高性能区块链系统原理及构建

陈鹏飞

中山大学 数据科学与计算机学院





区块链回顾



区块链系统中分布在全球各地的一个个节点;而这些节点可以简单理解为一台服务器服务器集群;在一个典型的P2P系统加入货币激励。



区块链回顾





区块链系统的基本角色: 物、人、事





投资





\$32m







\$50m





\$30m

- First funding rounds for three startups raised a total of \$102m
- Goldman Sachs acts as a VC in the Bitcoin's Block chain space
- Santander InnoVentures has created a \$100m investment fund,
- Visa joined a syndicate to invest \$30m in Chain

内部开发









- NASDAQ in-house platform described on p.40
- Citibank has created citicoin, a cryptocurrency being tested internally
- IBM is working on a blockchain platform for the Internet of Things

合作伙伴















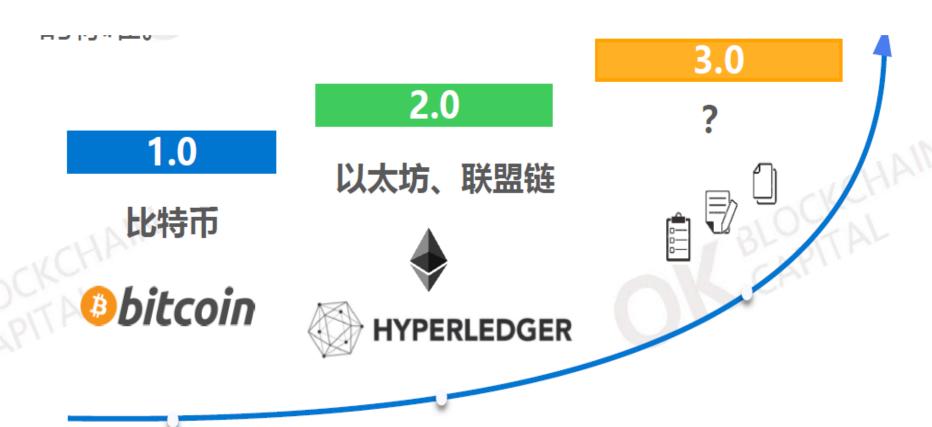
safello

- Ripple has attracted many partners for its real-time payments system
- UBS and Clearmatics partnered to create a securities settlement system
- Barclays partnered with Safello to deliver "proof-of-concepts"



区块链系统发展





区块链3.0: 侧链与跨链

OK BlockchainCapital



区块链展望



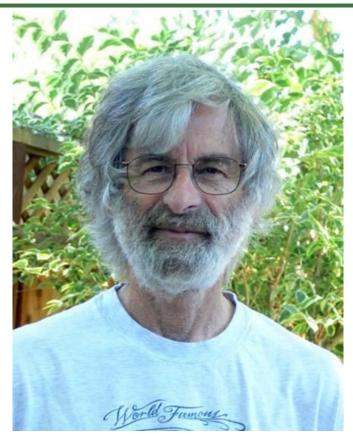




状态机复制







[PDF] Using Time Instead of Timeout for Fault-Tolerant Distributed Systems

https://www.microsoft.com/en-us/research/uploads/.../using-time-Copy.pdf ▼ 翻译此页

作者: L LAMPORT - 1984 - 被引用次数: 397 - 相关文章

A general method is described for implementing a distributed system with ... use of absolute times instead of timeouts, and can be considered an extension of.



Paxos!

状态机复制







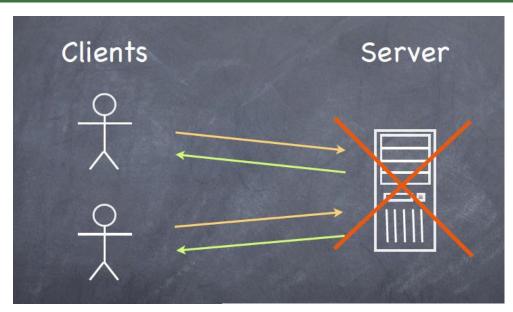
[PDF] Implementing Fault-Tolerant Services Using the State Machine ... https://www.cs.cornell.edu/fbs/publications/SMSurvey.pdf ▼ 翻译此页 作者: FB SCHNEIDER - 1990 - 被引用次数: 2384 - 相关文章

Implementing Fault-Tolerant Services. 303 state machine and running a replica on each of the processors in a distributed system. Provided each replica being ...

