



中国科学技术大学
University of Science and Technology of China

Raft—consensus algorithm

马子杰

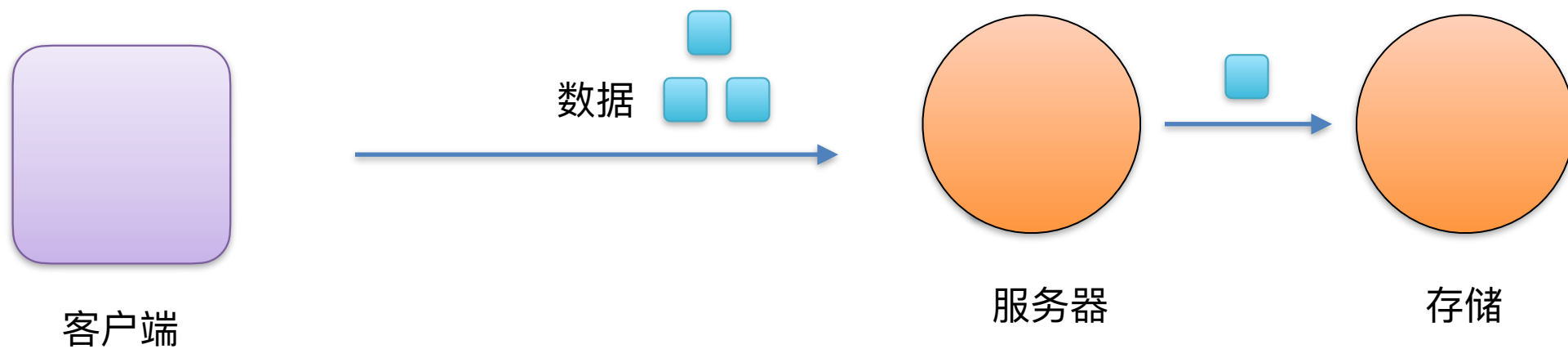


Introduction

数据备份问题



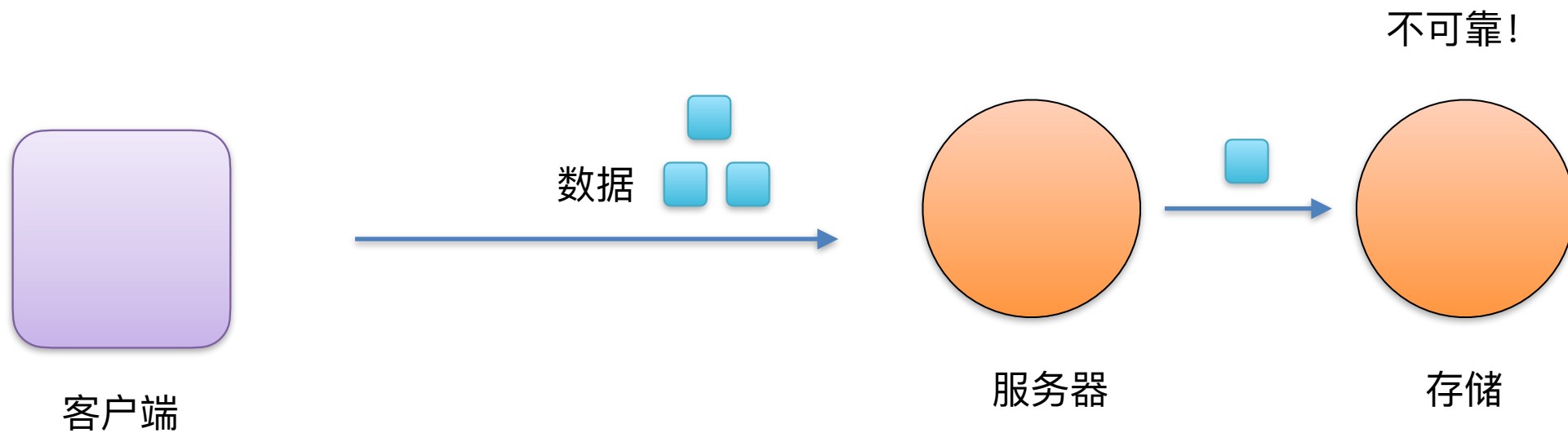
中国科学技术大学
University of Science and Technology of China



数据备份问题



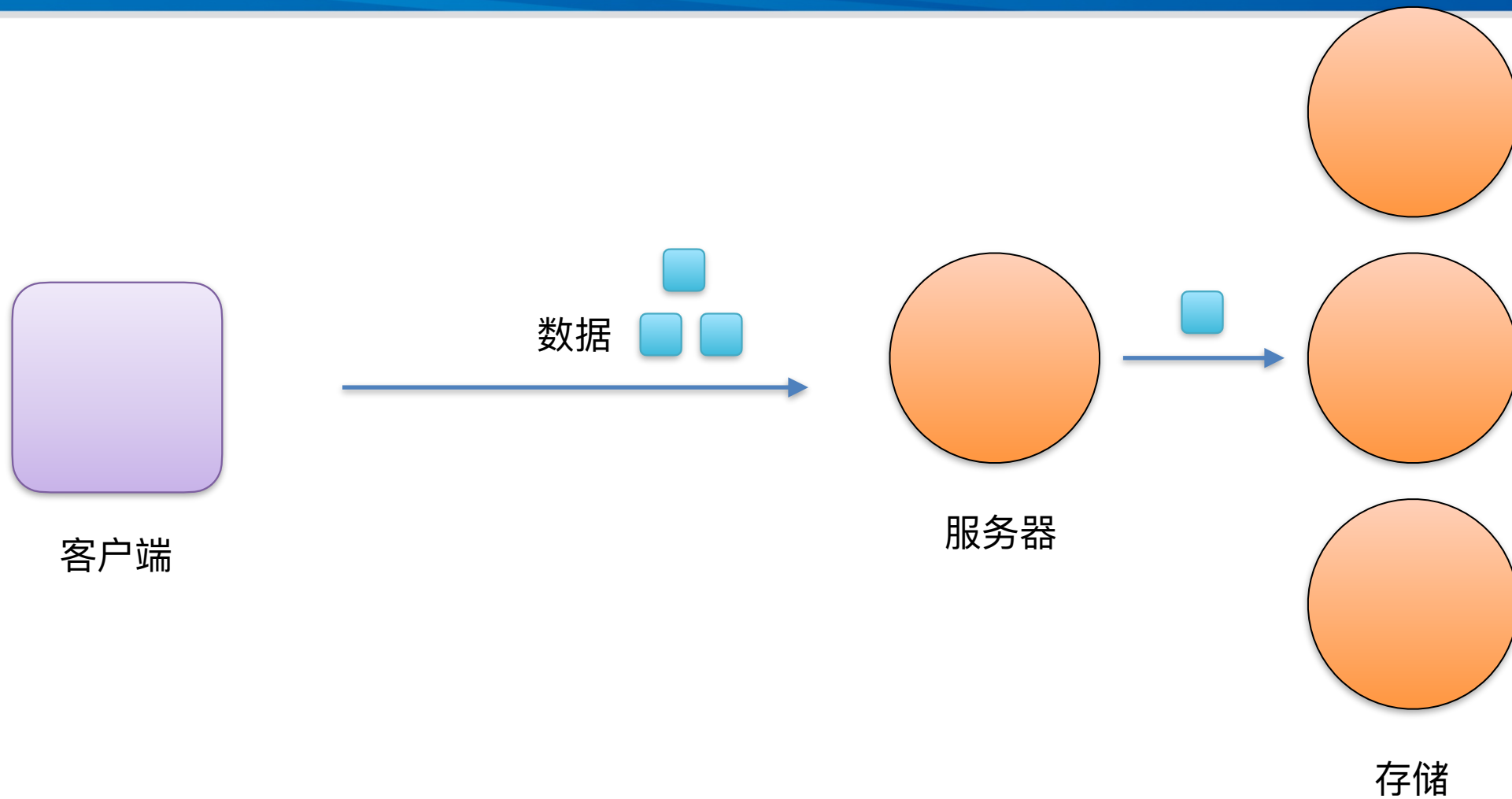
中国科学技术大学
University of Science and Technology of China



