

CSE 452

Distributed Systems

Partitioning  
aka Sharding  
aka Splitting the data

# Why build a distributed system?

- fault tolerance — replication
  - redundant computation + data
- performance
  - "edge computing" — move things close to users
  - split up the work / parallelism
  - "horizontal scaling"  
add more machines → get better performance

# KV Store

- distributed
- fault tolerant
- linearizable
- partitioning / sharding / splitting the data
- cross-partition transactions

# Partitioning

- split the data

- some keys in one partition  
some keys in another

get(k)  
put(k, v)  
append(k, v)

k <sub>1</sub>	k <sub>2</sub>	k <sub>3</sub>	k <sub>4</sub>	k <sub>5</sub>	k <sub>6</sub>
↓	↓	↓	↓	↓	↓
v <sub>1</sub>	v <sub>2</sub>	v <sub>3</sub>	v	v	v

# Static Partitioning

- ahead of time, decide which keys go to which partitions
- Cons
  - Stuck with choice of Key  $\rightarrow$  partition map
  - Changing workloads
  - # of partitions

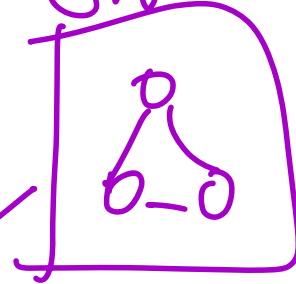
# Dynamic Partitioning

- move keys between partitions
- supports load balancing
  - lab 4: manual load balancing
- keep track of current mapping of keys  $\leftrightarrow$  partitions
- add/remove entire partitions

Admin



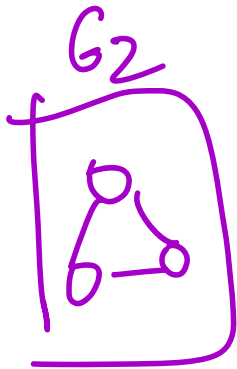
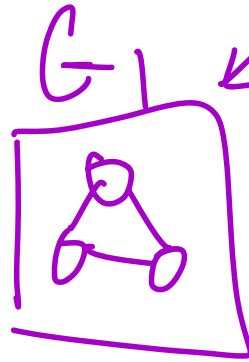
Shard Master



key  $\rightarrow$  group  
set of group

Client

you have  
key  $\parallel$



groups

# lab 4

part 1: Shard Master (not very distributed)

part 2: other groups that implement  
the sharded KV store  
- moving keys

part 3: multi key operations  
CROSS partition operations

# Shard Master

- Application: Command, Result
  - Query - reads current configuration
  - Join - add new group★
  - Leave - remove group★
  - Move - move a shard to specific group
- ★ = rebalance key space — move as few shards as possible

# Configuration: metadata

- Config num — monotonic

- Keys  $\rightarrow$  groups (Shards  $\leftrightarrow$  groups)

- set of groups

---

Shard = small set of keys

group = one paxos cluster available for KVStore ops