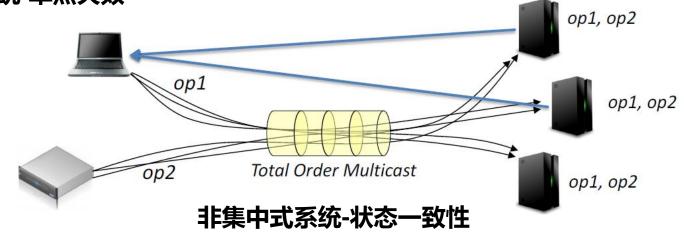


状态机复制





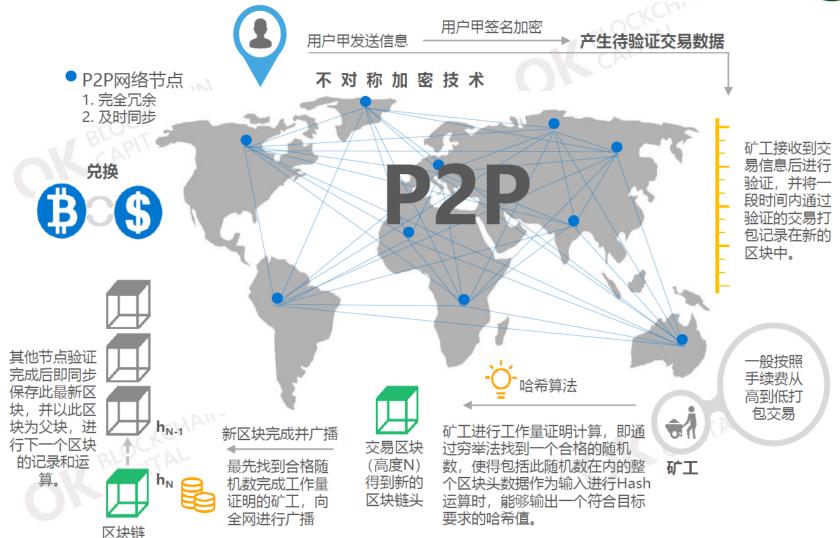
集中式系统-单点失效





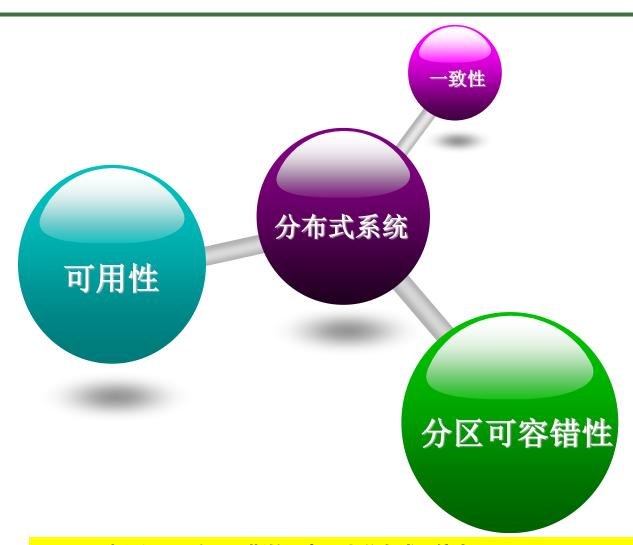
区块链交易系统





分布式系统"不可能三角"理论



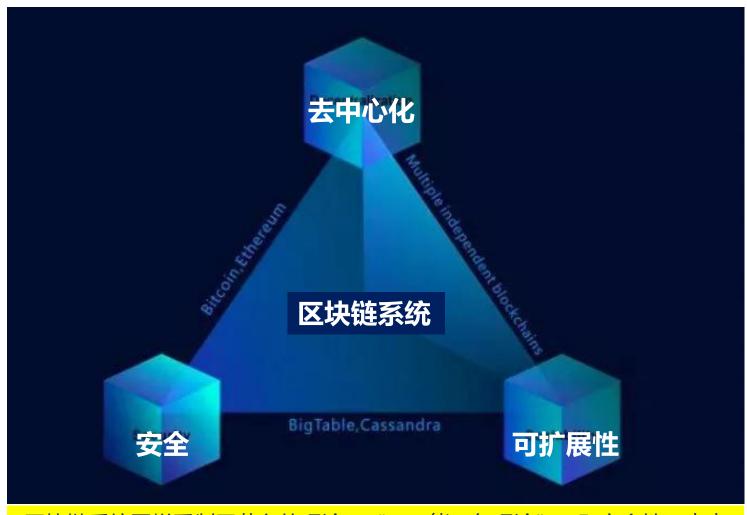


CAP原则又称CAP定理,指的是在一个分布式系统中,Consistency (一致性)、Availability (可用性)、Partition tolerance (分区容错性),三者不可得兼



区块链系统"不可能三角"理论





区块链系统同样受制于著名的理论: "不可能三角理论",即安全性、去中心化和可扩展性三者不可兼得



电子货币的交易性能



Cryptocurrencies Transaction Speeds Compared to Visa & Paypal





电子货币的交易性能



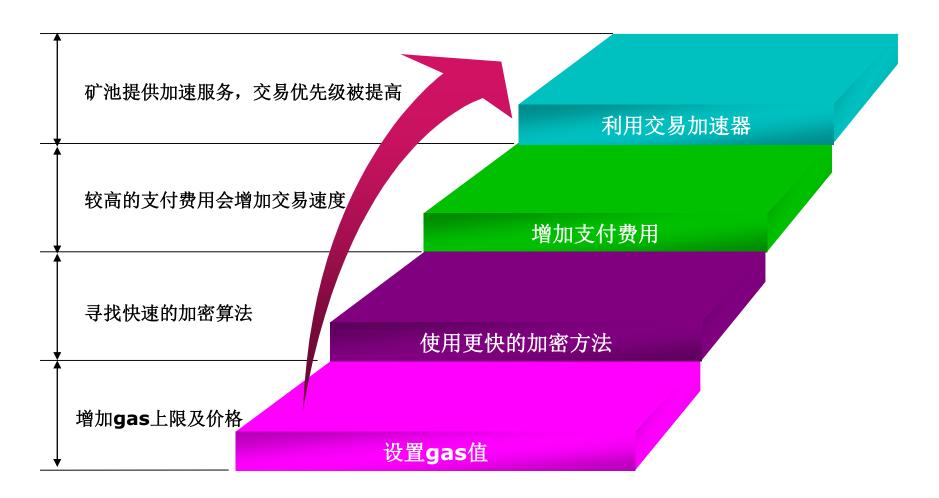
测量标准: 统计交易被确认的时间

电子货币	交易性能	电	子货币	交易性能
Bitcoin	单个交易时间平均78分钟;		Tron	单个交易时间平均15分钟;
Ripple	单个交易时间平均4秒;		Cardano	单个交易时间平均5分钟;
Litecoin	单个交易时间平均30分钟;		lota	单个交易时间平均3分钟;
以太坊	单个交易时间平均6分钟;	(2)	Zcash	单个交易时间平均15分钟;
EOS	单个交易时间平均1.5秒	M	Monero	单个交易时间平均30分钟;
Augur	单个交易时间平均6分钟;	151	Steem	单个交易时间平均3秒钟;
Dash	单个交易时间平均15分钟	A	Ardor	单个交易时间平均1分钟;



以太坊加快交易的方法

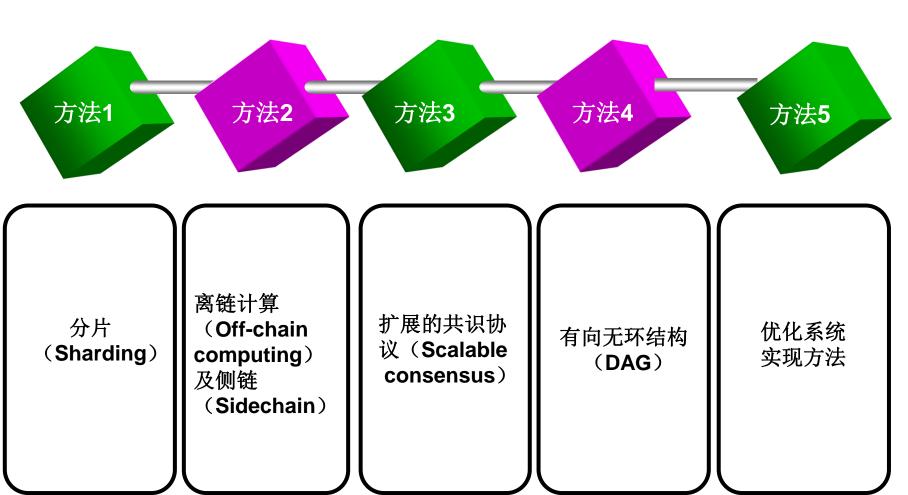






高性能高扩展性区块链系统原理







高性能高扩展性区块链系统原理





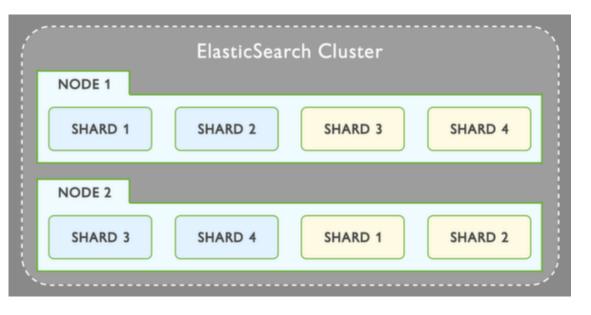
分区 (Sharding):

