不可能同时确保以下三个特性：强一致性（Consistency）、可用性（Availability）和分区容忍性（Partition），设计中往往需要弱化对某个特性的保证。

- **一致性：**所有节点在同一时刻能够看到相同的数据。
- **可用性：**在有限时间内，任何非失败节点都能应答请求；
- **分区容忍性：**分布式系统在遇到某节点或网络分区故障的时候，仍然能够对外提供满足一致性或可用性的服务。

一致性解读

- 对于客户端来说，一致性指的是并发访问时更新过的数据如何获取的问题。
- 从服务端来看，则是更新如何复制分布到整个系统，以保证数据最终一致。

规范来看，分布式系统达成一致的过程，应该满足：

- **可终止性（Termination）：**一致的结果在有限时间内能完成；意味着可以保障提供服务（Liveness）。
- **约同性（Agreement）：**不同节点最终完成决策的结果是相同的；
- **合法性（Validity）：**决策的结果必须是某个节点提出的提案。

CAP 原理认为，分布式系统最多只能保证三项特性中的两项特性。

比较直观地理解，当网络可能出现分区时候，系统是无法同时保证一致性和可用性的。要么，节点收到请求后因为没有得到其它节点的确认而不应答（牺牲可用性），要么节点只能应答非一致的结果（牺牲一致性）。

CAP 应用场景

既然 CAP 三种特性不可同时得到保障，则设计系统时候必然要弱化对某个特性的支持。

弱化一致性

对结果一致性不敏感的应用，可以允许在新版本上线后过一段时间才最终更新成功，期间不保证一致性。

例如网站静态页面内容、实时性较弱的查询类数据库等，简单分布式同步协议如 Gossip，以及 CouchDB、Cassandra 数据库等，都为此设计。

弱化可用性

对结果一致性很敏感的应用，例如银行取款机，当系统故障时候会拒绝服务。MongoDB、Redis、MapReduce 等为此设计。

Paxos、Raft 等共识算法，主要处理这种情况。在 Paxos 类算法中，可能存在着无法提供可用结果的情形，同时允许少数节点离线。

弱化分区容忍性

现实中，网络分区出现概率较小，但很难完全避免。

两阶段的提交算法，某些关系型数据库以及 ZooKeeper 主要考虑了这种设计。

实践中，网络可以通过双通道等机制增强可靠性，实现高稳定的网络通信。

共识算法

针对如上不稳定的情况，把出现故障（Crash 或 Fail-stop，即不响应）但不会伪造信息的情况称为“非拜占庭错误（Non-Byzantine Fault）”或“故障错误（Crash Fault）”；伪造信息恶意响应的情况称为“拜占庭错误”（Byzantine Fault），对应节点为拜占庭节点。显然，后者场景中因为存在“捣乱者”更难达成共识。

共识算法：

- BFT 类针对拜占庭错误有：PBFT 确定性算法、POW 工作量证明的概率算法等。
- CFT 类针对非拜占庭错误有：PAXOS 算法、Raft 算法（容错较好，处理快，保证不超过一半节点故障）

reference

常用共识算法总结（Paxos, Raft, PBFT, PoW, PoS, DPoS, Ripple）

[blockchain_guide](#)

2.1.3 Raft 算法

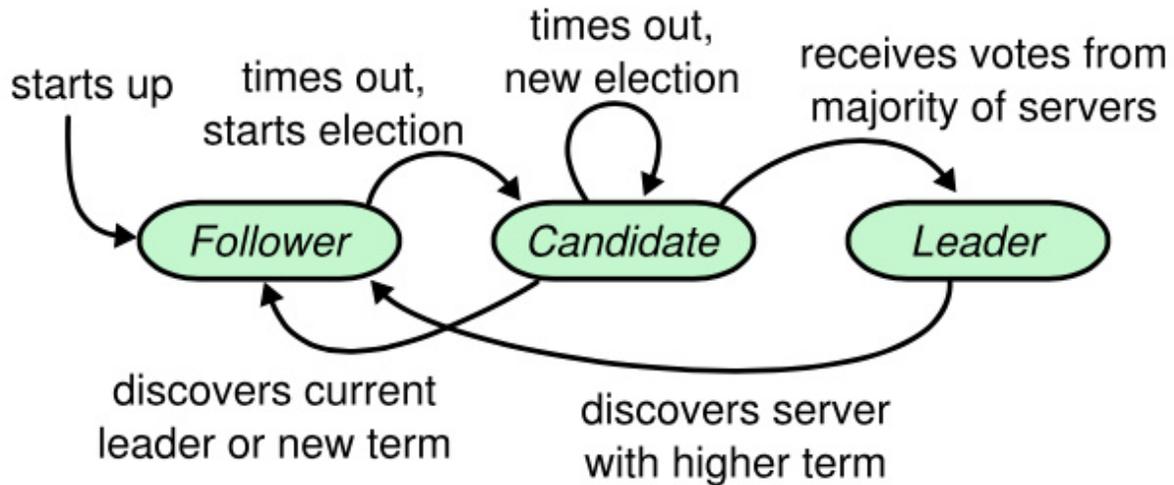
Raft 协议分为两个阶段

一、节点角色

角色定义

- 领导者 Leader
- 候选者 Candidate
- 追随者 Follower

角色变化



服务器状态。**Follower** 只响应来自其他服务器的请求。如果 **Follower** 接收不到消息，那么他就会变成 **Candidate** 并发起一次选举。获得集群中大多数选票的 **Candidate** 将成为 **Leader**。在一个任期内，**Leader** 一直都是 **Leader** 直到自己宕机了。

二、RPC 调用

AppendEntries RPC

由 **Leader** 负责调用来复制日志指令；也会用作 heartbeat

接收服务器需实现以下处理逻辑：

- 1、如果任期编号比本地小 $term < currentTerm$ ，则返回 `false`(5.1 节讨论)；
- 2、如果 **Follower** 在 `prevLogIndex` 索引处不包含 `prevLogTerm` 任期条目，则返回 `false`(5.3 节讨论)；
- 3、如果 **Follower** 有条目与新条目冲突(索引相同任期不同)，则删除该条目以及所有后续条目(5.3 节讨论)；

- 4、追加所有未在日志中的新条目；
- 5、如果 `leaderCommit > commitIndex`, 则将 `commitIndex` 更新为 `leaderCommit` 以及最后新条目索引两者的较小值；

RequestVote RPC

由候选人负责调用用来征集选票

接收服务器需实现以下处理逻辑：

- 1、如果 `term < currentTerm` 返回 `false` (5.2 节)
- 2、如果 `votedFor` 为空或者为 `candidateId`, 并且 **Candidate** 的日志至少和自己一样新, 那么就投票给他 (5.2 节, 5.4 节)

所有服务器需遵守的规则：

所有服务器：

- 如果 `commitIndex > lastApplied`, 那么就 `lastApplied` 加一, 并把 `log[lastApplied]` 条目应用到状态机中 (5.3 节)
- 如果接收到的 RPC 请求或响应中, 任期号 `T > currentTerm`, 那么就令 `currentTerm` 等于 `T`, 并切换状态为 **Follower** (5.1 节)

Follower (5.2 节)：

- 响应来自 **Candidate** 和 **Leader** 的请求
- 如果在超过选举超时时间的情况之前都没有收到 **Leader** 的心跳, 或者是 **Candidate** 请求投票的, 就自己变成 **Candidate**

Candidate (5.2 节)：

- 在转变成 **Candidate** 后就立即开始选举过程
 - 自增当前的任期号 (`currentTerm`)
 - 给自己投票
 - 重置选举超时计时器
 - 发送 **RequestVote RPC** 给其他所有服务器
- 如果接收到大多数服务器的选票, 那么就变成 **Leader**
- 如果接收到来自新的 **Leader** 的 **AppendEntries RPC**, 转变成 **Follower**
- 如果选举过程超时, 再次发起一轮选举

Leader:

- 一旦成为 **Leader**: 发送空的 **AppendEntries RPC** (心跳) 给其他所有的服务器; 在一定的空余时间之后不停的重复发送, 以阻止 **Follower** 超时 (5.2 节)
- 如果接收到来自客户端的请求: 附加条目到本地日志中, 在条目被应用到状态机后响应客户端 (5.3 节)
- 如果对于一个 **Follower**, 最后日志条目的索引值大于等于 `nextIndex`, 即 `lastIndex >= nextIndex`, 则发起 **AppendEntries RPC** 追加从 `nextIndex` 开始的日志条目:
 - 如果成功: 更新相应 **Follower** 的 `nextIndex` 和 `matchIndex`
 - 如果因为日志不一致而失败, 减少 `nextIndex` 重试, 直到成功。
- 如果存在一个满足 $N > \text{commitIndex}$ 的 N , 并且大多数 **Follower** 的 `matchIndex[i] ≥ N` 成立, 并且 `log[N].term == currentTerm` 成立, 那么令 `commitIndex` 等于这个 N (5.3 和 5.4 节)

阶段操作

选举阶段

Raft 采用 **心跳机制** 来触发 **领袖选举**。服务器启动后, 开始以 **Follower** 角色运行。只要它不断收到来自 **Leader** 或者 **Candidate** 的 RPC 请求, 便保持 **Follower** 状态不变。**Leader** 发送周期性心跳 (不带任何条目的 **AppendEntries RPC** 请求) 给所有 **Follower**, 以保持领导权。如果 **Follower** 超过一定时间 (选举超时时间) 没有收到任何通讯, 它便假设当前没有可见的 **Leader**, 进而发起新选举。

为发起选举, **Follower** 自增自身任期, 并转换为 **Candidate**。随后它为自己投票, 同时向其他服务器发起 **RequestVote RPC** 请求。**Candidate** 持续这个状态直到发生以下三种情形之一:

- **Candidate** 赢得选举, 成为 **Leader**, 随后它向其他所有服务器发送心跳信息, 建立领导权并阻止新的选举。
- 另一台机器以 **Leader** 身份连接 **Candidate**, 如果 **Leader** 领袖任期 `term` 至少与 **Candidate** 一样大, **Candidate** 便认为 **Leader** 合法, 并重回 **Follower** 状态。相反, **Candidate** 将拒绝 RPC 请求并继续保持 **Candidate** 状态。
- 超过设定时间但未产生获胜者, 开启一次新的选举——自增任期并开始下一轮 **RequestVote** 请求。

日志复制

Leader 被选举出来之后, 开始服务客户端请求。每个客户端请求包含一个可以被复制状态机执行的命令。**Leader** 将命令最为新条目追加到本地日志, 然后并行发起 **AppendEntries RPC** 请求往其他机器复制该条目。一旦条目安全完成复制 (已复制到过半数机器), **Leader** 将把条目应用到自己的状态机并将执行结果告知客户端。如果 **Follower** 节点宕机、响应缓慢或者网络丢包, **Leader** 将不断重试 **AppendEntries RPC** 请求 (就算已经响应客户端), 直到该节点日志同步。

解决问题

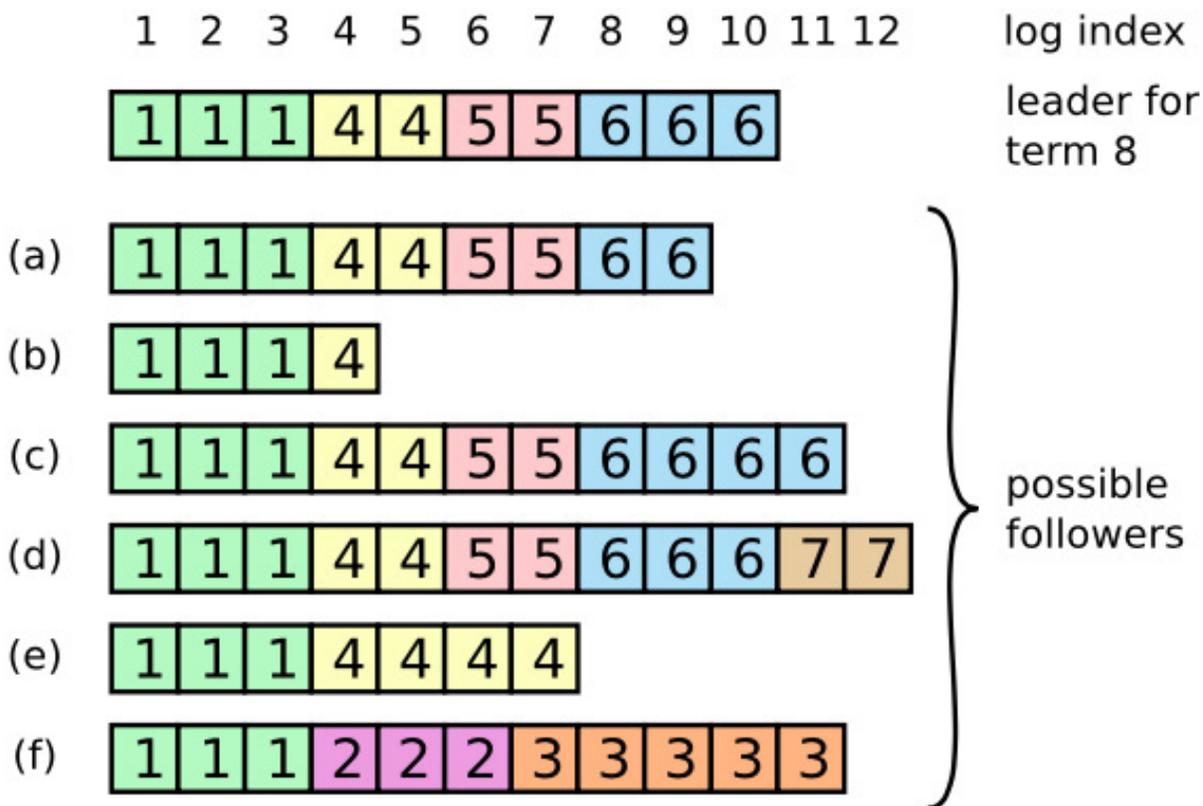
日志冲突解决

一致性检查

在发送 **AppendEntries RPC** 的时候，**Leader** 顺便带上前一日志条目的索引以及任期编号。**Follower** 节点如果找不到对应的条目，它将拒绝新条目。

一致性检查就像一个归纳步骤：一开始空的日志状态肯定是满足日志匹配特性的，然后一致性检查保护了日志匹配特性当日志扩展的时候。因此，每当 **AppendEntries RPC** 返回成功时，领导人就知道跟随者的日志一定是和自己相同的了。

不一致，Leader 强制 Follower 直接复制自己的日志



当顶部 **Leader** 节点启动后，从属节点日志状态存在各种可能：a-f。每个格子表示一个日志条目；数字代表任期。**Follower** 节点可能缺失部分日志条目 (a-b)；也可能包含多余的未提交日志；或者同时出现 (e-f)。举个例子，情景 f 可能是这样产生的：该服务器是任期 2 的 **Leader**，追加了几个条目，但是还没提交就宕机了；它迅速重启，在任期 3 继续担任 **Leader**，又追加了若干条目；在新条目提交前，服务器再次宕机，错过了接下来几个任期。

为了让 **Follower** 日志与自己保持一致，**Leader** 需要找到双方最后一个匹配的条目，从该点开始删除 **Follower** 所有不一致条目，并发送所有后续条目。这些操作根据 **AppendEntries RPC** 返回的一致性检查结果视情况执行。**Leader** 为每台 **Follower** 维护一个 `nextIndex` 变量，保存下一个发送给 **Follower** 的日志条目编号。当 **Leader** 刚开始服务时，将 `nextIndex` 设置为自己日志中下一个条目的索引（索引 11）。如果 **Follower** 日志与领袖不

一致，下一次 **AppendEntries RPC** 请求一致性检查将失败。请求被拒绝后，领袖将降低 `nextIndex` 并重试 **AppendEntries RPC** 请求。最终 `nextIndex` 将到达一个双方日志都匹配的点 (还可以用二分查找加快这个过程)。这时，**AppendEntries** 请求将成功，请求将删除属下所有有冲突条目并从领袖日志追加新条目 (如有)。一旦 **AppendEntries RPC** 请求成功，意味着属下日志与领袖一致，任期内后续时间也是如此。

选举约束

任何依赖领袖的共识算法，**Leader** 必须拥有所有已提交日志条目。

Raft 采用一种更简化的做法，选举保证新 **Leader** 必须包含先前任期所有已提交条目，无须传输缺失条目。这意味着，日志条目只有一个流向，从 **Leader** 流向 **Follower**，**Leader** 不会覆盖日志中已存在的条目。

Raft 通过选举过程阻止 **Candidate** 赢得选举，除非它拥有所有已提交日志。**Candidate** 想赢得选举必须与集群内大多数机器通讯，这意味着每个已提交条目必须出现在这些机器中至少一台。只要 **Candidate** 日志至少与其他大多数机器一样新 (`up-to-date`)，它便一定拥有所有已提交条目。**RequestVote RPC** 请求实现这个约束：该请求带有候选人日志信息，其他节点发现自己日志更新 (即更加新，`more up-to-date`) 则拒绝投票。

Raft 以日志最后条目的索引以及任期判断哪个日志更新。如果最后条目任期不同，则任期靠后的更新；如果最后条目任期相同索引不同，则索引大 (日志更长) 的更新。

2.1.4 PBFT 算法

raft 和 pbft 的最大容错节点数

pbft 算法的最大容错节点数量是 $(n-1) / 3$

raft 算法的最大容错节点数量是 $(n-1) / 2$

对于 **pbft** 算法，因为 **pbft** 算法的除了需要支持容错故障节点之外，还需要支持容错作恶节点。假设集群节点数为 N ，有问题的节点为 f 。有问题的节点中，可以既是故障节点，也可以是作恶节点，或者只是故障节点或者只是作恶节点。

- 第一种情况， f 个有问题节点既是故障节点，又是作恶节点，那么根据少数服从多数的原则，集群里正常节点只需要比 f 个节点再多一个节点，即 $f+1$ 个节点，确节点的数量就会比故障节点数量多，那么集群就能达成共识。也就是说这种情况支持的最大容错节点数量是 $(n-1) / 2$ 。
- 第二种情况，故障节点和作恶节点都是不同的节点。那么就会有 f 个问题节点和 f 个故障节点，当发现节点是问题节点后，会被集群排除在外，剩下 f 个故障节点，那么根据少数服从多数的原则，集群里正常节点只需要比 f 个节点再多一个节点，即 $f+1$ 个节点，确节点的数量就会比故障节点数量多，那么集群就能达成共识。所以，所有类型的节点数量加起来就是 $f+1$ 个正确节点， f 个故障节点和 f 个问题节点，即 $3f+1=n$ 。

算法基本流程

流程

pbft 算法的基本流程主要有以下四步：

- 客户端发送请求给主节点
- 主节点广播请求给其它节点，节点执行 **pbft** 算法的三阶段共识流程。
- 节点处理完三阶段流程后，返回消息给客户端。
- 客户端收到来自 $f+1$ 个节点的相同消息后，代表共识已经正确完成。

PREPARE 消息是收到 PRE-PREPARE 消息的节点发送出来的，primary 收到 REQUEST 消息后不会给自己发送 PRE-PREPARE 消息，也不会发送 PREPARE 消息，所以一个节点收到的消息数满足 $2f+1-1=2f$ 个就能满足没问题的节点数比有问题节点多了（包括自身节点）。COMMIT 阶段 primary 节点也会在收到 PREPARE 消息后发送 COMMIT 消息，所以收到的消息数满足 $2f+1$ 个就能满足没问题的节点数比有问题节点多了（包括自身节点）。

PRE-PREPARE 和 PREPARE 阶段保证了所有正常的节点对请求的处理顺序达成一致，它能够保证如果 $PREPARE(m, v, n, i)$ 是真的话， $PREPARE(m', v, n, j)$ 就一定是假的，其中 j 是任意一个正常节点的编号，只要 $D(m) \neq D(m')$ 。因为如果有 $3f+1$ 个节点，至少有 $f+1$ 个正常的节点发送了 PRE-PREPARE 和 PREPARE 消息，所以如果 $PREPARE(m', v, n, j)$ 是真的话，这些节点中就至少有一个节点发了不同的 PRE-PREPARE 或者 PREPARE 消息，这和它是正常的节点不一致。当然，还有一个假设是安全强度是足够的，能够保证 $m \neq m'$ 时， $D(m) \neq D(m')$ ， $D(m)$ 是消息 m 的摘要。

确定好了每个请求的处理顺序，怎么能保证按照顺序执行呢？网络消息都是无序到达的，每个节点达成一致的顺序也是不一样的，有可能在某个节点上 n 比 $n-1$ 先达成一致。其实每个节点都会把 PRE-PREPARE、PREPARE 和 COMMIT 消息缓存起来，它们都会有一个状态来标识现在处理的情况，然后再按顺序处理。而且序列号 n 在不同 view 中也是连续的，所以 $n-1$ 处理完了，处理 n 就好了。

客户端 C 向主节点发起请求，主节点 0 收到客户端请求

```
<Request,o,t,c>σc
```

- σc 是 client 对请求消息的签名
- o 是 operation
- c 是 client

Pre-prepare 阶段

主节点收到 $\langle \text{Request}, o, t, c \rangle \sigma c$ 之后，生成序列号 n ，广播 Pre-prepare 消息

```
<<Pre-Prepare,v,n,d>σp,m>
```

- v 是 view
- d 是消息的摘要
- n 是序列号
- σp 是 primary 对消息的签名
- m 是发给 backup 节点的消息

节点收到 pre-prepare 消息后，会有两种选择，一种是接受，一种是不接受。什么时候才不接受主节点发来的 pre-prepare 消息呢？一种典型的情况就是如果一个节点接受到了一条 pre-prepare 消息，消息里的 v 和 n 在之前收到里的消息是曾经出现过的，但是 d 和 m 却和之前的消息不一致，或者请求编号不在高低水位之间，这时候就会拒绝请求。拒绝的逻辑就是主节点不会发送两条具有相同的 v 和 n ，但 d 和 m 却不同的消息。

Prepare 阶段

节点验证只要没有问题，就发送 prepare 消息（不涉及消息的真实性）

节点同意请求后会向其它节点发送 prepare 消息。这里要注意一点，同一时刻不是只有一个节点在进行这个过程，可能有 n 个节点也在进行这个过程。因此节点是有可能收到其它节点发送的 prepare 消息的。

```
<Prepare, v, n, d, i> $\sigma$ i
```

- i 是节点 i 的编号
- σi 是节点 i 对 prepare 消息的签名

节点收到了 prepare 消息之后，先判断 v, n 等数据签名没问题，如果没问题，就把收到的消息写入 log, 包括 (pre-prepare, prepare 消息)

在一定时间范围内，如果收到超过 $2f$ 个不同节点的 prepare 消息，就代表 prepare 阶段已经完成。

Commit 阶段

节点进入 commit 阶段。向其它节点广播 commit 消息，同理，这个过程可能是有 n 个节点也在进行的。因此可能会收到其它节点发过来的 commit 消息。

```
<Commit, v, n, D(m), i> $\sigma$ i
```

- $D(m)$ 是 m 的信息摘要

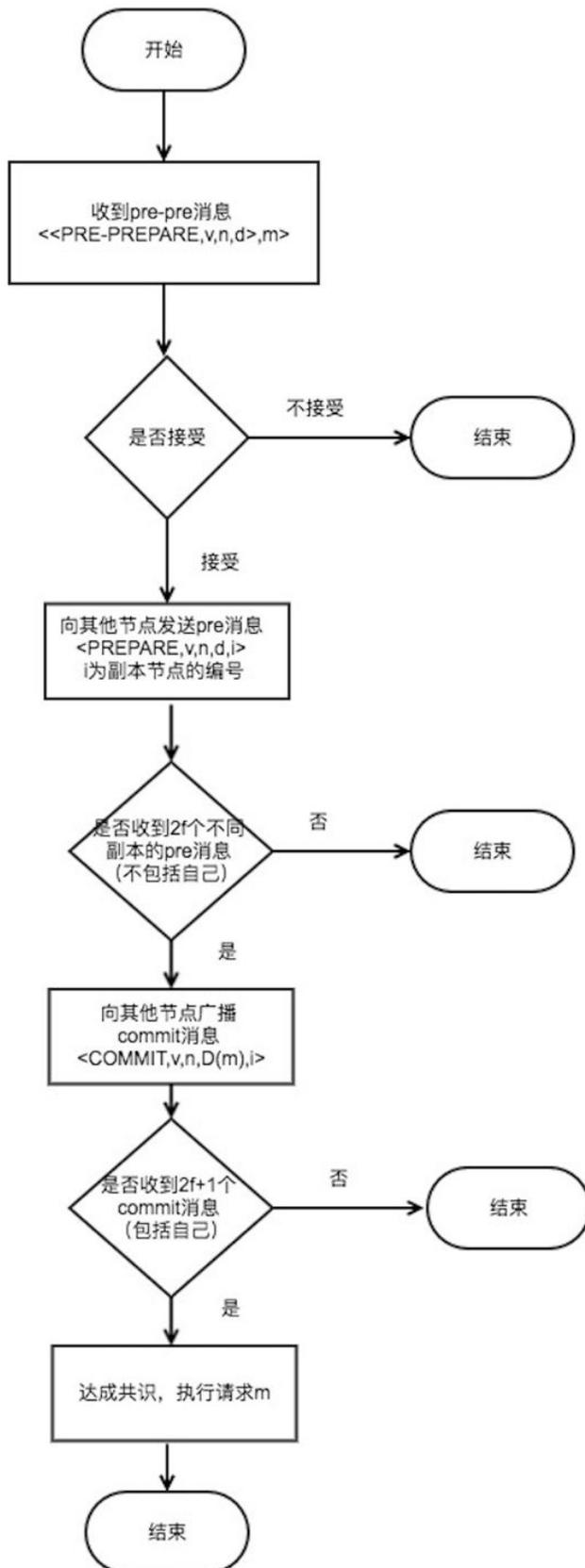
当收到 $2f+1$ 个 commit 消息后（包括自己），代表大多数节点已经进入 commit 阶段，这一阶段已经达成共识，于是节点就会执行请求，写入数据。

Reply

节点达成共识，会回复 Reply 消息给 client

```
<Reply, v, t, c, i, r> $\sigma$ i
```

- v 是 view
- t 是时间戳
- r 是节点 i 的回复
- σi 是节点对 Reply 消息的签名

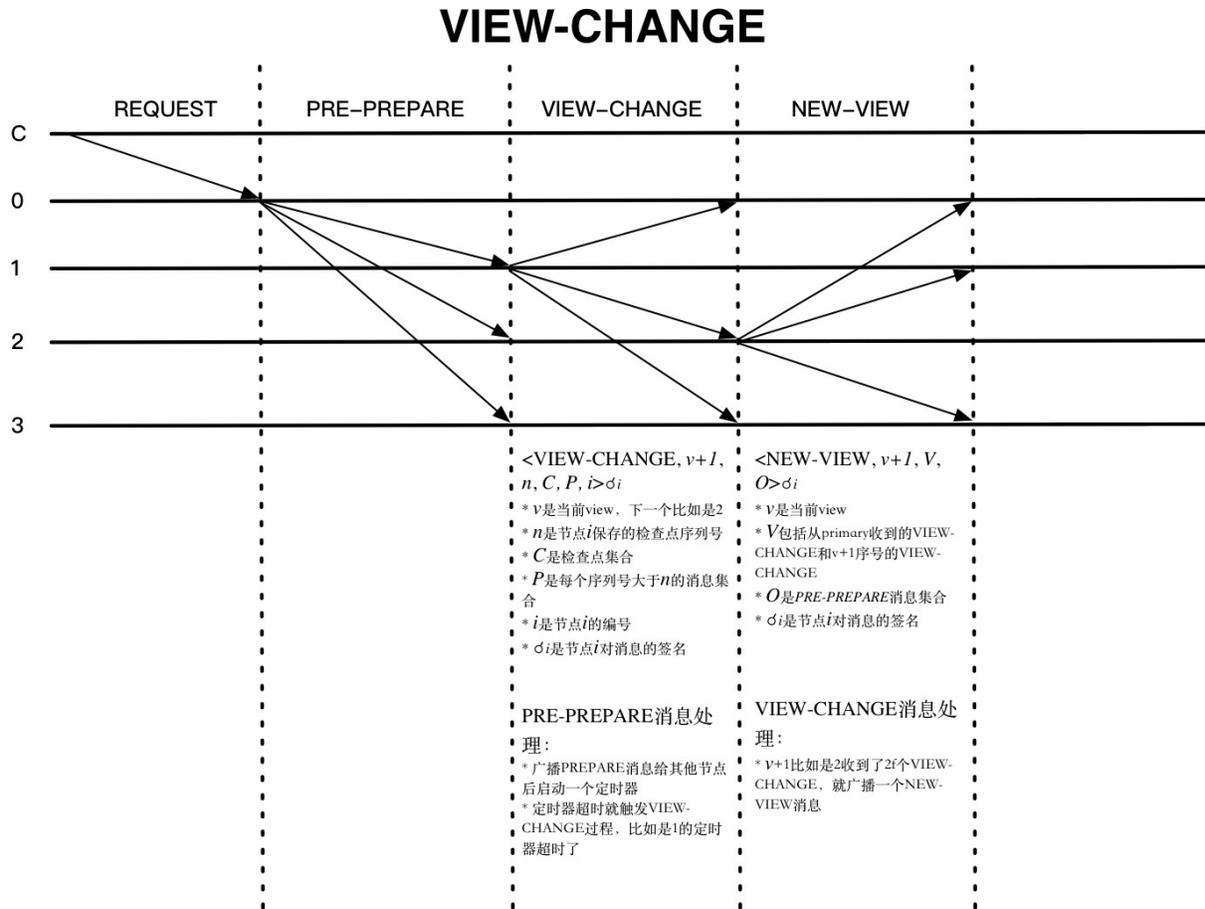


所有的流程如上图所示

ViewChange

ViewChange (视图更改) 事件

当主节点挂了（超时无响应）或者从节点集体认为主节点是问题节点时，就会触发 ViewChange 事件，ViewChange 完成后，视图编号将会加 1。



viewchange 会有三个阶段，分别是 view-change，view-change-ack 和 new-view 阶段。从节点认为主节点有问题时，会向其它节点发送 view-change 消息，当前存活的节点编号最小的节点将成为新的主节点。

```
<View-Change, v+1, C, P, Q, i>
```

- h 是节点 i 的最新 stable checkpoint
- C 是每个 checkpoint < 序号, digest > 的集合
- 集合 P 存储了之前 view 中到达 prepared 状态的请求。数据结构为 <n, d, v> 意为节点 i 收集到 view v 中序号为 n, digest 为 d 的请求的 prepared certificate, 并且节点 i 上的下一个 view 中就不能有相同序号的 prepared 请求。
- 集合 Q 存储了之前 view 中到达 pre-prepared 状态, 数据结构为 <n,d,v> 意为 i 已经 pre-prepare 了一个请求并且在下一 view 中这个请求没有相同序号的 pre-prepare 请求。

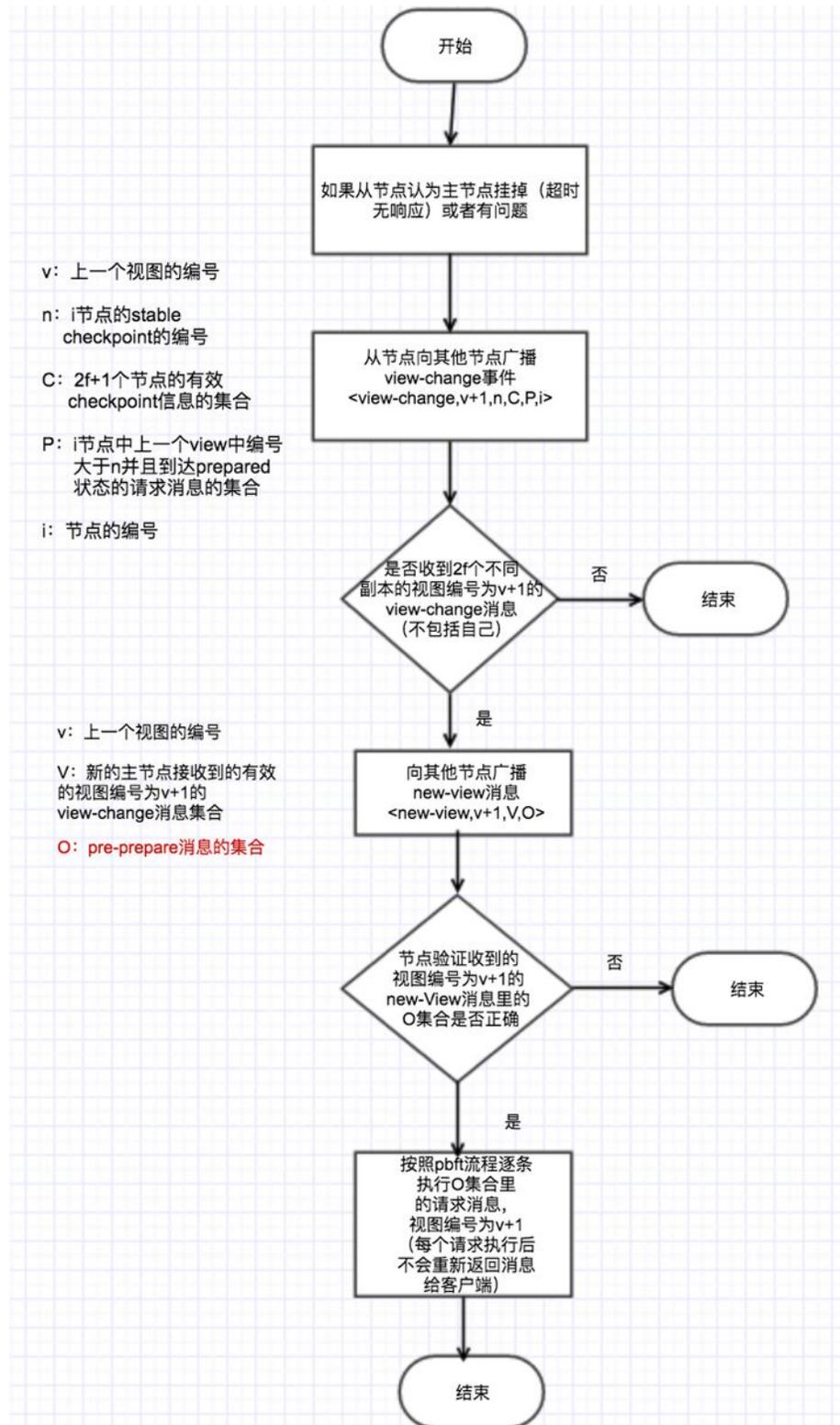
副本收集 VIEW-CHANGE 消息后发送 acknowledgments 以达成切换到 v+1 的共识

```
<View-Change-ACK, v+1, i, j, d>
```

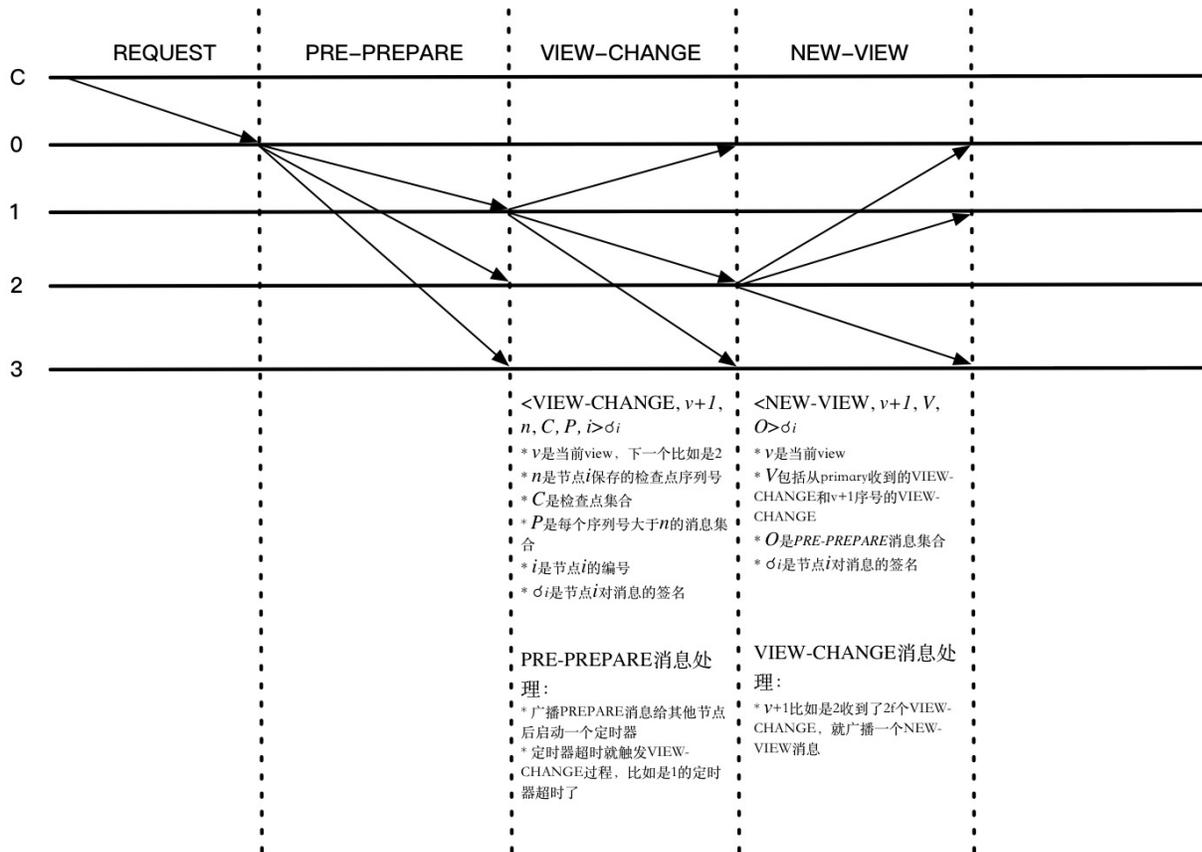
其中 i 是发送者的标识, d 是 acknowledged VIEW-CHANGE message 的 digest, j 是发送 VIEW-CHANGE 消息的节点。

当新的主节点收到 $2f$ 个其它节点的 view-change 消息, 则证明有足够多人的节点认为主节点有问题, 于是就会向其它节点广播 New-view 消息

对于主节点, 发送 new-view 消息后会继续执行上个视图未处理完的请求, 从 pre-prepare 阶段开始。其它节点验证 new-view 消息通过后, 就会处理主节点发来的 pre-prepare 消息, 这时执行的过程就是前面描述的 pbft 过程。到这时, 正式进入 v+1 (视图编号加 1) 的时代了。



VIEW-CHANGE



上图是发生 VIEW-CHANGE 的一种情况, 就是节点正常收到 PRE-PREPARE 消息以后都会启动一个定时器, 如果在设置的时间内都没有收到回复, 就会触发 VIEW-CHANGE, 该节点就不会再接收除 CHECKPOINT、VIEW-CHANGE 和 NEW-VIEW 等消息外的其他消息了。NEW-VIEW 是由新一轮的 primary 节点发送的, O 是不包含捎带的 REQUEST 的 PRE-PREPARE 消息集合, 计算方法如下:

primary 节点确定 V 中最新的稳定检查点序列号 $\min-s$ 和 PRE-PREPARE 消息中最大的序列号 $\max-s$ 对 $\min-s$ 和 $\max-s$ 之间每个序列号 n 都生成一个 PRE-PREPARE 消息。这可能有两种情况: P 的 VIEW-CHANGE 消息中至少存在一个集合, 序列号是 n 不存在上面的集合第一种情况, 会生成新的 PRE-PREPARE 消息

```
<PRE-PREPARE,  $v+1$ ,  $n$ ,  $d$ >
```

高低水位 (处理垃圾回收)

- checkpoint 就是当前节点处理的最新请求序号。
- stable checkpoint 就是大部分节点 ($2f+1$) 已经共识完成的最大请求序号。比如系统有 4 个节点, 三个节点都已经共识完了的请求编号是 213, 那么这个 213 就是 stable checkpoint 了。

为了不保留过去的所有 preprepare, prepare 以及 commit 等 log 信息, 所以我们每个一段时间设立一个检查点, 这个检查点需要得到全网的共识。这个经过全网共识的检查点被称为 stable checkpoint。在 stable checkpoint 前面的所有 log 信息就可以删除。

stable checkpoint 最大的目的就是减少内存的占用。

对于每个共识结点每 K 个请求，就设立一个 `checkpoint`，但由于每个共识结点共识过程是异步的，所以可能结点 1 刚刚处理了请求 1，但结点 2 已经处理了请求 k ，此时结点 2 设立 `checkpoint_1`，并广播出去（告诉大家我已经到 K 了， K 时候的状态是这样的）。此时，结点 2 也不会等其他结点（如结点 1）返回同样的 `checkpoint_1` 信息再继续处理请求。结点 2 一直向前处理交易，比如处理了 $2K$ 个交易后，他设立另一个 `checkpoint_2` 并广播出去。

结点 2 是不能一直狂奔的，这涉及到高低水位。因为如果主节点是坏的，它在给请求编号时故意选择了一个很大的编号，以至于超出了序号的范围，所以我们需要设置一个低水位（low water mark） h 和高水位（high water mark） H ，让主节点分配的编号在 h 和 H 之间，不能肆意分配。

```
低水位：h
高水位：H
间隔：L
即：H=h+L
```

如果我们的间隔设为 2 倍的 K ，低水位为最新的 `stable checkpoint_0`，其为 0；高水位为 $0+2K=2K$ ；现在结点 2 已经狂奔到 $2K$ 了，他不能向前了，所以他只能等待低水位 L 的前进，即他等到了 $2f$ 个 `checkpoint_1` 的确认信息。此时， L 会先前移动为 `stable_checkpoint_1`，其为 K ；高水位变成 $3K$ ，结点 2 可以继续前进了。

reference

PBFT

共识算法系列之一：raft 和 pbft 算法

2.1.5 POW 算法

掷骰子游戏

矿工节点挖矿

监控网络新的区块

所有的节点时刻监听传播到比特币网络的新区块。而这些新加入的区块对挖矿节点有着特殊的意义。矿工间的竞争以新区块的传播而结束，如同宣布谁是最后的赢家。对于矿工们来说，收到一个新区块进行验证意味着别人已经赢了，而自己则输了这场竞争。

节点同步区块后，去除交易池中的已成交交易，准备挖下一个块

接收并验证区块 277,315 后，节点会检查内存池中的全部交易，并移除已经在区块 277,315 中出现过的交易记录，确保任何留在内存池中的交易都是未确认的，等待被记录到新区块中。

收集交易，准备出块

收到新的交易后，验证交易，添加到自己的内存池（交易池），取出交易整合到一个候选区块中

创币交易

区块中的第一笔交易是笔特殊交易，称为创币交易或者 coinbase 交易。这个交易是出块节点 A 构造并用来奖励矿工节点 A 所做的贡献的。挖出区块获得的奖励金额是 coinbase 奖励（12.5 个全新的比特币）和区块中全部交易矿工费的总和。

与常规交易不同，创币交易没有输入，不消耗 UTXO。它只包含一个被称作 coinbase 的输入，仅仅用来创建新的比特币。创币交易有一个输出，支付到这个矿工的比特币地址。

创币交易需要 100 个区块确认高度，一般的需要 6 个区块确认高度

构建区块头

最后一个字段是 nonce，初始值为 0。

区块头完成全部的字段填充后，挖矿就可以开始进行了。挖矿的目标是找到一个使区块头哈希值小于难度目标的 nonce。挖矿节点通常需要尝试数十亿甚至数万亿个不同的 nonce 取值，直到找到一个满足条件的 nonce 值。

构建区块

比特币挖矿过程使用的是 SHA256 哈希函数。

用最简单的术语来说，挖矿就是重复计算区块头的哈希值，不断修改该参数，直到与哈希值匹配的一个过程。哈希函数的结果无法提前得知，也没有能得到一个特定哈希值的模式。哈希函数的这个特性意味着：得到哈希值的唯一方法是不断的尝试，每次随机修改输入，直到出现适当的哈希值。

维护链的方式：最长链，最长链又有 2 种方式，依据出块速度，依据块的难度

DPOS 节点作恶，不按排程产生区块的超级节点将在下一轮被投票剔除，被没收缴纳的保证金。

- 拒绝出块，使区块处于停滞状态，其它节点会补足
- 伪造交易，修改链上状态，其它节点无法通过校验

Pow

POW 算法

比特币的共识算法不适合于私有链和联盟链。其原因首先是它是一个最终一致性共识算法，不是一个强一致性共识算法。第二个原因是其共识效率低。提供共识效率又会牺牲共识协议的安全性。

工作证明必须产生小于目标的哈希值。更高的目标意味着找到低于目标的哈希是不太困难的。较低的目标意味着在目标下方找到哈希更难。目标和难度是成反比。

2.1.6 寻找一种易于理解的一致性算法（扩展版）

摘要

Raft 是一种为了管理复制日志的一致性算法。它提供了和 Paxos 算法相同的功能和性能，但是它的算法结构和 Paxos 不同，使得 Raft 算法更加容易理解并且更容易构建实际的系统。为了提升可理解性，Raft 将一致性算法分解成了几个关键模块，例如领导人选举、日志复制和安全性。同时它通过实施一个更强的一致性来减少需要考虑的状态的数量。从一个用户研究的结果可以证明，对于学生而言，Raft 算法比 Paxos 算法更加容易学习。Raft 算法还包括一个新的机制来允许集群成员的动态改变，它利用重叠的大多数来保证安全性。

1 介绍

一致性算法允许一组机器像一个整体一样工作，即使其中一些机器出现故障也能够继续工作下去。正因为如此，一致性算法在构建可信赖的大规模软件系统中扮演着重要的角色。在过去的 10 年里，Paxos 算法统治着一致性算法这一领域：绝大多数的实现都是基于 Paxos 或者受其影响。同时 Paxos 也成为了教学领域里讲解一致性问题时的示例。

但是不幸的是，尽管有很多工作都在尝试降低它的复杂性，但是 Paxos 算法依然十分难以理解。并且，Paxos 自身的算法结构需要进行大幅的修改才能够应用到实际的系统中。这些都导致了工业界和学术界都对 Paxos 算法感到十分头疼。

和 Paxos 算法进行过努力之后，我们开始寻找一种新的一致性算法，可以为构建实际的系统和教学提供更好的基础。我们的做法是不寻常的，我们的首要目标是可理解性：我们是否可以在实际系统中定义一个一致性算法，并且能够比 Paxos 算法以一种更加容易的方式来学习。此外，我们希望该算法方便系统构建者的直觉的发展。不仅一个算法能够工作很重要，而且能够显而易见的知道为什么能工作也很重要。

Raft 一致性算法就是这些工作的结果。在设计 Raft 算法的时候，我们使用一些特别的技巧来提升它的可理解性，包括算法分解（Raft 主要被分成了领导人选举，日志复制和安全三个模块）和减少状态机的状态（相对于 Paxos，Raft 减少了非确定性和服务器互相处于非一致性的方式）。一份针对两所大学 43 个学生的研究表明 Raft 明显比 Paxos 算法更加容易理解。在这些学生同时学习了这两种算法之后，和 Paxos 比起来，其中 33 个学生能够回答有关于 Raft 的问题。

Raft 算法在许多方面和现有的一致性算法都很相似（主要是 Oki 和 Liskov 的 Viewstamped Replication），但是它也有一些独特的特性：

- **强领导者**：和其他一致性算法相比，Raft 使用一种更强的领导能力形式。比如，日志条目只从领导者发送给其他的服务器。这种方式简化了对复制日志的管理并且使得 Raft 算法更加易于理解。
- **领导选举**：Raft 算法使用一个随机计时器来选举领导者。这种方式只是在任何一致性算法都必须实现的心跳机制上增加了一点机制。在解决冲突的时候会更加简单快捷。
- **成员关系调整**：Raft 使用一种共同一致的方法来处理集群成员变换的问题，在这种方法下，处于调整过程中的两种不同的配置集群中大多数机器会有重叠，这就使得集群在成员变换的时候依然可以继续工作。

我们相信，Raft 算法不论出于教学目的还是作为实践项目的基础都是要比 Paxos 或者其他一致性算法要优异的。它比其他算法更加简单，更加容易理解；它的算法描述足以实现一个现实的系统；它有好多开源的实现并且在很多公司里使用；它的安全性已经被证明；它的效率和其他算法比起来也不相上下。

接下来，这篇论文会介绍以下内容：复制状态机问题（第 2 节），讨论 Paxos 的优点和缺点（第 3 节），讨论我们为了可理解性而采取的方法（第 4 节），阐述 Raft 一致性算法（第 5-8 节），评价 Raft 算法（第 9 节），以及一些相关的工作（第 10 节）。

2 复制状态机

一致性算法是从复制状态机的背景下提出的（参考英文原文引用 37）。在这种方法中，一组服务器上的状态机产生相同状态的副本，并且在一些机器宕掉的情况下也可以继续运行。复制状态机在分布式系统中被用于解决很多容错的问题。例如，大规模的系统中通常都有一个集群领导者，像 GFS、HDFS 和 RAMCloud，典型应用就是一个独立的复制状态机去管理领导选举和存储配置信息并且在领导人宕机的情况下也要存活下来。比如 Chubby 和 ZooKeeper。

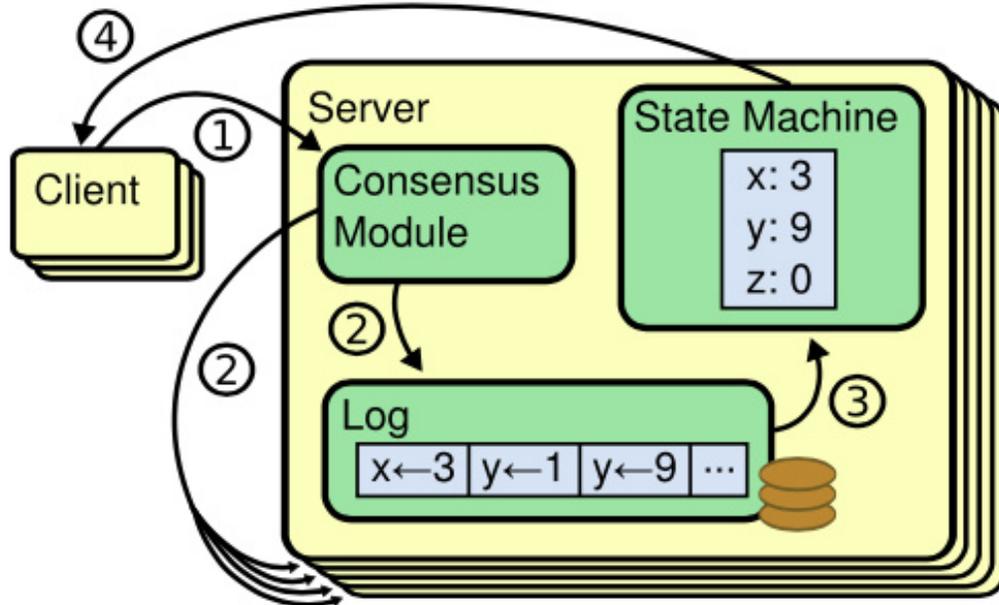


图 1：复制状态机的结构。一致性算法管理着来自客户端指令的复制日志。状态机从日志中处理相同顺序的相同指令，所以产生的结果也是相同的。

复制状态机通常都是基于复制日志实现的，如图 1。每一个服务器存储一个包含一系列指令的日志，并且按照日志的顺序进行执行。每一个日志都按照相同的顺序包含相同的指令，所以每一个服务器都执行相同的指令序列。因为每个状态机都是确定的，每一次执行操作都产生相同的状态和同样的序列。

保证复制日志相同就是一致性算法的工作了。在一台服务器上，一致性模块接收客户端发送来的指令然后增加到自己的日志中去。它和其他服务器上的一致性模块进行通信来保证每一个服务器上的日志最终都以相同的顺序包含相同的请求，尽管有些服务器会宕机。一旦指令被正确的复制，每一个服务器的状态机按照日志顺序处理他们，然后输出结果被返回给客户端。因此，服务器集群看起来形成一个高可靠的状态机。

实际系统中使用的一致性算法通常含有以下特性：

- 安全性保证（绝对不会返回一个错误的结果）：在非拜占庭错误情况下，包括网络延迟、分区、丢包、冗余和乱序等错误都可以保证正确。
- 可用性：集群中只要有大多数的机器可运行并且能够相互通信、和客户端通信，就可以保证可用。因此，一个典型的包含 5 个节点的集群可以容忍两个节点的失败。服务器被停止就认为是失败。他们当有稳定的存储的时候可以从状态中恢复回来并重新加入集群。
- 不依赖时序来保证一致性：物理时钟错误或者极端的消息延迟只有在最坏情况下才会导致可用性问题。
- 通常情况下，一条指令可以尽可能快的在集群中大多数节点响应一轮远程过程调用时完成。小部分比较慢的节点不会影响系统整体的性能。

3 Paxos 算法的问题

在过去的 10 年里，Leslie Lamport 的 Paxos 算法几乎已经成为一致性的代名词：Paxos 是在课程教学中最经常使用的算法，同时也是大多数一致性算法实现的起点。Paxos 首先定义了一个能够达成单一决策一致的协议，比如单条的复制日志项。我们把这一子集叫做单决策 Paxos。然后通过组合多个 Paxos 协议的实例来促进一系列决策的达成。Paxos 保证安全性和活性，同时也支持集群成员关系的变更。Paxos 的正确性已经被证明，在通常情况下也很高效。

不幸的是，Paxos 有两个明显的缺点。第一个缺点是 Paxos 算法特别的难以理解。完整的解释是出了名的不透明；通过极大的努力之后，也只有少数人成功理解了这个算法。因此，有了几次用更简单的术语来解释 Paxos 的尝试。尽管这些解释都只关注了单决策的子集问题，但依然很具有挑战性。在 2012 年 NSDI 的会议中的一次调查显示，很少有人对 Paxos 算法感到满意，甚至在经验老道的研究者中也是如此。我们自己也尝试去理解 Paxos；我们一直没能理解 Paxos 直到我们读了很多对 Paxos 的简化解释并且设计了我们自己的算法之后，这一过程花了近一年时间。

我们假设 Paxos 的不透明性来自它选择单决策问题作为它的基础。单决策 Paxos 是晦涩微妙的，它被划分成了两种没有简单直观解释和无法独立理解的情景。因此，这导致了很难建立起直观的感受为什么单决策 Paxos 算法能够工作。构成多决策 Paxos 增加了很多错综复杂的规则。我们相信，在多决策上达成一致性的问题（一份日志而不是单一的日志记录）能够被分解成其他的方式并且更加直接和明显。

Paxos 算法的第二个问题就是它没有提供一个足够好的用来构建一个现实系统的基础。一个原因是还没有一种被广泛认同的多决策问题的算法。Lamport 的描述基本上都是关于单决策 Paxos 的；他简要描述了实施多决策 Paxos 的方法，但是缺乏很多细节。当然也有很多具体化 Paxos 的尝试，但是他们都互相不一样，和 Paxos 的概述也不同。例如 Chubby 这样的系统实现了一个类似于 Paxos 的算法，但是大多数的细节并没有被公开。

而且，Paxos 算法的结构也不是十分易于构建实践的系统；单决策分解也会产生其他的结果。例如，独立的选择一组日志条目然后合并成一个序列化的日志并没有带来太多的好处，仅仅增加了不少复杂性。围绕着日志来设计一个系统是更加简单高效的；新日志条目以严格限制的顺序增添到日志中去。另一个问题是，Paxos 使用了一种对等的点对点的方式作为它的核心（尽管它最终提议了一种弱领导人的方法来优化性能）。在只有一个决策会被制定的简化世界中是很有意义的，但是很少有现实的系统使用这种方式。如果有一系列的决策需要被制定，首先选择一个领导人，然后让他去协调所有的决议，会更加简单快速。

因此，实际的系统中很少有和 Paxos 相似的实践。每一种实现都是从 Paxos 开始研究，然后发现很多实现上的难题，再然后开发了一种和 Paxos 明显不一样的结构。这样是非常费时和容易出错的，并且理解 Paxos 的难度使得这个问题更加糟糕。Paxos 算法在理论上被证明是正确可行的，但是现实的系统和 Paxos 差别是如此的大，以至于这些证明没有什么太大的价值。下面来自 Chubby 实现非常典型：

在 Paxos 算法描述和实现现实系统中间有着巨大的鸿沟。最终的系统建立在一种没有经过证明的算法之上。

由于以上问题，我们认为 Paxos 算法既没有提供一个良好的基础给实践的系统，也没有给教学很好的帮助。基于一致性问题在大规模软件系统中的重要性，我们决定看看我们是否可以设计一个拥有更好特性的替代 Paxos 的一致性算法。Raft 算法就是这次实验的结果。

4 为了可理解性的设计

设计 Raft 算法我们有几个初衷：它必须提供一个完整的实际的系统实现基础，这样才能大大减少开发者的工作；它必须在任何情况下都是安全的并且在大多数的情况下都是可用的；并且它的大部分操作必须是高效的。但是我们最重要也是最大的挑战是可理解性。它必须保证对于普遍的人群都可以十分容易的去理解。另外，它必须能够让人形成直观的认识，这样系统的构建者才能够现实中必然的扩展。

在设计 Raft 算法的时候，有很多的点需要我们在各种备选方案中进行选择。在这种情况下，我们评估备选方案基于可理解性原则：解释各个备选方案有多大的难度（例如，Raft 的状态空间有多复杂，是否有微妙的暗示）？对于一个读者而言，完全理解这个方案和暗示是否容易？

我们意识到对这种可理解性分析上具有高度的主观性；尽管如此，我们使用了两种通常适用的技术来解决这个问题。第一个技术就是众所周知的问题分解：只要有可能，我们就将问题分解成几个相对独立的，可被解

决的、可解释的和可理解的子问题。例如，Raft 算法被我们分成领导人选举，日志复制，安全性和角色改变几个部分。

我们使用的第二个方法是通过减少状态的数量来简化需要考虑的状态空间，使得系统更加连贯并且在可能的時候消除不确定性。特别的，所有的日志是不允许有空洞的，并且 Raft 限制了日志之间变成不一致状态的可能。尽管在大多数情况下我们都试图去消除不确定性，但是也有一些情况下不确定性可以提升可理解性。尤其是，随机化方法增加了不确定性，但是他们有利于减少状态空间数量，通过处理所有可能选择时使用相似的方法。我们使用随机化去简化 Raft 中领导人选举算法。

5 Raft 一致性算法

Raft 是一种用来管理章节 2 中描述的复制日志的算法。图 2 为了参考之用，总结这个算法的简略版本，图 3 列举了这个算法的一些关键特性。图中的这些元素会在剩下的章节逐一介绍。

Raft 通过选举一个高贵的领导人，然后给予他全部的管理复制日志的责任来实现一致性。领导人从客户端接收日志条目，把日志条目复制到其他服务器上，并且当保证安全性的时候告诉其他的服务器应用日志条目到他们的状态机中。拥有一个领导人大大简化了对复制日志的管理。例如，领导人可以决定新的日志条目需要放在日志中的什么位置而不需要和其他服务器商议，并且数据都从领导人流向其他服务器。一个领导人可以宕机，可以和其他服务器失去连接，这时一个新的领导人会被选举出来。

通过领导人的方式，Raft 将一致性问题分解成了三个相对独立的子问题，这些问题会在接下来的子章节中进行讨论：

- **领导选举**：一个新的领导人需要被选举出来，当现存的领导人宕机的时候（章节 5.2）
- **日志复制**：领导人必须从客户端接收日志然后复制到集群中的其他节点，并且强制要求其他节点的日志保持和自己相同。
- **安全性**：在 Raft 中安全性的关键是在图 3 中展示的状态机安全：如果有任何的服务器节点已经应用了一个确定的日志条目到它的状态机中，那么其他服务器节点不能在同一个日志索引位置应用一个不同的指令。章节 5.4 阐述了 Raft 算法是如何保证这个特性的；这个解决方案涉及到一个额外的选举机制（5.2 节）上的限制。

在展示一致性算法之后，这一章节会讨论可用性的一些问题和计时在系统的作用。

状态：

附加日志 RPC：

由领导人负责调用来复制日志指令；也会用作 heartbeat

接收者实现：

1. 如果 `term < currentTerm` 就返回 `false`（5.1 节）
2. 如果日志在 `prevLogIndex` 位置处的日志条目的任期号和 `prevLogTerm` 不匹配，则返回 `false`（5.3 节）
3. 如果已经存在的日志条目和新的产生冲突（索引值相同但是任期号不同），删除这一条和之后所有的（5.3 节）
4. 附加日志中尚未存在的任何新条目
5. 如果 `leaderCommit > commitIndex`，令 `commitIndex` 等于 `leaderCommit` 和新日志条目索引值中较小的一个

请求投票 RPC：

由候选人负责调用来征集选票（5.2 节）

接收者实现：

1. 如果 `term < currentTerm` 返回 `false`（5.2 节）

2. 如果 `votedFor` 为空或者为 `candidateId`, 并且候选人的日志至少和自己一样新, 那么就投票给他 (5.2 节, 5.4 节)

所有服务器需遵守的规则:

所有服务器:

- 如果 `commitIndex > lastApplied`, 那么就 `lastApplied` 加一, 并把 `log[lastApplied]` 应用到状态机中 (5.3 节)
- 如果接收到的 RPC 请求或响应中, 任期号 `T > currentTerm`, 那么就令 `currentTerm` 等于 `T`, 并切换状态为跟随者 (5.1 节)

跟随者 (5.2 节):

- 响应来自候选人和领导者的请求
- 如果在超过选举超时时间的情况之前都没有收到领导人的心跳, 或者是候选人请求投票的, 就自己变成候选人

候选人 (5.2 节):

- 在转变成候选人后就立即开始选举过程
 - 自增当前的任期号 (`currentTerm`)
 - 给自己投票
 - 重置选举超时计时器
 - 发送请求投票的 RPC 给其他所有服务器
- 如果接收到大多数服务器的选票, 那么就变成领导人
- 如果接收到来自新的领导人的附加日志 RPC, 转变成跟随者
- 如果选举过程超时, 再次发起一轮选举

领导人:

- 一旦成为领导人: 发送空的附加日志 RPC (心跳) 给其他所有的服务器; 在一定的空余时间之后不停的重复发送, 以阻止跟随者超时 (5.2 节)
- 如果接收到来自客户端的请求: 附加条目到本地日志中, 在条目被应用到状态机后响应客户端 (5.3 节)
- 如果对于一个跟随者, 最后日志条目的索引值大于等于 `nextIndex`, 那么: 发送从 `nextIndex` 开始的所有日志条目:
 - 如果成功: 更新相应跟随者的 `nextIndex` 和 `matchIndex`
 - 如果因为日志不一致而失败, 减少 `nextIndex` 重试
- 如果存在一个满足 `N > commitIndex` 的 `N`, 并且大多数的 `matchIndex[i] ≥ N` 成立, 并且 `log[N].term == currentTerm` 成立, 那么令 `commitIndex` 等于这个 `N` (5.3 和 5.4 节)

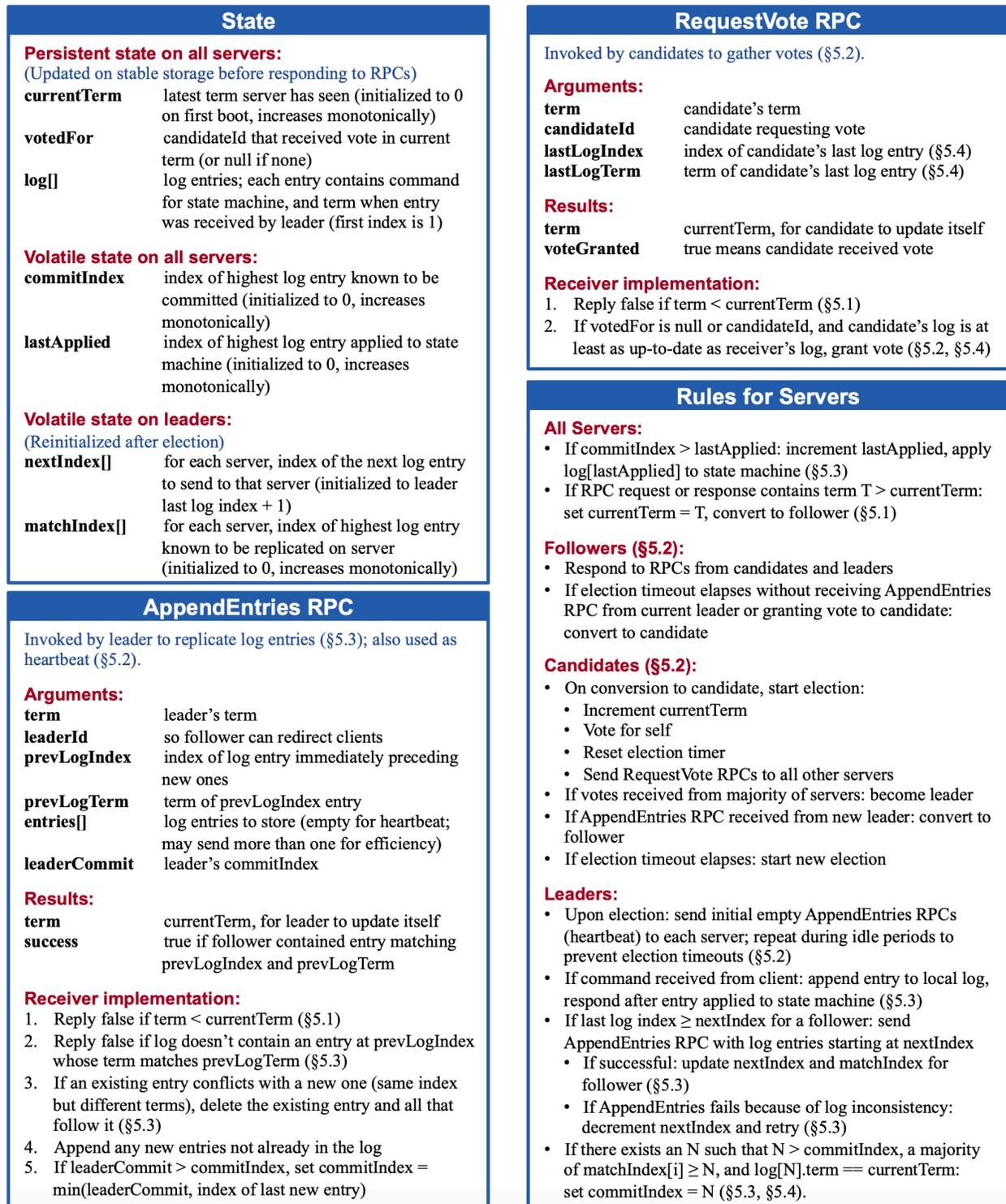


图 2: 一个关于 Raft 一致性算法的浓缩总结 (不包括成员变换和日志压缩)。

Election Safety: at most one leader can be elected in a given term. §5.2

Leader Append-Only: a leader never overwrites or deletes entries in its log; it only appends new entries. §5.3

Log Matching: if two logs contain an entry with the same index and term, then the logs are identical in all entries up through the given index. §5.3

Leader Completeness: if a log entry is committed in a given term, then that entry will be present in the logs of the leaders for all higher-numbered terms. §5.4

State Machine Safety: if a server has applied a log entry at a given index to its state machine, no other server will ever apply a different log entry for the same index. §5.4.3

图 3: Raft 在任何时候都保证以上的各个特性。

5.1 Raft 基础

一个 Raft 集群包含若干个服务器节点；通常是 5 个，这允许整个系统容忍 2 个节点的失效。在任何时刻，每一个服务器节点都处于这三个状态之一：领导人、跟随者或者候选人。在通常情况下，系统中只有一个领导人并且其他的节点全部都是跟随者。跟随者都是被动的：他们不会发送任何请求，只是简单的响应来自领导人或者候选人的请求。领导人处理所有的客户端请求（如果一个客户端和跟随者联系，那么跟随者会把请求重定向给领导人）。第三种状态，候选人，是用来在 5.2 节描述的选举新领导人时使用。图 4 展示了这些状态和他们之间的转换关系；这些转换关系会在接下来进行讨论。

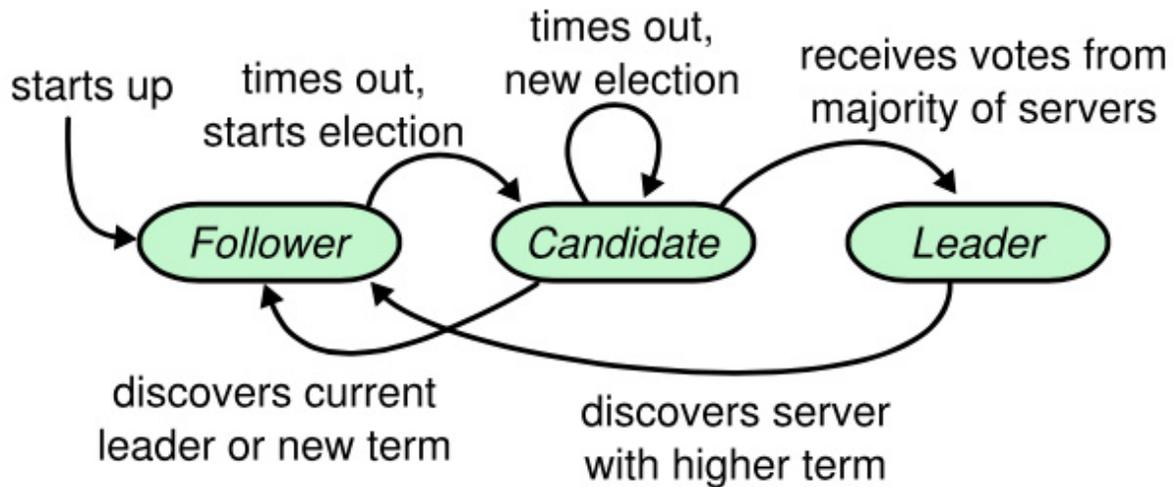


图 4: 服务器状态。跟随者只响应来自其他服务器的请求。如果跟随者接收不到消息, 那么他就会变成候选人并发起一次选举。获得集群中大多数选票的候选人将成为领导者。在一个任期内, 领导人一直都会是领导人直到自己宕机了。

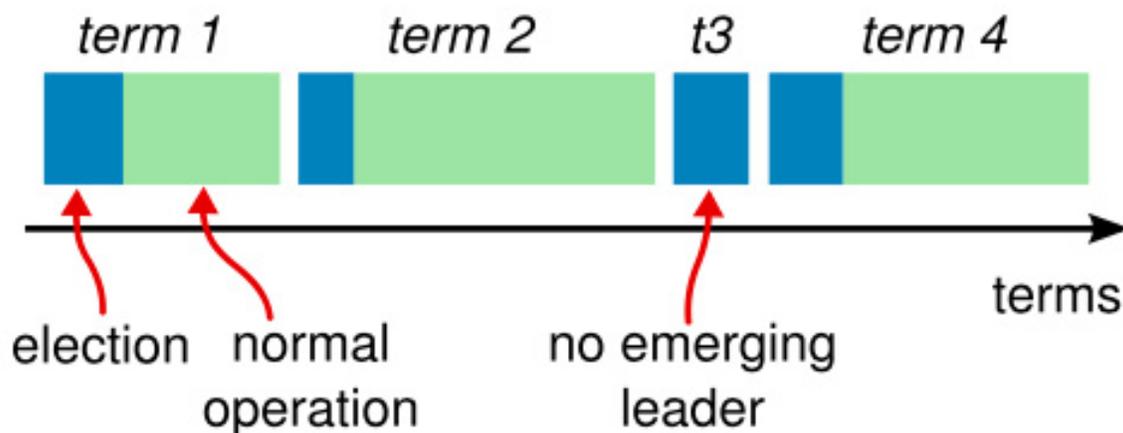


图 5: 时间被划分成一个个的任期, 每个任期开始都是一次选举。在选举成功后, 领导人会管理整个集群直到任期结束。有时候选举会失败, 那么这个任期就会没有领导人而结束。任期之间的切换可以在不同的时间不同的服务器上观察到。

Raft 把时间分割成任意长度的**任期**, 如图 5。任期用连续的整数标记。每一段任期从一次**选举**开始, 就像章节 5.2 描述的一样, 一个或者多个候选人尝试成为领导者。如果一个候选人赢得选举, 然后他就在接下来的任期内充当领导人的职责。在某些情况下, 一次选举过程会造成选票的瓜分。在这种情况下, 这一任期会以没有领导人结束; 一个新的任期 (和一次新的选举) 会很快重新开始。Raft 保证了在一个给定的任期内, 最多只有一个领导者。

不同的服务器节点可能多次观察到任期之间的转换, 但在某些情况下, 一个节点也可能观察不到任何一次选举或者整个任期全程。任期在 Raft 算法中充当逻辑时钟的作用, 这会允许服务器节点查明一些过期的信息比如陈旧的领导者。每一个节点存储一个当前任期号, 这一编号在整个时期内单调的增长。当服务器之间通信的时候会交换当前任期号; 如果一个服务器的当前任期号比其他人小, 那么他会更新自己的编号到较大的编号值。如果一个候选人或者领导者发现自己的任期号过期了, 那么他会立即恢复成跟随者状态。如果一个节点接收到一个包含过期的任期号的请求, 那么他会直接拒绝这个请求。

Raft 算法中服务器节点之间通信使用远程过程调用 (RPCs), 并且基本的一致性算法只需要两种类型的 RPCs。

请求投票 (RequestVote) RPCs 由候选人在选举期间发起 (章节 5.2), 然后附加条目 (AppendEntries) RPCs 由领导人发起, 用来复制日志和提供一种心跳机制 (章节 5.3)。第 7 节为了在服务器之间传输快照增加了第三种 RPC。当服务器没有及时的收到 RPC 的响应时, 会进行重试, 并且他们能够并行的发起 RPCs 来获得最佳的性能。

5.2 领导人选举

Raft 使用一种心跳机制来触发领导人选举。当服务器程序启动时, 他们都是跟随者身份。一个服务器节点继续保持着跟随者状态只要他从领导人或者候选者处接收到有效的 RPCs。领导者周期性的向所有跟随者发送心跳包 (即不包含日志项内容的附加日志项 RPCs) 来维持自己的权威。如果一个跟随者在一段时间里没有接收到任何消息, 也就是**选举超时**, 那么他就会认为系统中没有可用的领导者, 并且发起选举以选出新的领导者。

要开始一次选举过程, 跟随者先要增加自己的当前任期号并且转换到候选人状态。然后他会并行的向集群中的其他服务器节点发送请求投票的 RPCs 来给自己投票。候选人会继续保持着当前状态直到以下三件事情之一发生: (a) 他自己赢得了这次的选举, (b) 其他的服务器成为领导者, (c) 一段时间之后没有任何一个获胜的人。这些结果会分别在下面的段落里进行讨论。

当一个候选人从整个集群的大多数服务器节点获得了针对同一个任期号的选票, 那么他就赢得了这次选举并成为领导人。每一个服务器最多会对一个任期号投出一张选票, 按照先来先服务的原则 (注意: 5.4 节在投票上增加了一点额外的限制)。要求大多数选票的规则确保了最多只会有一个候选人赢得此次选举 (图 3 中的选举安全性)。一旦候选人赢得选举, 他就立即成为领导人。然后他会向其他的服务器发送心跳消息来建立自己的权威并且阻止新的领导人的产生。

在等待投票的时候, 候选人可能会从其他的服务器接收到声明它是领导人的附加日志项 RPC。如果这个领导人的任期号 (包含在此次的 RPC 中) 不小于候选人当前的任期号, 那么候选人会承认领导人合法并回到跟随者状态。如果此次 RPC 中的任期号比自己小, 那么候选人就会拒绝这次的 RPC 并且继续保持候选人状态。

第三种可能的结果是候选人既没有赢得选举也没有输: 如果有多个跟随者同时成为候选人, 那么选票可能会被瓜分以至于没有候选人可以赢得大多数人的支持。当这种情况发生的时候, 每一个候选人都会超时, 然后通过增加当前任期号来开始一轮新的选举。然而, 没有其他机制的话, 选票可能会被无限的重复瓜分。

Raft 算法使用随机选举超时时间的方法来确保很少会发生选票瓜分的情况, 就算发生也能很快的解决。为了阻止选票起初就被瓜分, 选举超时时间是从一个固定的区间 (例如 150-300 毫秒) 随机选择。这样可以把服务器都分散开以至于在大多数情况下只有一个服务器会选举超时; 然后他赢得选举并在其他服务器超时之前发送心跳包。同样的机制被用在选票瓜分的情况下。每一个候选人在开始一次选举的时候会重置一个随机的选举超时时间, 然后在超时时间内等待投票的结果; 这样减少了在新的选举中另外的选票瓜分的可能性。9.3 节展示了这种方案能够快速选出一个领导人。

