

区块链共识机制介绍

科普版本

刘鹏

2019.02.21

PoET

背景

- Intel在Sawtooth Lake项目中所采用的算法。
- 本质是利用intel的硬件安全可信环境

基本原理

- **Leader选举**：每个节点都从可信的SGX环境中获得一个随机等待时间，这个等待时间与节点贡献到网络中的资源反比，即贡献越大，随机等待时间越短。
- 随机等待时间最短的节点，会被选举成为**Leader**的节点，打包出块。
- 选举和验证均是基于intel的CPU中提供的指令，intel保证其CPU是完全可信的，所以整个出块和验证的过程均为可信的。（潜台词是依赖于第三方）
- SGX的发展也可能出现漏洞，因此该共识依赖于专用硬件的安全性。

Paxos

解决问题

- 使得参与分布式处理的每个参与者逐步达成一致意见。
- 无拜占庭节点的前提假设，大多用于联盟链，无法用于公链的共识。

主要思想

- 法定集合的性质，即两个法定集合一定存在交集。对，就这么简单。

通俗理解

- 第一阶段：因为存在多个“提议者”，如果都提意见，那么“接受者”接受谁的不接受谁的？太乱了。所以，要先明确哪个“提议者”是意见领袖有权提出提议，未来，“接受者”们就主要处理这个“提议者”的提议了。（目的是为了尽早形成多数派）
- 第二阶段：由上阶段选出的意见领袖提出提议，接受者反馈意见，如果多数接受者认可一个提议，那么提议就通过了。

Paxos

算法流程

A: Proposer选择一个提议编号n, 向所有Acceptor广播Prepare(n)请求

B: Acceptor接收到Prepare(n)的请求

若之前并未收到过任何提议, 则告诉其收到的编号最大的Proposer具有写入资格。

若已经写入过value, 且提议编号n比之前接受到的请求n都大, 则承诺不会接受比n小的提议, 且回复之前接受的提案中编号<n的最大的提议, 否则不予理会。

PAXOS

PREPARE (占坑)

COMMIT (投票)

