

Distributed Systems

An Intro

TYR

Definition

- A collection of computers running in a network
- communicate and coordinate towards a goal

Things to Consider

- Communication / coordination
- fault tolerance (software crash, hardware failure, network down)
- Scalability

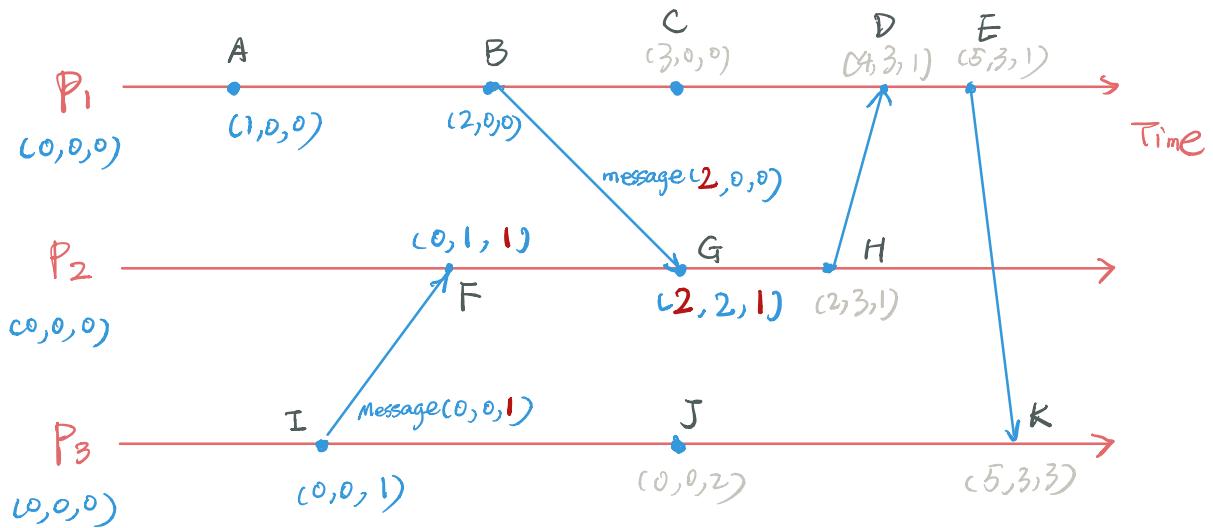
Concepts

- **Node**: an independent computer
- **Network**: the channel for communication
 - could be heterogeneous
 - communication has a cost!
- **Protocol**: rules for transmitting information

Comparison

	Solo	Cloud computing	fully distributed
global clock	strict	yes	no
environment	simple, controlled	managed, a little moving parts	nodes come and go
admin	single	single	no
security	easy to implement	harder	hardest
fault-tolerance	bad	good	best
protocol	n/a	private/easy to change	public/hard to change
peer discovery	n/a	static/dynamic (consul)	gossip /
upgrade	easy	rolling / blue-green	hard (consensus)
consensus	n/a	static/paxos/zab/raft	PBFT/PoX

Vector Timestamp



Causally-related

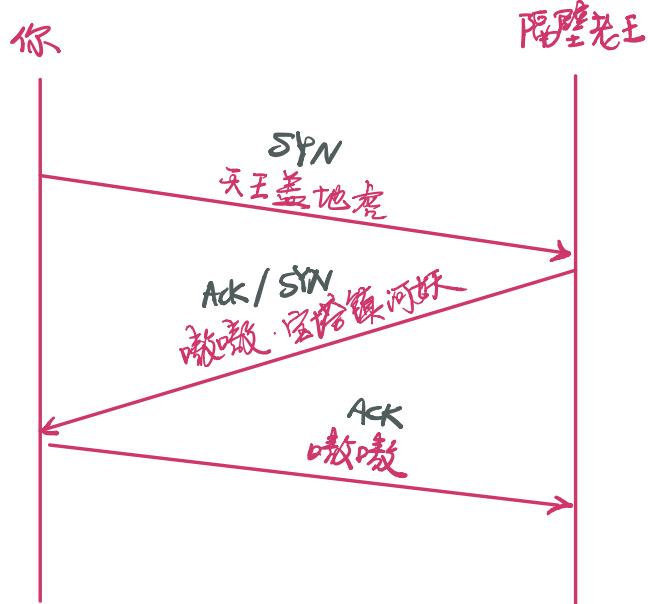
$$A \rightarrow B \quad (1, 0, 0) \prec (2, 0, 0)$$

