# Distributed databases building blocks

OREN EINI

WIZARD

HIBERNATING RHINOS
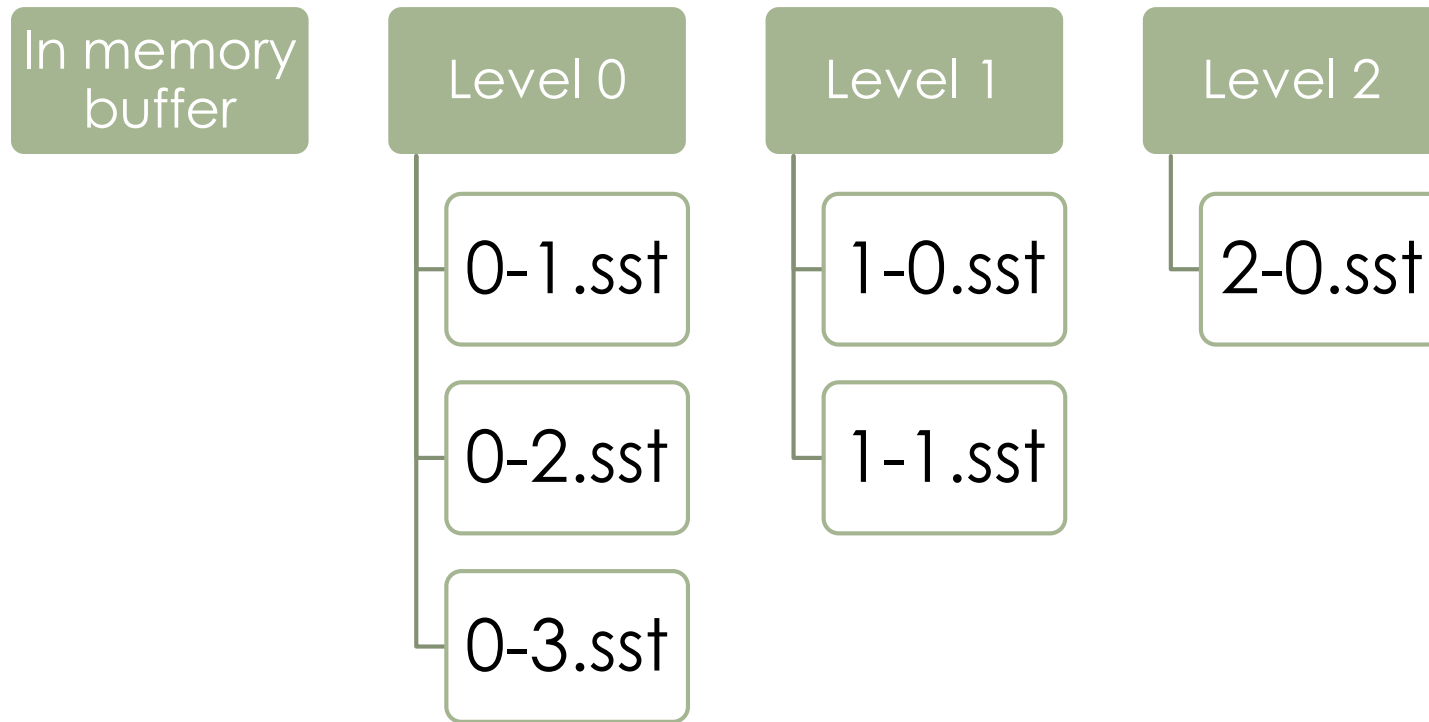
# I write databases for a living