如果未得到超过半数的Acceptor响应, 提议失败。

A: Proposer得到Acceptor的响应后

如果超过半数响应

如果所有Acceptor都未接受过value, 那么向所有的Acceptor发送自己的value和n

如果有部分的Acceptor接受过value, 那么从返回值中选择n最大的作为自己的value, 提议编号n不变。

B: Acceptor收到提议后, 如果该提议的n不等于自身之前保存的n, 不接受该请求, 如果相等, 则写入本地。

Raft

解决问题

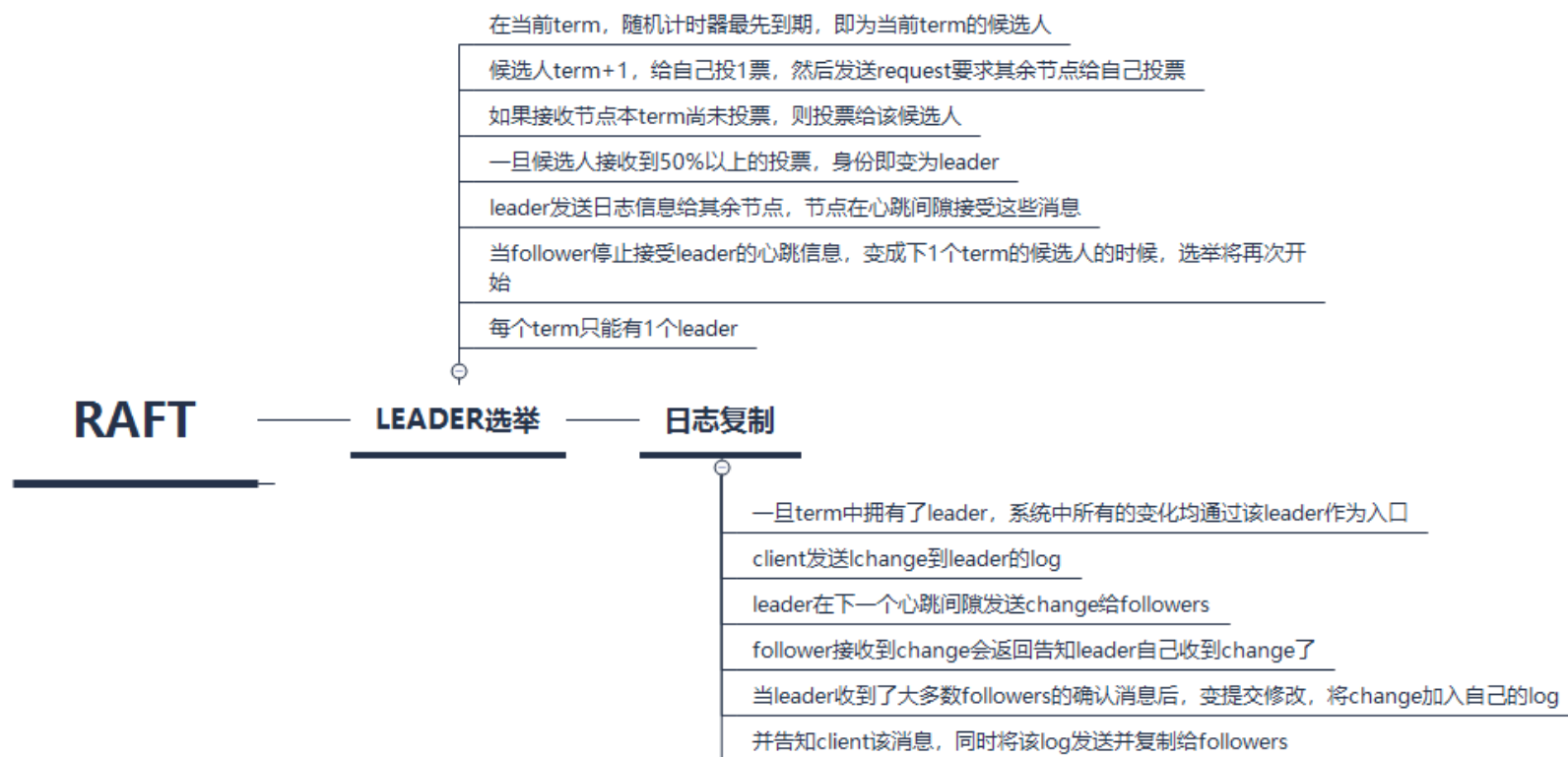
- 相比于Paxos，构建一个容易理解和工程实现的分布式协议。
- 无拜占庭问题的假设前提。

基本思想

- 在随机时间leader选举基础上叠加了各种补丁的方案。
- Log复制机制状态机
 - 首先，选出leader，leader主服务器负责接收客户端的命令请求，并将该指令以log日志的形式分发复制到其他服务器保存起来，在主服务器认为数据可靠有保障的时候（超过50%的服务器完成复制），再通知其他服务器（包括自己）将这条指令在各自的状态机中进行执行。
 - 为了区分日志的时序关系，Raft对日志的时序赋予了两个概念：Terms和序号，每当发送leader选举时，terms都会递增，而序号也是单向递增的。这两个概念是为了作为一个逻辑时钟存在。
 - <http://thesecretlivesofdata.com/raft/?spm=a2c4e.11153940.blogcont610948.8.30a799bfUMsBTx>