数据备份问题



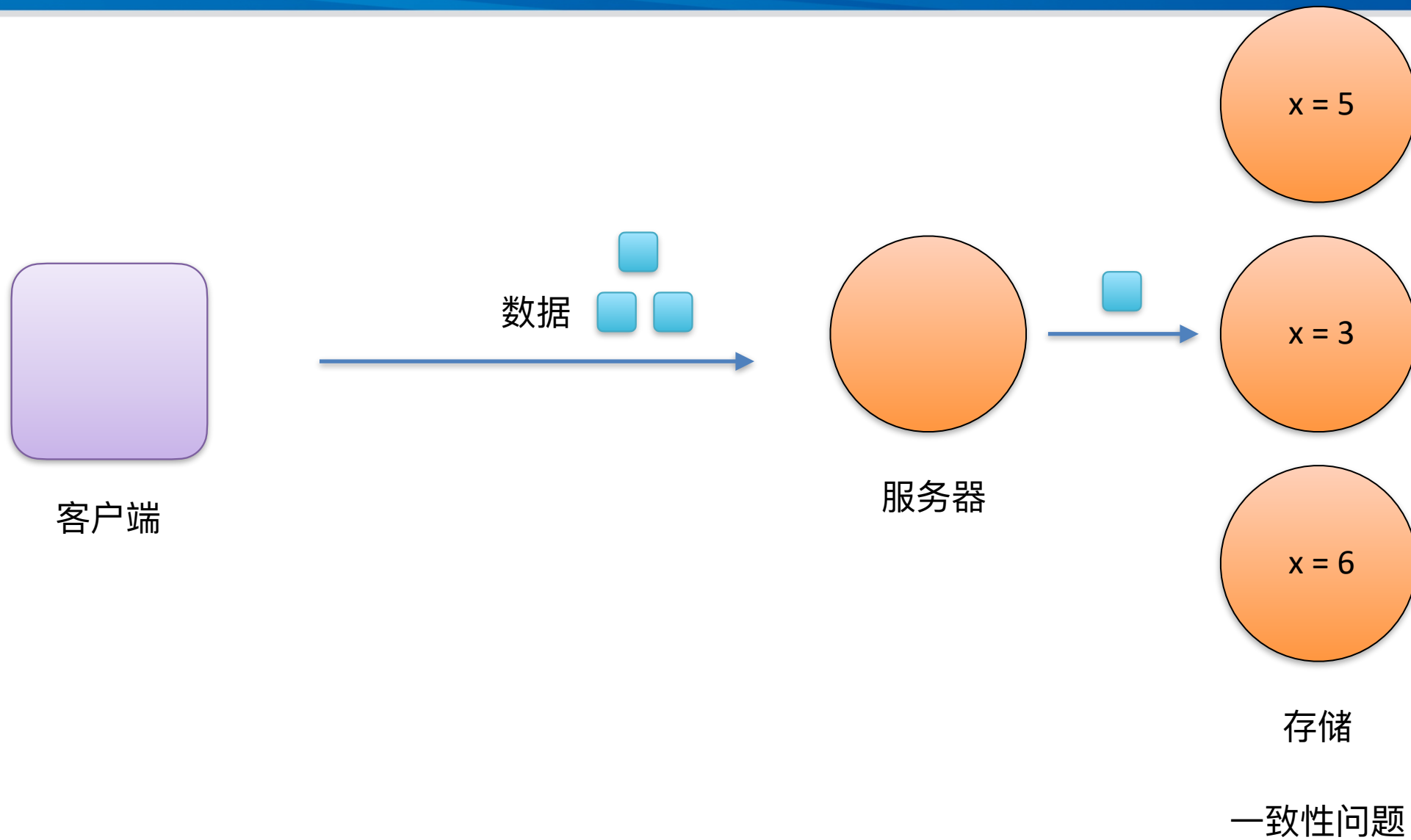
中国科学技术大学
University of Science and Technology of China



数据备份问题



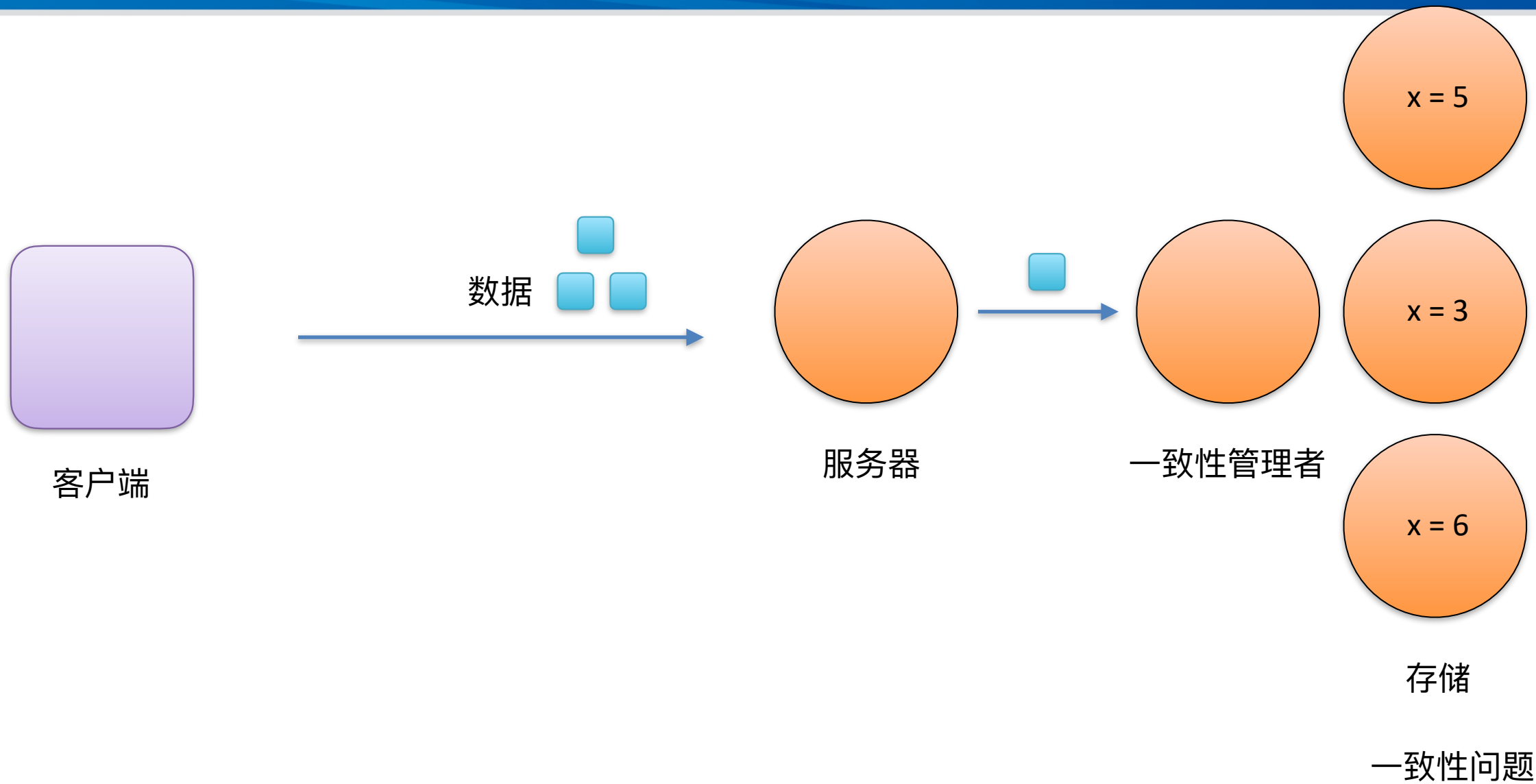
中国科学技术大学
University of Science and Technology of China



数据备份问题



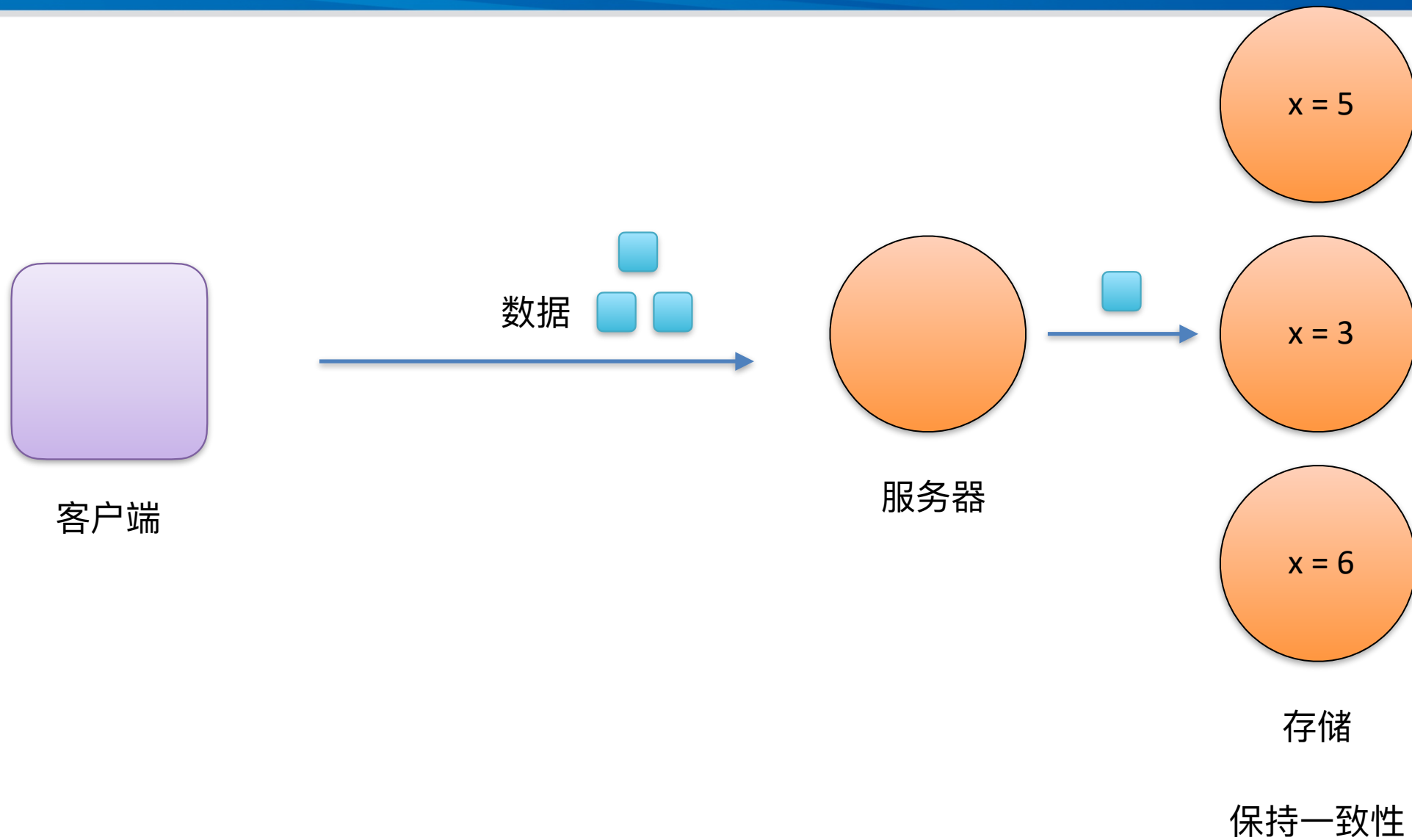
中国科学技术大学
University of Science and Technology of China