领导人选举这个例子, 体现了可理解性原则是如何指导我们进行方案设计的。起初我们计划使用一种排名系统: 每一个候选人都被赋予一个唯一的排名, 供候选人之间竞争时进行选择。如果一个候选人发现另一个候选人拥有更高的排名, 那么他就会回到跟随者状态, 这样高排名的候选人能够更加容易的赢得下一次选举。但是我们发现这种方法在可用性方面会有一点问题 (如果高排名的服务器宕机了, 那么低排名的服务器可能会超时并再次进入候选人状态。而且如果这个行为发生得足够快, 则可能会导致整个选举过程都被重置掉)。我们针对算法进行了多次调整, 但是每次调整之后都会有新的问题。最终我们认为随机重试的方法是更加明显和易于理解的。

5.3 日志复制

一旦一个领导人被选举出来，他就开始为客户端提供服务。客户端的每一个请求都包含一条被复制状态机执行的指令。领导人把这条指令作为一条新的日志条目附加到日志中去，然后并行的发起附加条目 RPCs 给其他的服务器，让他们复制这条日志条目。当这条日志条目被安全的复制（下面会介绍），领导人会应用这条日志条目到它的状态机中然后把执行的结果返回给客户端。如果跟随者崩溃或者运行缓慢，又或者网络丢包，领导人会不断的重复尝试附加日志条目 RPCs（尽管已经回复了客户端）直到所有的跟随者都最终存储了所有的日志条目。

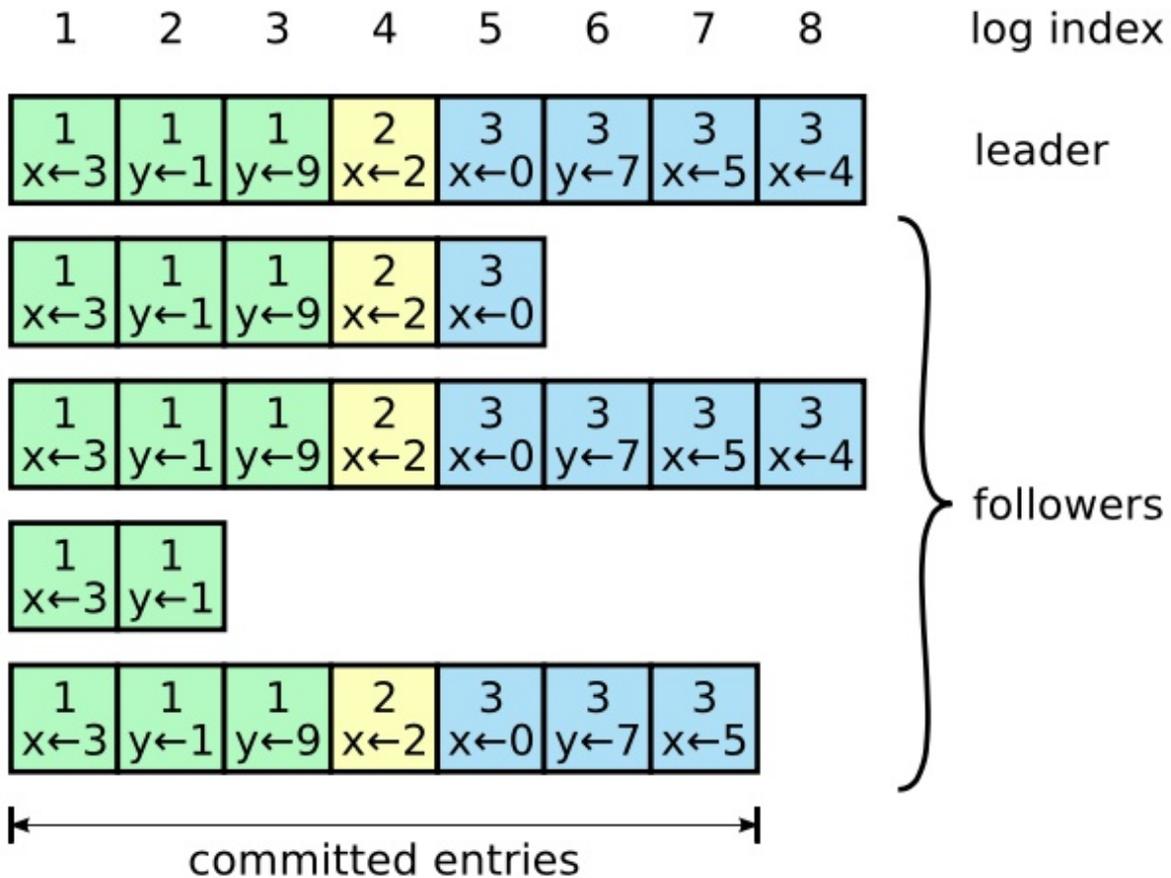


图 6: 日志由有序序号标记的条目组成。每个条目都包含创建时的任期号（图中框中的数字），和一个状态机需要执行的指令。一个条目当可以安全的被应用到状态机中去的时候，就认为是可以提交了。

日志以图 6 展示的方式组织。每一个日志条目存储一条状态机指令和从领导人收到这条指令时的任期号。日志中的任期号用来检查是否出现不一致的情况，同时也用来保证图 3 中的某些性质。每一条日志条目同时也都有一个整数索引值来表明它在日志中的位置。

领导人来决定什么时候把日志条目应用到状态机中是安全的；这种日志条目被称为**已提交**。Raft 算法保证所有已提交的日志条目都是持久化的并且最终会被所有可用的状态机执行。在领导人将创建的日志条目复制到大多数的服务器上的时候，日志条目就会被提交（例如在图 6 中的条目 7）。同时，领导人的日志中之前的所有日志条目也都会被提交，包括由其他领导人创建的条目。5.4 节会讨论某些当在领导人改变之后应用这条规则的隐晦内容，同时他也展示了这种提交的定义是安全的。领导人跟踪了最大的将会被提交的日志项的索引，并且索引值会被包含在未来的所有附加日志 RPCs（包括心跳包），这样其他的服务器才能最终知道领导人的提交位置。一旦跟随者知道一条日志条目已经被提交，那么他也会将这个日志条目应用到本地的状态机中（按照日志的顺序）。

我们设计了 Raft 的日志机制来维护一个不同服务器的日志之间的高层次的一致性。这么做不仅简化了系统的行为也使得更加可预计，同时他也是安全性保证的一个重要组件。Raft 维护着以下的特性，这些同时也组成了图 3 中的日志匹配特性：

- 如果在不同的日志中的两个条目拥有相同的索引和任期号，那么他们存储了相同的指令。
- 如果在不同的日志中的两个条目拥有相同的索引和任期号，那么他们之前的所有日志条目也全部相同。

第一个特性来自这样的一个事实，领导人最多在一个任期里在指定的一个日志索引位置创建一条日志条目，同时日志条目在日志中的位置也从来不会改变。第二个特性由附加日志 RPC 的一个简单的一致性检查所保证。在发送附加日志 RPC 的时候，领导人会把新的日志条目紧接着之前的条目的索引位置和任期号包含在里面。如果跟随者在它的日志中找不到包含相同索引位置和任期号的条目，那么他就会拒绝接收新的日志条目。一致性检查就像一个归纳步骤：一开始空的日志状态肯定是满足日志匹配特性的，然后一致性检查保护了日志匹配特性当日志扩展的时候。因此，每当附加日志 RPC 返回成功时，领导人就知道跟随者的日志一定是和自己相同的了。

在正常的操作中，领导人和跟随者的日志保持一致性，所以附加日志 RPC 的一致性检查从来不会失败。然而，领导人崩溃的情况会使得日志处于不一致的状态（老的领导人可能还没有完全复制所有的日志条目）。这种不一致问题会在领导人和跟随者的一系列崩溃下加剧。图 7 展示了跟随者的日志可能和新的领导人不同的方式。跟随者可能会丢失一些在新的领导人中有的日志条目，他也可能拥有一些领导人没有的日志条目，或者两者都发生。丢失或者多出日志条目可能会持续多个任期。

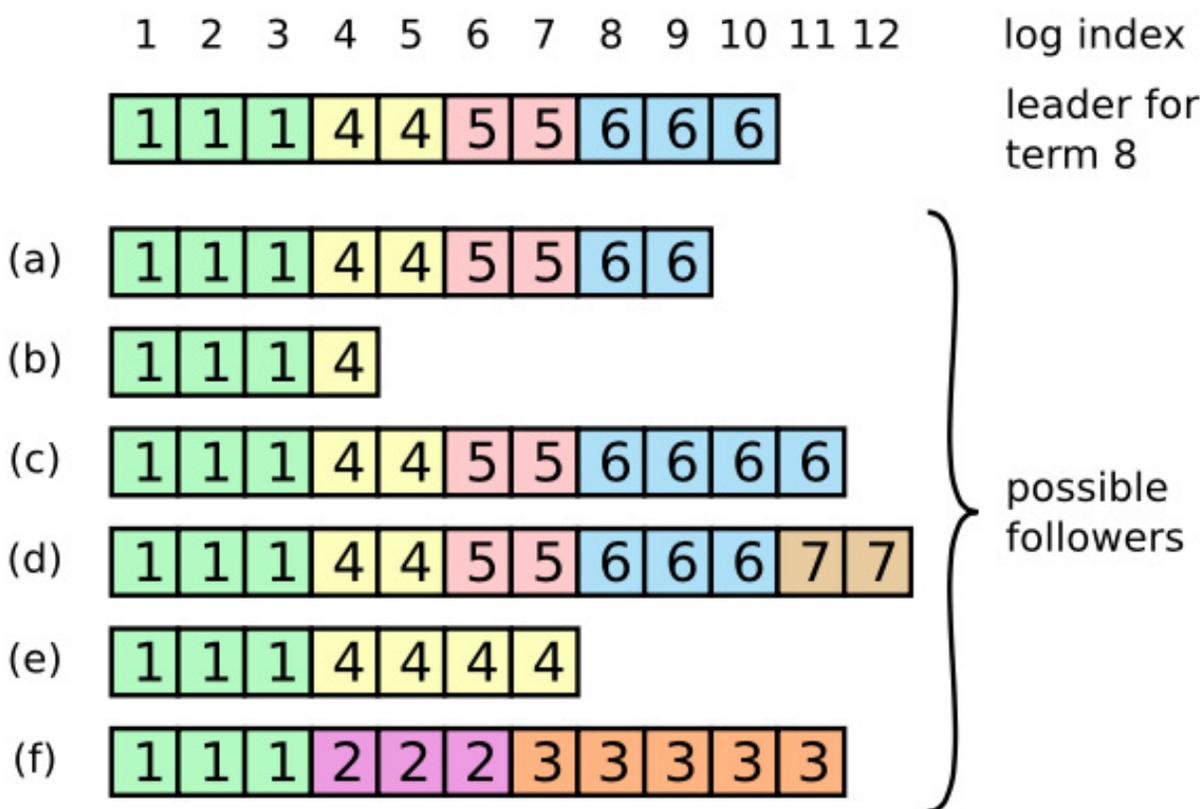


图 7：当一个领导人成功当选时，跟随者可能是任何情况 (a-f)。每一个盒子表示是一个日志条目；里面的数字表示任期号。跟随者可能会缺少一些日志条目 (a-b)，可能会有有一些未被提交的日志条目 (c-d)，或者两种情况都存在 (e-f)。例如，场景 f 可能会这样发生，某服务器在任期 2 的时候是领导人，已附加了一些日志条目到自己的日志中，但在提交之前就崩溃了；很快这个机器就被重启了，在任期 3 重新被选为领导人，并且又增加了一些日志条目到自己的日志中；在任期 2 和任期 3 的日志被提交之前，这个服务器又宕机了，并且在接下来的几个任期里一直处于宕机状态。

在 Raft 算法中，领导人处理不一致是通过强制跟随者直接复制自己的日志来解决的。这意味着在跟随者中的冲突的日志条目会被领导人的日志覆盖。5.4 节会阐述如何通过增加一些限制来使得这样的操作是安全的。

要使得跟随者的日志进入和自己一致的状态，领导人必须找到最后两者达成一致的地方，然后删除从那个点之后的所有日志条目，发送自己的日志给跟随者。所有的这些操作都在进行附加日志 RPCs 的一致性检查时完成。领导人针对每一个跟随者维护了一个 `nextIndex`，这表示下一个需要发送给跟随者的日志条目的索引地址。当一个领导人刚获得权力的时候，他初始化所有的 `nextIndex` 值为自己的最后一条日志的 `index` 加 1（图 7 中的 11）。如果一个跟随者的日志和领导人不一致，那么在下一轮的附加日志 RPC 时的一致性检查就会失败。在被跟随者拒绝之后，领导人就会减小 `nextIndex` 值并进行重试。最终 `nextIndex` 会在某个位置使得领导人和跟随者的日志达成一致。当这种情况发生，附加日志 RPC 就会成功，这时就会把跟随者冲突的日志条目全部删除并且加上领导人的日志。一旦附加日志 RPC 成功，那么跟随者的日志就会和领导人保持一致，并且在接下来的任期里一直继续保持。

如果需要的话，算法可以通过减少被拒绝的附加日志 RPCs 的次数来优化。例如，当附加日志 RPC 的请求被拒绝的时候，跟随者可以包含冲突的条目的任期号和自己存储的那个任期的最早的索引地址。借助这些信息，领导人可以减小 `nextIndex` 越过所有那个任期冲突的所有日志条目；这样就变成每个任期需要一次附加条目 RPC 而不是每个条目一次。在实践中，我们十分怀疑这种优化是否是必要的，因为失败是很少发生的并且也不大可能会有这么多不一致的日志。

通过这种机制，领导人在获得权力的时候就不需要任何特殊的操作来恢复一致性。他只需要进行正常的操作，然后日志就能自动的在回复附加日志 RPC 的一致性检查失败的时候自动趋于一致。领导人从来不会覆盖或者删除自己的日志（图 3 的领导人只附加特性）。

日志复制机制展示出了第 2 节中形容的一致性特性：Raft 能够接受，复制并应用新的日志条目只要大部分的机器是工作的；在通常的情况下，新的日志条目可以在一次 RPC 中被复制给集群中的大多数机器；并且单个的缓慢的跟随者不会影响整体的性能。

5.4 安全性

前面的章节里描述了 Raft 算法是如何选举和复制日志的。然而，到目前为止描述的机制并不能充分的保证每一个状态机会按照相同的顺序执行相同的指令。例如，一个跟随者可能会进入不可用状态同时领导人已经提交了若干的日志条目，然后这个跟随者可能会被选举为领导人并且覆盖这些日志条目；因此，不同的状态机可能会执行不同的指令序列。

这一节通过在领导选举的时候增加一些限制来完善 Raft 算法。这一限制保证了任何的领导人对于给定的任期号，都拥有了之前任期的所有被提交的日志条目（图 3 中的领导人完整特性）。增加这一选举时的限制，我们对于提交时的规则也更加清晰。最终，我们将展示对于领导人完整特性的简要证明，并且说明领导人是如何领导复制状态机的做出正确行为的。

5.4.1 选举限制

在任何基于领导人的一致性算法中，领导人都必须存储所有已经提交的日志条目。在某些一致性算法中，例如 `Viewstamped Replication`，某个节点即使是一开始并没有包含所有已经提交的日志条目，它也能被选为领导者。这些算法都包含一些额外的机制来识别丢失的日志条目并把他们传送给新的领导人，要么是在选举阶段要么在之后很快进行。不幸的是，这种方法会导致相当大的额外的机制和复杂性。Raft 使用了一种更加简单的方法，它可以保证所有之前的任期号中已经提交的日志条目在选举的时候都会出现在新的领导人中，不需要传送这些日志条目给领导人。这意味着日志条目的传送是单向的，只从领导人传给跟随者，并且领导人从不会覆盖自身本地日志中已经存在的条目。

Raft 使用投票的方式来阻止一个候选人赢得选举除非这个候选人包含了所有已经提交的日志条目。候选人为了赢得选举必须联系集群中的大部分节点，这意味着每一个已经提交的日志条目在这些服务器节点中肯定存在于至少一个节点上。如果候选人的日志至少和大多数的服务器节点一样新（这个新的定义会在下面讨论），那么他一定持有了所有已经提交的日志条目。请求投票 RPC 实现了这样的限制：RPC 中包含了候选人的日志信息，然后投票人会拒绝掉那些日志没有自己新的投票请求。

Raft 通过比较两份日志中最后一条日志条目的索引值和任期号定义谁的日志比较新。如果两份日志最后的条目的任期号不同，那么任期号大的日志更加新。如果两份日志最后的条目任期号相同，那么日志比较长的那个就更加新。

5.4.2 提交之前任期内的日志条目

如同 5.3 节介绍的那样，领导人知道一条当前任期内的日志记录是可以被提交的，只要它被存储到了大多数的服务器上。如果一个领导人在提交日志条目之前崩溃了，未来后续的领导人会继续尝试复制这条日志记录。然而，一个领导人不能断定一个之前任期里的日志条目被保存到大多数服务器上的时候就一定已经提交了。图 8 展示了一种情况，一条已经被存储到大多数节点上的老日志条目，也依然有可能会被未来的领导人覆盖掉。

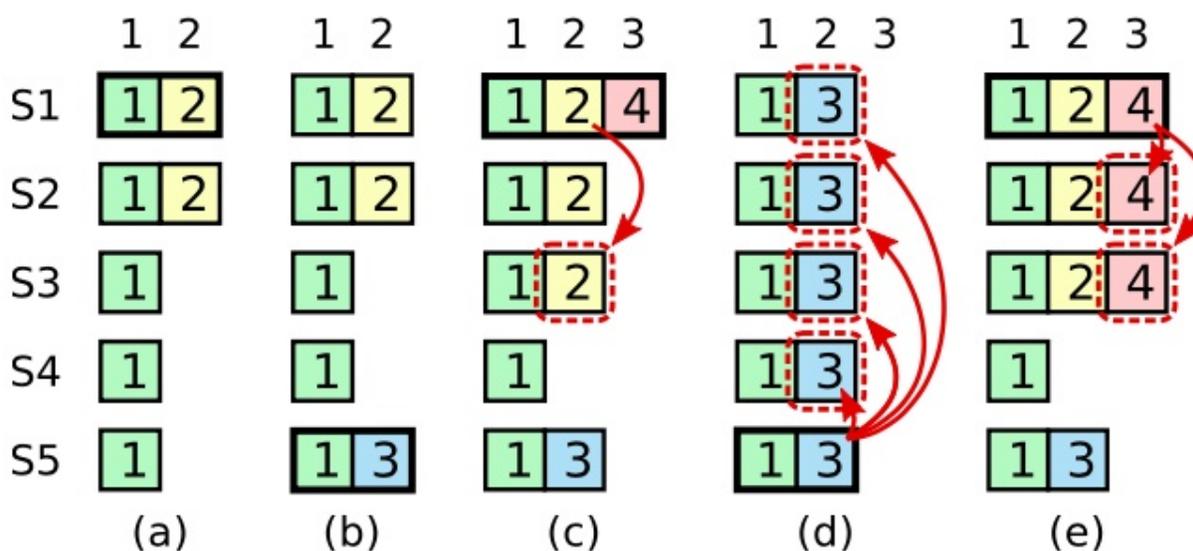


图 8：如图的时间序列展示了为什么领导人无法决定对老任期号的日志条目进行提交。在 (a) 中，S1 是领导者，部分的复制了索引位置 2 的日志条目。在 (b) 中，S1 崩溃了，然后 S5 在任期 3 里通过 S3、S4 和自己的选票赢得选举，然后从客户端接收了一条不一样的日志条目放在了索引 2 处。然后到 (c)，S5 又崩溃了；S1 重新启动，选举成功，开始复制日志。在这时，来自任期 2 的那条日志已经被复制到了集群中的大多数机器上，但是还没有被提交。如果 S1 在 (d) 中又崩溃了，S5 可以重新被选举成功（通过来自 S2、S3 和 S4 的选票），然后覆盖了他们在索引 2 处的日志。反之，如果在崩溃之前，S1 把自己主导的新任期里产生的日志条目复制到了大多数机器上，就如 (e) 中那样，那么在后面任期里面这些新的日志条目就会被提交（因为 S5 就不可能选举成功）。这样在同一时刻就同时保证了，之前的所有老的日志条目就会被提交。

为了消除图 8 里描述的情况，Raft 永远不会通过计算副本数目的方式去提交一个之前任期内的日志条目。只有领导人当前任期里的日志条目通过计算副本数目可以被提交；一旦当前任期的日志条目以这种方式被提交，那么由于日志匹配特性，之前的日志条目也都会被间接的提交。在某些情况下，领导人可以安全的知道一个老的日志条目是否已经被提交（例如，该条目是否存储到所有服务器上），但是 Raft 为了简化问题使用一种更加保守的方法。

当领导人复制之前任期里的日志时，Raft 会为所有日志保留原始的任期号，这在提交规则上产生了额外的复杂性。在其他的一致性算法中，如果一个新的领导人要重新复制之前的任期里的日志时，它必须使用当前新的任期号。Raft 使用的方法更加容易辨别出日志，因为它可以随着时间和日志的变化对日志维护着同一个任期编号。另外，和其他的算法相比，Raft 中的新领导人只需要发送更少日志条目（其他算法中必须在他们被提交之前发送更多的冗余日志条目来为他们重新编号）。

5.4.3 安全性论证

在给定了完整的 Raft 算法之后，我们现在可以更加精确的讨论领导人完整性特性（这一讨论基于 9.2 节的安全性证明）。我们假设领导人完整性特性是不存在的，然后我们推出矛盾来。假设任期 T 的领导人（领导人 T）在任期内提交了一条日志条目，但是这条日志条目没有被存储到未来某个任期的领导人的日志中。设大于 T 的最小任期 U 的领导人 U 没有这条日志条目。

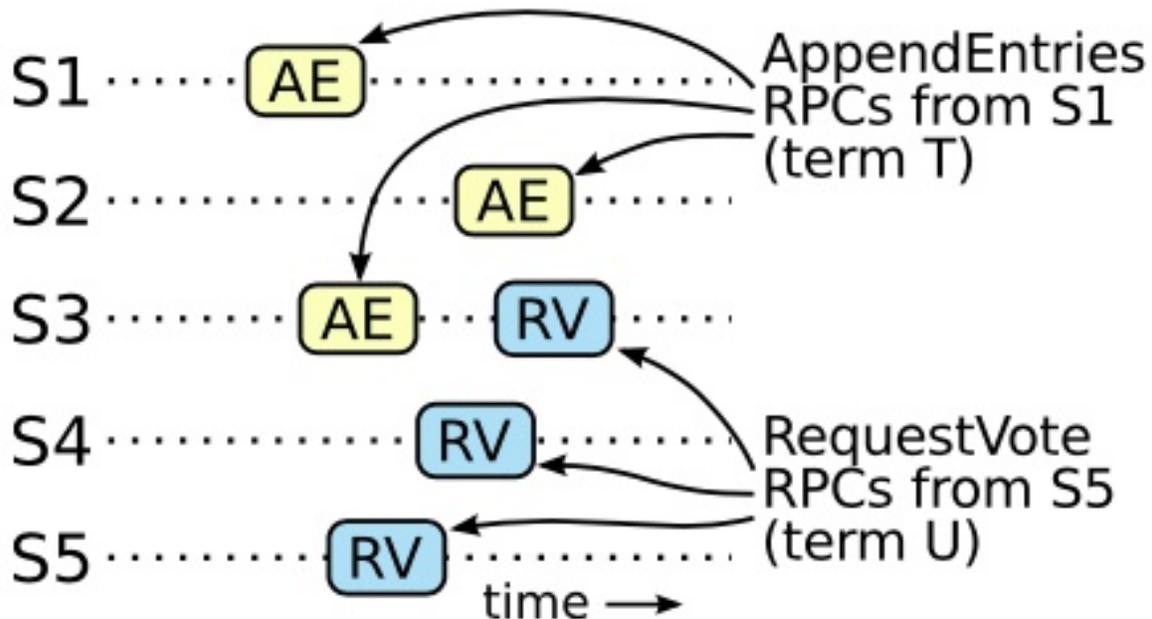


图 9: 如果 S1（任期 T 的领导者）提交了一条新的日志在它的任期内，然后 S5 在之后的任期 U 里被选举为领导人，然后至少会有一个机器，如 S3，既拥有来自 S1 的日志，也给 S5 投票了。

1. 在领导人 U 选举的时候一定没有那条被提交的日志条目（领导人从不会删除或者覆盖任何条目）。
2. 领导人 T 复制这条日志条目给集群中的大多数节点，同时，领导人 U 从集群中的大多数节点赢得了选票。因此，至少有一个节点（投票者、选民）同时接受了来自领导人 T 的日志条目，并且给领导人 U 投票了，如图 9。这个投票者是产生这个矛盾的关键。
3. 这个投票者必须在给领导人 U 投票之前先接受了从领导人 T 发来的已经被提交的日志条目；否则他就会拒绝来自领导人 T 的附加日志请求（因为此时他的任期号会比 T 大）。
4. 投票者在给领导人 U 投票时依然保存有这条日志条目，因为任何中间的领导人都包含该日志条目（根据上述的假设），领导人从不会删除条目，并且跟随者只有在和领导人冲突的时候才会删除条目。
5. 投票者把自己选票投给领导人 U 时，领导人 U 的日志必须和投票者自己一样新。这就导致了两者矛盾之一。
6. 首先，如果投票者和领导人 U 的最后一条日志的任期号相同，那么领导人 U 的日志至少和投票者一样长，所以领导人 U 的日志一定包含所有投票者的日志。这是另一处矛盾，因为投票者包含了那条已经被提交的日志条目，但是在上述的假设里，领导人 U 是不包含的。
7. 除此之外，领导人 U 的最后一条日志的任期号就必须比投票人大了。此外，他也比 T 大，因为投票人的最后一条日志的任期号至少和 T 一样大（他包含了来自任期 T 的已提交的日志）。创建了领导人 U 最后一条日志的之前领导人一定已经包含了那条被提交的日志（根据上述假设，领导人 U 是第一个不包含该日志条目的领导人）。所以，根据日志匹配特性，领导人 U 一定也包含那条被提交的日志，这里产生矛盾。
8. 这里完成了矛盾。因此，所有比 T 大的领导人一定包含了所有来自 T 的已经被提交的日志。

9. 日志匹配原则保证了未来的领导人也会同时包含被间接提交的条目，例如图 8 (d) 中的索引 2。

通过领导人完全特性，我们就能证明图 3 中的状态机安全特性，即如果服务器已经在某个给定的索引值应用了日志条目到自己的状态机里，那么其他的服务器不会应用一个不一样的日志到同一个索引值上。在一个服务器应用一条日志条目到他自己的状态机中时，他的日志必须和领导人的日志，在该条目和之前的条目上相同，并且已经被提交。现在我们来考虑在任何一个服务器应用一个指定索引位置的日志的最小任期；日志完全特性保证拥有更高任期号的领导人会存储相同的日志条目，所以之后的任期里应用某个索引位置的日志条目也会是相同的值。因此，状态机安全特性是成立的。

最后，Raft 要求服务器按照日志中索引位置顺序应用日志条目。和状态机安全特性结合起来看，这就意味着所有的服务器会应用相同的日志序列集到自己的状态机中，并且是按照相同的顺序。

5.5 跟随者和候选人崩溃

到目前为止，我们都只关注了领导人崩溃的情况。跟随者和候选人崩溃后的处理方式比领导人要简单的多，并且他们的处理方式是相同的。如果跟随者或者候选人崩溃了，那么后续发送给他们的 RPCs 都会失败。Raft 中处理这种失败就是简单的通过无限的重试；如果崩溃的机器重启了，那么这些 RPC 就会完整的成功。如果一个服务器在完成了 RPC，但是还没有响应的时候崩溃了，那么在他重新启动之后就会再次收到同样的请求。Raft 的 RPCs 都是幂等的，所以这样重试不会造成任何问题。例如一个跟随者如果收到附加日志请求但是他已经包含了这一日志，那么他就会直接忽略这个新的请求。

5.6 时间和可用性

Raft 的要求之一就是安全性不能依赖时间：整个系统不能因为某些事件运行的比预期快一点或者慢一点就产生了错误的结果。但是，可用性（系统可以及时的响应客户端）不可避免的要依赖于时间。例如，如果消息交换比服务器故障间隔时间长，候选人将没有足够长的时间来赢得选举；没有一个稳定的领导人，Raft 将无法工作。

领导人选举是 Raft 中对时间要求最为关键的方面。Raft 可以选举并维持一个稳定的领导人，只要系统满足下面的时间要求：

$$\text{广播时间 (broadcastTime)} \ll \text{选举超时时间 (electionTimeout)} \ll \text{平均故障间隔时间 (MTBF)}$$

在这个不等式中，广播时间指的是从一个服务器并行的发送 RPCs 给集群中的其他服务器并接收响应的平均时间；选举超时时间就是在 5.2 节中介绍的选举的超时时间限制；然后平均故障间隔时间就是对于一台服务器而言，两次故障之间的平均时间。广播时间必须比选举超时时间小一个量级，这样领导人才能够发送稳定的心跳消息来阻止跟随者开始进入选举状态；通过随机化选举超时时间的方法，这个不等式也使得选票瓜分的情况变得不可能。选举超时时间应该要比平均故障间隔时间小上几个数量级，这样整个系统才能稳定的运行。当领导人崩溃后，整个系统会大约相当于选举超时的时间里不可用；我们希望这种情况在整个系统的运行中很少出现。

广播时间和平均故障间隔时间是由系统决定的，但是选举超时时间是我们自己选择的。Raft 的 RPCs 需要接收方将信息持久化的保存到稳定存储中去，所以广播时间大约是 0.5 毫秒到 20 毫秒，取决于存储的技术。因此，选举超时时间可能需要在 10 毫秒到 500 毫秒之间。大多数的服务器的平均故障间隔时间都在几个月甚至更长，很容易满足时间的需求。

6 集群成员变化

到目前为止，我们都假设集群的配置（加入到一致性算法的服务器集合）是固定不变的。但是在实践中，偶尔是会改变集群的配置的，例如替换那些宕机的机器或者改变复制级别。尽管可以通过暂停整个集群，更新所有配置，然后重启整个集群的方式来实现，但是在更改的时候集群会不可用。另外，如果存在手工操作步骤，那么就会有操作失误的风险。为了避免这样的问题，我们决定自动化配置改变并且将其纳入到 Raft 一致性算法中来。

为了让配置修改机制能够安全，那么在转换的过程中不能够存在任何时间点使得两个领导人同时被选举成功在同一个任期里。不幸的是，任何服务器直接从旧的配置直接转换到新的配置的方案都是不安全的。一次性自动的转换所有服务器是不可能的，所以在转换期间整个集群存在划分成两个独立的大多数群体的可能性（见图 10）。

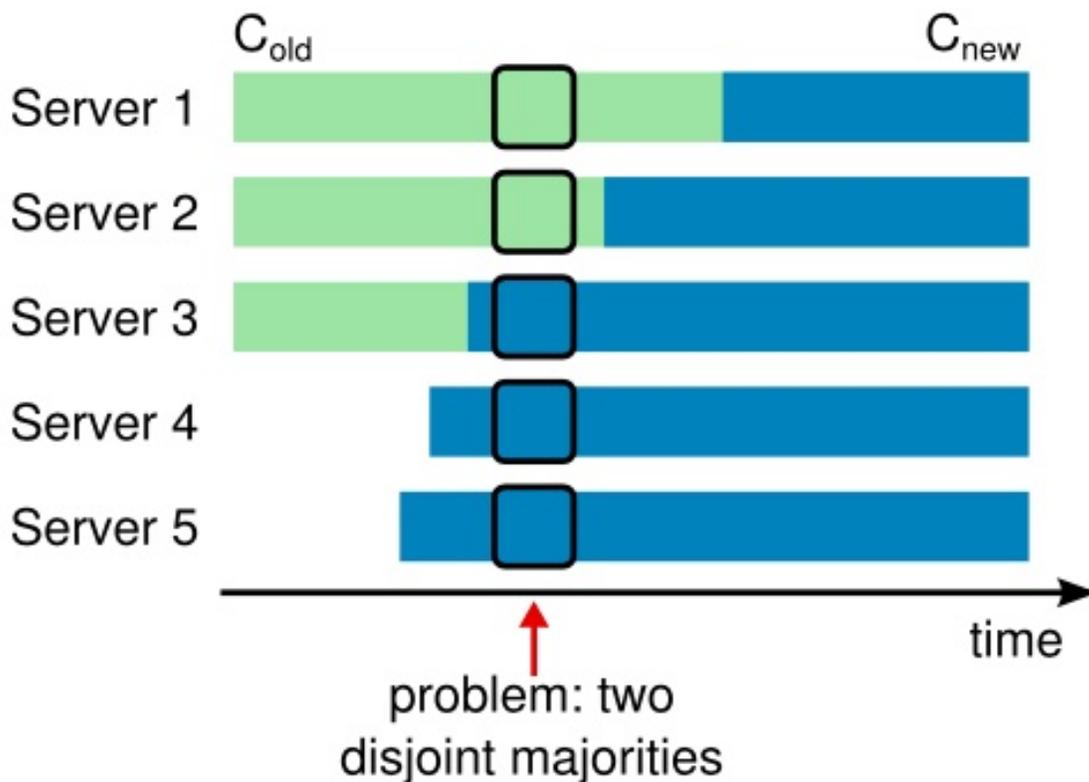


图 10：直接从一种配置转到新的配置是十分不安全的，因为各个机器可能在任何时候进行转换。在这个例子中，集群配额从 3 台机器变成了 5 台。不幸的是，存在这样的一个时间点，两个不同的领导人在同一个任期里都可以被选举成功。一个是通过旧的配置，一个通过新的配置。

为了保证安全性，配置更改必须使用两阶段方法。目前有很多种两阶段的实现。例如，有些系统在第一阶段停掉旧的配置所以集群就不能处理客户端请求；然后在第二阶段在启用新的配置。在 Raft 中，集群先切换到一个过渡的配置，我们称之为共同一致；一旦共同一致已经被提交了，那么系统就切换到新的配置上。共同一致是老配置和新配置的结合：

- 日志条目被复制给集群中新、老配置的所有服务器。
- 新、旧配置的服务器都可以成为领导人。
- 达成一致（针对选举和提交）需要分别在两种配置上获得大多数的支持。

共同一致允许独立的服务器在不影响安全性的前提下，在不同的时间进行配置转换过程。此外，共同一致可以让集群在配置转换的过程中依然响应客户端的请求。

集群配置在复制日志中以特殊的日志条目来存储和通信；图 11 展示了配置转换的过程。当一个领导人接收到一个改变配置从 C_{old} 到 C_{new} 的请求，他会为了共同一致存储配置（图中的 $C_{old,new}$ ），以前面描述的日志条目和副本的形式。一旦一个服务器将新的配置日志条目增加到它的日志中，他就会用这个配置来做出未来所有的决定（服务器总是使用最新的配置，无论他是否已经被提交）。这意味着领导人要使用 $C_{old,new}$ 的规则来决定日志条目 $C_{old,new}$ 什么时候需要被提交。如果领导人崩溃了，被选出来的新领导人可能是使用 C_{old} 配置也可能是 $C_{old,new}$ 配置，这取决于赢得选举的候选人是否已经接收到了 $C_{old,new}$ 配置。在任何情况下， C_{new} 配置在这一时期都不会单方面的做出决定。

一旦 $C_{old,new}$ 被提交，那么无论是 C_{old} 还是 C_{new} ，在没有经过他人批准的情况下都不可能做出决定，并且领导人完全特性保证了只有拥有 $C_{old,new}$ 日志条目的服务器才有可能被选举为领导人。这个时候，领导人创建一条关于 C_{new} 配置的日志条目并复制给集群就是安全的了。再者，每个服务器在见到新的配置的时候就会立即生效。当新的配置在 C_{new} 的规则下被提交，旧的配置就变得无关紧要，同时不使用新的配置的服务器就可以被关闭了。如图 11， C_{old} 和 C_{new} 没有任何机会同时做出单方面的决定；这保证了安全性。

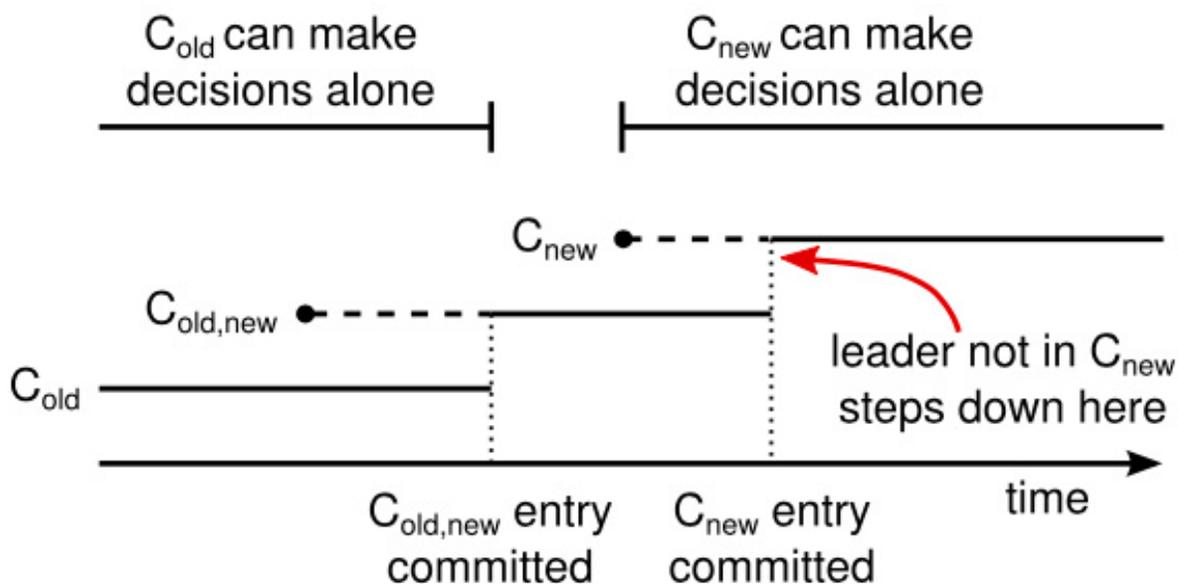


图 11: 一个配置切换的时间线。虚线表示已经被创建但是还没有被提交的条目，实线表示最后被提交的日志条目。领导人首先创建了 $C_{old,new}$ 的配置条目在自己的日志中，并提交到 $C_{old,new}$ 中 (C_{old} 的大多数和 C_{new} 的大多数)。然后他创建 C_{new} 条目并提交到 C_{new} 中的大多数。这样就不存在 C_{new} 和 C_{old} 可以同时做出决定的时间点。

在关于重新配置还有三个问题需要提出。第一个问题是，新的服务器可能初始化没有存储任何的日志条目。当这些服务器以这种状态加入到集群中，那么他们需要一段时间来更新追赶，这时还不能提交新的日志条目。为了避免这种可用性的间隔时间，Raft 在配置更新的时候使用了一种额外的阶段，在这个阶段，新的服务器以没有投票权身份加入到集群中来（领导人复制日志给他们，但是不考虑他们是大多数）。一旦新的服务器追赶上了集群中的其他机器，重新配置可以像上面描述的一样处理。

第二个问题是，集群的领导人可能不是新配置的一员。在这种情况下，领导人就会在提交了 C_{new} 日志之后退位（回到跟随者状态）。这意味着有这样的一段时间，领导人管理着集群，但是不包括他自己；他复制日志但是不把他自己算作是大多数之一。当 C_{new} 被提交时，会发生领导人过渡，因为这时是最早新的配置可以独立工作的时间点（将总是能够在 C_{new} 配置下选出新的领导人）。在此之前，可能只能从 C_{old} 中选出领导人。

第三个问题是，移除不在 C_{new} 中的服务器可能会扰乱集群。这些服务器将不会再接收到心跳，所以当选举超时，他们就会进行新的选举过程。他们会发送拥有新的任期号的请求投票 RPCs，这样会导致当前的领导人回退成跟随者状态。新的领导人最终会被选出来，但是被移除的服务器将会再次超时，然后这个过程会再次重复，导致整体可用性大幅降低。

为了避免这个问题，当服务器确认当前领导人存在时，服务器会忽略请求投票 RPCs。特别的，当服务器在

当前最小选举超时时间内收到一个请求投票 RPC，他不会更新当前的任期号或者投出选票。这不会影响正常的选举，每个服务器在开始一次选举之前，至少等待一个最小选举超时时间。然而，这有利于避免被移除的服务器扰乱：如果领导人能够发送心跳给集群，那么他就不会被更大的任期号废黜。

7 日志压缩

Raft 的日志在正常操作中不断增长，但是在实际的系统中，日志不能无限制的增长。随着日志不断增长，他会占用越来越多的空间，花费越来越多的时间来重置。如果没有一定的机制去清除日志里积累的陈旧的信息，那么会带来可用性问题。

快照是最简单的压缩方法。在快照系统中，整个系统的状态都以快照的形式写入到稳定的持久化存储中，然后到那个时间点之前的日志全部丢弃。快照技术被使用在 Chubby 和 ZooKeeper 中，接下来的章节会介绍 Raft 中的快照技术。

增量压缩的方法，例如日志清理或者日志结构合并树，都是可行的。这些方法每次只对一小部分数据进行操作，这样就分散了压缩的负载压力。首先，他们先选择一个已经积累的大量已经被删除或者被覆盖对象的区域，然后重写那个区域还活跃的对象，之后释放那个区域。和简单操作整个数据集合的快照相比，需要增加复杂的机制来实现。状态机可以实现 LSM tree 使用和快照相同的接口，但是日志清除方法就需要修改 Raft 了。

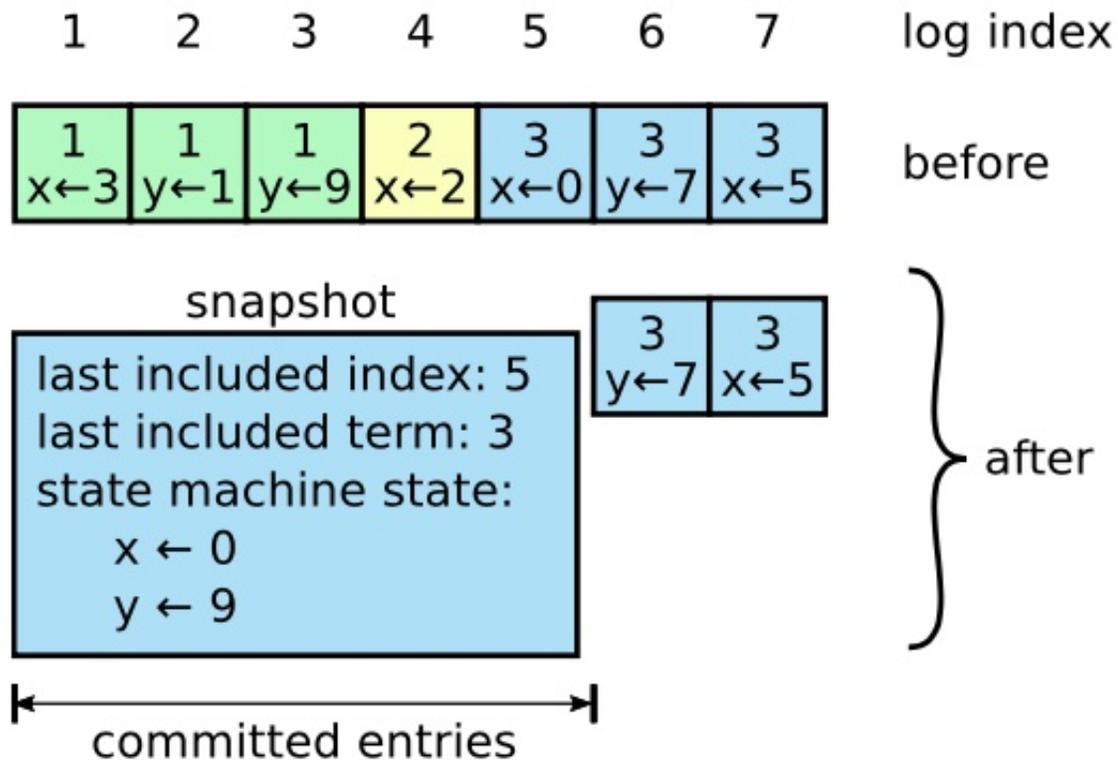


图 12: 一个服务器用新的快照替换了从 1 到 5 的条目，快照值存储了当前的状态。快照中包含了最后的索引位置和任期号。

图 12 展示了 Raft 中快照的基础思想。每个服务器独立的创建快照，只包括已经被提交的日志。主要的工作包括将状态机的状态写入到快照中。Raft 也包含一些少量的元数据到快照中：**最后被包含索引**指的是被快照取代的最后的条目在日志中的索引值（状态机最后应用的日志），**最后被包含的任期**指的是该条目的任期号。保留这些数据是为了支持快照后紧接着的第一个条目的附加日志请求时的一致性检查，因为这个条目需要前一日志条目的索引值和任期号。为了支持集群成员更新（第 6 节），快照中也将最后的一次配置作为最后一个条目存下来。一旦服务器完成一次快照，他就可以删除最后索引位置之前的所有日志和快照了。

尽管通常服务器都是独立的创建快照，但是领导人必须偶尔的发送快照给一些落后的跟随者。这通常发生在当领导人已经丢弃了下一条需要发送给跟随者的日志条目的时候。幸运的是这种情况不是常规操作：一个与领导人保持同步的跟随者通常都会有这个条目。然而一个运行非常缓慢的跟随者或者新加入集群的服务器（第 6 节）将不会有这个条目。这时让这个跟随者更新到最新的状态的方式就是通过网络把快照发送给他们。

安装快照 RPC：

由领导人调用以将快照的分块发送给跟随者。领导者总是按顺序发送分块。

接收者实现：

1. 如果 `term < currentTerm` 就立即回复
2. 如果是第一个分块（`offset` 为 0）就创建一个新的快照
3. 在指定偏移量写入数据
4. 如果 `done` 是 `false`，则继续等待更多的数据
5. 保存快照文件，丢弃具有较小索引的任何现有或部分快照
6. 如果现存的日志条目与快照中最后包含的日志条目具有相同的索引值和任期号，则保留其后的日志条目并进行回复
7. 丢弃整个日志
8. 使用快照重置状态机（并加载快照的集群配置）

InstallSnapshot RPC

Invoked by leader to send chunks of a snapshot to a follower.
Leaders always send chunks in order.

Arguments:

term	leader's term
leaderId	so follower can redirect clients
lastIncludedIndex	the snapshot replaces all entries up through and including this index
lastIncludedTerm	term of lastIncludedIndex
offset	byte offset where chunk is positioned in the snapshot file
data[]	raw bytes of the snapshot chunk, starting at offset
done	true if this is the last chunk

Results:

term	currentTerm, for leader to update itself
-------------	--

Receiver implementation:

1. Reply immediately if term < currentTerm
2. Create new snapshot file if first chunk (offset is 0)
3. Write data into snapshot file at given offset
4. Reply and wait for more data chunks if done is false
5. Save snapshot file, discard any existing or partial snapshot with a smaller index
6. If existing log entry has same index and term as snapshot's last included entry, retain log entries following it and reply
7. Discard the entire log
8. Reset state machine using snapshot contents (and load snapshot's cluster configuration)

图 13: 一个关于安装快照的简要概述。为了便于传输, 快照都是被成分块的; 每个分块都给了跟随者生命的迹象, 所以跟随者可以重置选举超时计时器。

在这种情况下领导人使用一种叫做安装快照的新的 RPC 来发送快照给太落后的跟随者；见图 13。当跟随者通过这种 RPC 接收到快照时，他必须自己决定对于已经存在的日志该如何处理。通常快照会包含没有在接收者日志中存在的信息。在这种情况下，跟随者丢弃其整个日志；它全部被快照取代，并且可能包含与快照冲突的未提交条目。如果接收到的快照是自己日志的前面部分（由于网络重传或者错误），那么被快照包含的条目将会被全部删除，但是快照后面的条目仍然有效，必须保留。

这种快照的方式背离了 Raft 的强领导人原则，因为跟随者可以在不知道领导人情况下创建快照。但是我们认为这种背离是值得的。领导人的存在，是为了解决在达成一致性的时候的冲突，但是在创建快照的时候，一致性已经达成，这时不存在冲突了，所以没有领导人也是可以的。数据依然是从领导人传给跟随者，只是跟随者可以重新组织他们的数据了。

我们考虑过一种替代的基于领导人的快照方案，即只有领导人创建快照，然后发送给所有的跟随者。但是这样做有两个缺点。第一，发送快照会浪费网络带宽并且延缓了快照处理的时间。每个跟随者都已经拥有了所有产生快照需要的信息，而且很显然，自己从本地的状态中创建快照比通过网络接收别人发来的要经济。第二，领导人的实现会更加复杂。例如，领导人需要发送快照的同时并行的将新的日志条目发送给跟随者，这样才不会阻塞新的客户端请求。

还有两个问题影响了快照的性能。首先，服务器必须决定什么时候应该创建快照。如果快照创建的过于频繁，那么就会浪费大量的磁盘带宽和其他资源；如果创建快照频率太低，他就要承受耗尽存储容量的风险，同时也增加了从日志重建的时间。一个简单的策略就是当日志大小达到一个固定大小的时候就创建一次快照。如果这个阈值设置的显著大于期望的快照的大小，那么快照对磁盘压力的影响就会很小了。

第二个影响性能的问题就是写入快照需要花费显著的一段时间，并且我们还不希望影响到正常操作。解决方案是通过写时复制的技术，这样新的更新就可以被接收而不影响到快照。例如，具有函数式数据结构的状态机天然支持这样的功能。另外，操作系统的写时复制技术的支持（如 Linux 上的 fork）可以被用来创建完整的状态机的内存快照（我们的实现就是这样的）。

8 客户端交互

这一节将介绍客户端是如何和 Raft 进行交互的，包括客户端如何发现领导人和 Raft 是如何支持线性化语义的。这些问题对于所有基于一致性的系统都存在，并且 Raft 的解决方案和其他的也差不多。

Raft 中的客户端发送所有请求给领导人。当客户端启动的时候，他会随机挑选一个服务器进行通信。如果客户端第一次挑选的服务器不是领导人，那么那个服务器会拒绝客户端的请求并且提供他最近接收到的领导人的信息（附加条目请求包含了领导人的网络地址）。如果领导人已经崩溃了，那么客户端的请求就会超时；客户端之后会再次重试随机挑选服务器的过程。

我们 Raft 的目标是要实现线性化语义（每一次操作立即执行，只执行一次，在他调用和收到回复之间）。但是，如上述，Raft 是可以执行同一条命令多次的：例如，如果领导人在提交了这条日志之后，但是在响应客户端之前崩溃了，那么客户端会和新的领导人重试这条指令，导致这条命令就被再次执行了。解决方案就是客户端对于每一条指令都赋予一个唯一的序列号。然后，状态机跟踪每条指令最新的序列号和相应的响应。如果接收到一条指令，它的序列号已经被执行了，那么就立即返回结果，而不重新执行指令。

只读的操作可以直接处理而不需要记录日志。但是，在不增加任何限制的情况下，这么做可能会冒着返回脏数据的风险，因为领导人响应客户端请求时可能已经被新的领导人作废了，但是他还不知道。线性化的读操作必须不能返回脏数据，Raft 需要使用两个额外的措施在不使用日志的情况下保证这一点。首先，领导人必须有关于被提交日志的最新信息。领导人完全特性保证了领导人一定拥有所有已经被提交的日志条目，但是在他任期开始的时候，他可能不知道那些是已经被提交的。为了知道这些信息，他需要在他的任期里提交一条日志条目。Raft 中通过领导人在任期开始的时候提交一个空白的没有任何操作的日志条目到日志中去来实现。第二，领导人在处理只读的请求之前必须检查自己是否已经被废黜了（他自己的信息已经变脏了如果一个更新的领导人被选举出来）。Raft 中通过让领导人在响应只读请求之前，先和集群中的大多数节点交换一次心跳信息来处理这个问题。可选的，领导人可以依赖心跳机制来实现一种租约的机制，但是这种方法依赖时间来保证安全性（假设时间误差是有界的）。

9 算法实现和评估

我们已经为 RAMCloud 实现了 Raft 算法作为存储配置信息的复制状态机的一部分，并且帮助 RAMCloud 协调故障转移。这个 Raft 实现包含大约 2000 行 C++ 代码，其中不包括测试、注释和空行。这些代码是开源的。同时也有大约 25 个其他独立的第三方的基于这篇论文草稿的开源实现，针对不同的开发场景。同时，很多公司已经部署了基于 Raft 的系统。

这一节会从三个方面来评估 Raft 算法：可理解性、正确性和性能。

9.1 可理解性

为了和 Paxos 比较 Raft 算法的可理解能力，我们针对高层次的本科生和研究生，在斯坦福大学的高级操作系统课程和加州大学伯克利分校的分布式计算课程上，进行了一次学习的实验。我们分别拍了针对 Raft 和 Paxos 的视频课程，并准备了相应的小测验。Raft 的视频讲课覆盖了这篇论文的所有内容除了日志压缩；Paxos 讲课包含了足够的资料来创建一个等价的复制状态机，包括单决策 Paxos，多决策 Paxos，重新配置和一些实际系统需要的性能优化（例如领导人选举）。小测验测试一些对算法的基本理解和解释一些边角的示例。每个学生都是看完第一个视频，回答相应的测试，再看第二个视频，回答相应的测试。大约有一半的学生先进行 Paxos 部分，然后另一半先进行 Raft 部分，这是为了说明两者从第一部分的算法学习中获得的表现和经验的差异。我们计算参加人员的每一个小测验的得分来看参与者是否在 Raft 算法上更加容易理解。

我们尽可能的使得 Paxos 和 Raft 的比较更加公平。这个实验偏爱 Paxos 表现在两个方面：43 个参加者中有 15 个人在之前有一些 Paxos 的经验，并且 Paxos 的视频要长 14%。如表格 1 总结的那样，我们采取了一些措施来减轻这种潜在的偏见。我们所有的材料都可供审查。

表 1: 考虑到可能会存在的偏见，对于每种情况的解决方法，和相应的材料。

参加者平均在 Raft 的测验中比 Paxos 高 4.9 分（总分 60，那么 Raft 的平均得分是 25.7，而 Paxos 是 20.8）；图 14 展示了每个参与者的得分。配置 t-检验（又称 student 's t-test）表明，在 95% 的可信度下，真实的 Raft 分数分布至少比 Paxos 高 2.5 分。

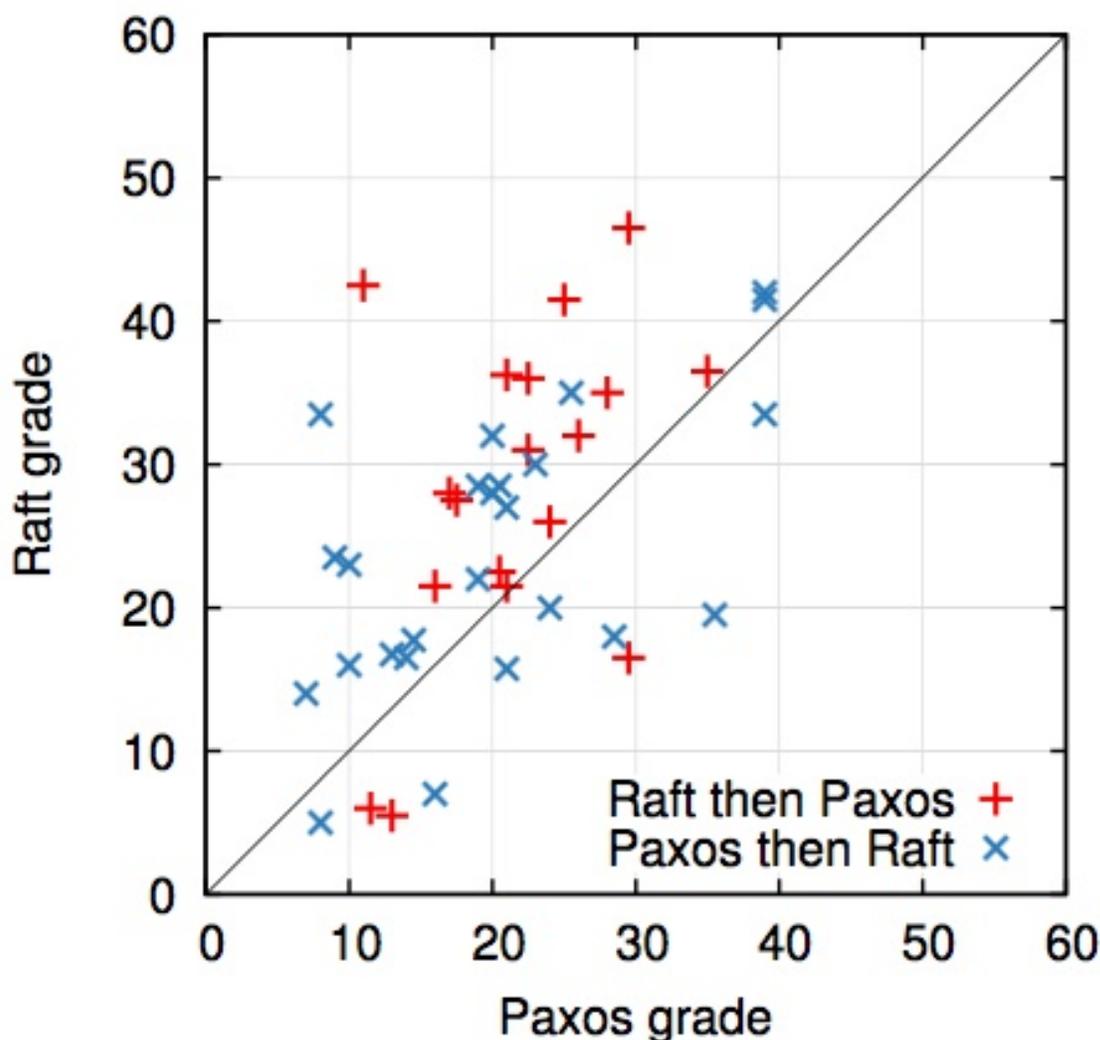


图 14: 一个散点图表示了 43 个学生在 Paxos 和 Raft 的小测验中的成绩。在对角线之上的点表示在 Raft 获得了更高分数的学生。

我们也建立了一个线性回归模型来预测一个新的学生的测验成绩，基于以下三个因素：他们使用的是哪个小测验，之前对 Paxos 的经验，和学习算法的顺序。模型预测，对小测验的选择会产生 12.5 分的差别。这显著的高于之前的 4.9 分，因为很多学生在之前都已经有了对于 Paxos 的经验，这相当明显的帮助 Paxos，对 Raft 就没什么太大影响了。但是奇怪的是，模型预测对于先进行 Paxos 小测验的人而言，Raft 的得分低了 6.3 分；虽然我们不知道为什么，这似乎在统计上是有意义的。

我们同时也在测验之后调查了参与者，他们认为哪个算法更加容易实现和解释；这个的结果在图 15 上。压倒性的结果表明 Raft 算法更加容易实现和解释（41 人中的 33 个）。但是，这种自己报告的结果不如参与者的成绩更加可信，并且参与者可能因为我们的 Raft 更加易于理解的假说而产生偏见。

图 15: 通过一个 5 分制的问题，参与者（左边）被问哪个算法他们觉得在一个高效正确的系统里更容易实现，右边被问哪个更容易向学生解释。

关于 Raft 用户学习有一个更加详细的讨论。

9.2 正确性

在第 5 节，我们已经制定了正式的规范，和对一致性机制的安全性证明。这个正式规范使用 TLA+ 规范语言使图 2 中总结的信息非常清晰。它长约 400 行，并作为证明的主题。同时对于任何想实现 Raft 的人也是十分有用的。我们通过 TLA 证明系统非常机械的证明了日志完全特性。然而，这个证明依赖的约束前提还没有被机械证明（例如，我们还没有证明规范的类型安全）。而且，我们已经写了一个非正式的证明关于状态机安全性是完备的，并且是相当清晰的（大约 3500 个词）。

9.3 性能

Raft 和其他一致性算法例如 Paxos 有着差不多的性能。在性能方面，最重要的关注点是，当领导人被选举成功时，什么时候复制新的日志条目。Raft 通过很少数量的消息包（一轮从领导人到集群大多数机器的消息）就达成了这个目的。同时，进一步提升 Raft 的性能也是可行的。例如，很容易通过支持批量操作和管道操作来提高吞吐量和降低延迟。对于其他一致性算法已经提出过很多性能优化方案；其中有很多也可以应用到 Raft 中来，但是我们暂时把这个问题放到未来的工作中去。

我们使用我们自己的 Raft 实现来衡量 Raft 领导人选举的性能并且回答两个问题。首先，领导人选举的过程收敛是否快速？第二，在领导人宕机之后，最小的系统宕机时间是多久？

图 16：发现并替换一个已经崩溃的领导人的时间。上面的图考察了在选举超时时间上的随机化程度，下面的图考察了最小选举超时时间。每条线代表了 1000 次实验（除了 150-150 毫秒只试了 100 次），和相应的确定的选举超时时间。例如，150-155 毫秒意思是，选举超时时间从这个区间范围内随机选择并确定下来。这个实验在一个拥有 5 个节点的集群上进行，其广播时延大约是 15 毫秒。对于 9 个节点的集群，结果也差不多。

为了衡量领导人选举，我们反复的使一个拥有五个节点的服务器集群的领导人宕机，并计算需要多久才能发现领导人已经宕机并选出一个新的领导人（见图 16）。为了构建一个最坏的场景，在每一的尝试里，服务器都有不同长度的日志，意味着有些候选人是没有成为领导人的资格的。另外，为了促成选票瓜分的情况，我们的测试脚本在终止领导人之前同步的发送了一次心跳广播（这大约和领导人在崩溃前复制一个新的日志给其他机器很像）。领导人均匀的随机的在心跳间隔里宕机，也就是最小选举超时时间的一半。因此，最小宕机时间大约就是最小选举超时时间的一半。

图 16 中上面的图表明，只需要在选举超时时间上使用很少的随机化就可以大大避免选票被瓜分的情况。在没有随机化的情况下，在我们的测试里，选举过程往往都需要花费超过 10 秒钟由于太多的选票瓜分的情况。仅仅增加 5 毫秒的随机化时间，就大大的改善了选举过程，现在平均的宕机时间只有 287 毫秒。增加更多的随机化时间可以大大改善最坏情况：通过增加 50 毫秒的随机化时间，最坏的完成情况（1000 次尝试）只要 513 毫秒。

图 16 中下面的图显示，通过减少选举超时时间可以减少系统的宕机时间。在选举超时时间为 12-24 毫秒的情况下，只需要平均 35 毫秒就可以选举出新的领导人（最长的一次花费了 152 毫秒）。然而，进一步降低选举超时时间的话就会违反 Raft 的时间不等式需求：在选举新领导人之前，领导人就很难发送完心跳包。这会导致没有意义的领导人改变并降低了系统整体的可用性。我们建议使用更为保守的选举超时时间，比如 150-300 毫秒；这样的时间不大可能导致没有意义的领导人改变，而且依然提供不错的可用性。

10 相关工作

已经有很多关于一致性算法的工作被发表出来，其中很多都可以归到下面的类别中：

- Lamport 关于 Paxos 的原始描述，和尝试描述的更清晰。
- 关于 Paxos 的更详尽的描述，补充遗漏的细节并修改算法，使得可以提供更加容易的实现基础。
- 实现一致性算法的系统，例如 Chubby, ZooKeeper 和 Spanner。对于 Chubby 和 Spanner 的算法并没有公开发表其技术细节，尽管他们都声称是基于 Paxos 的。ZooKeeper 的算法细节已经发表，但是和 Paxos 着实有着很大的差别。
- Paxos 可以应用的性能优化。
- Oki 和 Liskov 的 Viewstamped Replication (VR)，一种和 Paxos 差不多的替代算法。原始的算法描述和分布式传输协议耦合在了一起，但是核心的一致性算法在最近的更新里被分离了出来。VR 使用了一种基于领导人的方法，和 Raft 有很多相似之处。

Raft 和 Paxos 最大的不同之处就在于 Raft 的强领导特性：Raft 使用领导人选举作为一致性协议里必不可少的部分，并且将尽可能多的功能集中到了领导人身上。这样就可以使得算法更加容易理解。例如，在 Paxos 中，领导人选举和基本的一致性协议是正交的：领导人选举仅仅是性能优化的手段，而且不是一致性所必须要求的。但是，这样就增加了多余的机制：Paxos 同时包含了针对基本一致性要求的两阶段提交协议和针对领导人选举的独立的机制。相比较而言，Raft 就直接将领导人选举纳入到一致性算法中，并作为两阶段一致性的第一步。这样就减少了很多机制。

像 Raft 一样，VR 和 ZooKeeper 也是基于领导人的，因此他们也拥有一些 Raft 的优点。但是，Raft 比 VR 和 ZooKeeper 拥有更少的机制因为 Raft 尽可能的减少了非领导人的功能。例如，Raft 中日志条目都遵循着从领导人发送给其他人这一个方向：附加条目 RPC 是向外发送的。在 VR 中，日志条目的流动是双向的（领导人可以在选举过程中接收日志）；这就导致了额外的机制和复杂性。根据 ZooKeeper 公开的资料看，它的日志条目也是双向传输的，但是它的实现更像 Raft。

和上述我们提及的其他基于一致性的日志复制算法中，Raft 的消息类型更少。例如，我们数了一下 VR 和 ZooKeeper 使用的用来基本一致性需要和成员改变的消息数（排除了日志压缩和客户端交互，因为这些都比较独立且和算法关系不大）。VR 和 ZooKeeper 都分别定义了 10 中不同的消息类型，相对的，Raft 只有 4 中消息类型（两种 RPC 请求和对应的响应）。Raft 的消息都稍微比其他算法的要信息量大，但是都很简单。另外，VR 和 ZooKeeper 都在领导人改变时传输了整个日志；所以为了能够实践中使用，额外的消息类型就很必要了。

Raft 的强领导人模型简化了整个算法，但是同时也排斥了一些性能优化的方法。例如，平等主义 Paxos (EPaxos) 在某些没有领导人的情况下可以达到很高的性能。平等主义 Paxos 充分发挥了在状态机指令中的交换性。任何服务器都可以在一轮通信下就提交指令，除非其他指令同时被提出了。然而，如果指令都是并发的被提出，并且互相之间不通信沟通，那么 EPaxos 就需要额外的一轮通信。因为任何服务器都可以提交指令，所以 EPaxos 在服务器之间的负载均衡做的很好，并且很容易在 WAN 网络环境下获得很低的延迟。但是，他在 Paxos 上增加了非常明显的复杂性。

一些集群成员变换的方法已经被提出或者在其他的工作中被实现，包括 Lamport 的原始的讨论，VR 和 SMART。我们选择使用共同一致的方法因为他对一致性协议的其他部分影响很小，这样我们只需要很少的一些机制就可以实现成员变换。Lamport 的基于 α 的方法之所以没有被 Raft 选择是因为它假设在没有领导人的情况下也可以达到一致性。和 VR 和 SMART 相比较，Raft 的重新配置算法可以在不限制正常请求处理的情况下进行；相比较的，VR 需要停止所有的处理过程，SMART 引入了一个和 α 类似的方法，限制了请求处理的数量。Raft 的方法同时也需要更少的额外机制来实现，和 VR、SMART 比较而言。

11 结论

算法的设计通常会把正确性，效率或者简洁作为主要的目标。尽管这些都是很有意义的目标，但是我们相信，可理解性也是一样的重要。在开发者把算法应用到实际的系统中之前，这些目标没有一个会被实现，这些都会必然的偏离发表时的形式。除非开发人员对这个算法有着很深的理解并且有着直观的感觉，否则将会对他们而言很难在实现的时候保持原有期望的特性。

在这篇论文中，我们尝试解决分布式一致性问题，但是一个广为接受但是十分令人费解的算法 Paxos 已经困扰了无数学生和开发者很多年了。我们创造了一种新的算法 Raft，显而易见的比 Paxos 要容易理解。我们同时也相信，Raft 也可以为实际的实现提供坚实的基础。把可理解性作为设计的目标改变了我们设计 Raft 的方式；随着设计的进展，我们发现自己重复使用了一些技术，比如分解问题和简化状态空间。这些技术不仅提升了 Raft 的可理解性，同时也使我们坚信其正确性。

12 感谢

这项研究必须感谢以下人员的支持：Ali Ghodsi, David Mazie'res, 和伯克利 CS 294-91 课程、斯坦福 CS 240 课程的学生。Scott Klemmer 帮我们设计了用户调查，Nelson Ray 建议我们进行统计学的分析。在用户调查时使用的关于 Paxos 的幻灯片很大一部分是从 Lorenzo Alvisi 的幻灯片上借鉴过来的。特别的，非常感谢 DavidMazieres 和 Ezra Hoch，他们找到了 Raft 中一些难以发现的漏洞。许多人提供了关于这篇论文十分有用的反馈和用户调查材料，包括 Ed Bugnion, Michael Chan, Hugues Evrard, Daniel Giffin, Arjun Gopalan, Jon Howell, Vimalkumar Jeyakumar, Ankita Kejriwal, Aleksandar Kracun, Amit Levy, Joel Martin, Satoshi Matsushita, Oleg Pesok, David Ramos, Robbert van Renesse, Mendel Rosenblum, Nicolas Schiper, Deian Stefan, Andrew Stone, Ryan Stutsman, David Terei, Stephen Yang, Matei Zaharia 以及 24 位匿名的会议审查人员（可能有重复），并且特别感谢我们的领导人 Eddie Kohler。Werner Vogels 发了一条早期草稿链接的推特，给 Raft 带来了极大的关注。我们的工作由 Gigascale 系统研究中心和 Multiscale 系统研究中心给予支持，这两个研究中心由关注中心研究程序资金支持，一个是半导体研究公司的程序，由 STARnet 支持，一个半导体研究公司的程序由 MARCO 和 DARPA 支持，在国家科学基金会的 0963859 号批准，并且获得了来自 Facebook, Google, Mellanox, NEC, NetApp, SAP 和 Samsung 的支持。Diego Ongaro 由 Jungle 公司，斯坦福的毕业团体支持。

原文链接

https://github.com/maemual/raft-zh_cn/blob/master/raft-zh_cn.md

2.1.7 转载 PoW 挖矿算法原理及其在比特币、以太坊中的实现

PoW，全称 Proof of Work，即工作量证明，又称挖矿。大部分公有链或虚拟货币，如比特币、以太坊，均基于 PoW 算法，来实现其共识机制。即根据挖矿贡献的有效工作，来决定货币的分配。

比特币区块

比特币区块由区块头和该区块所包含的交易列表组成。区块头大小为 80 字节，其构成包括：

4 字节：版本号 32 字节：上一个区块的哈希值 32 字节：交易列表的 Merkle 根哈希值 4 字节：当前时间戳 4 字节：当前难度值 4 字节：随机数 Nonce 值

此 80 字节长度的区块头，即为比特币 Pow 算法的输入字符串。交易列表附加在区块头之后，其中第一笔交易为矿工获得奖励和手续费的特殊交易。

bitcoin-0.15.1 源码中区块头和区块定义：

```

class CBlockHeader
{
public:
    int32_t nVersion;                //版本号

    uint256 hashPrevBlock;          //上一个区块的哈希值

    uint256 hashMerkleRoot;        //交易列表的Merkle根哈希值

    uint32_t nTime;                //当前时间戳

    uint32_t nBits;                //
    ↪ 当前挖矿难度, nBits越小难度越大

    uint32_t nNonce;               //随机数Nonce值
    //其它代码略
};
class CBlock : public CBlockHeader
{
public:
    //交易列表
    std::vector<CTransactionRef> vtx;
    //其它代码略
};
//代码位置src/primitives/block.h

```

比特币 Pow 算法原理

Pow 的过程, 即为不断调整 Nonce 值, 对区块头做双重 SHA256 哈希运算, 使得结果满足给定数量前导 0 的哈希值的过程。其中前导 0 的个数, 取决于挖矿难度, 前导 0 的个数越多, 挖矿难度越大。

具体如下:

- 1、生成铸币交易, 并与其它所有准备打包进区块的交易组成交易列表, 生成 Merkle 根哈希值。
- 2、将 Merkle 根哈希值, 与区块头其它字段组成区块头, 80 字节长度的区块头作为 Pow 算法的输入。
- 3、不断变更区块头中的随机数 Nonce, 对变更后的区块头做双重 SHA256 哈希运算, 与当前难度的目标值做比对, 如果小于目标难度, 即 Pow 完成。

Pow 完成的区块向全网广播, 其他节点将验证其是否符合规则, 如果验证有效, 其他节点将接收此区块, 并附加在已有区块链之后。之后将进入下一轮挖矿。

bitcoin-0.15.1 源码中 Pow 算法实现:

```

UniValue generateBlocks(std::shared_ptr<CReserveScript> coinbaseScript, int nGenerate,
    ↪ uint64_t nMaxTries, bool keepScript)
{
    static const int nInnerLoopCount = 0x10000;
    int nHeightEnd = 0;
    int nHeight = 0;
    { // Don't keep cs_main locked
        LOCK(cs_main);
        nHeight = chainActive.Height();
        nHeightEnd = nHeight+nGenerate;
    }
    unsigned int nExtraNonce = 0;
    UniValue blockHashes(UniValue::VARR);
    while (nHeight < nHeightEnd)

```

(下页继续)

(续上页)

```

{
    std::unique_ptr<CBlockTemplate> pblocktemplate(BlockAssembler(Params()).
↳CreateNewBlock(coinbaseScript->reserveScript));
    if (!pblocktemplate.get())
        throw JSONRPCError(RPC_INTERNAL_ERROR, "Couldn't create new block");
    CBlock *pblock = &pblocktemplate->block;
    {
        LOCK(cs_main);
        IncrementExtraNonce(pblock, chainActive.Tip(), nExtraNonce);
    }
    //不断变更区块头中的随机数Nonce
    //对变更后的区块头做双重SHA256哈希运算
    //与当前难度的目标值做比对, 如果小于目标难度, 即Pow完成
    //uint64_t nMaxTries = 1000000;即重试100万次
    while (nMaxTries > 0 && pblock->nNonce < nInnerLoopCount && !
↳CheckProofOfWork(pblock->GetHash(), pblock->nBits, Params().GetConsensus())) {
        ++pblock->nNonce;
        --nMaxTries;
    }
    if (nMaxTries == 0) {
        break;
    }
    if (pblock->nNonce == nInnerLoopCount) {
        continue;
    }
    std::shared_ptr<const CBlock> shared_pblock = std::make_shared<const CBlock>
↳(*pblock);
    if (!ProcessNewBlock(Params(), shared_pblock, true, nullptr))
        throw JSONRPCError(RPC_INTERNAL_ERROR, "ProcessNewBlock, block not
↳accepted");
    ++nHeight;
    blockHashes.push_back(pblock->GetHash().GetHex());
    //mark script as important because it was used at least for one coinbase
↳output if the script came from the wallet
    if (keepScript)
    {
        coinbaseScript->KeepScript();
    }
    return blockHashes;
}
//代码位置src/rpc/mining.cpp

```

另附 bitcoin-0.15.1 源码中生成铸币交易和创建新块:

```

std::unique_ptr<CBlockTemplate> BlockAssembler::CreateNewBlock(const CScript&
↳scriptPubKeyIn, bool fMineWitnessTx)
{
    int64_t nTimeStart = GetTimeMicros();
    resetBlock();
    pblocktemplate.reset(new CBlockTemplate());
    if (!pblocktemplate.get())
        return nullptr;
    pblock = &pblocktemplate->block; // pointer for convenience
    pblock->vtx.emplace_back();
    pblocktemplate->vTxFees.push_back(-1); // updated at end