Raft

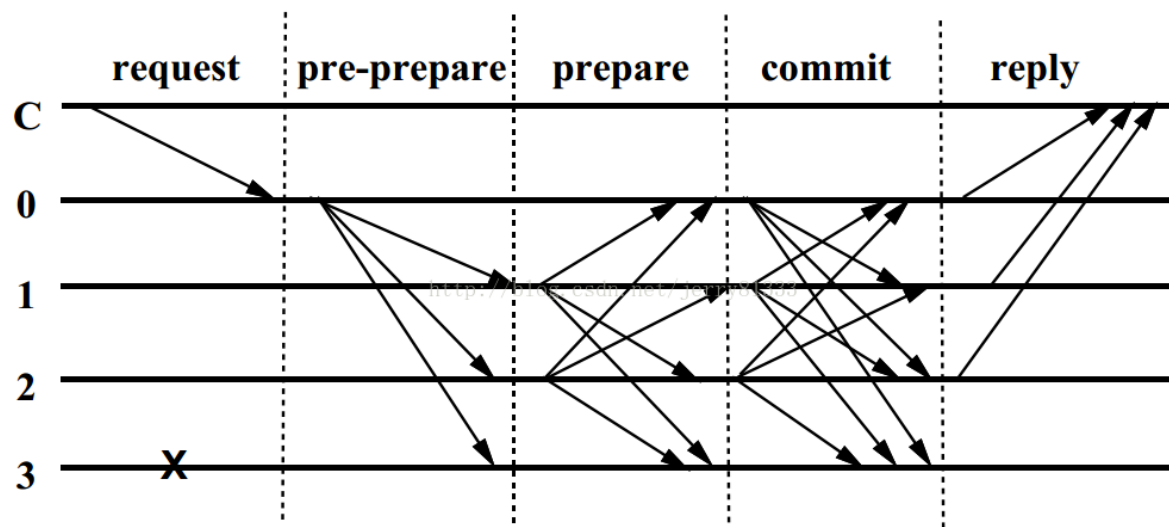
算法流程



PBFT

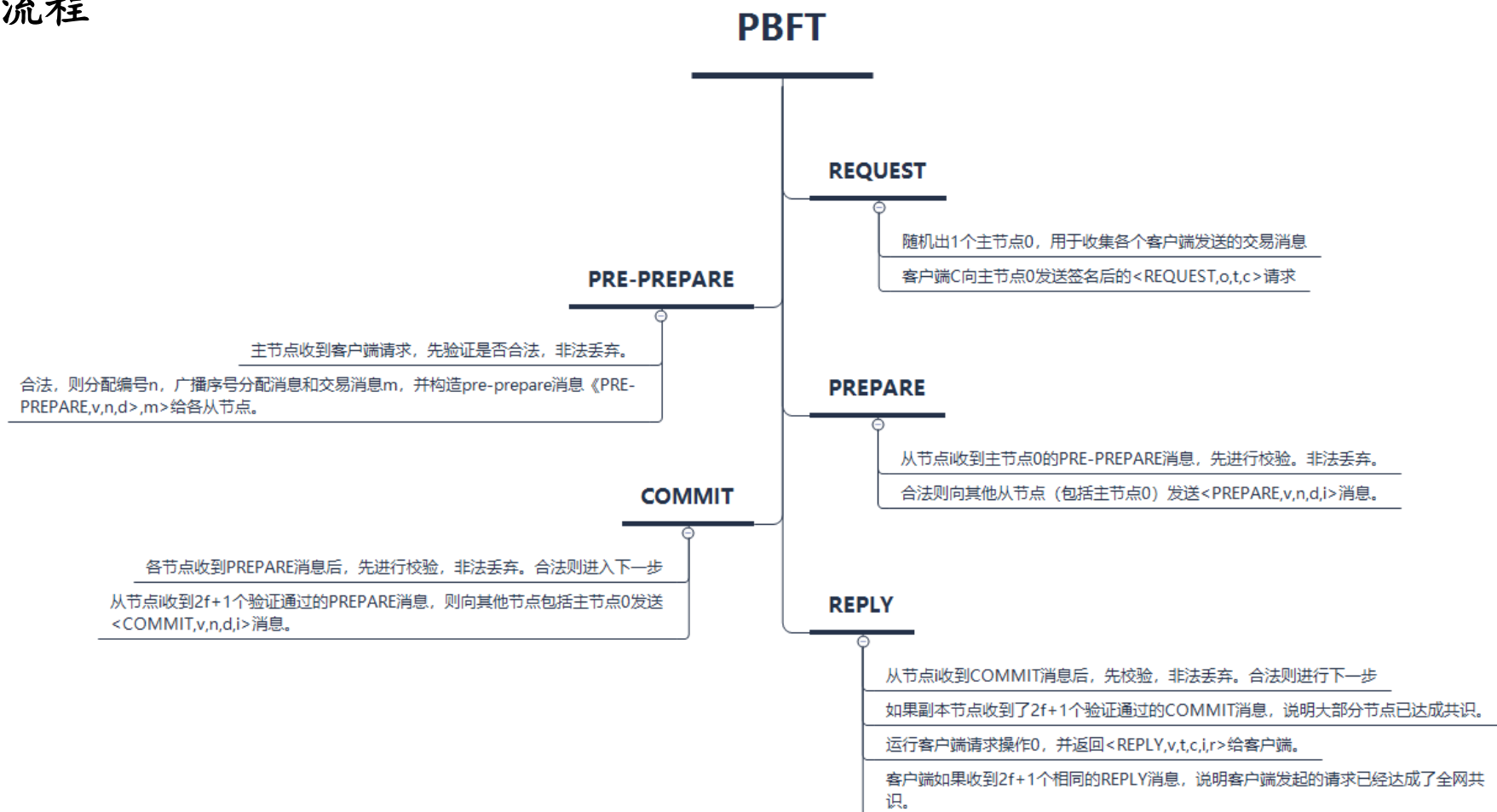
关键点

- PBFT (Practical Byzantine Fault Tolerance), 实用拜占庭容错算法。
- 解决了原始拜占庭容错算法效率不高的问题, 时间复杂度是 $O(n^2)$ 。
- 每个节点都需要与其他节点进行P2P共识同步, 节点增加, 性能下降很快。主要用于联盟链。