一种数据库技术,将大型数据分成更小、更快、更容易管理的部分



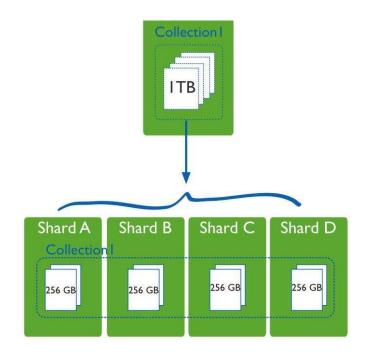
分片 (Sharding)





Elasticsearch sharding

■ 分片 (sharing) 其实是一种传统的数据库技术, 它将数据库分割成多个碎片并将这些碎片放置在 不同的服务器上,用于提高**性能和可用性**。

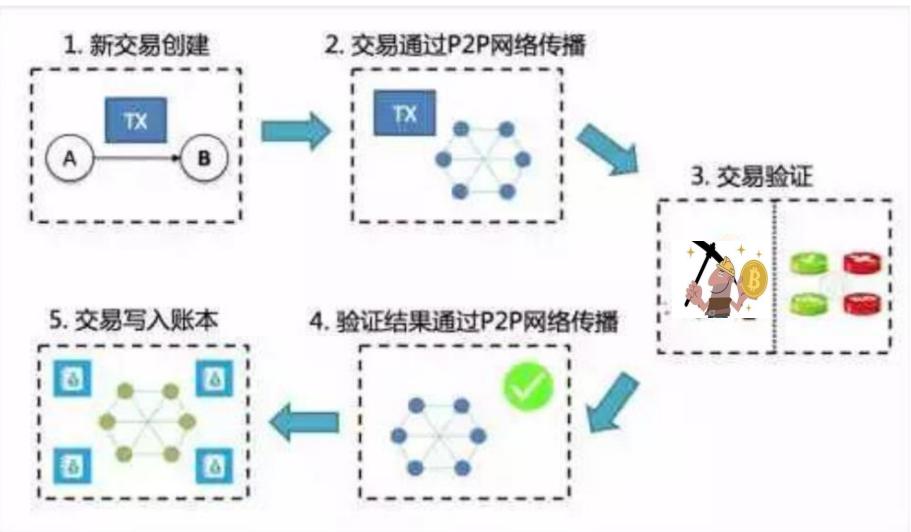


MongoDB sharding



回顾







区块分片 (Sharding)



- 传统区块链单一的区块打包队列;
- 将验证、打包的过程并行化;
- 交易量提高到干笔每秒;
- 交易费用减少;



分片策略



> 网络分片

区块链节点的网络被分割成不同的碎片,每个碎片都能形成独立的处理 过程并在不同的交易子集上达成共识。 需要保证安全性,方案: 随机性

> 交易分片

根据交易的某些属性划分碎片,比如根据交易记录的Hash值 (存在的问

题?) 或者根据交易账户;

> 状态分片

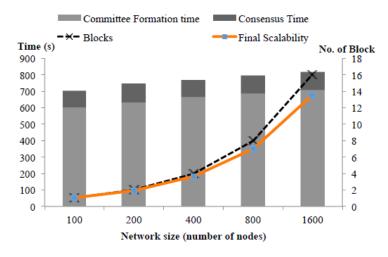
公共节点都承担着存储交易、智能合约和各种状态的负担,这可能使其在为了获得更大的存储空间而进行巨大的花费。状态分片让每个节点只负责托管自己的分片数据,而不是存储完整的区块链状态。存在的问题?



分片策略

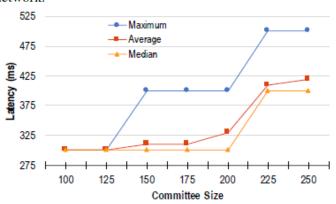
Luu L, Narayanan V, Zheng C, et al. A Secure Sharding Protocol For Open Blockchains[C] ACM Sigsac Conference on Computer and Communications Security. ACM, 2016:17-30.

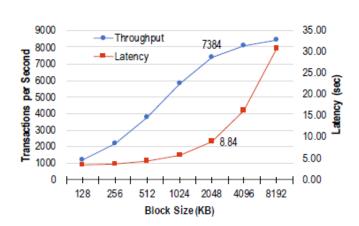
Fig



RapidChain: Scaling Blockchain via Full Sharding

Figure 1: ELASTICO scales up the throughput nearly linearly in the computation capacity of the network.



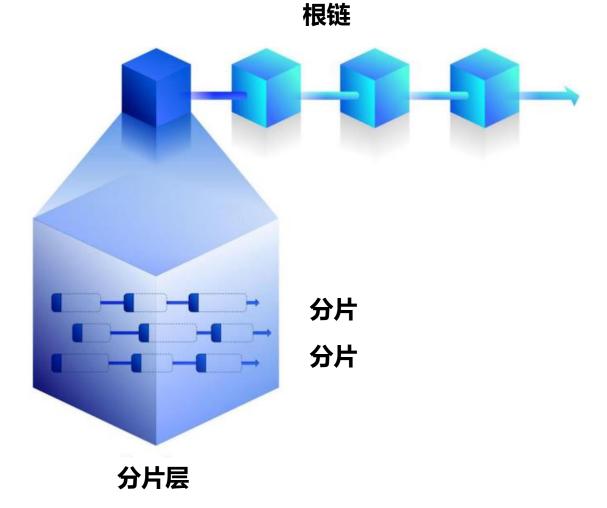




夸克链 (1)