```

(下页继续)

```

pblocktemplate->vTxSigOpsCost.push_back(-1); // updated at end
LOCK2(cs_main, mempool.cs);
CBlockIndex* pindexPrev = chainActive.Tip();
nHeight = pindexPrev->nHeight + 1;
//版本号
pblock->nVersion = ComputeBlockVersion(pindexPrev, chainparams.GetConsensus());
if (chainparams.MineBlocksOnDemand())
    pblock->nVersion = gArgs.GetArg("-blockversion", pblock->nVersion);
//当前时间戳
pblock->nTime = GetAdjustedTime();
const int64_t nMedianTimePast = pindexPrev->GetMedianTimePast();
nLockTimeCutoff = (STANDARD_LOCKTIME_VERIFY_FLAGS & LOCKTIME_MEDIAN_TIME_PAST
    ? nMedianTimePast
    : pblock->GetBlockTime());
fIncludeWitness = IsWitnessEnabled(pindexPrev, chainparams.GetConsensus()) &&
↪ fMineWitnessTx;
int nPackagesSelected = 0;
int nDescendantsUpdated = 0;
addPackageTxs(nPackagesSelected, nDescendantsUpdated);
int64_t nTime1 = GetTimeMicros();
nLastBlockTx = nBlockTx;
nLastBlockWeight = nBlockWeight;
//创建铸币交易
CMutableTransaction coinbaseTx;
coinbaseTx.vin.resize(1);
coinbaseTx.vin[0].prevout.SetNull();
coinbaseTx.vout.resize(1);
//挖矿奖励和手续费
coinbaseTx.vout[0].scriptPubKey = scriptPubKeyIn;
coinbaseTx.vout[0].nValue = nFees + GetBlockSubsidy(nHeight, chainparams.
↪ GetConsensus());
coinbaseTx.vin[0].scriptSig = CScript() << nHeight << OP_0;
//第一笔交易即为矿工获得奖励和手续费的特殊交易
pblock->vtx[0] = MakeTransactionRef(std::move(coinbaseTx));
pblocktemplate->vchCoinbaseCommitment = GenerateCoinbaseCommitment(*pblock,
↪ pindexPrev, chainparams.GetConsensus());
pblocktemplate->vTxFees[0] = -nFees;
LogPrintf("CreateNewBlock(): block weight: %u txs: %u fees: %ld sigops %d\n",
↪ GetBlockWeight(*pblock), nBlockTx, nFees, nBlockSigOpsCost);
//上一个区块的哈希值
pblock->hashPrevBlock = pindexPrev->GetBlockHash();
UpdateTime(pblock, chainparams.GetConsensus(), pindexPrev);
//当前挖矿难度
pblock->nBits = GetNextWorkRequired(pindexPrev, pblock, chainparams.
↪ GetConsensus());
//随机数Nonce值
pblock->nNonce = 0;
pblocktemplate->vTxSigOpsCost[0] = WITNESS_SCALE_FACTOR *
↪ GetLegacySigOpCount(*pblock->vtx[0]);
CValidationState state;
if (!TestBlockValidity(state, chainparams, *pblock, pindexPrev, false, false)) {
    throw std::runtime_error(strprintf("%s: TestBlockValidity failed: %s", __func_
↪ __, FormatStateMessage(state)));
}
int64_t nTime2 = GetTimeMicros();
LogPrint(BCLog::BENCH, "CreateNewBlock() packages: %.2fms (%d packages, %d
↪ updated descendants), validity: %.2fms (total %.2fms)\n", 0.001 * (nTime1 -
↪ nTimeStart), nPackagesSelected, nDescendantsUpdated, 0.001 * (nTime2 - nTime1),
↪ 0.001 * (nTime2 - nTimeStart));

```

(续上页)

```

    return std::move(pblocktemplate);
}
//代码位置src/miner.cpp

```

比特币挖矿难度计算

每创建 2016 个块后将计算新的难度，此后的 2016 个块使用新的难度。计算步骤如下：

- 1、找到前 2016 个块的第一个块，计算生成这 2016 个块花费的时间。即最后一个块的时间与第一个块的时间差。时间差不小于 3.5 天，不大于 56 天。
- 2、计算前 2016 个块的难度总和，即单个块的难度 x 总时间。
- 3、计算新的难度，即 2016 个块的难度总和/14 天的秒数，得到每秒的难度值。
- 4、要求新的难度，难度不低于参数定义的最小难度。

bitcoin-0.15.1 源码中计算挖矿难度代码如下：

```

//nFirstBlockTime即前2016个块的第一个块的时间戳
unsigned int CalculateNextWorkRequired(const CBlockIndex* pindexLast, int64_t_
↪nFirstBlockTime, const Consensus::Params& params)
{
    if (params.fPowNoRetargeting)
        return pindexLast->nBits;
    //计算生成这2016个块花费的时间
    int64_t nActualTimespan = pindexLast->GetBlockTime() - nFirstBlockTime;
    //不小于3.5天
    if (nActualTimespan < params.nPowTargetTimespan/4)
        nActualTimespan = params.nPowTargetTimespan/4;
    //不大于56天
    if (nActualTimespan > params.nPowTargetTimespan*4)
        nActualTimespan = params.nPowTargetTimespan*4;
    // Retarget
    const arith_uint256 bnPowLimit = UintToArith256(params.powLimit);
    arith_uint256 bnNew;
    bnNew.SetCompact(pindexLast->nBits);
    //计算前2016个块的难度总和
    //即单个块的难度*总时间
    bnNew *= nActualTimespan;
    //计算新的难度
    //即2016个块的难度总和/14天的秒数
    bnNew /= params.nPowTargetTimespan;
    //bnNew越小，难度越大
    //bnNew越大，难度越小
    //要求新的难度，难度不低于参数定义的最小难度
    if (bnNew > bnPowLimit)
        bnNew = bnPowLimit;
    return bnNew.GetCompact();
}
//代码位置src/pow.cpp

```

以太坊区块

以太坊区块由 Header 和 Body 两部分组成。

其中 Header 部分成员如下：ParentHash，父区块哈希 UncleHash，叔区块哈希，具体为 Body 中 Uncles 数组的 RLP 哈希值。RLP 哈希，即某类型对象 RLP 编码后做 SHA3 哈希运算。Coinbase，矿工地址。Root，StateDB 中 state Trie 根节点 RLP 哈希值。TxHash，Block 中 tx Trie 根节点 RLP 哈希值。ReceiptHash，Block 中 Receipt Trie 根节点的 RLP 哈希值。Difficulty，区块难度，即当前挖矿难度。Number，区块序号，即父区块 Number+1。GasLimit，区块内所有 Gas 消耗的理论上限，创建时指定，由父区块 GasUsed 和 GasLimit 计算得出。GasUsed，区块内所有 Transaction 执行时消耗的 Gas 总和。Time，当前时间戳。Nonce，随机数 Nonce 值。

有关叔区块：叔区块，即孤立的块。以太坊成块速度较快，导致产生孤块。以太坊会给发现孤块的矿工以回报，激励矿工在新块中引用孤块，引用孤块使主链更重。在以太坊中，主链是指最重的链。

有关 state Trie、tx Trie 和 Receipt Trie：state Trie，所有账户对象可以逐个插入一个 Merkle-PatriciaTrie(MPT) 结构中，形成 state Trie。tx Trie：Block 中 Transactions 中所有 tx 对象，逐个插入 MPT 结构中，形成 tx Trie。Receipt Trie：Block 中所有 Transaction 执行后生成 Receipt 数组，所有 Receipt 逐个插入 MPT 结构中，形成 Receipt Trie。

Body 成员如下：Transactions，交易列表。Uncles，引用的叔区块列表。

go-ethereum-1.7.3 源码中区块头和区块定义：

```

type Header struct {
    //父区块哈希
    ParentHash common.Hash
    //叔区块哈希
    UncleHash common.Hash
    //矿工地址
    Coinbase common.Address
    //StateDB中state Trie根节点RLP哈希值
    Root common.Hash
    //Block中tx Trie根节点RLP哈希值
    TxHash common.Hash
    //Block中Receipt Trie根节点的RLP哈希值
    ReceiptHash common.Hash
    Bloom Bloom
    //区块难度
    Difficulty *big.Int
    //区块序号
    Number *big.Int
    //区块内所有Gas消耗的理论上限
    GasLimit *big.Int
    //区块内所有Transaction执行时消耗的Gas总和
    GasUsed *big.Int
    //当前时间戳
    Time *big.Int
    Extra []byte
    MixDigest common.Hash
    //随机数Nonce值
    Nonce BlockNonce
}

type Body struct {
    //交易列表
    Transactions []*Transaction
    //引用的叔区块列表
    Uncles []*Header
}

//代码位置core/types/block.go

```

以太坊 Pow 算法原理

以太坊 Pow 算法可以表示为如下公式: $RAND(h, n) \leq M / d$

其中 $RAND()$ 表示一个概念函数, 代表一系列的复杂运算。其中 h 和 n 为输入, 即区块 Header 的哈希、以及 Header 中的 Nonce。 M 表示一个极大的数, 此处使用 $2^{256}-1$ 。 d , 为区块难度, 即 Header 中的 Difficulty。

因此在 h 和 n 确定的情况下, d 越大, 挖矿难度越大, 即为 Difficulty 本义。即不断变更 Nonce, 使 $RAND(h, n)$ 满足 $RAND(h, n) \leq M / d$, 即完成 Pow。

go-ethereum-1.7.3 源码中 Pow 算法实现:

```
func (ethash *Ethash) mine(block *types.Block, id int, seed uint64, abort chan struct
→{}), found chan *types.Block) {
    // Extract some data from the header
    var (
        header = block.Header()
        hash    = header.HashNoNonce().Bytes()
        //target, 即 M / d, 即 (2^256-1)/Difficulty
        target = new(big.Int).Div(maxUint256, header.Difficulty)
        number = header.Number.Uint64()
        dataset = ethash.dataset(number)
    )
    // Start generating random nonces until we abort or find a good one
    var (
        attempts = int64(0)
        nonce     = seed
    )
    logger := log.New("miner", id)
    logger.Trace("Started ethash search for new nonces", "seed", seed)
    for {
        select {
        case <-abort:
            // Mining terminated, update stats and abort
            logger.Trace("Ethash nonce search aborted", "attempts", nonce-seed)
            ethash.hashrate.Mark(attempts)
            return
        default:
            // We don't have to update hash rate on every nonce, so update after
            →after 2^X nonces
            attempts++
            if (attempts % (1 << 15)) == 0 {
                ethash.hashrate.Mark(attempts)
                attempts = 0
            }
            //hashimotoFull即RAND(h, n)所代表的一系列的复杂运算
            digest, result := hashimotoFull(dataset, hash, nonce)
            //result满足RAND(h, n) <= M / d
            if new(big.Int).SetBytes(result).Cmp(target) <= 0 {
                // Correct nonce found, create a new header with it
                header = types.CopyHeader(header)
                header.Nonce = types.EncodeNonce(nonce)
                header.MixDigest = common.BytesToHash(digest)
                // Seal and return a block (if still needed)
                select {
                case found <- block.WithSeal(header):
                    logger.Trace("Ethash nonce found and reported", "attempts", nonce-
                    →seed, "nonce", nonce)
                case <-abort:

```

(下页继续)

(续上页)

```

        logger.Trace("Ethash nonce found but discarded", "attempts", ↵
↵nonce-seed, "nonce", nonce)
    }
    return
}
//不断变更Nonce
nonce++
}
}
}
//代码位置 consensus/ethash/sealer.go

```

以太坊挖矿难度计算

以太坊每次挖矿均需计算当前区块难度。按版本不同有三种计算难度的规则，分别为：calcDifficultyByzantium (Byzantium 版)、calcDifficultyHomestead (Homestead 版)、calcDifficultyFrontier (Frontier 版)。此处以 calcDifficultyHomestead 为例。

计算难度时输入有：parent_timestamp：父区块时间戳 parent_diff：父区块难度 block_timestamp：当前区块时间戳 block_number：当前区块的序号

当前区块难度计算公式，即：

```

block_diff = parent_diff
+ (parent_diff / 2048 * max(1 - (block_timestamp - parent_timestamp) // 10, -99)
+ 2^((block_number // 100000) - 2)

```

其中//为整数除法运算符，a/b，即先计算 a/b，然后取不大于 a/b 的最大整数。

调整难度的目的，即为使挖矿时间保持在 10-19s 期间内，如果低于 10s 增大挖矿难度，如果大于 19s 将减小难度。另外，计算出的当前区块难度不应低于以太坊创世区块难度，即 131072。

go-ethereum-1.7.3 源码中计算挖矿难度代码如下：

```

func calcDifficultyHomestead(time uint64, parent *types.Header) *big.Int {
    // https://github.com/ethereum/EIPs/blob/master/EIPS/eip-2.mediawiki
    // algorithm:
    // diff = (parent_diff +
    //         (parent_diff / 2048 * max(1 - (block_timestamp - parent_timestamp) // ↵
↵10, -99))
    //         ) + 2^(periodCount - 2)
    bigTime := new(big.Int).SetUint64(time)
    bigParentTime := new(big.Int).Set(parent.Time)
    // holds intermediate values to make the algo easier to read & audit
    x := new(big.Int)
    y := new(big.Int)
    // 1 - (block_timestamp - parent_timestamp) // 10
    x.Sub(bigTime, bigParentTime)
    x.Div(x, big10)
    x.Sub(big1, x)
    // max(1 - (block_timestamp - parent_timestamp) // 10, -99)
    if x.Cmp(bigMinus99) < 0 {
        x.Set(bigMinus99)
    }
    // (parent_diff + parent_diff // 2048 * max(1 - (block_timestamp - parent_
↵timestamp) // 10, -99))

```

(下页继续)

(续上页)

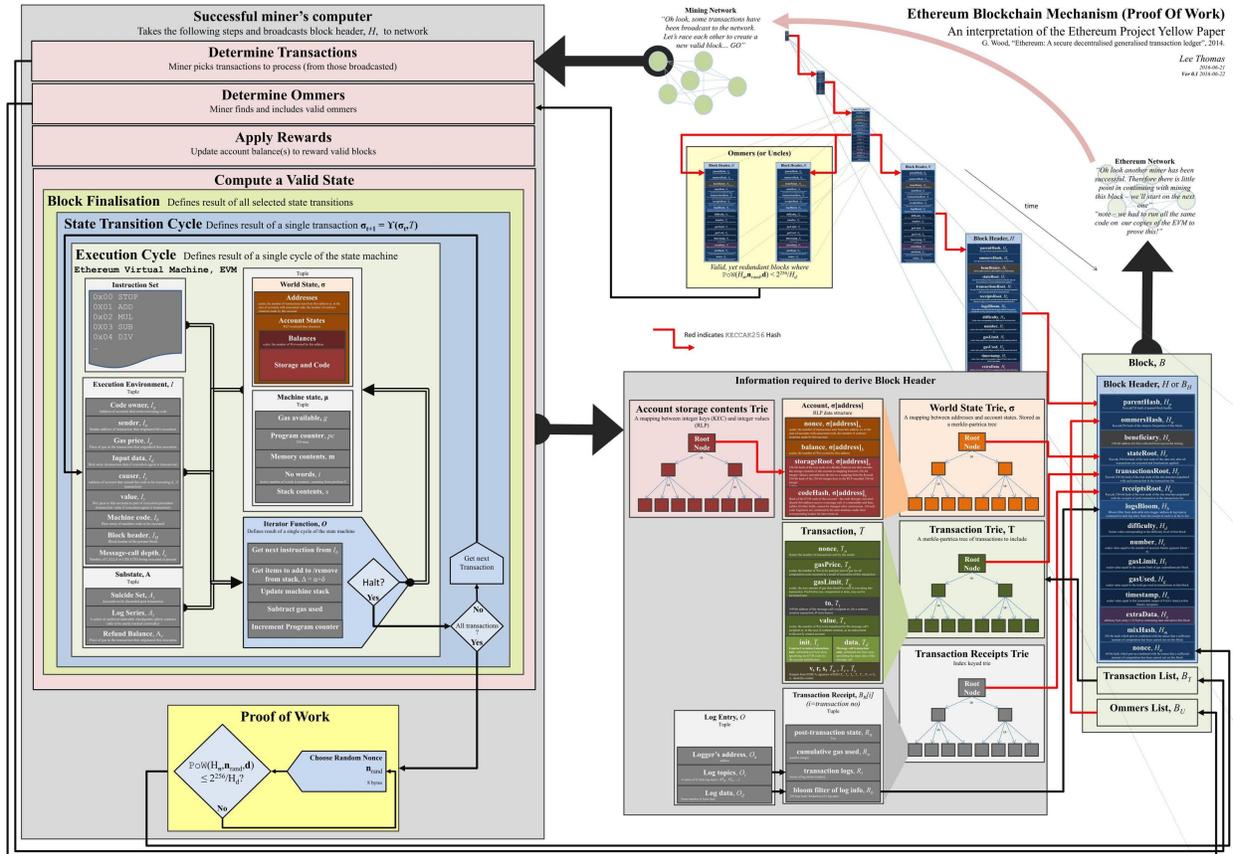
```
y.Div(parent.Difficulty, params.DifficultyBoundDivisor)
x.Mul(y, x)
x.Add(parent.Difficulty, x)
// minimum difficulty can ever be (before exponential factor)
if x.Cmp(params.MinimumDifficulty) < 0 {
    x.Set(params.MinimumDifficulty)
}
// for the exponential factor
periodCount := new(big.Int).Add(parent.Number, big1)
periodCount.Div(periodCount, expDiffPeriod)
// the exponential factor, commonly referred to as "the bomb"
// diff = diff + 2^(periodCount - 2)
if periodCount.Cmp(big1) > 0 {
    y.Sub(periodCount, big2)
    y.Exp(big2, y, nil)
    x.Add(x, y)
}
return x
}
//代码位置 consensus/ethash/consensus.go
```

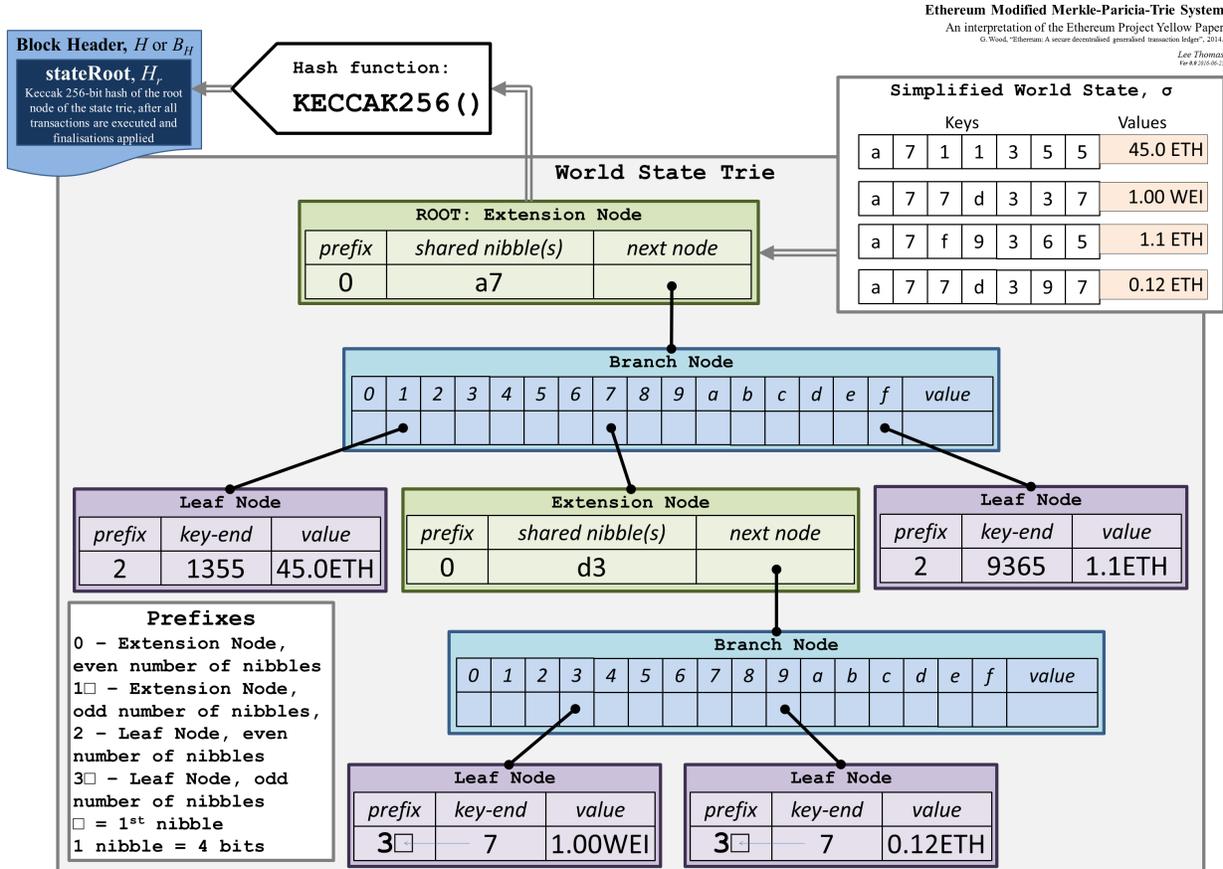
后记

Pow 算法概念简单，即工作端提交难以计算但易于验证的计算结果，其他节点通过验证这个结果来确信工作端完成了相当的工作量。但其缺陷也很明显：1、随着节点将 CPU 挖矿升级为 GPU、甚至矿机挖矿，节点数和算力已渐渐失衡；2、比特币等网络每秒需完成数百万亿次哈希计算，资源大量浪费。为此，业内提出了 Pow 的替代者如 PoS 权益证明算法，即要求用户拥有一定数量的货币，才有权参与确定下一个合法区块。另外，相对拥有 51% 算力，购买超过半数以上的货币难度更大，也使得恶意攻击更加困难。

2.2 ETH

2.2.1 ETH Overview





2.2.2 ETH SPV

什么是 SPV

轻钱包 (SPV) 是“Simplified Payment Verification” (简单支付验证) 的缩写。中本聪论文简要地提及了这一概念，指出：不运行完全节点也可验证支付，用户只需要保存所有的 block header 就可以了。用户虽然不能自己验证交易，但如果能够从区块链的某处找到相符的交易，他就可以知道网络已经认可了这笔交易，而且得到了网络的多少个确认。

ETH Hash 算法

- DAG 分层级的有向无环图，每 3 万个块生成一次，
 - 每个 epoch 的 seed 是存在依赖的，

```
Seed Hash[0] = Kec(32 byte 0)
Seed Hash[i] = Kec(32 byte i-1)
```

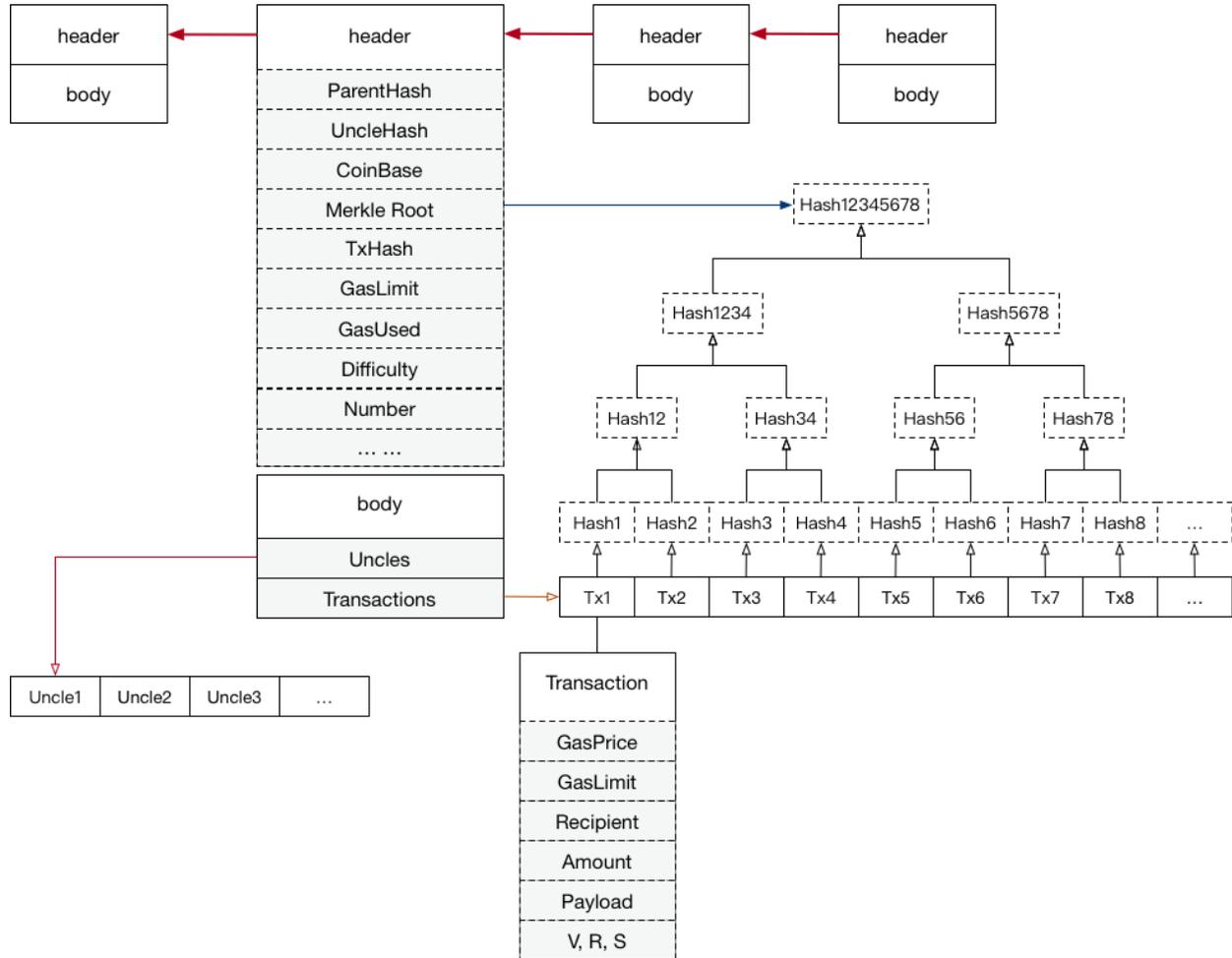
- Seed 生成 cache，初始值是 16MB，每个 epoch 增长 128K B，每个 cache 的单位是 64B
- cache 生成 Data，初始值是 1GB，每个 epoch 增长 8M。

ETH 数据结构

Block

以太坊中的一个区块由区块头 Header，以及交易列表，以及叔块的 header 集合三部分组成。

具体结构如下图



Block Header

```
// Header represents a block header in the Ethereum blockchain.
type Header struct {
    ParentHash    common.Hash    `json:"parentHash"      gencodec:"required"`
    UncleHash    common.Hash    `json:"sha3Uncles"      gencodec:"required"`
    Coinbase     common.Address `json:"miner"           gencodec:"required"`
    Root         common.Hash    `json:"stateRoot"       gencodec:"required"`
    TxHash       common.Hash    `json:"transactionsRoot" gencodec:"required"`
    ReceiptHash  common.Hash    `json:"receiptsRoot"    gencodec:"required"`
    Bloom        Bloom          `json:"logsBloom"       gencodec:"required"`
    Difficulty   *big.Int      `json:"difficulty"      gencodec:"required"`
}
```

(下页继续)

(续上页)

Number	*big.Int	`json:"number"`	gencodec:"required"`
GasLimit	uint64	`json:"gasLimit"`	gencodec:"required"`
GasUsed	uint64	`json:"gasUsed"`	gencodec:"required"`
Time	uint64	`json:"timestamp"`	gencodec:"required"`
Extra	[] byte	`json:"extraData"`	gencodec:"required"`
MixDigest	common.Hash	`json:"mixHash"`	
Nonce	BlockNonce	`json:"nonce"`	

Header 包括以下字段。

- ParentHash 用 Hp 表示。
 - 父节点的 hash 值, 每个区块有且只有一个父区块, 但允许某个时候存在多个子块 (分叉)。
- UncleHash 用 Ho 表示。
 - Block 结构体的成员 uncles 进行 RLP 编码后的 Hash 值, 这块是跟 GHOST 相关的,
- Coinbase 用 Hc 表示。矿工 address
- stateRoot 用 Hr 表示。
 - StateDB 中的 state Trie 的根节点的 RLP 哈希值。Block 中, 每个账户以 stateObject 对象表示, 账户以 Address 为唯一标示, 其信息在相关交易 Transaction 的执行中被修改。所有账户对象可以逐个插入一个 Merkle-PatriciaTrie (MPT) 结构里, 形成 state Trie。
- transactionsRoot 用 Ht 表示。
 - Block 中 tx Trie 的根节点的 RLP 哈希值。Block 的成员变量 transactions 中所有的 tx 对象, 被逐个插入一个 MPT 结构, 形成 tx Trie。
- receiptsRoot 用 He 表示。
- Block 中的 Receipt Trie 的根节点的 RLP 哈希值。Block 的所有 Transaction 执行完后会生成一个 Receipt 数组, 这个数组中的所有 Receipt 被逐个插入一个 MPT 结构中, 形成 Receipt Trie。
- logsBloom 用 Hb 表示。
 - Bloom 过滤器 (Filter), 用来快速判断一个参数 Log 对象是否存在于一组已知的 Log 集合中
- difficulty 用 Hd 表示。
 - 区块的难度。Block 的 Difficulty 由共识算法基于 parentBlock 的 Time 和 Difficulty 计算得出
- number 用 Hi 表示。区块号
- gasLimit 用 Hl 表示。区块的 gas 数量限制, 即区块中交易使用掉的 gas 值不应该超过该值。
- gasUsed 用 Hg 表示。区块使用掉的 gas 数量,
- timestamp 用 Hs 表示。时间戳,
- extraData 用 Hx 表示。额外的数据, 合法的交易对长度有限制,
- mixHash 用 Hm 表示。与 nonce 一起用作工作量证明,
- nonce 用 Hn 表示。PoW 中使用的随机数与 mixHash 一起用作工作量证明,

区块头校验

- 校验区块号和区块 hash 的符号化表示
- 难度计算的符号化表示
- 当区块编号为 0 的时候其难度值是固定好的，在这里用 D_0 表示，其值为 131072.
 - 对于其他区块，其难度值需要根据其父区块难度值以及一些其他因素，出块的间隔时间，区块编号等有关进行调节的，若小于 D_0 ，则难度值调整为 D_0 。
- 调节系数 the adjustment factor x 的定义。
 - 难度系数 (difficulty parameter) ς_2 的定义。该系数主要与出块间隔时间有关，当间隔大的时候，系数变大，难度也会相应变大，当间隔小的时候，系数变小，难度也会变小。使得区块链在整体上出块时间是趋于稳定的。其中 y 值根据父节点的 uncle 节点是否为空而有所区别，可以看出当父节点的 uncle 不为空的时候， y 值为 2，说明当前的分叉程度较大，适当调大难度，一定程度上会减少分叉。

- gasLimit 限制的符号化表示

$$\begin{aligned}
 H_1 &< P(H)_{H_1} + \left\lfloor \frac{P(H)_{H_1}}{1024} \right\rfloor \quad \wedge \\
 H_1 &> P(H)_{H_1} - \left\lfloor \frac{P(H)_{H_1}}{1024} \right\rfloor \quad \wedge \\
 H_1 &\geq 5000
 \end{aligned}$$

- 区块 gasLimit 必须大于等于 5000，和上个区块的 gasLimit 差值不超过上个块 gasLimit 除以 1024
- 时间戳的符号化表示
 - 当前区块的时间戳必须大于父区块的时间戳
- mixHash 和 nonce 相关符号化表示
- 总体而言：区块头验证的符号化表示如下

$$\begin{aligned}
V(H) \equiv n &\leq \frac{2^{256}}{H_d} \wedge m = H_m \quad \wedge \\
H_d &= D(H) \quad \wedge \\
H_g &\leq H_1 \quad \wedge \\
H_1 &< P(H)_{H_1} + \left\lfloor \frac{P(H)_{H_1}}{1024} \right\rfloor \quad \wedge \\
H_1 &> P(H)_{H_1} - \left\lfloor \frac{P(H)_{H_1}}{1024} \right\rfloor \quad \wedge \\
H_1 &\geq 5000 \quad \wedge \\
H_s &> P(H)_{H_s} \quad \wedge \\
H_i &= P(H)_{H_i} + 1 \quad \wedge \\
\|H_x\| &\leq 32
\end{aligned}$$

为什么用 SPV

以太坊区块的存储的信息主要包含两部分内容：Header 和 Body. Body 主要由交易列表和 Uncle Block 构成；Header 的内容更为丰富，主要有父区块的 hash，时间戳，挖矿的难度等。而 Header 中重要的是以太坊的 "三棵树"，对应于区块头中的 StateRoot，TransactionRoot 和 ReceiptRoot 三个哈希值

light client 只会同步 ETH 的 Block header data。由于没有具体的交易信息，只有 Block header data 无法验证某一个交易是否被确认，所以只能通过 Merkle 证明来判断一笔交易是否在现在的区块链交易列表中。

用户在验证一笔交易的时候就不需要知道其他交易的数据，只需要其他交易的交易哈希 (Txhash)，然后计算这笔交易的 hash 值，再逐步计算出 Merkle Root 并与接收到的 Block Header 的 Merkle Root 对比就可以验证交易是否存在了（需要用户提供交易的信息：区块块高，交易的 index）

light client 无法识别最长链（规范链）

如何维护一个 light client

- 计算 receiptsRoot 提供 tx 的校验
- 计算 header data 的校验

Reference

轻节点 SPV

go-ethereum 源码分析 - 区块结构 (1)

以太坊黄皮书详解（一）

2.2.3 Solidity 系列之基础

小知识

- 合约地址是通过调用地址 + 当前 nonce 计算出来的，所以合约地址是可以被提前预估的

```
func CreateAddress(b common.Address, nonce uint64) common.Address {
    data, _ := rlp.EncodeToBytes([]interface{}{b, nonce})
    return common.BytesToAddress(Keccak256(data)[12:])
}
```

不常见语法

习惯

- 变量命名以下划线开始，以区别全局变量。
- 私有函数的名字用下划线起始

函数修饰符

可见修饰符

由于 Solidity 有两种函数调用（内部调用不会产生实际的 EVM 调用或称为“消息调用”，而外部调用则会产生一个 EVM 调用），函数和状态变量有四种可见性类型。函数可以指定为 external, public, internal 或者 private。对于状态变量，不能设置为 external，默认是 internal。

- external 外部函数作为合约接口的一部分，意味着我们可以从其他合约和交易中调用。一个外部函数 f 不能从内部调用（即 f 不起作用，但 this.f() 可以）。当收到大量数据的时候，外部函数有时候会更有效率，因为数据不会从 calldata 复制到内存。
- public 函数是合约接口的一部分，可以在内部或通过消息调用。对于 public 状态变量，会自动生成一个 getter 函数。
- internal 这些函数和状态变量只能是内部访问（即从当前合约内部或从它派生的合约访问）。
- private 函数和状态变量仅在当前定义它们的合约中使用，并且不能被派生合约使用。

状态修饰符

- modifier 修饰符跟函数很类似，不过是用来修饰其他已有函数用的，在其他语句执行前，为它检查下先验条件。关键字 modifier 告诉编译器，这是个 modifier(修饰符)，而不是个 function(函数)。它不能像函数那样被直接调用，只能被添加到函数定义的末尾，用以改变函数的行为。

```
modifier onlyOwner() {
    require(msg.sender == owner);
    _; //修饰符的最后一行为 _;，表示修饰符调用结束后返回，并执行调用函数余下的部分
}

contract MyContract is Ownable {
    event LaughManiacally(string laughter);

    function likeABoss() external onlyOwner {
```

(下页继续)

(续上页)

```

    LaughManiacally("Muahahahaha");
  }
}

```

注意 likeABoss 函数上的 onlyOwner 修饰符。当你调用 likeABoss 时，首先执行 onlyOwner 中的代码，执行到 onlyOwner 中的 _; 语句时，程序再返回并执行 likeABoss 中的代码。

可见，尽管函数修饰符也可以应用到各种场合，但最常见的还是放在函数执行之前添加快速的 require 检查。

- view 意味着它只能读取数据不能更改数据

```

string greeting = "What's up dog";
function sayHello() public returns (string) {
    return greeting;
}

```

- pure 这个函数甚至都不访问应用里的数据。返回值完全取决于它的输入参数。

```

function _multiply(uint a, uint b) private pure returns (uint) {
    return a * b;
}

```

- payable 修饰符

计算

- 加法: $x + y$
- 减法: $x - y$,
- 乘法: $x * y$
- 除法: x / y
- 取模 / 求余: $x \% y$ (例如, $13 \% 5$ 余 3, 因为 13 除以 5, 余 3)
- 乘方操作 (如: x 的 y 次方) // 例如: `uint x = 5 ** 2; // equal to $5^2 = 25$`

预定义的全局变量和函数

当合约在 EVM 中执行时，它可以访问一组有限的全局对象，包括 block、msg 和 tx 对象。另外，以太坊把一些 EVM 字节码通过预定义函数的方式对外提供。

- 区块和交易对象属性

Block 的数据要有安全性考虑，不能盲目使用，因为其数据取决于矿工如何打包交易，出块

- `blockhash(uint blockNumber) returns (bytes32)`: 指定区块的哈希值，仅限于当前区块之前不超过 256 个区块，否则返回 0
- `block.chainid(uint)`: current chain id
- `block.coinbase(address payable)`: 当前区块的矿工地址
- `block.difficulty(uint)`: 当前区块工作量证明的难度
- `block.gaslimit(uint)`: current block gaslimit

- `block.number (uint)`: current block number
- `block.timestamp (uint)`: current block timestamp as seconds since unix epoch, 由矿工写入的当前区块的时间戳。
- `gasleft () returns (uint256)`: remaining gas
- `msg.data (bytes calldata)`: complete calldata, 调用合约时传入的数据。
- `msg.sender (address)`: sender of the message (current call)
- `msg.sig (bytes4)`: first four bytes of the calldata (i.e. function identifier)
- `msg.value (uint)`: number of wei sent with the message
- `tx.gasprice (uint)`: gas price of the transaction
- `tx.origin (address)`: sender of the transaction (**full call chain**)

此处要有安全性考虑, 使用 `tx.origin` 是调用链路交易的发起者, 不能盲目使用

• 计算和加密方法

- `addmod (uint x, uint y, uint k) returns (uint)` 计算 $(x + y) \% k$
- `mulmod (uint x, uint y, uint k) returns (uint)` 计算 $(x * y) \% k$
- `keccak256 (bytes memory) returns (bytes32)` 对数据做 Keccak-256 hash
- `sha256 (bytes memory) returns (bytes32)` 对数据做 sha256 hash
- `ripemd160 (bytes memory) returns (bytes20)` 对数据做 RIPEMD-160 hash
- `ecrecover (bytes32 hash, uint8 v, bytes32 r, bytes32 s) returns (address)`

`ecrecover` 通过签名和数据 hash 可以复原签名的地址, 如果还原出错, 返回全 0 地址

• 地址对象

- `<address>.balance (uint256)` balance of the Address in Wei
- `<address>.code (bytes memory)` code at the Address (can be empty), 指的是合约地址
- `<address>.codehash (bytes32)` the codehash of the Address 1
- `<address payable>.transfer (uint256 amount)`
send given amount of Wei to Address, reverts on failure, forwards 2300 gas stipend, not adjustable
- `<address>.call (bytes memory) returns (bool, bytes memory)`

一种底层 CALL 函数, 可以构建一个包含自定义数据的调用, 出错时会返回 false。警告: 这不是一种安全的调用, 调用接收方可以无意或有意地耗尽你的 gas, 导致合约因为 OOG 异常而停止, 总是要检查 call 的返回值。

此处要有安全性考虑, 使用 `transfer` 是其默认 2300 gas, 执行 fallback 函数有限, 但 `call` 会默认调用全部 gas

- `<address>.delegatecall (bytes memory) returns (bool, bytes memory)`
底层 DELEGATECALL with the given payload, returns success condition and return data, forwards all available gas, adjustable
- `<address>.staticcall (bytes memory) returns (bool, bytes memory)`
issue low-level STATICCALL with the given payload, returns success condition and return data, forwards all available gas, adjustable

- 合约相关

- **this** (current contract' s type) 当前合约, 可以转为 **Address** 类型
- **selfdestruct** (**address payable recipient**)

销毁当前合约, 将其资金发送到给定地址并结束执行。请注意, **selfdestruct** 有一些继承自 EVM 的特性:

- * **recipient** 如果是个合约, 接收合约的接收函数不会执行。

此处要有安全性考虑, **recipient** 合约必须接受 **eth**, 合约余额强制增加, 并且不执行 **fallback** 函数

- * 合约仅在交易结束时才真正被销毁, 而 **revert** 可能会“撤消”销毁。

selfdestruct 是唯一一个能在单个区块中变更无限个状态对象的操作码

其他所有的操作码都只能操作账户中的单个值或者存储树上的单个 **key**, 所以它们能变更多少固定大小的对象是有限制的 (通常, 调用一个操作码只能变更一个对象)。但是, **SELFDESTRUCT** 可以删除整棵存储树。

2.2.4 Solidity 系列之 EVM

以太坊 EVM 的特点

- EVM 是一种基于栈的虚拟机 (区别于基于寄存器的虚拟机), 用于编译、执行智能合约
- EVM 是图灵完备的 (图灵完备是指: 具有无限存储能力的通用物理机器或编程语言, 简单来说就是可以解决一切可计算的问题)
- EVM 是一个完全隔离的环境, 在运行期间不能访问网络、文件, 即使不同合约之间也有有限的访问权限
- 操作数栈调用深度为 `1024 params.CallCreateDepth 1024`
- 机器码长度一个字节, 最多可以有 256 个操作码

State Trie 全局状态树

Transactions Trie 区块级交易树

Receipts Trie 区块级数据树

Storage Trie 全局存储数

四个 Trie 都不在网络传输, 网络中只传输交易数据

三种情况

执行步骤

EVM 代码的执行环境

EVM 细节

- EVM 代码执行的实际 gas 消耗与其对内存 memory 的使用有关，并不是固定的。
- 鼓励最小化使用存储 storage，用 sstore 操作将非 0 值存储区域重置为 0 值，会获得实时的 gas 返还。
- 交易执行的最后会删除执行过程中接触过的所有“空账户”和自毁列表中的账户，这也会返还一定量的 gas。
- EVM 代码的执行必定会持续到一个正常终止或一个异常终止，但无法用代码直接触发一个异常终止。
- EVM 代码执行的异常终止会撤销当前交易中所有对状态的更改，但执行过程中所有消耗的 gas 不会返还。

消息调用

设 recipient 为消息调用的目标地址，sender 为消息调用的发起者地址，code address 为消息调用实际执行的代码所属的地址，我们有 4 种发起消息调用的方法

GHOST 算法

Greedy Heaviest-Observed Sub-Tree

分支选举协议

比特币是最长量分支 10 分钟，

以太坊是 16 秒所以不能最长量

执行模型-以太坊虚拟机

存储设计

费用设计

1. 区分临时存储 (Memory，存在于 VM 的每个实例中，并在 VM 执行结束后消失) 和永久存储 (Storage，存在于区块链状态层)；

Reference

- 以太坊虚拟机 EVM 执行原理
- 以太坊黄皮书精简解读

2.3 BTC

2.3.1 BTC 系列之 docker 节点

Docker 启动

testnet

目前 BTC 测试网数据，大约占据了 32G 硬盘存储空间，耗费大约 80 分钟

```
docker run --name btc-test-node -p 18332:18332 -d fpchen/btc-testnet:m1
```

regtest

```
docker run --name btc-dev-node -p 18443:18443 -d fpchen/btc-regtest:m1
```

程序启动

- 1、配置文件, 编辑 /etc/bitcoin/bitcoin.conf

```
daemon=1
server=1
rpcuser=test
rpcpassword=test
regtest=1
txindex=1
rpccallowip=0.0.0.0/0
discover=0
listen=0
```

- 2、启动本地私链

```
❏ bitcoind -conf=/etc/bitcoin/bitcoin.conf
Bitcoin Core starting
```

- 3、获取 bitcoind 监听的端口, 本地一般是 18843

```
❏ netstat --ip -lpa | grep bitcoind
tcp    0  0  localhost:18443      0.0.0.0:*          LISTEN      6768/bitcoind
```

- 4、发送请求, 测试服务是否运行

```

❏ curl --request POST \
  --user test:test \
  --data-binary '{"jsonrpc": "1.0", "id": "curltest", "method":
  ↪ "getblockchaininfo", "params": [] }' \
  -H 'content-type: applicaiton/json;' \
  http://127.0.0.1:18443/
{
  "result": {
    "chain": "regtest",
    "blocks": 303,
    "headers": 303,
    "bestblockhash":
  ↪ "5eba253aed009e25d347d59906c39bf47ae13ca7dc6ccd58c41edd549dd06214",
    "difficulty": 4.656542373906925e-10,
    "mediantime": 1600689629,
    "verificationprogress": 1,
    "initialblockdownload": false,
    "chainwork":
  ↪ "00000000000000000000000000000000000000000000000000000000000000260",
    "size_on_disk": 91959,
    "pruned": false,
    "softforks": {
      .....
    },
    "warnings": ""
  },
  "error": null,
  "id": "curltest"
}

```

- 5、导入私钥进钱包, 默认钱包名称为””

```

❏ bitcoin-cli -conf=/etc/bitcoin/bitcoin.conf -rpcwallet="" importprivkey
  ↪ "cUDFdzioB3SqbN9vutRTUrpw5EH9srrg6RPibacPo1fGHpfPKqL"

```

- 6、构造生成块, 需要生成 101 个块, 如果需要使用 coinbase 交易, 需要 100 个块后确认, 最好多于 100 个块

```

❏ bitcoin-cli -conf=/etc/bitcoin/bitcoin.conf generatetoaddress 101
  ↪ mvuTXVzk7n9QYHMhyUUfaPdeQ4QVwA2fmT

```

- 7、根据 hash 2d9ce904ac4d0df8dfec9a9a7e2c65883c357c37bceab044fb70f3515f9720c 查询块内交易

```

❏ bitcoin-cli -conf=/etc/bitcoin/bitcoin.conf getblock
  ↪ 2d9ce904ac4d0df8dfec9a9a7e2c65883c357c37bceab044fb70f3515f9720c
{
  "hash": "2d9ce904ac4d0df8dfec9a9a7e2c65883c357c37bceab044fb70f3515f9720c",
  "confirmations": 196,
  "strippedsize": 217,
  "size": 253,
  "weight": 904,
  "height": 108,
  "version": 536870912,
  "versionHex": "20000000",
  "merkleroot": "d0bce39e8a8a52ec58e71cb0acc28f21a2d54917e9c201619454d65078766ce4
  ↪ ",

```

(下页继续)

(续上页)

```

"tx": [
  "d0bce39e8a8a52ec58e71cb0acc28f21a2d54917e9c201619454d65078766ce4"
],
"time": 1600689597,
"mediantime": 1600689596,
"nonce": 0,
"bits": "207fffff",
"difficulty": 4.656542373906925e-10,
"chainwork": "00000000000000000000000000000000000000000000000000000000000000da",
"nTx": 1,
"previousblockhash":
↪ "642bf8c9135ef0fc37799a7348ecbc5d4e70a89cfce7e59586ad6fe34b64e093",
"nextblockhash":
↪ "1f767a69b1c168ba55ade509c2b5cca11d848672e331a91dc5035778b243c0a7"
}

```

- 8、根据交易 hash d0bce39e8a8a52ec58e71cb0acc28f21a2d54917e9c201619454d65078766ce4 查询交易

```

❏ bitcoin-cli -conf=/etc/bitcoin/bitcoin.conf gettxout_
↪ d0bce39e8a8a52ec58e71cb0acc28f21a2d54917e9c201619454d65078766ce4 0
{
  "bestblock": "5eba253aed009e25d347d59906c39bf47ae13ca7dc6ccd58c41edd549dd06214",
  "confirmations": 196, #coinbase 交易需要100个确认才能使用
  "value": 50.00000000,
  "scriptPubKey": {
    "asm": "OP_DUP OP_HASH160 a8cb707e4d0a5c6e690189bc0065a8f787aabced OP_
↪ EQUALVERIFY OP_CHECKSIG",
    "hex": "76a914a8cb707e4d0a5c6e690189bc0065a8f787aabced88ac",
    "reqSigs": 1,
    "type": "pubkeyhash",
    "addresses": [
      "mvuTXVzk7n9QYHMhyUUfaPdeQ4QVwA2fmT"
    ]
  },
  "coinbase": true
}

```

Docker 脚本

主要以测试网为主，私链无所谓配置

- bitcoin.conf

```

daemon=1
server=1
rpcuser=axontest
rpcpassword=axontest
txindex=1
testnet=1
fallbackfee=0.02
datadir=/data/bitcoin/bitcoin-data

[test]
rpccallwip=0.0.0.0/0

```

(下页继续)

(续上页)

```

rpcbind=0.0.0.0
rpcport=18332
addnode=47.100.162.210
addnode=114.214.228.185
# https://bitnodes.io/#join-the-network 查找可用节点同步数据

```

- Dockerfile

```

FROM ubuntu:18.04
ADD https://bitcoin.org/bin/bitcoin-core-0.20.0/bitcoin-0.20.0-x86_64-linux-gnu.
  tar.gz .
RUN tar -xzvf bitcoin-0.20.0-x86_64-linux-gnu.tar.gz -C ./
ADD entrypoint.sh ./
ADD bitcoin.conf /etc/bitcoin/bitcoin.conf
user root
RUN chmod 777 ./entrypoint.sh
RUN chmod go-w /etc/bitcoin/bitcoin.conf
RUN mkdir -p /data/bitcoin/bitcoin-data
RUN chown -R 755 /data/bitcoin/bitcoin-data

EXPOSE 18332
ENTRYPOINT ["/entrypoint.sh"]

```

- 启动脚本

```

#!/bin/bash

export PATH=/bitcoin-0.20.0/bin:$PATH

bitcoind --conf=/etc/bitcoin/bitcoin.conf
echo "wait 90 sec for bitcoind start" # 需要等一会儿确保服务启动
sleep 90

bitcoin-cli -conf=/etc/bitcoin/bitcoin.conf createwallet miner
bitcoin-cli -conf=/etc/bitcoin/bitcoin.conf -rpcwallet="miner" importprivkey
  "cURtxPqTGqaA5oLit5sMszceoEAbiLFsTRz7AHo23piqamtzbzav"

bitcoin-cli -conf=/etc/bitcoin/bitcoin.conf getnetworkinfo # network
bitcoin-cli -conf=/etc/bitcoin/bitcoin.conf getpeerinfo # peer info
bitcoin-cli -conf=/etc/bitcoin/bitcoin.conf getblockchaininfo #
  verificationprogress 代表区块同步进度

tail -f /data/bitcoin/bitcoin-data/testnet3/debug.log # 同步的日志所在位置

```

2.3.2 BTC 系列之账户脚本

UXTO

比特币完整节点跟踪所有可找到的和可使用的输出，称为“未花费的交易输出”（unspent transaction outputs），即 UTXO。所有 UTXO 的集合被称为 UTXO 集，目前有数百万个 UTXO。当新的 UTXO 被创建，UTXO 集就会变大，当 UTXO 被消耗时，UTXO 集会随着缩小。每一个交易都代表 UTXO 集的变化（状态转换）

```

{
  "version": 1,

```

(下页继续)

(续上页)

```

"locktime": 0,
"vin": [
  {
    "txid": "7957a35fe64f80d234d76d83a2a8f1a0d8149a41d81de548f0a65a8a999f6f18",
    "vout": 0,
    "scriptSig":
↪ "3045022100884d142d86652a3f47ba4746ec719bbfbd040a570b1deccbb6498c75c4ae24cb02204b9f039ff08df09cbe9
↪ 0484ecc0d46f1918b30928fa0e4ed99f16a0fb4fde0735e7ade8416ab9fe423cc5412336376789d172787ec3457eee41c0
↪ ",
    "sequence": 4294967295
  }
],
"vout": [
  {
    "value": 0.01500000,
    "scriptPubKey": "OP_DUP OP_HASH160 ab68025513c3dbd2f7b92a94e0581f5d50f654e7 OP_
↪ EQUALVERIFY OP_CHECKSIG"
  },
  {
    "value": 0.08450000,
    "scriptPubKey": "OP_DUP OP_HASH160 7f9b1a7fb68d60c536c2fd8aeaa53a8f3cc025a8 OP_
↪ EQUALVERIFY OP_CHECKSIG",
  }
]
}

```

交易输出 (创世交易先有输出)

```

"vout": [
  {
    "value": 0.01500000,
    "scriptPubKey": "OP_DUP OP_HASH160 ab68025513c3dbd2f7b92a94e0581f5d50f654e7 OP_
↪ EQUALVERIFY
OP_CHECKSIG"
  },
  {
    "value": 0.08450000,
    "scriptPubKey": "OP_DUP OP_HASH160 7f9b1a7fb68d60c536c2fd8aeaa53a8f3cc025a8 OP_
↪ EQUALVERIFY OP_CHECKSIG",
  }
]

```

交易输入

交易输入将 UTXO (通过引用) 标记为将被消费, 并通过解锁脚本提供所有权证明。

输入的第一部分是一个指向 UTXO 的指针, 通过指向 UTXO 被记录在区块链中所在的交易的哈希值和序列号来实现

```

"vin": [
  {
    "txid": "7957a35fe64f80d234d76d83a2a8f1a0d8149a41d81de548f0a65a8a999f6f18",
    "vout": 0,

```

(下页继续)

(续上页)

```

"scriptSig" :
↪ "3045022100884d142d86652a3f47ba4746ec719bbfbd040a570b1deccbb6498c75c4ae24cb02204b9f039ff08df09cbe9
↪ 0484ecc0d46f1918b30928fa0e4ed99f16a0fb4fde0735e7ade8416ab9fe423cc5412336376789d172787ec3457eee41c0
↪ ",
  "sequence": 4294967295
}
]

```

- 一个交易 ID，引用包含正在使用的 UTXO 的交易
- 一个输出索引 (vout)，用于标识来自该交易的哪个 UTXO 被引用 (第一个为零)
- 一个 scriptSig (解锁脚本)，满足放置在 UTXO 上的条件，解锁它用于支出
- 一个序列号 (稍后讨论)

在 Alice 的交易中，输入指向的交易 ID 是：

```
7957a35fe64f80d234d76d83a2a8f1a0d8149a41d81de548f0a65a8a999f6f18
```

输出索引是 0 (即由该交易创建的第一个 UTXO)，那个输出 (如上: value) 就可以被引用了。

如果不检索输入中引用的 UTXO，则不知道它们的值。因此，在单个交易中计算交易费用的简单操作，实际上涉及多个交易的多个步骤和数据。

比特币交易脚本和脚本语言

脚本构建 (锁定与解锁)

- **锁定脚本**是一个放置在输出上面的花费条件：它指定了今后花费这笔输出必须要满足的条件。由于锁定脚本往往含有一个公钥或比特币地址 (公钥哈希值)，在历史上它曾被称为脚本公钥 (scriptPubKey)。
- **解锁脚本**是一个“解决”或满足被锁定脚本在一个输出上设定的花费条件的脚本，它将允许输出被消费。解锁脚本是每一笔比特币交易输入的一部分，而且往往含有一个由用户的比特币钱包 (通过用户的私钥) 生成的数字签名。由于解锁脚本常常包含一个数字签名，因此它曾被称作 ScriptSig。

每一个比特币验证节点会通过同时执行锁定和解锁脚本来验证一笔交易。每个输入都包含一个解锁脚本，并引用了之前存在的 UTXO。

最常见类型的比特币交易 (P2PKH: 对公钥哈希的付款) 的解锁和锁定脚本的示例

P2PKH (Pay-to-Public-Key-Hash)

```
OP_DUP OP_HASH160 <Cafe Public Key Hash> OP_EQUALVERIFY OP_CHECKSIG
```

脚本中的 Cafe Public Key Hash 即为咖啡馆的比特币地址，但该地址不是基于 Base58Check 编码上述锁定脚本相应的解锁脚本是：

```
<Cafe Signature> <Cafe Public Key>
```

将两个脚本结合起来可以形成如下组合验证脚本：

```
<Cafe Signature> <Cafe Public Key> OP_DUP OP_HASH160
<Cafe Public Key Hash> OP_EQUALVERIFY OP_CHECKSIG
```

只有当解锁脚本得到了咖啡馆的有效签名，交易执行结果才会被通过 (结果为真)

余额计算

为了构建“总接收”数量，区块链浏览器首先解码比特币地址的 Base58Check 编码，以检索编码在地址中的 Bob 的公钥的 160 位哈希值。然后，区块链浏览器搜索交易数据库，使用包含 Bob 公钥哈希的 P2PKH 锁定脚本寻找输出。通过总结所有输出的值，浏览器可以产生接收的总值。

区块链浏览器将当前未被使用的输入保存为一个分离的数据库——UTXO 集。为了维护这个数据库，区块链浏览器必须监视比特币网络，添加新创建的 UTXO，并在已被使用的 UTXO 出现在未经确认的交易中时，实时地删除它们。这是一个复杂的过程，不但要实时地跟踪交易在网络上的传播，同时还要保持与比特币网络的共识，确保在正确的链上。有时区块链浏览器未能保持同步，导致其对 UTXO 集的跟踪扫描不完整或不正确。

地址

以公钥 K 为输入，计算其 SHA256 哈希值，并以此结果计算 RIPEMD160 哈希值，得到一个长度为 160 位 (20 字节) 的数字: $A = \text{RIPEMD160}(\text{SHA256}(K))$

公式中， K 是公钥， A 是生成的比特币地址。

脚本类型

- 不同的锁定脚本依赖的源数据生成。对应关系如下
- 地址的表现形式

实际上我们看到的地址是经过 Base58 编码或 bech32 编码的地址，地址中添加了地址类型和校验码所以不同的地址开头会呈现不同的数值。生成锁定脚本的时候要解码数据并且去掉地址类型码。

不同地址的标志位对照表如下：

锁定脚本的结构

- P2PKH:

是利用公钥的哈希，相对来说比较简单，现在新版本的钱包应用的越来越少。

OP_DUP OP_HASH160 OP_CODE_LEN ADDR OP_EQUALVERIFY OP_CHECKSIG

对应的字节表示 0x76 0xa9 0x14 20 字节地址 0x88 0xac

P2PKH 长度: 25 字节

OP_CODE_LEN: 0x14 即地址 ADDR 的长度

ADDR: 是由公钥生成的 20 字节

- P2SH:

利用了脚本的哈希，并且锁定脚本相较于 P2PKH 更加简洁。由于脚本更加复杂，可以根据需要设计不同的逻辑控制，因此应用会被广泛应用。锁定脚本中表示简单，但是赎回脚本和签名处变的复杂了。

OP_HASH160 OP_CODE_LEN ADDR OP_EQUAL

对应的字节表示: 0xa9 0x14 20 字节地址 0x87

P2SH 长度: 23

OP_CODE_LEN: 0x14 即地址 ADDR 的长度

ADDR: 是由脚本经过 HASH160 生成的 20 字节

- **P2WPKH:**

是隔离见证中公钥地址的表示，格式主要是为了与非隔离见证的锁定脚本区别。旧钱包是不支持的。

VER OPCODE_LEN ADDR

对应的字节表示 0x00 0x14 20 字节地址

P2WPKH 长度: 22

0x00: 版本号, 固定 0

OPCODE_LEN: 0x14 即地址 ADDR 的长度

ADDR: 是由公钥生成的 20 字节

- **P2WSH:**

是隔离见证中脚本地址的表示，格式主要是为了与非隔离见证的锁定脚本区别。旧钱包是不支持的。因为可以用到复杂脚本，因此可以做复杂逻辑应用。

VER OPCODE_LEN ADDR

对应的字节表示 0x00 0x20 32 字节地址

P2WSH 长度: 34

0x00: 版本号, 固定 0

OPCODE_LEN: 0x20 即地址 ADDR 的长度

- **P2SH-P2WPKH:**

P2SH-P2WPKH 是对 P2WPKH 进行 HASH160 得到 20 字节地址，此地址可在 P2SH 中应用，与普通的 P2SH 地址没有任何区别。

格式见 P2SH

- **P2SH-P2WSH**

P2SH-P2WSH 是对 P2WSH 进行 HASH160 得到 20 字节地址，此地址可在 P2SH 中应用，与普通的 P2SH 地址没有任何区别。

格式见 P2SH

比特币地址类型:

比特币地址类型分为 3 种格式: legacy p2sh-segwit bech32

- legacy 类型实际上就是取公钥或脚本的 HASH160 值得到 20 位字节地址。
- p2sh-segwit 对应生成 P2SH-P2WPKH 和 P2SH-P2WSH 中的地址
- bech32 对应生成 P2WPKH 和 P2WSH 中的地址

隔离见证中的公钥类型

隔离见证中生成地址用到的公钥必须是压缩的公钥，否则可能会无法赎回此地址的币。

隔离见证地址用长度 32 或 20 来判定地址是由脚本或公钥生成的，0（版本）开头表示是隔离见证地址。

脚本中的地址和钱包中的地址区别

我们在 UI 上看到的地址是带有地址类型码并且经过 Bash58 或 Bech32 编码后的地址的字符串，因此与实际地址长度并不符合。当我们传入此地址后，会解码并去除类型码，并根据类型码生成不同的 ScriptPubkey (P2PKH P2SH P2WSH P2WPKH P2SH-P2WSH P2SH-P2WPKH)。

reference

比特币锁定脚本及地址解析

2.3.3 BTC 系列之标准与非标准交易

操作码

标准交易

如果交易的锁定和解锁脚本与一小部分被认为是安全的模板相匹配，则可以被网络接受。这是 `*isStandard ()` 和 `isStandardTx ()` * 测试，通过它的交易称为标准交易。

- 如果所有输出（锁定脚本）仅使用标准交易表单，则 `*isStandard ()` * 函数将为 TRUE。
- 如果所有输入（解锁脚本）仅根据其花费的输出使用标准交易表单，则 `*isStandardTx ()` * 函数将提供 TRUE。

定义和检查标准交易背后的主要原因是通过广播有害交易来防止某人攻击比特币。此外，这些检查还使用户无法创建交易，而这会使将来增加新的交易功能更加困难。有七种标准的交易类型。

P2PKH (pay to public key hash)

- 锁定脚本

```
OP_DUP OP_HASH160 <Cafe Public Key Hash> OP_EQUALVERIFY OP_CHECKSIG
```

Cafe Public Key Hash 即为咖啡馆的地址公钥经过 OP_HASH160 操作后的数据

- 解锁脚本

```
<Cafe Signature> <Cafe Public Key>
```

- 解锁花费 utxo 是的计算过程如下

```
<Cafe Signature> <Cafe Public Key> OP_DUP OP_HASH160 <Cafe Public Key Hash> OP_
↪EQUALVERIFY OP_CHECKSIG
```

只有当解锁脚本得到了咖啡馆的有效签名，交易执行结果才会被通过（结果为 true. 被 push 到栈中）

P2PK (pay to public key)

- 锁定脚本

```
<PUBLIC KEY A> OP_CHECKSIG
```

PUBLIC KEY A 即为咖啡馆的地址公钥

- 解锁脚本

```
<A Signature>
```

- 解锁花费 utxo 是的计算过程如下

```
<A Signature> <PUBLIC KEY A> OP_CHECKSIG
```

只有当解锁脚本得到了咖啡馆的有效签名，交易执行结果才会被通过（结果为 true. 被 push 到栈中）

multi sign

多重签名脚本设置了一个条件，其中在脚本中记录了 N 个公钥，并且这些签名中至少 M 个必须用于解锁交易。这也称为 N 对 M 方案。

其中 N 是密钥总数， M 是验证所需签名的下限阈值。当前比特币核心实现的最大 M 为 15

- 锁定脚本

```
M <PUBLIC KEY 1> <PUBLIC KEY 2> ... <PUBLIC KEY N> N OP_CHECKMULTISIG
```

- 解锁脚本

```
OP_0 <SIGNATURE 1> <SIGNATURE 2> ... <SIGNATURE M>
```

- 解锁花费 utxo 是的计算过程如下。以 3-2 为例

```
OP_0 <SIGNATURE 1> <SIGNATURE 2> 2 <PUBLIC KEY 1> <PUBLIC KEY 2> <PUBLIC KEY 3> 3_
↪ OP_CHECKMULTISIG
```

请注意，前缀 OP_0 是因为在最初的实现中存在错误需要 CHECKMULTISIG：由于这个错误，需要在栈上多了一个说法。CHECKMULTISIG 只是将其视为占位符。

OP_RETURN

将数据输出交易被用来不相关的比特币付款存储数据，任何带有 OP_RETURN 的输出都证明是不可花费的，无解锁脚本

- 锁定脚本

```
OP_RETURN <DATA>
```

P2SH (pay to script hash)

在 input，也就是锁定脚本中，包含解锁脚本的 hash。

例如，我们可 hash **5-2 多重签名交易**。P2SH 等效交易不是“支付给该 5 键多重签名脚本”，而是“支付给具有此哈希的脚本”。因此，该脚本仅存储 20 字节的哈希值，而不存储五个 pubkey（使用压缩形式大约为 180 字节）

- 锁定脚本

```
OP_HASH160 <2-OF-5 多重签名脚本hash> OP_EQUAL
```

- 解锁脚本

```
<SIG1> <SIG2> <2-OF-5 多重签名脚本>
```

- 解锁花费 utxo 是的计算过程如下。提前构造好赎回脚本 (2-OF-5 多重签名脚本)

```
<赎回脚本> OP_HASH160 <赎回脚本hash> OP_EQUAL
```

P2WPKH (pay to witness public key hash)

- 锁定脚本

```
OP_0 <PUBLIC KEY A HASH>
```

- 解锁脚本 (witness 字段)

```
<SIGNATURE S> <PUBLIC KEY A>
```

P2WSH (pay to witness script hash)

非标准交易

交易通过比特币核心参考实现中的 `*isStandard ()` 和 `isStandardTx ()` * 函数进行验证。如果它们没有通过这样的测试，则将它们直接丢弃。但是，有矿工会放弃这些严格检查，这种交易虽说不符合 Bitcoin Core 强制执行的标准，由于矿工不坚持，其实也可以上链。非标准交易使用更复杂的脚本形式，常见于链上挑战或报备 bug。它们的不同在于非标准的 input 或 output 脚本。

正确验证非标准交易会增加创建未来交易的难度，原因有两个：

1. 某些脚本可能会损害网络。比如引入复杂验证脚本（需要验证 5, 6 个小时，那节点同步交易之后，出现的空窗期会导致很多孤块产生）
2. 一些脚本可能会使将来的升级变得更加困难。

Pay to Public Key Hash 0 [2011]

它对应于 P2PKH 的一种失真，不同之处在于，它不是公共密钥的哈希，而是一个 0 值。

锁定脚本的格式为 `OP_DUP OP_HASH 160 0 OP_EQUALVERIFY OP_CHECKSIG`。

这些交易是不可花费的，因为作为 P2PKH 交易，矿工为了验证它们，需要与锁定脚本中的哈希相对应的公钥，以及生成相应签名的键。但是，HASH 160 返回一个 20 字节长的哈希：因此，没有任何传递给哈希函数的键可以返回 0。一个 Bitcointalk 线程 [16] 表示此偏差主要由 MtGox 执行。产出代表 2609.36304319 BTC 的价值，当时约为 8 000 美元（如今约为 2000 万美元）。

P2PKH NOP [2011, 2014]

此交易与 P2PKH 相同，唯一的区别是锁定脚本中包含 NOP（不执行任何操作的操作）：

`OP_DUP OP_HASH 160 <HASHPUBKEY> OP_EQUALVERIFY OP_CHECKSIG OP_NOP`。

此交易可能用于测试 OP_NOP 运算符。可以使用与 P2PKH 相同的脚本将其解锁。

OnlyHash [2011-2014]

该 *OnlyHash* 交易在区块链的最大量非标准类。这些交易在锁定脚本中包含哈希，通常是文件的哈希，用于使用区块链作为分类账来注册文档。因此，区块链系统的安全性和弹性可以被数字票据服务，证券交易所证书和智能合约之类的应用程序使用。阻止脚本对应于 `<HASH-OF-SOMETHING>`。这些交易是在引入 OP_RETURN 之前使用的，OP_RETURN 可以用于相同的目的。

P2Pool 错误 [2012]

这些交易是由于错误 P2Pool 中的 192012 年 2 月 2 日至 2012 年 4 月 1 日之间使用 20 个挖矿工具。添加了以下脚本，而不是普通的 P2PKH 锁定脚本：`OP_IFDUP OP_IF OP_2SWAP OP_VERIFY OP_2OVER OP_DEPTH` 该脚本没有任何意义，甚至无效，因为 OP_IF 没有被相应的 OP_ENDIF 关闭，可以将其解锁。

OP_CHECKLOCKTIMEVERIFY OP_DROP (CLTV) [2012]

在这种情况下，锁定脚本的格式为：`<DATA> OP_CHECKLOCKTIMEVERIFY OP_DROP`。所述 OP_CHECKLOCKTIMEVERIFY 操作者使交易无效的，如果在堆栈的顶部的元件比越大 *nLockTime* 一个交易的字段。在实践中，使用 OP_CHECKLOCKTIMEVERIFY 可以使资金在未来的某个特定时刻被证明无法使用，即，这是一种将资金冻结到未来某个特定日期的方法。因此，通过对照上述规则检查 `<DATA>` 元素，我们可以通过使用将 TRUE 插入堆栈的解锁脚本来验证交易。

OP_MIN OP_EQUAL [2012]

这些交易相关的脚本 `OP_MIN 3 OP_EQUAL`，其仅需要对应于方程两个数 $X \leq Y \wedge X = 3$ 进行验证。因此，为了解锁，可以将解锁脚本准备为 `3 4`。我们可以将此交易视为如何使用方程式生成比特币交易的证明。这意味着任何人都可以在没有任何私钥的情况下轻松解锁此类交易。

Pay to Hash (P2H) [2012-2015]

这些交易是简单的 P2SH 交易的简化；我们有一个与 P2SH 相同的阻止脚本，不同之处在于内部哈希不引用兑换脚本，而是一个十六进制字符串的哈希：`OP_HASH 160 <HASH 160 OFSOMETHING> OP_EQUALVERIFY`。有两种变体，其中仅散列运算符的类型发生变化：一种对应于 **HASH 256**，而另一种对应于 **SHA256**。两个阻止脚本是 `OP_HASH 256 <HASH 256 OFSOMETHING> OP_EQUAL` 和 `OP_SHA 256 <SHA 256 OFSOMETHING> OP_EQUAL`。我们可以将这些交易视为网络中的“竞争”，以在交易中找到哈希值的正确值。

P2H 不能被视为哈希时间锁定合同 (HTLC)。HTLC 本质上是一种付款方式，其中两个人同意一种财务安排，其中一方将向另一方支付一定数量的加密货币。但是，接收方只有一定的时间才能接受付款，否则，钱将退还给发送方。取而代之的是，P2H 交易与 HTLC 不同，因为它生成了可以被接收方接受且不受任何时间限制的支付。

这些输出只能通过提供一旦通过加密函数进行哈希处理（称为 *hashlock*）等于给定哈希值的数据来使用。验证步骤很简单：阻止脚本中哈希的十六进制字符串足以花费交易。

UnLocked (UL) [2015]

此交易有一个空的锁定脚本：只需将 **TRUE** 作为解锁脚本即可将其解锁。几乎所有这些交易的总金额为 0 BTC，即它们是无价值的交易。这样的交易除了可以用作交易费之外，还可以用作向矿工捐赠资金的一种方式：进行此类交易的任何矿工还可以包括汇款到他们控制的地址后，再增加。

OP_RETURN ERROR [2016-2017]

这些交易与 OP_RETURN 相同，不同之处在于脚本代码中有错误。该代码要求将比实际代码本身更多的操作码压入堆栈。例如，OP_RETURN 脚本可能要求在代码中插入下一个 40 个字节，但是如果只有 28 个字节，则执行失败。该交易显然是由于编程错误引起的。锁定脚本为：`OP_RETURN ERROR`（返回错误）。

OP_2 OP_3 错误 [2017-2018]

这些交易类似于 OP_RETURN 错误，但脚本代码中没有 OP_RETURN。由于 OP_RETURN ERROR，他们的代码要求将比代码本身实际字节高一些的字节压入堆栈。这些交易（例如 OP_RETURN ERROR）可能与实现错误有关。这是锁定脚本：`OP_2 OP_3 ERROR`（返回错误）。

支付脚本哈希

正如我们之前介绍的，P2SH 交易在其锁定脚本中包含脚本的哈希（称为“赎回脚本”）。在区块链中，有 149,410,668 个 P2SH 交易，其中花费了 140,620,401 个 (94,12%)。因此，我们决定分析用过的 P2SH 交易的解锁脚本中的赎回脚本的内容。

P2SH 中的非标准交易

在 P2SH 中，我们发现 111122 个非标准交易（占总数的 0.08%）；这个数量是我们“简单地”分析区块链（0.02%）时获得的数量的四倍。我们发现了几类新的非标准交易，它们与以前的交易有所不同。

OP_CHECKLOCKTIMEVERIFY OP_DROP (CLVT)

我们已经发现了五种不同类型的交易，它们利用了 CLVT 运算符，正如我们在第 7 节中已经看到的那样，使得交易在特定日期之前被证明是不可花费的。从本质上讲，它允许用户创建一个比特币交易，其输出仅在将来的某个时候可用。CLTV 对于正常运行的支付渠道（例如，闪电网络）是必需的。这些渠道实际上是一系列“脱链”交易，它们受益于典型的链上交易的所有安全性，并具有一些额外的好处

- 其中一个等于区块链中已经存在的一组交易：`<DATA> OP_CHECKLOCKTIMEVERIFY OP_DROP`
- 第二个是用 CLVT 运算符锁定的 P2PK，这是锁定脚本：`<DATA> OP_CHECKLOCKTIMEVERIFY OP_DROP <PUBLIC KEY A> OP_CHECKSIG` 该脚本意味着无法使用签名（与公钥相关）在脚本中的日期之前花费此交易。
- 第三个是用 CLVT 运算符锁定的 P2PKH，这是脚本：`<DATA> OP_CHECKLOCKTIMEVERIFY OP_DROP OP_DUP OP_HASH 160 <PUBLIC KEY A HASH> OP_EQUAL OP_CHECKSIG.`
- 然后我们有一个类 `<DATA> OP_CHECKLOCKTIMEVERIFY OP_DROP 1 OP_ADD 2 OP_EQUAL`，其仅需要对应于方程式的数目 $x + 1 = 2 \wedge x = 1$ ：它因此等待直到日期已过期。
- 最后一个是 P2PKH 变体，它也需要进一步的哈希才能解锁。锁定脚本为 `<DATA> OP_CHECKLOCKTIMEVERIFY OP_DROP OP_SHA 256 <SHA 256 OFSOMETHING> OP_DUP OP_HASH 160 <PUBLIC KEY A HASH> OP_EQUAL OP_CHECKSIG。”`

OP_DROP

该交易允许在区块链中存储一些数据，而不会使交易变得不可花费。实际上，所有这些交易都以 `<DATA> OP_DROP` 开头，其中 `<DATA>` 是我们放入块链中的内容，`OP_DROP` 是将数据从堆栈中删除以便与标准交易相同的运算符。

我们发现了使用 `OP_DROP` 的四种不同类型的交易。

- 第一种类型不需要任何解锁，脚本是：`<DATA> OP_DROP 1。`
- 2-2 多重签名类型 `<DATA> OP_DROP 2 <PUBLIC KEY A> <PUBLIC KEY B> 2 OP_CHECKMULTISIG`，它需要像正常的 2-2 多重签名一样将两个签名解锁。
- 带有 `OP_DROP` 的 P2PKH，由 `<DATA> OP_DROP OP_DUP OP_HASH 160 <PUBLIC KEY A HASH> OP_EQUAL OP_CHECKSIG` 标识。同样，对于 P2PKH 交易，此交易仅需要签名。
- 最后一个是带有 P2PK 的 `OP_DROP`：`<DATA> OP_DROP <PUBLIC KEY A> OP_CHECKSIG。`

最常用的类型是 2-2 多重签名，几乎有 25,000 次出现。我们可以说这些交易只是伪装的标准输出，使用 `OP_DROP` 运算符添加验证期间丢弃的数据

Reference

- [An Analysis of Non-standard Transactions](#)

2.3.4 BTC 系列之隔离见证

为什么需要隔离见证

消除交易的延展性

“Transaction Malleability”是指在 Transaction 被比特币网络确认之前，txid 可能被别人（攻击者）修改。

举例来说，A 提交了一个 Transaction（其功能是转移一个 BTC 给 B），记为 tx1，这个交易的 Input 的 scriptSig 字段中包含了 A 的签名。这个 Transaction 在打包确认前，攻击者稍微调整一下 scriptSig（比如增加无关紧要的指令，scriptSig 的关键信息 A 的签名是不能变的，因为不能影响到这个交易信息的验证），这时会得到一个新 txid（因为 scriptSig 参与了 txid 的生成），记为 tx2。现在 tx1 和 tx2 都在等待打包确认状态。如果恰好 tx2 被先打包（这很少出现，但可以伪造更多的 tx3, tx4 来增加被先打包的概率），那么 tx1 在验证时就会失败。因为交易 Input 中引用的 UTXO 已经被 tx2 花费了。这样，A 成功转移了一个 BTC 给 B，但 txid 却不是 A 当初提交的那个 tx1 了。当然如要 tx1 被先打包，那么 tx2 会验证失败，这里 txid 还是 A 当初提交的那个 tx1。