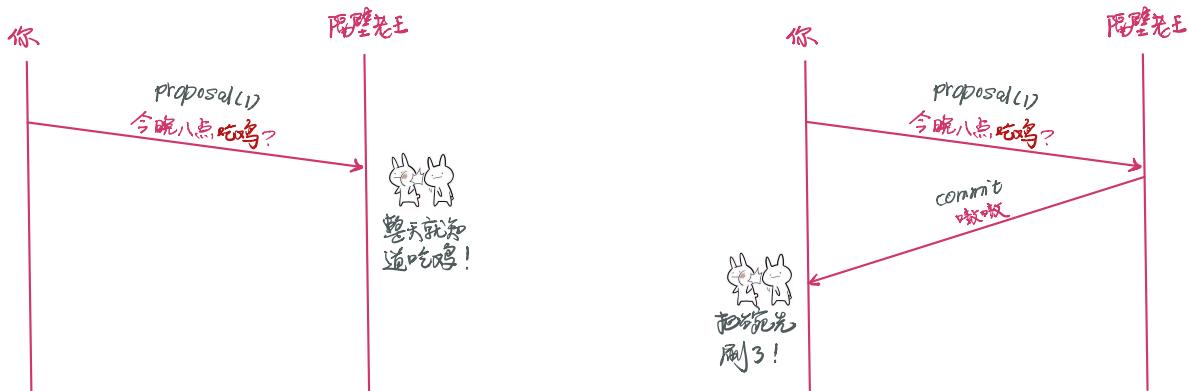
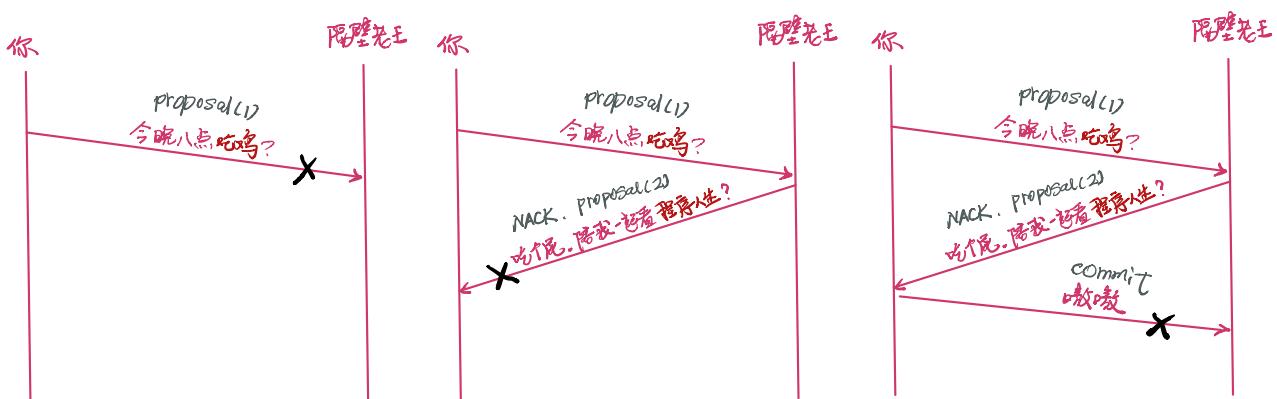
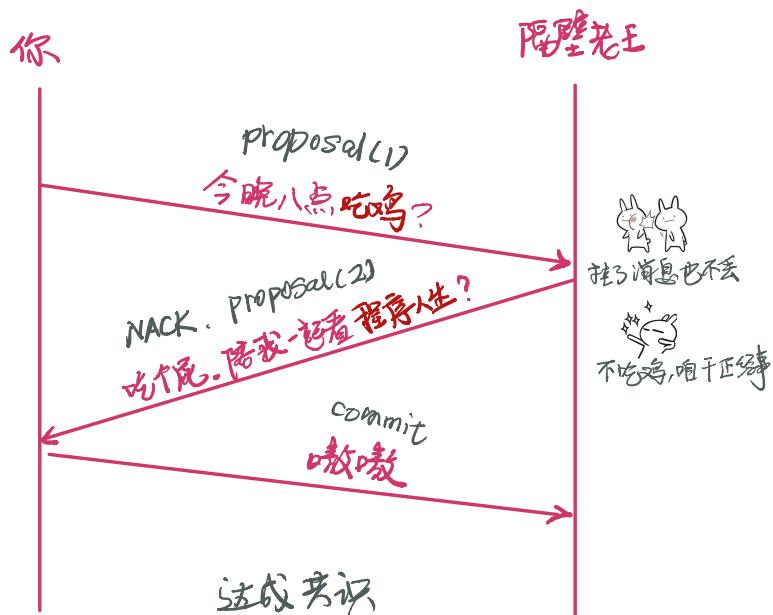
$$B \rightarrow G \quad (2, 0, 0) \prec (2, 2, 1)$$