系统架构



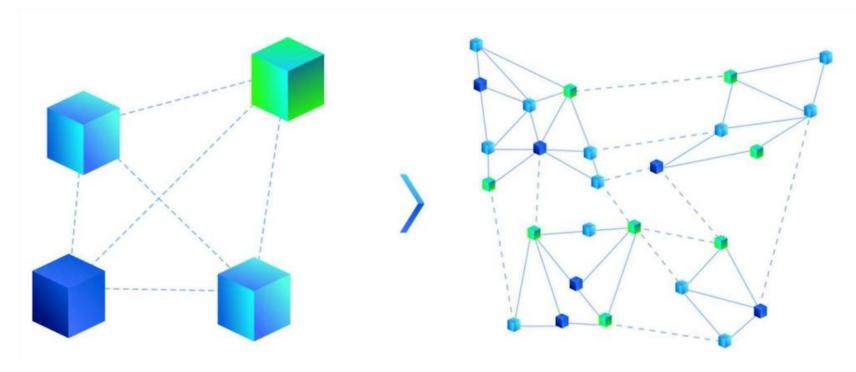
其中每条分片处理一个事务的子集合,而根链则通过在根链中包含分区的区块头来确认分区。



夸克链 (2)

SIN LINES UNITED

抗中心化横向扩展性



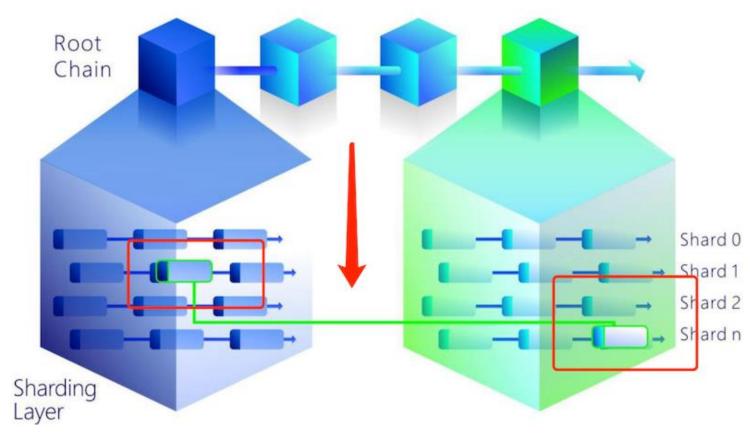
QuarkChain 网络的横向扩展性,其中四个超全节点(左)被四个节点集群(右)取代,其中每个集群中的节点彼此诚实。(实线表示诚实的连接,虚线表示不可靠的连接)



夸克链 (2)



高效、安全的分片交易



只要跨分片交易被根链确认, 交易就可以确认



高性能高扩展性区块链系统原理





离链计算/侧链(Off-chain computing):

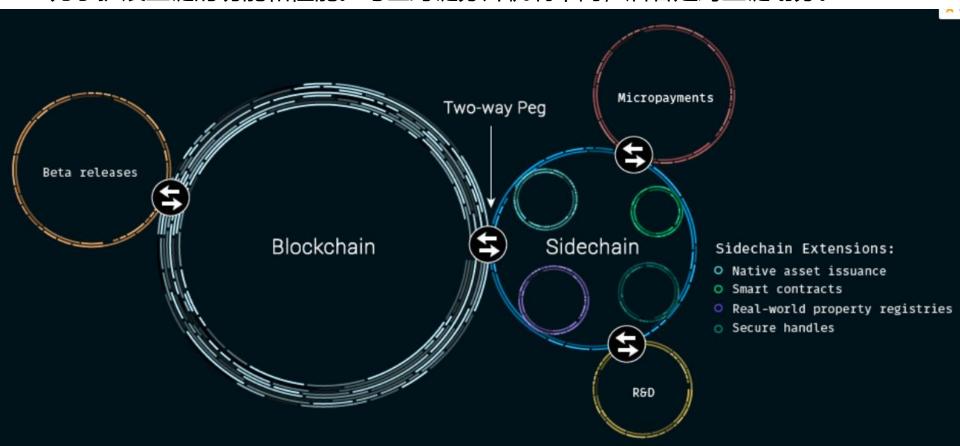
在主链之下进行交易,中间交 易不进入主链确认,待最后一 笔交易完成后回归到主链。



侧链



侧链的概念是通过双向锚定实现了主链与侧链之间的价值转移,侧链的目的是 为了扩展主链的功能和性能。与区跨链分片机制不同,后者是对主链划分。





区块链3.0: 侧链与跨链

侧链要解决的关键问题



1.跨链交易验证问题

链与链之间如何建立一个信任机制,验证跨链之间的交易数据?解决方案有:

a. 公证人机制, b. 区块头Oracle + SPV简易验证

2. 跨链事务管理问题

跨链交易包含多个子交易,这些子交易构成了一个事务。跨链的事务管理 又分为两个子问题:

2.1 如何确定子交易是否被最终确认,永不回滚?

解决方案:

- a. 等待足够多确认,b. 区块纠缠,c. DPoS/xBFT
- 2.2 如何保证交易的原子性? 所有子交易要么都成功, 要么都失败。

解决方案有: 哈希时间锁

3. 锁定资产管理问题

资产跨链转移时,锁定资产如何管理?

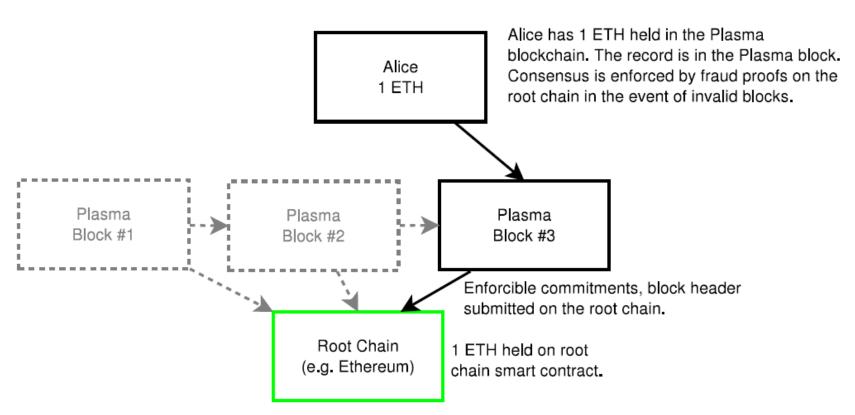
解决方案有:

a. 单一托管人b. 联盟托管人c. 智能合约托管



Plasma侧链





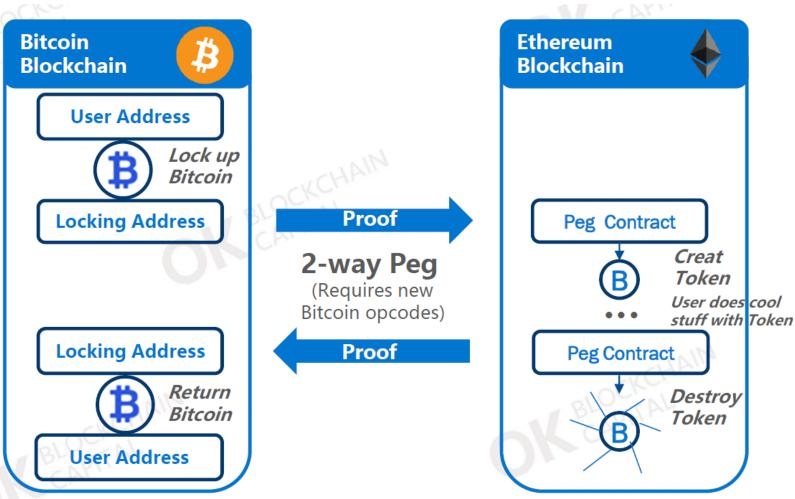
Plasma blockchains are a chain within a blockchain.



Plasma: Scalable Autonomous Smart Contracts

不同链之间的信息交换









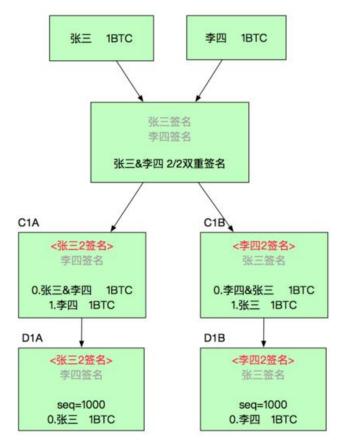
- ▶ 基于微支付通道 (双向支付通道)
- ▶ 两种类型的交易合约——序列到期可撤销合约RSMC (Revocable Sequence Maturity Contract) 和哈希时间锁定合约HTLC (Hashed Timelock Contract)

Commitment Tx

资金池 Funding Tx 共同可见

➤ RSMC解决了通道中货币单向流动问题,HTLC解决了货币跨节点传递的问题 引入sequence

引入sequence,阻止后续交易进块(D1A),给出一个实施惩罚窗口期。

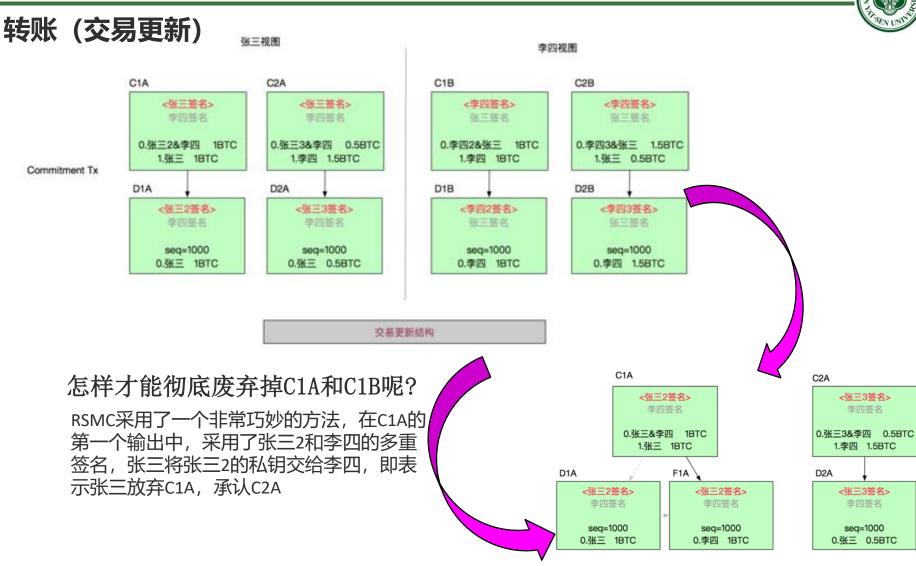


RSMC交易结构图

建立微 (双向) 支付通道





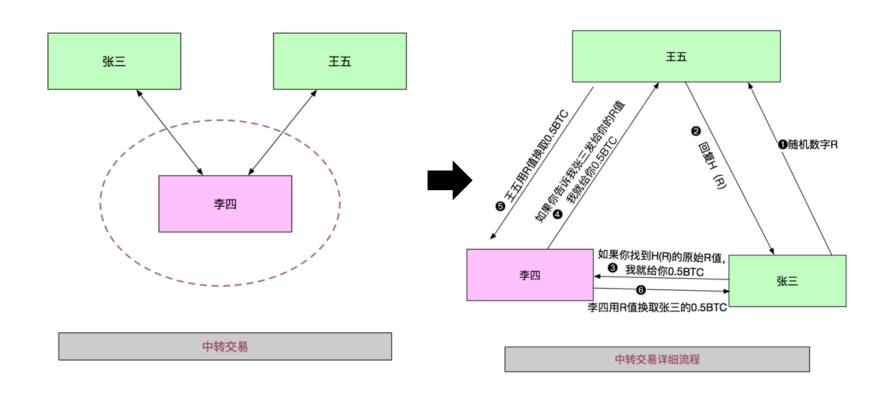




交易更新过程

SIN YATE SEN UNITE

中转交易



张三通过李四与王五进行中转交易





闪电网络本质

RSMC通过巧妙的 设置Commitment TX 的多重签名输出,以 的多重签名输出,以 及sequence的延迟 块形成惩罚窗口期, 解决了在微支付通动 中的货币单向流动向 断形LC则解决了货币 节点的传递。

闪电网络是一种链下交易,只有关键的环节才会发送到比特币主链,非关键环节则在链下(闪电网络上)进行计算、运行,这样大大降低了主链负荷。



构建侧链



利用以太坊构建侧链 (下回分解……)



高性能高扩展性区块链系统原理





扩展的共识协议(Scalable consensus):