Transaction Malleability 为什么会带来问题呢？我们以例子来说明，假设 A 开了一个交易所，而 B 是该交易所的用户。B 以前存入了 100 个 BTC，现在向 A 提出提币申请，A 向 B 转移 100 个 BTC，记下 txid 为 tx1，B 利用 Transaction Malleability，伪造了另外一个交易 tx2。这时恰好 tx2 被先打包了，tx1 失败。B 对 A 说“我没有收到 100 个 BTC”，A 去查询他当时提交的 txid（即 tx1），发现确实 tx1 没有打包，然后再发起一个新交易，向 B 转移 100 BTC。这里，其实 A 已经向 B 转移了 200 BTC。

当然，这种情况是可以避免的。如果 A 不根据 txid 来判断是否转账成功，而是查询它钱包地址的余额变化，那么他是可以发现已经转移过 100 个 BTC 给用户 B 了（只是 txid 变了而已）。

比特币交易平台 Mt.Gox，曾经由于 Transaction Malleability 而损失了比特币，参见 [Bitcoin Transaction Malleability and MtGox](#)

采用隔离见证时，由于 scriptSig 信息总是为空，见证数据保存在另外的字段（witness）中，修改见证数据，也不会对 txid 造成变化，所以消除了 Transaction Malleability。

优化交易存储

见证数据通常在交易的总大小中占了很大比重。更复杂的脚本，如用于多签或支付通道的脚本体积会非常大。在某些情况下，这些脚本会占到交易数据的大多数（超过 75%）。通过将见证数据从交易中移出，隔离见证提高了比特币的可扩展性。节点可以在验证签名后删减见证数据。见证数据不再需要发送到所有节点，也不需要被所有节点存储在磁盘上。

签名验证优化

隔离见证升级签名函数（CHECKSIG，CHECKSIGG 等），以减少算法的计算复杂度。在引入隔离见证之前，用于生成签名的算法需要大量与交易大小成比例的哈希操作。相对于签名操作的数量数据哈希计算复杂度增加到 $O(n^2)$ ，给验证签名的所有节点带来了巨大的计算负担。有了隔离见证，算法的复杂度降低到 $O(n)$ 。

离线签名改进

隔离见证签名包含由签名的哈希中的每个输入引用的值（金额）。以前，离线签名设备（如硬件钱包）必须在签名交易之前验证每个输入的金额。这通常是通过流式传输大量先前作为输入引用的交易的数据来实现的。由于金额现在是已签名的提交哈希的一部分，因此离线设备不需要以前的交易。如果金额不匹配（由被入侵的在线系统篡改），签名将是无效的。

隔离见证交易案例

从某个地址往 mno5h3qJZ3cT1TtfJZVErzFkF2BJQtkJMh 转账 0.01 BTC

- p2pkh.

```

bitcoin-cli getrawtransaction
↳2b95053f3f0f59d358e1afb12bd8153245405954e8dcd6eca8b82daf02adae71 2
{
  "txid": "2b95053f3f0f59d358e1afb12bd8153245405954e8dcd6eca8b82daf02adae71",
  "hash": "2b95053f3f0f59d358e1afb12bd8153245405954e8dcd6eca8b82daf02adae71",
  "version": 2,
  "size": 226,
  "vsize": 226,
  "weight": 904,
  "locktime": 0,
  "vin": [
    {
      "txid": "04c718ae8f34024f0b0650cf2f8ebd40c63214e5fa8836206bedca26650f29c2",
      "vout": 0,
      "scriptSig": {
        "asm":
↳"3045022100b84617263695d31407d6ae230d3b629d071e65d936ee946c8adf84304163cb2d022045435bf766cfcf5
↳029ba616506337ff7f0ccaac35938af4ba614b9853d670c4a45d830d7d38de7823",
        "hex":
↳"483045022100b84617263695d31407d6ae230d3b629d071e65d936ee946c8adf84304163cb2d022045435bf766cfcf5
↳"
      },
      "sequence": 4294967295
    }
  ],
  "vout": [
    {
      "value": 0.01000000,
      "n": 0,
      "scriptPubKey": {
        "asm": "OP_DUP OP_HASH160 4fd5b5cbe9d25a7cbbef9847d1a6e97d8a076975 OP_
↳EQUALVERIFY OP_CHECKSIG",
        "hex": "76a9144fd5b5cbe9d25a7cbbef9847d1a6e97d8a07697588ac",
        "reqSigs": 1,
        "type": "pubkeyhash",
        "addresses": [
          "mno5h3qJZ3cT1TtfJZVErzFkF2BJQtkJMh"
        ]
      }
    },
    {
      "value": 0.98997750,
      "n": 1,

```

(下页继续)

(续上页)

```

    "scriptPubKey": {
      "asm": "OP_DUP OP_HASH160 db9d0eb300a248a494f55ec1abb140d5d580216f OP_
↪EQUALVERIFY OP_CHECKSIG",
      "hex": "76a914db9d0eb300a248a494f55ec1abb140d5d580216f88ac",
      "reqSigs": 1,
      "type": "pubkeyhash",
      "addresses": [
        "n1YARc5gEby2WPf81yPUzBDaCHTvRbtTus"
      ]
    }
  ],
  "hex":
↪"0200000001c2290f6526caed6b203688fae51432c640bd8e2fcf50060b4f02348fae18c7040000000006b483045022
↪",
  "blockhash": "047ec3252f08ace282415d5d25ac1ff9eb9dfb55b73612888ee595c55cee6b00",
  "confirmations": 395,
  "time": 1617772888,
  "blocktime": 1617772888
}

```

- p2wpkh

```

bitcoin-cli getrawtransaction
↪3dbdf66064a5dd8414f4afea80f2cd318991fef28fb9dbbf491b0e7f58d6226e 2
{
  "txid": "3dbdf66064a5dd8414f4afea80f2cd318991fef28fb9dbbf491b0e7f58d6226e",
  "hash": "03dcdd25e721f4090ce0123917b2db94019bd1290cfb3af128e712a993f3f1e6",
  "version": 2,
  "size": 225,
  "vsize": 144,
  "weight": 573,
  "locktime": 0,
  "vin": [
    {
      "txid": "c320be2f30ab7e35975133b30adb8b804e580f674dda78d34b199b7da915df97",
      "vout": 0,
      "scriptSig": {
        "asm": "",
        "hex": ""
      },
      "txinwitness": [
↪"304402201e9df1d606a9fcbd159b62f1fd0675ba31092ef9d8faf8afa09b48c765d40c74022016f272575b03e22c1
↪",
        "029ba616506337ff7f0ccaac35938af4ba614b9853d670c4a45d830d7d38de7823"
      ],
      "sequence": 4294967295
    }
  ],
  "vout": [
    {
      "value": 0.01000000,
      "n": 0,
      "scriptPubKey": {
        "asm": "OP_DUP OP_HASH160 4fd5b5cbe9d25a7cbbef9847d1a6e97d8a076975 OP_
↪EQUALVERIFY OP_CHECKSIG",

```

(下页继续)

(续上页)

```

    "hex": "76a9144fd5b5cbe9d25a7cbbef9847d1a6e97d8a07697588ac",
    "reqSigs": 1,
    "type": "pubkeyhash",
    "addresses": [
      "mno5h3qJZ3cT1TtfJZVErzFkF2BJQtkJMh"
    ]
  },
  {
    "value": 0.98997740,
    "n": 1,
    "scriptPubKey": {
      "asm": "0 db9d0eb300a248a494f55ec1abb140d5d580216f",
      "hex": "0014db9d0eb300a248a494f55ec1abb140d5d580216f",
      "reqSigs": 1,
      "type": "witness_v0_keyhash",
      "addresses": [
        "bcrt1qmwwsavcq5fy2f984tmq6hv2q6h2cqgt06jy455"
      ]
    }
  }
],
"hex":
↪ "020000000010197df15a97d9b194bd378da4d670f584e808bdb0ab3335197357eab302fbe20c300000000fffff
↪ ",
"blockhash": "208069b14e551149519671300e80419f944439cfc76c2a99dc16fd5a99502db0",
"confirmations": 302,
"time": 1617772939,
"blocktime": 1617772939
}

```

隔离见证规范

数据结构

在交易中新增了一个数据结构：witness，下面展示新旧数据结构

```

old: [nVersion] [input num] [tx inputs] [output num] [tx outputs] [nLockTime]

new: [nVersion] [marker] [flag] [input num] [tx inputs] [output num] [tx outputs]_
↪ [witness] [nLockTime]

```

以上面的交易数据为例，具体解析如下：

- P2PKH

```
[nVersion] [input num] [tx inputs] [output num] [tx outputs] [nLockTime]
```

- rawtransaction

```
0200000001c2290f6526caed6b203688fae51432c640bd8e2fcf50060b4f02348fae18c7040000000006b48304502
```

- 解析具体字段

```

02000000          // nVersion
01                // input number
c2290f6526caed6b203688fae51432c640bd8e2fcf50060b4f02348fae18c704 // 小端序的
↳previous tx hash
00000000          // tx vin index
6b                // 解锁脚本长度 10进制 107
483045022100b84617263695d31407d6ae230d3b629d071e65d936ee946c8adf84304163cb2d022045435bf766cf

ffffffff          // sequence 通常为0xFFFFFFFF。除非交易的锁定时间>
↳0, 否则不相关
02                // output num
40420f00000000001976a9144fd5b5cbe9d25a7cbbef9847d1a6e97d8a07697588ac
f695e605000000001976a914db9d0eb300a248a494f55ec1abb140d5d580216f88ac
00000000          // nLockTime 定义交易有效（或最小块高）的最早时间

```

• P2WPKH

```

[nVersion] [marker] [flag] [input num] [tx inputs] [output num] [tx outputs]
↳[witness][nLockTime]

```

- rawtransaction

```

02000000000010197df15a97d9b194bd378da4d670f584e808bdb0ab3335197357eab302fbe20c300000000ffff

```

- 解析具体字段

```

02000000          // nVersion
00                // marker 必须是 0
01                // flag 必须是 1
01
97df15a97d9b194bd378da4d670f584e808bdb0ab3335197357eab302fbe20c3 // 小端序的
↳previous tx hash
00000000          // tx vin index
00                // 解锁脚本长度为0, 因为在 witness 字段
ffffffff          // sequence 通常为0xFFFFFFFF。除非交易的锁定时间>
↳0, 否则不相关
02                // output num
40420f00000000001976a9144fd5b5cbe9d25a7cbbef9847d1a6e97d8a07697588ac
ec95e60500000000160014db9d0eb300a248a494f55ec1abb140d5d580216f
02                // txinwitness 有2部分
47                // txinwitness 签名部分的长度. 10进制 71
304402201e9df1d606a9fcbd159b62f1fd0675ba31092ef9d8faf8afa09b48c765d40c74022016f272575b03e22c
21                // txinwitness 公钥部分的长度. 10进制 33
029ba616506337ff7f0ccaac35938af4ba614b9853d670c4a45d830d7d38de7823
00000000          // nLockTime 定义交易有效（或最小块高）的最早时间

```

锁定脚本与验证脚本

P2WPKH

• 对应字段

```

witness:         <signature> <pubkey>
scriptSig:       (empty)

```

(下页继续)

(续上页)

```
scriptPubKey: 0 <20-byte-key-hash>
              (0x0014{20-byte-key-hash})
```

以上述的交易为例

```
witness:      [
  ↪ "304402201e9df1d606a9fcbd159b62f1fd0675ba31092ef9d8faf8afa09b48c765d40c74022016f272575b03e22c1
  ↪ ",
      "029ba616506337ff7f0ccaac35938af4ba614b9853d670c4a45d830d7d38de7823"
    ]
scriptSig:    (empty)
scriptPubKey: 0014db9d0eb300a248a494f55ec1abb140d5d580216f
```

- 实际验证的脚本

```
<signature> <pubkey> CHECKSIG
```

txid 与 hash

每个交易存在 hash 和 txid. txid 不对 witness 数据也进行 double sha256, 这样 witness 数据即便是修改了, 也无法影响 txid.

- txid

```
double sha256([nVersion] [input num] [tx inputs] [output num] [tx outputs]
  ↪ [nLockTime])
```

- hash

```
double sha256([nVersion] [marker] [flag] [input num] [tx inputs] [output num] [tx
  ↪ outputs] [witness] [nLockTime])
```

区块变相扩容

原本 Block 的大小限制为 1,000,000 字节 (1MB) .

BIP 141 提出一种新的计算方式: weight. 新的 Block 限制为 4M weight.

```
weight = (tx_size - witness_size) * 3 + tx_size < 4M
```

具体对应字段如下所示:

换句话说, witness_size 包含了 marker flag witness 1 byte 就对应到 1 weight

如果是 non-witness program 则 1 byte 就占了 4 weight.

例如上面例子: 从某个地址往 mno5h3qJZ3cT1TtFJZVErzFkF2BJQtkJMh 转账 0.01 BTC.

- p2pkh 计算 weight

```
(226 - 0) * 3 + 226 = 904
```

- p2wpkh 计算 weight

```
weight = (tx_size[225] - marker[1] - flag[1] - txinwitness[1 + 1 + 73 + 1 + 33])_
↪+ tx_size[225]
(225 - 1 - 1 - 1 - 1 - 71 - 1 - 33) + 225 = 573
```

按照以前的算法，Block 大小限制为 1,000,000 字节，最多可以放下

- $1,000,000 / 226 = 4424$ 笔 p2pkh 交易.
- $1,000,000 / 225 = 4444$ 笔 p2wpkh 交易. (当然，以前也没有这种交易)

但现在，Block 的 weight 限制为 4 M weight。

- $4,000,000 / 904 = 4424$ 笔 p2pkh 交易. (对于以前的交易存储限制其实并无影响)
- $4,000,000 / 573 = 6980$ 笔 p2wpkh 交易.

相同的目的（转账 0.01BTC）可以再一个块内多存储 2 千多笔交易，变相的扩容 1.5 倍了。

Reference

- [Bitcoin Transactions](#)
- [BIP-0141](#)

2.3.5 BTC 系列之数据存储

ETH 可以在交易构造之时直接使用 data 字段存储数据。BTC 存储数据有不同的方式

背景介绍

- dust: 构造交易时，一般认为当交易费用高于 1/3 交易价值时，即可称作“Dust”或尘埃交易，目前而言，尘埃交易是指交易价值低于 546 satoshis 比特币 (即 0.000000546 BTC) 的交易。
- 标准交易和非标准交易

在比特币中，最常见的交易形式为“鲍勃付爱丽丝 (Bob pays Alice)”，它基于“付给公共密钥哈希 (P2PKH)”脚本，该脚本通过发送公钥和由以下人员创建的数字签名来解决：相应的私钥。P2PKH 交易只是许多标准类别之一：如果交易通过了 *Bitcoin Core IsStandard ()* 和 *IsStandardTx ()* 测试，则该交易为标准交易。但是，通过创建临时脚本来锁定 (和解锁) 交易，还可以生成非标准的脚本。交易，但仍可以广播和挖掘。

OP_RETURN

op_return 存储的最大字节数设置为 83 个字节 (默认为 80 个字节，外加 3 个字节的开销)。

但一个 BTC 标准交易只能有一个 op_return 的输出。

OP_DROP

op_drop 运算符删除顶部堆栈元素。

P2SH 脚本

构造一个赎回脚本，锁定脚本为其 hash。解锁脚本需要用这个赎回脚本去验证锁定脚本，同时保障赎回脚本通过。

- 具体脚本

```

赎回脚本 <OP_DUP OP_HASH160 0xca4159109e41b78863d23ceb5493e770cc53c6f4 OP_
↳EQUALVERIFY OP_CHECKSIG>
解锁脚本 <SIGN> <PUBKEY> <赎回脚本>
锁定脚本 OP_HASH160 ebcc108d534d6bab12f9f657d29e9c982ea68453 OP_EQUAL

```

- 解锁过程

- 两个脚本经由两步实现组合。首先，将赎回脚本与锁定脚本比对以确认其与哈希是否匹配:

```

<OP_DUP OP_HASH160 0xca4159109e41b78863d23ceb5493e770cc53c6f4 OP_EQUALVERIFY_
↳OP_CHECKSIG> OP_HASH160 ebcc108d534d6bab12f9f657d29e9c982ea68453 OP_EQUAL

```

- 假如赎回脚本与哈希匹配，解锁脚本会被执行以释放赎回脚本

```

<SIGN> <PUBKEY> OP_DUP OP_HASH160 0xca4159109e41b78863d23ceb5493e770cc53c6f4_
↳OP_EQUALVERIFY OP_CHECKSIG

```

OP_DROP 存储数据

利用 P2SH 脚本，先构造一个由 OP_DROP 和 P2PKH 组合成的赎回脚本。

```

赎回脚本 <OP_DROP OP_DUP OP_HASH160 <PUBKEY> OP_EQUALVERIFY OP_CHECKSIG>
解锁脚本 <SIGN> <PUBKEY> <DATA> <赎回脚本>
锁定脚本 OP_HASH160 <RedeemScriptHash> OP_EQUAL

```

具体解锁过程如下：

- 首先，将赎回脚本与锁定脚本比对以确认其与哈希是否匹配

```

<OP_DROP OP_DUP OP_HASH160 <PUBKEY> OP_EQUALVERIFY OP_CHECKSIG> OP_HASH160
↳<RedeemScriptHash> OP_EQUAL

```

- 执行解锁脚本. 上链数据，存储在 Data 字段，这样第一个字段就会被丢弃

```

<SIGN> <PUBKEY> <Data > OP_DROP OP_DUP OP_HASH160 <PUBKEY Hash> OP_EQUALVERIFY_
↳OP_CHECKSIG

```

难题和进度

OP_RETURN

每笔交易 unlock 2 条 burn tx 记录。因为 burn tx hash 长度为 32，能容纳 2 笔。

目前构造的交易 size 大小为 449 byte，按照 100 聪左右的手续费率，大概每笔 unlock 需要 0.00044900 BTC (蛮贵的)

P2SH + OP_DROP

看论文大概能存储 1500 左右字节的数据

- 多签 unlock 地址会变更，目前尚未测试多签 + op_drop

目前本地构造的交易，赎回脚本和锁定脚本 OK 了，但是解锁脚本尚未通过签名验证

OP_DROP + P2PKH

解锁脚本	<SIGN>	<PUBKEY>	<DATA>
锁定脚本	OP_DROP	OP_DUP	OP_HASH160 <PUBKEY> OP_EQUALVERIFY OP_CHECKSIG

属于非标准交易，目前尚未找到使矿工节点忽略标准检查，打包非标准交易的方法。

Reference

- [Data Insertion in Bitcoin's Blockchain](#)
- [An Analysis of Non-standard Transactions](#)

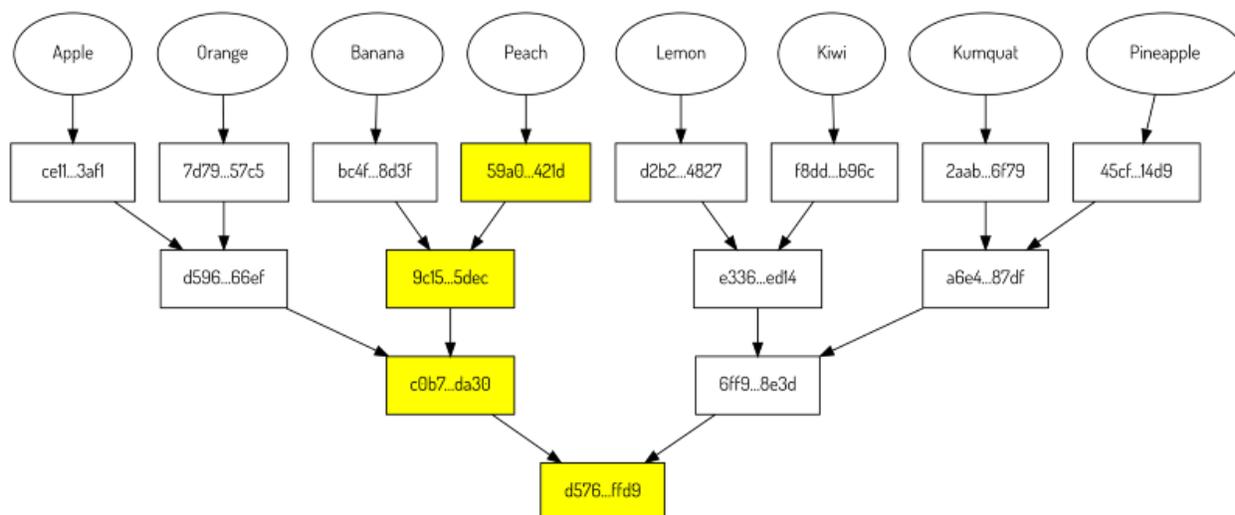
2.4 Crypto

2.4.1 默克尔树

默克尔证明

默克尔路径

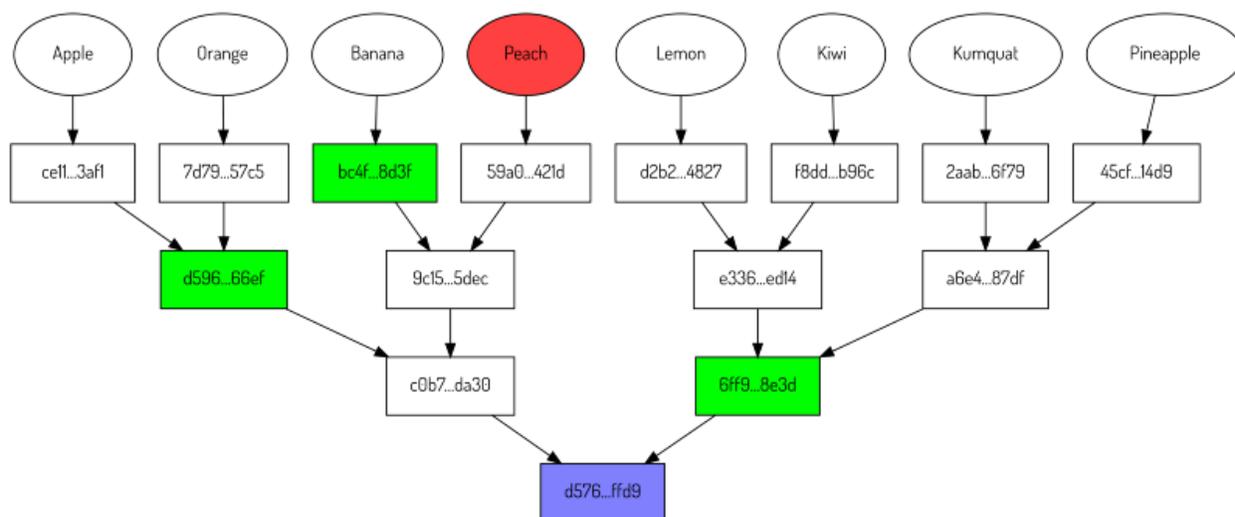
默克尔路径指的是某个输入值到默克尔根节点之间所有哈希值的集合。下图显示了输入值“Peach”的默克尔路径：



-“Peach” 的默克尔路径-

默克尔证明

默克尔证明指的是不需要知道一个数据集中的其他值就能证明某个值属于这个集合。



-默克尔证明-

默克尔证明需要三样东西：输入值（红色标记）、中间支哈希值（绿色标记）和默克尔根节点（蓝色标记）。每个输入值对应的中间支哈希值集合各不相同。

区块链系统经常会用到默克尔证明，证明某个数据集合内存在某个输入值，这样就不需要将整个数据集合都存储在区块链上了。假设一个以太坊合约内有一个白名单列表，只允许列表内的账户购买以太坊。如果将白名单内每个账户信息都存储在区块链上，势必要付出很高的成本。在这种情况下，只需要创建一个默克尔树，再将根节点存储在区块链上即可。

例如，如果将根节点存储在一个智能合约上，这个智能合约很容易就能证明某个账户包含在白名单内：这个账户需提供中间支哈希值（合约所有者通过某种链下方式提供给账户持有者），智能合约将这个账户的哈希值依次与中间支哈希值进行哈希计算。如果最后得出的结果与默克尔根节点一致的话，就证明这个账户确实存在在白名单里。

请注意最后两张图中默克尔路径和默克尔证明的哈希值之间的关系。在同一棵树的同一个层级中，默克尔证明的哈希值与默克尔路径的哈希值是相互关联的。由此可见，默克尔证明能够重塑输入值的默克尔路径，这就是为什么最终结果是默克尔根节点的原因。

至此，可以看出默克尔证明具有以下特征：

- 在链上存储默克尔证明所需的空间远远小于直接存储输入值所需的空间
- 在链上公开存储默克尔证明也不会暴露整个输入值集合
- 要证明某个输入值集合内是否存在某个值，验证默克尔证明的成本低于核对整个输入值集合的成本

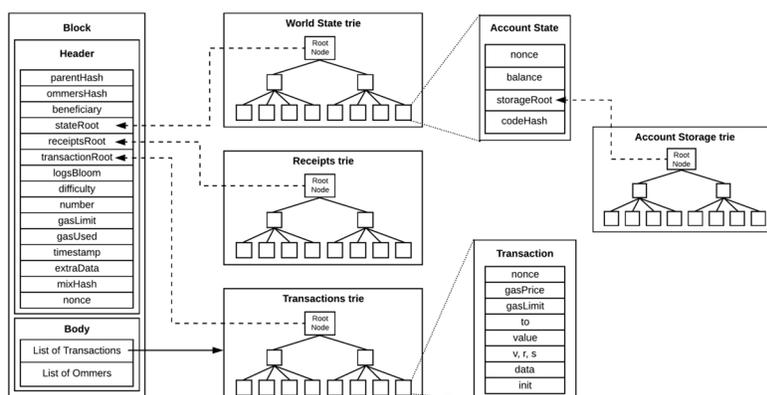
ETH 中的默克尔实现

1. **世界状态树**包括了从地址到账户状态之间的映射。世界状态树的根节点哈希值由区块保存（在 `stateRoot` 字段），它标示了区块创建时的当前状态。整个网络中只有一个世界状态树。
2. **账户存储树**保存了与某一智能合约相关的数据信息。由账户状态保存账户存储树的根节点哈希值（在 `storageRoot` 字段）。每个账户都有一个账户存储树。
3. **交易树**包含了一个区块中的所有交易信息。由区块头（在 `transactionsRoot` 区域）保存交易树的根节点哈希值。每个区块都有一棵交易树。
4. **交易收据树**包含了一个区块中所有交易的收据信息。同样由区块头（在 `receiptsRoot` 区域）保存交易收据树的根节点哈希值；每个区块都有对应的交易收据树。

我们今天讨论的对象有：

1. **世界状态**：以太坊这台分布式计算机的硬盘。它是从地址到账户状态的映射。
2. **账户状态**：保存着每个以太坊账户的状态信息。账户状态同样保存着账户状态树的 `storageRoot`，后者包含了该账户的存储数据。
3. **交易**：标示了系统中的状态转移。它可以是资金的转移、消息调用或是合约的部署。
4. **区块**：包括对前序区块（`parentHash`）的链接，并且保存了当执行时会在系统中产生新状态的交易。区块同时保存了 `stateRoot`、`transactionRoot`、`receiptsRoot`、世界状态树的根节点哈希、交易树以及对应的交易收据树。

我想用一张图来表示文中提及的各种概念信息。



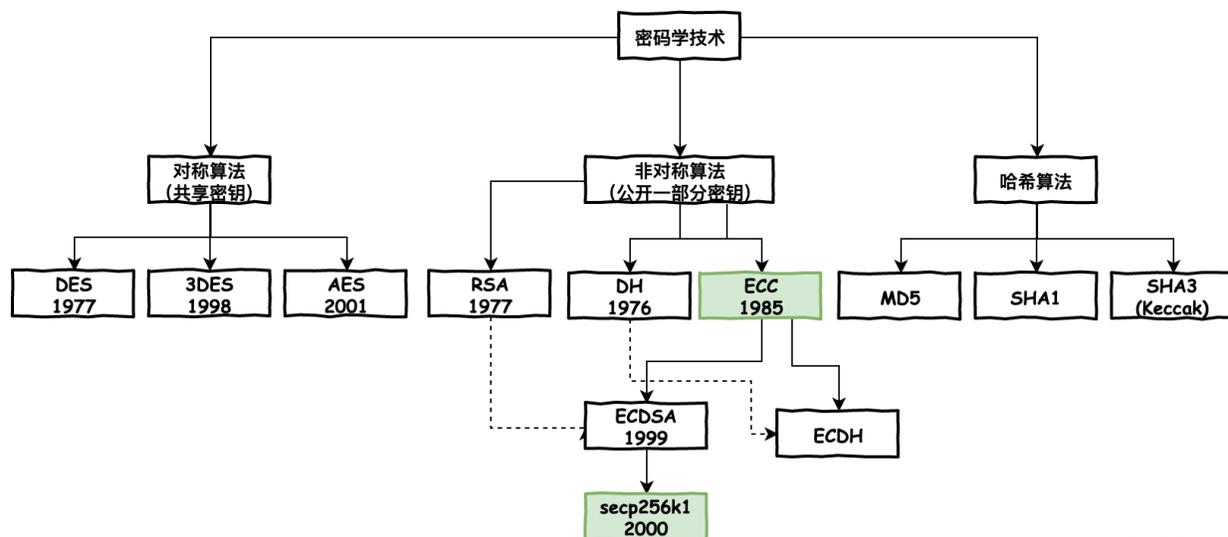
-区块、交易、账户状态对象以及以太坊的默克尔树-

Reference

- 什么是 Merkle Pollard

2.4.2 签名

签名算法类别



以太坊技术与实现

简单理解椭圆曲线

- 模乘运算：已知 2 个大质数 (x, y) ，相乘算出结果比较容易，例如 $2003 * 4003 = 8018009$ ，
- 模除运算：但已知相乘结果 8018009，推算出 (x, y) 成本是要高的。

在椭圆曲线运算中，质数的模乘运算较为简单，但模除运算（反向运算）几乎不可能。（离散对数难题，笔者望而却步）。

ETH 和 BTC 的签名算法：椭圆曲线算法 secp256k1。

签名中最重要的是 checksig，bitcoin 单签执行一定包括 signature pubkey checksig，即针对签名结果进行验签

椭圆曲线数字签名算法 (ECDSA)

主要信息

- 用户的密钥对: (d, Q) ; (d 为私钥, Q 为公钥)
- 待签名的信息: M ;
- 签名: $Signature(M) = (r, s)$

签名过程:

已知私钥 d , 和待签名的消息 M , 计算其签名 (r, s)

- 1、根据 ECC 算法随机生成一个临时密钥对 (k, R) , $R=(xR, yR)$ 。(曲线上 R 点 (x, y) 的大整数)
- 2、令 $r = xR \bmod n$, 如果 $r=0$, 则返回步骤 1. 此处可计算得到 r
- 3、计算 $H = \text{Hash}(M)$
- 4、按照数据类型转换规则, 将 H 转化为一个 big endian 的整数 e
- 5、 $s = k^{-1} (e + rd) \bmod n$, 若 $s=0$, 则返回步骤 1
- 6、输出的 $S=(r,s)$ 即为签名

可参考代码 [threshold-signature-demo](#)

```
def sign(private_key: int, message: bytes) -> tuple:
    """Create ECDSA signature (r, s)"""
    e = hash_to_int(message)
    r, s = 0, 0
    while not r or not s:
        k = random.randrange(1, curve.n)
        k_x, _ = scalar_multiply(k, curve.g)
        r = k_x % curve.n
        s = ((e + r * private_key) * modular_multiplicative_inverse(k, curve.n))
    ↪ % curve.n
    return r, s
```

验证过程:

已知公钥 Q , 和待签名的消息 M , 签名 (r, s) , 验证签名是否有效

- 1、计算 $H = \text{Hash}(M)$
- 2、按照数据类型转换规则, 将 H 转化为一个 big endian 的整数 e
- 3、计算 $u1 = es^{-1} \bmod n, u2 = rs^{-1} \bmod n$
- 4、计算 $R = (xR, yR) = u1G + u2Q$, 如果 $R =$ 零点, 则验证该签名无效
- 5、令 $v = xR \bmod n$
- 6、若 $v == r$, 则签名有效, 若 $v \neq r$, 则签名无效

可参考代码. [threshold-signature-demo](#)

```
def verify_signature(public_key: tuple, message: bytes, signature: tuple) -> bool:
    """Verify signature with public key and message"""
    e = hash_to_int(message) # 计算 message 的 hash, 并且转换为 BigInteger e
    r, s = signature # 签名其实就是 r, s 值
    w = modular_multiplicative_inverse(s, curve.n) # 计算 s^{-1}
    u1 = (w * e) % curve.n # u1 = e * s^{-1} mod n
    u2 = (w * r) % curve.n # u2 = r * s^{-1} mod n
    x, _ = add(scalar_multiply(u1, curve.g), scalar_multiply(u2, public_key))
    return r == (x % curve.n)
```

恢复过程 Recover

已知消息 M 和签名 (r, s) ，恢复计算出公钥 Q 。

- 1、计算 $H = \text{Hash}(M)$
- 2、按照数据类型转换规则，将 H 转化为一个 big endian 的整数 e
- 3、计算 $R=(x, y)$ ，其中 $x=r, r+n, r+2n\cdots$ ，代入椭圆曲线方程计算获得 R
- 4、计算 $u1 = -er^{-1} \bmod n, u2 = sr^{-1} \bmod n$
- 5、计算公钥 $Q= (x', y') = u1 \cdot G + u2 \cdot R$

ETH 黄皮书

$$\text{ECDSAPUBKEY}(pr \in B32) \equiv pu \in B64$$

$$\text{ECDSASIGN}(e \in B32, pr \in B32) \equiv (v \in B1, r \in B32, s \in B32)$$

$$\text{ECDSARECOVER}(e \in B32, v \in B1, r \in B32, s \in B32) \equiv pu \in B64$$

Reference

- [threshold-signature-demo](#)
- 椭圆曲线数字签名算法
- 比特币系统采用的公钥密码学方案和 ECDSA 签名算法介绍——第一部分：原理
- 一个数字引发的探索——ECDSA 解析
- 以太坊技术与实现

2.5 Wallet

2.5.1 加密与签名

加密算法

椭圆曲线加密 (ECC: Elliptic Curve Cryptography)

1. 原理性质

椭圆曲线加密法是一种基于离散对数问题的非对称加密法，可以用对椭圆曲线上的点进行加法或乘法运算来表达。下图是一个椭圆曲线的示例，类似于比特币所用的曲线。

secp256k1

$$y^2 = (x^3 + 7) \text{ over } (F_p)$$

$$y^2 \bmod p = (x^3 + 7) \bmod p$$

上述 $\text{mod } p$ (素数 p 取模) 表明该曲线是在素数阶 p 的有限域内, 也写作 F_p , 其中 $p = 2^{256} - 2^{32} - 2^9 - 2^8 - 2^7 - 2^6 - 2^4 - 1$, 这是个非常大的素数

数字签名 (ECDSA)

BTC 数字签名算法是椭圆曲线数字签名算法 (Elliptic Curve Digital Signature Algorithm) 或 ECDSA。

ECDSA 用于脚本函数 `OP_CHECKSIG`, `OP_CHECKSIGVERIFY`, `OP_CHECKMULTISIG` 和 `OP_CHECKMULTISIGVERIFY`。

数字签名在比特币中有三种用途。

- 第一, 签名证明私钥的所有者, 即资金所有者, 已经授权支出这些资金
- 第二, 授权证明是不可否认的 (不可否认性)
- 第三, 签名证明交易 (或交易的具体部分) 在签字之后没有也不能被任何人修改。

通常使用椭圆曲线算法生成密钥对

- 比特币密钥长度: 256 位
- 公钥哈希值 = `RIMPED160(SHA256(公钥))`
- 比特币地址 = `1+Base58(0+ 公钥哈希值 + 校验码)`
- 校验码 = 前四字节 (`SHA256(SHA256(0+ 公钥哈希值))`)

签名

- 发送方使用 HASH 算法计算数据的 HASH 值
- 发送方使用本方的私钥加密 HASH 值, 得到签名
- 接收方使用 HASH 算法计算数据的 HASH 值
- 接收方使用发送方的公钥解密签名得到发送的 HASH 值
- 比较两个 HASH 值的一致性

数字签名如何工作

数字签名是一种由两部分组成的数学方案:

- 第一部分是使用私钥 (签名密钥) 从消息 (交易) 创建签名的算法;
- 第二部分是允许任何人验证签名的算法, 给定消息和公钥。

创建签名

在比特币的 ECDSA 算法的实现中, 被签名的“消息”是交易, 或更确切地说是交易中特定数据子集的哈希值 (参见签名哈希类型 (SIGHASH))。签名密钥是用户的私钥, 结果是签名:

$$\text{Sig} = F\{\text{sig}\}(F\{\text{hash}\}(m), \quad dA)$$

- dA 是签名私钥
- m 是交易 (或其部分)
- F_{hash} 是散列函数
- F_{sig} 是签名算法

- Sig 是结果签名

函数 `Fsig` 产生由两个值组成的签名 `Sig`，通常称为 `R` 和 `S`：

验证签名

要验证签名，必须有签名（`R` 和 `S`）、序列化交易和公钥（对应于用于创建签名的私钥）。本质上，签名的验证意味着“只有生成此公钥的私钥的所有者，才能在此交易上产生此签名。”

签名验证算法采用消息（交易或其部分的哈希值）、签名者的公钥和签名（`R` 和 `S` 值），如果签名对该消息和公钥有效，则返回 `TRUE` 值。

ECDSA 数学

签名由产生由两个值 `R` 和 `S` 组成的签名的数学函数 `Fsig` 创建。

签名算法首先生成一个 `ephemeral`（临时）私公钥对。在涉及签名私钥和交易哈希的变换之后，该临时密钥对用于计算 `R` 和 `S` 值。

临时密钥对基于随机数 `k`，用作临时私钥。从 `k`，我们生成相应的临时公钥 `P`（以 $P = k * G$ 计算，与派生比特币公钥相同）；参见 `[pubkey]` 部分）。数字签名的 `R` 值则是临时公钥 `P` 的 `x` 坐标。

加密常用库

JavaScript :

Elliptic

支持以下曲线

- `secp256k1 ==ECDSA==`
- `p192`
- `p224`
- `p256`
- `p384`
- `p521`
- `curve25519 ==ECDH==`
- `ed25519 ==EdDSA==`

Eccrypto

Only `secp256k1 curve ==ECDSA== ==ECDH== ==ECIES==`

tiny-secp256k1

crypto-browserify

实现以下：

- createHash (sha1, sha224, sha256, sha384, sha512, md5, rmd160)
- createHmac (sha1, sha224, sha256, sha384, sha512, md5, rmd160)
- pbkdf2
- pbkdf2Sync
- randomBytes
- pseudoRandomBytes
- createCipher (aes)
- createDecipher (aes)
- createDiffieHellman
- createSign (rsa, ecdsa)
- createVerify (rsa, ecdsa)
- createECDH (secp256k1)
- publicEncrypt/privateDecrypt (rsa)
- privateEncrypt/publicDecrypt (rsa)

2.5.2 多重签名

多重签名脚本设置了一个条件，其中 N 个公钥被记录在脚本中，并且至少有 M 个必须提供签名来解锁资金。这也称为 M - N 方案，其中 N 是密钥的总数， M 是验证所需的签名的数量。例如， $2/3$ 的多重签名是三个公钥被列为潜在签名人，至少有 2 个有效的签名才能花费资金。

设置 M - N 多重签名条件的锁定脚本的一般形式是：

```
M <Public Key 1> <Public Key 2> ... <Public Key N> N CHECKMULTISIG
```

M 是花费输出所需的签名的数量， N 是列出的公钥的总数。设置 2 到 3 多重签名条件的锁定脚本如下所示：

```
2 <Public Key A> <Public Key B> <Public Key C> 3 CHECKMULTISIG
```

上述锁定脚本可由含有签名和公钥的脚本予以解锁：或者由 3 个存档公钥中的任意 2 个相一致的私钥签名组合予以解锁。两个脚本组合将形成一个验证脚本：

```
<Signature B> <Signature C> 2 <Public Key A> <Public Key B> <Public Key C> 3  
↪CHECKMULTISIG
```

上述例子中相应的设置条件即为：解锁脚本是否含有 3 个公钥中的任意 2 个相对应的私钥的有效签名。

reference

椭圆曲线密码学相关概念与开源实现精通比特币 - 密钥和地址

2.5.3 钱包常用知识

钱包登录

- 单链钱包：只支持一条链的钱包，比方说比特币钱包。
- 多链钱包：支持所有币种的钱包。
- 单账户钱包：每种币只支持一个账户。
- 多账户钱包：每种币支持多个账户

钱包登录的方案三种：

钱包密码

可用于转账确认，调用合约等，以及辅助 keystore 登录

助记词

私钥是 64 位长度的十六进制的字符，不利于记录且容易记错，且每个账号对应一个私钥，多个账号就有多个私钥，不易管理，所以用算法将一串随机数转化为了一串 12 ~ 24 个容易记住的单词，方便保存记录。

通过助记词可以获取相关联的多个私钥，但是通过其中一个私钥是不能获取助记词的，因此助记词 ≠ 私钥。

根据助记词推算种子的算法是 **PBKDF2**，使用的哈希函数是 Hmac-SHA512，其中，输入是助记词的 UTF-8 编码，并设置 Key 为 mnemonic+password，循环 2048 次，得到最终的 64 字节种子。

KeyStore

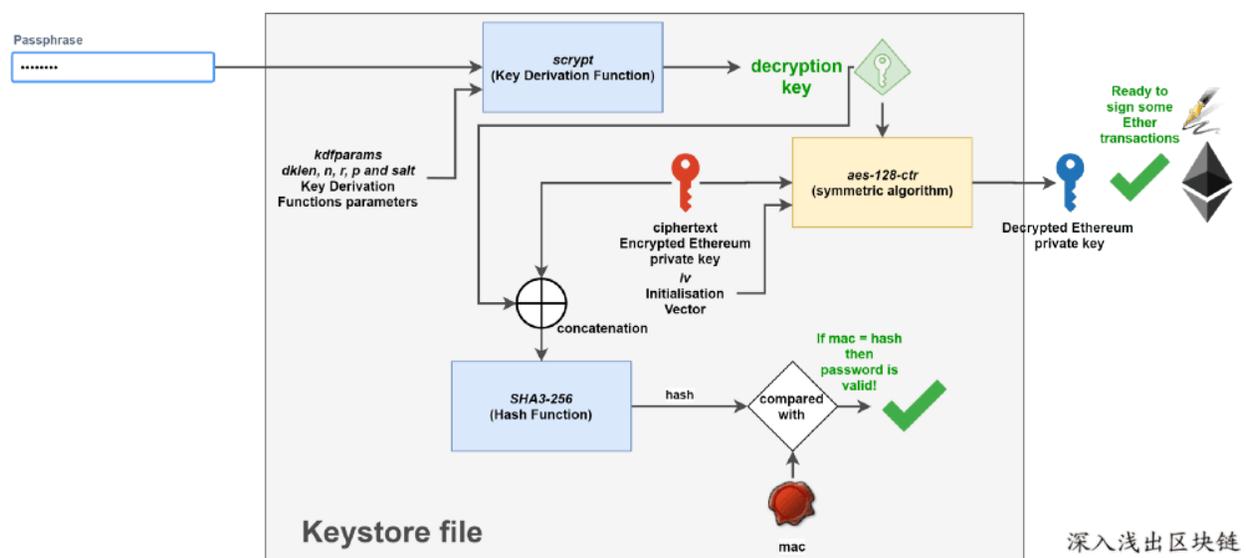
```
{
  "address": "856e604698f79cef417aab...",
  "crypto": {
    "cipher": "aes-128-ctr",
    "ciphertext": "13a3ad2135bef1ff228e399dfc8d7757eb4bb1a81d1b31....",
    "cipherparams": {
      "iv": "92e7468e8625653f85322fb3c..."
    },
    "kdf": "scrypt",
    "kdfparams": {
      "dklen": 32,
      "n": 262144,
      "p": 1,
      "r": 8,
      "salt": "3ca198ce53513ce01bd651aee54b16b6a...."
    },
    "mac": "10423d837830594c18a91097d09b7f2316..."
  },
  "id": "5346bac5-0a6f-4ac6-baba-e2f3ad464f3f",
  "version": 3
}
```

- address: 账号地址
- version: Keystore 文件的版本，目前为第 3 版，也称为 V3 Keystore。
- salt 是一段随机的盐
- dk_len 是输出的哈希值的长度
- n 是 CPU/Memory 开销值，越高的开销值，计算就越困难。
- r 表示块大小
- p 表示并行度。
- id : uuid
- crypto: 加密推倒的相关配置:
 1. cipher 是用于加密以太坊私钥的对称加密算法。用的是 aes-128-ctr 。
 2. cipherparams 是 aes-128-ctr 加密算法需要的参数。在这里，用到的唯一的参数 iv。
 3. ciphertext 是加密算法输出的密文，也是将来解密时的需要的输入。
 4. kdf: 指定使用哪一个算法，这里使用的是 scrypt。
 5. kdfparams: scrypt 函数需要的参数
 6. mac: 用来校验密码的正确性， $mac = sha3(DK [16:32], ciphertext)$

Keystore 文件的产生：

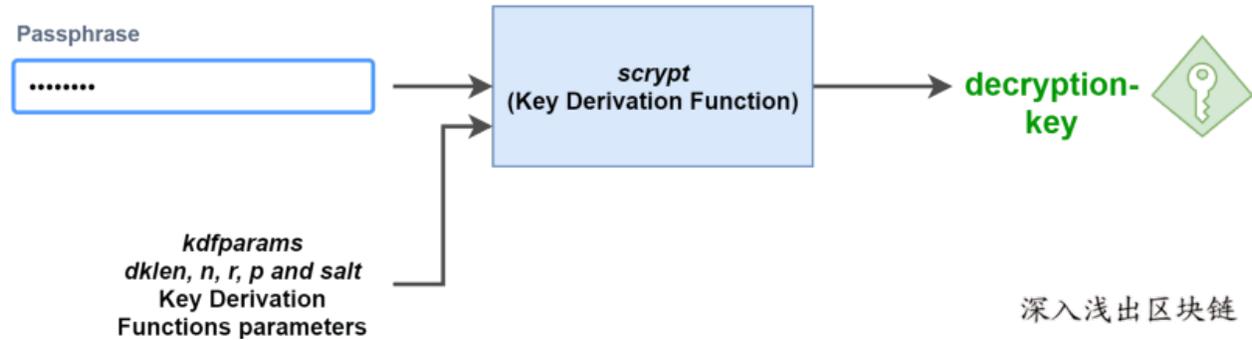
- 1、使用 scrypt 函数（根据密码和相应的参数）生成密钥
- 2、使用上一步生成的密钥 + 账号私钥 + 参数进行对称加密。
- 3、把相关的参数和输出的密文保存为以上格式的 JSON 文件

整体流程



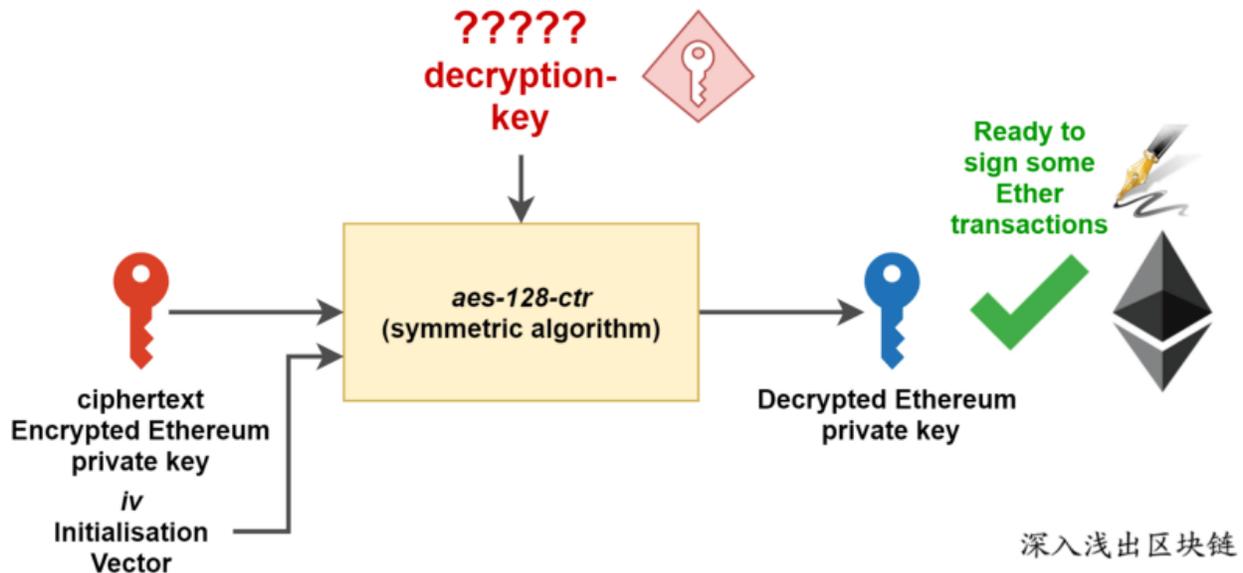
在 Keystore 中，用的是 **Scrypt** 算法，用一个公式来表示的话，派生的 Key 生成方程为：

$$DK = \text{Scrypt}(\text{salt}, \text{dk_len}, n, r, p)$$



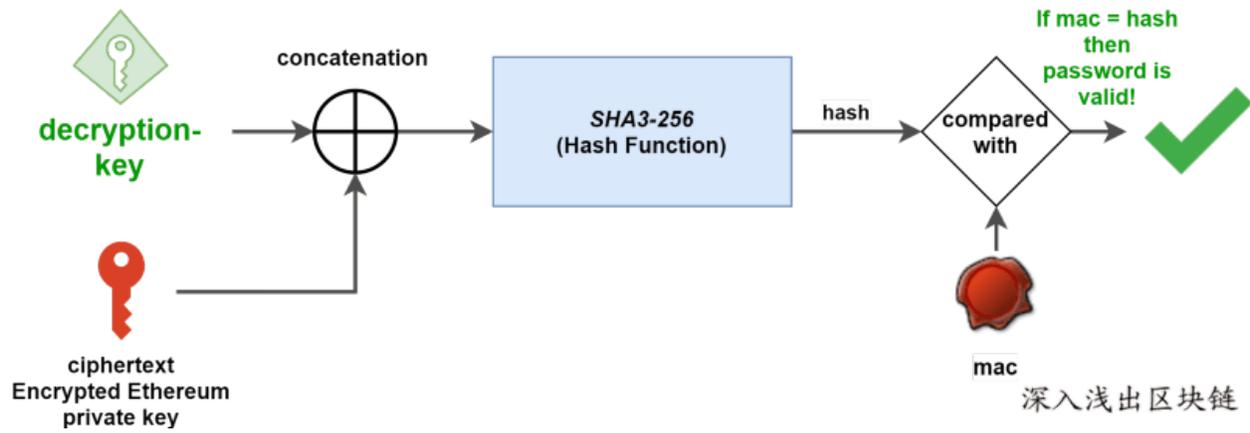
校验 Keystore

当我们在使用 Keystore 文件来还原私钥时，依然是使用 kdf 生成一个密钥，然后用密钥对 ciphertext 进行解密，其过程如下：



无论使用说明密码，来进行这个操作，都会生成一个私钥，但是最终计算的以太坊私钥到底是不是正确的，却不得而知。

这就是 keystore 文件中 mac 值的作用。**mac 值是 kdf 输出和 ciphertext 密文进行 SHA3-256 运算的结果，显然密码不同，计算的 mac 值也不同，因此可以用来检验密码的正确性。**

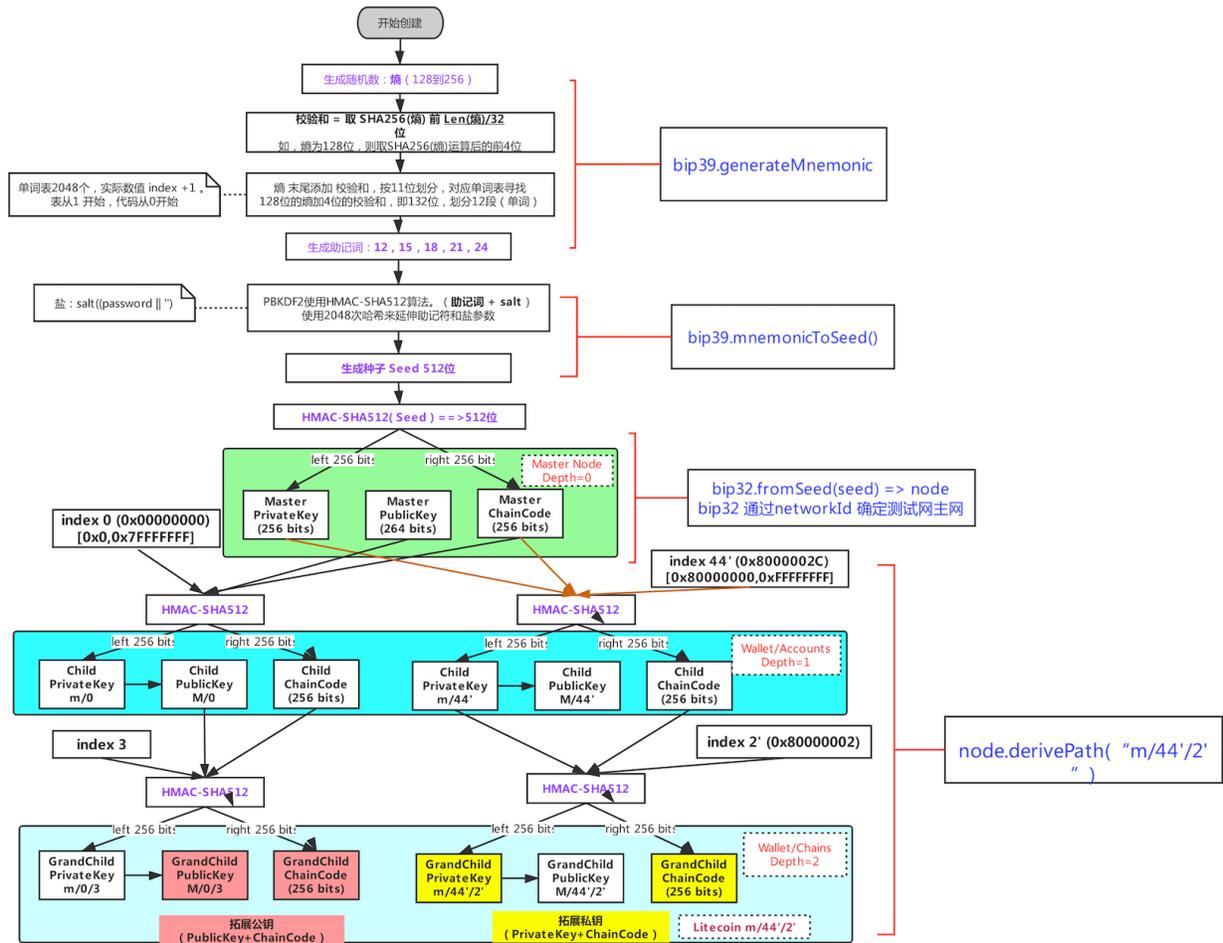


reference

账号 Keystore 文件导入导出精通比特币-钱包

2.5.4 BIP39、BIP44、BIP32 协议

HD 钱包导入流程



BIP39

熵每次都可以得到不同的助记词

```

CS = ENT / 32
checksum = SHA256(entropy)
MS = (ENT + CS) / 11
    
```

ENT: 熵长度。CS: 校验长度。MS: 助记词长度。checksum: 校验和校验和: 对熵取 SHA256 哈希结果前 CS 位

创建助记词

1. 创建一个 128 到 256 位的随机序列（熵）。
2. 提出 SHA256 哈希前几位（熵长/32），就可以创造一个随机序列的校验和。
3. 将校验和添加到随机序列的末尾。
4. 将序列划分为包含 11 位的不同部分。
5. 将每个包含 11 位部分的值与一个已经预先定义 2048 个单词的字典做对应。
6. 生成的有顺序的单词组就是助记码。

助记词生成种子（seed）

创建助记词之后的 7-9 步是：

1. PBKDF2 密钥延伸函数的第一个参数是从步骤 6 生成的助记符。
2. PBKDF2 密钥延伸函数的第二个参数是盐。由字符串常数“助记词”与可选的用户提供的密码字符串连接组成。
3. PBKDF2 使用 HMAC-SHA512 算法，使用 2048 次哈希来延伸助记符和盐参数，产生一个 512 位的值作为其最终输出。这个 512 位的值就是种子。

图 5-7 显示了从助记词如何生成种子

BIP32

分层路径

从 Seed 创建 HDWallet

HD 钱包从单个根种子（root seed）中创建，为 128 到 256 位的随机数。常见的通过助记词可以确定种子 Root Seed。

HD 钱包的所有的确定性都衍生自这个根种子。任何兼容 HD 钱包的根种子也可重新创造整个 HD 钱包。

- 父密钥（没有压缩过的椭圆曲线推导的私钥或公钥 ECDSA uncompressed key）
- 链码作为熵（chain code 256 bits）
- 子代索引序号（index 32 bits）

钱包安全的核心在私钥，而公钥则比较容易被找到，如果子节点生成过程只依赖父节点公钥和子节点序号，那么黑客拿到父节点公钥之后就能复原出所有子节点、孙节点的公钥，这样就会破坏隐私性，CKD 里面引入的 Chain Code 则是在整个子节点派生过程中引入确定的随机数，为 HD 钱包的隐私性增加了一重保障。

正常衍生子密钥

强化衍生子密钥

- 图上标识了索引号码根据正常衍生和硬化衍生不同，索引的范围不同，对于正常衍生的索引号范围为 $[0x0, 0x7FFFFFFF]$ ，而硬化衍生的索引号范围为 $[0x80000000, 0xFFFFFFFF]$
- 硬化衍生的索引号太长，一般为了便于阅读，都是会将索引号右上角加上撇号，譬如：0x80000000 记为 0'，0x80000001 记为 1'，以此类推。

针对扩展密钥的学习，可以看到有三种生成规则，分别如下：

- Private parent key -> private child key 即，从父私钥和父链码计算生成子私钥和子链码。用公式表示就是：

$$\text{CKDpriv}((\text{kpar}, \text{cpar}), i) \rightarrow (\text{ki}, \text{ci})$$

- 检查是否 $i \geq 2^{31}$ (子私钥)。如果是 (硬化的子密钥)：让 $I = \text{HMAC-SHA512}(\text{Key} = \text{cpar}, \text{Data} = 0x00 \parallel \text{ser256}(\text{kpar}) \parallel \text{ser32}(i))$ (注意：0x00 将私钥补齐到 33 字节长。) 如果不是 (普通的普通子密钥)：让 $I = \text{HMAC-SHA512}(\text{Key} = \text{cpar}, \text{Data} = \text{serP}(\text{point}(\text{kpar})) \parallel \text{ser32}(i))$ 。
- 将 I 分为两个 32 字节序列，IL 和 IR。
- 返回的子密钥 ki 是 $\text{parse256}(\text{IL}) + \text{kpar} \pmod{n}$ 。
- 返回的链码 ci 是 IR。
- 如果 $\text{parse256}(\text{IL}) \geq n$ 或 $\text{ki} = 0$ ，则生成的密钥无效，并且应继续下一个 i 值。(注：概率低于 1/2127)
- Public parent key -> public child key 即，从父公钥和父链码计算生成子公钥和子链码。公式表示如下：

$$\text{CKDpub}((\text{Kpar}, \text{cpar}), i) \rightarrow (\text{Ki}, \text{ci})$$

- 检查是否 $i \geq 2^{31}$ (子密钥是否是硬化密钥) 如果是 (硬化子密钥)：返回失败如果不是 (普通子密钥)：让 $I = \text{HMAC-SHA512}(\text{Key} = \text{cpar}, \text{Data} = \text{serP}(\text{Kpar}) \parallel \text{ser32}(i))$ 。
- 将 I 分为两个 32 字节序列，IL 和 IR。
- 返回的子密钥 Ki 是 $\text{point}(\text{parse256}(\text{IL})) + \text{Kpar}$ 。
- 返回的链码 ci 是 IR。
- 如果 $\text{parse256}(\text{IL}) \geq n$ 或 Ki 是无限远的点，则生成的密钥无效，并且应继续下一个 i 值。
- Private parent key -> public child key
即，从父私钥和父链码计算生成子公钥和子链码。公式表示如下：

$$N((k, c)) \rightarrow (K, c)$$

- 返回的密钥 K 是 $\text{point}(k)$ 。
- 返回的链码 c 只是传递的链码。

要计算父私钥的公用子密钥：

- $N(\text{CKDpriv}((\text{kpar}, \text{cpar}), i))$ (总是工作)。
- $\text{CKDpub}(N(\text{kpar}, \text{cpar}), i)$ (仅适用于非硬化子密钥)。

它们等价的事实是使非硬化密钥有用 (可以在不知道任何私钥的情况下导出给定父密钥的子公钥)，以及它们与硬密钥的区别。不总是使用非硬化键 (更有用) 的原因是安全性；后面可以了解更详细的信息。

BIP44

BIP32 路径中的 5 个层次定义:

```
m / purpose' / coin_type' / account' / change / address_index
```

路径中的撇号表示使用了经过硬化 (hardern) 处理的 BIP32 派生。

1. purpose': 固定值 44'(或 0x8000002C), 代表是 BIP44
2. coin_type': 这个代表的是币种, 可以兼容很多种币, 比如 BTC 是 0', ETH 是 60' 硬化的派生在这个级别上使用。(所以右上角有' 标记) 币种, 代表一个主节点 (种子) 可用于无限数量的独立加密货币, 如比特币等。此级别为每个加密货币创建一个单独的子树, 避免重用已经在其它链上存在的地址。开发人员可以为他们的项目注册未使用的号码
3. Account: 一类 coin 下能够有多个账户, 就相当于你的人民币会存在多张银行卡以顺序递增的方式从索引 0 开始编号。这个数字在 BIP32 派生中用作子索引。
4. Change: 表明为外部链 0 或内部链 1: 常量 0 用于外部链, 常量 1 用于内部链 (也称为更改地址)。外部链用于在钱包外部可见的地址 (例如用于接收付款)。内部链用于不在钱包外部可见的地址, 用于返回交易更改。一个是用来接收地址一个是用来创造找零地址
5. Index: 被 HD 钱包衍生的真正可用的地址是第四层级的子级, 就是第五层级的树的 "address_index"。比如, 第三个层级的主账户如 M/44'/0'/1' 收到比特币真正支付的地址是 M/44'/0'/0'/0/2

最后我们简要介绍一下 BIP44 中的账户发现算法:

1. 推导出第一个账户节点 (生成编号初始值设为 0)
2. 推导出这个账户的外部链
3. 依次扫描外部链上的地址, 如果连续二十个地址都没有交易记录, 停止扫描
4. 如果外部链上没有发现交易, 退出
5. 如果存在交易, 将生成编号的值增加 1, 跳转到第一步

一句话概括下 BIP44 就是: 给 BIP32 的分层路径定义规范

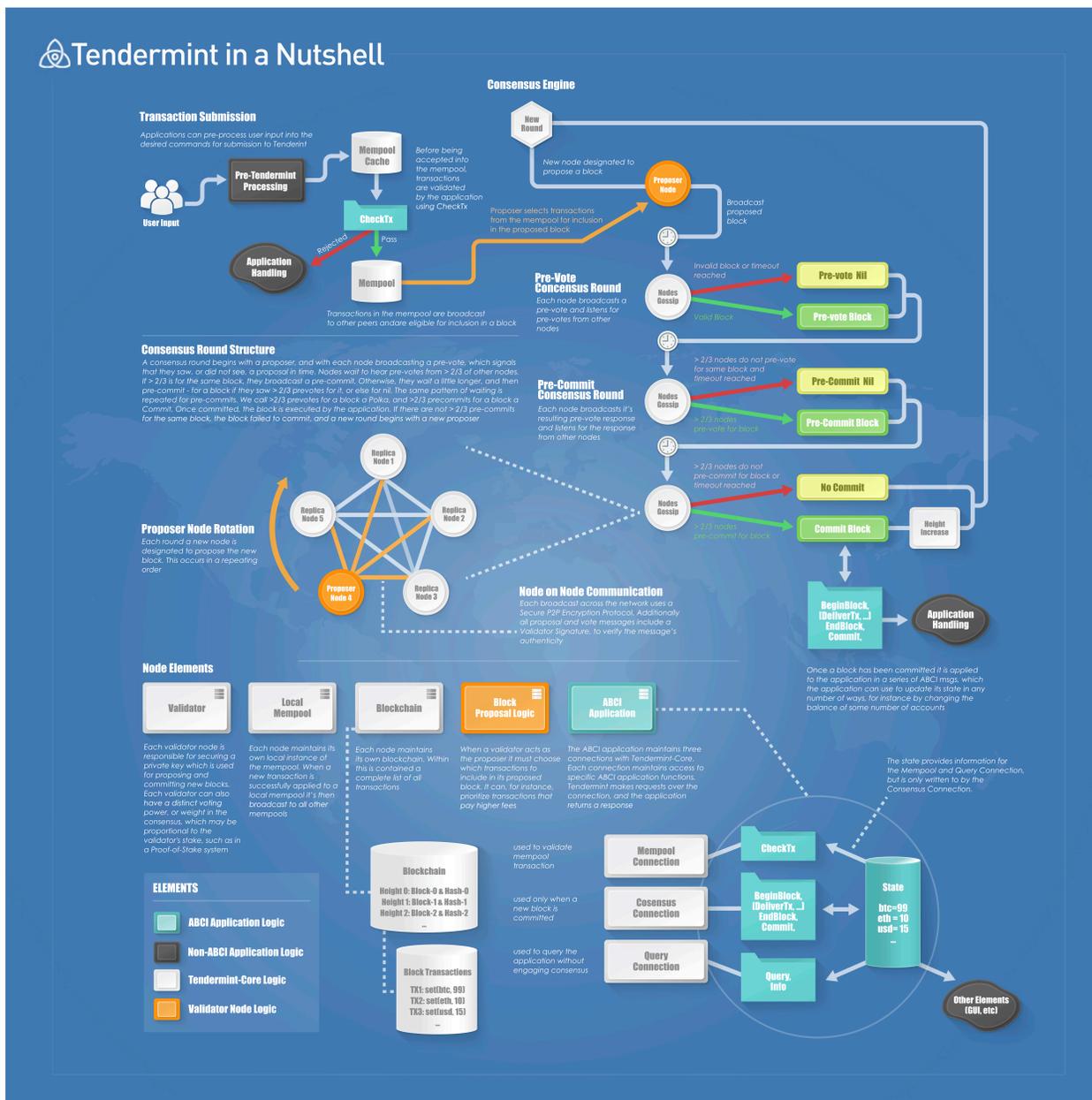
reference

Ethereum HD Wallet (虚拟货币钱包) -BIP32、BIP39、BIP44理解开发 HD 钱包涉及的 BIP32、BIP44、BIP39 基于 BIP-32 和 BIP-39 规范生成 HD 钱包 (分层确定性钱包) 分层确定性钱包 HD Wallet 剖析: 设计和实现

2.6 Cosmos

2.6.1 Tendermint

架构流程



client 通过 RPC 接口 broadcast_tx_commit 提交交易；

mempool 调用 ABCI 接口 CheckTx 用于校验交易的有效性，比如交易序号、发送者余额等，同时订阅交易执行后的事件并等待监听。

共识从 mempool 中获取交易开始共识排序，

协议中有两个角色：

- 验证人：协议中的角色或者节点，不同的验证者在投票过程中具备不同的权力（vote power）。
- 提议人：由验证人轮流产生。

验证人轮流对交易的区块提议并对提议的区块投票。区块被提交到链上，且每个区块就是一个区块高度。但区块也有可能提交失败，这种情况下协议将选择下一个验证人在相同高度上提议一个新块，重新开始投票。

从图中可以看到，成功提交一个区块，必须经过两阶段的投票，称为 pre-vote 和 pre-commit。当超过 2/3 的验证人在同一轮提议中对同一个块进行了 pre-commit 投票，那么这个区块才会被提交。

由于离线或者网络延迟等原因，可能造成提议人提议区块失败。这种情况在 Tendermint 中也是允许的，因为验证人会在进入下一轮提议之前等待一定时间，用于接收提议人提议的区块。

假设少于三分之一的验证人是拜占庭节点，Tendermint 能够保证验证人永远不会在同一高度重复提交区块而造成冲突。为了做到这一点，Tendermint 引入了锁定机制，一旦验证人预投票了一个区块，那么该验证人就会被锁定在这个区块。然后：

1. 该验证人必须在预提交的区块进行预投票。
2. 当前一轮预提议和预投票没成功提交区块时，该验证人就会被解锁，然后进行对新块的下一轮预提交。

Commit 阶段先执行 BeginBlock, DeliverTxAsync, EndBlock, 然后再执行 Commit

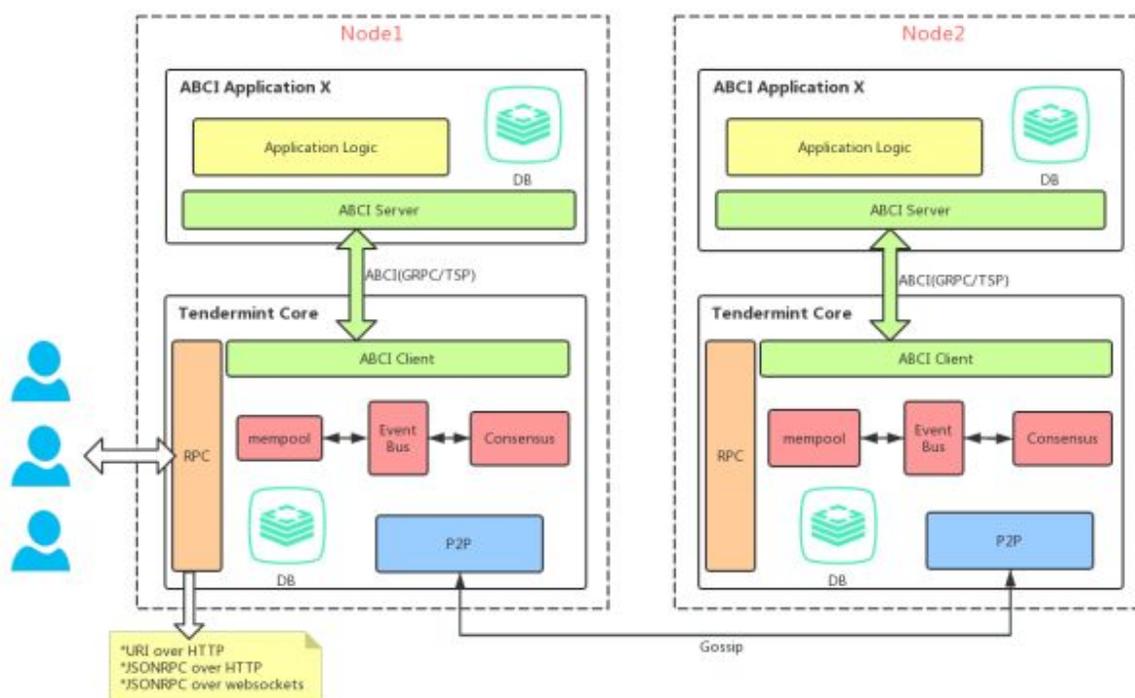
PBFT 相同点：

1. 同属 BFT 体系。
2. 抗 1/3 拜占庭节点攻击。
3. 三阶段提交，第一阶段广播交易（区块），后两阶段广播签名（确认）。
4. 两者都需要达到法定人数才能提交块。

与 PBFT 不同点：

1. Tendermint 与 PBFT 的区别主要是在超过 1/3 节点为拜占庭节点的情况下。当拜占庭节点数量在验证者数量的 1/3 和 2/3 之间时，PBFT 算法无法提供保证，使得攻击者可以将任意结果返回给客户端。而 Tendermint 共识模型认为必须超过 2/3 数量的 precommit 确认才能提交块。举个例子，如果 1/2 的验证者是拜占庭节点，Tendermint 中这些拜占庭节点能够阻止区块的提交，但他们自己也无法提交恶意块。而在 PBFT 中拜占庭节点却是可以提交块给客户端。
2. 另一个不同点在于拜占庭节点概念不同，PBFT 指的是节点数，而 Tendermint 代表的是节点的权益数，也就是投票权力。
3. 最后一点，PBFT 需要预设一组固定的验证人，而 Tendermint 是通过要求超过 2/3 法定人数的验证人员批准会员变更，从而支持验证人的动态变化。

关于 P2P 网络



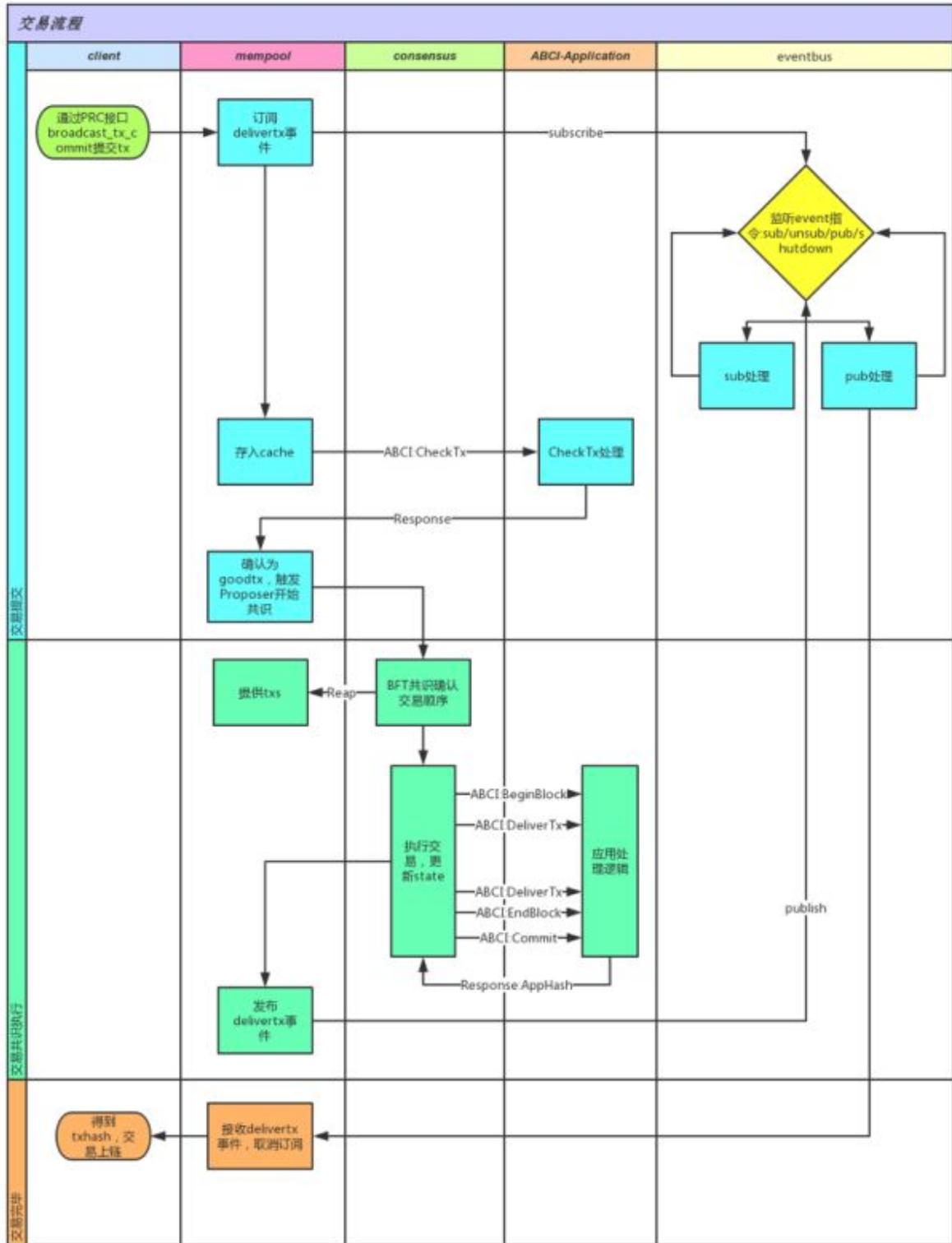
- P2P: 采用 gossip 算法，默认端口是 46656。
- RPC: 区块链对外接口，默认端口是 46657。支持三种访问方式：URI over HTTP、JSONRPC over HTTP、JSONRPC over websockets。