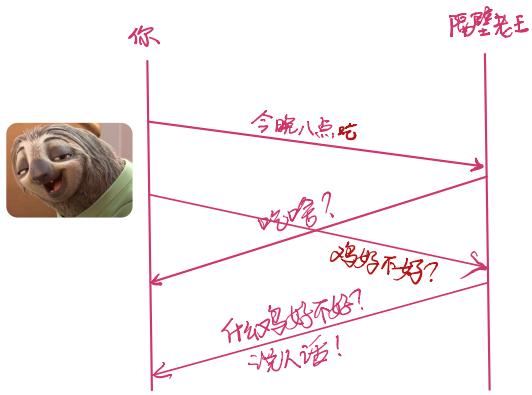
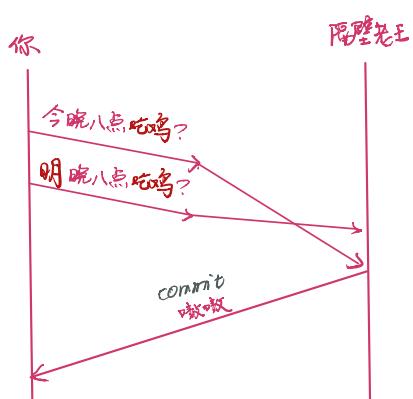
Concurrent

$$C \& G \quad (3, 0, 0) \parallel (2, 2, 1)$$

$$I \& C \quad (0, 0, 1) \parallel (3, 0, 0)$$

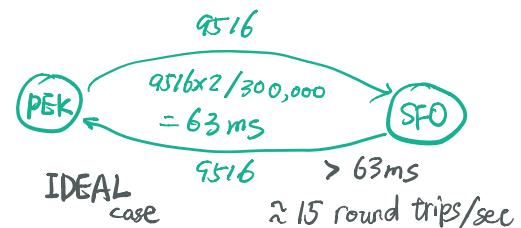






Prejudice

- Network is reliable
 - Can we rely on network to deliver the message?
- Network is homogeneous
 - wired, wireless, cellular, satellite
 - speed varies a lot
- Latency is neglectable
 - $L = \text{distance} / C$