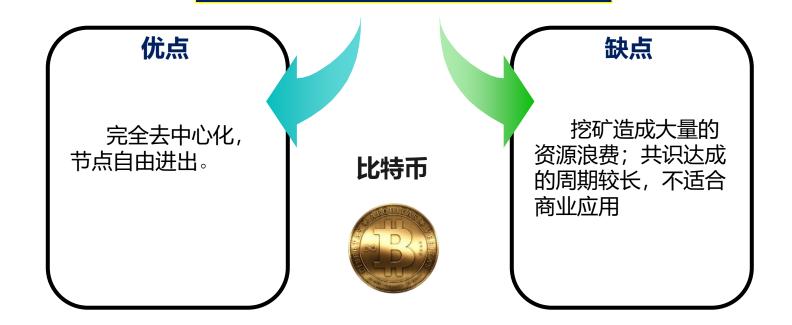
采用新的可扩展的共识协议



PoW工作量证明



计算出一个满足规则的随机数,即获 得本次记账权,发出本轮需要记录的 数据,全网其它节点验证后一起存储

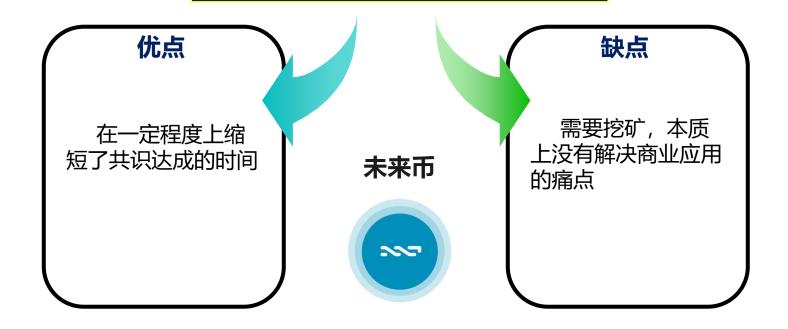




PoS 权益证明



Pow的一种升级共识机制;根据每个节点 所占代币的比例和时间;等比例的降低挖 矿难度,从而加快找随机数的速度。

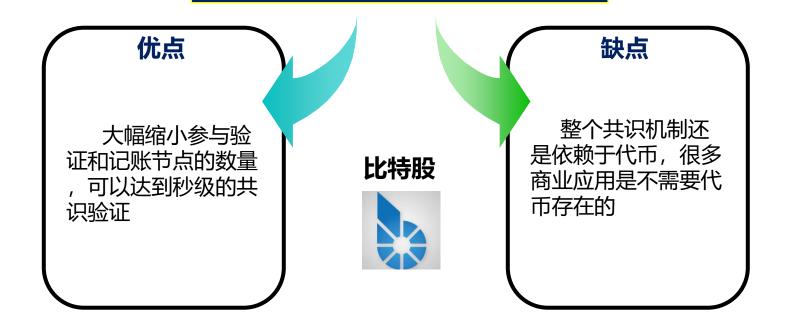




DPoS 权益证明



类似于董事会投票,持币者投出 一定数量的节点,代理他们进行验证 和记账。

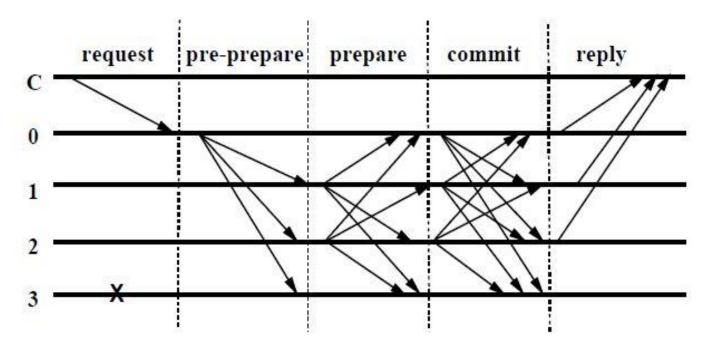




PBFT拜占庭容错算法



- ▶ 当节点发现leader作恶时,通过算法选举其他的replica为leader;
- ▶ leader通过pre-prepare 消息把它选择的 value广播给其他replica节点,其他的 replica节点如果接受则发送 prepare,如果失败则不发送;
- ▶ 一旦2f个节点接受prepare消息,则节点发送commit消息;
- ➤ 当2f+1个节点接受commit消息后,代表该value值被确定





PBFT拜占庭容错算法



消息传递的一致性算法,算法经过三个阶段达成一致性,这些阶段可能因为失败而重复进行

优点

- 系统运转可以脱离币的存在,安全性与稳定性由业务相关方保证;
- ▶ 共识的时延大约 在2~5秒钟;
- > 共识效率高;

超级账本



HYPERLEDGER

缺点

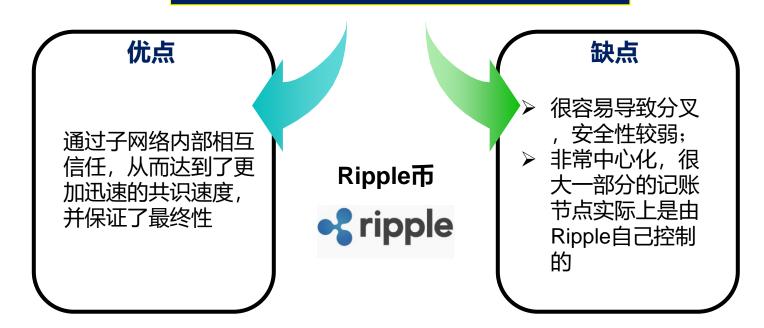
- 其封闭性(节点数目提前确定并互相联通);
- ▶ 高性能开销 (O(N^2)的消息量):
- 连接节点数量有限,无法用于公链;



Ripple共识协议



PCA对交易分两个阶段完成,第一阶段是达成交易集(即未打包进入区块的合法交易的集合)的共识,第二阶段是对新生成的区块进行提议,最终形成被共识过的区块。





共识机制比较



共识机制	优势	劣势
POW	1. 安全稳定,节点自由度高 2. 去中心化程度高,节点系统 开放	1.扩展性弱,性能低 2.没有最终性 3.造成硬件设备浪费
POS	 1. 能源耗费少 2. 去中心化程度较高,节点系统开放 	1.实现过程复杂 2. 存在安全漏洞
DPOS	1.能源耗费少 2. 性能高 3. 具备最终性	去中心化程度弱,节点系统相对封闭
BFT	 1.性能较高 2. 具备最终性 3. 安全性好 	1.去中心化程度弱, 节点系统 封闭 2. 容错率低



DAG





有向无环结构 (DAG)



DAG



1.单元:区块链组成单元是Block(区块),

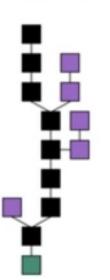
DAG组成单元是TX(交易);

2.**拓扑**:区块链是由Block区块组成的单链,只能按出块时间同步依次写入,好像单核单线程

CPU; DAG是由交易单元组成的网络,可以 异步并发写入交易,好像多核多线程CPU;