- Rhino DHT
- RavenDB
  - Voron

- This is hard ☺
  - Also fun
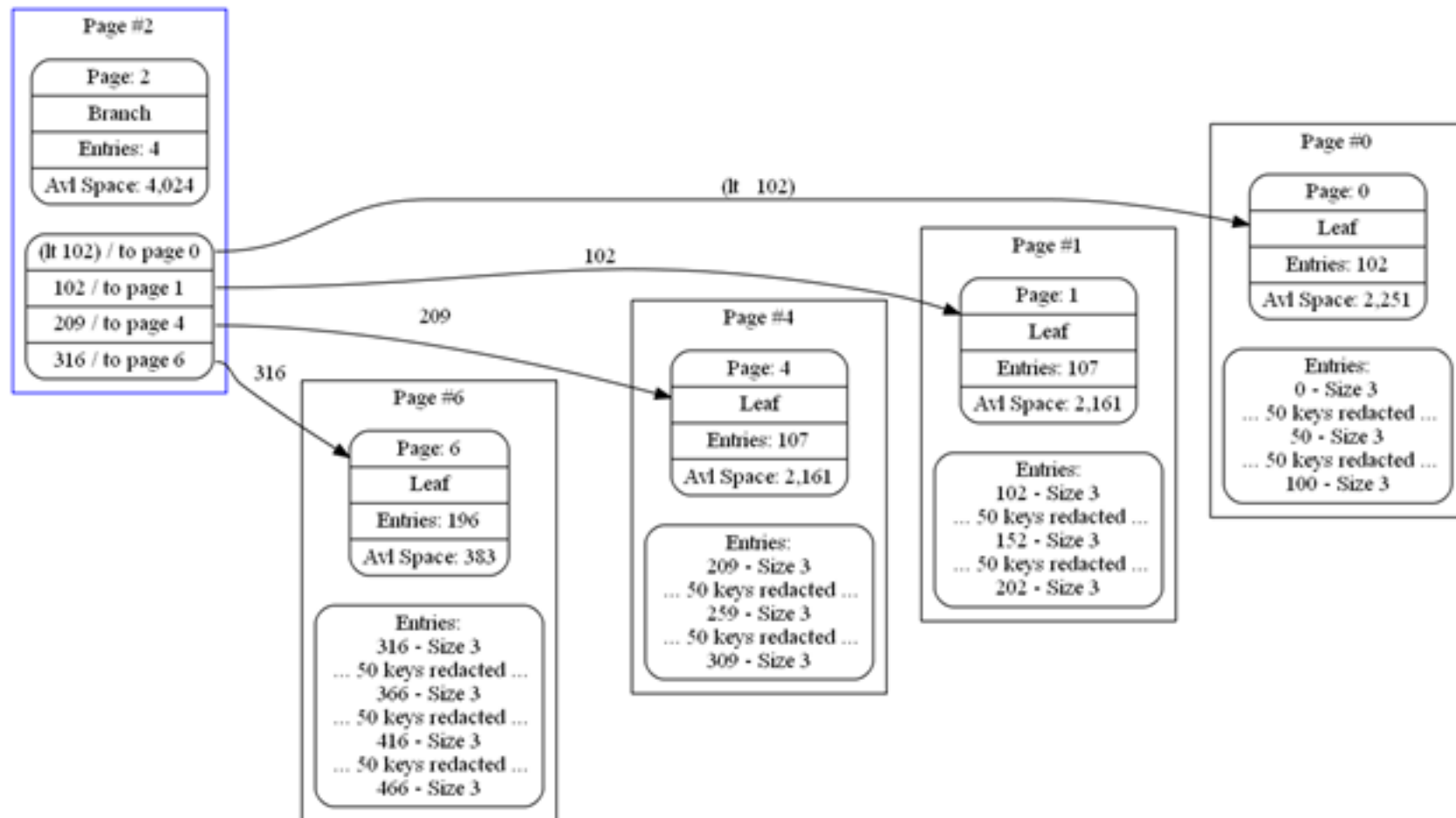
# Storing data

- What do we need?

- File.Open ?