Tendermint 的 P2P 网络协议借鉴了比特币的对等发现协议，更准确地说，Tendermint 是采用了 BTCD 的 P2P 地址簿 (Address Book) 机制。当连接建立后，新节点将自身的 Address 信息 (包含 IP、Port、ID 等) 发送给相邻节点，相邻节点接收到信息后加入到自己的地址簿，再将此条 Address 信息，转播给它的相邻节点。

此外为了保证节点之间数据传输的安全性，Tendermint 采用了基于 Station-to-Station 协议的认证加密方案，此协议是一种密钥协商方案，基于经典的 DH 算法，并提供相互密钥和实体认证。大致的流程如下：

1. 每一个节点都必须生成一对 ED25519 密钥对作为自己的 ID。
2. 当两个节点建立起 TCP 连接时，两者都会生成一个临时的 ED25519 密钥对，并把临时公钥发给对方。
3. 两个节点分别将自己的私钥和对方的临时公钥相乘，得到共享密钥。这个共享密钥对称加密密钥。
4. 将两个临时公钥以一定规则进行排序，并将两个临时公钥拼接起来后使用 Ripemd160 进行哈希处理，后面填充 4 个 0，这样可以得到一个 24 字节的随机数。
5. 得到的随机数作为加密种子，但为了保证相同的随机数不会被相同的私钥使用两次，我们将随机数最后一个 bit 置为 1，这样就得到了两个随机数，同时约定排序更高的公钥使用反转过的随机数来加密自己的消息，而另外一个用于解密对方节点的消息。
6. 使用排序的临时公钥拼接起来，并进行 SHA256 哈希，得到一个挑战码。
7. 每个节点都使用自己的私钥对挑战码进行签名，并将自己的公钥和签名发给其它节点校验。
8. 校验通过之后，双方的认证就验证成功了。后续的通信就使用共享密钥和随机数进行加密，保护数据的安全。

Tendermint workflow



上图简单描述了 Tendermint 的工作流。大致为：

1. client 通过 RPC 接口 `broadcast_tx_commit` 提交交易；
2. mempool 调用 ABCI 接口 `CheckTx` 用于校验交易的有效性，比如交易序号、发送者余额等，同时订阅交易执行后的事件并等待监听。
3. 共识从 mempool 中获取交易开始共识排序，打包区块，确定之后依次调用 ABCI 相关接口更新当前的世界状态，并触发事件。
4. 最终将交易信息返回 client。

ABCI 接口

ABCI 接口可以分为三类：信息查询、交易校验以及共识相关处理。而 Tendermint Core 作为 ABCI Client 在启动时，会与 ABCI Server 建立三个连接，分别用于这三类接口消息的处理。

在 Tendermint Core 与 Application 交互的所有消息类型中，有 3 种主要的消息类型：

1. `CheckTx` 消息用于验证交易。Tendermint Core 中的 mempool 通过此消息校验交易的合法性，通过之后才会将交易广播给其它节点。
2. `DeliverTx` 消息是应用的主要工作流程，通过此消息真正执行交易，包括验证交易、更新应用程序的状态。
3. `Commit` 消息通知应用程序计算当前的世界状态，并存在下一区块头中。

Reference

- 深度剖析区块链跨链技术 Cosmos(上篇)

2.7 Rust

2.7.1 基础语法

实例声明

变量

- 不可变变量 `let`：代码可读性强；每次需要创建实例，开销较大；
- 可变变量 `let mut`：复用变量/地址空间；

常量

- 声明方式 `const`：不允许用 `mut`；必须注明值的类型；常量只能被设置为常量表达式；

```
const MAX_POINTS: u32 = 100_000;
```

```
let mut guess = String::new();  
// 创建了一个可变变量，当前它绑定到一个新的 String 空实例上
```

::new 那一行的 :: 语法表明 new 是 String 类型的一个 **关联函数** (*associated function*)。关联函数是针对类型实现的，在这个例子中是 String，而不是 String 的某个特定实例。一些语言中把它称为 **静态方法** (*static method*)。

new 函数创建了一个新的空字符串 (String 类的实例)，很多类型上有 new 函数

数据类型

标量 (scalar)

复合 (compound)

元组 (tuple)

- 元组长度固定：一旦声明，其长度不会增大或缩小。

```
fn main() {
    //创建了一个包含多个类型元素的元组，并绑定到 tup 变量上。
    let tup = (500, 6.4, 1);

    // 接着使用了 let 和一个模式将 tup 分成了三个不同的变量，x、y 和 z。这叫做
    ↪ 解构 (destructuring)
    let (x, y, z) = tup;

    println!("The value of y is: {}", y);

    let tuple: (i32, f64, u8) = (500, 6.4, 1);

    // 也可以使用点号 (.) 后跟值的索引来直接访问它们
    println!("The value of x is: {}", tuple.0);
}
```

数组 (array)

- 数组固定长度：一旦声明，它们的长度不能增长或缩小。

```
fn main() {
    // [元素类型; 元素数量] i32: 标量类型; 5个元素
    let a: [i32; 5] = [1, 2, 3, 4, 5];

    // 数组是一整块分配在栈上的内存。可以使用索引来访问数组的元素
    let first = a[0];
    println!("first = {} ", first);
}
```

函数调用

- 每个参数需要注明类型，返回值用 `->` 表示

```
fn main() {
    let x = plus_one(5);
    println!("The value of x is: {}", x);
}

fn plus_one(x: i32) -> i32 {
    // return x + 1; 显式的 x + 1
    x + 1
    // x + 1; 会报错，相当于执行 return 空
}
```

循环控制

- Rust 有三种循环: `loop`、`while` 和 `for`

```
fn main() {
    let mut a: [i32; 4] = [0;4];
    let mut i: i32 = 1; let mut count = 0;

    // result 表达式: 会返回 loop 过程中的返回值
    let result = loop {
        if i > 10 {
            break;
        }
        i += i; a[count] = i; count += 1;
    };
    result; // 执行刚刚的代码块

    // while 循环
    while i < 100 {
        i += i;
    }
    println!("i is {}", i);

    // for 循环遍历:
    for element in a.iter() {
        println!("iter value is: {}", element);
    }

    // 循环遍历 rev: 反转 range
    for number in (1..4).rev() {
        println!("rev {}", number);
    }
}
```

2.7.2 变量与数据

所有权

所有权规则

- 1、rust 中的每个值都有一个被称为**所有者** (owner) 的变量
- 2、值在任意时刻有且只有一个所有者
- 3、当**所有者** (变量) 离开作用域, 这个值会被丢弃

内存与分配

- String 类型, 为了支持一个可变, 可增长的文本片段, 需要在堆上分配一块在编译时未知大小的内存来存放内容。这意味着:
 - 必须在运行时向操作系统请求内存。
 - 需要一个当我们处理完 String 时将内存返回给操作系统的方法。
- 当变量离开作用域, Rust 为我们调用一个特殊的函数。这个函数叫做 drop, 在这里 String 的作者可以放置释放内存的代码。Rust 在结尾的 } 处自动调用 drop。

变量与数据交互的方式 (一): 移动

- Rust 永远也不会自动创建数据的“深拷贝”。因此, 任何**自动**的复制可以被认为对运行时性能影响较小。
- String 会申请内存资源, 复制时涉及到内存 (值) 的所有权转移, 参见所有权规则 1

```
fn move_test(){
    // 像整型这样的在编译时已知大小的类型被整个存储在栈上
    let stack_str = "stack";
    let copy_stack = stack_str;
    println!("{}, hello!, copy {}", stack_str, copy_stack);

    // String 类型被分配到堆上, 所以能够存储在编译时未知大小的文本。
    let heap_str = String::from("heap");

    // String实例 heap_str 中指向堆存储的指针作废, move_pointer_
    →复制或者说移动刚刚的指针
    // 拷贝指针、长度和容量而不拷贝数据
    let move_pointer = heap_str;
    println!("{0}, hello!; move stack pointer, {0} data no change", move_pointer);
    →
}
```

变量与数据交互的方式 (二): 克隆

- 会重新生成一个

```
fn clone_test() {
    let heap_str = String::from("hello");

    // clone 会深度复制 String 中堆上的数据，而不仅仅是栈上的数据
    let heap_copy = heap_str.clone();
    println!("heap_str = {}, copy = {}", heap_str, heap_copy);
}
```

只在栈上的数据: 拷贝

- 任何简单标量值的组合可以是 Copy 的，不需要分配内存或某种形式资源的类型是 Copy 的
 - 所有整数类型，比如 u32。
 - 布尔类型，bool，它的值是 true 和 false。
 - 所有浮点数类型，比如 f64。
 - 字符类型，char。
 - 元组，当且仅当其包含的类型也都是 Copy 的时候。比如，(i32, i32) 是 Copy 的，但 (i32, String) 就不是。

所有权转移

- 将值赋给另一个变量时移动它。当持有堆中数据值的变量离开作用域时，其值将通过 drop 被清理掉，除非数据被移动为另一个变量所有。
- 函数参数可转移所有权: (copy) / (move)。目前来看取决于传递的类型是标量值还是占据内存资源的值
- 返回值可转移所有权

```
// gives_ownership 将返回值move给调用它的函数
fn gives_ownership() -> String {
    let some_string = String::from("hello"); // some_string 进入作用域。
    some_string                               // 返回 some_string_
    ↪ 并移出给调用的函数
}

// takes_and_gives_back 将传入字符串并返回该值
fn takes_and_gives_back(a_string: String) -> String { // a_string 进入作用域
    "change".to_owned() // 返回并move给调用的函数
} // a_string 移出作用域并被丢弃。

fn makes_copy(some_integer: i32) { // some_integer 进入作用域
    println!("{}", some_integer);
} // 这里，some_integer 移出作用域。不会有特殊操作

fn ownership_test() {
    let s1 = gives_ownership(); // gives_ownership 将返回值move 给s1

    let s2 = String::from("hello"); // s2 进入作用域
```

(下页继续)

(续上页)

```

// s2 被move到函数中,它也将返回值移给 s3。
// 根据规则二: s2已经没有所有权,不能使用
let s3 = takes_and_gives_back(s2);

println!("s1 = {}, s3 = {}", s1, s3);

let x = 5; // x 进入作用域

makes_copy(x); // x 应该移动到函数里,但 i32 是 Copy
→的,所以在后面可继续使用 x
} // 根据规则三: 这里, s3 移出作用域并被丢弃。s2 也移出作用域,但已被移走,
// 所以什么也不会发生。s1 移出作用域并被丢弃

```

引用

规则

- 在任意给定时间,要么只能有一个可变引用,要么只能有多个不可变引用。
- 引用必须总是有效的。

实例

- 可变引用/

```

fn reference_test() {
    let s1 = String::from("hello");

    let len = calculate_length(&s1);
    println!("The length of '{}' is {}.", s1, len);

    // 在特定作用域中的特定数据只能有一个可变引用。可以通过代码块来变通约束条件
    {
        let r1 = &mut s1;

    } // r1 在这里离开了作用域,所以我们完全可以创建一个新的引用

    let r2 = &mut s1;
}

fn calculate_length(s: &String) -> usize { // s 是对 String 的引用;
→获取引用作为函数参数称为 借用 (borrowing)
    s.len()
} // 这里, s 离开了作用域。但因为它并不拥有引用值的所有权,所以什么也不会发生

```


3.1 链表

3.1.1 实现链表的排序

使用 $O(n \log n)$ 的时间复杂度对链表进行排序

```
/* 链表排序
 * 1、每一个节点是一个有序链表，当前节点与排好序的链表进行合并
 * 2、合并两个有序链表，只需要两个当前指针进行比大小即可
 * 3、 $O(n \log n)$  需要进行分治，合并是 $O(n)$ ，分治才能 $O(\log n)$ 
 * 4、两个快慢指针，快指针跳两下，慢指针跳一下，达到对半分治的效果
 */
func sortList(head *ListNode) *ListNode {
    if head == nil || head.Next == nil { return head }

    slow, fast := head, head
    var prev *ListNode

    for fast != nil && fast.Next != nil { // fast 比 slow 快一倍，即 slow 是半点
        prev = slow
        slow = slow.Next
        fast = fast.Next.Next
    }
    prev.Next = nil

    l1 := sortList(head)
    l2 := sortList(slow)
    return mergeTwoLists(l1, l2)
}

func mergeTwoLists(l1 *ListNode, l2 *ListNode) *ListNode {
    var tmp = &ListNode{}
}
```

(下页继续)

```

var head = tmp
for l1!=nil && l2!=nil{
    if(l1.Val < l2.Val){
        tmp.Next = l1
        l1 = l1.Next
        tmp = tmp.Next
    }else{
        tmp.Next = l2
        l2 = l2.Next
        tmp = tmp.Next
    }
}
if(l1 != nil){
    tmp.Next = l1
}else if(l2 != nil){
    tmp.Next = l2
}
return head.Next
}

```

3.1.2 实现高效的单向链表逆序输出

```

// 如何实现一个高效的单向链表逆序输出？
import "fmt"

type node struct {
    data int
    next *node
}

func reverse(head *node) {
    if head == nil || head.next == nil {
        return
    }
    var prev *node
    var next *node
    pcur := head
    for pcur != nil {
        if pcur.next == nil { //此时为最后一个节点
            pcur.next = prev
            break
        }
        next = pcur.next //取下一个节点
        pcur.next = prev // 将当前节点挂在下一个节点 (pcur) 的后面
        prev = pcur // 当前节点
        pcur = next //将下一个节点指给 pcur
    }
    tmp := pcur
    for tmp != nil {
        fmt.Print(tmp.data, " ")
        tmp = tmp.next
    }
}

func main() {

```

(续上页)

```

var head *node = &node{data: 0, next: nil} //指针指向资源区
var varHead *node = head
fmt.Println("before:")
for i := 1; i < 6; i++ { // 构建链表
    var tmp node = node{data: i, next: nil}
    fmt.Print(tmp.data, " ")
    varHead.next = &tmp
    varHead = varHead.next
}
fmt.Println()
fmt.Println("after:")
reverse(head)
}

```

3.2 Tree

3.2.1 验证一棵树是否为二叉搜索树

```

func isValidBST(root *TreeNode) bool {
    if(root == nil) {return true}
    if(root.Left !=nil && (getRight(root.Left) >= root.Val)) {return false}
    if(root.Right !=nil && (getLeft(root.Right) <= root.Val)) {return false}
    return (isValidBST(root.Left) && isValidBST(root.Right))
}
func getRight(root *TreeNode) int {
    if(root.Right == nil){return root.Val }
    return getRight(root.Right)
}
func getLeft(root *TreeNode) int {
    if(root.Left == nil){return root.Val }
    return getLeft(root.Left)
}
}

```

3.2.2 平衡二叉树

//二叉树，其中每个节点的左和右子树的高度相差不超过1。

```

func isBalanced(root *TreeNode) bool {
    if(root == nil){ return true }
    if(isBalanced(root.Left) && isBalanced(root.Right)){
        return WithBranch(getHeight(root.Left) - getHeight(root.Right)) <=1
    }
    return false
}
func getHeight(root *TreeNode) int{
    if(root == nil) { return 0 }
    l := getHeight(root.Left) +1
    r := getHeight(root.Right) +1
    if(l <= r){return r}
    return l
}
}

```

(下页继续)

(续上页)

```
func WithBranch(n int) int {  
    if n < 0 {  
        return -n  
    }  
    return n  
}
```

4.1 Mysql

4.1.1 索引

索引用于：

- 快速找出匹配一个 WHERE 子句的行；
- 当执行联结时，从其他表检索行；
- 对特定的索引列找出 MAX() 或 MIN() 值；
- 如果排序或分组在一个可用键的最左面前缀上进行（例如，ORDER BY key_part_1,key_part_2），排序或分组一个表。
- 如果所有键值部分跟随 DESC，键以倒序被读取。
- 在一些情况中，一个查询能被优化来检索值，不用咨询数据文件。

如果对某些表的所有使用的列是数字型的并且构成某些键的最左面前缀，为了更快，值可以从索引树被检索出来。

索引原理

索引要求做到减少 IO，加速查询

- **在表中有大量数据的前提下，创建索引速度会很慢**

在没有数据的情况下先建索引或者说目录快，还是已经存在好多的数据了，然后再去建索引，哪个快，肯定是没有数据的时候快，因为如果已经有了很多数据了，你再去根据这些数据建索引，是不是要将数据全部遍历一遍，然后根据数据建立索引。

- **在索引创建完毕后，对表的查询性能会幅度提升，但是写性能会降低**

索引建立好之后再添加数据快，还是没有索引的时候添加数据快，索引是用来干什么的，是用来加速查询的，那对你写入数据会有什么影响，肯定是慢一些了，因为你但凡加入一些新的数据，都需要把索引或者说书的目录重新做一个，**所以索引虽然会加快查询，但是会降低写入的效率。**

磁盘 IO 与预读

磁盘读取数据靠的是机械运动，每次读取数据花费的时间可以分为**寻道时间、旋转延迟、传输时间**三个部分。

- 寻道时间指的是磁臂移动到指定磁道所需要的时间，主流磁盘一般在 5ms 以下；
- 旋转延迟就是我们经常听说的磁盘转速，比如一个磁盘 7200 转/min，表示每分钟能转 7200 次，也就是说 1 秒钟能转 120 次，旋转延迟就是 $1/120/2 = 4.17\text{ms}$ ，也就是半圈的时间（这里有两个时间：平均寻道时间，受限于目前的物理水平，大概是 5ms 的时间，找到磁道了，还需要找到你数据存在的那个点，寻点时间，这寻点时间的一个平均值就是半圈的时间，这个半圈时间叫做平均延迟时间，那么平均延迟时间加上平均寻道时间就是你找到一个数据所消耗的平均时间，大概 9ms，其实机械硬盘慢主要是慢在这两个时间上了，当找到数据然后把数据拷贝到内存的时间是非常短暂的，和光速差不多了）；
- 传输时间指的是从磁盘读出或将数据写入磁盘的时间，一般在零点几毫秒，相对于前两个时间可以忽略不计

那么访问一次磁盘的时间，即一次磁盘 IO 的时间约等于 $5+4.17 = 9\text{ms}$ 左右。

计算机硬件延迟的对比图：

L1 cache reference	0.5 ns
Branch mispredict	5 ns
L2 cache reference	7 ns
Mutex lock/unlock	25 ns
Main memory reference	100 ns
Compress 1K bytes with Zippy	3,000 ns
Send 2K bytes over 1 Gbps network	20,000 ns
Read 1 MB sequentially from memory	250,000 ns
Round trip within same datacenter	500,000 ns
Disk seek	10,000,000 ns
Read 1 MB sequentially from disk	20,000,000 ns
Send packet CA->Netherlands->CA	150,000,000 ns

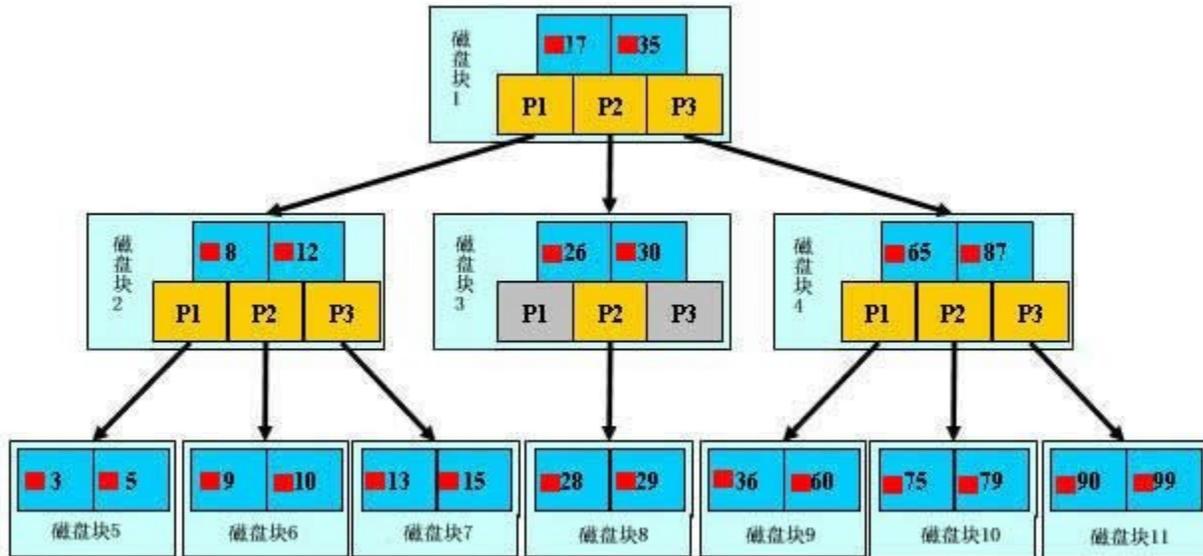
考虑到磁盘 IO 是非常高昂的操作，计算机操作系统做了一些优化，当一次 IO 时，不光把当前磁盘地址的数据，而是把相邻的数据也都读取到内存缓冲区内，因为局部预读性原理告诉我们，当计算机访问一个地址的数据的时候，与其相邻的数据也会很快被访问到。每一次 IO 读取的数据我们称之为`一页 (page)`。具体一页有多大数据跟操作系统有关，一般为 4k 或 8k，也就是我们读取一页内的数据时候，实际上才发生了一次 IO，这个理论对于索引的数据结构设计非常有帮助。

索引的数据结构

红黑树等数据结构也可以用来实现索引，但是文件系统及数据库系统普遍采用 B-/+Tree 作为索引结构。

一般来说，索引本身也很大，不可能全部存储在内存中，因此索引往往以索引文件的形式存储在磁盘中。这样的话，索引查找过程中就要产生磁盘 I/O 消耗，索引的结构组织要尽量减少查找过程中磁盘 I/O 的存取次数，最好是常数数量级。

详解 b+ 树



如上图，是一颗 b+ 树，关于 b+ 树的定义可以参见 B+ 树，这里只说一些重点，浅蓝色的块我们称之为一个磁盘块，可以看到每个磁盘块包含几个数据项（深蓝色所示）和指针（黄色所示），如磁盘块 1 包含数据项 17 和 35，包含指针 P1、P2、P3，P1 表示小于 17 的磁盘块，P2 表示在 17 和 35 之间的磁盘块，P3 表示大于 35 的磁盘块。真实的数据存在于叶子节点即 3、5、9、10、13、15、28、29、36、60、75、79、90、99。非叶子节点只不存储真实的数据，只存储指引搜索方向的数据项，如 17、35 并不真实存在于数据表中。

b+ 树的查找过程

如图所示，如果要查找数据项 29，那么首先会把磁盘块 1 由磁盘加载到内存，此时发生一次 IO，在内存中用二分查找确定 29 在 17 和 35 之间，锁定磁盘块 1 的 P2 指针，内存时间因为非常短（相比磁盘的 IO）可以忽略不计，通过磁盘块 1 的 P2 指针的磁盘地址把磁盘块 3 由磁盘加载到内存，发生第二次 IO，29 在 26 和 30 之间，锁定磁盘块 3 的 P2 指针，通过指针加载磁盘块 8 到内存，发生第三次 IO，同时内存中做二分查找找到 29，结束查询，总计三次 IO。真实的情况是，3 层的 b+ 树可以表示上百万的数据，如果上百万的数据查找只需要三次 IO，性能提高将是巨大的，如果没有索引，每个数据项都要发生一次 IO，那么总共需要百万次的 IO，显然成本非常非常高。

b+ 树性质

1**. 索引字段要尽量的小 **: 通过上面的分析, 我们知道 IO 次数取决于 b+ 数的高度 h 或者说层级, 这个高度或者层级就是你每次查询数据的 IO 次数, 假设当前数据表的数据为 N , 每个磁盘块的数据项的数量是 m , 则有 $h = \lceil \log_{m+1} N \rceil$, 当数据量 N 一定的情况下, m 越大, h 越小; 而 $m = \text{磁盘块的大小} / \text{数据项的大小}$, 磁盘块的大小也就是一个数据页的大小, 是固定的, 如果数据项占的空间越小, 数据项的数量越多, 树的高度越低。这就是为什么每个数据项, 即索引字段要尽量的小, 比如 `int` 占 4 字节, 要比 `bigint8` 字节少一半。这也是为什么 b+ 树要求把真实的数据放到叶子节点而不是内层节点, 一旦放到内层节点, 磁盘块的数据项会大幅度下降, 导致树增高。当数据项等于 1 时将会退化线性表。

2. 索引的最左匹配特性: 简单来说就是你的数据来了以后, 从数据块的左边开始匹配, 在匹配右边的。

我们继续学下面的内容。当 b+ 树的数据项是复合的数据结构, 比如 `(name,age,sex)` 的时候, b+ 树是按照从左到右的顺序来建立搜索树的, 比如当 `(张三,20,F)` 这样的数据来检索的时候, b+ 树会优先比较 `name` 来确定下一步的所搜方向, 如果 `name` 相同再依次比较 `age` 和 `sex`, 最后得到检索的数据; 但当 `(20,F)` 这样的没有 `name` 的数据来的时候, b+ 树就不知道下一步该查哪个节点, 因为建立搜索树的时候 `name` 就是第一个比较因子, 必须先根据 `name` 来搜索才能知道下一步去哪里查询。比如当 `(张三,F)` 这样的数据来检索时, b+ 树可以用 `name` 来指定搜索方向, 但下一个字段 `age` 的缺失, 所以只能把名字等于张三的数据都找到, 然后再匹配性别是 F 的数据了, 这个是非常重要的性质, 即索引的最左匹配特性。

建索引的几大原则

1. 最左前缀匹配原则, 非常重要的原则, mysql 会一直向右匹配直到遇到范围查询 (`>`、`<`、`between`、`like`) 就停止匹配, 比如 `a = 1 and b = 2 and c > 3 and d = 4` 如果建立 `(a,b,c,d)` 顺序的索引, `d` 是用不到索引的, 如果建立 `(a,b,d,c)` 的索引则都可以用到, `a,b,d` 的顺序可以任意调整。

2. `=` 和 `in` 可以乱序, 比如 `a = 1 and b = 2 and c = 3` 建立 `(a,b,c)` 索引可以任意顺序, mysql 的查询优化器会帮你优化成索引可以识别的形式。

3. 尽量选择区分度高的列作为索引, 区分度的公式是 $\text{count}(\text{distinct col}) / \text{count}(*)$, 表示字段不重复的比例, 比例越大我们扫描的记录数越少, 唯一键的区分度是 1, 而一些状态、性别字段可能在大数据面前区分度就是 0, 那可能有人会问, 这个比例有什么经验值吗? 使用场景不同, 这个值也很难确定, 一般需要 `join` 的字段我们都要求是 0.1 以上, 即平均 1 条扫描 10 条记录。

4. 索引列不能参与计算, 保持列“干净”, 比如 `from_unixtime(create_time) = '2014-05-29'` 就不能使用到索引, 原因很简单, b+ 树中存的都是数据表中的字段值, 但进行检索时, 需要把所有元素都应用函数才能比较, 显然成本太大。所以语句应该写成 `create_time = unix_timestamp('2014-05-29')`。

5. 尽可能的扩展索引, 不要新建索引。比如表中已经有 `a` 的索引, 现在要加 `(a,b)` 的索引, 那么只需要修改原来的索引即可。

MySQL 索引管理

功能

```
#1. 索引的功能就是加速查找
#2. mysql中的primary↵
↵key, unique, 联合唯一也都是索引, 这些索引除了加速查找以外, 还有约束的功能
```

MySQL 常用的索引

普通索引 INDEX: 加速查找

唯一索引:

- 主键索引 PRIMARY KEY: 加速查找+约束 (不为空、不能重复)
- 唯一索引 UNIQUE: 加速查找+约束 (不能重复)

联合索引:

- PRIMARY KEY(id,name): 联合主键索引
- UNIQUE(id,name): 联合唯一索引
- INDEX(id,name): 联合普通索引

参考如下:

MySQL 索引原理及慢查询优化

MySQL 索引

4.1.2 gorm

创建连接

```
db, err := gorm.Open("mysql", "user:password@/dbname?charset=utf8&parseTime=True&
↳ loc=Local")
```

数据模型定义

表名, 列名如何对应结构体

在 Gorm 中, 表名是结构体名的复数形式, 列名是字段名的蛇形小写。

即, 如果有一个 Order 表, 那么如果你定义的结构体名为: User, gorm 会默认表名为 orders 而不是 order。

```
CREATE TABLE `orders` (
  `order_id` varchar(3) NOT NULL,
  `status` bigint(20) DEFAULT NULL,
  PRIMARY KEY (`order_id`),
  KEY `idx_orders_status` (`status`),
) ENGINE=InnoDB DEFAULT CHARSET=utf8;
```

那么对应的结构体定义如下:

```
type Order struct {
  OrderId      string `gorm:"PRIMARY_KEY;type:varchar(3)" json:"orderId"`
  Status       int64  `gorm:"index;" json:"status"` //索引
}
```

如何全局禁用表名复数呢?

可以在创建数据库连接的时候设置如下参数:

```
// 全局禁用表名复数
db.SingularTable(true) // 如果设置为true, `User` 的默认表名为 `user`,
↳使用 `TableName` 设置的表名不受影响
```

这样的话, 表名默认即为结构体的首字母小写形式。

生成表

当数据库需要频繁更新结构时, 代码与数据库难以保持一致是烦人的问题. 而 go 的 gorm 库, 有 Auto Migration 功能, 可以根据 go 里的 struct tag 自动更新数据库结构, 非常方便。

```
// 自动生成表
db.AutoMigrate(&Order{})
```

查询

```
query := db.Model(Order{}).Where("order_id = ? AND status = ?", order_id, status)
var count int
err := db.Model(&Order{}).Where(&Order{order_id: id, status: status}).Count(&count).
↳Error
```

先用 db.Model() 选择一个表, 再用 db.Where() 构造查询条件, 后面可以使用 db.Count() 计算数量, 如果要获取对象, 可以使用 db.Find(&Likes) 或者只需要查一条记录 db.First(&Order)

事务

要在事务中执行一组操作, 一般流程如下。

```
func CreateAnimals(db *gorm.DB) err {
    tx := db.Begin()
    // 注意, 一旦你在一个事务中, 使用tx作为数据库句柄

    if err := tx.Create(&Animal{Name: "Giraffe"}).Error; err != nil {
        tx.Rollback()
        return err
    }
    // 在事务中做一些数据库操作 (从这一点使用 'tx', 而不是 'db')
    if err := tx.Create(&Animal{Name: "Lion"}).Error; err != nil {
        tx.Rollback()
        return err
    }

    tx.Commit()
    return nil
}
```

日志

Gorm 有内置的日志记录器支持，默认情况下，它会打印发生的错误。

```
// 启用 Logger，显示详细日志
db.LogMode(true)

// 禁用日志记录器，不显示任何日志
db.LogMode(false)

// 调试单个操作，显示此操作的详细日志
db.Debug().Where("name = ?", "xiaoming").First(&User{})
```

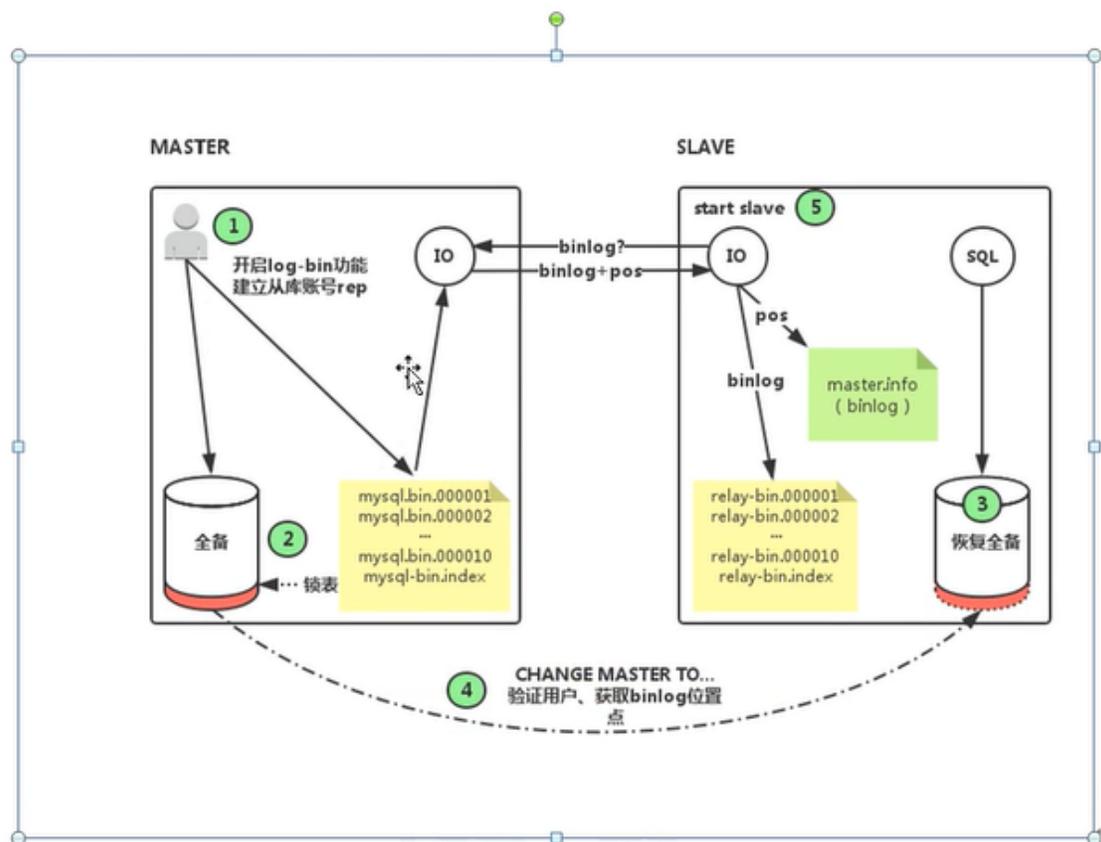
参考如下：

Go orm 框架 gorm 学习

4.1.3 使用 Docker 搭建 MySQL 主从复制

主从复制原理

· 复制的工作过程



该过程的第一部分就是 master 记录二进制日志。在每个事务更新数据完成之前，master 在二日志记录这些改变。MySQL 将事务串行的写入二进制日志，即使事务中的语句都是交叉执行的。在事件写入二进制日志完成后，master 通知存储引擎提交事务。

复制原理 mysql 的复制原理大致如下。

- 主库开启 bin-log 二进制日志功能，并建立 slave 账号，并授权从库连接主库，从库通过 change master 得到主库的相关同步信息。mysql 主库在事务提交时会把数据库变更作为事件 Events 记录在二进制文件 binlog 中；mysql 主库上的 sys_binlog 控制 binlog 日志刷新到磁盘。
- 从库然后连接主库进行验证，主库产生的新数据会导入到 bin-log 二进制文件中，同时主库会开启 IO 线程，从库也会开启 IO 线程以及 SQL 线程，从库中的 IO 线程与主库的 IO 线程连接一旦主库数据有所变更则从库将变更的数据复制到 relay-bin（中继日志）中。SQL 线程会读取 relay log（中继日志）文件中的日志，并解析成具体操作，至此整个同步过程完成。

主从复制操作

流程命令集

```
# 1、搭建部署 mysql
docker run --name mysql1 -p 33061:3306 -e MYSQL_ROOT_PASSWORD=123 -d mysql:5.7 --
↳character-set-server=utf8mb4 --collation-server=utf8mb4_unicode_ci
docker run --name mysql2 -p 33062:3306 -e MYSQL_ROOT_PASSWORD=123 -d mysql:5.7 --
↳character-set-server=utf8mb4 --collation-server=utf8mb4_unicode_ci

# 2、主机配置
## 2.1、进入mysql
mysql -u root -h 127.0.0.1 -P 33061 -p

## 2.2、主机配置 配置一个从机登录用户
GRANT REPLICATION SLAVE ON *.* TO 'user'@'%' IDENTIFIED BY '123';

## 2.3、进入docker容器
docker exec -it mysql1 bash

## 2.4、需要修改 MySQL 的配置，开启 binlog
vi /etc/mysql/mysql.conf.d/mysqld.cnf

# 3、从机配置
## 3.1、进入mysql
mysql -u root -h 127.0.0.1 -P 33061 -p

## 3.2、从配置 配置一个从机登录用户
change master to master_host='192.168.66.114',master_port=33061,master_user='user',
↳master_password='123',master_log_file='binlog.000001',master_log_pos=154;

## 3.3、slave 启动或停止
start slave;
stop slave;

## 3.4、查看slave配置正确与否
show slave status\G;

####Slave_IO_Running 和 Slave_SQL_Running ，这两个的值必须为 Yes。
```

配置相关文件

主机 MySQL 配置文件 mysqld.cnf

- binlog 的开启，需要修改 MySQL 的配置，因此，我们需要进入到容器内部去执行。
docker exec -it mysql1 /bin/bash
- MySQL 配置文件/etc/mysql/mysql.conf.d/mysqld.cnf 追加以下内容

```
# 第一行表示配置 binlog 的位置，理论上 binlog 可以放在任意位置，但是该位置，MySQL
↳一定要有操作权限。
log-bin=/var/lib/mysql/binlog
```

(下页继续)

(续上页)

```
#server-id 表示集群中，每个实例的一标识符
server-id=1

# 不同步哪些数据库
binlog-ignore-db = mysql
binlog-ignore-db = test
binlog-ignore-db = information_schema

# 只同步哪些数据库，除此之外，其他不同步
binlog-do-db = game
```

```
mysql> show master status;
+-----+-----+-----+-----+-----+
| File           | Position | Binlog_Do_DB   | Binlog_Ignore_DB | Executed_Gtid_Set |
+-----+-----+-----+-----+-----+
| binlog.000001 |      154 | cmdb,db1,db2,db3 |                   |                   |
+-----+-----+-----+-----+-----+
1 row in set (0.00 sec)
```

File 和 Position 需要记着，这两个标记了二进制日志的起点位置，在从机的配置中将使用到这两个参数。

从机配置

- show slave status;
 - Slave_IO_State: 当前 I/O 线程的状态
 - Master_Log_File: 当前同步的主服务器的二进制文件
 - Read_Master_Log_Pos: 当前同步的主服务器的二进制文件的偏移量，单位为字节，如图中为已经同步了 12.9M(13630580/1024/1024) 的内容
 - Relay_Master_Log_File: 当前中继日志同步的二进制文件
 - Slave_IO_Running: 从服务器中 I/O 线程的运行状态，YES 为运行正常
 - Slave_SQL_Running: 从服务器中 SQL 线程的运行状态，YES 为运行正常
 - Exec_Master_Log_Pos: 表示同步完成的主服务器的二进制日志偏移量
 - Seconds_Behind_Master: 表示从服务器数据比主服务器落后的持续时长

参考链接

使用 Docker 搭建 MySQL 主从复制

4.2 Redis

4.2.1 redis 使用规范

1、核心业务数据和边缘业务数据分离

这两类数据分开为不同的缓存实例，避免边缘业务影响核心业务服务；

2、数据缓存按业务功能分离

不能将众多的业务服务依赖同一缓存实例，而应该按业务进行拆分（尤其是核心的高并发的数据应单独的缓存服务实例）。

3、缓存 key 命名规范

包括三个维度来命名（命名空间、数据分类名、变量），三者之间采用冒号“:”隔开

4、提升访问性能（pipeline）

多个操作命令发送到服务端，客户端不依赖每个操作命令返回的结果（即可批量化）时，可以使用 pipeline 机制减少 tcp 往返次数，提升 redis 访问性能；

5、实现优雅删除（unlink）

del 指令会同步释放对象内存，大部分情况下没有问题，当遇到比较大的 key 时（如上千万记录的 list、hash 等），则该指令会造成服务卡顿；因此 redis 从 4.0 版本开始发布了 unlink 指令来解决该问题。unlink 用于替代 del 指令，解决大 key 删除卡顿问题。在执行 unlink 指令时 redis 第一步将 key 从 redis 空间摘除此 key 的引用，第二步会在一个异步线程中释放此 key 的内存（因是异步的，不会对工作线程造成影响）。

6、安全遍历（scan）

线上应避免使用 keys 指令，keys 会遍历 redis 中所有的 key，当 key 的数据达到千万时，执行该指令会导致服务卡顿；可以采用 scan 指令带上 limit 参数来替代 keys。

redis 是单线程的。keys 指令会导致线程阻塞一段时间，线上服务会停顿，直到指令执行完毕，服务才能恢复。这个时候可以使用 scan 指令，scan 指令可以无阻塞的提取出指定模式的 key 列表，但是会有一定的重复概率，在客户端做一次去重就可以了，但是整体所花费的时间会比直接用 keys 指令长。

7、预防缓存击穿（缓存默认值）

当查询的数据在 Redis 数据库本身不存在时，每次请求一般都会查询数据库。恶意用户会利用此场景对数据库进行攻击从而影响服务。可以采用下发的方式进行预防：

从数据库查询不到数据时，在缓存中设置一个特殊的默认值并设置一个较短的过期时间。

常见攻击

- 缓存击穿当我们缓存 key 设置过期时间，恰巧在这一刻这个 key 在某一时刻被高并发的访问，把所有的请求都打到了 DB 中这就可能会导致 DB 挂了。

缓存击穿的话，设置热点数据永远不过期。或者加上互斥锁就能搞定了

或者使用分布锁，使一个线程去重新回原数据，其他请求返回空值，UI 友好提示。

- 缓存穿透使用大量不存在的 key 发布访问，由于缓存中找不到，故把所有请求全都打到 DB 中，使用 DB 压力巨增，导致 DB 挂掉。

缓存穿透我会在接口层增加校验，比如用户鉴权校验，参数做校验，不合法的参数直接代码 Return，比如：id 做基础校验，id <=0 的直接拦截等。

1. 使用 Boolm 过滤器，2. 使用 Hystrix 第三方组件进行限流、熔断，从而保护服务器。

- 缓存雪崩雪崩指的是多个 key 同时失效，当高并发的请求过来所有请求都打到数据库导致挂了

处理缓存雪崩简单，在批量往 Redis 存数据的时候，把每个 Key 的失效时间都加个随机值就好了，这样可以保证数据不会在同一时间大面积失效

```
setRedis (Key, value, time + Math.random() * 10000) ;
```

或者设置热点数据永远不过期，有更新操作就更新缓存就好了

8、缓存失效（加锁）

高并发时，如果缓存失效，会有大量的请求打到数据库，对数据库造成瞬间冲击，影响服务质量，因此在查询数据库加载到缓存时需要加锁，以预防数据库瞬间流量冲击。

4.2.2 redis 应用场景

一、分布式锁

二、数据存储

redis 可用作数据缓存，包括五种类型：

- **String:**

Strings 数据结构是简单的 key-value 类型，value 其实不仅是 String，也可以是数字。

常用命令: set,get,decr,incr,mget 等。

应用场景：String 是最常用的一种数据类型，普通的 key/ value 存储都可以归为此类。即可以完全实现目前 Memcached 的功能，并且效率更高。还可以享受 Redis 的定时持久化，操作日志及 Replication 等功能。除了提供与 Memcached 一样的 get、set、incr、decr 等操作外，Redis 还提供了下面一些操作：

获取字符串长度

往字符串 append 内容

设置和获取字符串的某一段内容

设置及获取字符串的某一位 (bit)

批量设置一系列字符串的内容

实现方式：String 在 redis 内部存储默认就是一个字符串，被 redisObject 所引用，当遇到 incr,decr 等操作时会转成数值型进行计算，此时 redisObject 的 encoding 字段为 int。

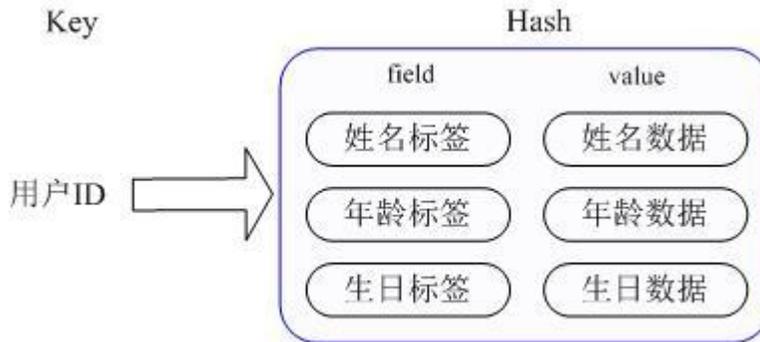
扩容过程:

当 len 长度小于 1M 时, 扩容时双倍扩容

当 len 长度大于 1M 时, 每次扩容 1M 容量

- **Hash:**

Redis 的 Hash 实际是内部存储的 Value 为一个 HashMap, 并提供了直接存取这个 Map 成员的接口, 如下图:



也就是说, Key 仍然是用户 ID, value 是一个 Map, 这个 Map 的 key 是成员的属性名, value 是属性值, 这样对数据的修改和存取都可以直接通过其内部 Map 的 Key(Redis 里称内部 Map 的 key 为 field), 也就是通过 key(用户 ID) + field(属性标签) 就可以操作对应属性数据了, 既不需要重复存储数据, 也不会带来序列化和并发修改控制的问题。很好的解决了问题。

常用命令: hget,hset,hgetall 等。

HSET 命令用来给字段赋值, 而 HGET 命令用来获得字段的值。用法如下:

```
HSET key field value
HGET key field
redis> HSET car price 500
(integer) 1
redis> HSET car name BMW
(integer) 1
redis> HGET car name
"BMW"

HMSET key field1 value1 field2 value2
redis> HMGET car price name
1) "500"
2) "BMW"
```

HSET 命令的方便之处在于不区分插入和更新操作, 这意味着修改数据时不用事先判断字段是否存在来决定要执行的是插入操作 (update) 还是更新操作 (insert)。当执行的是插入操作时 (即之前字段不存在) HSET 命令会返回 1, 当执行的是更新操作时 (即之前字段已经存在) HSET 命令会返回 0。更进一步, 当键本身不存在时, HSET 命令还会自动建立它。

Redis 提供了接口 (hgetall) 可以直接取到全部的属性数据, 但是如果内部 Map 的成员很多, 那么涉及到遍历整个内部 Map 的操作, 由于 Redis 单线程模型的原因, 这个遍历操作可能会比较耗时, 而另其它客户端的请求完全不响应, 这点需要格外注意。

实现方式:

上面已经说到 Redis Hash 对应 Value 内部实际就是一个 HashMap, 实际这里会有 2 种不同实现, 这个 Hash 的成员比较少时 Redis 为了节省内存会采用类似一维数组的方式来紧凑存储, 而不会采用真正

的 HashMap 结构，对应的 value redisObject 的 encoding 为 zipmap，当成员数量增大时会自动转成真正的 HashMap，此时 encoding 为 ht。

- **List:**

常用命令：lpush,rpush,lpop,rpop,lrange 等。

应用场景：

Redis list 的应用场景非常多，也是 Redis 最重要的数据结构之一，比如 twitter 的关注列表，粉丝列表等都可以用 Redis 的 list 结构来实现。

Lists 就是链表，相信略有数据结构知识的人都应该能理解其结构。使用 Lists 结构，我们可以轻松地实现最新消息排行等功能。Lists 的另一个应用就是消息队列，可以利用 Lists 的 PUSH 操作，将任务存在 Lists 中，然后工作线程再用 POP 操作将任务取出进行执行。Redis 还提供了操作 Lists 中某一段的 api，你可以直接查询，删除 Lists 中某一段的元素。

实现方式：

Redis list 的实现为一个双向链表，即可以支持反向查找和遍历，更方便操作，不过带来了部分额外的内存开销，Redis 内部的很多实现，包括发送缓冲队列等也都是用的这个数据结构。

- **Set**

常用命令：sadd,spop,smembers,sunion 等。

应用场景：

Redis set 对外提供的功能与 list 类似是一个列表的功能，特殊之处在于 set 是可以自动排重的，当你需要存储一个列表数据，又不希望出现重复数据时，set 是一个很好的选择，并且 set 提供了判断某个成员是否在一个 set 集合内的重要接口，这个也是 list 所不能提供的。

Sets 集合的概念就是一堆不重复值的组合。利用 Redis 提供的 Sets 数据结构，可以存储一些集合性的数据，比如在微博应用中，可以将一个用户所有的关注人存在一个集合中，将其所有粉丝存在一个集合。Redis 还为集合提供了求交集、并集、差集等操作，可以非常方便的实现如共同关注、共同喜好、二度好友等功能，对上面的所有集合操作，你还可以使用不同的命令选择将结果返回给客户端还是存集到一个新的集合中。

实现方式：

set 的内部实现是一个 value 永远为 null 的 HashMap，实际就是通过计算 hash 的方式来快速排重的，这也是 set 能提供判断一个成员是否在集合内的原因。

- **Sorted Set**

常用命令：zadd,zrange,zrem,zcard 等

使用场景：

Redis sorted set 的使用场景与 set 类似，区别是 set 不是自动有序的，而 sorted set 可以通过用户额外提供一个优先级 (score) 的参数来为成员排序，并且是插入有序的，即自动排序。当你需要一个有序的并且不重复的集合列表，那么可以选择 sorted set 数据结构，比如 twitter 的 public timeline 可以以发表时间作为 score 来存储，这样获取时就是自动按时间排好序的。

另外还可以用 Sorted Sets 来做带权重的队列，比如普通消息的 score 为 1，重要消息的 score 为 2，然后工作线程可以选择按 score 的倒序来获取工作任务。让重要的任务优先执行。

实现方式：

Redis sorted set 的内部使用 HashMap 和跳跃表 (SkipList) 来保证数据的存储和有序，HashMap 里放的是成员到 score 的映射，而跳跃表里存放的是所有的成员，排序依据是 HashMap 里存的 score，使用跳跃表的结构可以获得比较高的查找效率，并且在实现上比较简单。

参考如下:

Redis 之数据存储结构

Redis 面试题

三、消息队列

一般使用 list 结构作为队列，rpush 生产消息，lpop 消费消息。当 lpop 没有消息的时候，要适当 sleep 一会再重试。

```

#无限循环读取任务队列中的内容
loop
task=RPOP queue
if task
#如果任务队列中有任务则执行它
execute( task)
else
#如果没有则等待1秒以免过于频繁地请求数据
wait 1 second

```

当任务队列中没有任务时消费者每秒都会调用一次 RPOP 命令查看是否有新任务。如果可以实现一旦有新任务加入任务队列就通知消费者就好了。其实借助 BRPOP 命令就可以实现这样的需求。

BRPOP 命令和 RPOP 命令相似，唯一的区别是当列表中没有元素时 BRPOP 命令会一直阻塞住连接，在没有消息的时候，它会阻塞住直到有新元素加入。

```

loop
#如果任务队列中没有新任务，BRPOP 命令会一直阻塞，不会执行execute()。
task=BRPOP queue, 0
#返回值是一个数组（见下介绍），数组第二个元素是我们需要的任务。
execute( task[1])

```

BRPOP 命令接收两个参数，第一个是键名，第二个是超时时间，单位是秒。当超过了此时间仍然没有获得新元素的话就会返回 nil。上例中超时时间为“0”，表示不限制等待的时间，即如果没有新元素加入列表就会永远阻塞下去。

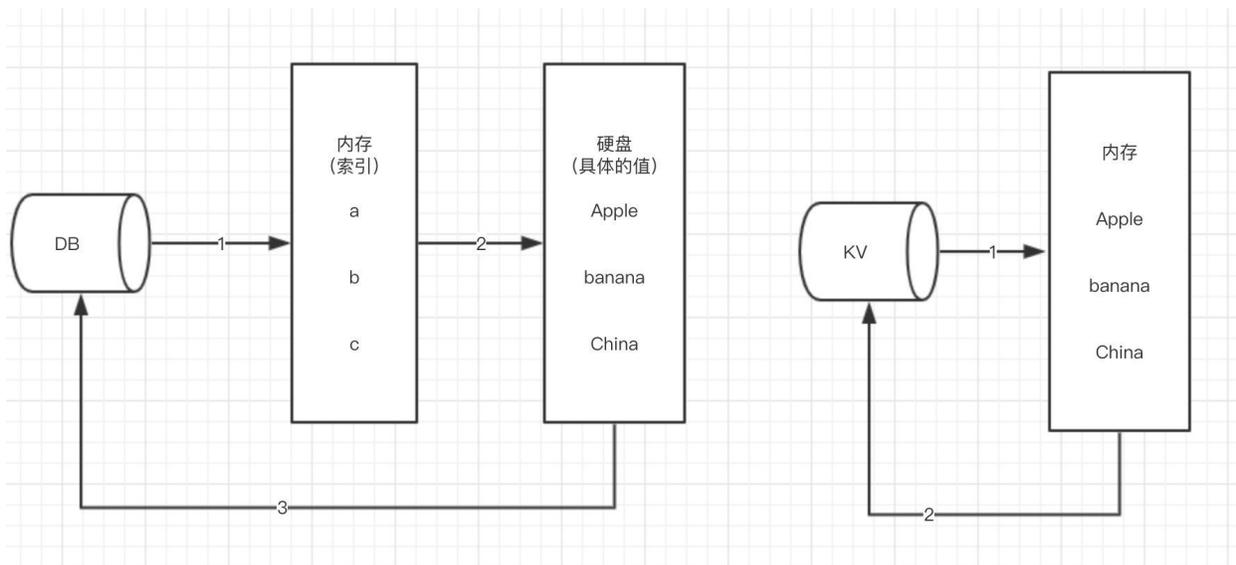
四、订阅推送

Pub/Sub 从字面上理解就是发布 (Publish) 与订阅 (Subscribe)，在 Redis 中，你可以设定对某一个 key 值进行消息发布及消息订阅，当一个 key 值上进行了消息发布后，所有订阅它的客户端都会收到相应的消息。这一功能最明显的用法就是用作实时消息系统，比如普通的即时聊天，群聊等功能。

“发布/订阅”模式中包含两种角色，分别是发布者和订阅者。订阅者可以订阅一个或若干个频道 (channel)，而发布者可以向指定的频道发送消息，所有订阅此频道的订阅者都会收到此消息。

pub/sub 在消费者下线的情况下，生产的消息会丢失。

4.2.3 redis 面试



Redis 采用的是基于内存的采用的是单进程单线程模型的 KV 数据库，由 C 语言编写，官方提供的数据库是可以达到 100000+ 的 QPS（每秒内查询次数）。

- 完全基于内存，绝大部分请求是纯粹的内存操作，非常快速。它的数据存在内存中，类似于 **HashMap**，**HashMap** 的优势就是查找和操作的时间复杂度都是 $O(1)$ ；
- 数据结构简单，对数据操作也简单，**Redis** 中的数据结构是专门进行设计的；
- 采用单线程，避免了不必要的上下文切换和竞争条件，也不存在多进程或者多线程导致的切换而消耗 CPU，不用去考虑各种锁的问题，不存在加锁释放锁操作，没有因为可能出现死锁而导致的性能消耗；
- 使用多路 I/O 复用模型，非阻塞 IO，将速度优势发挥到最大，也提供了较简单的计算功能；
- 使用底层模型不同，它们之间底层实现方式以及与客户端之间通信的应用协议不一样，**Redis** 直接自己构建了 VM 机制，因为一般的系统调用系统函数的话，会浪费一定的时间去移动和请求；
- **Redis** 提供了事务的功能，可以保证一串命令的原子性，中间不会被任何操作打断

Redis 虽然是单线程的，但是可以通过开多个 **Redis** 实例利用多核机器

Redis 持久化数据和缓存怎么做扩容？

- 如果 **Redis** 被当做缓存使用，使用一致性哈希实现动态扩容扩容。

shared 一致性哈希采用以下方案：

1. **Redis** 服务器节点划分：将每台服务器节点采用 hash 算法划分为 160 个虚拟节点 (可以配置划分权重)
2. 将划分虚拟节点采用 **TreeMap** 存储
3. 对每个 **Redis** 服务器的物理连接采用 **LinkedHashMap** 存储
4. 对 **Key** or **KeyTag** 采用同样的 hash 算法，然后从 **TreeMap** 获取大于等于键 hash 值得节点，取最近节点存储；当 **key** 的 hash 值大于虚拟节点 hash 值得最大值时，存入第一个虚拟节点

sharded 采用的 hash 算法：**MD5** 和 **MurmurHash** 两种；默认采用 64 位的 **MurmurHash** 算法；有兴趣的可以研究下，**MurmurHash** 是一种高效，低碰撞的 hash 算法；参考地址：

<http://blog.csdn.net/yfkiss/article/details/7337382>

<https://sites.google.com/site/murmurhash/>

- 如果 Redis 被当做一个持久化存储使用，必须使用固定的 keys-to-nodes 映射关系，节点的数量一旦确定不能变化。否则的话（即 Redis 节点需要动态变化的情况），必须使用可以在运行时进行数据再平衡的一套系统，而当前只有 Redis 集群可以做到这样。

redis 的一些其他特点

- Redis 是单进程单线程的 redis 利用队列技术将并发访问变为串行访问，消除了传统数据库串行控制的开销
- 读写分离模型通过增加 Slave DB 的数量，读的性能可以线性增长。为了避免 Master DB 的单点故障，集群一般都会采用两台 Master DB 做双机热备，所以整个集群的读和写的可用性都非常高。读写分离架构的缺陷在于，不管是 Master 还是 Slave，每个节点都必须保存完整的数据，如果在数据量很大的情况下，集群的扩展能力还是受限于单个节点的存储能力，而且对于 Write-intensive 类型的应用，读写分离架构并不适合
- 数据分片模型为了解决读写分离模型的缺陷，可以将数据分片模型应用进来。可以将每个节点看成都是独立的 master，然后通过业务实现数据分片。结合上面两种模型，可以将每个 master 设计成由一个 master 和多个 slave 组成的模型。
- Redis 的回收策略
 - volatile-lru: 从已设置过期时间的数据集 (server.db[i].expires) 中挑选最近最少使用的数据淘汰
 - volatile-ttl: 从已设置过期时间的数据集 (server.db[i].expires) 中挑选将要过期的数据淘汰
 - volatile-random: 从已设置过期时间的数据集 (server.db[i].expires) 中任意选择数据淘汰
 - allkeys-lru: 从数据集 (server.db[i].dict) 中挑选最近最少使用的数据淘汰
 - allkeys-random: 从数据集 (server.db[i].dict) 中任意选择数据淘汰
 - no-eviction (驱逐): 禁止驱逐数据

注意这里的 6 种机制，volatile 和 allkeys 规定了是对已设置过期时间的数据集淘汰数据还是从全部数据集淘汰数据，后面的 lru、ttl 以及 random 是三种不同的淘汰策略，再加上一种 no-eviction 永不回收的策略。使用策略规则：如果数据呈现幂律分布，也就是一部分数据访问频率高，一部分数据访问频率低，则使用 allkeys-lru 如果数据呈现平等分布，也就是所有的数据访问频率都相同，则使用 allkeys-random

Redis 虽然是一种内存型数据库，一旦服务器进程退出，数据库的数据就会丢失，为了解决这个问题 Redis 提供了两种持久化的方案，将内存中的数据保存到磁盘中，避免数据的丢失。

RDB 持久化

redis 提供了 RDB 持久化的功能，这个功能可以将 redis 在内存中的的状态保存到硬盘中，它可以**手动执行**，也可以再 redis.conf 中配置，**定期执行**。

RDB 持久化产生的 RDB 文件是一个**经过压缩**的二进制文件，这个文件被保存在硬盘中，redis 可以通过这个文件还原数据库当时的状态。

RDB 的创建与载入

RDB 文件可以通过两个命令来生成：

- SAVE：阻塞 redis 的服务器进程，直到 RDB 文件被创建完毕。
- BGSAVE：派生 (fork) 一个子进程来创建新的 RDB 文件，记录接收到 BGSAVE 当时的数据库状态，父进程继续处理接收到的命令，子进程完成文件的创建之后，会发送信号给父进程，而与此同时，父进程处理命令的同时，通过轮询来接收子进程的信号。

而 RDB 文件的载入一般情况是自动的，redis 服务器启动的时候，redis 服务器再启动的时候如果检测到 RDB 文件的存在，那么 redis 会自动载入这个文件。

如果服务器开启了 AOF 持久化，那么服务器会优先使用 AOF 文件来还原数据库状态。

RDB 是通过保存键值对来记录数据库状态的，采用 copy on write 的模式，每次都是全量的备份。

自动保存间隔

BGSAVE 可以在不阻塞主进程的情况下完成数据的备份。可以通过 redis.conf 中设置多个自动保存条件，只要有一个条件被满足，服务器就会执行 BGSAVE 命令。

```
# 以下配置表示的条件：  
# 服务器在900秒之内被修改了1次  
save 900 1  
# 服务器在300秒之内被修改了10次  
save 300 10  
# 服务器在60秒之内被修改了10000次  
save 60 10000
```

AOF 持久化

AOF 持久化 (Append-Only-File)，与 RDB 持久化不同，AOF 持久化是通过保存 Redis 服务器执行的写状态来记录数据库的。

具体来说，RDB 持久化相当于备份数据库状态，而 AOF 持久化是备份数据库接收到的命令，所有被写入 AOF 的命令都是以 redis 的协议格式来保存的。

在 AOF 持久化的文件中，数据库会记录下所有变更数据库状态的命令，除了指定数据库的 select 命令，其他的命令都是来自 client 的，这些命令会以追加 (append) 的形式保存到文件中。

服务器配置中有一项 appendfsync，这个配置会影响服务器多久完成一次命令的记录：

- always：将缓存区的内容总是即时写到 AOF 文件中。
- everysec：将缓存区的内容每隔一秒写入 AOF 文件中。
- no：写入 AOF 文件中的操作由操作系统决定，一般而言为了提高效率，操作系统会等待缓存区被填满，才会开始同步数据到磁盘。

redis 默认实用的是 everysec。

redis 在载入 AOF 文件的时候，会创建一个虚拟的 client，把 AOF 中每一条命令都执行一遍，最终还原回数据库的状态，它的载入也是自动的。在 RDB 和 AOF 备份文件都有的情况下，redis 会优先载入 AOF 备份文件

AOF 文件可能会随着服务器运行的时间越来越大，可以利用 AOF 重写的功能，来控制 AOF 文件的大小。AOF 重写功能会首先读取数据库中现有的键值对状态，然后根据类型使用一条命令来替代前的键值对多条命令。

AOF 重写功能有大量写入操作，所以 redis 才用子进程来处理 AOF 重写。这里带来一个新的问题，由于处理重新的是子进程，这意味着如果主线程的数据在此时被修改，备份的数据和主库的数据将会有不一致的情况发生。因此 redis 还设置了一个 AOF 重写缓冲区，这个缓冲区在子进程被创建开始之后开始使用，这个期间，所有的命令会被存两份，一份在 AOF 缓存空间，一份在 AOF 重写缓冲区，当 AOF 重写完成之后，子进程发送信号给主进程，通知主进程将 AOF 重写缓冲区的内容添加到 AOF 文件中。

相关配置

```
#AOF 和 RDB 持久化方式可以同时启动并且无冲突。
#如果AOF开启，启动redis时会加载aof文件，这些文件能够提供更好的保证。
appendonly yes

# 只增文件的文件名称。（默认是appendonly.aof）
# appendfilename appendonly.aof
#redis支持三种不同的写入方式：
#
# no:不调用，之等待操作系统来清空缓冲区当操作系统要输出数据时。很快。
# always: 每次更新数据都写入仅增日志文件。慢，但是最安全。
# everysec: 每秒调用一次。折中。
appendfsync everysec

# 设置为yes表示rewrite期间对新写操作不fsync,暂时存在内存中,等rewrite完成后再写入。
→官方文档建议如果你有特殊的情况可以配置为'yes'。但是配置为'no'是最为安全的选择。
no-appendfsync-on-rewrite no

# 自动重写只增文件。
#_
→redis可以自动盲从的调用‘BGREWRITEAOF’来重写日志文件，如果日志文件增长了指定的百分比。
# 当前AOF文件大小是上次日志重写得到AOF文件大小的二倍时，自动启动新的日志重写过程。
auto-aof-rewrite-percentage 100
#_
→当前AOF文件启动新的日志重写过程的最小值，避免刚刚启动Reids时由于文件尺寸较小导致频繁的重写。
auto-aof-rewrite-min-size 64mb
```

对比

- AOF 更安全，可将数据及时同步到文件中，但需要较多的磁盘 IO，AOF 文件尺寸较大，文件内容恢复相对较慢，也更完整。
- RDB 持久化，安全性较差，它是正常时期数据备份及 master-slave 数据同步的最佳手段，文件尺寸较小，恢复数度较快。

作者：whthomas 链接：<https://www.jianshu.com/p/bedec93e5a7b> 来源：简书著作权归作者所有。商业转载请联系作者获得授权，非商业转载请注明出处。

4.2.4 redis 分布式锁

分布式锁

分布式锁一般有三种实现方式：1. 数据库乐观锁；2. 基于 Redis 的分布式锁；3. 基于 ZooKeeper 的分布式锁。本篇博客将介绍第二种方式，基于 Redis 实现分布式锁。虽然网上已经有各种介绍 Redis 分布式锁实现的博客，然而他们的实现却有着各种各样的问题，为了避免误人子弟，本篇博客将详细介绍如何正确地实现 Redis 分布式锁。

可靠性

首先，为了确保分布式锁可用，我们至少要确保锁的实现同时满足以下四个条件：

1. **互斥性**在任意时刻，只有一个客户端能持有锁。
2. **不会发生死锁**即使有一个客户端在持有锁的期间崩溃而没有主动解锁，也能保证后续其他客户端能加锁。
3. **具有容错性**只要大部分的 Redis 节点正常运行，客户端就可以加锁和解锁。
4. **解铃还须系铃人**加锁和解锁必须是同一个客户端，客户端自己不能把别人加的锁给解了。

代码实现

组件依赖

首先我们要通过 Maven 引入 Jedis 开源组件，在 pom.xml 文件加入下面的代码：

```
<dependency>
  <groupId>redis.clients</groupId>
  <artifactId>jedis</artifactId>
  <version>2.9.0</version>
</dependency>
```

加锁代码

Talk is cheap, show me the code。先展示代码，再带大家慢慢解释为什么这样实现：

```
public class RedisTool {

    private static final String LOCK_SUCCESS = "OK";
    private static final String SET_IF_NOT_EXIST = "NX";
    private static final String SET_WITH_EXPIRE_TIME = "PX";

    /**
     * 尝试获取分布式锁
     * @param jedis Redis客户端
     * @param lockKey 锁
     * @param requestId 请求标识
     * @param expireTime 超期时间
     * @return 是否获取成功
     */
    public static boolean tryGetDistributedLock(Jedis jedis, String lockKey, String
↵requestId, int expireTime) {
```

(下页继续)

(续上页)

```

        String result = jedis.set(lockKey, requestId, SET_IF_NOT_EXIST, SET_WITH_
↳EXPIRE_TIME, expireTime);

        if (LOCK_SUCCESS.equals(result)) {
            return true;
        }
        return false;
    }
}

```

可以看到，我们加锁就一行代码：`jedis.set(String key, String value, String nxxx, String expx, int time)`，这个 `set()` 方法一共有五个形参：

- 第一个为 `key`，我们使用 `key` 来当锁，因为 `key` 是唯一的。
- 第二个为 `value`，我们传的是 `requestId`，很多童鞋可能不明白，有 `key` 作为锁不就够了吗，为什么还要用到 `value`？原因就是我们在上面讲到可靠性时，分布式锁要满足第四个条件**解铃还须系铃人**，通过给 `value` 赋值为 `requestId`，我们就知道这把锁是哪个请求加的了，在解锁的时候就可以有依据。`requestId` 可以使用 `UUID.randomUUID().toString()` 方法生成。
- 第三个为 `nxxx`，这个参数我们填的是 `NX`，意思是 `SET IF NOT EXIST`，即当 `key` 不存在时，我们进行 `set` 操作；若 `key` 已经存在，则不做任何操作；
- 第四个为 `expx`，这个参数我们传的是 `PX`，意思是我们要给这个 `key` 加一个过期的设置，具体时间由第五个参数决定。
- 第五个为 `time`，与第四个参数相呼应，代表 `key` 的过期时间。

总的来说，执行上面的 `set()` 方法就只会导致两种结果：1. 当前没有锁（`key` 不存在），那么就进行加锁操作，并对锁设置个有效期，同时 `value` 表示加锁的客户端。2. 已有锁存在，不做任何操作。

心细的童鞋就会发现了，我们的加锁代码满足我们**可靠性**里描述的三个条件。首先，`set()` 加入了 `NX` 参数，可以保证如果已有 `key` 存在，则函数不会调用成功，也就是只有一个客户端能持有锁，满足互斥性。其次，由于我们对锁设置了过期时间，即使锁的持有者后续发生崩溃而没有解锁，锁也会因为到了过期时间而自动解锁（即 `key` 被删除），不会发生死锁。最后，因为我们将 `value` 赋值为 `requestId`，代表加锁的客户端请求标识，那么在客户端在解锁的时候就可以进行校验是否是同一个客户端。由于我们只考虑 Redis 单机部署的场景，所以容错性我们暂不考虑。

解锁代码

还是先展示代码，再带大家慢慢解释为什么这样实现：

```

public class RedisTool {

    private static final Long RELEASE_SUCCESS = 1L;

    /**
     * 释放分布式锁
     * @param jedis Redis客户端
     * @param lockKey 锁
     * @param requestId 请求标识
     * @return 是否释放成功
     */
    public static boolean releaseDistributedLock(Jedis jedis, String lockKey, String_
↳requestId) {

```

(下页继续)

(续上页)

```

String script = "if redis.call('get', KEYS[1]) == ARGV[1] then return redis.
↪call('del', KEYS[1]) else return 0 end";
Object result = jedis.eval(script, Collections.singletonList(lockKey), ↪
↪Collections.singletonList(requestId));

    if (RELEASE_SUCCESS.equals(result)) {
        return true;
    }
    return false;
}
}

```

可以看到，我们解锁只需要两行代码就搞定了！第一行代码，我们写了一个简单的 Lua 脚本代码，上一次见到这个编程语言还是在《黑客与画家》里，没想到这次居然用上了。第二行代码，我们将 Lua 代码传到 `jedis.eval()` 方法里，并使参数 `KEYS[1]` 赋值为 `lockKey`，`ARGV[1]` 赋值为 `requestId`。`eval()` 方法是将 Lua 代码交给 Redis 服务端执行。

那么这段 Lua 代码的功能是什么呢？其实很简单，首先获取锁对应的 `value` 值，检查是否与 `requestId` 相等，如果相等则删除锁（解锁）。那么为什么要使用 Lua 语言来实现呢？因为要确保上述操作是原子性的。关于非原子性会带来什么问题，会在锁已经不属于当前客户端的时候会解除他人加的锁。那么为什么执行 `eval()` 方法可以确保原子性，源于 Redis 的特性：

简单来说，就是在 `eval` 命令执行 Lua 代码的时候，Lua 代码将被当成一个命令去执行，并且直到 `eval` 命令执行完成，Redis 才会执行其他命令。

脚本

```

//加锁
conn.Do("SET", lockKey, locker, "PX", expiredInMS, "NX")

// 解锁脚本
var unlockScript = redis.NewScript(1, `
    if redis.call("get", KEYS[1]) == ARGV[1]
    then
        return redis.call("del", KEYS[1])
    else
        return 0
    end
`)

var unlockScriptWithState = redis.NewScript(1, `
    if redis.call("get", KEYS[1]) == ARGV[1]
    then
        redis.call("set", ARGV[2], ARGV[3])
        return redis.call("del", KEYS[1])
    else
        return -1
    end
`)

```

原文链接

Redis 分布式锁的正确实现方式 (Java 版)

4.3 Pulsar

4.3.1 pulsar-go

客户端

Pulsar 协议 url 使用 pulsar scheme 来指定被连接的集群，默认端口为 6650。以下是 localhost 的示例：

```
pulsar://localhost:6650
```

生产环境的 Pulsar 集群 URL 类似这样：

```
pulsar://pulsar.us-west.example.com:6650
```

创建客户端

```
import (
    "log"
    "runtime"

    "github.com/apache/pulsar/pulsar-client-go/pulsar"
)

func main() {
    client, err := pulsar.NewClient(pulsar.ClientOptions{
        URL: "pulsar://localhost:6650",
        OperationTimeoutSeconds: 5,
        MessageListenerThreads: runtime.NumCPU(),
    })

    if err != nil {
        log.Fatalf("Could not instantiate Pulsar client: %v", err)
    }
}
```

Producers

```
producer, err := client.CreateProducer(pulsar.ProducerOptions{
    Topic: "my-topic",
})

if err != nil {
    log.Fatalf("Could not instantiate Pulsar producer: %v", err)
}

defer producer.Close()
```

(下页继续)

(续上页)

```
msg := pulsar.ProducerMessage{
    Payload: []byte("Hello, Pulsar"),
}

if err := producer.Send(msg); err != nil {
    log.Fatalf("Producer could not send message: %v", err)
}
```

Producer operations

Consumers

消费者订阅一个或多个主题，并听取在该主题/这些主题上产生的传入消息。

```
msgChannel := make(chan pulsar.ConsumerMessage)

consumerOpts := pulsar.ConsumerOptions{
    Topic:           "my-topic",
    SubscriptionName: "my-subscription-1",
    Type:           pulsar.Exclusive,
    MessageChannel: msgChannel,
}

consumer, err := client.Subscribe(consumerOpts)

if err != nil {
    log.Fatalf("Could not establish subscription: %v", err)
}

defer consumer.Close()

for cm := range msgChannel {
    msg := cm.Message

    fmt.Printf("Message ID: %s", msg.ID())
    fmt.Printf("Message value: %s", string(msg.Payload()))

    consumer.Ack(msg)
}
```

Consumer operations

Readers

Pulsar readers process messages from Pulsar topics

readers 与 consumers 不同，因为使用读取器，您需要明确指定要从流中开始的消息（另一方面，consumers 会自动以最新的未确认消息开始）。

```
reader, err := client.CreateReader(pulsar.ReaderOptions{
    Topic: "my-golang-topic",
```

(下页继续)

(续上页)

```
    StartMessageId: pulsar.LatestMessage,  
  })
```


5.1 基础知识

5.1.1 数组

Go 语言中数组是值语义。一个数组变量即表示整个数组，它并不是隐式的指向第一个元素的指针（比如 C 语言的数组），而是一个完整的值。当一个数组变量被赋值或者被传递的时候，实际上会复制整个数组。如果数组较大的话，数组的赋值也会有较大的开销。为了避免复制数组带来的开销，可以传递一个指向数组的指针，但是数组指针并不是数组。

```
var a = [...]int{1, 2, 3} // a 是一个数组
var b = &a                // b 是指向数组的指针

fmt.Println(a[0], a[1])  // 打印数组的前2个元素
fmt.Println(b[0], b[1])  // 通过数组指针访问数组元素的方式和数组类似

for i, v := range b {    // 通过数组指针迭代数组的元素
    fmt.Println(i, v)
}
```

其实数组指针类型除了类型和数组不同之外，通过数组指针操作数组的方式和通过数组本身的操作类似，而且数组指针赋值时只会拷贝一个指针。

5.1.2 字符串

Go 语言字符串的底层结构在 `reflect.StringHeader` 中定义：

```
type StringHeader struct {
    Data uintptr
    Len  int
}
```

`uintptr` 是 `golang` 的内置类型，是能存储指针的整型

字符串结构由两个信息组成：第一个是字符串指向的底层字节数组，第二个是字符串的字节长度。字符串其实是一个结构体，因此字符串的赋值操作也就是 `reflect.StringHeader` 结构体的复制过程。

字符串虽然不是切片，但是支持切片操作，不同位置的切片底层也访问的同一块内存数据（因为字符串是只读的，相同的字符串面值常量通常是对应同一个字符串常量）：

5.1.3 切片

`reflect.SliceHeader`：

```
type SliceHeader struct {
    Data uintptr
    Len  int
    Cap  int
}
```

`uintptr` 是 `golang` 的内置类型，是能存储指针的整型

只要是切片的底层数据指针、长度和容量没有发生变化的话，对切片的遍历、元素的读取和修改都和数组是一样的。在对切片本身赋值或参数传递时，和数组指针的操作方式类似，只是复制切片头信息（`reflect.SliceHeader`），并不会复制底层的数据。对于类型，和数组的最大不同是，切片的类型和长度信息无关，只要是相同类型元素构成的切片均对应相同的切片类型。

切片操作并不会复制底层的数据。底层的数组会被保存在内存中，直到它不再被引用。但是有时候可能会因为一个小的内存引用而导致底层整个数组处于被使用的状态，这会延迟自动内存回收器对底层数组的回收。

切片中的底层数组部分是通过隐式指针传递（指针本身依然是传值的，但是指针指向的却是同一份的数据） 所以被调用函数是可以通过指针修改掉调用参数切片中的数据。

除了数据之外，切片结构还包含了切片长度和切片容量信息，这 2 个信息也是传值的。

如果被调用函数中修改了 `Len` 或 `Cap` 信息的话，就无法反映到调用参数的切片中，这时候我们一般会通过返回修改后的切片来更新之前的切片。这也是为何内置的 `append` 必须要返回一个切片的原因。