数据备份问题



中国科学技术大学
University of Science and Technology of China





Raft

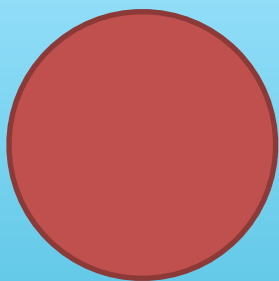
- 共识协议
- 实现了强一致性
- 强领导者协议
- 分布式存储系统
- 假设有 $2N+1$ 个节点的情况下，至多在 N 个节点宕机的情况下，对外正常服务，保持高可用性

- 线性一致性 OceanBase Spanner
- 顺序一致性 Zookeeper
- 因果一致性 CockroachDB
- 最终一致性 DNS -> 写后读一致性 单调读一致性

Raft



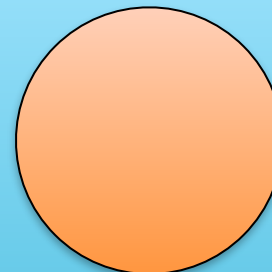
角色



Leader

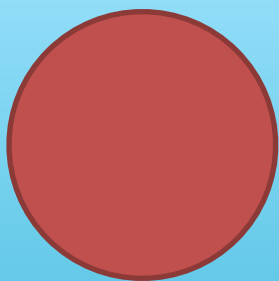


Follower



Candidate

角色



Leader

1. 负责接收命令
2. 将命令写入日志
3. 发送日志到其他主机中
4. 日志提交 返回结果到客户端
5. 空闲中发送心跳

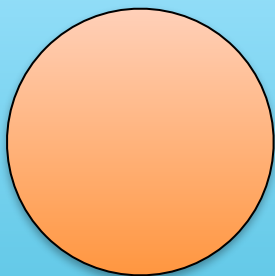
角色



Follower

1. 负责接收日志
2. 定时收到心跳
3. 未收到心跳，转变为Candidate

角色



Candidate

1. 负责进行选举
2. 向每台主机发送投票请求
3. 收到大多数主机的投票，变成 Leader
4. 收到Leader心跳之后，变化为 Follower

选举过程



A:Follower



B:Follower



C:Follower

选举过程

我没有收到心跳，我要
变成Candidate!!!