3. 粒度: 区块链每个区块单元记录多个用户的 多笔交易, DAG每个单元记录单个用户交易。

Blockchain







高性能高扩展性区块链系统原理





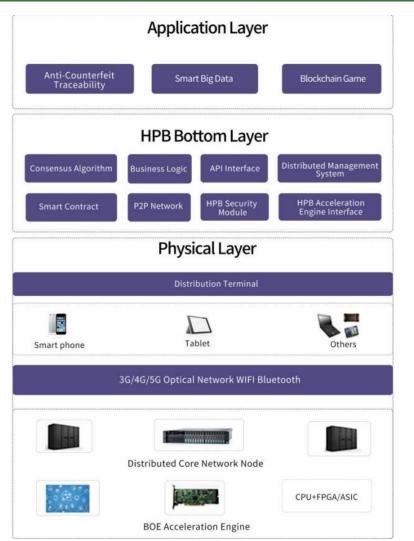
系统优化:

通过优化系统 架构、代码结构、 数据存储效率、 网络发送效率等, 提高区块链性能

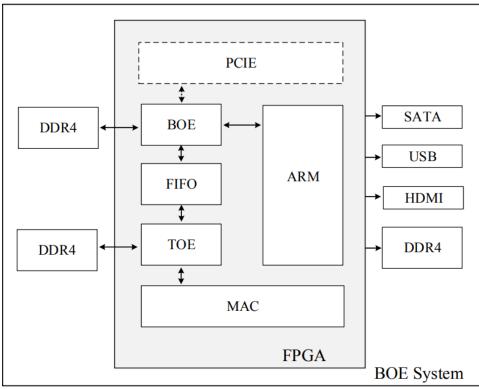


HPB区块链





硬件加速



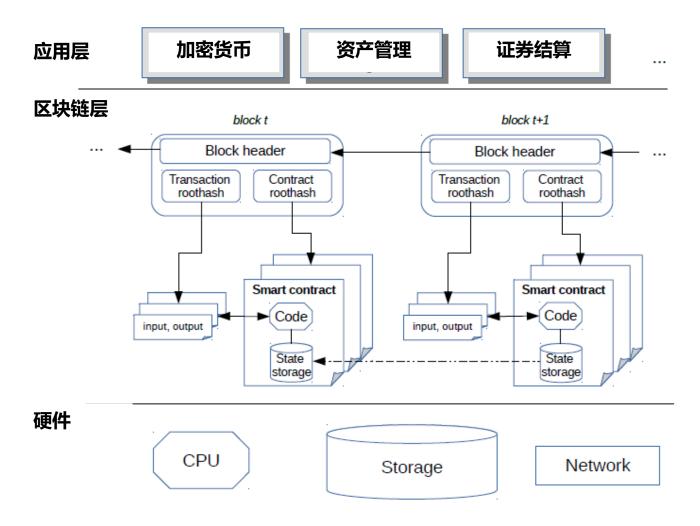
Blockchain Offload Engine (BOE)



区块链测试床



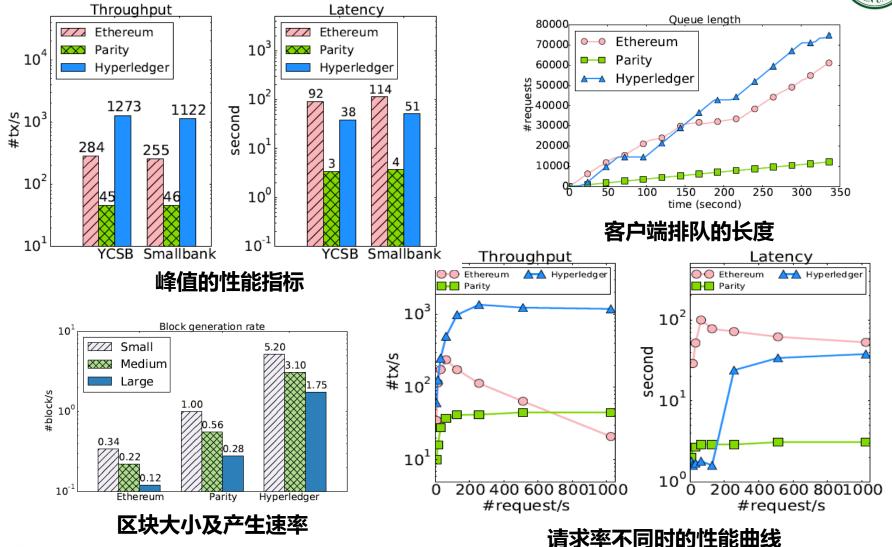
BLOCKBENCH: http://www.comp.nus.edu.sg/~dbsystem/blockbench/





区块链测试床







区块链测试床



	Layer	Examples
1	View 可视化界面	UTXO (BTC), Balance (ETH),
	- 1元 1元 FEI	Segments (Indigo)
2	Application 应用程序	Contracts (ETH), Processes (Indigo)
3	Storage 存储	DHT, Databases, CAN
4	Execution 执行环境	Compilers, VMs, Dockers
5	Consensus 共识协议	PoW, PoS, PBFT
6	Data Model 数据模型	Blocks, Transactions, Indexing
7	P2P Propagation P2P传播	Relay Network, Fast Relay Network,
	· < 1H	Falcon, Fibre
8	Side Plane	Sidechains, Payment Channels,
	侧链	Hubs, Oracles, Bridges to other
		networks
9	Network 网络	TCP/IP, UDP

区块链系统细粒度划分



区块链系统性能指标



9	Peak Memory Usage 峰值内存利用率	Execution	Data Model, Execution	RAM utilization (both on- and off-chain) for segment creation
10	执行时间	Execution	Execution	Time to create a segment, i.e., • Off-Chain: CPU Cycles behind an agent/method execution. • On-Chain: CPU cycles for running the Validation Code (as we do in Indigo wrt the Manifest). This could be compared with the CPU cycles for a similar problem-statement solved by a Smart Contract (ETH).
11	Availability 可用性	Reliability 可靠性	Network, P2P Propagation, Consensus	The ratio of uptime to total time.





区块链测试标准的意义



	Comparison Type	Example
1.	Standards-Based 提供标准	How do our results compare with an industry-accepted standard?
2.	Inter-Stack Benchmark 不同平台比较	How does the performance of technology stack A compare with technology stack B, as a solution to the same high-level problem?
3.	Intra-Stack Benchmark 平台内优化	How does the performance of a given technology stack change as it configured with different options?