5.1.4 Make 和 New

虽然 `make` 和 `new` 都是能够用于初始化数据结构，但是它们两者能够初始化的结构类型却有着较大的不同；`make` 在 `Go` 语言中只能用于初始化语言中的基本类型：

```
slice := make([]int, 0, 100)
hash := make(map[int]bool, 10)
ch := make(chan int, 5)
```

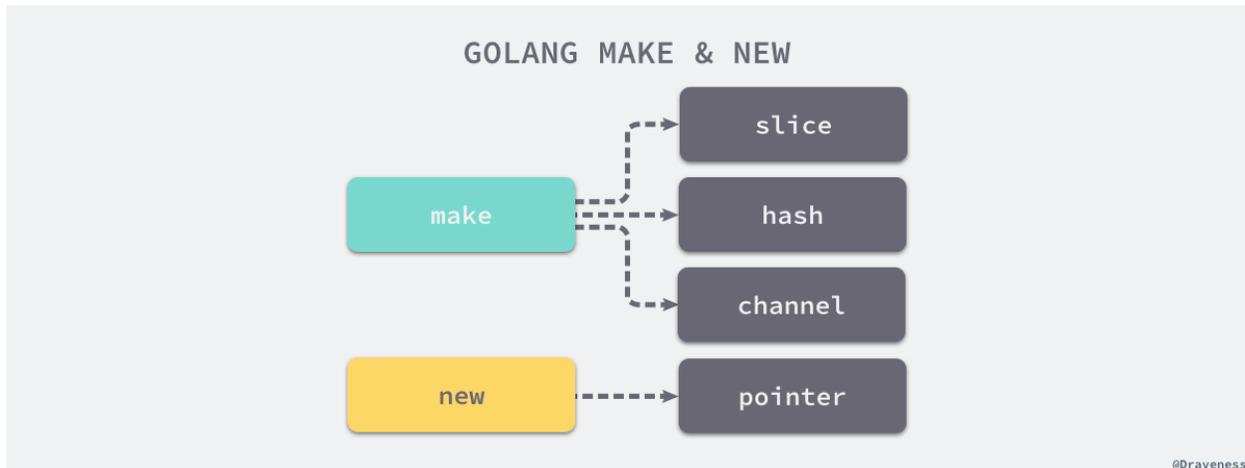
- `slice` 是一个包含 `data`、`cap` 和 `len` 的**结构体**；
- `hash` 是一个指向 `hmap` 结构体的**指针**；

- `ch` 是一个指向 `hchan` 结构体的指针；

而另一个用于初始化数据结构的关键字 `new` 的作用其实就非常简单了，它只是接收一个类型作为参数然后返回一个指向这个类型的指针：

```
i := new(int)
var v int
i := &v
```

上述代码片段中的两种不同初始化方法其实是等价的，它们都会创建一个指向 `int` 零值的指针。

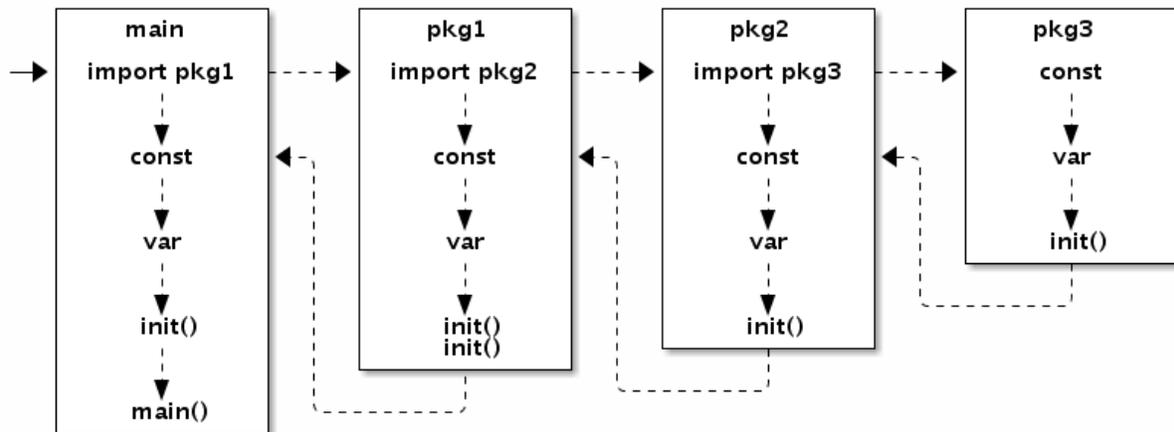


到了这里我们对 Go 语言中这两种不同关键字的使用也有了一定的了解：`make` 用于创建切片、哈希表和管道等内置数据结构，`new` 用于分配并创建一个指向对应类型的指针。

5.1.5 函数初始化

Go 语言中的函数有具名和匿名之分：具名函数一般对应于包级的函数，是匿名函数的一种特例，当匿名函数引用了外部作用域中的变量时就成了闭包函数，闭包函数是函数式编程语言的核心。

Go 语言程序的初始化和执行总是从 `main.main` 函数开始的。但是如果 `main` 包导入了其它的包，则会按照顺序将它们包含进 `main` 包里（这里的导入顺序依赖具体实现，一般可能是以文件名或包路径名的字符串顺序导入）。如果某个包被多次导入的话，在执行的时候只会导入一次。当一个包被导入时，如果它还导入了其它的包，则先将其它的包包包含进来，然后创建和初始化这个包的常量和变量，再调用包里的 `init` 函数，如果一个包有多个 `init` 函数的话，调用顺序未定义（实现可能是以文件名的顺序调用），同一个文件内的多个 `init` 则是以出现的顺序依次调用（`init` 不是普通函数，可以定义有多个，所以也不能被其它函数调用）。最后，当 `main` 包的所有包级常量、变量被创建和初始化完成，并且 `init` 函数被执行后，才会进入 `main.main` 函数，程序开始正常执行。



5.1.6 函数调用

Go 语言中对于整型和数组类型的参数都是值传递的，也就是在调用函数时会对内容进行拷贝，需要注意的是如果当前数组的大小非常的大，这种直接复制传值的方式就会对性能造成比较大的影响。

函数可以直接或间接地调用自己，也就是支持递归调用。Go 语言函数的递归调用深度逻辑上没有限制，函数调用的栈是不会出现溢出错误的，因为 Go 语言运行时会根据需要动态地调整函数栈的大小。每个 goroutine 刚启动时只会分配很小的栈（4 或 8KB，具体依赖实现），根据需要动态调整栈的大小，栈最大可以达到 GB 级（依赖具体实现，在目前的实现中，32 位体系结构为 250MB,64 位体系结构为 1GB）

```
func f(x int) *int {
    return &x
}

func g() int {
    x = new(int)
    return *x
}
```

第一个函数直接返回了函数参数变量的地址——这似乎是不可以的，因为如果参数变量在栈上的话，函数返回之后栈变量就失效了，返回的地址自然也应该失效了。但是 Go 语言的编译器和运行时比我们聪明的多，它会保证指针指向的变量在合适的地方。第二个函数，内部虽然调用 new 函数创建了 *int 类型的指针对象，但是依然不知道它具体保存在哪里。

任何可以通过函数参数修改调用参数的情形，都是因为函数参数中显式或隐式传入了指针参数。

5.1.7 方法

Go 语言的方法却是关联到类型的，这样可以在编译阶段完成方法的静态绑定。一个面向对象的程序会用方法来表达其属性对应的操作，这样使用这个对象的用户就不需要直接去操作对象，而是借助方法来做一些事情。

对于给定的类型，每个方法的名字必须是唯一的，同时方法和函数一样也不支持重载。

5.1.8 接口

Go 的接口类型是对其它类型行为的抽象和概括；因为接口类型不会和特定的实现细节绑定在一起。Go 语言的接口类型是延迟绑定，可以实现类似虚函数的多态功能。

Go 语言对于接口类型的转换则非常的灵活。对象和接口之间的转换、接口和接口之间的转换都可能是隐式的转换。

```
var (
    a io.ReadCloser = (*os.File)(f) // 隐式转换, *os.File 满足 io.ReadCloser 接口
    b io.Reader      = a             // 隐式转换, io.ReadCloser 满足 io.Reader 接口
    c io.Closer      = a             // 隐式转换, io.ReadCloser 满足 io.Closer 接口
    d io.Reader      = c.(io.Reader) // 显式转换, io.Closer 不满足 io.Reader 接口
)
```

interface 应用场景

类型转换

类型推断可将接口变量还原为原始类型，或用来判断是否实现了某个更具体的接口类型。

```
type data int

func(d data)String()string{
    return fmt.Sprintf("data:%d",d)
}

func main() {
    var d data=15
    var x interface{} =d

    if n,ok:=x.(fmt.Stringer);ok{ // 转换为更具体的接口类型
        fmt.Println(n)
    }

    if d2,ok:=x.(data);ok{ // 转换回原始类型
        fmt.Println(d2)
    }

    e:=x.(error) // 错误:main.data is not error
    fmt.Println(e)
}

data:15
data:15
panic:interface conversion:main.data is not error:missing method Error
```

实现多态功能

多态功能是 interface 实现的重要功能，也是 Golang 中的一大行为特色，其多态功能一般要结合 Go method 实现，作为函数参数可以容易的实现多态功能。

5.1.9 指针和引用

结构体和指针在参数的传递过程中：如果传递的是结构体，那么在传递参数时依然会对结构体中的全部内容进行拷贝，而传递指针时复制的其实也是指针的内容 - 地址

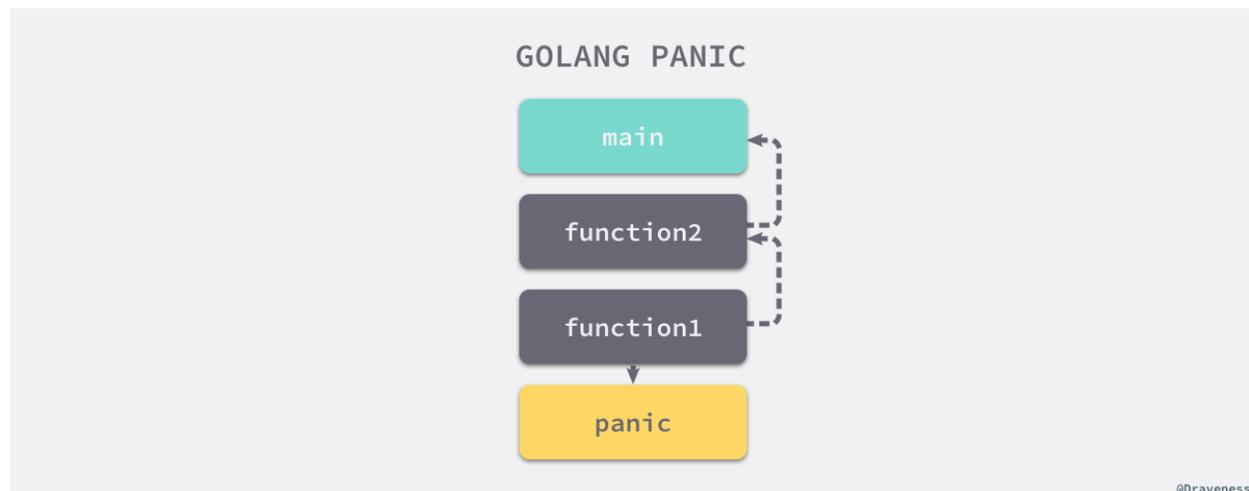
5.1.10 函数调用

Go 通过堆栈的方式对函数的参数和返回值进行传递和接受，在调用函数之前会在栈上为返回值分配合适的内存空间，随后按照入参从右到左按顺序压栈，被调用方接受参数时会对参数进行拷贝后再进行计算，返回值最终会被放置到调用者预留好的栈空间上，Go 语言函数调用的原理可以总结成以下的几条规则：

- 1、通过堆栈传递参数，入栈的顺序是从右到左；
- 2、函数返回值通过堆栈传递并由调用者预先分配内存空间；
- 3、调用函数时都是传值，接收方会对入参进行复制再计算；

5.1.11 panic 和 recover

panic 和 recover 两个关键字其实都是 Go 语言中的内置函数，panic 能够改变程序的控制流，当一个函数调用执行 panic 时，它会立刻停止执行函数中其他的代码，而是会运行其中的 defer 函数，执行成功后会返回到调用方。



对于上层调用方来说，调用导致 panic 的函数其实与直接调用 panic 类似，所以也会执行所有的 defer 函数并返回到它的调用方，这个过程会一直进行直到当前 Goroutine 的调用栈中不包含任何的函数，这时整个程序才会崩溃，这个『恐慌过程』不仅会被显式的调用触发，还会由于运行期间发生错误而触发。

然而 panic 导致的『恐慌』状态其实可以被 defer 中的 recover 中止，recover 是一个只在 defer 中能够发挥作用的函数，在正常的控制流程中，调用 recover 会直接返回 nil 并且没有任何的作用，但是如果当前的 Goroutine 发生了『恐慌』，recover 其实就能够捕获到 panic 抛出的错误并阻止『恐慌』的继续传播

5.1.12 参考资料:

golang 语言实现

5.2 开发常用

5.2.1 context 包

Go 语言中的每一个请求的都是通过一个单独的 Goroutine 进行处理的，HTTP/RPC 请求的处理器往往都会启动新的 Goroutine 访问数据库和 RPC 服务，我们可能会创建多个 Goroutine 来处理一次请求，而 Context 的主要作用就是在不同的 Goroutine 之间同步请求特定的数据、取消信号以及处理请求的截止日期。

Context 结构

```
// A Context carries a deadline, cancelation signal, and request-scoped values
// across API boundaries. Its methods are safe for simultaneous use by multiple
// goroutines.
type Context interface {
    // Done returns a channel that is closed when this Context is canceled
    // or times out.
    Done() <-chan struct{}

    // Err indicates why this context was canceled, after the Done channel
    // is closed.
    Err() error

    // Deadline returns the time when this Context will be canceled, if any.
    Deadline() (deadline time.Time, ok bool)

    // Value returns the value associated with key or nil if none.
    Value(key interface{}) interface{}
}
```

- Deadline 方法需要返回当前 Context 被取消的时间，也就是完成工作的截止日期；
- Done 方法需要返回一个 Channel，这个 Channel 会在当前工作完成或者上下文被取消之后关闭，多次调用 Done 方法会返回同一个 Channel；
- Err 方法会返回当前 Context 结束的原因，它只会在 Done 返回的 Channel 被关闭时才会返回非空的值；
 - 如果当前 Context 被取消就会返回 Canceled 错误；
 - 如果当前 Context 超时就会返回 DeadlineExceeded 错误；
- Value 方法会从 Context 中返回键对应的值，对于同一个上下文来说，多次调用 Value 并传入相同的 Key 会返回相同的结果，这个功能可以用来传递请求特定的数据；

Context 的实现方法

根 context

Context 虽然是个接口，但是并不需要使用方实现：Background() TODO

```

var (
    background = new(emptyCtx)
    todo = new(emptyCtx)
)
func Background() Context {
    //
    ↪主要用于main函数、初始化以及测试代码中，作为Context这个树结构的最顶层的Context，也就是根Context，它
    return background
}
func TODO() Context { //只在不确定时使用 context.TODO()
    return todo
}

```

一般在代码中，开始上下文的时候都是以这两个作为最顶层的 parent context，然后再衍生出子 context。这些 Context 对象形成一棵树：当一个 Context 对象被取消时，继承自它的所有 Context 都会被取消。

继承 context

一般情况下是以 context.Background() 做为根节点

```

//传递一个父Context作为参数，返回子Context，以及一个取消函数用来取消Context。
func WithCancel(parent Context) (ctx Context, cancel CancelFunc)

//
↪会多传递一个截止时间参数，意味着到了这个时间点，会自动取消Context，也可以提前通过取消函数进行取消。
func WithDeadline(parent Context, deadline time.Time) (Context, CancelFunc)

//超时自动取消，是多少时间后自动取消Context的意思。也可以提前通过取消函数进行取消。
func WithTimeout(parent Context, timeout time.Duration) (Context, CancelFunc)

// 为了生成一个绑定了一个键值对数据的Context，这个绑定的数据可以通过Context.
↪Value方法访问到，这是我们实际用经常要用到的技巧，一般我们想要通过上下文来传递数据时，可以通过这个方法
func WithValue(parent Context, key, val interface{}) Context

```

参考如下

[golang 服务器开发利器 context 用法详解](#)

[Golang Context 深入理解](#)

5.2.2 定时器

time.Timer

Timer 的结构定义

```
type Timer struct {
    C <-chan Time
    r runtimeTimer
}
```

其他方法

- `func After(d Duration) <-chan Time { return NewTimer(d).C }`

根据源码可以看到 `After` 直接是返回了 `Timer` 的 `channel`，这种就可以做超时处理。比如我们有这样一个需求：我们写了一个爬虫，爬虫在 `HTTP GET` 一个网页的时候可能因为网络的原因就一直等待着，这时候就需要做超时处理，比如只请求五秒，五秒以后直接丢掉不请求这个网页了，或者重新发起请求。

```
Get("http://baidu.com/")
func Get(url string) {
    response := make(chan string)
    response = http.Request(url)
    select {
    case html :=<- response:
        println(html)
    case <-time.After(time.Second * 5):
        println("超时处理")
    }
}
```

可以从代码中体现出来，如果五秒到了，网页的请求还没有下来就是执行超时处理，因为 `Timer` 的内部会是帮你在你设置的时间长度后自动向 `Timer.C` 中写入当前时间。

其实也可以写成这样：

```
func Get(url string) {
    response := make(chan string)
    response = http.Request(url)
    timeout := time.NewTimer(time.Second * 3)
    select {
    case html :=<- response:
        println(html)
    case <-timeout.C:
        println("超时处理")
    }
}
```

- `func (t *Timer) Reset(d Duration) bool`//强制的修改 `timer` 中规定的时间，`Reset` 会先调用 `stopTimer` 再调用 `startTimer`，类似于废弃之前的定时器，重新启动一个定时器，`Reset` 在 `Timer` 还未触发时返回 `true`；触发了或 `Stop` 了，返回 `false`。
- `func (t *Timer) Stop() bool`// 如果定时器还未触发，`Stop` 会将其移除，并返回 `true`；否则返回 `false`；后续再对该 `Timer` 调用 `Stop`，直接返回 `false`。
- `func AfterFunc(d Duration, f func()) *Timer` // 在时间 `d` 后自动执行函数 `f`

```
func main() {
    f := func() {fmt.Println("I Love You!")}
    time.AfterFunc(time.Second*2, f)
    time.Sleep(time.Second * 4)
}
```

time.Ticker

其实 Ticker 就是一个重复版本的 Timer，它会重复的在时间 d 后向 Ticker 中写数据

- func NewTicker(d Duration) *Ticker // 新建一个 Ticker
- func (t *Ticker) Stop() // 停止 Ticker
- func Tick(d Duration) <-chan Time // Ticker.C 的封装

参考

Golang time.Timer and time.Ticker

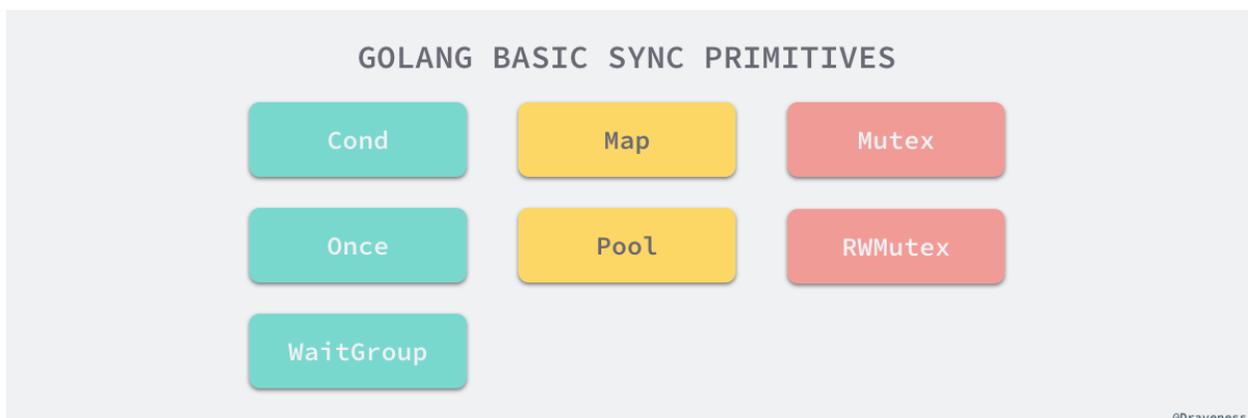
5.3 同步原语与锁

当提到并发编程、多线程编程时，我们往往都离不开『锁』这一概念，Go 语言作为一个原生支持用户态进程 Goroutine 的语言，也一定会为开发者提供这一功能，锁的主要作用就是保证多个线程或者 Goroutine 在访问同一片内存时不会出现混乱的问题，锁其实是一种并发编程中的同步原语（Synchronization Primitives）。

在这一节中我们就会介绍 Go 语言中常见的同步原语 Mutex、RWMutex、WaitGroup、Once 和 Cond 以及扩展原语 ErrGroup、Semaphore 和 SingleFlight 的实现原理，同时也会涉及互斥锁、信号量等并发编程中的常见概念。

5.3.1 1. 基本原语

Go 语言在 sync 包中提供了用于同步的一些基本原语，包括常见的互斥锁 Mutex 与读写互斥锁 RWMutex 以及 Once、WaitGroup。



这些基本原语的主要作用是提供较为基础的同步功能，我们应该使用 Channel 和通信来实现更加高级的同步机制，我们在这一节中并不会介绍标准库中全部的原语，而是会介绍其中比较常见的 Mutex、RWMutex、Once、WaitGroup 和 Cond，我们并不会涉及剩下两个用于存取数据的结构体 Map 和 Pool

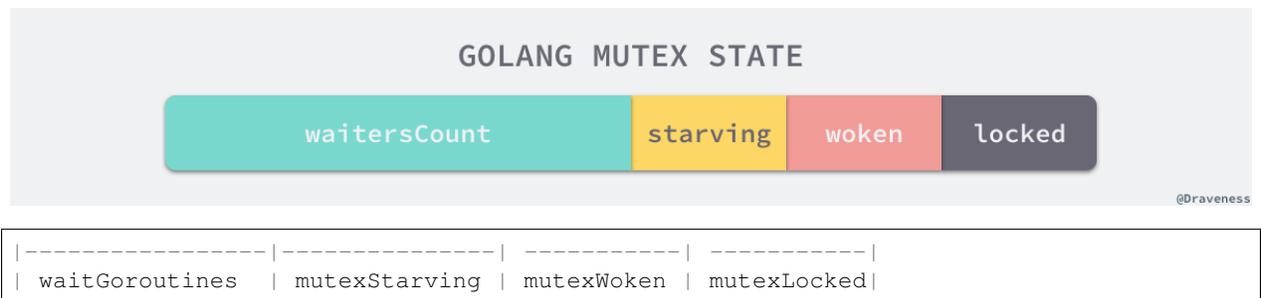
1.1. Mutex

Go 语言中的互斥锁在 sync 包中，它由两个字段 state 和 sema 组成，state 表示当前互斥锁的状态，而 sema 真正用于控制锁状态的信号量，这两个加起来只占 8 个字节空间的结构体就表示了 Go 语言中的互斥锁。

```
type Mutex struct {
    state int32
    sema  uint32
}
```

状态

互斥锁的状态是用 int32 来表示的，但是锁的 state 并不是互斥的，它的最低三位分别表示 mutexLocked、mutexWoken 和 mutexStarving，剩下的位置都用来表示当前有多少个 Goroutine 等待互斥锁被释放：

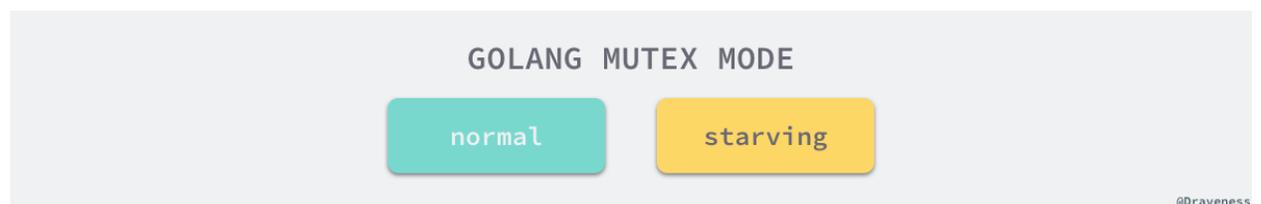


互斥锁在被创建出来时，所有的状态位的默认值都是 0，当互斥锁被锁定时 `mutexLocked` 就会被置成 1、当互斥锁被在正常模式下被唤醒时 `mutexWoken` 就会被置成 1、`mutexStarving` 用于表示当前的互斥锁进入了状态，最后的几位是在当前互斥锁上等待的 Goroutine 个数。

饥饿模式

Mutex 在可能会进入饥饿模式，饥饿模式主要功能就是保证互斥锁的获取的『公平性』(Fairness)。

互斥锁可以同时处于两种不同的模式，也就是正常模式和饥饿模式，在正常模式下，所有锁的等待者都会按照先进先出的顺序获取锁，但是如果一个刚刚被唤起的 Goroutine 遇到了新的 Goroutine 进程也调用了 Lock 方法时，大概率会获取不到锁，为了减少这种情况的出现，防止 Goroutine 被『饿死』，一旦 Goroutine 超过 1ms 没有获取到锁，它就会将当前互斥锁切换饥饿模式。



在饥饿模式中，互斥锁会被直接交给等待队列最前面的 Goroutine，新的 Goroutine 在这时不能获取锁、也不会进入自旋的状态，它们只会在队列的末尾等待，如果一个 Goroutine 获得了互斥锁并且它是队列中最末尾的协程或者它等待的时间少于 1ms，那么当前的互斥锁就会被切换回正常模式。

相比于饥饿模式，正常模式下的互斥锁能够提供更好地性能，饥饿模式的主要作用就是避免一些 Goroutine 由于陷入等待无法获取锁而造成较高的尾延时，这也是对 Mutex 的一个优化。

加锁

互斥锁 Mutex 的加锁是靠 Lock 方法完成的，最新的 Go 语言源代码中已经将 Lock 方法进行了简化，方法的主干只保留了最常见、简单并且快速的情况；当锁的状态是 0 时直接将 mutexLocked 位置成 1：

```
func (m *Mutex) Lock() {
    if atomic.CompareAndSwapInt32(&m.state, 0, mutexLocked) {
        return
    }
    m.lockSlow()
}
```

5.4 RPC

RPC (Remote Procedure Call)，即远程过程调用，是一个分布式系统间通信的必备技术。RPC 最核心要解决的问题就是在分布式系统间，如何执行另外一个地址空间上的函数、方法，就仿佛在本地图用一样。

rpc 是远端过程调用，其调用协议通常包含传输协议和编码协议

5.4.1 传输 (Transport)

TCP 协议是 RPC 的基石，一般来说通信是建立在 TCP 协议之上的，而且 RPC 往往需要可靠的通信，因此不采用 UDP。RPC 传输的 message 也就是 TCP body 中的数据，这个 message 也同样可以包含 header+body。body 也经常叫做 payload。TCP 协议栈存在端口的概念，端口是进程获取数据的渠道。

5.4.2 I/O 模型 (I/O Model)

做一个高性能 /scalable 的 RPC，需要能够满足：

- 服务端尽可能多的处理并发请求
- 同时尽可能短的处理完毕。

Socket I/O 可以看做是二者之间的桥梁，如何更好地协调二者，去满足前面说的两点要求，有一些模式 (pattern) 是可以应用的。RPC 框架可选择的 I/O 模型严格意义上有 5 种，这里不讨论基于信号驱动的 I/O (Signal Driven I/O)。它们分别是：

- 传统的阻塞 I/O (Blocking I/O)
- 非阻塞 I/O (Non-blocking I/O)
- I/O 多路复用 (I/O multiplexing)
- 异步 I/O (Asynchronous I/O)

5.4.3 参考如下:

理解 REST 和 RPC

RPC

Go 语言的 RPC 框架有两个比较有特色的设计:

- 一个是 RPC 数据打包时可以通过插件实现自定义的编码和解码;
- 另一个是 RPC 建立在抽象的 `io.ReadWriteCloser` 接口之上的

5.5 Protobuf3

5.5.1 [转]Protobuf3 语法指南(<https://colobu.com/2017/03/16/Protobuf3-language-guide/>)

语言指南 [译]Protobuf 语法指南 中文出处是 proto2 的译文, proto3 的英文出现后在原来基础上增改了, 水平有限, 还请指正

这个指南描述了如何使用 Protocol buffer 语言去描述你的 protocol buffer 数据, 包括 .proto 文件符号和如何从 .proto 文件生成类。包含了 proto2 版本的 protocol buffer 语言: 对于老版本的 proto3 符号, 请见 [Proto2 Language Guide](#) (以及中文译本, 抄了很多这里的感谢下老版本的翻译者)

本文是一个参考指南——如果要查看如何使用本文中描述的多个特性的循序渐进的例子, 请在教程中查找需要的语言的教程。

5.5.2 定义一个消息类型

先来看一个非常简单的例子。假设你想定义一个“搜索请求”的消息格式, 每一个请求含有一个查询字符串、你感兴趣的查询结果所在的页数, 以及每一页多少条查询结果。可以采用如下的方式来定义消息类型的 .proto 文件了:

```
syntax = "proto3";
message SearchRequest {
  string query = 1;
  int32 page_number = 2;
  int32 result_per_page = 3;
}
```

- 文件的第一行指定了你正在使用 proto3 语法: 如果你没有指定这个, 编译器会使用 proto2。这个指定语法行必须是文件的非空非注释的第一个行。
- SearchRequest 消息格式有 3 个字段, 在消息中承载的数据分别对应于每一个字段。其中每个字段都有一个名字和一种类型。

指定字段类型

在上面的例子中，所有字段都是标量类型：两个整型（`page_number` 和 `result_per_page`），一个 `string` 类型（`query`）。当然，你也可以为字段指定其他的合成类型，包括枚举（enumerations）或其他消息类型。

分配标识号

正如你所见，在消息定义中，每个字段都有唯一的一个数字标识符。这些标识符是用来在消息的二进制格式中识别各个字段的，一旦开始使用就不能够再改变。注：[1,15] 之内的标识号在编码的时候会占用一个字节。[16,2047] 之内的标识号则占用 2 个字节。所以应该为那些频繁出现的消息元素保留 [1,15] 之内的标识号。切记：要为将来有可能添加的、频繁出现的标识号预留一些标识号。

最小的标识号可以从 1 开始，最大到 $2^{29} - 1$, or 536,870,911。不可以使用其中的 [19000 - 19999]（从 `FieldDescriptor::kFirstReservedNumber` 到 `FieldDescriptor::kLastReservedNumber`）的标识号，Protobuf 协议实现中对这些进行了预留。如果非要在 .proto 文件中使用这些预留标识号，编译时就会报警。同样你也不能使用早期保留的标识号。

指定字段规则

所指定的消息字段修饰符必须是如下之一：

- `singular`：一个格式良好的消息应该有 0 个或者 1 个这种字段（但是不能超过 1 个）。
- `repeated`：在一个格式良好的消息中，这种字段可以重复任意多次（包括 0 次）。重复的值的顺序会被保留。

在 proto3 中，`repeated` 的标量域默认情况下使用 `packed`。

你可以了解更多的 `packed` 属性在 Protocol Buffer 编码

添加更多消息类型

在一个 .proto 文件中可以定义多个消息类型。在定义多个相关的消息的时候，这一点特别有用——例如，如果想定义与 `SearchResponse` 消息类型对应的回复消息格式的话，你可以将它添加到相同的 .proto 文件中，如：

```
message SearchRequest {
  string query = 1;
  int32 page_number = 2;
  int32 result_per_page = 3;
}
message SearchResponse {
  ...
}
```

添加注释

向 .proto 文件添加注释，可以使用 C/C++/Java 风格的双斜杠 (`//`) 语法格式，如：

```
message SearchRequest {
  string query = 1;
  int32 page_number = 2; // Which page number do we want?
  int32 result_per_page = 3; // Number of results to return per page.
}
```

保留标识符 (Reserved)

如果你通过删除或者注释所有域，以后的用户在更新这个类型的时候可能重用这些标识号。如果你使用旧版本加载相同的.proto 文件会导致严重的问题，包括数据损坏、隐私错误等等。现在有一种确保不会发生这种情况的方法就是为字段 tag (reserved name 可能会 JSON 序列化的问题) 指定 reserved 标识符，protocol buffer 的编译器会警告未来尝试使用这些域标识符的用户。

```
message Foo {
  reserved 2, 15, 9 to 11;
  reserved "foo", "bar";
}
```

注：不要在同一行 reserved 声明中同时声明域名字和 tag number。

从.proto 文件生成了什么？

当用 protocol buffer 编译器来运行.proto 文件时，编译器将生成所选择语言的代码，这些代码可以操作在.proto 文件中定义的消息类型，包括获取、设置字段值，将消息序列化到一个输出流中，以及从一个输入流中解析消息。

- 对 C++ 来说，编译器会为每个.proto 文件生成一个.h 文件和一个.cc 文件，.proto 文件中的每一个消息有一个对应的类。
- 对 Java 来说，编译器为每一个消息类型生成了一个.java 文件，以及一个特殊的 Builder 类（该类是用来创建消息类接口的）。
- 对 Python 来说，有点不太一样——Python 编译器为.proto 文件中的每个消息类型生成一个含有静态描述符的模块，该模块与一个元类（metaclass）在运行时（runtime）被用来创建所需的 Python 数据访问类。
- 对 go 来说，编译器会为每个消息类型生成了一个.pb.go 文件。
- 对于 Ruby 来说，编译器会为每个消息类型生成了一个.rb 文件。
- javaNano 来说，编译器输出类似域 java 但是没有 Builder 类
- 对于 Objective-C 来说，编译器会为每个消息类型生成了一个 pbobjc.h 文件和 pbobjcm 文件，.proto 文件中的每一个消息有一个对应的类。
- 对于 C# 来说，编译器会为每个消息类型生成了一个.cs 文件，.proto 文件中的每一个消息有一个对应的类。你可以从如下的文档链接中获取每种语言更多 API(proto3 版本的内容很快就公布)。API Reference

5.5.3 标量数值类型

一个标量消息字段可以含有一个如下的类型——该表格展示了定义于.proto 文件中的类型，以及与之对应的、在自动生成的访问类中定义的类型：

你可以在文章Protocol Buffer 编码中，找到更多“序列化消息时各种类型如何编码”的信息。

1. 在 java 中，无符号 32 位和 64 位整型被表示成他们的整型对应形式，最高位被储存在标志位中。
2. 对于所有的情况，设定值会执行类型检查以确保此值是有效。
3. 64 位或者无符号 32 位整型在解码时被表示成为 ilong，但是在设置时可以使用 int 型值设定，在所有的情况下，值必须符合其设置其类型的要求。
4. python 中 string 被表示成在解码时表示成 unicode。但是一个 ASCIIstring 可以被表示成 str 类型。
5. Integer 在 64 位的机器上使用，string 在 32 位机器上使用

5.5.4 默认值

当一个消息被解析的时候，如果被编码的信息不包含一个特定的 `singular` 元素，被解析的对象锁对应的域被设置位一个默认值，对于不同类型指定如下：

- 对于 `string`，默认是一个空 `string`
- 对于 `bytes`，默认是一个空的 `bytes`
- 对于 `bool`，默认是 `false`
- 对于数值类型，默认是 `0`
- 对于枚举，默认是第一个定义的枚举值，必须为 `0`；
- 对于消息类型（`message`），域没有被设置，确切的消息是根据语言确定的，详见 `generated code guide`

对于可重复域的默认值是空（通常情况下是对应语言中空列表）。

注：对于标量消息域，一旦消息被解析，就无法判断域释放被设置为默认值（例如，例如 `boolean` 值是否被设置为 `false`）还是根本没有被设置。你应该在定义你的消息类型时非常注意。例如，比如你不应该定义 `boolean` 的默认值 `false` 作为任何行为的触发方式。也应该注意如果一个标量消息域被设置为标志位，这个值不应该被序列化传输。

查看 `generated code guide` 选择你的语言的默认值的工作细节。

5.5.5 枚举

当需要定义一个消息类型的时候，可能想为一个字段指定某“预定义值序列”中的一个值。例如，假设要为每一个 `SearchRequest` 消息添加一个 `corpus` 字段，而 `corpus` 的值可能是 `UNIVERSAL`, `WEB`, `IMAGES`, `LOCAL`, `NEWS`, `PRODUCTS` 或 `VIDEO` 中的一个。其实可以很容易地实现这一点：通过向消息定义中添加一个枚举（`enum`）并且为每个可能的值定义一个常量就可以了。

在下面的例子中，在消息格式中添加了一个叫做 `Corpus` 的枚举类型——它含有所有可能的值——以及一个类型为 `Corpus` 的字段：

```
message SearchRequest {
  string query = 1;
  int32 page_number = 2;
  int32 result_per_page = 3;
  enum Corpus {
    UNIVERSAL = 0;
    WEB = 1;
    IMAGES = 2;
    LOCAL = 3;
    NEWS = 4;
    PRODUCTS = 5;
    VIDEO = 6;
  }
  Corpus corpus = 4;
}
```

如你所见，`Corpus` 枚举的第一个常量映射为 `0`：每个枚举类型必须将其第一个类型映射为 `0`，这是因为：

- 必须有一个 `0` 值，我们可以用这个 `0` 值作为默认值。
- 这个零值必须为第一个元素，为了兼容 `proto2` 语义，枚举类的第一个值总是默认值。

你可以通过将不同的枚举常量指定相同的值。如果这样做你需要将 `allow_alias` 设定位 `true`，否则编译器会在别名的地方产生一个错误信息。

```

enum EnumAllowingAlias {
  option allow_alias = true;
  UNKNOWN = 0;
  STARTED = 1;
  RUNNING = 1;
}
enum EnumNotAllowingAlias {
  UNKNOWN = 0;
  STARTED = 1;
  // RUNNING = 1; // Uncommenting this line will cause a compile error inside Google_
  ↪and a warning message outside.
}

```

枚举常量必须在 32 位整型值的范围内。因为 enum 值是使用可变编码方式的，对负数不够高效，因此不推荐在 enum 中使用负数。如上例所示，可以在一个消息定义的内部或外部定义枚举——这些枚举可以在 .proto 文件中的任何消息定义里重用。当然也可以在一个消息中声明一个枚举类型，而在另一个不同的消息中使用它——采用 `MessageType.EnumType` 的语法格式。

当对一个使用了枚举的 .proto 文件运行 protocol buffer 编译器的时候，生成的代码中将有一个对应的 enum（对 Java 或 C++ 来说），或者一个特殊的 EnumDescriptor 类（对 Python 来说），它被用来在运行时生成的类中创建一系列的整型值符号常量（symbolic constants）。

在反序列化的过程中，无法识别的枚举值会被保存在消息中，虽然这种表示方式需要依据所使用语言而定。在那些支持开放枚举类型超出指定范围之外的语言中（例如 C++ 和 Go），为识别的值会被表示成所支持的整型。在使用封闭枚举类型的语言中（Java），使用枚举中的一个类型来表示未识别的值，并且可以使用所支持整型来访问。在其他情况下，如果解析的消息被序列号，未识别的值将保持原样。

关于如何在你的应用程序的消息中使用枚举的更多信息，请查看所选择的语言 [generated code guide](#)。

5.5.6 使用其他消息类型

你可以将其他消息类型用作字段类型。例如，假设在每一个 SearchResponse 消息中包含 Result 消息，此时可以在相同的 .proto 文件中定义一个 Result 消息类型，然后在 SearchResponse 消息中指定一个 Result 类型的字段，如：

```

message SearchResponse {
  repeated Result results = 1;
}
message Result {
  string url = 1;
  string title = 2;
  repeated string snippets = 3;
}

```

导入定义

在上面的例子中，Result 消息类型与 SearchResponse 是定义在同一文件中的。如果想要使用的消息类型已经在其他 .proto 文件中已经定义过了呢？你可以通过导入（importing）其他 .proto 文件中的定义来使用它们。要导入其他 .proto 文件的定义，你需要在你的文件中添加一个导入声明，如：

```
import "myproject/other_protos.proto";
```

默认情况下你只能使用直接导入的 .proto 文件中的定义。然而，有时候你需要移动一个 .proto 文件到一个新的位置，可以不直接移动 .proto 文件，只需放入一个伪 .proto 文件在老的位置，然后使用 `import public` 转向新的位置。`import public` 依赖性会通过任意导入包含 `import public` 声明的 proto 文件传递。例如：

```
// 这是新的proto
// All definitions are moved here

// 这是旧的proto
// 这是所有客户端正在导入的包
import public "new.proto";
import "other.proto";

// 客户端proto
import "old.proto";
// 现在你可以使用新旧两种包的proto定义了。
```

通过在编译器命令行参数中使用 `-I/--proto_path` 编译器会在指定目录搜索要导入的文件。如果没有给出标志，编译器会搜索编译命令被调用的目录。通常你只要指定 `proto_path` 标志为你的工程根目录就好。并且指定好导入的正确名称就好。

使用 proto2 消息类型

在你的 proto3 消息中导入 proto2 的消息类型也是可以的，反之亦然，然后 proto2 枚举不可以在 proto3 的标识符中使用（如果仅仅在 proto2 消息中使用是可以的）。

5.5.7 嵌套类型

你可以在其他消息类型中定义、使用消息类型，在下面的例子中，Result 消息就定义在 SearchResponse 消息内，如：

```
message SearchResponse {
  message Result {
    string url = 1;
    string title = 2;
    repeated string snippets = 3;
  }
  repeated Result results = 1;
}
```

如果你想在它的父消息类型的外部重用这个消息类型，你需要以 `Parent.Type` 的形式使用它，如：

```
message SomeOtherMessage {
  SearchResponse.Result result = 1;
}
```

当然，你也可以将消息嵌套任意多层，如：

```
message Outer { // Level 0
  message MiddleAA { // Level 1
    message Inner { // Level 2
      int64 ival = 1;
      bool  booly = 2;
    }
  }
  message MiddleBB { // Level 1
    message Inner { // Level 2
      int32 ival = 1;
      bool  booly = 2;
    }
  }
}
```

(下页继续)

(续上页)

```

    }
  }
}

```

5.5.8 更新一个消息类型

如果一个已有的消息格式已无法满足新的需求——如，要在消息中添加一个额外的字段——但是同时旧版本写的代码仍然可用。不用担心！更新消息而不破坏已有代码是非常简单的。在更新时只要记住以下的规则即可。

- 不要更改任何已有的字段的数值标识。
- 如果你增加新的字段，使用旧格式的字段仍然可以被你新产生的代码所解析。你应该记住这些元素的默认值这样你的新代码就可以以适当的方式和旧代码产生的数据交互。相似的，通过新代码产生的消息也可以被旧代码解析：只不过新的字段会被忽视掉。注意，未被识别的字段会在反序列化的过程中丢弃掉，所以如果消息再被传递给新的代码，新的字段依然是不可用的（这和 proto2 中的行为是不同的，在 proto2 中未定义的域依然会随着消息被序列化）
- 非 required 的字段可以移除——只要它们的标识号在新的消息类型中不再使用（更好的做法可能是重命名那个字段，例如在字段前添加“OBSOLETE_”前缀，那样的话，使用的.proto 文件的用户将来就不会无意中重新使用了那些不该使用的标识号）。
- int32, uint32, int64, uint64, 和 bool 是全部兼容的，这意味着可以将这些类型中的一个转换为另外一个，而不会破坏向前、向后的兼容性。如果解析出来的数字与对应的类型不相符，那么结果就像在 C++ 中对它进行了强制类型转换一样（例如，如果把一个 64 位数字当作 int32 来读取，那么它就会被截断为 32 位的数字）。
- sint32 和 sint64 是互相兼容的，但是它们与其他整数类型不兼容。
- string 和 bytes 是兼容的——只要 bytes 是有效的 UTF-8 编码。
- 嵌套消息与 bytes 是兼容的——只要 bytes 包含该消息的一个编码过的版本。
- fixed32 与 sfixed32 是兼容的，fixed64 与 sfixed64 是兼容的。
- 枚举类型与 int32, uint32, int64 和 uint64 相兼容（注意如果值不相兼容则会被截断），然而在客户端反序列化之后他们可能会有不同的处理方式，例如，未识别的 proto3 枚举类型会被保留在消息中，但是他的表示方式会依照语言而定。int 类型的字段总会保留他们的

5.5.9 Any

Any 类型消息允许你在没有指定他们的.proto 定义的情况下使用消息作为一个嵌套类型。一个 Any 类型包括一个可以被序列化 bytes 类型的任意消息，以及一个 URL 作为一个全局标识符和解析消息类型。为了使用 Any 类型，你需要导入 `import google/protobuf/any.proto`。

```

import "google/protobuf/any.proto";
message ErrorStatus {
  string message = 1;
  repeated google.protobuf.Any details = 2;
}

```

对于给定的消息类型的默认类型 URL 是 `type.googleapis.com/packageName.messageName`。

不同语言的实现会支持动态库以线程安全的方式去帮助封装或者解封装 Any 值。例如在 java 中，Any 类型会有特殊的 `pack()` 和 `unpack()` 访问器，在 C++ 中会有 `PackFrom()` 和 `UnpackTo()` 方法。

```
// Storing an arbitrary message type in Any.
NetworkErrorDetails details = ...;
ErrorStatus status;
status.add_details()->PackFrom(details);
// Reading an arbitrary message from Any.
ErrorStatus status = ...;
for (const Any& detail : status.details()) {
    if (detail.Is<NetworkErrorDetails>()) {
        NetworkErrorDetails network_error;
        detail.UnpackTo(&network_error);
        ... processing network_error ...
    }
}
```

目前，用于 Any 类型的动态库仍在开发之中如果你已经很熟悉proto2 语法，使用 Any 替换扩展。

5.5.10 Oneof

如果你的消息中有很多可选字段，并且同时至多一个字段会被设置，你可以加强这个行为，使用 oneof 特性节省内存。

Oneof 字段就像可选字段，除了它们会共享内存，至多一个字段会被设置。设置其中一个字段会清除其它字段。你可以使用 case() 或者 WhichOneof() 方法检查哪个 oneof 字段被设置，看你使用什么语言了。

使用 Oneof

为了在.proto 定义 Oneof 字段，你需要在名字前面加上 oneof 关键字，比如下面例子的 test_oneof:

```
message SampleMessage {
    oneof test_oneof {
        string name = 4;
        SubMessage sub_message = 9;
    }
}
```

然后你可以增加 oneof 字段到 oneof 定义中. 你可以增加任意类型的字段, 但是不能使用 repeated 关键字.

在产生的代码中, oneof 字段拥有同样的 getters 和 setters, 就像正常的可选字段一样. 也有一个特殊的方法来检查到底那个字段被设置. 你可以在相应的语言API 指南中找到 oneof API 介绍.

Oneof 特性

- 设置 oneof 会自动清楚其它 oneof 字段的值. 所以设置多次后, 只有最后一次设置的字段有值.

```
SampleMessage message;
message.set_name("name");
CHECK(message.has_name());
message.mutable_sub_message(); // Will clear name field.
CHECK(!message.has_name());
```

- 如果解析器遇到同一个 oneof 中有多个成员, 只有最会一个会被解析成消息.
- oneof 不支持 repeated.
- 反射 API 对 oneof 字段有效.

- 如果使用 C++, 需确保代码不会导致内存泄漏. 下面的代码会崩溃, 因为 sub_message 已经通过 set_name() 删除了

```
SampleMessage message;
SubMessage* sub_message = message.mutable_sub_message();
message.set_name("name");           // Will delete sub_message
sub_message->set_...                 // Crashes here
```

- 在 C++ 中, 如果你使用 Swap() 两个 oneof 消息, 每个消息, 两个消息将拥有对方的值, 例如在下面的例子中, msg1 会拥有 sub_message 并且 msg2 会有 name。

```
SampleMessage msg1;
msg1.set_name("name");
SampleMessage msg2;
msg2.mutable_sub_message();
msg1.swap(&msg2);
CHECK(msg1.has_sub_message());
CHECK(msg2.has_name());
```

向后兼容性问题

当增加或者删除 oneof 字段时一定要小心. 如果检查 oneof 的值返回 None/NOT_SET, 它意味着 oneof 字段没有被赋值或者在一个不同的版本中赋值了. 你不会知道是哪种情况, 因为没有办法判断如果未识别的字段是一个 oneof 字段。

Tag 重用问题:

- **将字段移入或移除 oneof**: 在消息被序列化或者解析后, 你也许会失去一些信息 (有些字段也许会被清除)
- **删除一个字段或者加入一个字段**: 在消息被序列化或者解析后, 这也许会清除你现在设置的 oneof 字段
- **分离或者融合 oneof**: 行为与移动常规字段相似。

5.5.11 Map

如果你希望创建一个关联映射, protocol buffer 提供了一种快捷的语法:

```
map<key_type, value_type> map_field = N;
```

其中 key_type 可以是任意 Integer 或者 string 类型 (所以, 除了 floating 和 bytes 的任意标量类型都是可以的) value_type 可以是任意类型。

例如, 如果你希望创建一个 project 的映射, 每个 Project 使用一个 string 作为 key, 你可以像下面这样定义:

```
map<string, Project> projects = 3;
```

- Map 的字段可以是 repeated。
- 序列化后的顺序和 map 迭代器的顺序是不确定的, 所以你不要期望以固定顺序处理 Map
- 当为 .proto 文件产生生成文本格式的时候, map 会按照 key 的顺序排序, 数值化的 key 会按照数值排序。
- 从序列化中解析或者融合时, 如果有重复的 key 则后一个 key 不会被使用, 当从文本格式中解析 map 时, 如果存在重复的 key。

生成 map 的 API 现在对于所有 proto3 支持的语言都可用了, 你可以从 [API 指南](#) 找到更多信息。

向后兼容性问题

map 语法序列化后等同于如下内容，因此即使是不支持 map 语法的 protocol buffer 实现也是可以处理你的数据的：

```
message MapFieldEntry {
  key_type key = 1;
  value_type value = 2;
}
repeated MapFieldEntry map_field = N;
```

5.5.12 Package

当然可以为 .proto 文件新增一个可选的 package 声明符，用来防止不同的消息类型有命名冲突。如：

```
package foo.bar;
message Open { ... }
```

在其他的消息格式定义中可以使用包名 + 消息名的方式来定义域的类型，如：

```
message Foo {
  ...
  required foo.bar.Open open = 1;
  ...
}
```

包的声明符会根据使用语言的不同影响生成的代码。

- 对于 C++，产生的类会被包装在 C++ 的命名空间中，如上例中的 Open 会被封装在 foo::bar 空间中；- 对于 Java，包声明符会变为 java 的一个包，除非在 .proto 文件中提供了一个明确有 java_package；
- 对于 Python，这个包声明符是被忽略的，因为 Python 模块是按照其在文件系统中的位置进行组织的。
- 对于 Go，包可以被用做 Go 包名称，除非你显式的提供一个 option go_package 在你的 .proto 文件中。
- 对于 Ruby，生成的类可以被包装在内置的 Ruby 名称空间中，转换成 Ruby 所需的大小写样式（首字母大写；如果第一个符号不是一个字母，则使用 PB_ 前缀），例如 Open 会在 Foo::Bar 名称空间中。
- 对于 javaNano 包会使用 Java 包，除非你在你的文件中显式的提供一个 option java_package。
- 对于 C# 包可以转换为 PascalCase 后作为名称空间，除非你在你的文件中显式的提供一个 option csharp_namespace，例如，Open 会在 Foo.Bar 名称空间中

包及名称的解析

Protocol buffer 语言中类型名称的解析与 C++ 是一致的：首先从最内部开始查找，依次向外进行，每个包会被看作是其父类包的内部类。当然对于 (foo.bar.Baz) 这样以 “.” 分隔的意味着是从最外围开始的。

ProtocolBuffer 编译器会解析 .proto 文件中定义的所有类型名。对于不同语言的代码生成器会知道如何来指向每个具体的类型，即使它们使用了不同的规则。

5.5.13 定义服务 (Service)

如果想要将消息类型用在 RPC(远程方法调用) 系统中, 可以在 .proto 文件中定义一个 RPC 服务接口, protocol buffer 编译器将会根据所选择的不同语言生成服务接口代码及存根。如, 想要定义一个 RPC 服务并具有一个方法, 该方法能够接收 SearchRequest 并返回一个 SearchResponse, 此时可以在 .proto 文件中进行如下定义:

```
service SearchService {
  rpc Search (SearchRequest) returns (SearchResponse);
}
```

最直观的使用 protocol buffer 的 RPC 系统是 gRPC, 一个由谷歌开发的语言和平台中的开源的 PRC 系统, gRPC 在使用 protocol buffer 时非常有效, 如果使用特殊的 protocol buffer 插件可以直接为您从 .proto 文件中产生相关的 RPC 代码。

如果你不想使用 gRPC, 也可以使用 protocol buffer 用于自己的 RPC 实现, 你可以从 [proto2 语言指南](#) 中找到更多信息

还有一些第三方开发的 PRC 实现使用 Protocol Buffer。参考 [第三方插件 wiki](#) 查看这些实现的列表。

5.5.14 JSON 映射

Proto3 支持 JSON 的编码规范, 使他更容易在不同系统之间共享数据, 在下表中逐个描述类型。

如果 JSON 编码的数据丢失或者其本身就是 null, 这个数据会在解析成 protocol buffer 的时候被表示成默认值。如果一个字段在 protocol buffer 中表示为默认值, 体会在转化成 JSON 的时候编码的时候忽略掉以节省空间。具体实现可以提供在 JSON 编码中可选的默认值。

5.5.15 选项

定义 .proto 文件时能够标注一系列的 option。Option 并不改变整个文件声明的含义, 但却能够影响特定环境下处理方式。完整的可用选项可以在 [google/protobuf/descriptor.proto](#) 找到。

一些选项是文件级别的, 意味着它可以作用于最外范围, 不包含在任何消息内部、enum 或服务定义中。一些选项是消息级别的, 意味着它可以用在消息定义的内部。当然有些选项可以作用在域、enum 类型、enum 值、服务类型及服务方法中。到目前为止, 并没有一种有效的选项能作用于所有的类型。

如下就是一些常用的选项:

- `java_package` (文件选项): 这个选项表明生成 java 类所在的包。如果在 .proto 文件中没有明确的声明 `java_package`, 就采用默认的包名。当然了, 默认方式产生的 java 包名并不是最好的方式, 按照应用名称倒序方式进行排序的。如果不需要产生 java 代码, 则该选项将不起任何作用。如:

```
option java_package = "com.example.foo";
```

- `java_outer_classname` (文件选项): 该选项表明想要生成 Java 类的名称。如果在 .proto 文件中没有明确的 `java_outer_classname` 定义, 生成的 class 名称将会根据 .proto 文件的名称采用驼峰式的命名方式进行生成。如 (`foo_bar.proto` 生成的 java 类名为 `FooBar.java`), 如果不生成 java 代码, 则该选项不起任何作用。如:

```
option java_outer_classname = "Ponycopter";
```

- `optimize_for` (文件选项): 可以被设置为 `SPEED`, `CODE_SIZE`, 或者 `LITE_RUNTIME`。这些值将通过如下的方式影响 C++ 及 java 代码的生成:
 - `SPEED` (default): protocol buffer 编译器将通过在消息类型上执行序列化、语法分析及其他通用的操作。这种代码是最优的。

- `CODE_SIZE`: protocol buffer 编译器将会产生最少量的类, 通过共享或基于反射的代码来实现序列化、语法分析及各种其它操作。采用该方式产生的代码将比 `SPEED` 要少得多, 但是操作要相对慢些。当然实现的类及其对外的 API 与 `SPEED` 模式都是一样的。这种方式经常用在一些包含大量的 .proto 文件而且并不盲目追求速度的应用中。
- `LITE_RUNTIME`: protocol buffer 编译器依赖于运行时核心类库来生成代码 (即采用 `libprotobuf-lite` 替代 `libprotobuf`)。这种核心类库由于忽略了一些描述符及反射, 要比全类库小得多。这种模式经常在移动手机平台应用多一些。编译器采用该模式产生的方法实现与 `SPEED` 模式不相上下, 产生的类通过实现 `MessageLite` 接口, 但它仅仅是 `Message` 接口的一个子集。

```
option optimize_for = CODE_SIZE;
```

- `cc_enable_arenas`(文件选项): 对于 C++ 产生的代码启用 arena allocation
- `objc_class_prefix`(文件选项): 设置 Objective-C 类的前缀, 添加到所有 Objective-C 从此 .proto 文件产生的类和枚举类型。没有默认值, 所使用的前缀应该是苹果推荐的 3-5 个大写字母, 注意 2 个字节的前缀是苹果所保留的。
- `deprecated`(字段选项): 如果设置为 `true` 则表示该字段已经被废弃, 并且不应该在新的代码中使用。在大多数语言中没有实际的意义。在 java 中, 这回变成 `@Deprecated` 注释, 在未来, 其他语言的代码生成器也许会在字标识符中产生废弃注释, 废弃注释会在编译器尝试使用该字段时发出警告。如果字段没有被使用你也不希望有新用户使用它, 尝试使用保留语句替换字段声明。

```
int32 old_field = 6 [deprecated=true];
```

自定义选项

ProtocolBuffers 允许自定义并使用选项。该功能应该属于一个高级特性, 对于大部分人是用不到的。如果你的确希望创建自己的选项, 请参看 [Proto2 Language Guide](#)。注意创建自定义选项使用了拓展, 拓展只在 `proto3` 中可用。

5.5.16 生成访问类

可以通过定义好的 .proto 文件来生成 Java, Python, C++, Ruby, JavaNano, Objective-C, 或者 C# 代码, 需要基于 .proto 文件运行 protocol buffer 编译器 `protoc`。如果你没有安装编译器, 下载安装包并遵照 `README` 安装。对于 Go, 你还需要安装一个特殊的代码生成器插件。你可以通过 [GitHub](#) 上的 `protobuf` 库找到安装过程

通过如下方式调用 protocol 编译器:

```
protoc --proto_path=IMPORT_PATH --cpp_out=DST_DIR --java_out=DST_DIR --python_out=DST_
↳DIR --go_out=DST_DIR --ruby_out=DST_DIR --javanano_out=DST_DIR --objc_out=DST_DIR --
↳csharp_out=DST_DIR path/to/file.proto
```

- `IMPORT_PATH` 声明了一个 .proto 文件所在的解析 `import` 具体目录。如果忽略该值, 则使用当前目录。如果有多个目录则可以多次调用 `--proto_path`, 它们将会顺序的被访问并执行导入。 `-I=IMPORT_PATH` 是 `--proto_path` 的简化形式。
- 当然也可以提供一个或多个输出路径:
 - `--cpp_out` 在目标目录 `DST_DIR` 中产生 C++ 代码, 可以在 [C++ 代码生成参考](#) 中查看更多。
 - `--java_out` 在目标目录 `DST_DIR` 中产生 Java 代码, 可以在 [Java 代码生成参考](#) 中查看更多。
 - `--python_out` 在目标目录 `DST_DIR` 中产生 Python 代码, 可以在 [Python 代码生成参考](#) 中查看更多。
 - `--go_out` 在目标目录 `DST_DIR` 中产生 Go 代码, 可以在 [GO 代码生成参考](#) 中查看更多。
 - `--ruby_out` 在目标目录 `DST_DIR` 中产生 Ruby 代码, 参考正在制作中。

- `--javanano_out` 在目标目录 `DST_DIR` 中生成 `JavaNano`，`JavaNano` 代码生成器有一系列的选项用于定制自定义生成器的输出：你可以通过生成器的[README](#)查找更多信息，`JavaNano` 参考正在制作中。
- `--objc_out` 在目标目录 `DST_DIR` 中产生 `Object` 代码，可以在[Objective-C 代码生成参考](#)中查看更多。
- `--csharp_out` 在目标目录 `DST_DIR` 中产生 `Object` 代码，可以在[C# 代码生成参考](#)中查看更多。
- `--php_out` 在目标目录 `DST_DIR` 中产生 `Object` 代码，可以在[PHP 代码生成参考](#)中查看更多。

作为一个方便的拓展，如果 `DST_DIR` 以 `.zip` 或者 `.jar` 结尾，编译器会将输出写到一个 `ZIP` 格式文件或者符合 `JAR` 标准的 `.jar` 文件中。注意如果输出已经存在则会被覆盖，编译器还没有智能到可以追加文件。

- 你必须提议一个或多个 `.proto` 文件作为输入，多个 `.proto` 文件可以只指定一次。虽然文件路径是相对于当前目录的，每个文件必须位于其 `IMPORT_PATH` 下，以便每个文件可以确定其规范的名称。

6.1 Docker

6.1.1 Docker 基础

Docker 镜像: image

Docker 镜像是一个特殊的文件系统，它提供容器运行时所需的程序、库、资源、配置等文件，还包含了一些为运行时准备的一些配置参数（如匿名卷、环境变量、用户等）。镜像中不包含任何动态数据，其内容在构建之后也不会被改变，Docker 引擎基于 Docker 镜像、提供运行时所需的配置参数就可以运行一个 Docker 容器。Docker 镜像是一个虚拟的概念，设计时采用了分层存储的架构，由多层文件系统（如基础镜像）联合组成。

镜像构建时，会一层层构建，前一层是后一层的基础。每一层构建完就不会再发生改变，后一层上的任何改变只发生在自己这一层。比如，删除前一层文件的操作，实际不是真的删除前一层的文件，而是仅在当前层标记为该文件已删除。在最终容器运行的时候，虽然不会看到这个文件，但是实际上该文件会一直跟随镜像。因此，在构建镜像的时候，需要额外小心，每一层尽量只包含该层需要添加的东西，任何额外的东西应该在该层构建结束前清理掉。

分层存储的特征还使得镜像的复用、定制变的更为容易。甚至可以用之前构建好的镜像作为基础层，然后进一步添加新的层，以定制自己所需的内容，构建新的镜像。

容器运行时，会在镜像文件系统层上方再添加一个读写层，容器运行时在容器读写层进行文件读写操作，当容器停止时，docker 引擎会删除容器读写层。

查看一个镜像的层次结构，可通过 `docker history <image-id>` 查看：

容器: container

容器是镜像运行时的实体 (镜像实例化)。容器可以被创建、启动、停止、删除、暂停等。

容器的实质是进程，但与直接在宿主执行的进程不同，容器进程运行于自己的独立的命名空间。因此容器可以拥有自己的 `root` 文件系统、自己的网络配置、自己的进程空间，甚至自己的用户 ID 空间。

每一个容器运行时，是以镜像为基础层，在其上创建一个当前容器的存储层。容器存储层的生存周期和容器一样，因此，容器不应该向其存储层内写入任何数据，容器存储层要保持无状态化。所有的文件写入操作，都应该使用数据卷 (Volume)、或者绑定宿主目录，在这些位置的读写会跳过容器存储层，直接对宿主（或网络存储）发生读写，其性能和稳定性更高。

常用命令

```
# 显示本机镜像
docker images
docker image ls
# 编译镜像
docker build -t cfp/node
# 拉取镜像
docekr pull cfp/node-1:latest
# 推送镜像
docekr push cfp/node-1:latest
# 删除本地镜像
docker rmi cfp/node-1:latest
# 镜像tag
docker tag cfp/node-1:latest cfp/node-1:v1.0
# 显示本地容器
docker ps -a
# 删除本地容器
docker rm full10
# 进入容器
docker exec -it full10 /bin/sh
# 进入容器
docker attach full10
# 启动、停止、重启容器
docker start/stop/restart
# 查看容器log
docker logs -f full10
# 启动容器
docker run cfp/node ls /go/src/github.com/cosmos/launch
# 使用容器生成镜像
docker commit -a "author" -m "message" full10 image:new
```

- docker exec 与 attache 的区别
 - docker exec -it 分配伪终端和 stdin，跟正常的 console 一样执行命令
 - docker attach attach 到一个已经运行的容器的 stdin，然后进行命令执行。但是，如果从这个 stdin 中 exit(如 `ctl+c`)，会导致容器的停止
- docker commit 生成镜像
 - 将容器的存储层保存下来成为镜像的一层。就是在原有镜像的基础上，再叠加容器的存储层，并构成新的镜像。以后我们运行这个新镜像的时候，就会拥有原有容器最后的文件变化。
 - docker commit 生成镜像时黑箱操作，除了制作镜像的人知道执行过什么命令、怎么生成的镜像，别人根本无从得知。难用且不易维护。

- docker commit 生成镜像易包含其他无用修改，使镜像臃肿
- 通常不建议使用此方式构建镜像，推荐使用 Dockerfile
- 可用于特殊场景，如被入侵后保存现场等

6.1.2 Docker 应用

docker 运行 mysql

- 获取镜像

```
# 从 docker hub 的仓库中拉取 mysql 镜像
sudo docker pull mysql
```

- 运行一个 mysql 实例

```
docker run --name fp-mysql -p 23306:3306 -e MYSQL_ROOT_PASSWORD=123456 -d mysql:5.
↪6

5b6bf6f629bfe46b4c8786b555d8db1947680138b2de1f268f310a15ced7247a

run          运行一个容器
--name       后面是这个镜像的名称
-p 3306:3306 ↪
↪表示在这个容器中使用3306端口(第二个)映射到本机的端口号也为3306(第一个)
-d          表示使用守护进程运行，即服务挂在后台
```

- 装一个 mysql-client，访问 mysql 数据库

```
sudo apt-get install mysql-client-core-5.6

# mysql 命令访问服务器
mysql -h 192.168.95.4 -P 3306 -uroot -p 123456
```

•

6.1.3 Docker file

Docker file 分析

Dockerfile 是一个文本文件，用于构建生成 Docker 镜像，其内以**一个镜像为基础**，在基础镜像之上通过一条指令 (Instruction) 加工、定制，最终产出自己想要的镜像。每一条指令都基于前一层依赖的镜像层构建一层新的镜像层，因此每一条指令的内容，就是描述该层应当如何构建。

```
FROM cfp/build as builder

# Set working directory for the build
WORKDIR $GOPATH/src/github.com/fpchan/cfp

# Add cfp source files
COPY . .

# build cfp
```

(下页继续)

```

RUN make install
RUN make launchcmd

# FROM cfp/node
FROM cfp/build

# Set working directory for the build
WORKDIR $GOPATH/src/github.com/cosmos/launch

ARG LAUNCH_REPOSITORY=https://github.com/cfp/launch
ARG LAUNCH_BRANCH=dev
ARG SEED_NODE_NUM=1
ARG VAL_NODE_NUM=1
ARG FULL_NODE_NUM=1

RUN mkdir -p $GOPATH/src/github.com/cosmos \
  && cd $GOPATH/src/github.com/cosmos \
  && git clone $LAUNCH_REPOSITORY -b $LAUNCH_BRANCH

COPY --from=builder $GOPATH/bin $GOPATH/bin
COPY docker/launch/genfile.sh .
COPY docker/launch/start.sh .

EXPOSE 26656 26657 26658 26659 6060

# generate genesis file
RUN $GOPATH/src/github.com/cosmos/launch/genfile.sh ${SEED_NODE_NUM} ${VAL_NODE_NUM}
  &→ ${FULL_NODE_NUM}
ENTRYPOINT $GOPATH/src/github.com/cosmos/launch/start.sh
docker build \
--build-arg SEED_NODE_NUM=1 \
--build-arg VAL_NODE_NUM=4 \
--build-arg FULL_NODE_NUM=10 \
-f Dockerfile -t cfp/node-1:rr .