A:Follower



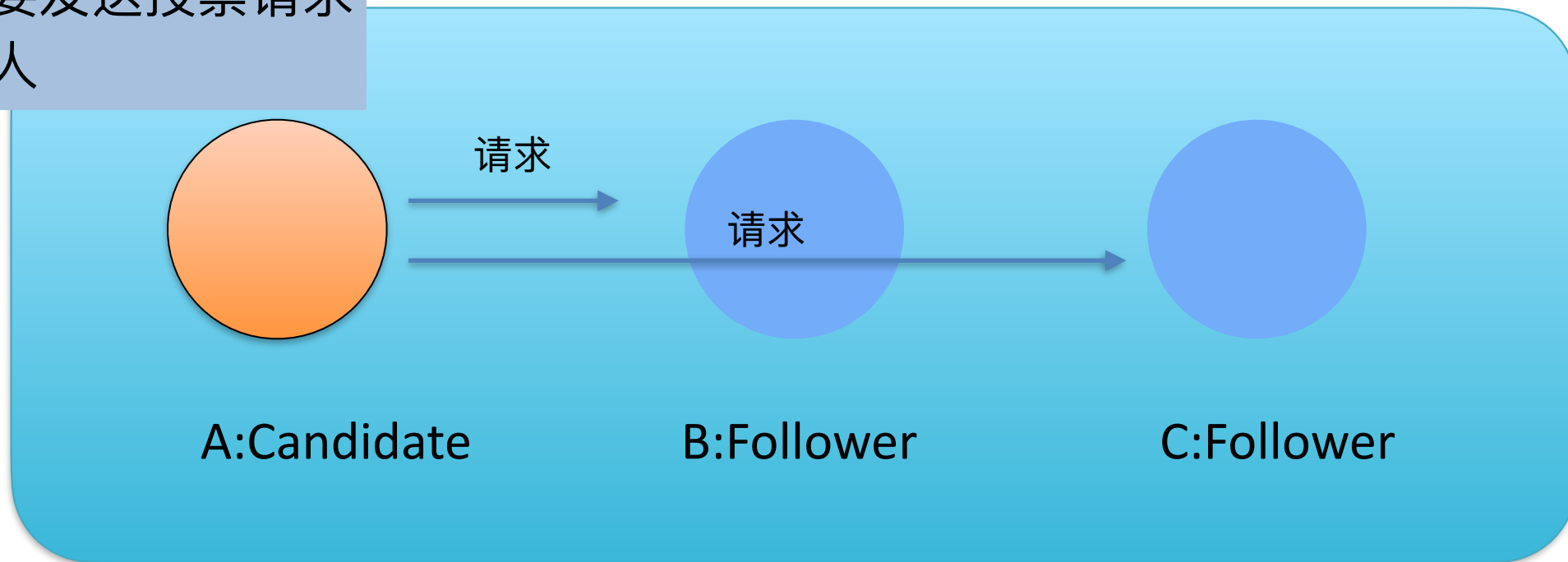
B:Follower



C:Follower

我给自己先投票
之后我要发送投票请求
给其他人

选举过程



选举过程

我自己有一票
B投给我一票
我收到2票
我要成为Leader



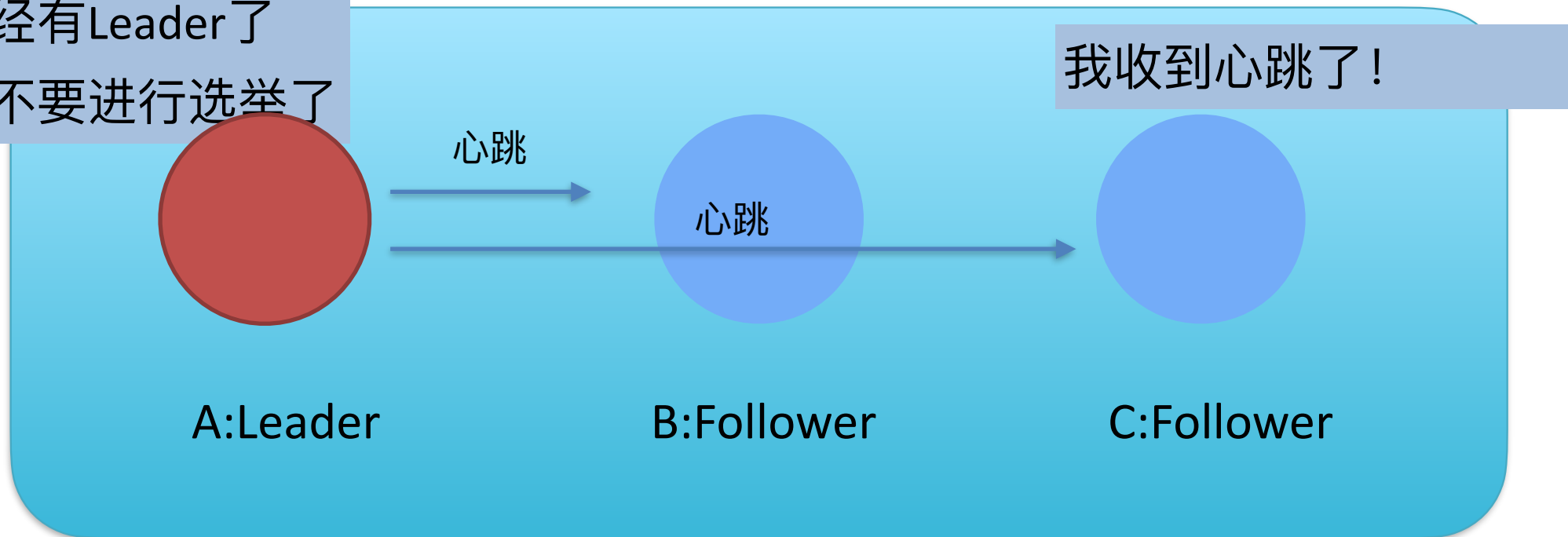
领导者选举



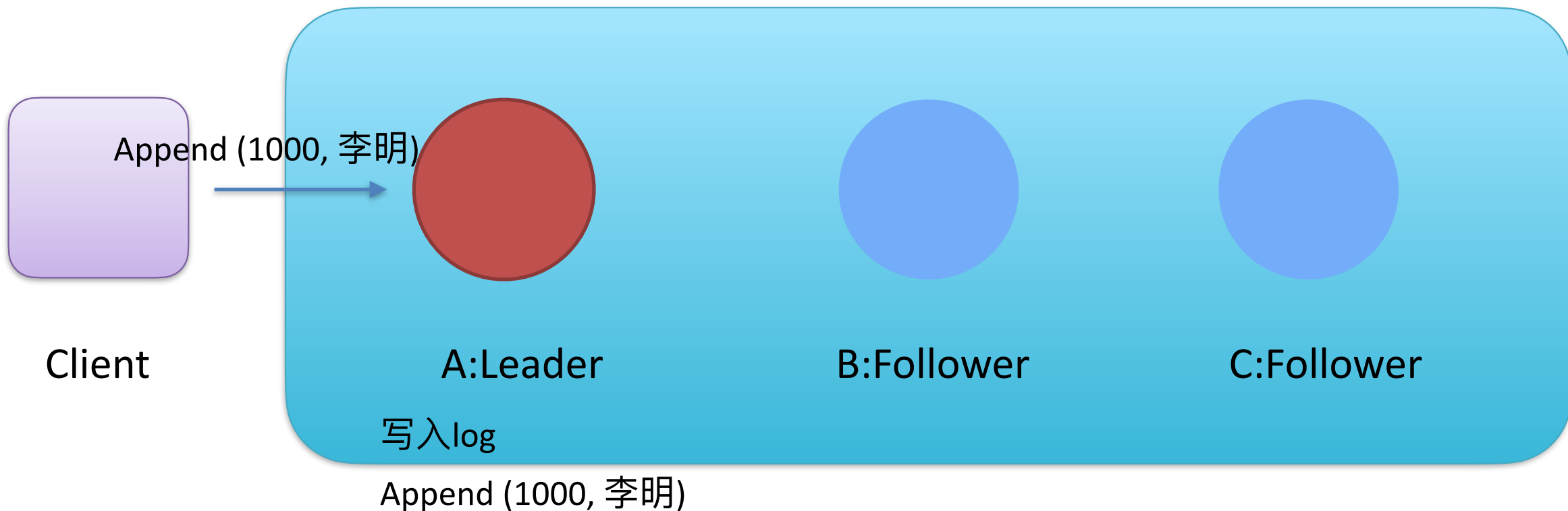
中国科学技术大学
University of Science and Technology of China