优化性能、完善系统



超级账本测试性能



Test	Name	Succ	Fail	Send Rate	Max Latency	Min Latency	Avg Latency	75%ile Latency	Throughput
1	simpleOperations	4000	0	200 tps	1.56 s	0.27 s	0.85 s	1.09 s	196 tps
2	simpleOperations	10000	0	499 tps	14.19 s	1.83 s	8.81 s	11.48 s	388 tps
3	simpleOperations	20000	0	950 tps	50.14 s	4.05 s	33.80 s	37.76 s	370 tps
4	query	2000	0	100 tps	0.03 s	0.00 s	0.01 s	0.01 s	100 tps
5	query	4000	0	200 tps	0.03 s	0.00 s	0.01 s	0.01 s	200 tps
6	query	10000	0	499 tps	0.06 s	0.00 s	0.01 s	0.01 s	498 tps
7	query	20000	0	1000 tps	0.46 s	0.00 s	0.08 s	0.10 s	999 tps

resource consumption

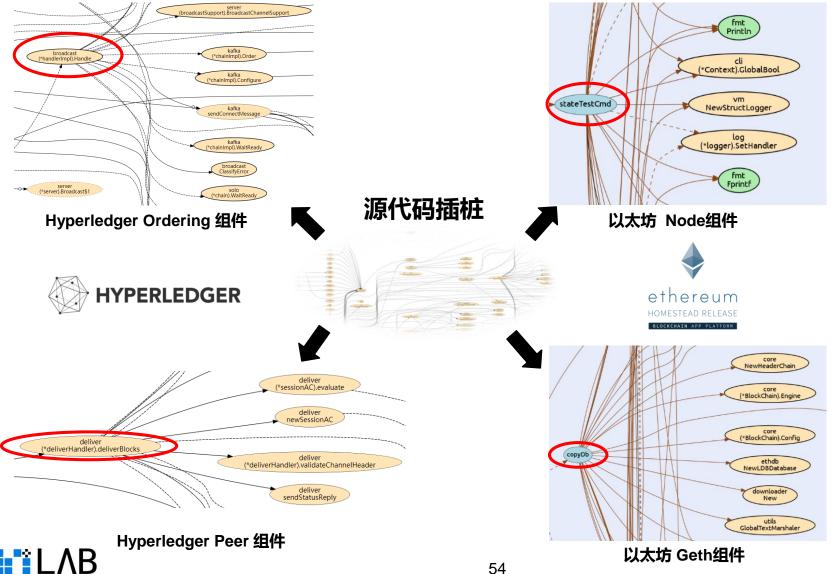
TYPE	NAME	Memory(max)	Memory(avg)	CPU(max)	CPU(avg)	Traffic In	Traffic Out
Process	node local-client.js(avg)	114.3MB	110.0MB	37.36%	28.85%	-	-
Docker	dev-peer1.org2.example.cole-v0	6.3MB	6.1MB	8.93%	7.28%	4.8MB	2.6MB
Docker	dev-peer1.org1.example.cole-v0	6.2MB	6.0MB	8.85%	7.40%	4.8MB	2.7MB
Docker	dev-peer0.org2.example.cole-v0	6.0MB	5.8MB	6.04%	5.07%	3.3MB	1.9MB



区块链系统静态分析

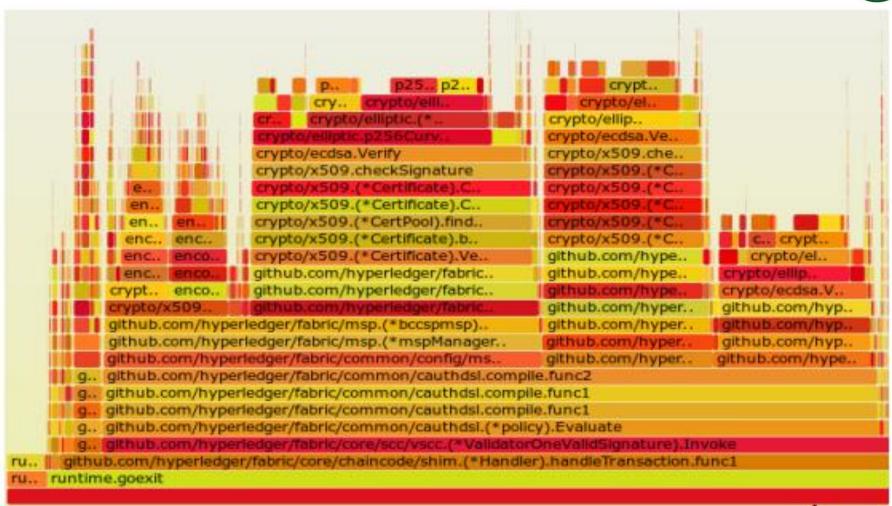
WWW.INPLUSLAB.COM





区块链系统动态分析





Hyperledger的MSP服务的函数调用次数及调用栈



区块链系统动态分析



[36m2018-05-19 08:17:40.491 UTC [orderer/common/broadcast] Handle -> DEBU 48e[0m [channel: mychannel] Broadcast is processing normal message from 172.17.0.1:51660 with txid '322a0a44e60d011f44aae05c708a6e8d07b890de22646d9ea868454a8eea0cf5' of type ENDORSER TRANSACTION

[36m2018-05-19 08:17:40.491 UTC [policies] Evaluate -> DEBU 48f[0m == Evaluating *policies.implicitMetaPolicy Policy /Channel/Writers == with txid '322a0a44e60d011f44aae05c708a6e8d07b890de22646d9ea868454a8eea0cf5' of type ENDORSER_TRANSACTION

[36m2018-05-19 08:17:40.491 UTC [policies] Evaluate -> DEBU 490[0m This is an implicit meta policy, it will trigger other policy evaluations, whose failures may be benign with txid '322a0a44e60d011f44aae05c708a6e8d07b890de22646d9ea868454a8eea0cf5' of type ENDORSER TRANSACTION

[36m2018-05-19 08:17:40.491 UTC [policies] Evaluate -> DEBU 491[0m == Evaluating *policies.implicitMetaPolicy Policy /Channel/Application/Writers == with txid '322a0a44e60d011f44aae05c708a6e8d07b890de22646d9ea868454a8eea0cf5' of type ENDORSER_TRANSACTION

Hyperledger Fabric 运行时产生的日志

Blockchain peer batch Consensus transactions



高性能区块链系统构建



BaaS (Blockchain as a Service)

应用层(Dapp)	天气预 智	能交 医疗(通 康	建 消费习 惯分析	地理环 境	生物信息	视频推 荐	+ 更多 ···	
应用框架层	Tensorflow Hadoop	SparkML Caffe	Theano HTCondor	H2O SystemM	_	充计 :合约	自定义AI服 务 + 更多	
区块链平台	以太坊		当级账本 HYPERLEDGER		Parity		Ripple	0-0 A 0-0 0-0
Kubernetes 云平台 基础设施	部署存储	消度		9络	+ 更多		云运维 Operator)	





谢谢!