- The disk goes round & round, round & round, round & round

# Log Structure Merge

| In memory buffer | Level 0 | Level 1 | Level 2 |
|---|---|---|---|
| | 0-1.sst | 1-0.sst | 2-0.sst |
| | 0-2.sst | 1-1.sst | |
| | 0-3.sst | | |

# B+Trees

# Concurrency & Isolations

**Tx 1**
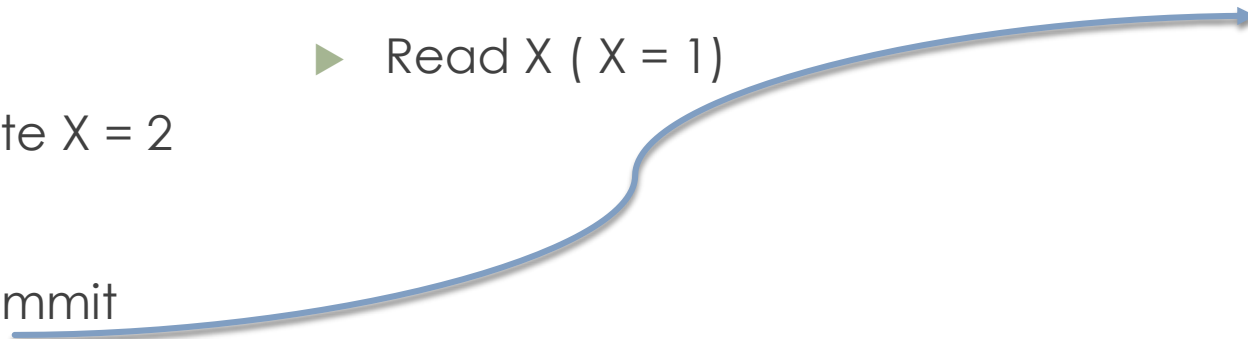
▶ Begin Tx

▶ Read X (X = 1)

▶ Write X = 2

▶ Commit

**Tx 2**

▶ Begin Tx

▶ Read X ( X = 1)

▶ Read X ( X = ?)