PBFT

算法流程



CONFLUX

要解决的问题

- 比特币的性能和效率问题；
- 现有对比特币的改进方式简单分为 增加区块容量 和 提高出块率，但会带来更多分叉，引发安全性和资源浪费的问题。

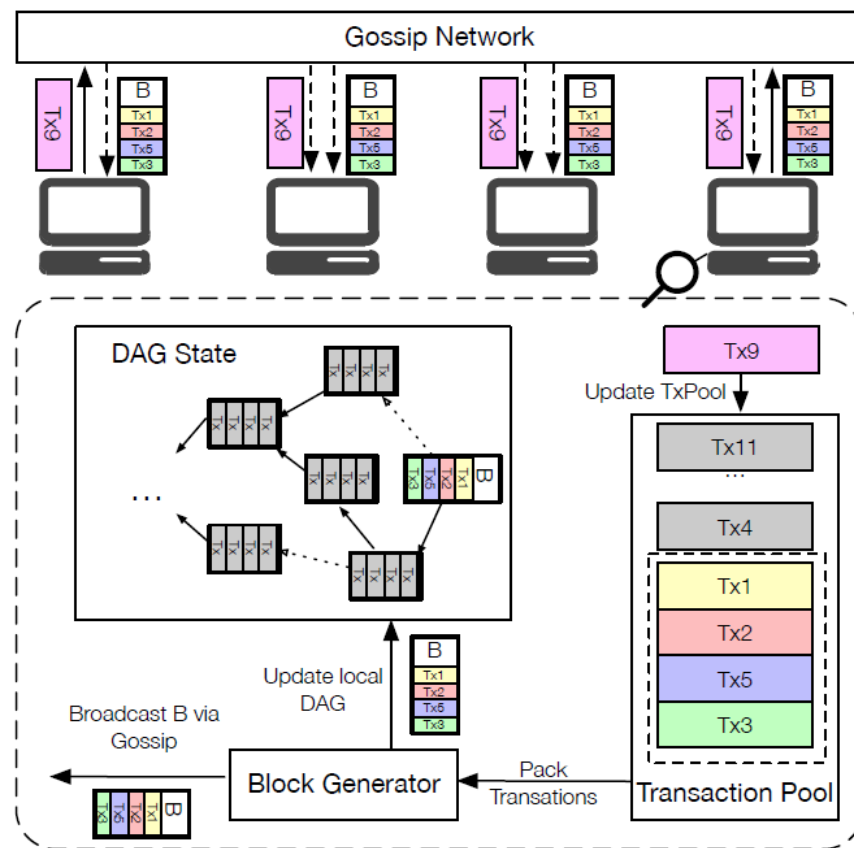
核心思想

- 采用DAG组织区块，并发的处理交易和区块，延迟确定事务的全局顺序。
- 主链不一定是最长链，防止双花。
- 通过DAG的组织 and 交易的全局排序，将分叉的无效区块变成有效区块，提升了有效区块的比例，从而提升系统的吞吐率，减少浪费。