我是Leader了
我要发送我的心跳
证明已经有Leader了
其他人不要进行选举了

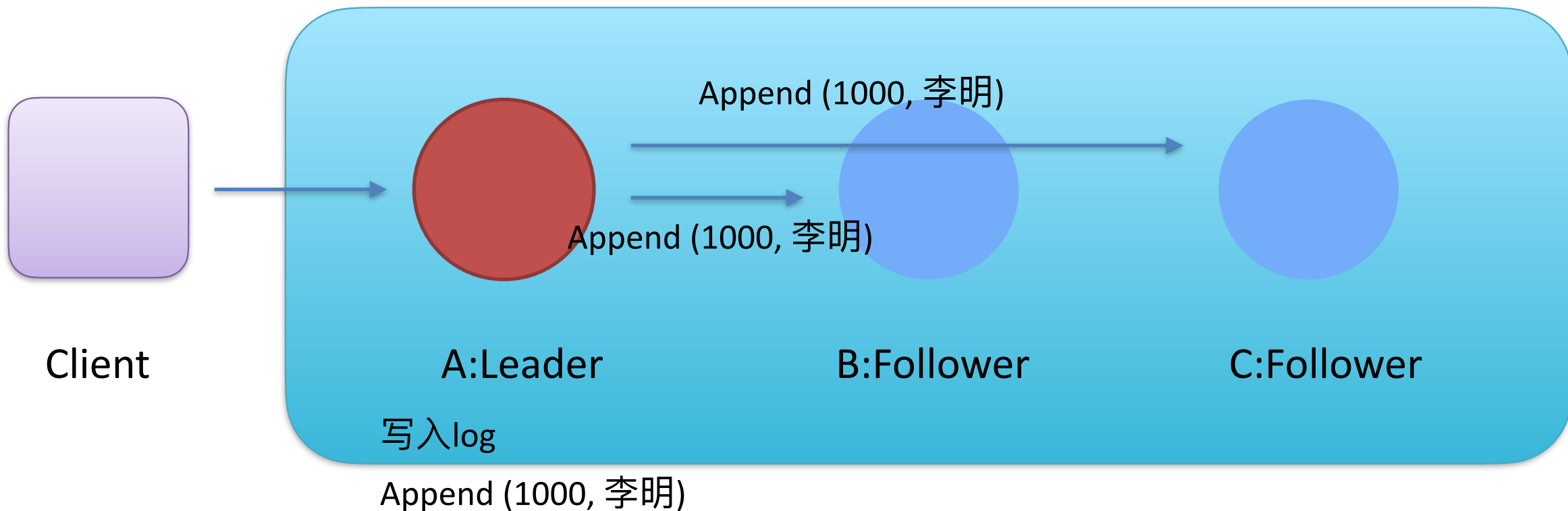
选举过程



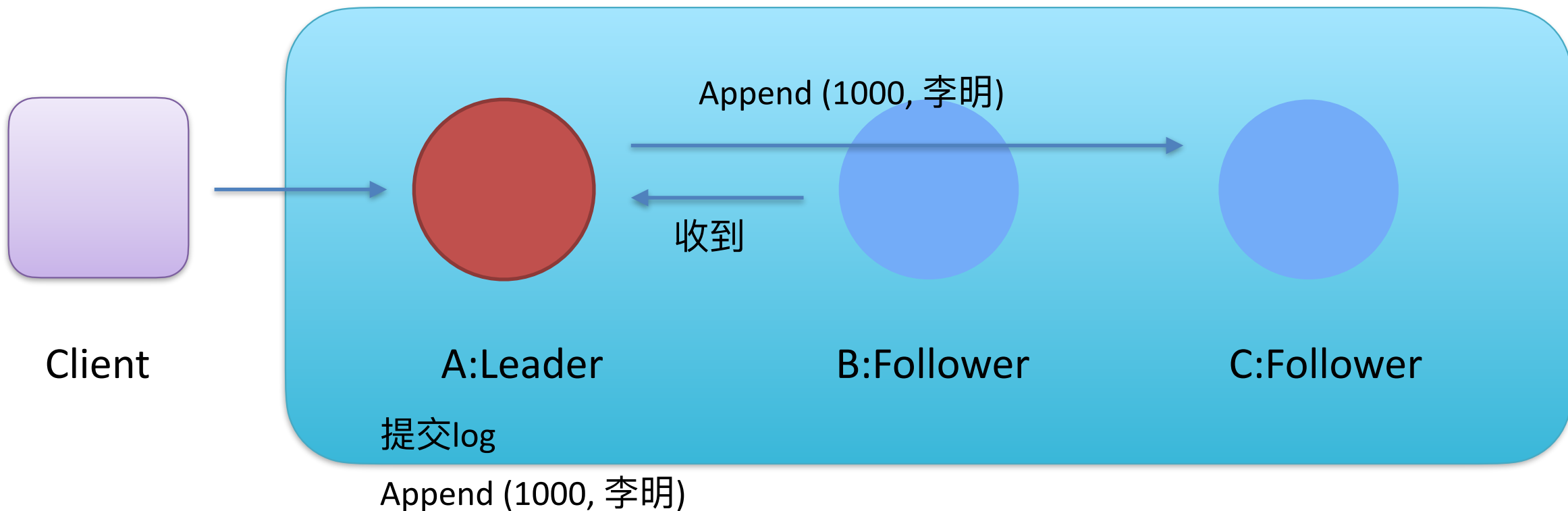
日志复制阶段



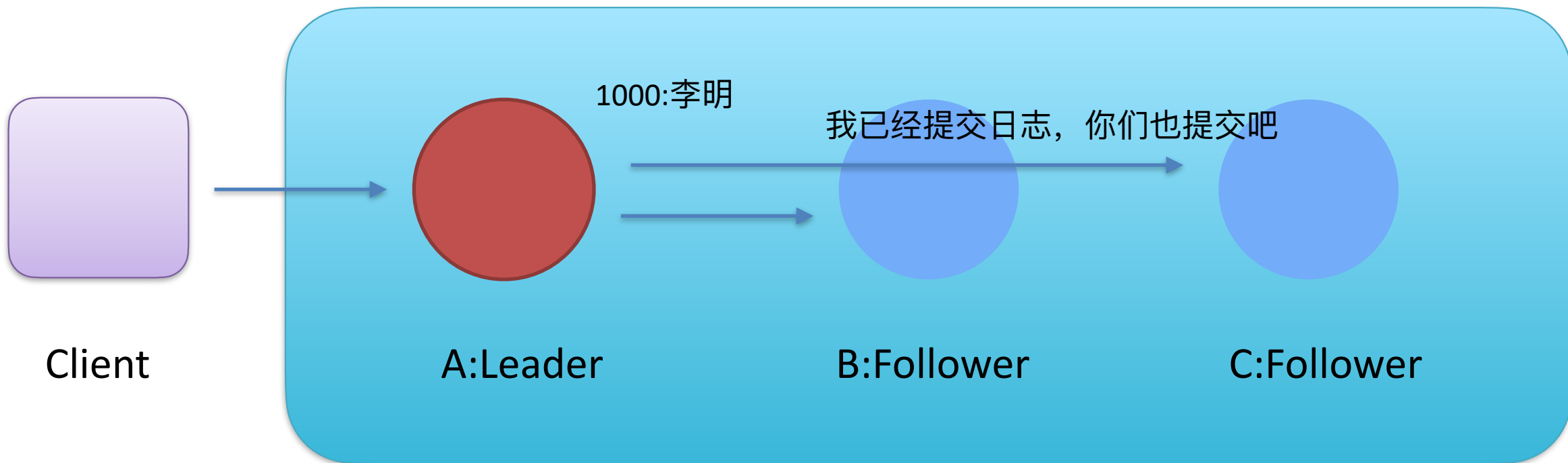
日志复制阶段



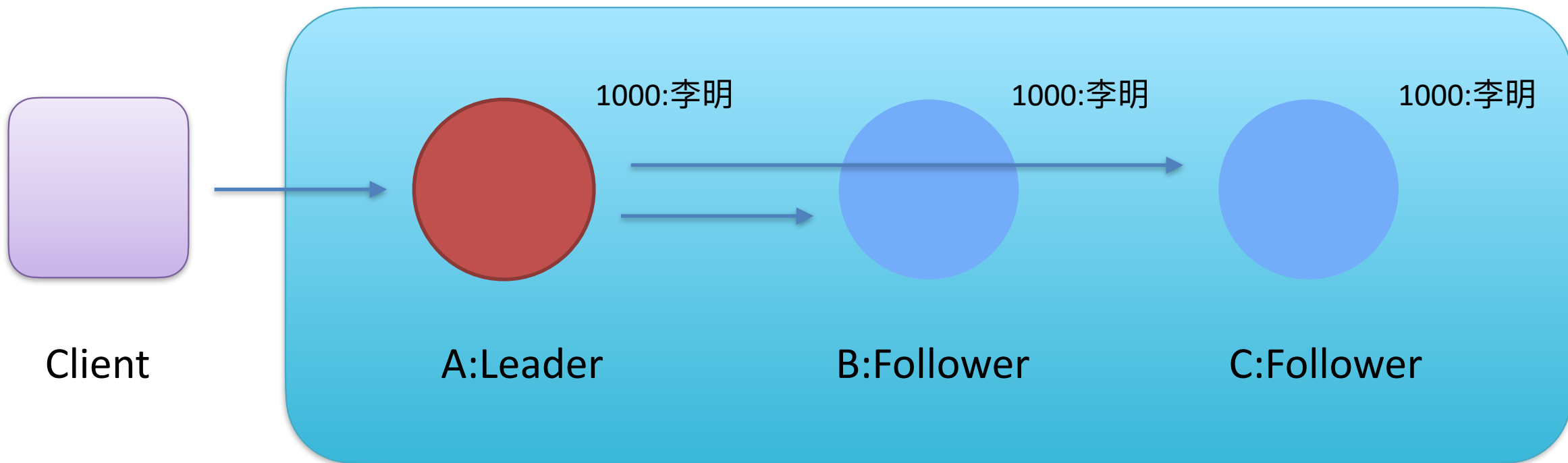
日志复制阶段