```

- FROM 指定基础镜像
- COPY 从构建上下文目录中 <源路径> 的文件/目录复制到新的一层的镜像内的 <目标路径> 位置
 - 构建上下文目录是指 docker build 指定的目录，通常使用 . 即当前目录，操作文件不可以超出上下文目录
 - COPY [--chown=<user>:<group>] <源路径>... <目标路径>
 - COPY [--chown=<user>:<group>] ["<源路径 1>", ... "<目标路径>"]
- ADD 与 COPY 基本一致，但如果 <源路径> 是 URL 或者压缩文件，ADD 可以自动下载和解压
- RUN 指定要执行的命令
 - shell 格式: RUN <命令>
 - exec 格式: RUN ["可执行文件", "参数 1", "参数 2"]
- CMD 启动容器时执行的命令，格式同 RUN
 - 在运行时可以替代，docker run 中执行指定的指令或 docker-compose 中的 command
 - 在指令格式上，一般推荐使用 exec 格式，这类格式在解析时会被解析为 JSON 数组，因此一定要使用双引号 "，而不要使用单引号

- ENTRYPOINT 同 CMD 含义相同，格式同 RUN
 - 当同时执行 CMD 和 ENTRYPOINT 后，CMD 不再是直接的运行命令，二是将其作为参数传给 ENTRYPOINT，变成 `<ENTRYPOINT> "<CMD>"`
 - ENTRYPOINT 比 CMD 的优势：
 - * 如果指定 CMD ["curl", "-s", "https://ip.cn"]，执行 `docker run {images}` 会输出 IP，如果想加入参数 -i，则执行 `docker run {images} -i` 会报错
 - * 如果指定 ENTRYPOINT ["curl", "-s", "https://ip.cn"]，同样的执行 `docker run {images} -i` 则没问题，相当于执行 `curl -s https://ip.cn -i`
- ENV 环境变量
- ARG 构建参数
- EXPOSE 声明容器提供的服务端口
- WORKDIR 工作目录
- USER 指定当前用户

Docker file 最佳实践

- 使用 .dockerignore 文件
- 使用多阶段构建
 - 多阶段构建，允许在一个 Dockerfile 中使用多个 FROM 指令，可直接将前置阶段、或其他镜像中的资源文件复制到当前构建阶段，并最终只产生最后一个镜像，大大减少了镜像大小及个数
- 避免安装不必要的包
- 一个容器只运行一个进程
- 镜像层数尽可能少
- 构建缓存
 - 如果你不想在构建过程中使用缓存，你可以在 `docker build` 命令中使用 `--no-cache=true` 选项
- 将多个 RUN 指令合并为一个
 - Dockerfile 中的每个指令都会创建一个新的镜像层。
 - 镜像层将被缓存和复用
 - 当 Dockerfile 的指令修改了，复制的文件变化了，或者构建镜像时指定的变量不同了，对应的镜像层缓存就会失效
 - 某一层的镜像缓存失效之后，它之后的镜像层缓存都会失效
 - 镜像层是不可变的，如果我们再某一层中添加一个文件，然后在下一层中删除它，则镜像中依然会包含该文件 (只是这个文件在 Docker 容器中不可见了)
- 基础镜像的标签不要用 latest
 - 当镜像更新时，latest 标签会指向不同的镜像，这时构建镜像有可能失败
- 每个 RUN 指令后删除多余文件
- 选择合适的基础镜像 (alpine 版本最好)
 - alpine 是一个极小化的 Linux 发行版，只有 4MB，这让它非常适合作为基础镜像

- COPY 与 ADD 优先使用前者
 - COPY 指令非常简单，仅用于将文件拷贝到镜像中。ADD 相对来讲复杂，可以用于下载远程文件以及解压压缩包
- 合理调整 COPY 与 RUN 的顺序
 - 应该把变化最少的部分放在 Dockerfile 的前面，这样可以充分利用镜像缓存。

6.1.4 Docker Compose

Docker Compose 介绍

Docker Compose 是 Docker 官方编排（orchestration）项目之一，负责快速的部署分布式应用。

它允许用户通过一个 Docker-compose.yml 模板文件（YAML 格式）来定义一组相关联的应用容器为一个项目（project）。

- 服务 (Service): 一个应用的容器，实际上可以包括若干运行相同镜像的容器实例。
- 项目 (Project): 由一组关联的应用容器组成的一个完整业务单元，在 Docker-compose.yml 文件中定义。
- Compose 的默认管理对象是项目，通过子命令对项目中的一组容器进行便捷地生命周期管理。
- Compose 项目由 Python 编写，实现上调用了 Docker 服务提供的 API 来对容器进行管理。

常用命令

大部分命令跟 docker 下的命令功能相同

up

尝试自动完成包括构建镜像，（重新）创建服务，启动服务，并关联服务相关容器

```
docker-compose -f localhost.yml up -d
```

- -f 指定使用的 Compose 模板文件，默认为 docker-compose.yml，可以多次指定
- -d 后台启动并运行所有的容器
- 默认网络为 bridge，可通过--x-network-driver DRIVER 设置
- 默认，如果服务容器已经存在，docker-compose up 将会尝试停止容器，然后重新创建（保持使用 volumes-from 挂载的卷），以保证新启动的服务匹配 docker-compose.yml 文件的最新内容。
- 如果用户不希望容器被停止并重新创建，可以使用 docker-compose up --no-recreate。这样将只会启动处于停止状态的容器，而忽略已经运行的服务。
- 如果用户只想重新部署某个服务，可以使用 docker-compose up --no-deps -d <SERVICE_NAME> 来重新创建服务并后台停止旧服务，启动新服务，并不会影响到其所依赖的服务。

down

停止 up 命令所启动的容器，并移除网络

```
docker-compose -f localhost.yml down
```

stop

停止已经处于运行状态的容器，但不删除。通过 `docker-compose start` 可以再次启动这些容器

```
docker-compose -f localhost.yml stop nginx
```

- 如果不指定服务名，即停止 yml 文件中的所有服务容器

start

启动已经存在的服务容器

```
docker-compose -f localhost.yml start nginx
```

- 如果不指定服务名，即启动 yml 文件中的所有服务容器

其余脚本命令

```
# 登录到nginx容器中
docker-compose exec nginx bash

# 重新启动nginx容器
docker-compose restart nginx

# 暂停nginx容器
docker-compose pause nginx

# 恢复nginx容器
docker-compose unpause nginx

# 删除容器（删除前必须关闭容器）
docker-compose rm nginx

# 停止nginx容器
docker-compose stop nginx

# 启动nginx容器
docker-compose start nginx

# 在php-fpm中不启动关联容器，并容器执行php -v 执行完成后删除容器
docker-compose run --no-deps --rm php-fpm php -v

# 构建镜像
docker-compose build nginx

# 不带缓存的构建。
```

(下页继续)

```

docker-compose build --no-cache nginx

# 验证 (docker-compose.
→yml) 文件配置, 当配置正确时, 不输出任何内容, 当文件配置错误, 输出错误信息。
docker-compose config -q

# 查看nginx的日志
docker-compose logs nginx

# 查看nginx的实时日志
docker-compose logs -f nginx

# 以json的形式输出nginx的docker日志
docker-compose events --json nginx

```

docker-compose file

分解案例

```

##服务基于已经存在的镜像
services:
  web:
    image: hello-world

##服务基于dockerfile
build: /path/to/build/dir
build: ./dir
build:
  context: ../
  dockerfile: path/of/Dockerfile

build: ./dir
image: webapp:tag

##command command命令可以覆盖容器启动后默认执行的命令
command: bundle exec thin -p 3000
command: [bundle, exec, thin, -p, 3000]

##container_name
## Compose 的容器名称格式是: <项目名称><服务名称><序号>
##_
→虽然可以自定义项目名称、服务名称, 但是如果你想完全控制容器的命名, 可以使用这个标签指定
container_name: app

##depends_on depends_on解决了容器的依赖、启动先后的问题
version: '2'
services:
  web:
    build: .
    depends_on:
      - db
      - redis

```

(下页继续)

(续上页)

```

redis:
  image: redis
db:

##dns
dns: 8.8.8.8
dns:
- 8.8.8.8
- 9.9.9.9

dns_search: example.com
dns_search:
- dc1.example.com
- dc2.example.com

##tmpfs 挂载临时目录到容器内部，与run的参数一样效果
tmpfs: /run
tmpfs:
- /run
- /tmp

##environment_
↪设置镜像变量，它可以保存变量到镜像里，也就是说启动的容器也会包含这些变量设置
environment:
  RACK_ENV: development
  SHOW: 'true'
  SESSION_SECRET:

environment:
- RACK_ENV=development
- SHOW=true
- SESSION_SECRET

##expose 用于指定暴露的端口，但是只是作为参考，端口映射的话还得ports标签
expose:
- "3000"
- "8000"

##external_links
# 在使用Docker的过程中，我们会有许多单独使用docker_
↪run启动的容器，为了使Compose能够连接这些不在docker-compose.
↪yaml中定义的容器，我们需要一个特殊的标签，就是external_
↪links，它可以让Compose项目里面的容器连接到那些项目配置外部的容器（前提是外部容器中必须至少有一个容

external_links:
- redis_1
- project_db_1:mysql
- project_db_1:postgresql

##extra_hosts 添加主机名的标签，就是往容器内部/etc/hosts文件中添加一些记录
extra_hosts:
- "somehost:162.242.195.82"
- "otherhost:50.31.209.229"

##labels 向容器添加元数据，和Dockerfile的lable指令一个意思

```

(下页继续)

```

labels:
  com.example.description: "Accounting webapp"
  com.example.department: "Finance"
  com.example.label-with-empty-value: ""
labels:
  - "com.example.description=Accounting webapp"
  - "com.example.department=Finance"
  - "com.example.label-with-empty-value"

##links
# 解决容器连接问题，与docker的-link一样的效果，会连接到其他服务中的容器，
→使用的别名将会自动在服务容器中的/etc/hosts里创建
links:
  - db
  - db:database
  - redis

##ports
# 用作端口映射
# 使用HOST:CONTAINER格式或者只是指定容器的端口，宿主机会随机映射端口
ports:
  - "3000"
  - "8000:8000"
  - "49100:22"
  - "127.0.0.1:8001:8001"
#_
→当使用HOST:CONTAINER格式来映射端口时，如果你使用的容器端口小于60你可能会得到错误的结果，因为YAML将

##security_opt
#_
→为每个容器覆盖默认的标签。简单说来就是管理全部服务的标签，比如设置全部服务的user标签值为USER
security_opt:
  - label:user:USER
  - label:role:ROLE

##volumes
#_
→挂载一个目录或者一个已经存在的数据卷容器，可以直接使用[HOST:CONTAINER]这样的格式，或者使用[HOST:CONTAINER]
# compose的数据卷指定路径可以是相对路径，使用 . 或者 ... 来指定性对目录

volumes:
  // 只是指定一个路径，Docker 会自动在创建一个数据卷（这个路径是容器内部的）。
  - /var/lib/mysql

  // 使用绝对路径挂载数据卷
  - /opt/data:/var/lib/mysql

  // 以 Compose 配置文件为中心的相对路径作为数据卷挂载到容器。
  - ./cache:/tmp/cache

  // 使用用户的相对路径（~/ 表示的目录是 /home/<用户目录>/ 或者 /root/）。
  - ~/configs:/etc/configs/:ro

  // 已经存在的命名的数据卷。

```

(续上页)

```

- datavolume:/var/lib/mysql

# 如果你不使用宿主机的路径, 你可以指定一个volume_driver。
volume_driver: mydriver

##volumes_from
#_
→从其它容器或者服务挂载数据卷, 可选的参数是:ro或者:rw, 前者表示容器只读, 后者表示容器对数据卷是可读

volumes_from:
- service_name
- service_name:ro
- container:container_name
- container:container_name:rw

##network_mode
# 网络模式, 与docker client的-net参数类似, 只是相对多了一个service:[service name]的格式
network_mode: "bridge"
network_mode: "host"
network_mode: "none"
network_mode: "service:[service name]"
network_mode: "container:[container name/id]"

##networks 加入指定网络
services:
  some-service:
    networks:
      - some-network
      - other-network

```

Compose file 版本和兼容性

6.2 Apllo

6.2.1 Apollo 配置中心介绍

Jason Song edited this page on 13 Jul · 34 revisions

- 1、What is Apollo
- 2、Why Apollo
- 3、Apollo at a glance
- 4、Apollo in depth
- 5、Contribute to Apollo

6.2.2 1、What is Apollo

1.1 背景

随着程序功能的日益复杂，程序的配置日益增多：各种功能的开关、参数的配置、服务器的地址……

对程序配置的期望值也越来越高：配置修改后实时生效，灰度发布，分环境、分集群管理配置，完善的权限、审核机制……

在这样的大环境下，传统的通过配置文件、数据库等方式已经越来越无法满足开发人员对配置管理的需求。

Apollo 配置中心应运而生！

1.2 Apollo 简介

Apollo（阿波罗）是携程框架部门研发的开源配置管理中心，能够集中化管理应用不同环境、不同集群的配置，配置修改后能够实时推送到应用端，并且具备规范的权限、流程治理等特性。

Apollo 支持 4 个维度管理 Key-Value 格式的配置：

1. application (应用)
2. environment (环境)
3. cluster (集群)
4. namespace (命名空间)

同时，Apollo 基于开源模式开发，开源地址：<https://github.com/ctripcorp/apollo>

1.2 配置基本概念

既然 Apollo 定位于配置中心，那么在这里有必要先简单介绍一下什么是配置。

按照我们的理解，配置有以下几个属性：

- **配置是独立于程序的只读变量**
 - 配置首先是独立于程序的，同一份程序在不同的配置下会有不同的行为。
 - 其次，配置对于程序是只读的，程序通过读取配置来改变自己的行为，但是程序不应该去改变配置。
 - 常见的配置有：DB Connection Str、Thread Pool Size、Buffer Size、Request Timeout、Feature Switch、Server Urls 等。
- **配置伴随应用的整个生命周期**
 - 配置贯穿于应用的整个生命周期，应用在启动时通过读取配置来初始化，在运行时根据配置调整行为。
- **配置可以有多种加载方式**
 - 配置也有很多种加载方式，常见的有程序内部 hard code，配置文件，环境变量，启动参数，基于数据库等
- **配置需要治理**
 - 权限控制
 - * 由于配置能改变程序的行为，不正确的配置甚至能引起灾难，所以对配置的修改必须有比较完善的权限控制

- 不同环境、集群配置管理
 - * 同一份程序在不同的环境（开发，测试，生产）、不同的集群（如不同的数据中心）经常需要有不同的配置，所以需要有良好的环境、集群配置管理
- 框架类组件配置管理
 - * 还有一类比较特殊的配置 - 框架类组件配置，比如 CAT 客户端的配置。
 - * 虽然这类框架类组件是由其他团队开发、维护，但是运行时是在业务实际应用内的，所以本质上可以认为框架类组件也是应用的一部分。
 - * 这类组件对应的配置也需要有比较完善的管理方式。

6.2.3 2、Why Apollo

正是基于配置的特殊性，所以 Apollo 从设计之初就立志于成为一个有治理能力的配置发布平台，目前提供了以下的特性：

- **统一管理不同环境、不同集群的配置**
 - Apollo 提供了一个统一界面集中式管理不同环境（environment）、不同集群（cluster）、不同命名空间（namespace）的配置。
 - 同一份代码部署在不同的集群，可以有不同的配置，比如 zookeeper 的地址等
 - 通过命名空间（namespace）可以很方便地支持多个不同应用共享同一份配置，同时还允许应用对共享的配置进行覆盖
- **配置修改实时生效（热发布）**
 - 用户在 Apollo 修改完配置并发布后，客户端能实时（1 秒）接收到最新的配置，并通知到应用程序
- **版本发布管理**
 - 所有的配置发布都有版本概念，从而可以方便地支持配置的回滚
- **灰度发布**
 - 支持配置的灰度发布，比如点了发布后，只对部分应用实例生效，等观察一段时间没问题后再推给所有应用实例
- **权限管理、发布审核、操作审计**
 - 应用和配置的管理都有完善的权限管理机制，对配置的管理还分为了编辑和发布两个环节，从而减少人为的错误。
 - 所有的操作都有审计日志，可以方便地追踪问题
- **客户端配置信息监控**
 - 可以在界面上方便地看到配置在被哪些实例使用
- **提供 Java 和 .Net 原生客户端**
 - 提供了 Java 和 .Net 的原生客户端，方便应用集成
 - 支持 Spring Placeholder, Annotation 和 Spring Boot 的 ConfigurationProperties，方便应用使用（需要 Spring 3.1.1+）
 - 同时提供了 Http 接口，非 Java 和 .Net 应用也可以方便地使用
- **提供开放平台 API**

- Apollo 自身提供了比较完善的统一配置管理界面，支持多环境、多数据中心配置管理、权限、流程治理等特性。不过 Apollo 出于通用性考虑，不会对配置的修改做过多限制，只要符合基本的格式就能保存，不会针对不同的配置值进行针对性的校验，如数据库用户名、密码，Redis 服务地址等
- 对于这类应用配置，Apollo 支持应用方通过开放平台 API 在 Apollo 进行配置的修改和发布，并且具备完善的授权和权限控制

• 部署简单

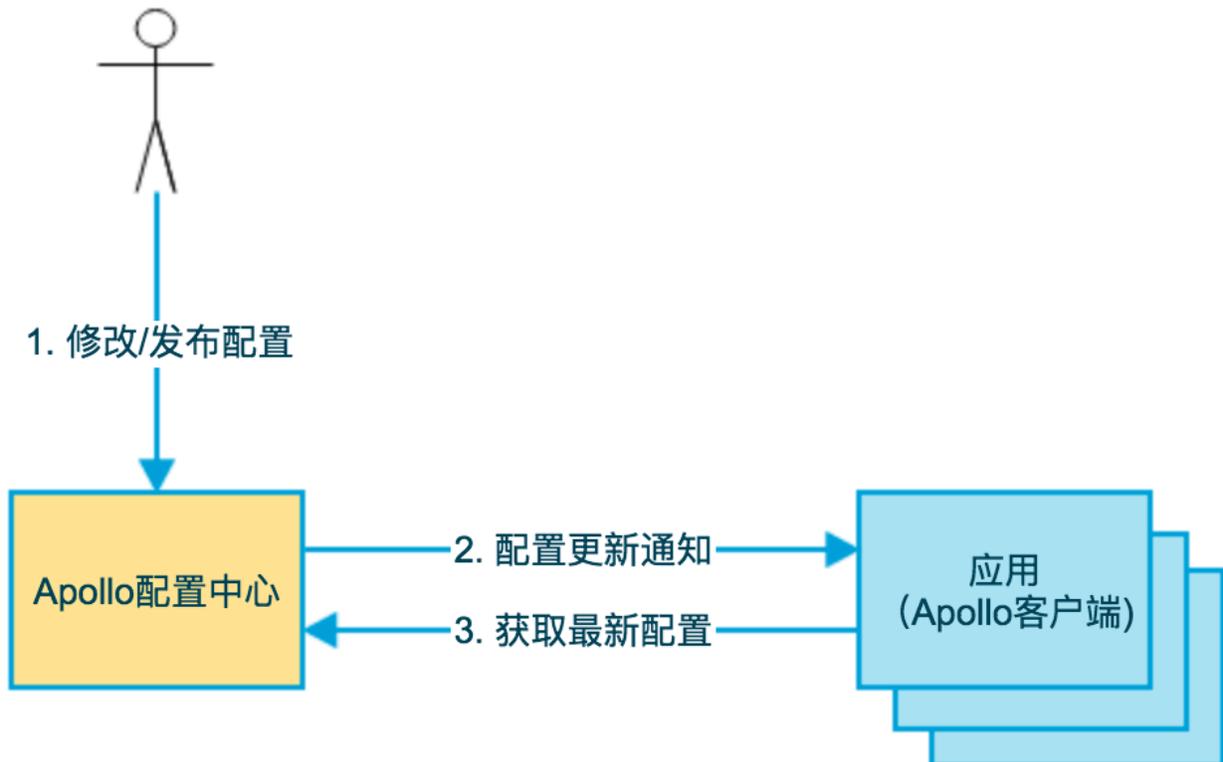
- 配置中心作为基础服务，可用性要求非常高，这就要求 Apollo 对外部依赖尽可能地少
- 目前唯一的外部依赖是 MySQL，所以部署非常简单，只要安装好 Java 和 MySQL 就可以让 Apollo 跑起来
- Apollo 还提供了打包脚本，一键就可以生成所有需要的安装包，并且支持自定义运行时参数

6.2.4 3、Apollo at a glance

3.1 基础模型

如下即是 Apollo 的基础模型：

1. 用户在配置中心对配置进行修改并发布
2. 配置中心通知 Apollo 客户端有配置更新
3. Apollo 客户端从配置中心拉取最新的配置、更新本地配置并通知到应用



3.2 界面概览

Apollo 配置中心

搜索项目(AppId、项目名) 帮助 song_s

私有

▼ application properties 发布 回滚 发布历史 授权 灰度

表格 文本 更改历史 实例列表

发布状态	Key	Value	备注	最后修改人	最后修改时间	操作
已发布	timeout	3000		song_s	2017-02-16 13:24:58	✎ ✕
已发布	kibana.uri	http://1.1.1.2:5600		song_s	2016-11-25 20:57:27	✎ ✕
已发布	elastic.document.type	biz1		song_s	2017-01-11 19:14:06	✎ ✕
已发布	elastic.cluster.name	es-cluster		song_s	2016-10-18 19:57:29	✎ ✕
已发布	elastic.cluster	2.2.2.2:9300,4.4.4.4:9300		zhanglea	2016-12-08 14:19:43	✎ ✕
已发布	page.size	20		song_s	2016-12-27 14:58:56	✎ ✕
已发布	zookeeper.address	10.1.12.2		song_s	2016-10-19 11:33:50	✎ ✕

关联

▼ FX.apollo properties 发布 回滚 发布历史 授权 灰度

表格 文本 更改历史 实例列表

覆盖的配置 filter by key ...

发布状态	Key	Value	备注	最后修改人	最后修改时间	操作
已发布	servers	3.3.3.3,4.4.4.4		song_s	2017-02-16 13:26:27	✎ ✕

公共的配置 (AppId:100003173, Cluster:default) filter by key ...

Key	Value	备注	最后修改人	最后修改时间	操作
batch	2000	样例项目会使用到, 勿删。	song_s	2017-02-16 13:27:07	✕
servers	1.1.1.1,2.2.2.2	样例项目会使用到, 勿删。	song_s	2016-10-12 14:03:34	

环境列表

FAT

default 集群

dev 集群

UAT

PRO

default 集群

SHAOY 集群

SHAIQ 集群

项目信息

AppId: 100004458

应用名: apollo-demo

部门: 框架(FX)

负责人: song_s

邮箱: song_s@ctrip.com

管理项目

添加集群

添加Namespace

上图是 Apollo 配置中心中一个项目的配置首页

- 在页面左上方的环境列表模块展示了所有的环境和集群，用户可以随时切换。
- 页面中央展示了两个 namespace(application 和 FX.apollo) 的配置信息，默认按照表格模式展示、编辑。用户也可以切换到文本模式，以文件形式查看、编辑。
- 页面上可以方便地进行发布、回滚、灰度、授权、查看更改历史和发布历史等操作

3.3 添加/修改配置项

用户可以通过配置中心界面方便的添加/修改配置项，更多使用说明请参见[应用接入指南](#)

Apollo 配置中心 应用ID/应用名 帮助 song_s

环境列表

- FAT
- UAT
- PRO
- default
- SHAOY
- SHAJQ

项目信息

AppId: 100004458
 应用名: apollo-demo
 部门: 框架(FX)
 负责人: song_s
 负责人Email: song_s@ctrip.com

application properties 发布 回滚 发布历史 授权 同步配置

表格 文本 更改历史 实例列表 过滤配置 + 新增配置

Key	Value	备注	最后修改人	最后修改时间	操作
request.timeout	200	请求超时时间 (毫秒)	song_s	2016-10-18 19:56:26	✎ ✕
kibana.url	http://1.1.1.2:5601		song_s	2016-10-18 19:57:29	✎ ✕
elastic.document.type	biz		song_s	2016-10-18 19:57:29	✎ ✕
elastic.cluster.name	es-cluster		song_s	2016-10-18 19:57:29	✎ ✕
elastic.cluster	2.2.2.2:9300,3.3.3.3,4.4.4.4:9300		song_s	2016-10-18 19:57:29	✎ ✕
page.size	10		song_s	2016-10-18 19:57:29	✎ ✕
zookeeper.address	10.1.12.1		song_s	2016-10-18 19:57:29	✎ ✕

输入配置信息：

Apollo 配置中心 应用ID/应用名 帮助 song_s

环境列表

- FAT
- UAT
- PRO
- default
- SHAOY
- SHAJQ

项目信息

AppId: 100004458
 应用名: apollo-demo
 部门: 框架(FX)
 负责人: song_s

修改配置项 ✕

Key: request.timeout

* Value:

Comment:

提交

Key	Value	备注	最后修改人	最后修改时间	操作
elastic.cluster	2.2.2.2:9300,3.3.3.3,4.4.4.4:9300		song_s	2016-10-18 19:57:29	✎ ✕
page.size	10		song_s	2016-10-18 19:57:29	✎ ✕
zookeeper.address	10.1.12.1		song_s	2016-10-18 19:57:29	✎ ✕

3.4 发布配置

通过配置中心发布配置：

Apollo 配置中心 应用ID/应用名 Go 帮助 song_s

环境列表

- FAT
- UAT
- PRO
- default 集群
- SHAOY 集群
- SHAJQ 集群

项目信息

AppId: 100004458
 应用名: apollo-demo
 部门: 框架(FX)
 负责人: song_s

application properties 有修改 1 发布 回滚 发布历史 授权 同步配置

表格 文本 更改历史 实例列表 过滤配置 新增配置

Key	Value	备注	最后修改人	最后修改时间	操作
request.timeout	100	请求超时时间 (毫秒)	song_s	2016-10-18 19:59:51	编辑 删除
kibana.url	http://1.1.1.2:5601		song_s	2016-10-18 19:57:29	编辑 删除
elastic.document.type	biz		song_s	2016-10-18 19:57:29	编辑 删除
elastic.cluster.name	es-cluster		song_s	2016-10-18 19:57:29	编辑 删除
elastic.cluster	2.2.2.2:9300,3.3.3.3,4.4.4.4:9300		song_s	2016-10-18 19:57:29	编辑 删除
page.size	10		song_s	2016-10-18 19:57:29	编辑 删除
zookeeper.address	10.1.12.1		song_s	2016-10-18 19:57:29	编辑 删除

填写发布信息：

Apollo 配置中心 应用ID/应用名 Go 帮助 song_s

环境列表

- FAT
- UAT
- PRO
- default 集群
- SHAOY 集群
- SHAJQ 集群

项目信息

AppId: 100004458
 应用名: apollo-demo
 部门: 框架(FX)
 负责人: song_s

发布

Changes	Key	Old Value	New Value	最后修改人	最后修改时间
	request.timeout	200	100	song_s	2016-10-18 19:59:51

* Release Name:

Comment:

关闭 发布

3.5 客户端获取配置 (Java API 样例)

配置发布后, 就能在客户端获取到了, 以 Java 为例, 获取配置的示例代码如下。Apollo 客户端还支持和 Spring 整合, 更多客户端使用说明请参见[Java 客户端使用指南](#)和[.Net 客户端使用指南](#)。

```
Config config = ConfigService.getAppConfig();
Integer defaultRequestTimeout = 200;
Integer requestTimeout = config.getIntProperty("requestTimeout",
    ↪ defaultRequestTimeout);
```

3.6 客户端监听配置变化

通过上述获取配置代码，应用就能实时获取到最新的配置了。

不过在某些场景下，应用还需要在配置变化时获得通知，比如数据库连接的切换等，所以 Apollo 还提供了监听配置变化的功能，Java 示例如下：

```
Config config = ConfigService.getAppConfig();
config.addChangeListener(new ConfigChangeListener() {
    @Override
    public void onChange(ConfigChangeEvent changeEvent) {
        for (String key : changeEvent.changedKeys()) {
            ConfigChange change = changeEvent.getChange(key);
            System.out.println(String.format(
                "Found change - key: %s, oldValue: %s, newValue: %s, changeType: %s",
                change.getPropertyName(), change.getOldValue(),
                change.getNewValue(), change.getChangeType()));
        }
    }
});
```

3.7 Spring 集成样例

Apollo 和 Spring 也可以很方便地集成，只需要标注 `@EnableApolloConfig` 后就可以通过 `@Value` 获取配置信息：

```
@Configuration
@EnableApolloConfig
public class AppConfig {}
```

```
@Component
public class SomeBean {
    //timeout的值会自动更新
    @Value("${request.timeout:200}")
    private int timeout;
}
```

6.2.5 4、Apollo in depth

通过上面的介绍，相信大家已经对 Apollo 有了一个初步的了解，并且相信已经覆盖到了大部分的使用场景。接下来会主要介绍 Apollo 的 cluster 管理（集群）、namespace 管理（命名空间）和对应的配置获取规则。

4.1 Core Concepts

在介绍高级特性前，我们有必要先来了解一下 Apollo 中的几个核心概念：

1. application (应用)

- 这个很好理解，就是实际使用配置的应用，Apollo 客户端在运行时需要知道当前应用是谁，从而可以去获取对应的配置
- 每个应用都需要有唯一的身份标识 -- `appId`，我们认为应用身份是跟着代码走的，所以需要在代码中配置，具体信息请参见[Java 客户端使用指南](#)。

2. environment (环境)

- 配置对应的环境，Apollo 客户端在运行时需要知道当前应用处于哪个环境，从而可以去获取应用的配置
- 我们认为环境和代码无关，同一份代码部署在不同的环境就应该能够获取到不同环境的配置
- 所以环境默认是通过读取机器上的配置（server.properties 中的 env 属性）指定的，不过为了开发方便，我们也支持运行时通过 System Property 等指定，具体信息请参见[Java 客户端使用指南](#)。

3. cluster (集群)

- 一个应用下不同实例的分组，比如典型的可以按照数据中心分，把上海机房的应用实例分为一个集群，把北京机房的应用实例分为另一个集群。
- 对不同的 cluster，同一个配置可以有不一样的值，如 zookeeper 地址。
- 集群默认是通过读取机器上的配置（server.properties 中的 idc 属性）指定的，不过也支持运行时通过 System Property 指定，具体信息请参见[Java 客户端使用指南](#)。

4. namespace (命名空间)

- 一个应用下不同配置的分组，可以简单地把 namespace 类比为文件，不同类型的配置存放在不同的文件中，如数据库配置文件，RPC 配置文件，应用自身的配置文件等
- 应用可以直接读取到公共组件的配置 namespace，如 DAL，RPC 等
- 应用也可以通过继承公共组件的配置 namespace 来对公共组件的配置做调整，如 DAL 的初始数据库连接数

4.2 自定义 Cluster

【本节内容仅对应用需要对不同集群应用不同配置才需要，如没有相关需求，可以跳过本节】

比如我们有应用在 A 数据中心和 B 数据中心都有部署，那么如果希望两个数据中心的配置不一样的话，我们可以通过新建 cluster 来解决。

4.2.1 新建 Cluster

新建 Cluster 只有项目的管理员才有权限，管理员可以在页面左侧看到“添加集群”按钮。



点击后就进入到集群添加页面，一般情况下可以按照数据中心来划分集群，如 SHAJQ、SHAOY 等。

不过也支持自定义集群，比如可以为 A 机房的某一台机器和 B 机房的某一台机创建一个集群，使用一套配置。

创建集群

* 应用AppId 100003173

* 集群名称

SHAJQ

(部署集群如:SHAJQ,SHAOY 或自定义集群如:SHAJQ-xx,SHAJQ-yy)

* 选择环境

DEV

提交

4.2.2 在 Cluster 中添加配置并发布

集群添加成功后,就可以为该集群添加配置了,首先需要按照下图所示切换到 SHAJQ 集群,之后配置添加流程和3.2 添加/修改配置项一样,这里就不再赘述了。

DEV

default

SHAJQ

i 应用信息

AppId: 100003173

应用名: apollo-portal

部门: 框架(FX)

application

表格 **文本** **更改历史**

Key	Value
requestTimeout	200

4.2.3 指定应用实例所属的 Cluster

Apollo 会默认使用应用实例所在的数据中心作为 cluster，所以如果两者一致的话，不需要额外配置。

如果 cluster 和数据中心不一致的话，那么就需要通过 System Property 方式来指定运行时 cluster：

- -Dapollo.cluster=SomeCluster
- 这里注意 apollo.cluster 为全小写

4.3 自定义 Namespace

【本节仅对公共组件配置或需要多个应用共享配置才需要，如没有相关需求，可以跳过本节】

如果应用有公共组件（如 hermes-producer，cat-client 等）供其它应用使用，就需要通过自定义 namespace 来实现公共组件的配置。

4.3.1 新建 Namespace

以 hermes-producer 为例，需要先新建一个 namespace，新建 namespace 只有项目的管理员才有权限，管理员可以在页面左侧看到“添加 Namespace”按钮。

项目信息  

AppId: 100003806

应用名: hermes-portal

部门: 框架(FX)

负责人: jhliang

邮箱: jhliang@ctrip.com

 管理项目

 添加集群

 添加Namespace

点击后就进入 namespace 添加页面，Apollo 会把应用所属的部门作为 namespace 的前缀，如 FX。

新建Namespace

关联已存在的Namespace 创建新的Namespace

应用ID 100003806

* 名称 FX.Hermes.Producer

备注

4.3.2 关联到环境和集群

Namespace 创建完，需要选择在哪些环境和集群下使用

新建Namespace

关联已存在的Namespace 创建新的Namespace

应用ID 100003806

* 选择集群 环境 集群

DEV default

* namespace

4.3.3 在 Namespace 中添加配置项

接下来在这个新建的 namespace 下添加配置项

添加配置项 ✕

*** Key**

*** Value**

Comment

*** 选择集群**

<input type="checkbox"/>	环境	集群
<input checked="" type="checkbox"/>	DEV	default
<input type="checkbox"/>	FAT	default

添加完成后就能在 FX.Hermes.Producer 的 namespace 中看到配置。

FX.Hermes.Producer

Key	Value
sender.batchsize	500

4.3.4 发布 namespace 的配置

发布 ×

Changes:	Key	Old Value	New Value	最后修改人	最后修改时间
	sender.batchsize		500	apollo	2016-06-15 11:24:13

* Release Name:

Comment:

4.3.5 客户端获取 Namespace 配置

对自定义 namespace 的配置获取，稍有不同，需要程序传入 namespace 的名字。Apollo 客户端还支持和 Spring 整合，更多客户端使用说明请参见[Java 客户端使用指南](#)和[.Net 客户端使用指南](#)。

```
Config config = ConfigService.getConfig("FX.Hermes.Producer");
Integer defaultSenderBatchSize = 200;
Integer senderBatchSize = config.getIntProperty("sender.batchsize", ↵
↵defaultSenderBatchSize);
```

4.3.6 客户端监听 Namespace 配置变化

```
Config config = ConfigService.getConfig("FX.Hermes.Producer");
config.addChangeListener(new ConfigChangeListener() {
    @Override
    public void onChange(ConfigChangeEvent changeEvent) {
        System.out.println("Changes for namespace " + changeEvent.getNamespace());
        for (String key : changeEvent.changedKeys()) {
            ConfigChange change = changeEvent.getChange(key);
            System.out.println(String.format(
                "Found change - key: %s, oldValue: %s, newValue: %s, changeType: %s",
                change.getPropertyName(), change.getOldValue(),
                change.getNewValue(), change.getChangeType()));
        }
    }
});
```

4.3.7 Spring 集成样例

```
@Configuration
@EnableApolloConfig("FX.Hermes.Producer")
public class AppConfig {}
```

```
@Component
public class SomeBean {
    //timeout的值会自动更新
    @Value("${request.timeout:200}")
    private int timeout;
}
```

4.4 配置获取规则

【本节仅当应用自定义了集群或 namespace 才需要，如无相关需求，可以跳过本节】

在有了 cluster 概念后，配置的规则就显得重要了。

比如应用部署在 A 机房，但是并没有在 Apollo 新建 cluster，这个时候 Apollo 的行为是怎样的？

或者在运行时指定了 cluster=SomeCluster，但是并没有在 Apollo 新建 cluster，这个时候 Apollo 的行为是怎样的？

接下来就来介绍一下配置获取的规则。

4.4.1 应用自身配置的获取规则

当应用使用下面的语句获取配置时，我们称之为获取应用自身的配置，也就是应用自身的 application namespace 的配置。

```
Config coig = ConfigService.getAppConfig();
```

对这种情况的配置获取规则，简而言之如下：

1. 首先查找运行时 cluster 的配置（通过 apollo.cluster 指定）
2. 如果没有找到，则查找数据中心 cluster 的配置
3. 如果还是没有找到，则返回默认 cluster 的配置

图示如下：



所以如果应用部署在 A 数据中心，但是用户没有在 Apollo 创建 cluster，那么获取的配置就是默认 cluster (default) 的。

如果应用部署在 A 数据中心，同时在运行时指定了 SomeCluster，但是没有在 Apollo 创建 cluster，那么获取的配置就是 A 数据中心 cluster 的配置，如果 A 数据中心 cluster 没有配置的话，那么获取的配置就是默认 cluster (default) 的。

4.4.2 公共组件配置的获取规则

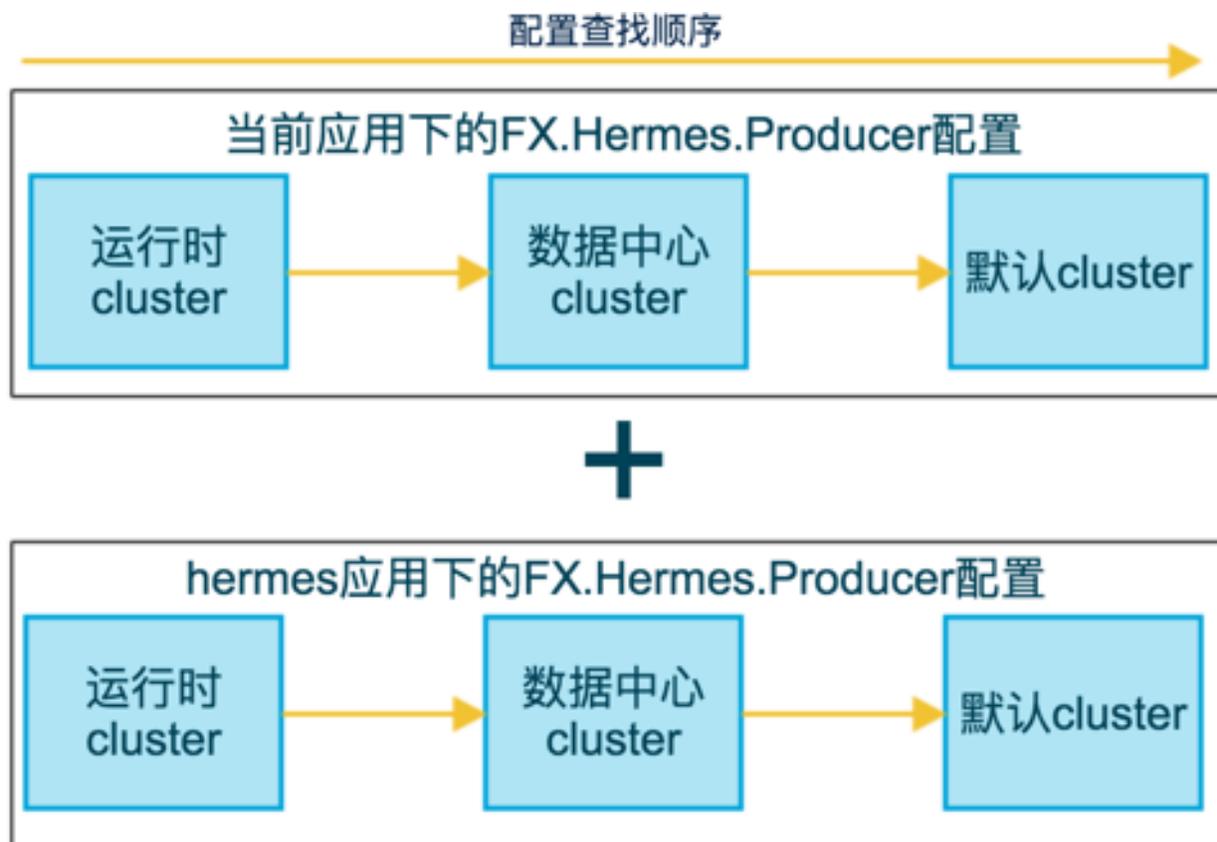
以 `FX.Hermes.Producer` 为例，hermes producer 是 hermes 发布的公共组件。当使用下面的语句获取配置时，我们称之为获取公共组件的配置。

```
Config config = ConfigService.getConfig("FX.Hermes.Producer");
```

对这种情况的配置获取规则，简而言之如下：

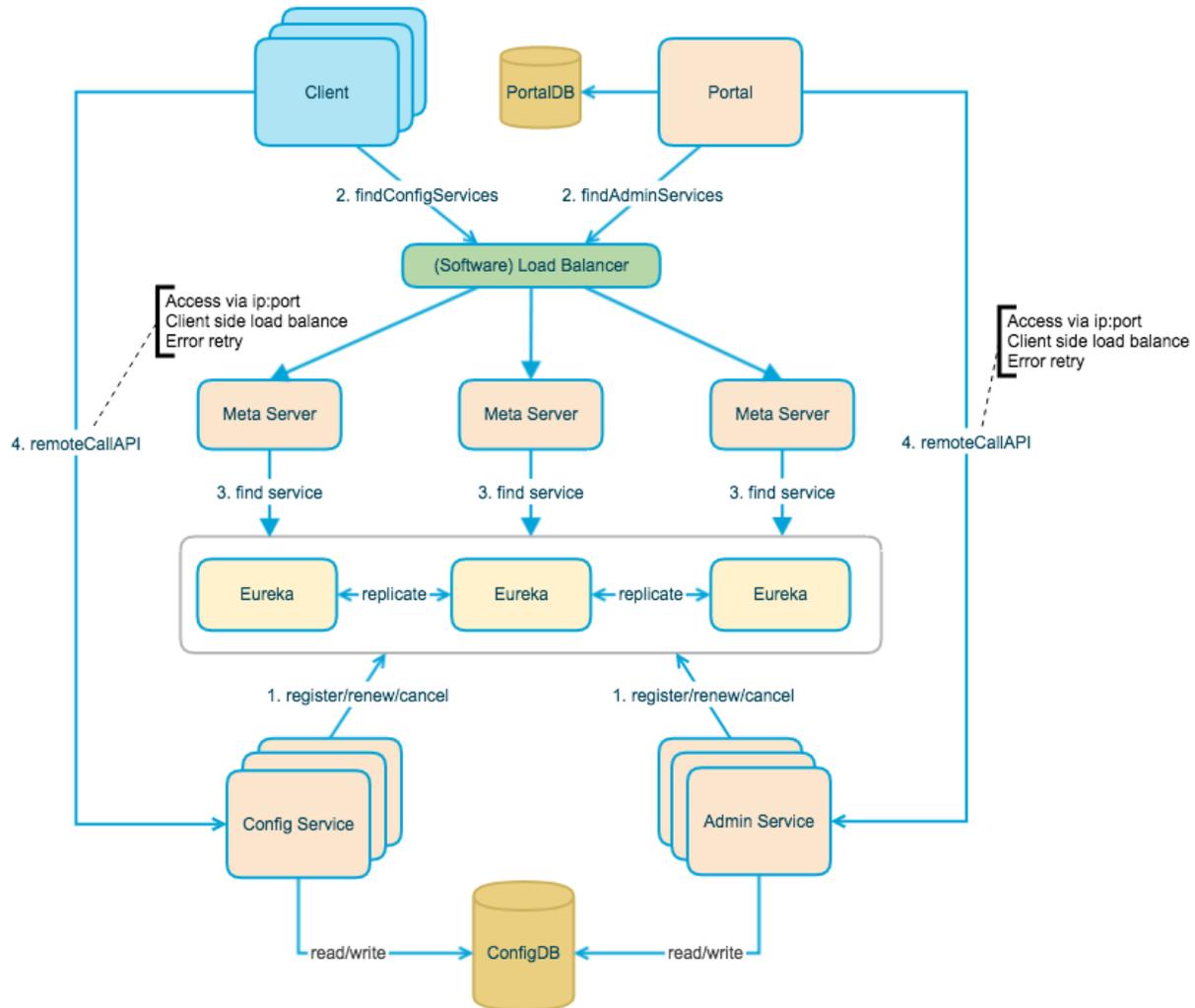
1. 首先获取当前应用下的 `FX.Hermes.Producer` namespace 的配置
2. 然后获取 hermes 应用下 `FX.Hermes.Producer` namespace 的配置
3. 上面两部分配置的并集就是最终使用的配置，如有 key 一样的部分，以当前应用优先

图示如下：



通过这种方式，就实现了对框架类组件的配置管理，框架组件提供方提供配置的默认值，应用如果有特殊需求，可以自行覆盖。

4.5 总体设计



上图简要描述了 Apollo 的总体设计，我们可以从下往上看：

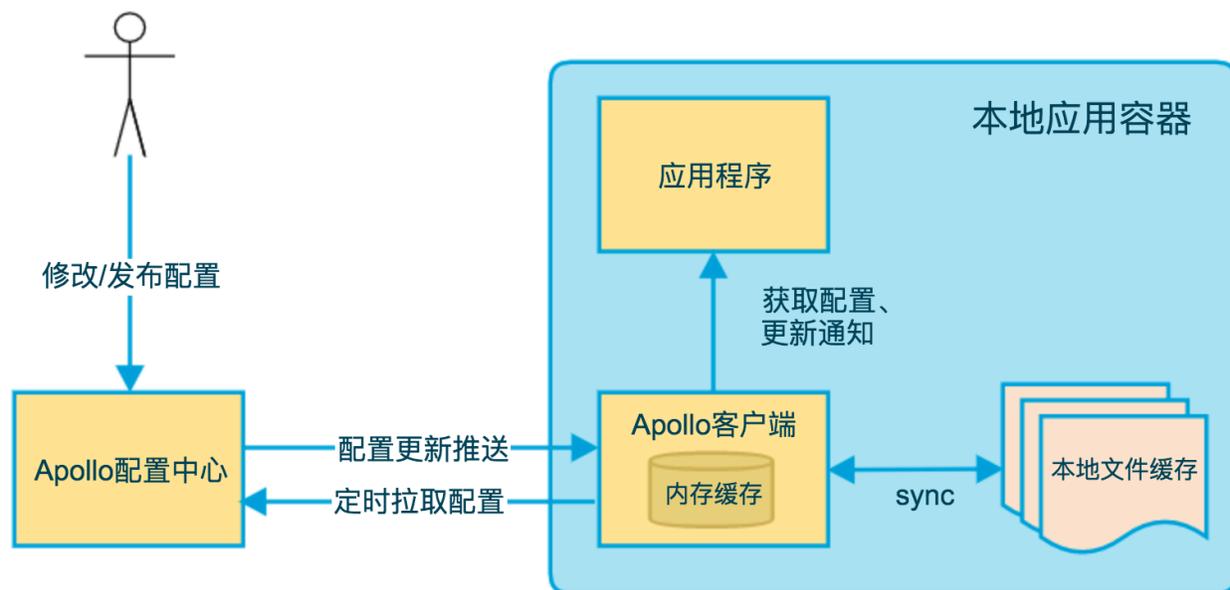
- Config Service 提供配置的读取、推送等功能，服务对象是 Apollo 客户端
- Admin Service 提供配置的修改、发布等功能，服务对象是 Apollo Portal（管理界面）
- Config Service 和 Admin Service 都是多实例、无状态部署，所以需要将自己注册到 Eureka 中并保持心跳
- 在 Eureka 之上我们架了一层 Meta Server 用于封装 Eureka 的服务发现接口
- Client 通过域名访问 Meta Server 获取 Config Service 服务列表（IP+Port），而后直接通过 IP+Port 访问服务，同时在 Client 侧会做 load balance、错误重试
- Portal 通过域名访问 Meta Server 获取 Admin Service 服务列表（IP+Port），而后直接通过 IP+Port 访问服务，同时在 Portal 侧会做 load balance、错误重试
- 为了简化部署，我们实际上会把 Config Service、Eureka 和 Meta Server 三个逻辑角色部署在同一个 JVM 进程中

4.5.1 Why Eureka

为什么我们采用 Eureka 作为服务注册中心，而不是使用传统的 zk、etcd 呢？我大致总结了一下，有以下几方面的原因：

- 它提供了完整的 Service Registry 和 Service Discovery 实现
 - 首先是提供了完整的实现，并且也经受住了 Netflix 自己的生产环境考验，相对使用起来会比较省心。
- 和 Spring Cloud 无缝集成
 - 我们的项目本身就使用了 Spring Cloud 和 Spring Boot，同时 Spring Cloud 还有一套非常完善的开源代码来整合 Eureka，所以使用起来非常方便。
 - 另外，Eureka 还支持在我们应用自身的容器中启动，也就是说我们的应用启动完之后，既充当了 Eureka 的角色，同时也是服务的提供者。这样就极大的提高了服务的可用性。
 - **这一点是我们选择 Eureka 而不是 zk、etcd 等的主要原因，为了提高配置中心的可用性和降低部署复杂度，我们需要尽可能地减少外部依赖。**
- Open Source
 - 最后一点是开源，由于代码是开源的，所以非常便于我们了解它的实现原理和排查问题。

4.6 客户端设计



上图简要描述了 Apollo 客户端的实现原理：

1. 客户端和服务端保持了一个长连接，从而能第一时间获得配置更新的推送。
2. 客户端还会定时从 Apollo 配置中心服务端拉取应用的最新配置。
 - 这是一个 fallback 机制，为了防止推送机制失效导致配置不更新
 - 客户端定时拉取会上报本地版本，所以一般情况下，对于定时拉取的操作，服务端都会返回 304 - Not Modified

- 定时频率默认为每 5 分钟拉取一次，客户端也可以通过在运行时指定 System Property: `apollo.refreshInterval` 来覆盖，单位为分钟。
3. 客户端从 Apollo 配置中心服务端获取到应用的最新配置后，会保存在内存中
 4. 客户端会把从服务端获取到的配置在本地文件系统缓存一份
 - 在遇到服务不可用，或网络不通的时候，依然能从本地恢复配置
 5. 应用程序从 Apollo 客户端获取最新的配置、订阅配置更新通知

4.6.1 配置更新推送实现

前面提到了 Apollo 客户端和服务端保持了一个长连接，从而能第一时间获得配置更新的推送。

长连接实际上我们是通过 Http Long Polling 实现的，具体而言：

- 客户端发起一个 Http 请求到服务端
- 服务端会保持住这个连接 60 秒
 - 如果在 60 秒内有客户端关心的配置变化，被保持住的客户端请求会立即返回，并告知客户端有配置变化的 namespace 信息，客户端会据此拉取对应 namespace 的最新配置
 - 如果在 60 秒内没有客户端关心的配置变化，那么会返回 Http 状态码 304 给客户端
- 客户端在收到服务端请求后会立即重新发起连接，回到第一步

考虑到会有数万客户端向服务端发起长连，在服务端我们使用了 `async servlet(Spring DeferredResult)` 来服务 Http Long Polling 请求。

4.7 可用性考虑

配置中心作为基础服务，可用性要求非常高，下面的表格描述了不同场景下 Apollo 的可用性：

6.2.6 5、Contribute to Apollo

Apollo 从开发之初就是以开源模式开发的，所以也非常欢迎有兴趣、有余力的朋友一起加入进来。

服务端开发使用的是 Java，基于 Spring Cloud 和 Spring Boot 框架。客户端目前提供了 Java 和 .Net 两种实现。

Github 地址：<https://github.com/ctripcorp/apollo>

欢迎大家发起 Pull Request!

- 设计文档
 - Apollo 配置中心介绍
 - Apollo 配置中心设计
 - Apollo 核心概念之 “Namespace”
- 部署文档
 - Quick Start
 - Docker 方式部署 Quick Start
 - 分布式部署指南
 - Apollo 源码解析（全）

- 开发文档
 - Apollo 开发指南
 - Code Styles
 - * Eclipse Code Style
 - * IntelliJ Code Style
 - Portal 实现用户登录功能
 - 邮件模板样例
- 系统使用文档
 - Apollo 使用指南
 - Java 客户端使用指南
 - .Net 客户端使用指南
 - Go、Python、NodeJS、PHP 等客户端使用指南
 - 其它语言客户端接入指南
 - Apollo 开放平台接入指南
 - Apollo 使用场景和示例代码
- FAQ
 - 常见问题回答
 - 部署 & 开发遇到的常见问题

6.3 Eureka

6.3.1 1、背景介绍

Eureka是 Netflix 开源的一款提供服务注册和发现的产品。

其官方文档中对自己的定义是：

Eureka is a REST (Representational State Transfer) based service that is primarily used in the AWS cloud for locating services for the purpose of load balancing and failover of middle-tier servers. We call this service, the Eureka Server. Eureka also comes with a Java-based client component, the Eureka Client, which makes interactions with the service much easier. The client also has a built-in load balancer that does basic round-robin load balancing.

我们在研发 Apollo 配置中心时 (ctripcorp/apollo)，考虑到配置中心是基础服务，有非常高的可用性要求，为了更好的支持服务动态扩容、缩容、失效剔除等特性，所以就选择了使用 Eureka 来提供服务注册和发现功能。

本着“当你选择一款开源产品后，你就应当对它负责，既要信任它又要挑战它”的原则，我花了一些时间比较深入的研究了 Eureka 的实现细节（好在 Eureka 的实现短小精悍，通读源码也没花太多时间），今天就来详细介绍一下。

6.3.2 2、Why Eureka?

那么为什么我们在项目中使用了 Eureka 呢？我大致总结了一下，有以下几方面的原因：

1) 它提供了完整的 **Service Registry** 和 **Service Discovery** 实现

首先是提供了完整的实现，并且也经受住了 Netflix 自己的生产环境考验，相对使用起来会比较省心。

2) 和 **Spring Cloud** 无缝集成

我们的项目本身就使用了 Spring Cloud 和 Spring Boot，同时 Spring Cloud 还有一套非常完善的开源代码来整合 Eureka，所以使用起来非常方便。

另外，Eureka 还支持在我们应用自身的容器中启动，也就是说我们的应用启动完之后，既充当了 Eureka 的角色，同时也是服务的提供者。这样就极大的提高了服务的可用性。

这一点是我们选择 **Eureka** 而不是 **zk**、**etcd** 等的主要原因，为了提高配置中心的可用性和降低部署复杂度，我们需要尽可能地减少外部依赖。

3) Open Source

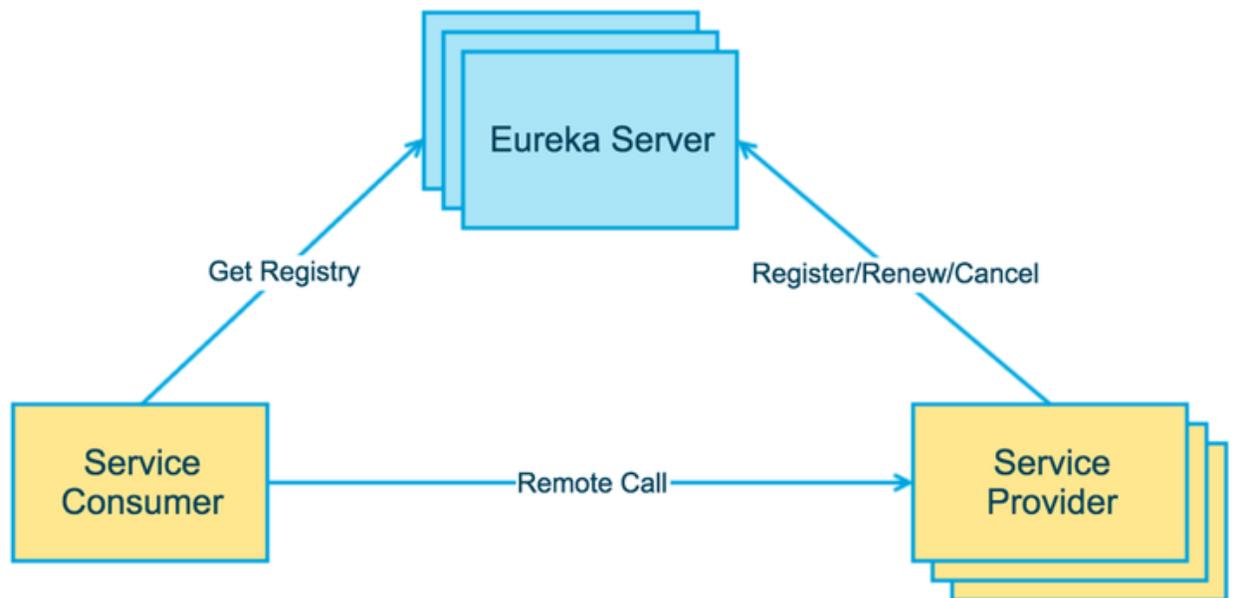
最后一点是开源，由于代码是开源的，所以非常便于我们了解它的实现原理和排查问题。

6.3.3 3、Dive into Eureka

相信大家看到这里，已经对 Eureka 有了一个初步的认识，接下来我们就来深入了解下它吧 ~

3.1 Overview

3.1.1 Basic Architecture



图

1

上图简要描述了 Eureka 的基本架构，由 3 个角色组成：

Eureka Server

提供服务注册和发现

Service Provider

服务提供方

将自身服务注册到 Eureka，从而使服务消费方能够找到

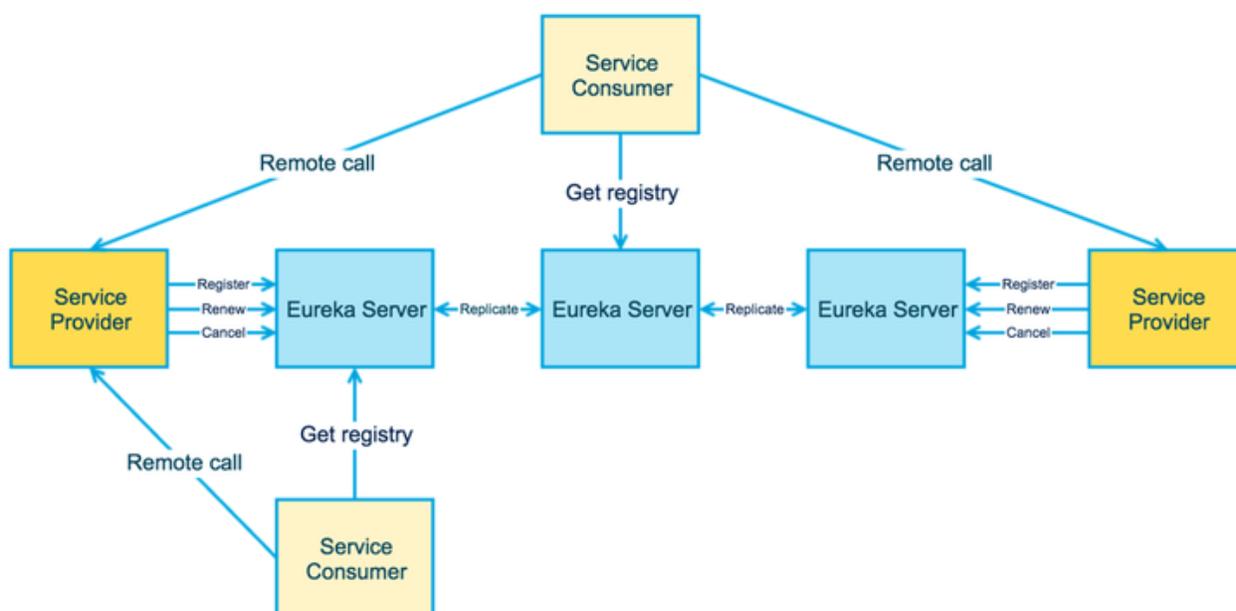
Service Consumer

服务消费方

从 Eureka 获取注册服务列表，从而能够消费服务

需要注意的是，上图中的 3 个角色都是逻辑角色。在实际运行中，这几个角色甚至可以是同一个实例，比如在我们项目中，Eureka Server 和 Service Provider 就是同一个 JVM 进程。

3.1.2 More in depth



图

2

上图更进一步的展示了 3 个角色之间的交互。

1. Service Provider 会向 Eureka Server 做 Register（服务注册）、Renew（服务续约）、Cancel（服务下线）等操作。
2. Eureka Server 之间会做注册服务的同步，从而保证状态一致
3. Service Consumer 会向 Eureka Server 获取注册服务列表，并消费服务

3.2 Demo

为了给大家一个更直观的印象，我们可以通过一个简单的 demo 来实际运行一下，从而对 Eureka 有更好的了解。

3.2.1 Git Repository

Git 仓库：[nobodyiam/spring-cloud-in-action](#)

这个项目使用了 Spring Cloud 相关类库，包括：

- Spring Cloud Config
- Spring Cloud Eureka (Netflix)

3.2.2 准备工作

Demo 项目使用了 Spring Cloud Config 做配置，所以第一步先在本地启动 Config Server。

由于项目基于 Spring Boot 开发，所以直接运行 `com.nobodyiam.spring.cloud.in.action.config.ConfigServerApplication` 即可。

3.2.3 Eureka Server Demo

Eureka Server 的 Demo 模块名是：`eureka-server`。

3.2.3.1 Maven 依赖

`eureka-server` 是一个基于 Spring Boot 的 Web 应用，我们首先需要做的就是要在 pom 中引入 Spring Cloud Eureka Server 的依赖。

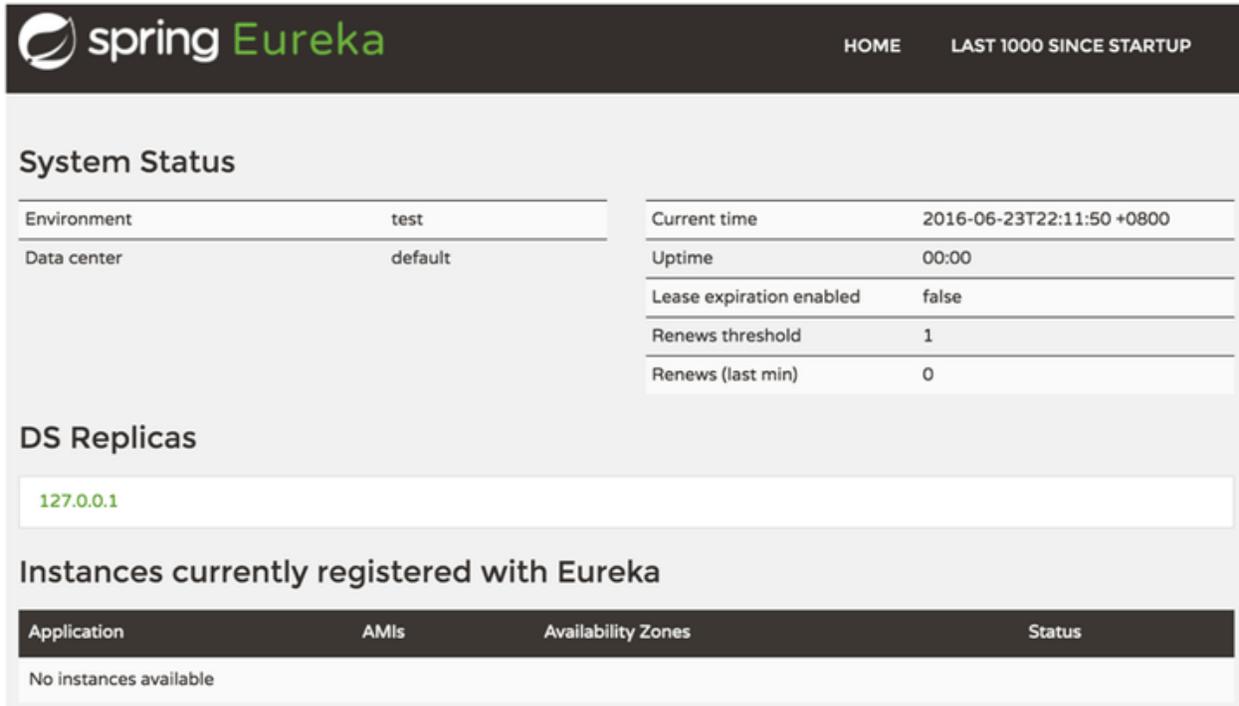
```
<dependency>
<groupId>org.springframework.cloud</groupId>
<artifactId>spring-cloud-starter-eureka-server</artifactId><version>1.2.0.
↪RELEASE</version>
</dependency>
```

3.2.3.2 启用 Eureka Server

启用 Eureka Server 非常简单，只需要加上 `@EnableEurekaServer` 即可。

```
@EnableEurekaServer
@SpringBootApplication
public class EurekaServiceApplication {
public static void main(String[] args){
SpringApplication.run(EurekaServiceApplication.class, args);
}
}
```

做完以上配置后，启动应用，Eureka Server 就开始工作了！
启动完，打开<http://localhost:8761>，就能看到启动成功的画面了。



图

3

3.2.4 Service Provider and Service Consumer Demo

Service Provider 的 Demo 模块名是：reservation-service。
Service Consumer 的 Demo 模块名是：reservation-client。

3.2.4.1 Maven 依赖

reservation-service 和 reservation-client 都是基于 Spring Boot 的 Web 应用，我们首先需要做的就是要在 pom 中引入 Spring Cloud Eureka 的依赖。

```
<dependency>
<groupId>org.springframework.cloud</groupId>
<artifactId>spring-cloud-starter-eureka</artifactId>    <version>1.2.0.
↪RELEASE</version>
</dependency>
```

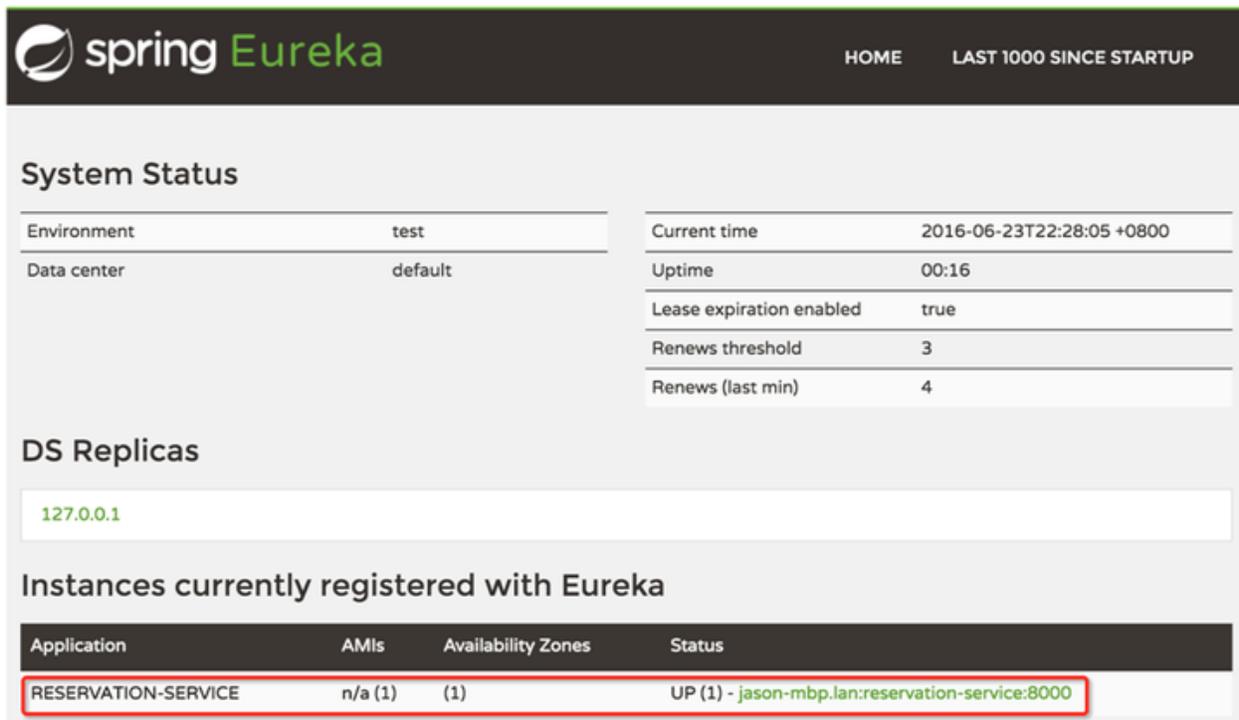
3.2.4.2 启动 Service Provider

启用 Service Provider 非常简单，只需要加上 @EnableDiscoveryClient 即可。

```
@EnableDiscoveryClient
@SpringBootApplication
public class ReservationServiceApplication {
    public static void main(String[] args) {
        ↪newSpringApplicationBuilder(ReservationServiceApplication.class)
        .run(args);
    }
}
```

做完以上配置后，启动应用，Server Provider 就开始工作了！

启动完，打开<http://localhost:8761>，就能看到服务已经注册到 Eureka Server 了。



The screenshot shows the Spring Eureka dashboard. At the top, there's a navigation bar with 'HOME' and 'LAST 1000 SINCE STARTUP'. Below that is the 'System Status' section, which contains two tables. The first table shows 'Environment: test' and 'Data center: default'. The second table shows 'Current time: 2016-06-23T22:28:05 +0800', 'Uptime: 00:16', 'Lease expiration enabled: true', 'Renews threshold: 3', and 'Renews (last min): 4'. Below this is the 'DS Replicas' section, which shows '127.0.0.1'. The main section is 'Instances currently registered with Eureka', which contains a table with the following data:

Application	AMIs	Availability Zones	Status
RESERVATION-SERVICE	n/a (1)	(1)	UP (1) - jason-mbp.lan:reservation-service:8000

图

4

3.2.4.3 启动 Service Consumer

启动 Service Consumer 其实和 Service Provider 一样，因为本质上 Eureka 提供的客户端是不区分 Provider 和 Consumer 的，一般情况下，Provider 同时也会是 Consumer。

```
@EnableDiscoveryClient
@SpringBootApplication
public class ReservationClientApplication {
    @Bean
    CommandLineRunner runner(DiscoveryClient dc) {
        return args -> {
            dc.getInstances("reservation-service")
        }
    }
}
```

(下页继续)

(续上页)

```

.forEach(si -> System.out.println(String.format (
"Found %s %s:%s", si.getServiceId(), si.getHost(), si.getPort())));
};
}
public static void main(String[] args) {    SpringApplication.
    ↪run(ReservationClientApplication.class, args);
}
}

```

上述代码中通过 `dc.getInstances(“reservation-service”)` 就能获取到当前所有注册的 reservation-service 服务。

3.3 Eureka Server 实现细节

看了前面的 demo，我们已经初步领略到了 Spring Cloud 和 Eureka 的强大之处，通过短短几行配置就实现了服务注册和发现！

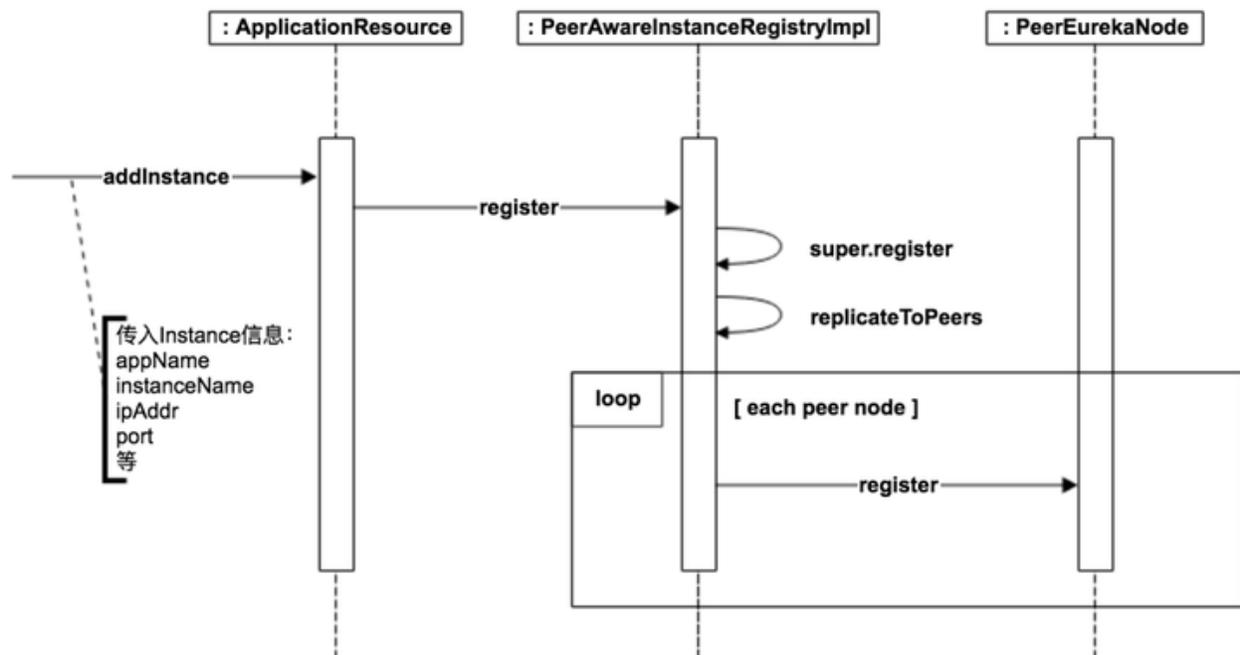
相信大家一定想了解 Eureka 是如何实现的吧，所以接下来我们继续 Dive！首先来看下 Eureka Server 的几个对外接口实现。

3.3.1 Register

首先来看 Register (服务注册)，这个接口会在 Service Provider 启动时被调用来实现服务注册。同时，当 Service Provider 的服务状态发生变化时（如自身检测认为 Down 的时候），也会调用来更新服务状态。

接口实现比较简单，如下图所示。

1. ApplicationResource 类接收 Http 服务请求，调用 PeerAwareInstanceRegistryImpl 的 register 方法
2. PeerAwareInstanceRegistryImpl 完成服务注册后，调用 replicateToPeers 向其它 Eureka Server 节点 (Peer) 做状态同步（异步操作）



图

注册的服务列表保存在一个嵌套的 hash map 中：

- 第一层 hash map 的 key 是 app name，也就是应用名字
- 第二层 hash map 的 key 是 instance name，也就是实例名字

以3.2.4.2中的截图为例，RESERVATION-SERVICE 就是 app name，jason-mbp.lan:reservation-service:8000 就是 instance name。

Hash map 定义如下：

```
private final ConcurrentHashMap<String, Map<String, Lease>> registry =new ConcurrentHashMap<String, Map<String, Lease>>();
```

3.3.2 Renew

Renew（服务续约）操作由 Service Provider 定期调用，类似于 heartbeat。主要是用来告诉 Eureka Server Service Provider 还活着，避免服务被剔除掉。接口实现如下图所示。

可以看到，接口实现方式和 register 基本一致：首先更新自身状态，再同步到其它 Peer。

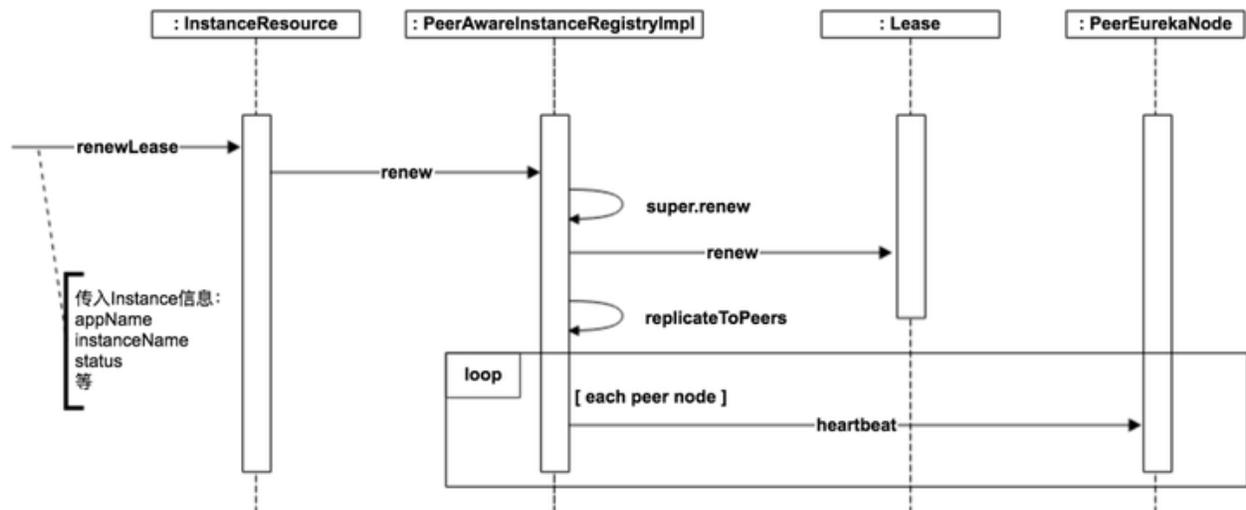
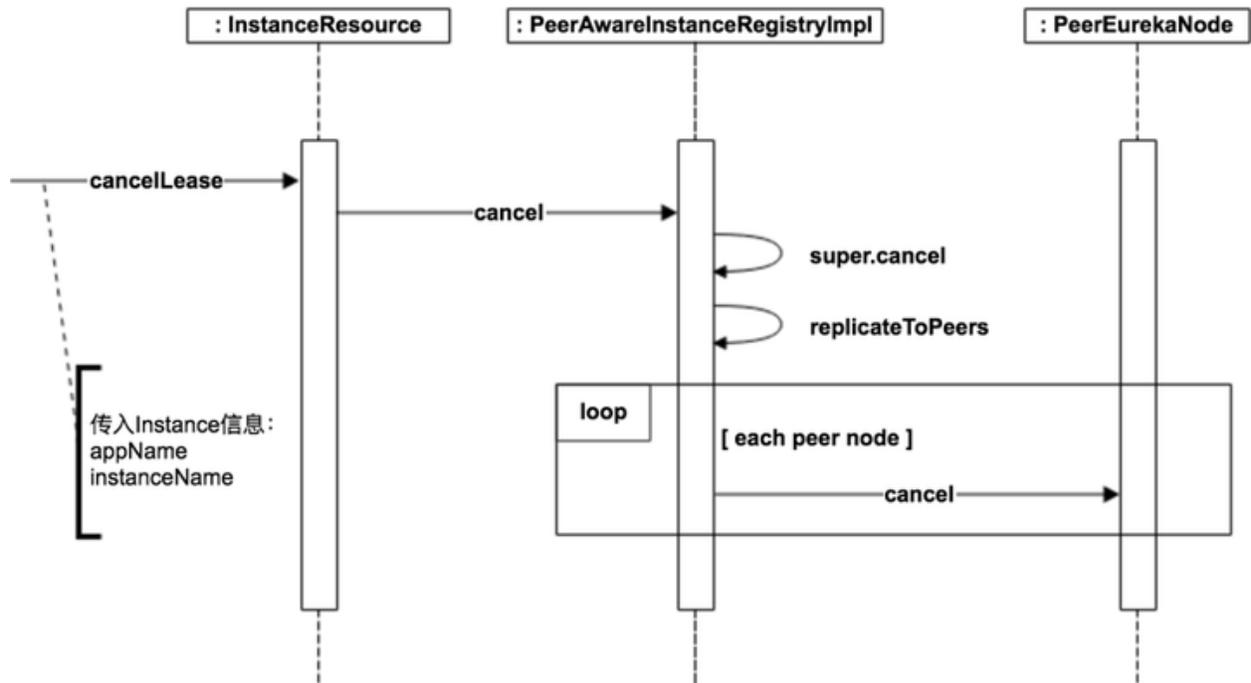


图 6

3.3.3 Cancel

Cancel（服务下线）一般在 Service Provider shut down 的时候调用，用来把自身的服务从 Eureka Server 中删除，以防客户端调用不存在的服务。接口实现如下图所示。



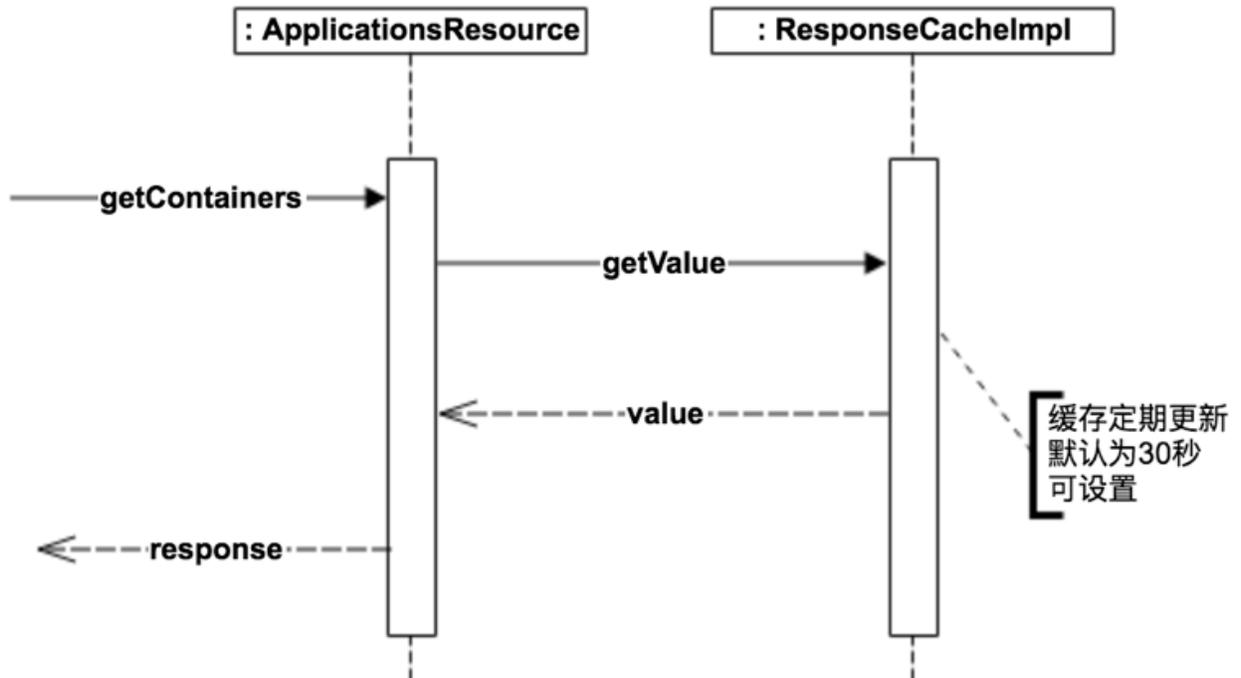
图

7

3.3.4 Fetch Registries

Fetch Registries 由 Service Consumer 调用，用来获取 Eureka Server 上注册的服务。

为了提高性能，服务列表在 Eureka Server 会缓存一份，同时每 30 秒更新一次。



图

8

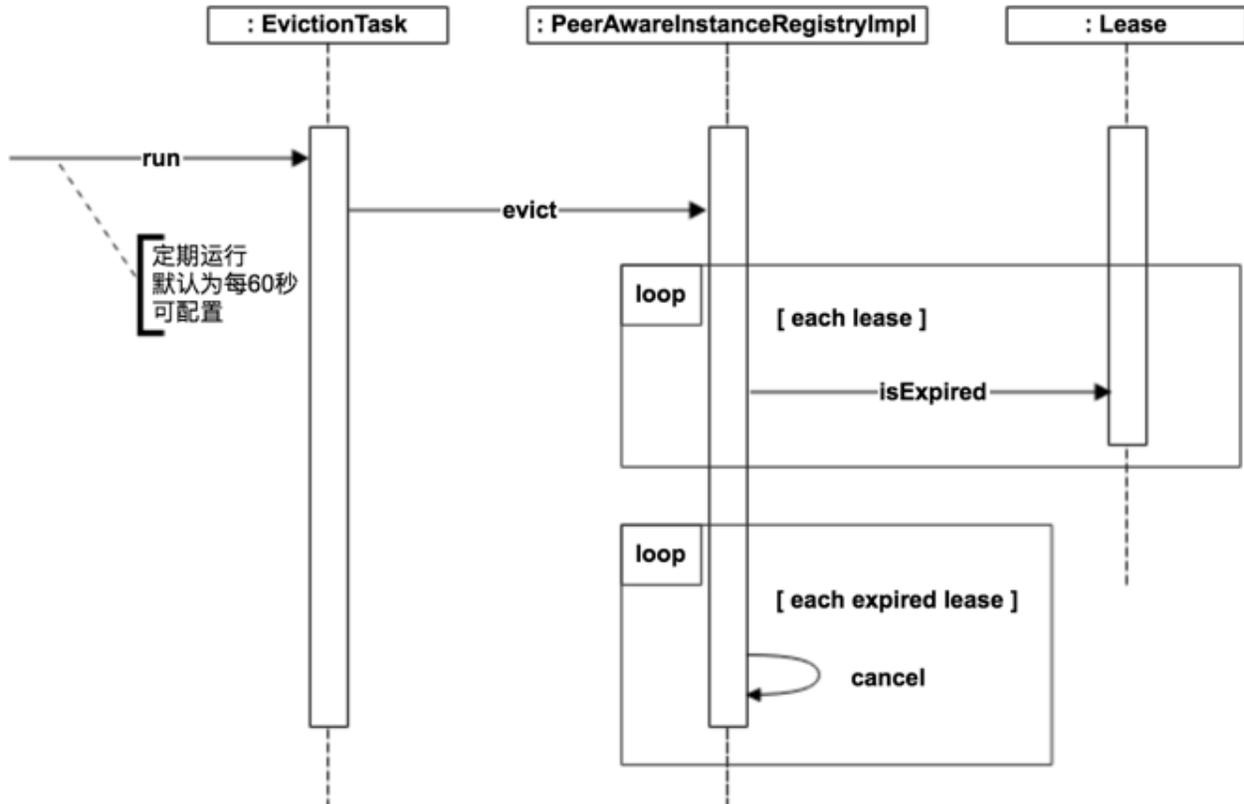
3.3.5 Eviction

Eviction（失效服务剔除）用来定期（默认为每 60 秒）在 Eureka Server 检测失效的服务，检测标准就是超过一定时间没有 Renew 的服务。

默认失效时间为 90 秒，也就是如果有服务超过 90 秒没有向 Eureka Server 发起 Renew 请求的话，就会被当做失效服务剔除掉。

失效时间可以通过 `eureka.instance.leaseExpirationDurationInSeconds` 进行配置，定期扫描时间可以通过 `eureka.server.evictionIntervalTimerInMs` 进行配置。