CONFLUX

总体架构

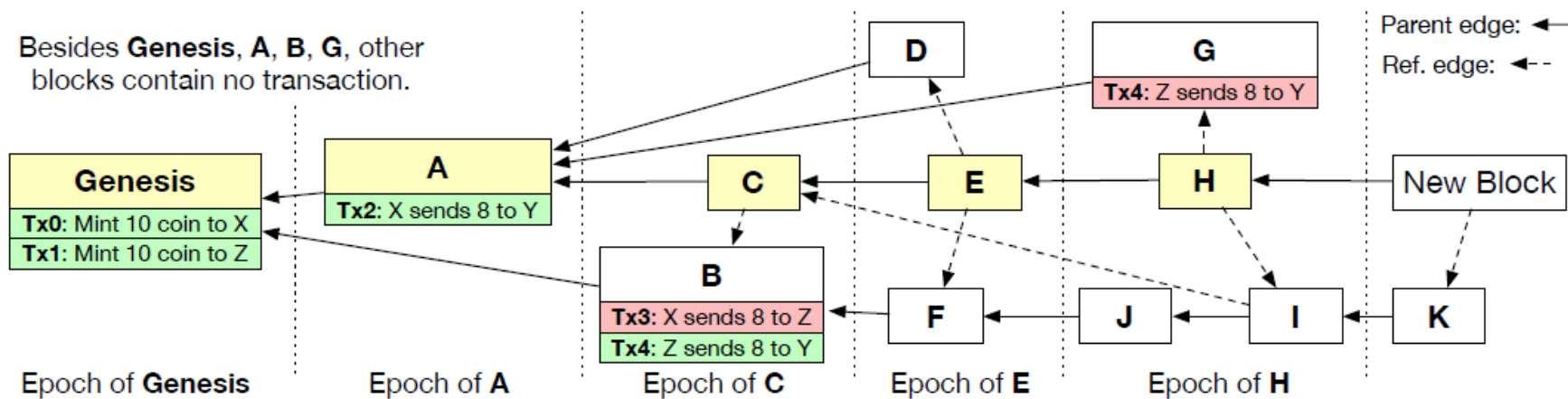
- 和比特币一样，采用Gossip实现P2P网络交互，实现交易或者区块的广播。
- 虚线部分表达了1个全节点的工作流程：node接收到1个Tx后，将Tx放入TxPool，Block generator模块将交易池中的交易进行打包。
- 打包后，将新的Block更新至本地DAG状态，将新的Block广播至全网，并根据DAG状态规定的区块顺序对该Block中的交易进行删除。
- 对于其他的全节点，如果接收到其他节点广播的block，会直接更新至本地DAG状态，并将该Block中包含的Tx从本地TxPool中删除，以避免重复交易打包。



CONFLUX

区块DAG

- Conflux的区块之间靠 父链接 和 引用链接 连接在一起，上图中实线表示 父链接，虚线表示 引用链接。链接在一起的区块关系，会形成DAG。
- Conflux中组成DAG的区块会先确定一条主链（Pivot Chian），在这个基础上再确定所有区块的先后顺序。



CONFLUX

确定主链

- $G = \langle B, g, P, E \rangle$, 其中, B 是DAG中的所有区块, g 是Genesis block, P 是映射函数 (每个区块可以通过该函数获取父区块), E 是所有区块的“父链接”和“引用链接”的集合。
- 主链上下一个区块的选取规则为 子区块个数多的、或者在子区块个数相等时区块Hash小的区块。(子区块的计算只包括父链接, 不包括引用链接)

划分Epoch

- 主链Pivot Chain中的每一个block对应一个epoch, 其他非Pivot Chain上的block也会被分配到相应的chain中。
- 分配的原则为, 被主链上的区块“链接“, 且没有被之前的区块”链接“的区块, 即属于该区块Epoch。

Input : The local state $G = \langle B, g, P, E \rangle$ and a starting block $b \in B$

Output: The last block in the pivot chain for the subtree of b in G

```
1 if Child( $G, b$ ) =  $\emptyset$  then
2   | return  $b$ 
3 else
4   |  $s \leftarrow \perp$ 
5   |  $w \leftarrow -1$ 
6   | for  $b' \in \text{Child}(G, b)$  do
7   |   |  $w' \leftarrow |\text{Subtree}(G, b')|$ 
8   |   | if  $w' > w$  or  $w' = w$  and Hash( $b'$ ) < Hash( $s$ )
9   |   |   | then
10  |   |   |   |  $w \leftarrow w'$ 
11  |   |   |   |  $s \leftarrow b'$ 
12  | return Pivot( $G, s$ )
```

CONFLUX

确定区块顺序

- 按照Epoch进行定序，在同一个epoch中，不同的blocks按照parental tree的拓扑关系来排序。

确定交易顺序

- 在区块顺序确定的前提下，前一个区块中的交易在后一个区块中的交易前面。
- 同一区块的交易，按照区块中交易顺序进行排序。
- 当交易产生冲突时，第一个交易有效，后续冲突交易都无效。