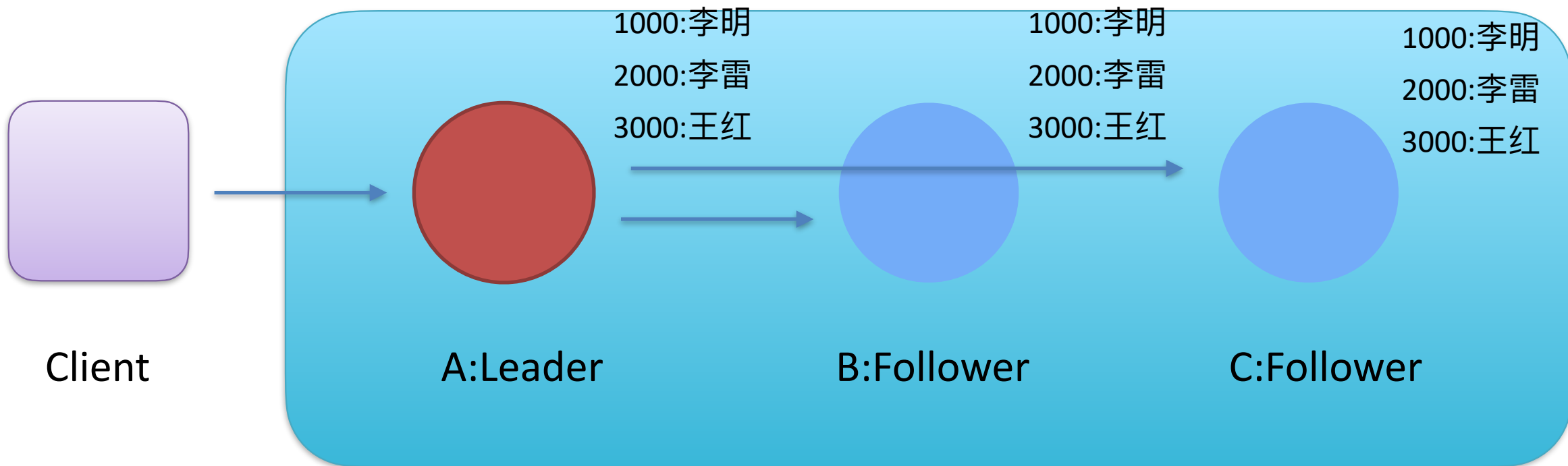
日志复制阶段



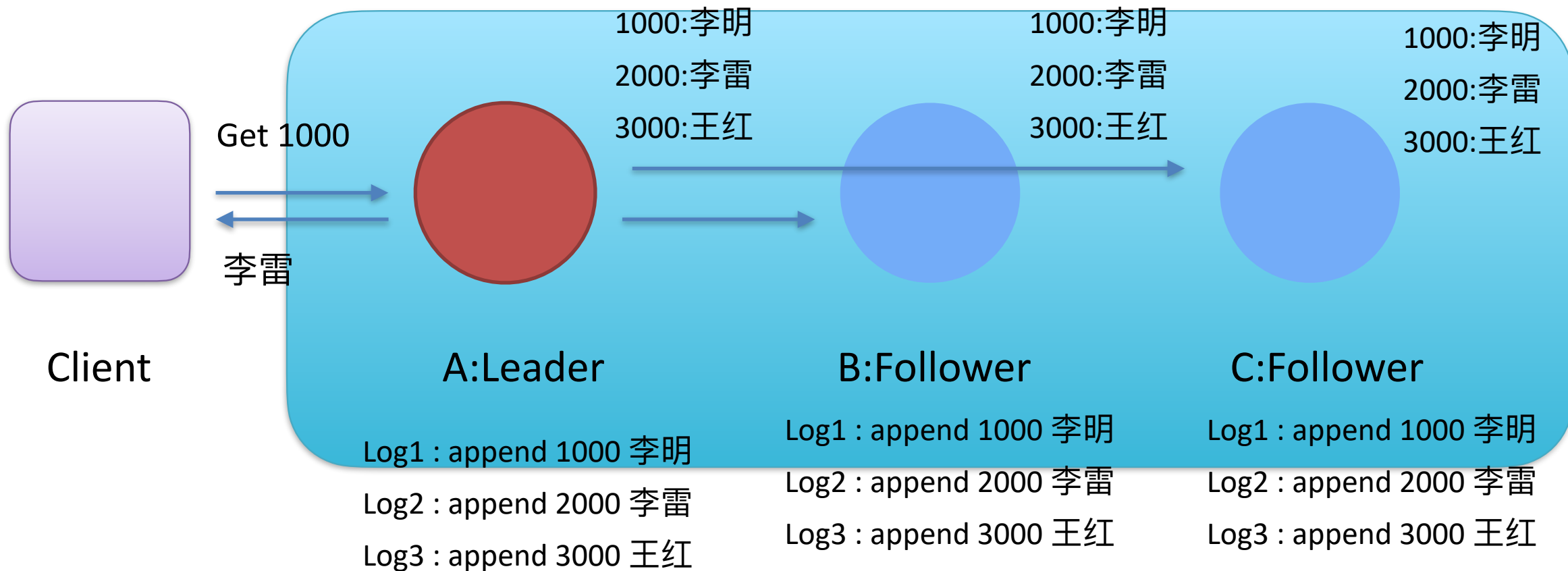
日志复制阶段



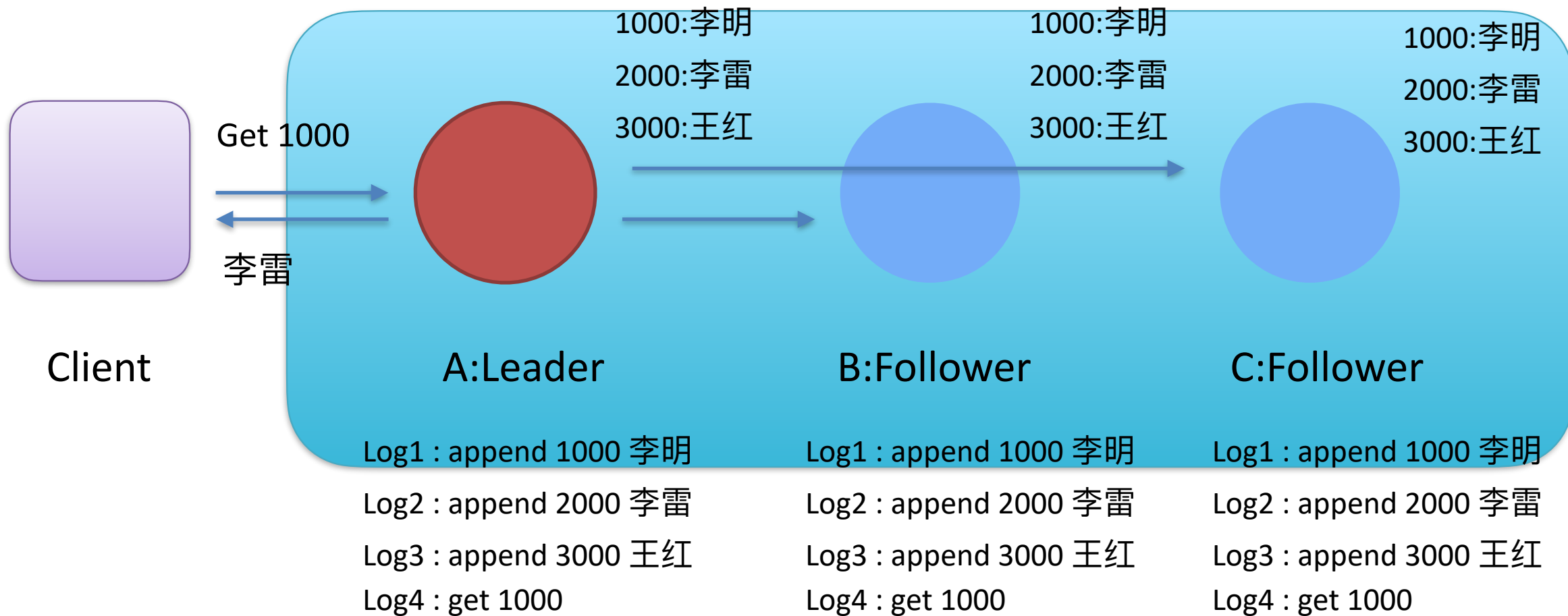
日志复制阶段



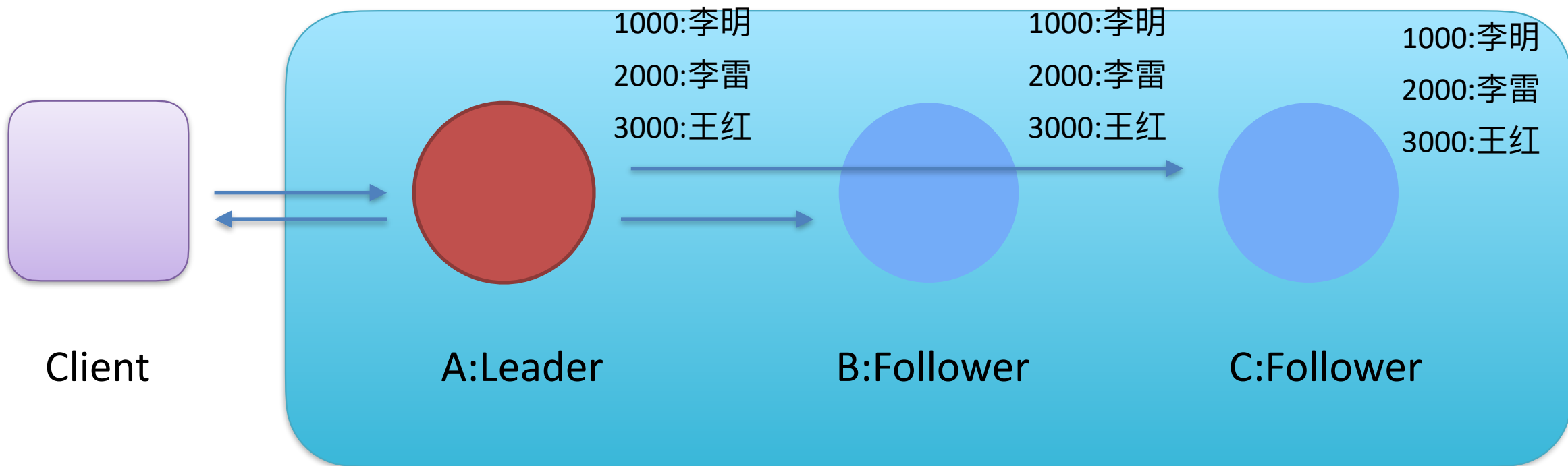
日志复制阶段



日志复制阶段

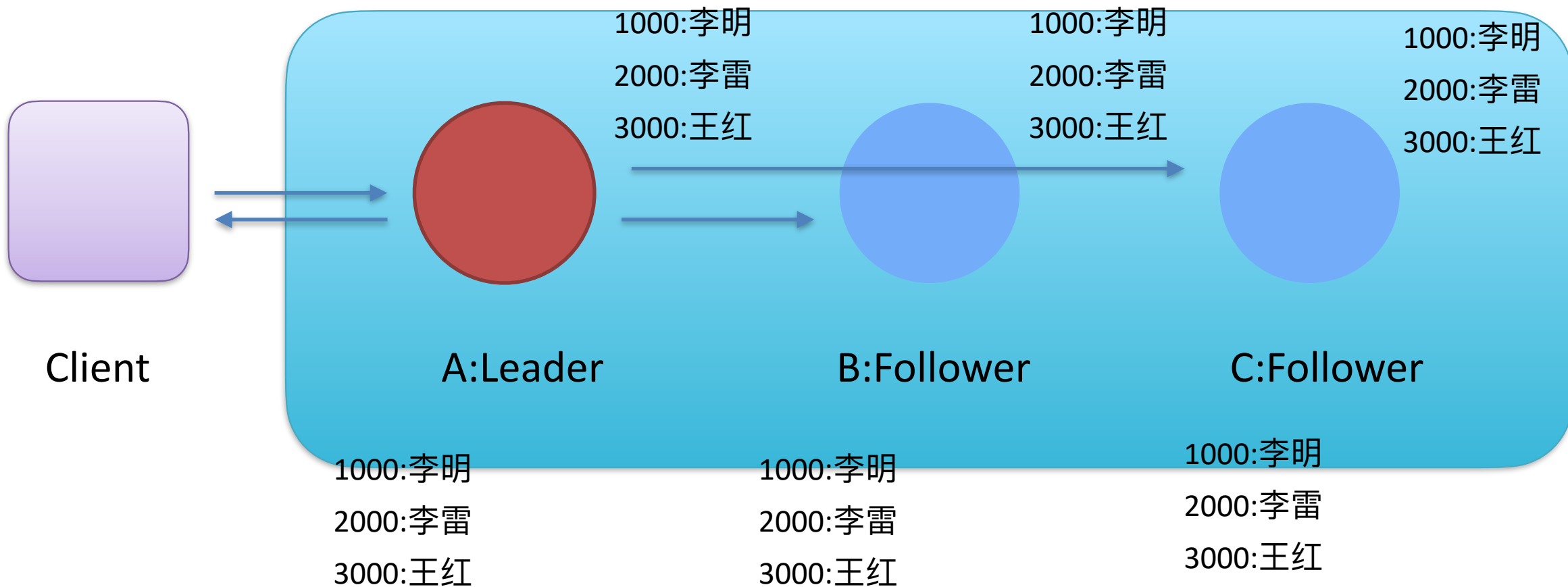


日志复制阶段