**Tx 3**
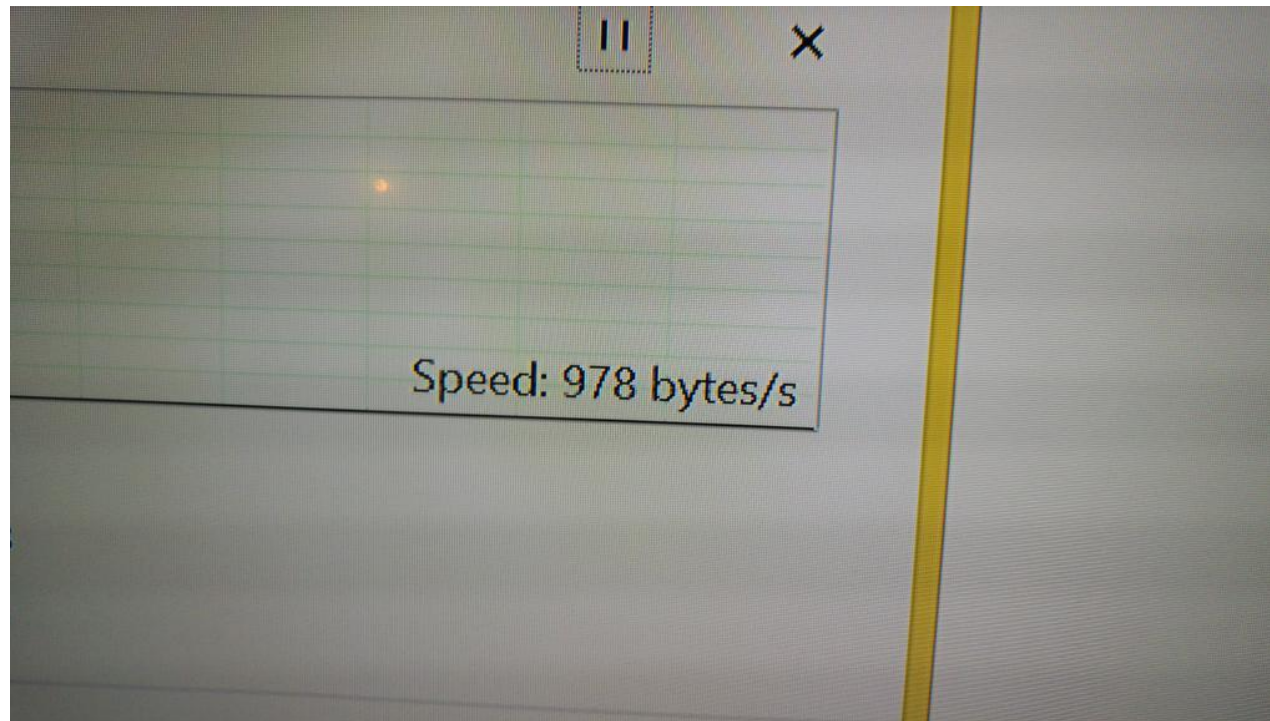
▶ Begin Tx
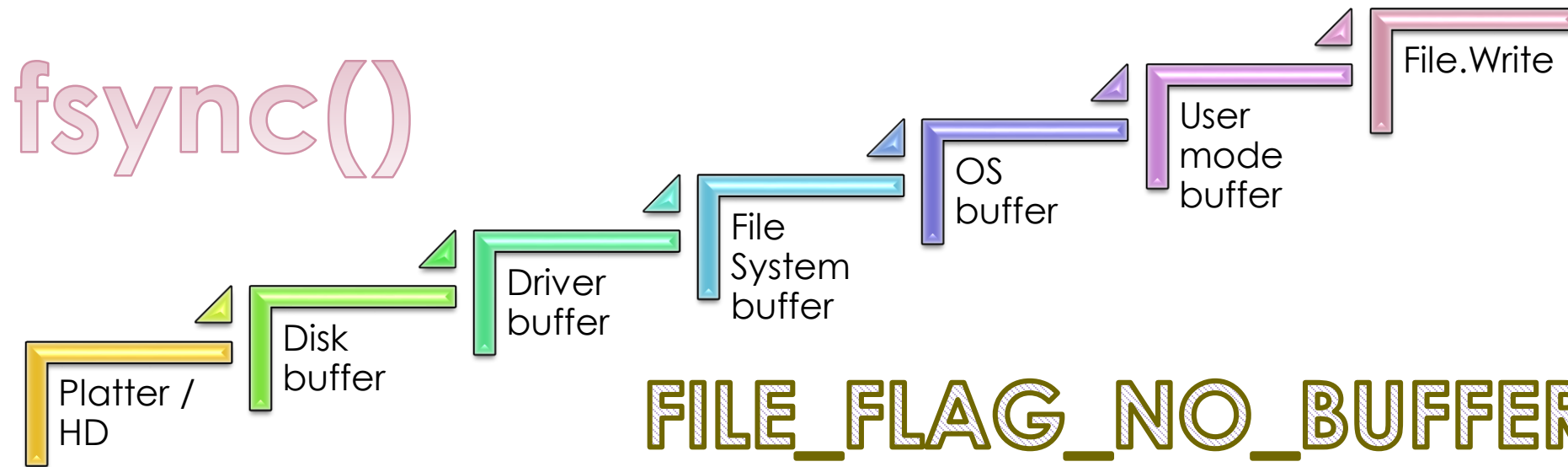
▶ Read X ( X = 2)

# Implementing concurrency

Locks

MVCC

# I/O IS SLOW

# Durability / transactions / fsync

fsync()