Input : The local state $G = \langle B, g, P, E \rangle$ and a block a

Output: A list of blocks $L = b_1 \circ b_2 \circ \dots \circ b_n$, where

$b_1 = g$ and $\forall 1 \leq i \leq n, b_i \in B$

```
1  $a' \leftarrow P(a)$ 
2 if  $a' = \perp$  then
3   return  $a$ 
4  $L \leftarrow \text{ConfluxOrder}(G, a')$ 
5  $B_\Delta \leftarrow \text{Past}(G, a) - \text{Past}(G, a')$ 
6 while  $B_\Delta \neq \emptyset$  do
7    $G' \leftarrow \langle B_\Delta, g, P, E \rangle$ 
8    $B'_\Delta \leftarrow \{x \mid |\text{Before}(G', x)| = 0\}$ 
9   Sort all blocks in  $B'_\Delta$  in order as  $b'_1, b'_2, \dots, b'_k$ 
10  such that  $\forall 1 \leq i < j \leq k, \text{Hash}(b'_i) < \text{Hash}(b'_j)$ 
11   $L \leftarrow L \circ b'_1 \circ b'_2 \circ \dots \circ b'_k$ 
12   $B_\Delta \leftarrow B_\Delta - B'_\Delta$ 
13 return  $L$ 
```

ALGORAND

背景

- Silvio Micali在2016年提出，75页的长文发表在Arxiv。
- 为了解决POW共识的算力浪费、扩展性差、TPS低、易分叉等问题。

优势

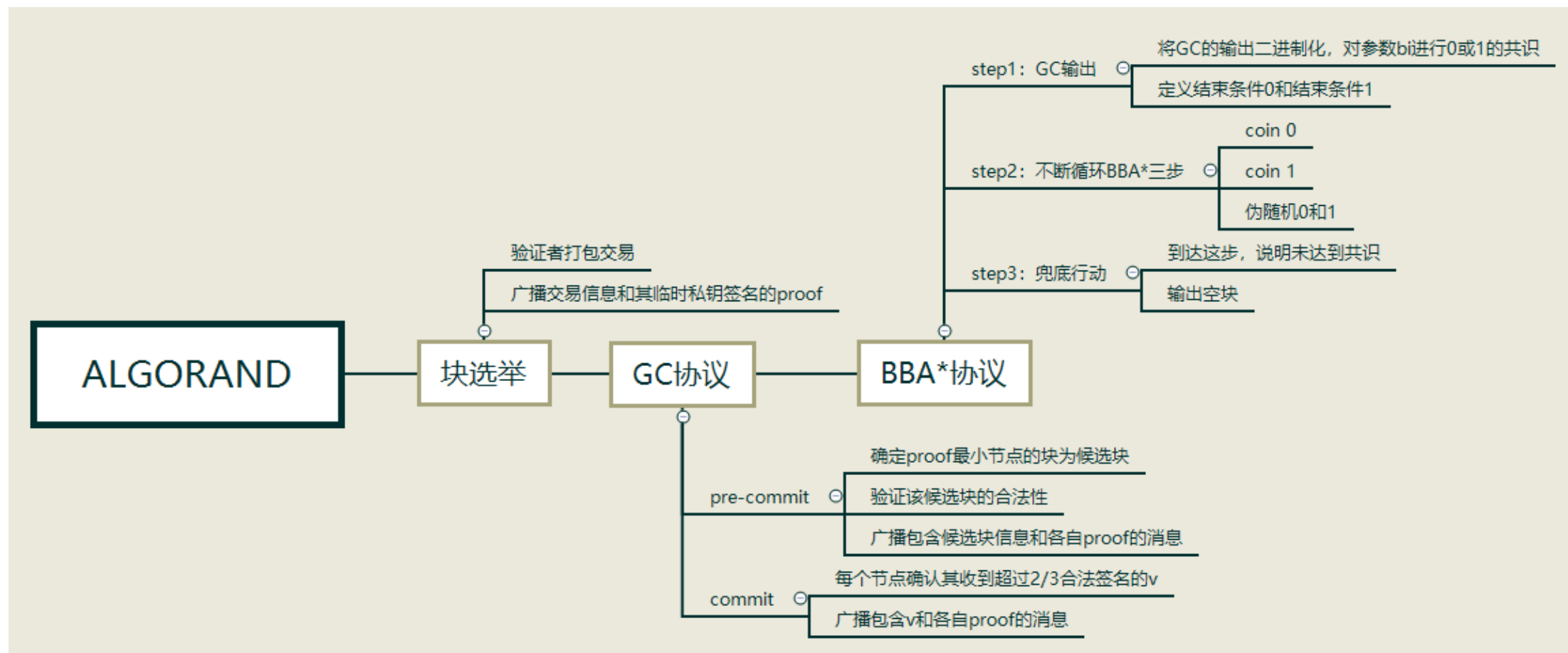
- 能耗低：1/1500的节点执行几秒钟的运算
- 民主化：不会出现矿池
- 分叉概率极低： 10^{-18}
- 可扩展性好