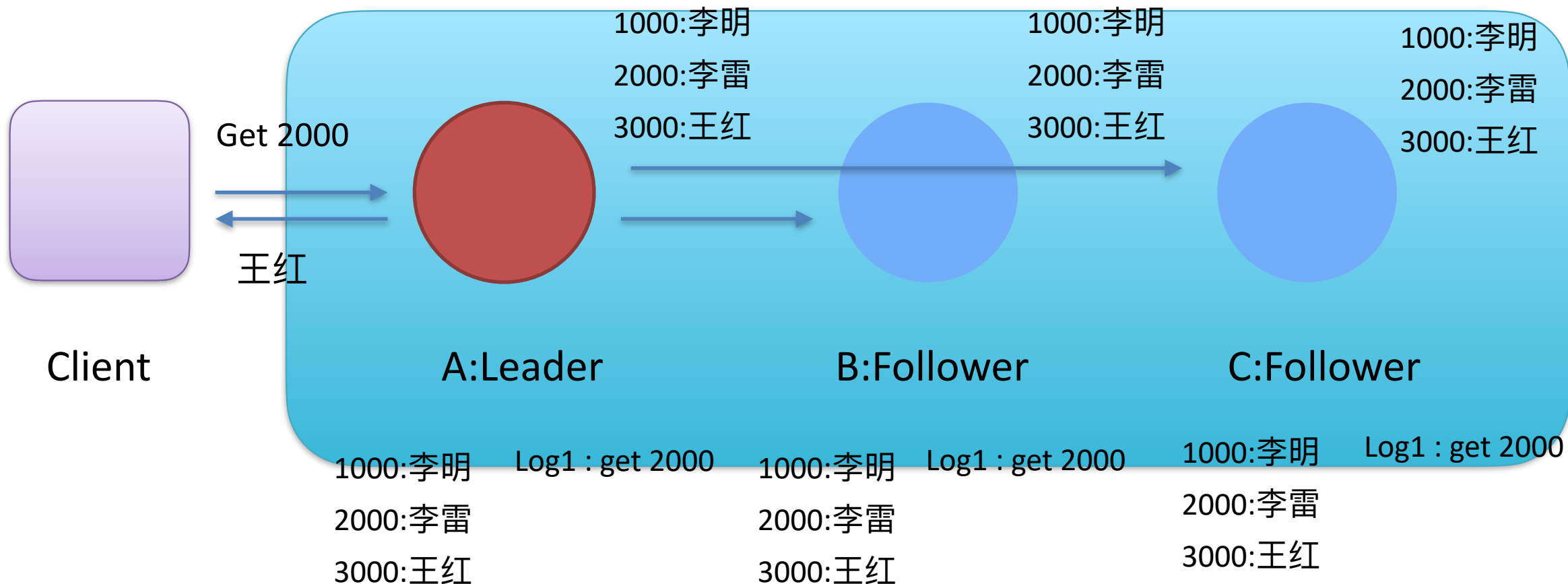


Log太多了 导致我的Log集合存不下!!!

日志复制阶段

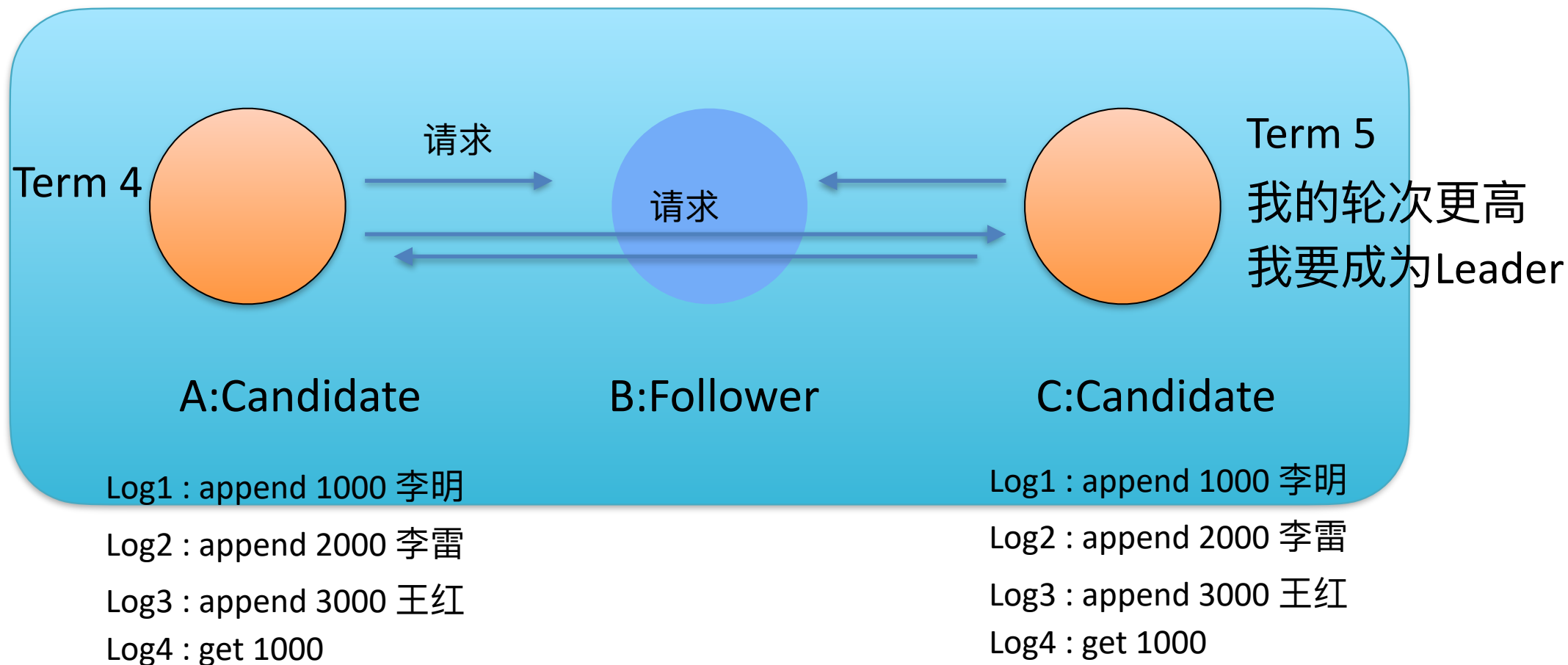


日志复制阶段



- Term在Raft中是一个时间概念
- 每个term最多只有一个leader
- 如果leader宕机，term增加，其他的follower会重新进行选举
- 在Raft中，term越大，代表数据越新
- Term小的主机会听从term大的主机