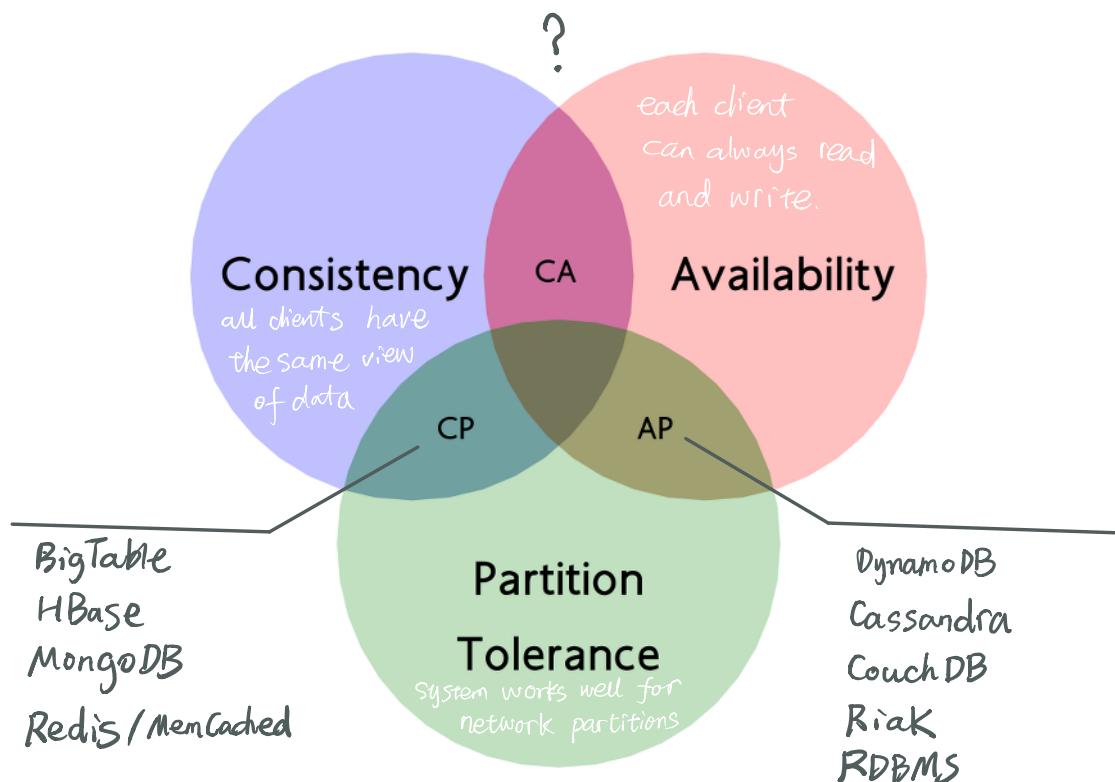
Latency Numbers Every Programmer Should Know

<ul style="list-style-type: none"> ■ 1 ns ■ L1 cache reference: 0.5 ns ■ Branch mispredict: 5 ns ■ L2 cache reference: 7 ns ■ Mutex lock/unlock: 25 ns ■ = ■ 100 ns 	<ul style="list-style-type: none"> ■ Main memory reference: 100 ns ■ = ■ 1 μs ■ Compress 1 KB with Zippy: 3 μs ■ = ■ 10 μs 	<ul style="list-style-type: none"> ■ Send 1 KB over 1 Gbps network: 10 μs ■ SSD random read (1Gb/s SSD): 150 μs ■ Read 1 MB sequentially from memory: 250 μs ■ Round trip in same datacenter: 500 μs ■ = ■ 1 ms 	<ul style="list-style-type: none"> ■ Read 1 MB sequentially from SSD: 1 ms ■ Disk seek: 10 ms ■ Read 1 MB sequentially from disk: 20 ms ■ = ■ 100 ms ■ Packet roundtrip CA to Netherlands: 150 ms
--	--	---	---

Source: <https://gist.github.com/2841832>

- Bandwidth is not a problem
 - upload speed is pretty limited for home users
- Network is secure
- Topology won't change
 - Node could come and go

CAP theorem

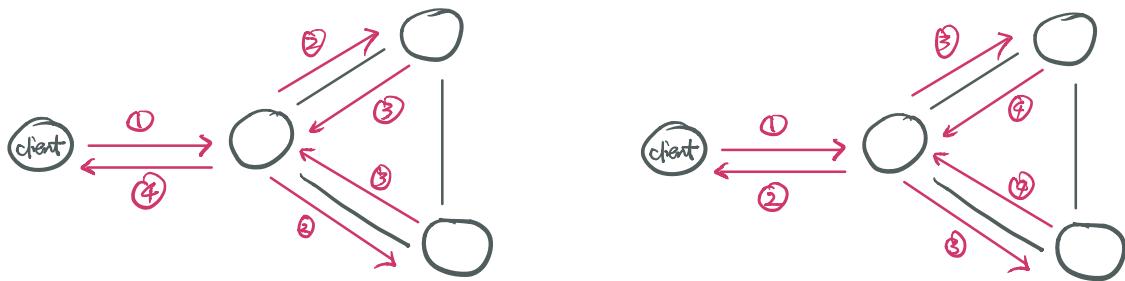


Consistency Models

- Linearizability
- Sequential consistency
- Causal consistency
- Eventual strong consistency
- Eventual consistency
- Monotonic read consistency
- Monotonic write consistency
- Read-your-writes consistency
- Writes follows reads consistency
- Serializability

Partition and Replication

- Split data to multiple nodes
- Copy data over to ensure consistency
 - sync : consistency over performance
 - async : performance over consistency



Algorithms

- Gossip
- Consistent hashing