

Raft in Theory and Practice

Building a Distributed Consistent Database



关注公众号回复help, 可获取更多经典学习 资料和文档,电子书



About me

- Tang Liu
- Chief Architect of PingCAP
- TiDB/TiKV
- Open sources: LedisDB, go-mysql, etc...



Agenda

- Why Raft
- What is Raft
- The optimization
- Multi-Raft



Why



One Node





One Node





Replication





Async Replication





Sync Replication





Sync Replication





Quorum , 2 N + 1





Failover





Brain Split





Linearizability





Consensus Algorithm

- Agreement on shared state
- Recovery from server failures autonomously
 - Minority failure: no problem
 - Majority failure: lose availability, retain consistency
- Strong consistency (linearizability)



Raft



Keypoint

- Replicated Log for replicated state machine
- Leader, Follower and Candidate
- Term and Leader election
- Log Replication
- Membership Change



Replicated State Machine





State

• Leader

- Only one Leader
- Elected by the majority of the peers
- Handles all the client requests
- Follower
 - Receives Replicated Logs from the Leader
- Candidate
 - For campaign



Term









Log Replication









Membership Change





Membership Change









Membership Change





Optimization



Pre-Vote





Pipeline











Writing to the Leader in parallel





Non-Voter





Multi-Raft



How to scale?

- Logical split
- Just Split && Move







Scale-out (initial state)





Scale-out (add new node)





Scale-out (balancing)





Scale-out (balancing)











Questions?



Oracle 2017 技术嘉年华交流