缺点

- 强同步网络的假设
- 步骤太多，且每一步都要全网广播，导致时间过长

ALGORAND

算法流程



ALGORAND

关键技术

1、BBA*协议

2、种子参数

选取完全无法预测的种子参数，从而保证不被敌手影响，上一轮的种子会参与下一轮投票用户的生成。

3、秘密抽签和秘密资格

所有参与共识投票的用户都是本地得知自己的身份，投票后身份才被暴露，虽然敌手可以马上腐蚀他们，但是消息已经无法撤回。另外，消息生成后，用于签名的临时秘钥会被丢弃，敌手即使在该轮控制了他们，也无法生成合法消息。

4、用户可替换

BA*的每一步，SV都会轮换，且每一步的SV是相互独立的，因此可扩展性很强。

ALGORAND

关键技术

5、密码学抽签（VRF函数）

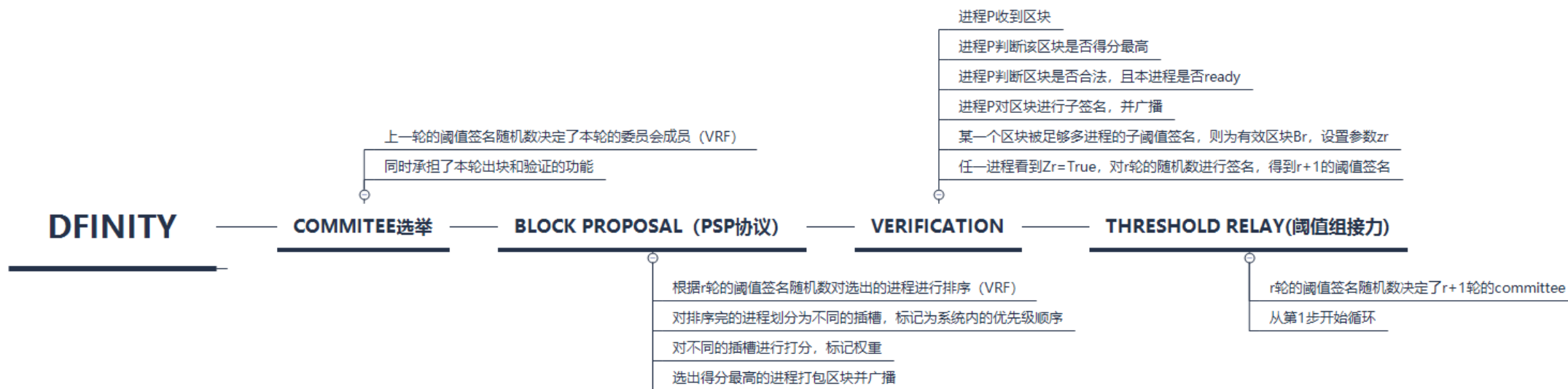
每一个节点持有私钥SK，并公开其验证公钥VK。如果一个用户想知道他是否被选中在参与出块r的委员会中，他讲行动如下：

- 调用 $\text{Evaluate}(\text{SK}, Q_{r-1}) \rightarrow (H(\sigma), \sigma)$, Q_{r-1} 是系统中对所有人公开的随机种子, $\sigma = \text{SIG}_i(r, s, Q_{r-1})$;
- 检查 $H(\sigma)$ 是否在值域 $[0, P]$ 范围内, 该值域取决于用户拥有的权益比重。

如果检查通过, 那么用户就被选中为SV成员, 且可被他人通过一组值 $(Q_{r-1}, \text{VK}, H(\sigma), \sigma)$ 来进行验证。也就是说, 用户被选择为SV成员的过程, 是随机不可预测的, 但又是可被他人验证的。