File.Write

User
mode
buffer

OS
buffer

File
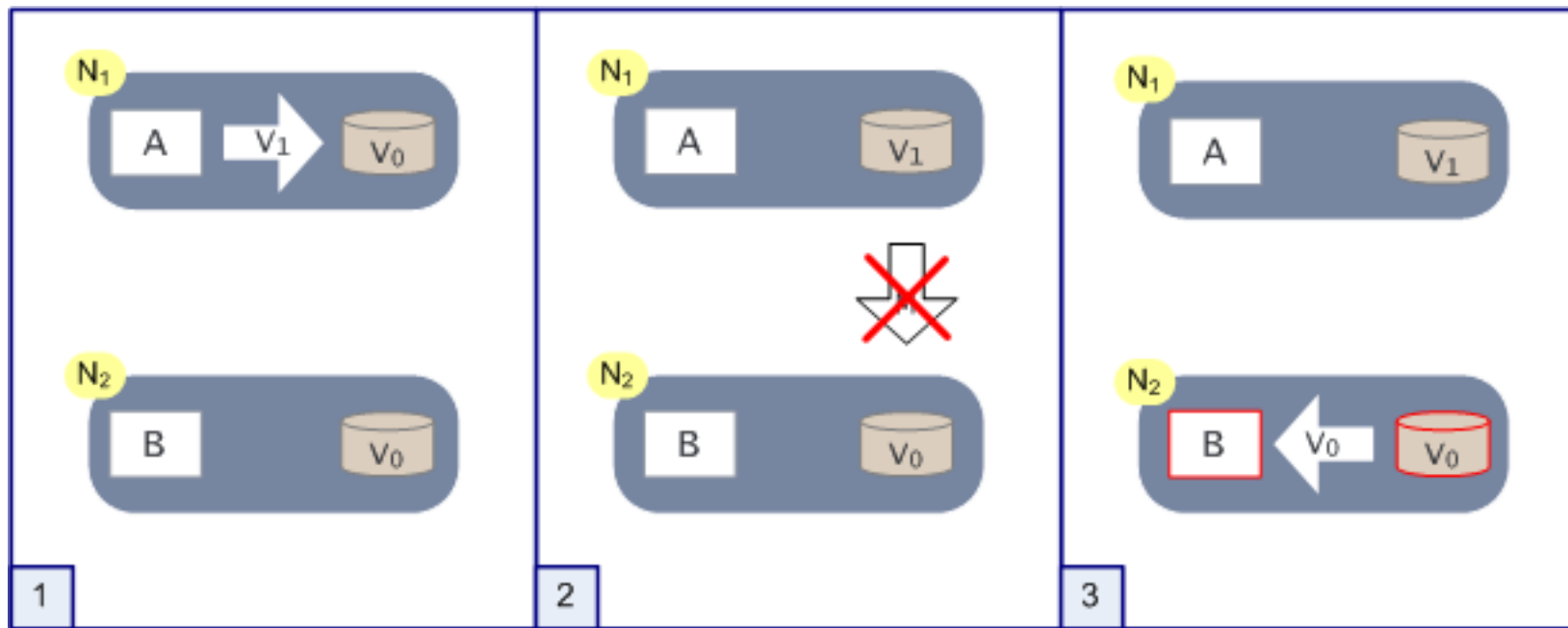System
buffer

Driver
buffer

Disk
buffer

Platter /
HD

O_SYNC

FILE_FLAG_NO_BUFFERING |
FILE_FLAG_WRITE_THROUGH

# The curse of the single node…

# What about the distributed part?

# Before that…

- What is the distribution model?
  - Consensus?
  - Collaborative?
  - Repair?
- Partition model

# Consensus

Paxos ???