接口实现逻辑见下图：



图

9

3.3.6 How Peer Replicates

在前面的 Register、Renew、Cancel 接口实现中，我们看到了都会有 `replicateToPeers` 操作，这个就是用来做 Peer 之间的状态同步。

通过这种方式，Service Provider 只需要通知到任意一个 Eureka Server 后就能保证状态会在所有的 Eureka Server 中得到更新。

具体实现方式其实很简单，就是接收到 Service Provider 请求的 Eureka Server，把请求再次转发到其它的 Eureka Server，调用同样的接口，传入同样的参数，除了会在 header 中标记 `isReplication=true`，从而避免重复的 replicate。

Peer 之间的状态是采用异步的方式同步的，所以不保证节点间的状态一定是一致的，不过基本能保证最终状态是一致的。

结合服务发现的场景，实际上也并不需要节点间的状态强一致。在一段时间内（比如 30 秒），节点 A 比节点 B 多一个服务实例或少一个服务实例，在业务上也是完全可以接受的（Service Consumer 侧一般也会实现错误重试和负载均衡机制）。

所以按照 CAP 理论，Eureka 的选择就是放弃 C，选择 AP。

3.3.7 How Peer Nodes are Discovered

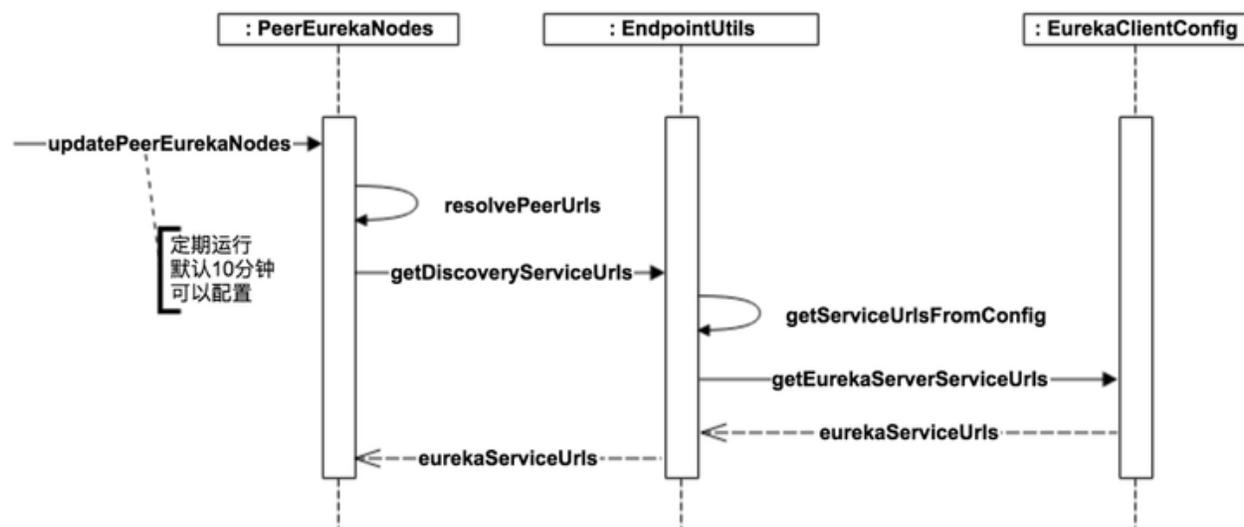
那大家可能会有疑问，Eureka Server 是怎么知道有多少 Peer 的呢？

Eureka Server 在启动后会调用 `EurekaClientConfig.getEurekaServerServiceUrls` 来获取所有的 Peer 节点，并且会定期更新。定期更新频率可以通过 `eureka.server.peerEurekaNodesUpdateIntervalMs` 配置。

这个方法的默认实现是从配置文件读取，所以如果 Eureka Server 节点相对固定的话，可以通过在配置文件中配置来实现。

如果希望能更灵活的控制 Eureka Server 节点，比如动态扩容/缩容，那么可以 override `getEurekaServerServiceUrls` 方法，提供自己的实现，比如我们的项目中会通过数据库读取 Eureka Server 列表。

具体实现如下图所示：



图

10

3.3.8 How New Peer Initializes

最后再来看一下一个新的 Eureka Server 节点加进来，或者 Eureka Server 重启后，如何做初始化，从而能够正常提供服务。

具体实现如下图所示，简而言之就是启动时把自己当做是 Service Consumer 从其它 Peer Eureka 获取所有服务的注册信息。然后对每个服务，在自己这里执行 Register，`isReplication=true`，从而完成初始化。

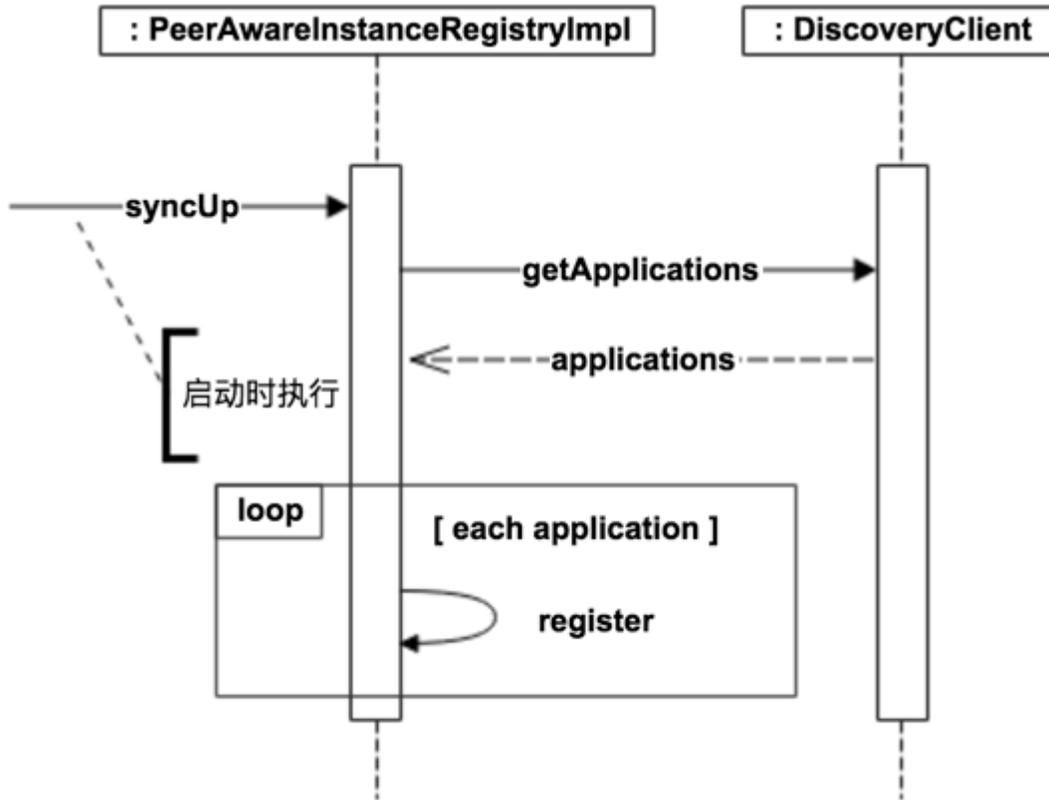


图 11

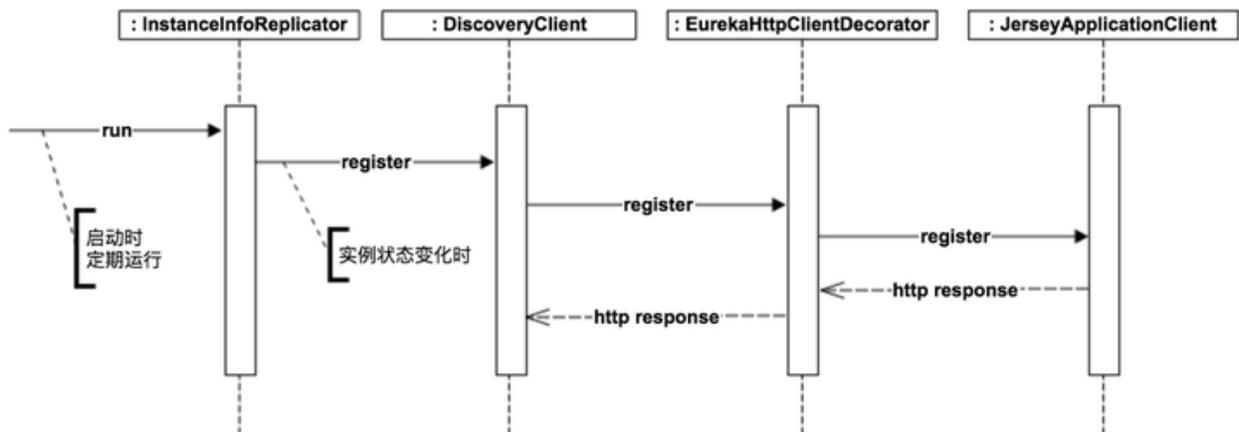
3.4 Service Provider 实现细节

现在来看下 Service Provider 的实现细节，主要就是 Register、Renew、Cancel 这 3 个操作。

3.4.1 Register

Service Provider 要对外提供服务，一个很重要的步骤就是把自己注册到 Eureka Server 上。

这部分的实现比较简单，只需要在启动时和实例状态变化时调用 Eureka Server 的接口注册即可。需要注意的是，需要确保配置 `eureka.client.registerWithEureka=true`。



图

3.4.2 Renew

Renew 操作会在 Service Provider 端定期发起，用来通知 Eureka Server 自己还活着。这里有两个比较重要的配置需要注意一下：

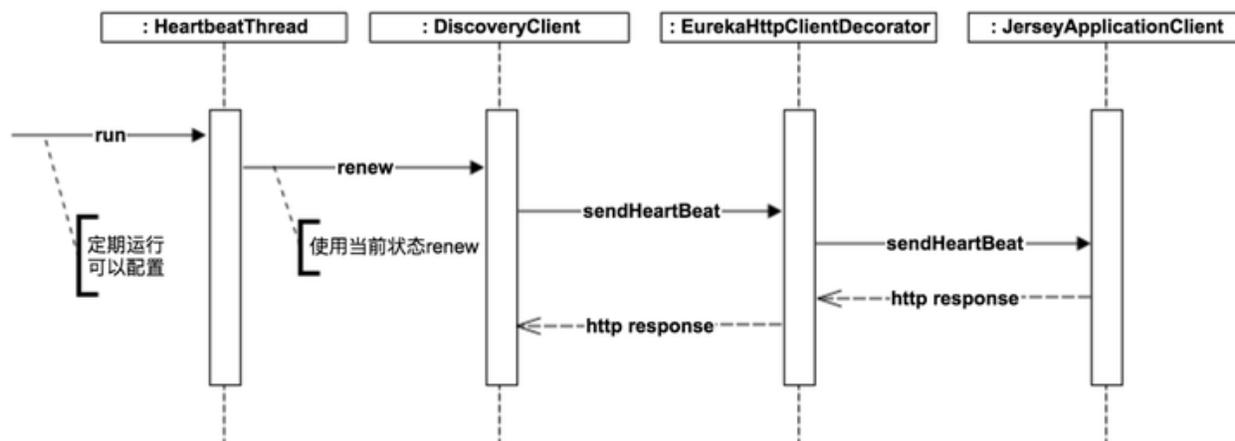
1. `instance.leaseRenewalIntervalInSeconds`

Renew 频率。默认是 30 秒，也就是每 30 秒会向 Eureka Server 发起 Renew 操作。

1. `instance.leaseExpirationDurationInSeconds`

服务失效时间。默认是 90 秒，也就是如果 Eureka Server 在 90 秒内没有接收到来自 Service Provider 的 Renew 操作，就会把 Service Provider 剔除。

具体实现如下：



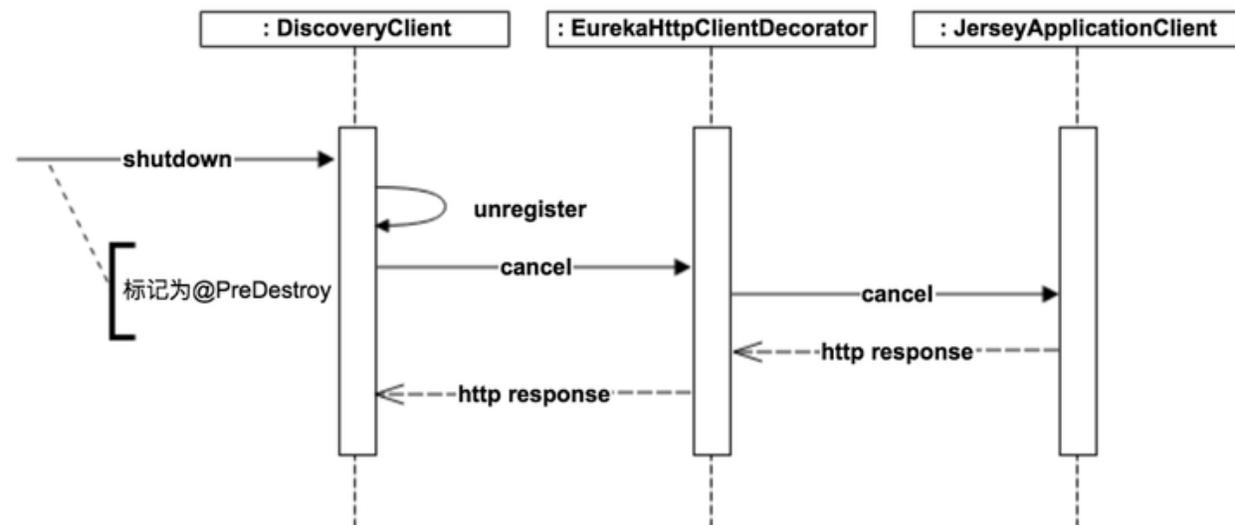
图

13

3.4.3 Cancel

在 Service Provider 服务 shut down 的时候，需要及时通知 Eureka Server 把自己剔除，从而避免客户端调用已经下线的服务。

逻辑本身比较简单，通过对方法标记 `@PreDestroy`，从而在服务 shut down 的时候会被触发。



图

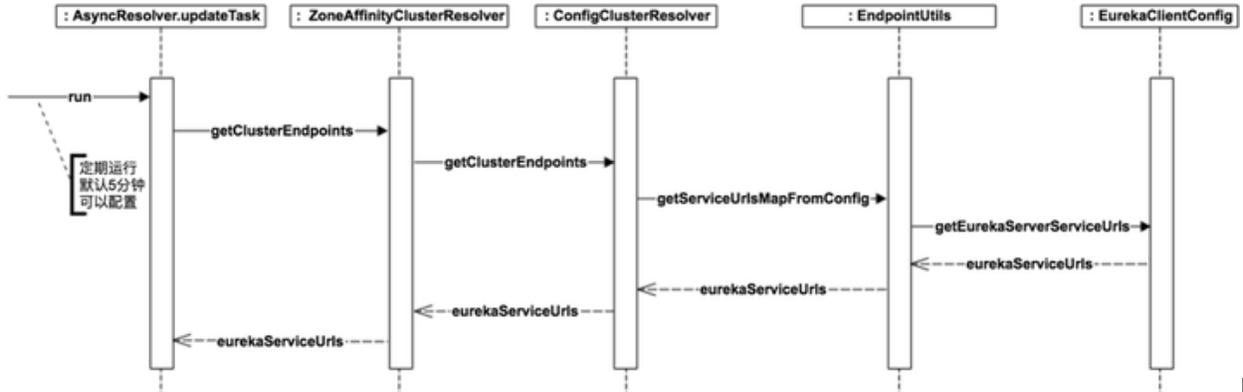
14

3.4.4 How Eureka Servers are Discovered

这里大家疑问又来了，Service Provider 是怎么知道 Eureka Server 的地址呢？

其实这部分的主体逻辑和 3.3.7 How Peer Nodes are Discovered 几乎是一样的。

也是默认从配置文件读取，如果需要更灵活的控制，可以通过 override `getEurekaServerServiceUrls` 方法来提供自己的实现。定期更新频率可以通过 `eureka.client.eurekaServiceUrlPollIntervalSeconds` 配置。



图

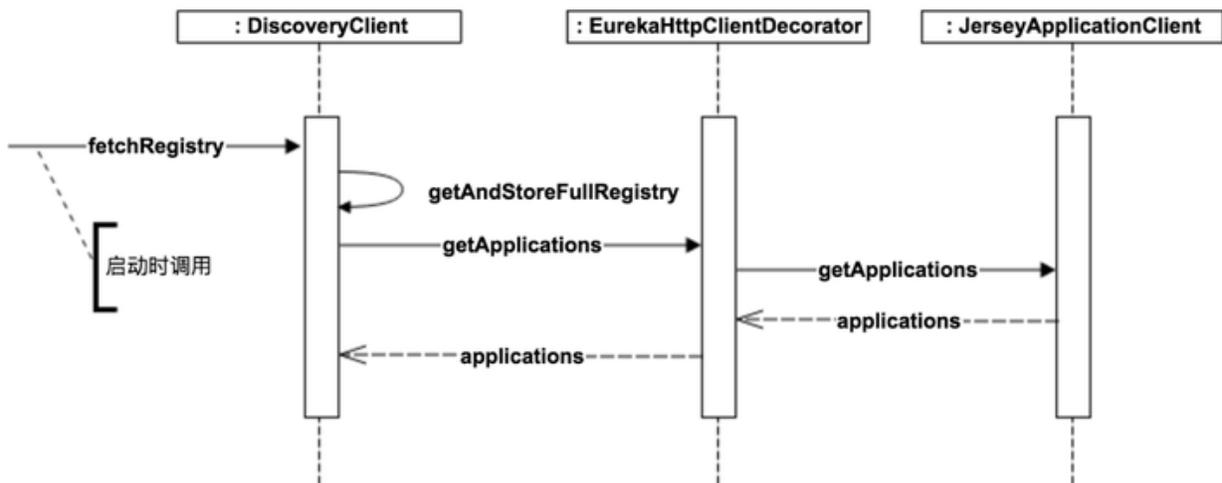
15

3.5 Service Consumer 实现细节

Service Consumer 这块的实现相对就简单一些，因为它只涉及到从 Eureka Server 获取服务列表和更新服务列表。

3.5.1 Fetch Service Registries

Service Consumer 在启动时会从 Eureka Server 获取所有服务列表，并在本地缓存。需要注意的是，需要确保配置 `eureka.client.shouldFetchRegistry=true`。

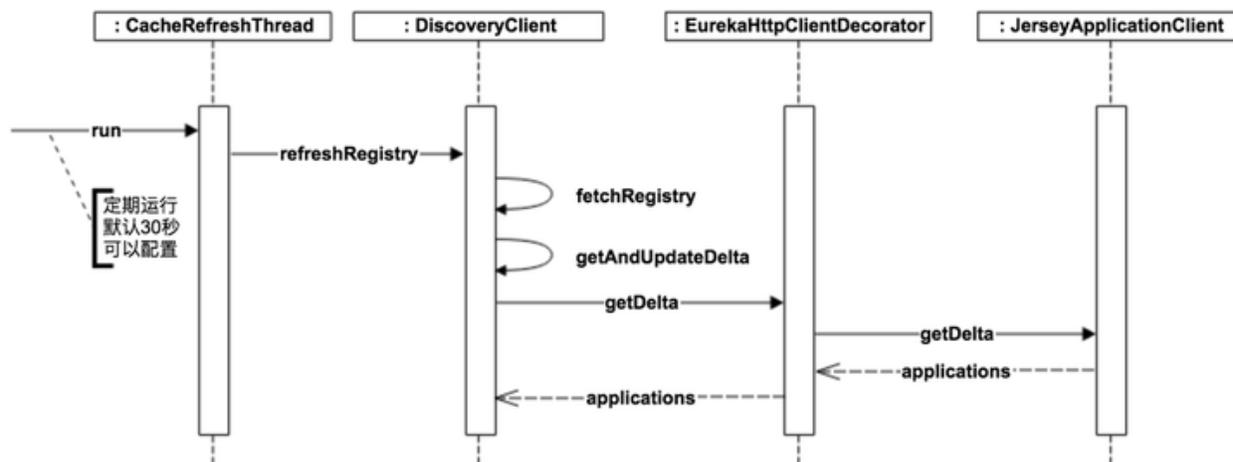


图

16

3.5.2 Update Service Registries

由于在本地有一份缓存，所以需要定期更新，定期更新频率可以通过 `eureka.client.registryFetchIntervalSeconds` 配置。



图

17

3.5.3 How Eureka Servers are Discovered

Service Consumer 和 Service Provider 一样，也有一个如何知道 Eureka Server 地址的问题。

其实由于 Service Consumer 和 Service Provider 本质上是同一个 Eureka 客户端，所以这部分逻辑是一样的，这里就不再赘述了。详细信息见 3.4.4 节。

6.4 Log

6.4.1 日志查看

查找关键日志 grep

Linux 查找文件内容的常用命令方法

从文件内容查找匹配指定字符串的行：`$ grep "被查找的字符串" 文件名`

从文件内容查找与正则表达式匹配的行：`$ grep -e "正则表达式" 文件名`

查找时不区分大小写：`$ grep -i "被查找的字符串" 文件名`

查找匹配的行数：`$ grep -c "被查找的字符串" 文件名`

从文件内容查找不匹配指定字符串的行：`$ grep -v "被查找的字符串" 文件名`

从根目录开始查找所有扩展名为 `.log` 的文本文件，并找出包含 "ERROR" 的行 `find / -type f -
↪name "*.log" | xargs grep "ERROR"`

如果需要查找的内容包含特殊符号，比如 `$` 等等，`grep` 要加参数 `find ./ -name "*.php" | xargs
↪grep -F '要查找的内容'`

linux 查看日志文件内容命令 tail、cat、tac、head、echo

```
tail -f test.log
```

你会看到屏幕不断有内容被打印出来. 这时候中断第一个进程 Ctrl-C,

linux 如何显示一个文件的某几行 (中间几行)

从第 3000 行开始, 显示 1000 行. 即显示 3000~3999 行

```
cat filename | tail -n +3000 | head -n 1000
```

显示 1000 行到 3000 行

```
cat filename | head -n 3000 | tail -n +1000
```

注意两种方法的顺序

tail -n 1000: 显示最后 1000 行 tail -n +1000: 从 1000 行开始显示, 显示 1000 行以后的 head -n 1000: 显示前面 1000 行

用 sed 命令

```
sed -n '5,10p' filename
```

这样你就可以只查看文件的第 5 行到第 10 行。

例: cat mylog.log | tail -n 1000 输出 mylog.log 文件最后一千行

cat 主要有三大功能:

- 1. 一次显示整个文件。\$ cat filename
- 2. 从键盘创建一个文件。\$ cat > filename 只能创建新文件, 不能编辑已有文件.
- 3. 将几个文件合并为一个文件: \$cat file1 file2 > file 参数: -n 或 --number 由 1 开始对所有输出的行数编号 -b 或 --number-nonblank 和 -n 相似, 只不过对于空白行不编号 -s 或 --squeeze-blank 当遇到有连续两行以上的空白行, 就代换为一行的空白行 -v 或 --show-nonprinting 例:

```
把 textfile1 的档案内容加上行号后输入 textfile2 这个档案里
cat -n textfile1 > textfile2
```

```
把 textfile1 和 textfile2 的档案内容加上行号 (空白行不加) 之后将内容附加到 textfile3 里。
```

```
cat -b textfile1 textfile2 >> textfile3
```

```
把 test.txt 文件扔进垃圾箱, 赋空值 test.txt
```

```
cat /dev/null > /etc/test.txt
```

注意: >意思是创建, >>是追加。千万不要弄混了。

tac (反向列示)

tac 是将 cat 反写过来，所以他的功能就跟 cat 相反，cat 是由第一行到最后一行连续显示在萤幕上，而 tac 则是由最后一行到第一行反向在萤幕上显示出来！

在 Linux 中 echo 命令用来在标准输出上显示一段字符，比如：echo "the echo command test!"

这个就会输出 "the echo command test!" 这一行文字！

echo "the echo command test!">a.sh 这个就会在 a.sh 文件中输出 "the echo command test!" 这一行文字！该命令的一般格式为：echo [-n] 字符串其中选项 n 表示输出文字后不换行；字符串能加引号，也能不加引号。用 echo 命令输出加引号的字符串时，将字符串原样输出；用 echo 命令输出不加引号的字符串时，将字符串中的各个单词作为字符串输出，各字符串之间用一个空格分割。

6.5 开发踩坑

1、程序的可用性：

程序的设计很多常见的算法都是空间换时间。增加空间复杂度，减少时间复杂度，提高用户体验度。比如参数过滤，使用 bloom 过滤器，底层需要增加一个非常大的数组进行存储错误的参数值和可能错误的参数值。但是能够过滤大部分的错误参数。

2、错误处理：

程序设计需要对一些考虑到的边界问题进行处理，需要尽可能减少错误率。

例如数据库键重复，数据库底层会爆出键重复的异常，上层虽然有捕捉，但反馈结果是报错。这种设计是不合理的，这种应该在数据的入口就进行排查。程序设计的错误不仅仅是用户级别，也在监控级别。不仅要尽量减少用户反馈的错误，也要尽量减少监控能检测的错误。

真的错误是指程序无法处理的错误，比如系统奔溃。程序的错误一定需要提前考虑并进行处理。

3、技术选型

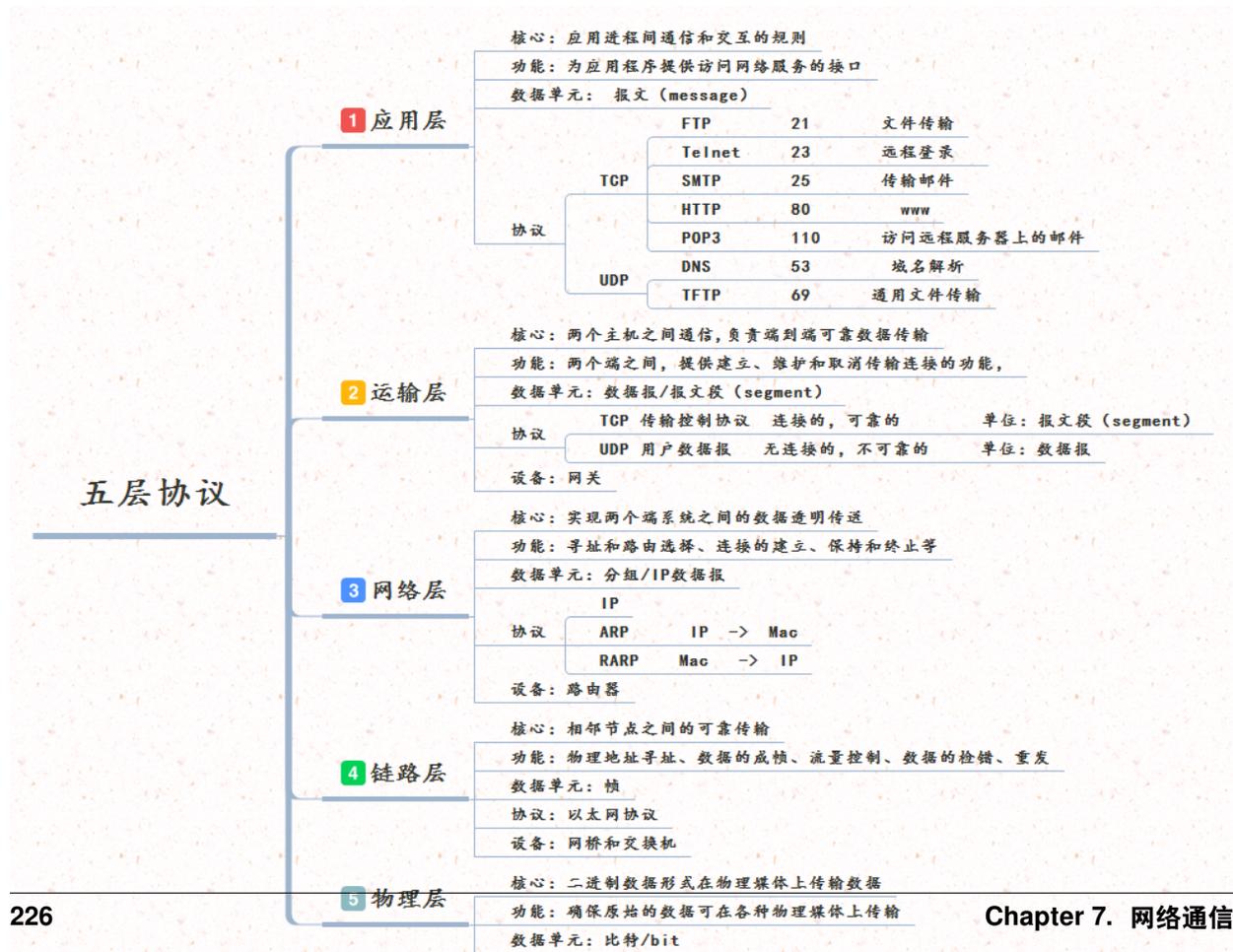
- 场景适配

4、方案设计

- 需求切入点
- 需求切割

7.1 网络协议

7.1.1 计算机网络



7.1.2 HTTP 协议

请求数据

请求码

- 1xx: 指示信息--表示请求已接收, 继续处理
- 2xx: 成功--表示请求已被成功接收、理解、接受
- 3xx: 重定向--要完成请求必须进行更进一步的操作
- 4xx: 客户端错误--请求有语法错误或请求无法实现
- 5xx: 服务器端错误--服务器未能实现合法的请求

常见请求码

200 OK	//客户端请求成功
301 Moved Permanently	//请求永久重定向
302 Moved Temporarily	//请求临时重定向
400 Bad Request	//客户端请求有语法错误, 不能被服务器所理解
401 Unauthorized	//请求未经授权, 这个状态代码必须和WWW-
↪Authenticate报头域一起使用	
403 Forbidden	//服务器收到请求, 但是拒绝提供服务
404 Not Found	//请求资源不存在, eg: 输入了错误的URL
500 Internal Server Error	//服务器发生不可预期的错误
503 Server Unavailable	//服务器当前不能处理客户端的请求, 一段时间后可能恢复正常

请求方法

GET	请求指定的页面信息, 并返回实体主体。
HEAD	类似于get请求, 只不过返回的响应中没有具体的内容, 用于获取报头
POST	↪向指定资源提交数据进行处理请求(例如提交表单或者上传文件)。数据被包含在请求体中。POST请求可能会导↪或已有资源的修改。
PUT	从客户端向服务器传送的数据取代指定的文档的内容。
DELETE	请求服务器删除指定的页面。
CONNECT	HTTP/1.1协议中预留给能够将连接改为管道方式的代理服务器。
OPTIONS	允许客户端查看服务器的性能。
TRACE	回显服务器收到的请求, 主要用于测试或诊断。

HTTP 工作原理

HTTP 协议定义 Web 客户端如何从 Web 服务器请求 Web 页面, 以及服务器如何把 Web 页面传送给客户端。HTTP 协议采用了请求/响应模型。客户端向服务器发送一个请求报文, 请求报文包含请求的方法、URL、协议版本、请求头部和请求数据。服务器以一个状态行作为响应, 响应的内容包括协议的版本、成功或者错误代码、服务器信息、响应头部和响应数据。

以下是 HTTP 请求/响应的步骤:

1. 客户端连接到 Web 服务器一个 HTTP 客户端, 通常是浏览器, 与 Web 服务器的 HTTP 端口 (默认为 80) 建立一个 TCP 套接字连接。例如, <http://www.luffycity.com>。

2. 发送 HTTP 请求通过 TCP 套接字, 客户端向 Web 服务器发送一个文本的请求报文, 一个请求报文由请求行、请求头部、空行和请求数据 4 部分组成。
3. 服务器接受请求并返回 HTTP 响应 Web 服务器解析请求, 定位请求资源。服务器将资源副本写到 TCP 套接字, 由客户端读取。一个响应由状态行、响应头部、空行和响应数据 4 部分组成。
4. 释放连接 TCP 连接若 connection 模式为 close, 则服务器主动关闭 TCP 连接, 客户端被动关闭连接, 释放 TCP 连接; 若 connection 模式为 keepalive, 则该连接会保持一段时间, 在该时间内可以继续接收请求;
5. 客户端浏览器解析 HTML 内容客户端浏览器首先解析状态行, 查看表明请求是否成功的状态代码。然后解析每一个响应头, 响应头告知以下为若干字节的 HTML 文档和文档的字符集。客户端浏览器读取响应数据 HTML, 根据 HTML 的语法对其进行格式化, 并在浏览器窗口中显示。

例如: 在浏览器地址栏键入 URL, 按下回车之后会经历以下流程:

1. 浏览器向 DNS 服务器请求解析该 URL 中的域名所对应的 IP 地址;
2. 解析出 IP 地址后, 根据该 IP 地址和默认端口 80, 和服务器建立 TCP 连接;
3. 浏览器发出读取文件 (URL 中域名后面部分对应的文件) 的 HTTP 请求, 该请求报文作为 TCP 三次握手的第三个报文的数据发送给服务器;
4. 服务器对浏览器请求作出响应, 并把对应的 html 文本发送给浏览器;
5. 释放 TCP 连接;
6. 浏览器将该 html 文本并显示内容;

参考如下:

[HTTP 协议超级详解](#)

[HTTP 详解](#)

7.1.3 websocket

建立在 TCP 协议之上

WebSocket 协议建立在 http 协议的基础之上, 因此二者有很多类似之处。事实上, 在使用 websocket 协议时, 浏览器与服务端最开始建立的还是 http 连接, 之后再讲协议从 http 转换成 websocket, 协议转换的过程称之为握手 (handshake), 表示服务端与客户端都同意建立 websocket 协议。

请求

```
GET ws://server.example.com/ws HTTP/1.1
Host: server.example.com
Upgrade: websocket
Connection: Upgrade
Sec-WebSocket-Key: x3JJHMbDL1EzLkh9GBhXDw==
Sec-WebSocket-Protocol: chat, superchat
Sec-WebSocket-Version: 13
Origin: http://example.com
```

Upgrade 字段表示将通信协议从 HTTP/1.1 转向该字段指定的协议。Connection 字段表示浏览器通知服务器, 如果可以的话, 就升级到 WebSocket 协议。Origin 字段用于提供请求发出的域名, 供服务器验证是否许可的范围内 (服务器也可以不验证)。Sec-WebSocket-Key 则是用于握手协议的密钥, 是 Base64 编码的 16 字节随机字符串。

- Payload data: (x + y) Byte。” extension data” 后接 “application data”。

参考如下:

WebSocket 浅析

Netty-Websocket 根据 URL 路由, 分发机制的实现

netty

WebSocket 协议深入探究

学习 WebSocket 协议—从顶层到底层的实现原理 (修订版)

基于 Redis 以及 WebSocket 的一个实时消息推送系统

7.2 IO 通信

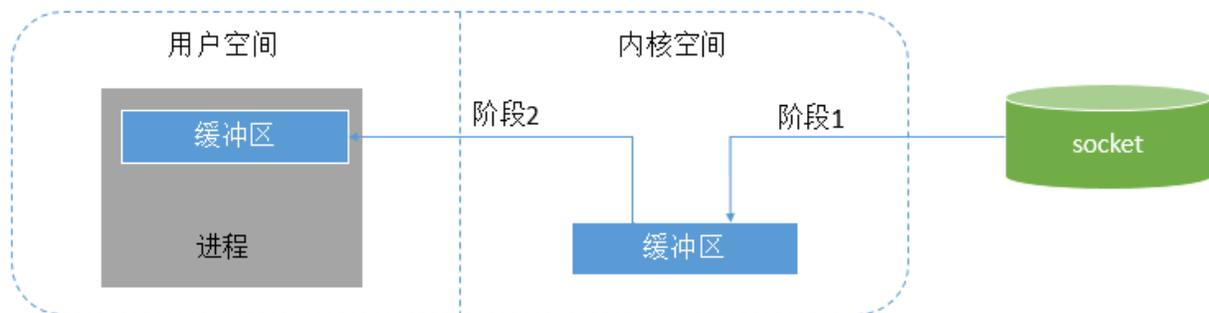
7.2.1 网络 I/O 模型

数据操作

对于一个 network IO (以 read 举例), 它会涉及到两个系统对象: 一个是调用这个 IO 的进程, 另一个就是系统内核 (kernel)。当一个 read 操作发生时, 它会经历两个阶段:

阶段 1: 等待数据准备 (Waiting for the data to be ready)

阶段 2: 将数据从内核拷贝到进程中 (Copying the data from the kernel to the process)



用户空间是常规进程所在区域。JVM 就是常规进程, 驻守于用户空间。用户空间是非特权区域: 比如, 在该区域执行的代码就不能直接访问硬件设备。

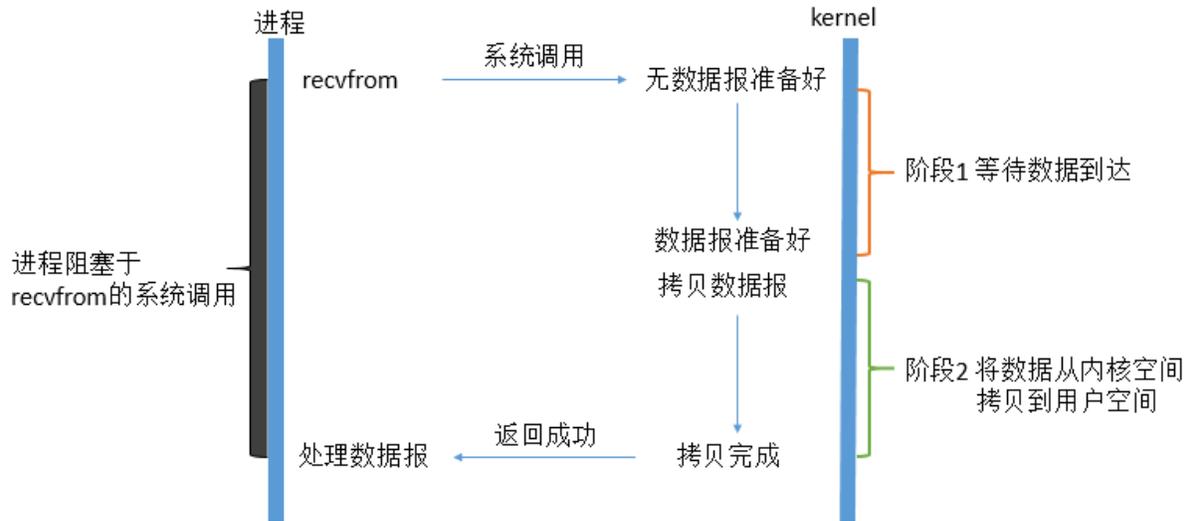
内核空间是操作系统所在区域。内核代码有特别的权力: 它能与设备控制器通讯, 控制着用户区域进程的运行状态, 等等。最重要的是, 所有 I/O 都直接 (如这里所述) 或间接通过内核空间。

当进程请求 I/O 操作的时候, 它执行一个系统调用将控制权移交给内核。C/C++ 程序员所熟知的底层函数 open()、read()、write() 和 close() 要做的无非就是建立和执行适当的系统调用。当内核以这种方式被调用, 它随即采取任何必要步骤, 找到进程所需数据, 并把数据传送到用户空间内的指定缓冲区。内核试图对数据进行高速缓存或预读取, 因此进程所需数据可能已经在内核空间里了。如果是这样, 该数据只需简单地拷贝出来即可。如果数据不在内核空间, 则进程被挂起, 内核着手把数据读进内存。

1、阻塞式 I/O: blocking IO

在 linux 中，默认情况下所有的 socket 都是 blocking，一个典型的读操作流程大概是这样：

- 第一步通常涉及等待数据从网络中到达。当所有等待数据到达时，它被复制到内核中的某个缓冲区。
- 第二步就是把数据从内核缓冲区复制到应用程序缓冲区。

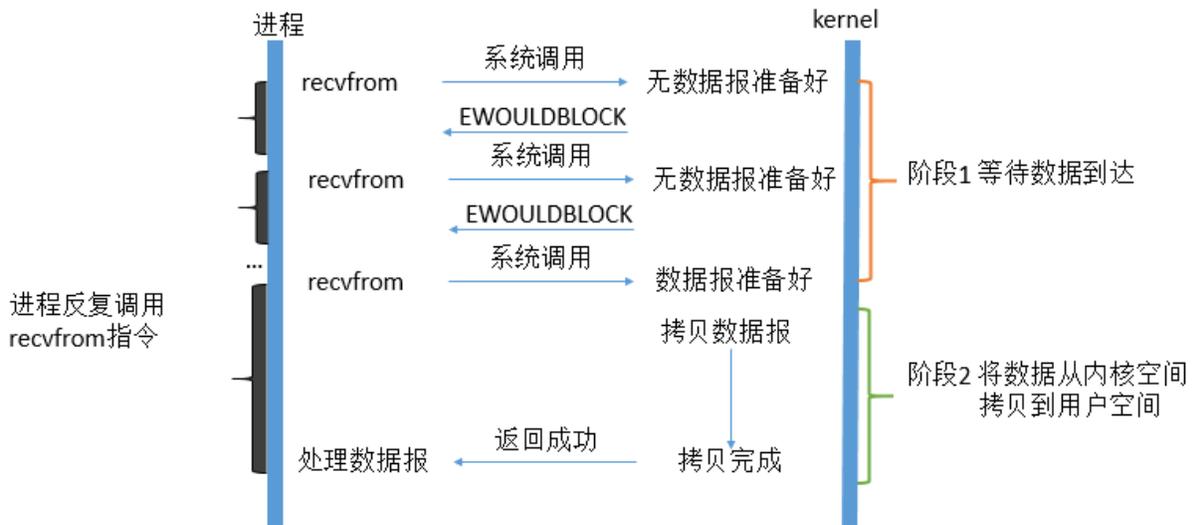


当用户进程调用了 `recvfrom` 这个系统调用，kernel 就开始了 IO 的第一个阶段：准备数据。对于 network io 来说，很多时候数据在一开始还没有到达（比如，还没有收到一个完整的 UDP 包），这个时候 kernel 就要等待足够的数据到来。而在用户进程这边，整个进程会被阻塞。当 kernel 一直等到数据准备好了，它就会将数据从 kernel 中拷贝到用户内存，然后 kernel 返回结果，用户进程才解除 block 的状态，重新运行起来。

所以，blocking IO 的特点就是在 IO 执行的两个阶段都被 block 了。

2、非阻塞式 I/O: nonblocking IO

linux 下，可以通过设置 socket 使其变为 non-blocking。当对一个 non-blocking socket 执行读操作时，流程是这个样子：



从图中可以看出，当用户进程发出 read 操作时，如果 kernel 中的数据还没有准备好，那么它并不会 block 用户进程，而是立刻返回一个 error。从用户进程角度讲，它发起一个 read 操作后，并不需要等待，而是马上就得到了一个结果。用户进程判断结果是一个 error 时，它就知道数据还没有准备好，于是它可以再次发送 read 操作。一旦 kernel 中的数据准备好了，并且又再次收到了用户进程的 system call，那么它马上就将数据拷贝到了用户内存，然后返回。

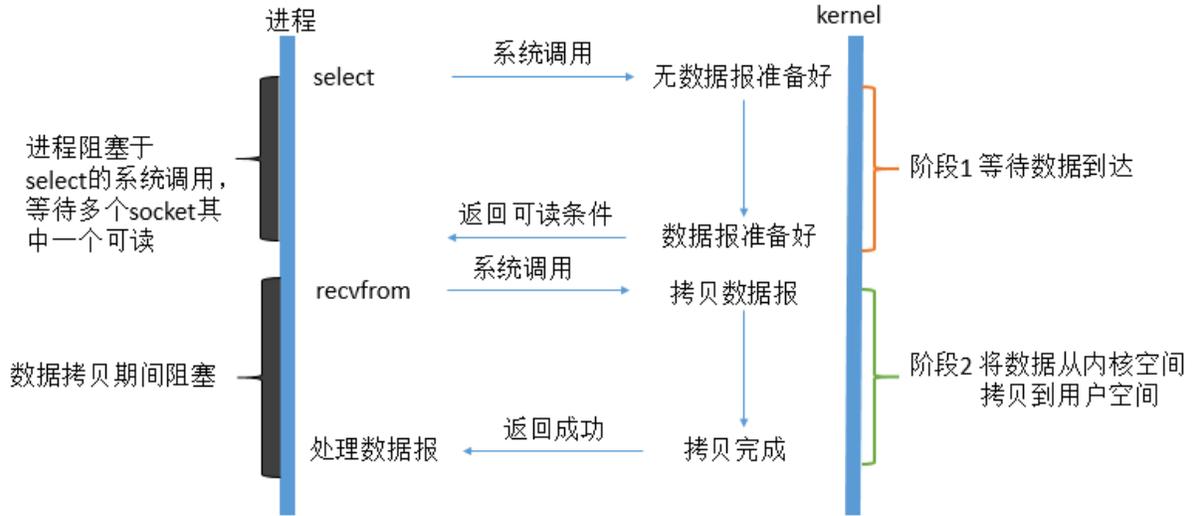
所以，用户进程第一个阶段不是阻塞的，需要不断的主动询问 kernel 数据好了没有；第二个阶段依然总是阻塞的。

3、I/O 复用 (select, poll, epoll...): IO multiplexing

IO 复用同非阻塞 IO 本质一样，不过利用了新的 select 系统调用，由内核来负责本来是请求进程该做的轮询操作。看似比非阻塞 IO 还多了一个系统调用开销，不过因为可以支持多路 IO，才算提高了效率。

基于内核，建立在 epoll 或者 kqueue 上实现，I/O 多路复用最大的优势是用户可以在一个线程内同时处理多个 Socket 的 I/O 请求。通过一个线程监听全部的 TCP 连接，有任何事件发生就通知用户态处理即可。

它的基本原理就是 select /epoll 这个 function 会不断的轮询所负责的所有 socket，当某个 socket 有数据到达了，就通知用户进程。它的流程如图：



当用户进程调用了 `select`，那么整个进程会被 `block`，而同时，`kernel` 会“监视”所有 `select` 负责的 `socket`，当任何一个 `socket` 中的数据准备好了，`select` 就会返回。这个时候用户进程再调用 `read` 操作，将数据从 `kernel` 拷贝到用户进程

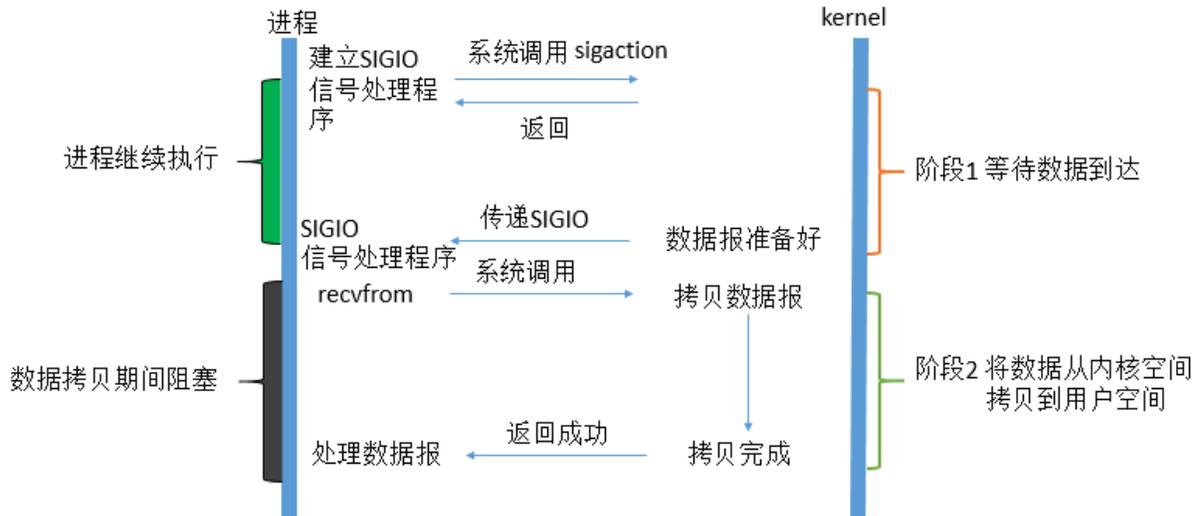
这里需要使用两个 `system call` (`select` 和 `recvfrom`)，而 `blocking IO` 只调用了一个 `system call` (`recvfrom`)。但是，用 `select` 的优势在于它可以同时处理多个 `connection`。（多说一句。所以，如果处理的连接数不是很高的话，使用 `select/epoll` 的 `web server` 不一定比使用 `multi-threading + blocking IO` 的 `web server` 性能更好，可能延迟还更大。

select/epoll 的优势并不是对于单个连接能处理得更快，而是在于能处理更多的连接。

实际中，对于每一个 `socket`，一般都设置成为 `non-blocking`，但是，如上图所示，整个用户的 `process` 其实是一直被 `block` 的。只不过 `process` 是被 `select` 这个函数 `block`，而不是被 `socket IO` 给 `block`。

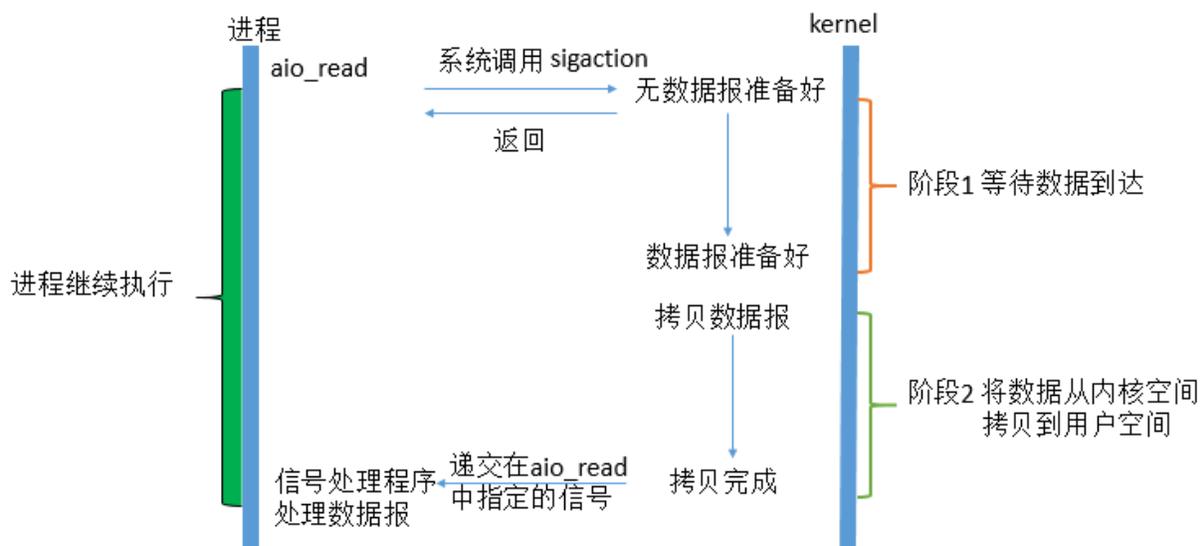
4、信号驱动式 I/O (SIGIO)：signal driven IO

用的很少



5、异步 I/O (POSIX 的 aio_ 系列函数): asynchronous IO

这类函数的工作机制是告知内核启动某个操作，并让内核在整个操作（包括将数据从内核拷贝到用户空间）完成后通知我们。如图：



用户进程发起 read 操作之后，立刻就可以开始去做其它的事。而另一方面，从 kernel 的角度，当它受到一个 asynchronous read 之后，首先它会立刻返回，所以不会对用户进程产生任何 block。然后，kernel 会等待数据准备完成，然后将数据拷贝到用户内存，当这一切都完成之后，kernel 会给用户进程发送一个 signal，告诉它 read 操作完成了。在这整个过程中，进程完全没有被 block。

I/O 多路复用往往对应 Reactor 模式，异步 I/O 往往对应 Proactor。

总结

前四种 I/O 模型都是同步 I/O 操作，他们的区别在于第一阶段，而他们的第二阶段是一样的：在数据从内核复制到应用缓冲区期间（用户空间），进程阻塞于 recvfrom 调用。

参考

Unix 五种 IO 模型

7.3 通信框架

7.3.1 Netty 实现 WebSocket

处理 WebSocket frame

WebSockets 在“帧”里面来发送数据，其中每一个都代表了一个消息的一部分。一个完整的消息可以利用了多个帧。WebSocket “Request for Comments” (RFC) 定义了六中不同的 frame; Netty 给他们每个都提供了一个 POJO 实现

- CloseWebSocketFrame

- PingWebSocketFrame
- PongWebSocketFrame
- TextWebSocketFrame
- BinaryWebSocketFrame

其中 Close、Ping 以及 Pong 称之为控制帧，Close 关闭帧很容易理解，客户端如果接受到了就关闭连接，客户端也可以发送关闭帧给服务端。Ping 和 Pong 是 websocket 里的心跳，用来保证客户端是在线的。Text 和 Binary 表示这个传输的帧是文本类型还是二进制类型，二进制类型传输的数据可以是图片或者语音之类的。

如何实现消息推送

之前的案例，还是通过客户端发送请求，服务端进行响应的方式进行交互。WebSocket 协议还可以实现服务端主动向客户端推送消息的功能。

主动推送需要考虑的问题包括：

1、服务端要保存用户唯一标记(假设为 `userId`) 与其对应的 `SocketChannel` 的对应关系，以便之后根据这个唯一标记给用户推送消息。

对应关系的保存很简单，用一个 Map 来记录即可。保存这种对应关系的最佳时机是 websocket 连接刚刚建立的时候，因为这种对应关系只需要保存一次，之后就可以拿着这个 `userId` 找到对应 `SocketChannel` 给用户推送消息。以之前的代码为例，很明显这个操作，应该放在 `handleHttpRequest` 方法中处理，因为这个方法在一个 websocket 连接的生命周期只会调用一次。

我们知道，当连接建立时，我们在服务端就能获取到客户端对应的 `SocketChannel`，我们如果获取到 `userId` 呢？WebSocket 协议与 HTTP 协议一样，都支持在 url 中添加参数。如：

```
String userId=params.get("userId").get(0);
//保存userId和Channel的对应关系
map.put(userId,ctx.channel());
```

2、通过 `userId` 给用户推送消息

通常消息推送是由运营人员根据业务需要，从一个后台管理界面，过滤出符合特定条件的用户来推送消息，因此需要推送的 `userId` 和推送内容都是运营指定的。推送服务端只需要根据 `userId` 找到对应的 `Channel`，将推送消息推送给浏览器即可。代码片段如下：

```
//运营人员制定要推送的userId和消息
String userId="123456";
String msg="xxxxxxxxxxxxxxxxxxxxxxxx";

//根据userId找到对应的channel，并将消息写出
Channel channel = map.get(userId);
channel.writeAndFlush(new TextWebSocketFrame(msg));
```

一般公司的推送的消息量都比较大，所以需要推送的消息一般会放到一个外部消息队列中，如 `rocketmq`，`kafka`。推送服务端通过消费这些队列，来给用户推送消息。

3、客户端 `ack` 接受到的消息

一般情况下，我们会每一个推送的消息都设置一个唯一的 `msgId`。客户端在接受推送的消息后，需要对这条消息进行 `ack`，也就是告诉服务端自己接收到了这条消息，`ack` 时需要将 `msgId` 带回来。

如果我们只使用 websocket 做推送服务，那么之前代码中的 `handleWebSocketFrame` 方法，要处理的主要就是这类 `ack` 消息。

4、离线消息

如果当给某个用户推送消息的时候，其并不在线，可以将消息保存下来，一般都是保存到一个缓存服务器中(如 redis)，每次当一个用户与服务端建立连接的时候，首先检查缓存服务器中有没有其对应的离线消息，如果有直接取出来发送给用户。

5、消息推送记录的存储

服务端推送完成消息记录之后，还需要将这些消息存储下来，以支持一个用户查看自己曾经受到过的推送消息。

6、疲劳度控制

在实现推送功能时，我们通常还会做一个疲劳度控制的功能，也就是限制给一个用户每天推送消息的数量，以免频繁的推送消息给用户造成不好的体验。疲劳度控制可以从多个维度进行设计。例如用户一天可以接受到消息总数，根据业务划分的某种特定类型的消息的一天可以接受到的消息总数等。

参考如下

Netty 实现 WebSocket 聊天功能

WebSocket 协议

7.4 Nginx

7.4.1 Nginx

Nginx 的模块根据其功能基本上可以分为以下几种类型：

- **event module:** 搭建了独立于操作系统的事件处理机制的框架，及提供了各具体事件的处理。包括 ngx_events_module, ngx_event_core_module 和 ngx_epoll_module 等。Nginx 具体使用何种事件处理模块，这依赖于具体的操作系统和编译选项。
- **phase handler:** 此类型的模块也被直接称为 handler 模块。主要负责处理客户端请求并产生待响应内容，比如 ngx_http_static_module 模块，负责客户端的静态页面请求处理并将对应的磁盘文件准备为响应内容输出。
- **output filter:** 也称为 filter 模块，主要是负责对输出的内容进行处理，可以对输出进行修改。例如，可以实现对输出的所有 html 页面增加预定义的 footer 一类的工作，或者对输出的图片的 URL 进行替换之类的工作。
- **upstream:** upstream 模块实现反向代理的功能，将真正的请求转发到后端服务器上，并从后端服务器上读取响应，发回客户端。upstream 模块是一种特殊的 handler，只不过响应内容不是真正由自己产生的，而是从后端服务器上读取的。
- **load-balancer:** 负载均衡模块，实现特定的算法，在众多的后端服务器中，选择一个服务器出来作为某个请求的转发服务器。

example:

```
upstream getsvr {
    #server 13.231.34.83:1317;#请求rest地址
    server 192.168.13.125:26659;#请求rest地址
}

server {
    listen 7777;
    location / {
        add_header Access-Control-Allow-Origin 'http://192.168.71.154:5200';
    }
}
```

↔#需要设置前端部署的地址 (下页继续)

(续上页)

```

        add_header Access-Control-Allow-Methods 'POST, GET, OPTIONS';
        add_header Access-Control-Allow-Headers 'X-Requested-With,Content-
↪Type,origin,content-type,accept,authorization,Action,Module,access-control-
↪allow-origin,app-type,timeout,devid';
        add_header Access-Control-Allow-Credentials true;
        if ($request_method = 'OPTIONS') {
            return 204;
        }
        proxy_pass http://getsvr;
    }
}

```

Nginx 使用一个多进程模型来对外提供服务，其中一个 master 进程，多个 worker 进程。master 进程负责管理 Nginx 本身和其他 worker 进程。

所有实际上的业务处理逻辑都在 worker 进程。worker 进程中有一个函数，执行无限循环，不断处理收到的来自客户端的请求，并进行处理，直到整个 Nginx 服务被停止。

worker 进程中，`ngx_worker_process_cycle()` 函数就是这个无限循环的处理函数。在这个函数中，一个请求的简单处理流程如下：

- 操作系统提供的机制（例如 `epoll`, `kqueue` 等）产生相关的事件。
- 接收和处理这些事件，如是接受到数据，则产生更高层的 `request` 对象。
- 处理 `request` 的 `header` 和 `body`。
- 产生响应，并发送回客户端。
- 完成 `request` 的处理。
- 重新初始化定时器及其他事件。

请求的处理流程

为了让大家更好的了解 Nginx 中请求处理过程，我们以 HTTP Request 为例，来做一下详细地说明。

从 Nginx 的内部来看，一个 HTTP Request 的处理过程涉及到以下几个阶段。

- 初始化 HTTP Request（读取来自客户端的数据，生成 HTTP Request 对象，该对象含有该请求所有的信息）。
- 处理请求头。
- 处理请求体。
- 如果有的话，调用与此请求（URL 或者 Location）关联的 `handler`。
- 依次调用各 `phase handler` 进行处理。

在这里，我们需要了解一下 `phase handler` 这个概念。`phase` 字面的意思，就是阶段。所以 `phase handlers` 也就好理解了，就是包含若干个处理阶段的一些 `handler`。

在每一个阶段，包含有若干个 `handler`，再处理到某个阶段的时候，依次调用该阶段的 `handler` 对 HTTP Request 进行处理。

通常情况下，一个 `phase handler` 对这个 `request` 进行处理，并产生一些输出。通常 `phase handler` 是与定义在配置文件中的某个 `location` 相关联的。

一个 `phase handler` 通常执行以下几项任务：

- 获取 `location` 配置。

- 产生适当的响应。
- 发送 response header。
- 发送 response body。

当 Nginx 读取到一个 HTTP Request 的 header 的时候, Nginx 首先查找与这个请求关联的虚拟主机的配置。如果找到了这个虚拟主机的配置, 那么通常情况下, 这个 HTTP Request 将会经过以下几个阶段的处理 (phase handlers):

- NGX_HTTP_POST_READ_PHASE: 读取请求内容阶段
- NGX_HTTP_SERVER_REWRITE_PHASE: Server 请求地址重写阶段
- NGX_HTTP_FIND_CONFIG_PHASE: 配置查找阶段:
- NGX_HTTP_REWRITE_PHASE: Location 请求地址重写阶段
- NGX_HTTP_POST_REWRITE_PHASE: 请求地址重写提交阶段
- NGX_HTTP_PREACCESS_PHASE: 访问权限检查准备阶段
- NGX_HTTP_ACCESS_PHASE: 访问权限检查阶段
- NGX_HTTP_POST_ACCESS_PHASE: 访问权限检查提交阶段
- NGX_HTTP_TRY_FILES_PHASE: 配置项 try_files 处理阶段
- NGX_HTTP_CONTENT_PHASE: 内容产生阶段
- NGX_HTTP_LOG_PHASE: 日志模块处理阶段

在内容产生阶段, 为了给一个 request 产生正确的响应, Nginx 必须把这个 request 交给一个合适的 content handler 去处理。如果这个 request 对应的 location 在配置文件中被明确指定了一个 content handler, 那么 Nginx 就可以通过对 location 的匹配, 直接找到这个对应的 handler, 并把这个 request 交给这个 content handler 去处理。这样的配置指令包括像, perl, flv, proxy_pass, mp4 等。

如果一个 request 对应的 location 并没有直接有配置的 content handler, 那么 Nginx 依次尝试:

- 如果一个 location 里面有配置 random_index on, 那么随机选择一个文件, 发送给客户端。
- 如果一个 location 里面有配置 index 指令, 那么发送 index 指令指明的文件, 给客户端。
- 如果一个 location 里面有配置 autoindex on, 那么就发送请求地址对应的服务端路径下的文件列表给客户端。
- 如果这个 request 对应的 location 上有设置 gzip_static on, 那么就查找是否有对应的.gz 文件存在, 有的话, 就发送这个给客户端 (客户端支持 gzip 的情况下)。
- 请求的 URI 如果对应一个静态文件, static module 就发送静态文件的内容到客户端。内容产生阶段完成以后, 生成的输出会被传递到 filter 模块去进行处理。filter 模块也是与 location 相关的。所有的 filter 模块都被组织成一条链。输出会依次穿越所有的 filter, 直到有一个 filter 模块的返回值表明已经处理完成。

这里列举几个常见的 filter 模块, 例如:

- server-side includes。
- XSLT filtering。
- 图像缩放之类的。
- gzip 压缩。

在所有的 filter 中, 有几个 filter 模块需要关注一下。按照调用的顺序依次说明如下:

- write: 写输出到客户端, 实际上是写到连接对应的 socket 上。
- postpone: 这个 filter 是负责 subrequest 的, 也就是子请求的。

- copy: 将一些需要复制的 buf(文件或者内存) 重新复制一份然后交给剩余的 body filter 处理。

参考如下：深入浅出 NginxNginx 入门指南

7.4.2 nginx 代理跨域

nginx 反向代理接口跨域

跨域原理：同源策略是浏览器的安全策略，不是 HTTP 协议的一部分。服务器端调用 HTTP 接口只是使用 HTTP 协议，不会执行 JS 脚本，不需要同源策略，也就不存在跨越问题。

实现思路：通过 nginx 配置一个代理服务器（域名与 domain1 相同，端口不同）做跳板机，反向代理访问 domain2 接口，并且可以顺便修改 cookie 中 domain 信息，方便当前域 cookie 写入，实现跨域登录。

nginx 具体配置：

```
#proxy服务器
server {
    listen      81;
    server_name www.domain1.com;

    location / {
        proxy_pass      http://www.domain2.com:8080; #反向代理
        proxy_cookie_domain www.domain2.com www.domain1.com; #修改cookie里域名
        index index.html index.htm;

        # 当用webpack-dev-
        ↪server等中间件代理接口访问nginx时，此时无浏览器参与，故没有同源限制，下面的跨域配置可不启用
        add_header Access-Control-Allow-Origin http://www.domain1.com;
        ↪#当前端只跨域不带cookie时，可为*
        add_header Access-Control-Allow-Credentials true;
    }
}
```

配置

```
location / {
    add_header Access-Control-Allow-Origin *;
    add_header Access-Control-Allow-Headers "Origin, X-Requested-With, Content-Type, ↵
    ↪Accept";
    add_header Access-Control-Allow-Methods "GET, POST, OPTIONS";
}
```

1. Access-Control-Allow-Origin

服务器默认是不被允许跨域的。给 Nginx 服务器配置 Access-Control-Allow-Origin * 后，表示服务器可以接受所有的请求源（Origin），即接受所有跨域的请求。

2. Access-Control-Allow-Headers

是为了防止出现以下错误:

```
Request header field Content-Type is not allowed by Access-Control-Allow-Headers in
↪preflight response.
```

这个错误表示当前请求 Content-Type 的值不被支持。其实是我们发起了”application/json” 的类型请求导致的。

3. Access-Control-Allow-Methods

是为了防止出现以下错误:

```
Content-Type is not allowed by Access-Control-Allow-Headers in preflight response.
```

发送”预检请求”时, 需要用到方法 OPTIONS, 所以服务器需要允许该方法。

预检请求 (preflight request)

跨域资源共享 (CORS) 标准新增了一组 HTTP 首部字段, 允许服务器声明哪些源站有权限访问哪些资源。另外, 规范要求, 对那些可能对服务器数据产生副作用的 HTTP 请求方法 (特别是 GET 以外的 HTTP 请求, 或者搭配某些 MIME 类型的 POST 请求), 浏览器必须首先使用 OPTIONS 方法发起一个预检请求 (preflight request), 从而获知服务端是否允许该跨域请求。服务器确认允许之后, 才发起实际的 HTTP 请求。在预检请求的返回中, 服务器端也可以通知客户端, 是否需要携带身份凭证 (包括 Cookies 和 HTTP 认证相关数据)。

Content-Type 不属于以下 MIME 类型的, 都属于预检请求:

```
application/x-www-form-urlencoded
multipart/form-data
text/plain
```

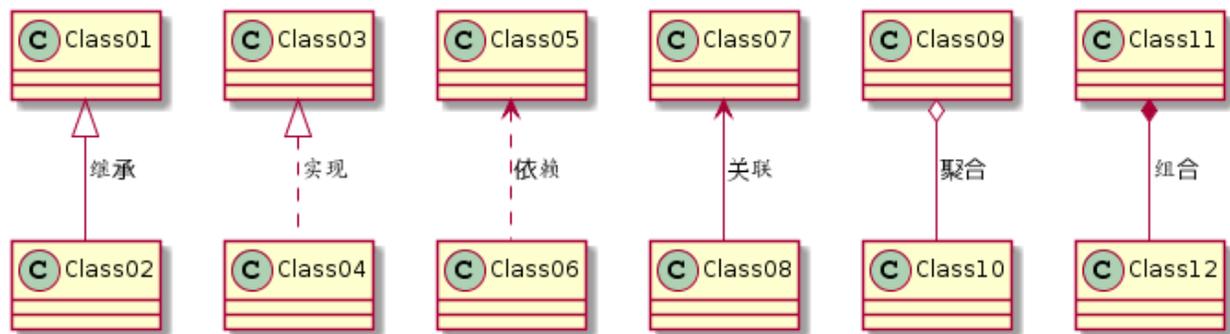
所以 application/json 的请求会在正式通信之前, 增加一次”预检”请求, 这次”预检”请求会带上头部信息 Access-Control-Request-Headers: Content-Type

8.1 类图关系

- **泛化 (generalization)**: 表示 is-a 的关系, 是对象之间耦合度最大的一种关系, 子类继承父类的所有细节。直接使用语言中的继承表达。**带三角箭头的实线表示, 箭头从子类指向父类**
- **实现 (Realization)**: 在类图中就是接口和实现的关系。三角箭头的虚线表示, **箭头从实现类指向接口**。
- **依赖 (Dependency)**: 是临时性的关联。代码中一般指由局部变量、函数参数、返回值建立的对于其他对象的调用关系。一个类调用被依赖类中的某些方法而得以完成这个类的一些职责。**带箭头的虚线表示, 箭头从使用类指向被依赖的类**
- **** 关联 (Association) ****: 对象之间一种引用关系, 比如客户类与订单类之间的关系。这种关系通常使用类的属性表达。关联又分为一般关联、聚合关联与组合关联。后两种在后面分析。**使用带箭头的实线表示, 箭头从使用类指向被关联的类。可以是单向和双向**。
- **聚合 (Aggregation)**: 表示 has-a 的关系, 是一种不稳定的包含关系。较强于一般关联, 有整体与局部的关系, 并且没有了整体, 局部也可单独存在。如公司和员工的关系, 公司包含员工, 但如果公司倒闭, 员工依然可以换公司。**在类图使用空心的菱形表示, 菱形从局部指向整体**
- **组合 (Composition)**: 表示 contains-a 的关系, 是一种强烈的包含关系。组合类负责被组合类的生命周期。是一种更强的聚合关系。部分不能脱离整体存在。如公司和部门的关系, 没有了公司, 部门也不能存在了; **在类图使用实心的菱形表示, 菱形从局部指向整体**。

8.2 PUML 图形

```
@startuml
Class01 <|-- Class02 : 继承
Class03 <|.. Class04 : 实现
Class05 <.. Class06 : 依赖
Class07 <-- Class08 : 关联
Class09 o-- Class10 : 聚合
Class11 *-- Class12 : 组合
@enduml
```



8.2.1 PUML 图像网站

8.3 单例模式

8.3.1 涉及的 JVM 与 CPU

Java 虚拟机在执行该对象创建语句时具体做了哪些事情，我们简单概括为 3 步：

- 1 在栈内存中创建 singleton 变量，在堆内存中开辟一个空间用来存放 Singleton 实例对象，该空间会得到一个随机地址，假设为 0x0001。
- 2 对 Singleton 实例对象初始化。
- 3 将 singleton 变量指向该对象，也就是将该对象的地址 0x0001 赋值给 singleton 变量，此时 singleton 就不为 null 了。

CPU 在执行指令时并不是无节操随意打乱顺序，而是有一定的原则可寻的，这个原则也叫先行发生原则 (happens-before)

先行发生原则 (happens-before)：

- 1 程序次序规则：一个线程内，按照代码顺序，书写在前面的操作先行发生于书写在后面的操作
- 2 锁定规则：一个 unlock 操作先行发生于后面对同一个锁额 lock 操作
- 3 volatile 变量规则：对一个变量的写操作先行发生于后面对这个变量的读操作
- 4 传递规则：如果操作 A 先行发生于操作 B，而操作 B 又先行发生于操作 C，则可以得出操作 A 先行发生于操作 C
- 5 线程启动规则：Thread 对象的 start() 方法先行发生于此线程的每一个一个动作
- 6 线程中断规则：对线程 interrupt() 方法的调用先行发生于被中断线程的代码检测到中断事件的发生

- 7 线程终结规则：线程中所有的操作都先行发生于线程的终止检测，我们可以通过 `Thread.join()` 方法结束、`Thread.isAlive()` 的返回值手段检测到线程已经终止执行
- 8 对象终结规则：一个对象的初始化完成先行发生于他的 `finalize()` 方法的开始

8.3.2 单例模式介绍

所谓单例，就是整个程序有且仅有一个实例。该类负责创建自己的对象，同时确保只有一个对象被创建。在 Java，一般常用在工具类的实现或创建对象需要消耗资源。

特点

- 类构造器私有
- 持有自己类型的属性
- 对外提供获取实例的静态方法

饿汉模式

线程安全，比较常用，但容易产生垃圾，因为一开始就初始化

```
public class Singleton {
    private static Singleton instance = new Singleton();
    private Singleton (){}
    public static Singleton getInstance() {
        return instance;
    }
}
```

直接调用 `Singleton.getInstance()` 来创建对象

只要这个类一加载完，就会创建实例对象，显得很着急，像一个恶汉一样，所以称之为恶汉式单例模式。

由于饿汉模式会带来一个问题，类加载的时候静态变量便申请了空间，如果没有用到的话，而且复杂的话，会造成大量的资源浪费，所以用一种懒加载的方式来创建对象，等需要用到这个对象的时候再创建。

懒汉模式

线程不安全，延迟初始化，严格意义上不是不是单例模式

```
public class Singleton {
    private static Singleton instance; // 静态类型
    private Singleton (){}

    public static Singleton getInstance() {
        if (instance == null) {
            instance = new Singleton();
        }
        return instance;
    }
}
```

直接调用 `Singleton.getInstance()` 来创建对象

由于只有当用到时才创建对象，比较懒，我们称之为懒汉式单例模式。

双重锁模式

线程安全，延迟初始化。这种方式采用双锁机制，安全且在多线程情况下能保持高性能。

```
public class Singleton {
    private volatile static Singleton singleton;

    private Singleton() {
        System.out.print("实例化对象\n");
    }

    public static Singleton getInsatnce() {
        if (singleton == null) { //-----1
            synchronized (Singleton.class) { //-----2
                if (singleton == null) {
                    singleton = new Singleton();
                }
            }
        }
        return singleton;
    }
}
```

这里 `synchronized` 的用法和上面的有所不同，上面我们用 `synchronized` 来修饰方法，表示给整个方法上了把锁，我们称之为同步方法。这里我们只给语句 1 和语句 2 加了把锁，我们称这种结构为同步代码块，同步代码块 `synchronized` 后面的括号中需要一个对象，可以任意，这里我们用了 `Singleton` 的类对象 `Singleton.class`。可以看到我们在方法中进行了两次对象是否为空的判断，一次在同步代码块外面，一次在里面。因此称之为双重检验锁模式（Double Checked Locking Pattern）。

双重检查模式，进行了两次的判断，第一次是为了避免不要的实例，第二次是为了进行同步，避免多线程问题。由于 `singleton=new Singleton()` 对象的创建在 JVM 中可能会进行重排序，在多线程访问下存在风险，使用 `volatile` 修饰 `singleton` 实例变量有效，解决该问题。

静态内部类单例模式

```
public class Singleton {
    private Singleton(){}
    public static Singleton getInstance(){
        return Inner.instance;
    }
    private static class Inner {
        private static final Singleton instance = new Singleton();
    }
}
```

只有第一次调用 `getInstance` 方法时，虚拟机才加载 `Inner` 并初始化 `instance`，只有一个线程可以获得对象的初始化锁，其他线程无法进行初始化，保证对象的唯一性。目前此方式是所有单例模式中最推荐的模式，但具体还是根据项目选择。

我们在 `Singleton` 类内部定义一个 `Inner` 静态内部类，里面有一个对象实例，当然由于 Java 虚拟机自身的特性，只有调用该静态内部类时才会创建该对象，从而实现了单例的延迟加载，同样由于虚拟机的特性，该模式是线程安全的，并且后续读取该单例对象时不会进行同步加锁的操作，提高了性能。

枚举单例模式

- 枚举类隐藏了私有的构造器。
- 枚举类的域是相应类型的一个实例对象

```
public enum Singleton {
    INSTANCE

    // doSomething 该实例支持的行为

    // 可以省略此方法，通过Singleton.INSTANCE进行操作
    public static Singleton get Instance() {
        return Singleton.INSTANCE;
    }
}
```

参考如下：

Java 大白话讲解设计模式之 -- 单例模式