DEFINITY

算法流程

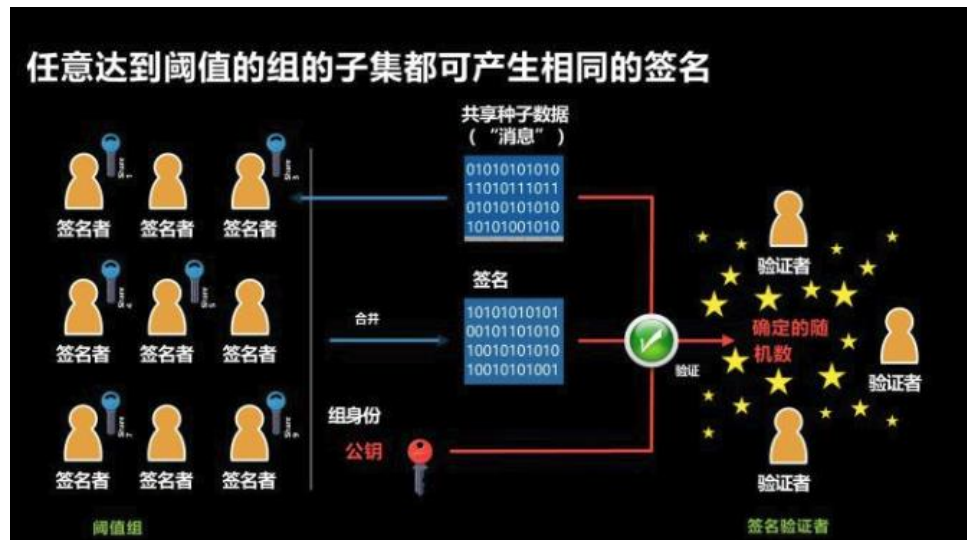


DFINITY

关键技术

1、阈值签名（基于BLS算法）

- 一个由公钥识别的阈值组仅可在给定的种子数据上产生单个唯一有效的签名输出。
- 这个组具有容错能力，任何 \geq 阈值的子集将组签名分片合成组签名，即可完成群签。
- 所得到的阈值签名可以由任何具有该组公钥和种子数据的人来验证。
- 这个签名是确定性产生的随机数。
- 给定阈值组的公钥和种子数据，验证者就可以在无需执行共识协议的情况下立即就随机数达成一致。

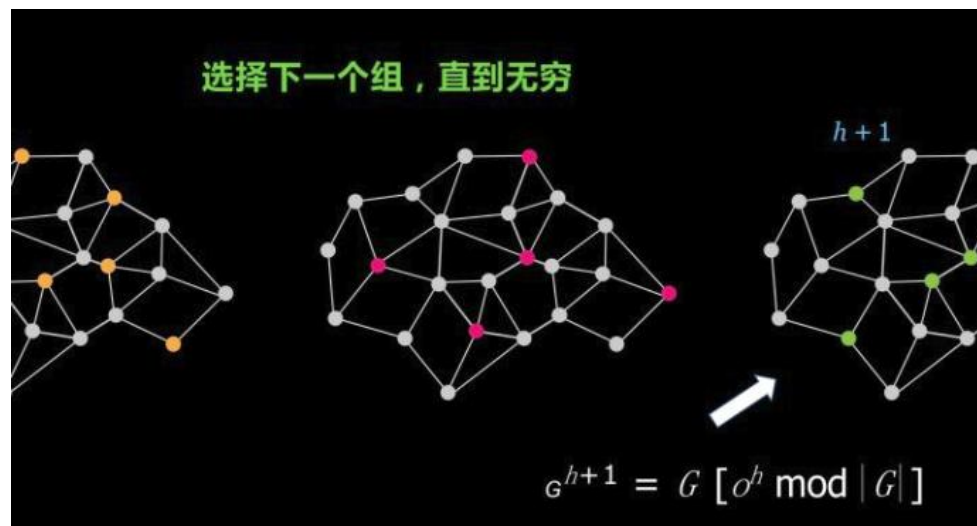


DFINITY

关键技术

2、阈值接力

- 高度为 h 的阈值组对高度 $h-1$ 的阈值签名进行签名，生成随机数。
- 该随机数选择 $h+1$ 高度的阈值组，这个过程称为阈值接力。
- 按此可产生一个去中心化的可验证随机函数（VRF）
- 该随机数序列是 确定的、不可操纵的、又是可验证的。
- 当前组达到阈值后，发布下一个值，因此又是不可预测的。



ACP

算法流程