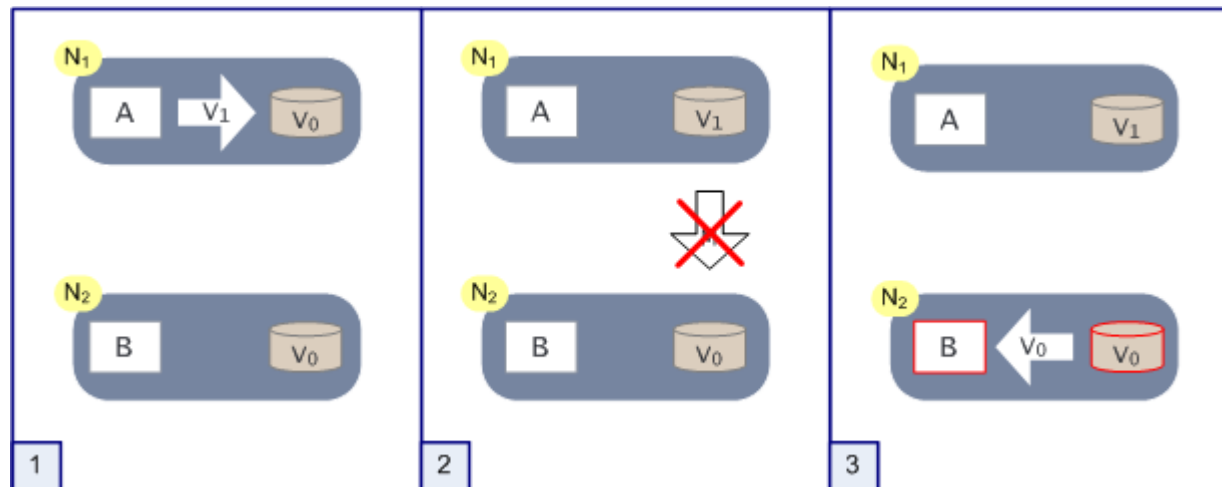# Raft

# Consensus – log of operations

- Set x = 1
- Set y = 2
- Set n = 1
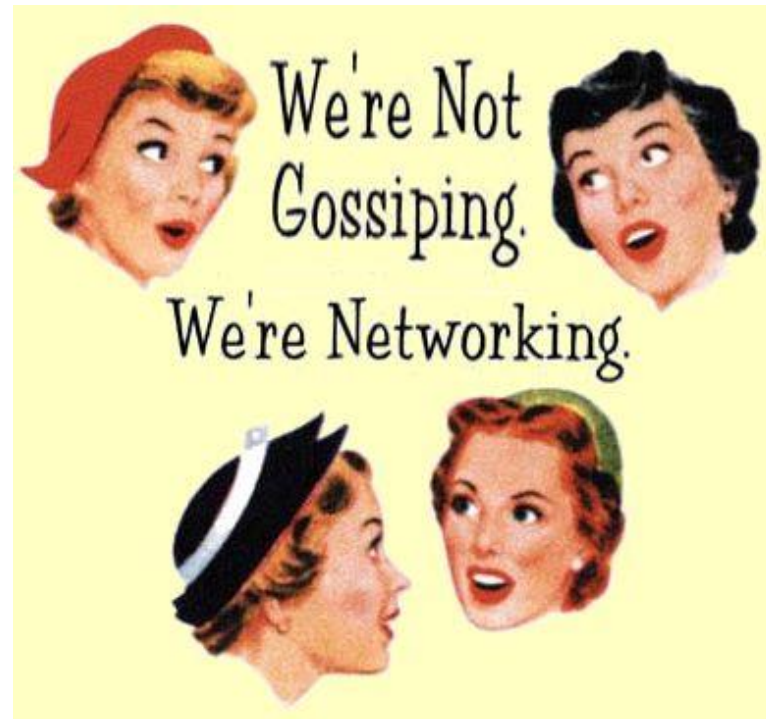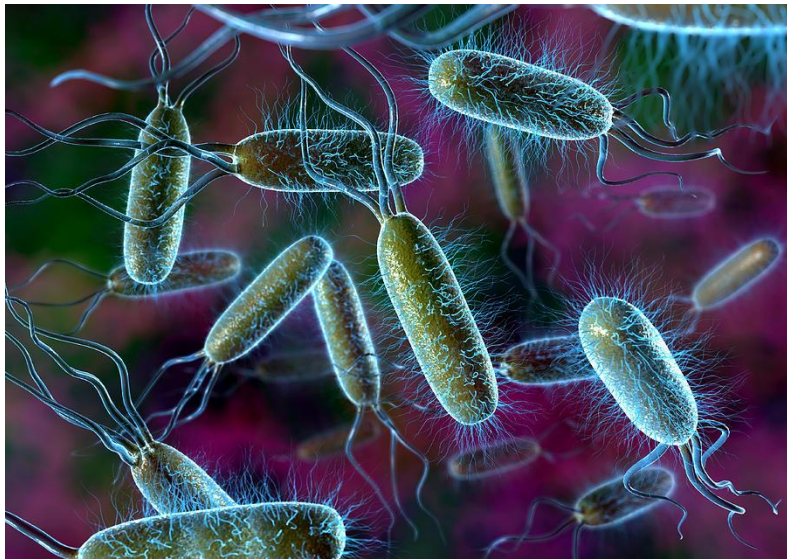- Set x = y + n

# Problems?

# Collaborations (master / master)

- Allow writes on any node

- Conflict-free replicated data types ?