case 1 term不相等

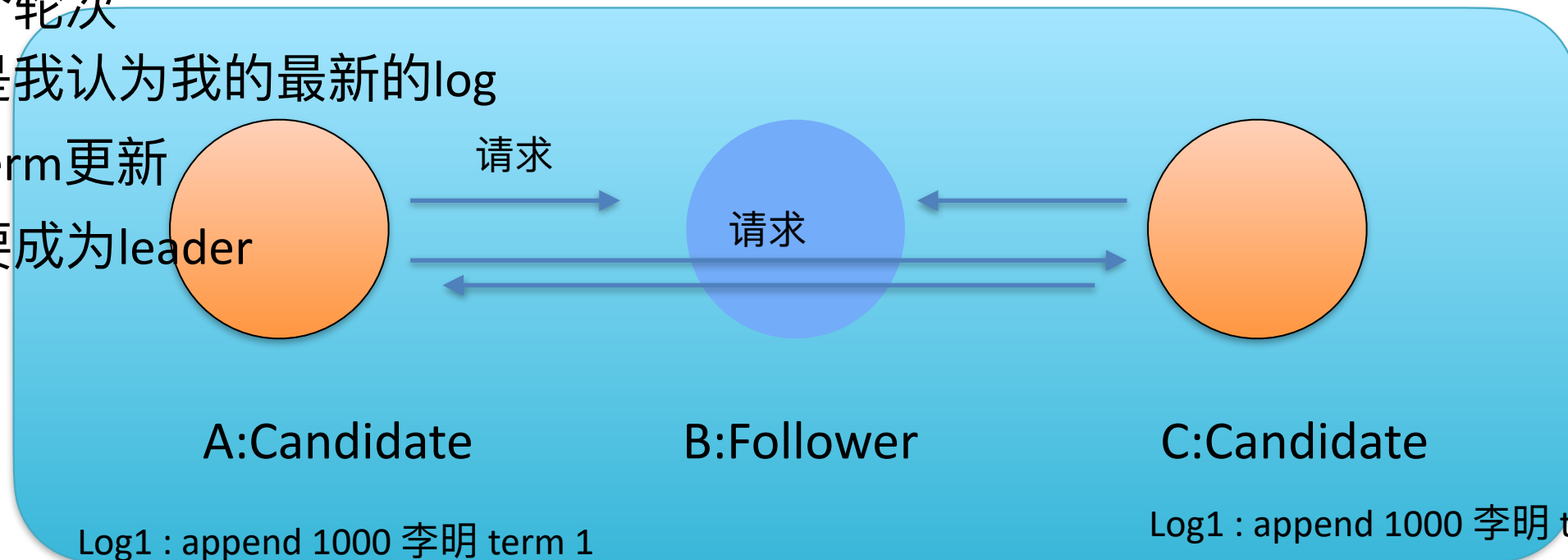


case 2 term相等

虽然咱们两个处于
一个轮次

但是我认为我的最新的log
的term更新

我要成为leader



Log1 : append 1000 李明 term 1

Log2 : append 2000 李雷 term 2

Log3 : append 3000 王红 term 3

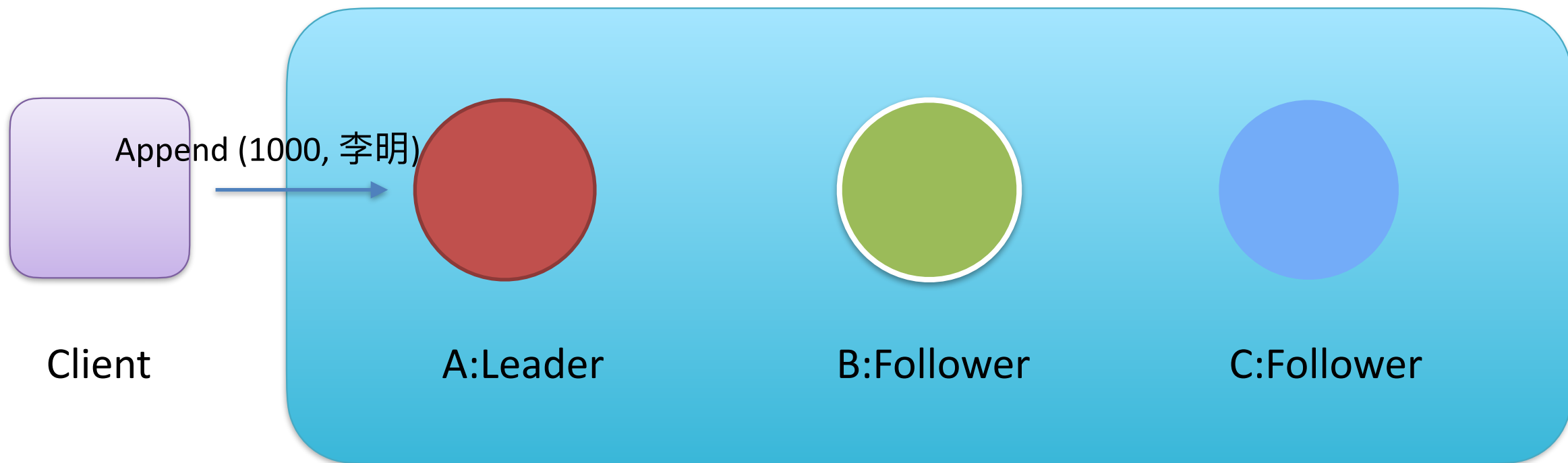
Log4 : get 1000 term 6

Log1 : append 1000 李明 term 1

Log2 : append 2000 李雷 term 2

Log3 : append 3000 王红 term 3

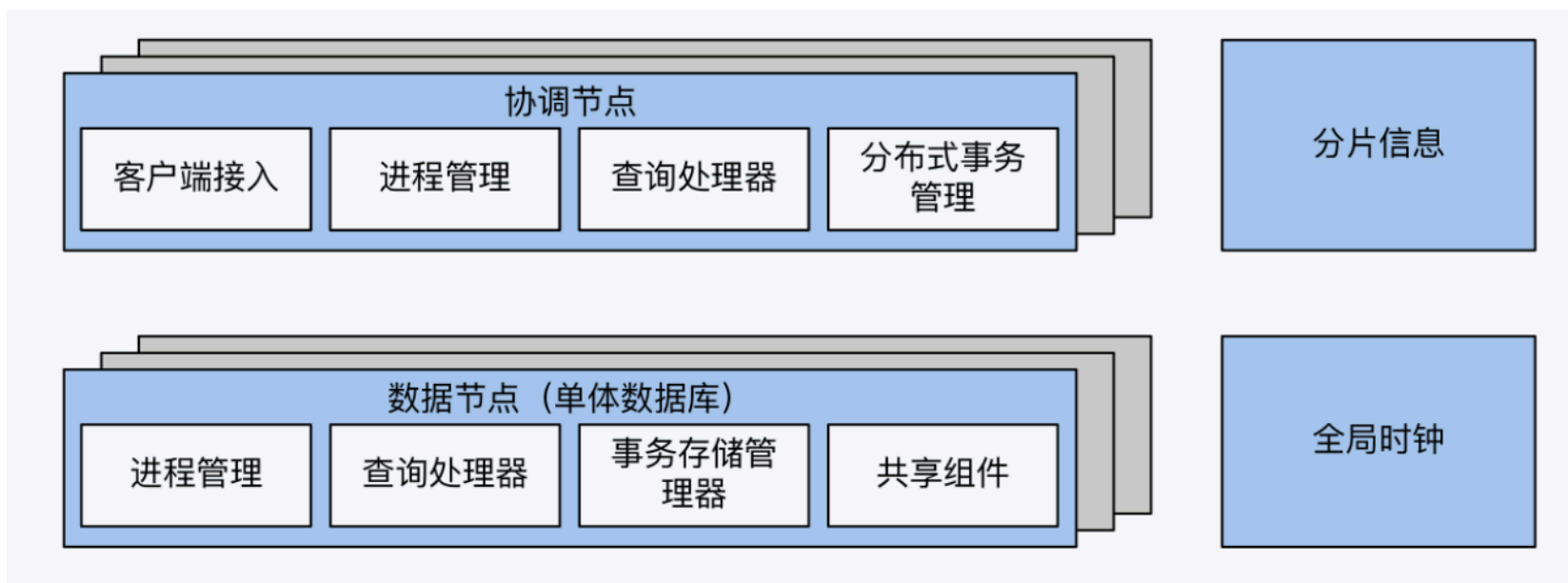
case 3 日志不统一



case 3 日志不统一

	index 11	index 12	index 13	Index 14
A:Leader Term5	Term 4	Term 4	Term 5	Term 5
B	Term 4	Term 4	Term 5	
C	Term 4	Term 4	Term5	
D	Term 4	Term 4	Term 4	

– 传统SQL (proxy) 与NewSQL

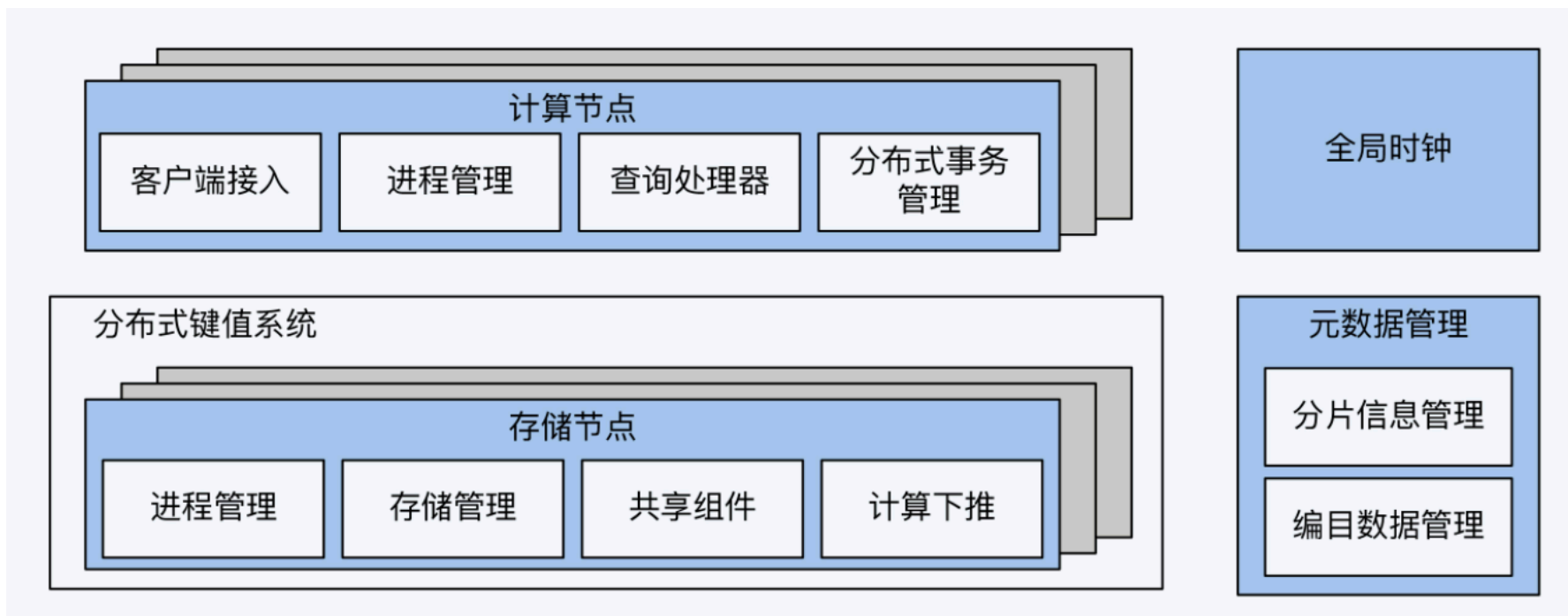


深入理解Raft



中国科学技术大学
University of Science and Technology of China

- Multi-raft-group
- 传统SQL (proxy) 与NewSQL





- 如何理解实现了强一致性
- Consul针对raft实现了三个一致性模型
- Default(无法防止脑裂)
- Consistent(强一致性)
- Stale(可以读到过时数据)



中国科学技术大学
University of Science and Technology of China

Q & A