- Merges

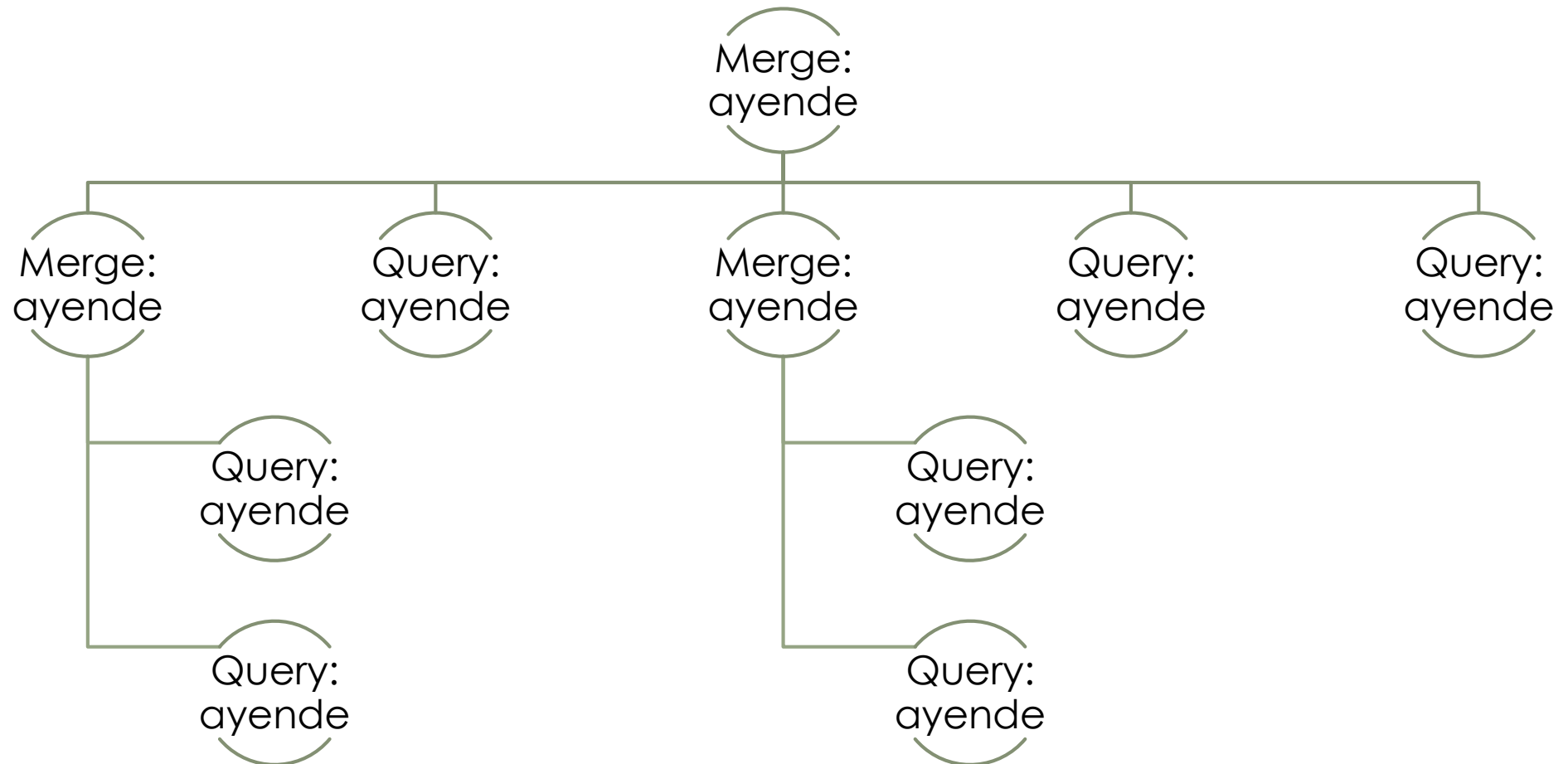# Sharded

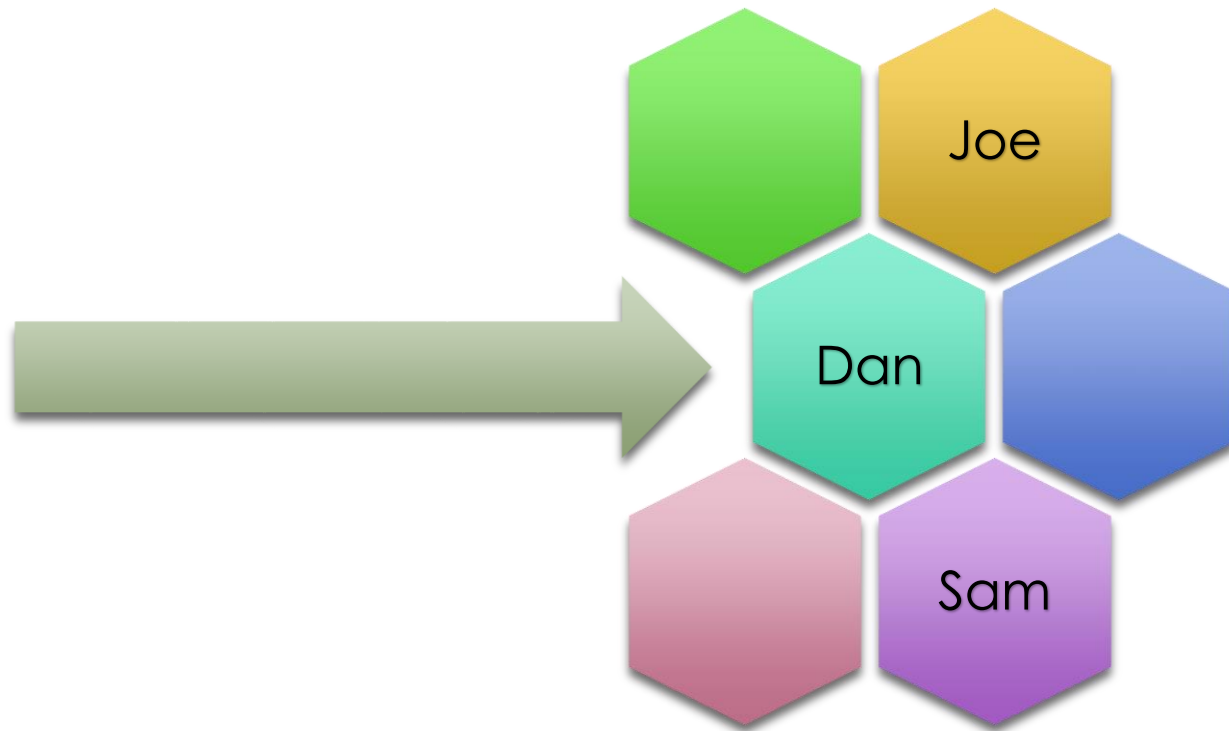- Some data on some node
- Multiple consensus groups

# Gossip




We're Not Gossiping. We're Networking.

# Data model

- So now you know how to store data (single node, multiple nodes)
- What do you *do* with this?

# Entire dataset

# Portioned data

# Imagine a banking system…

- Core data model:
  - Accounts
  - Account operations (incoming, outgoing, fees, etc)


- New rule changes – gossip
- Reporting

# Questions?