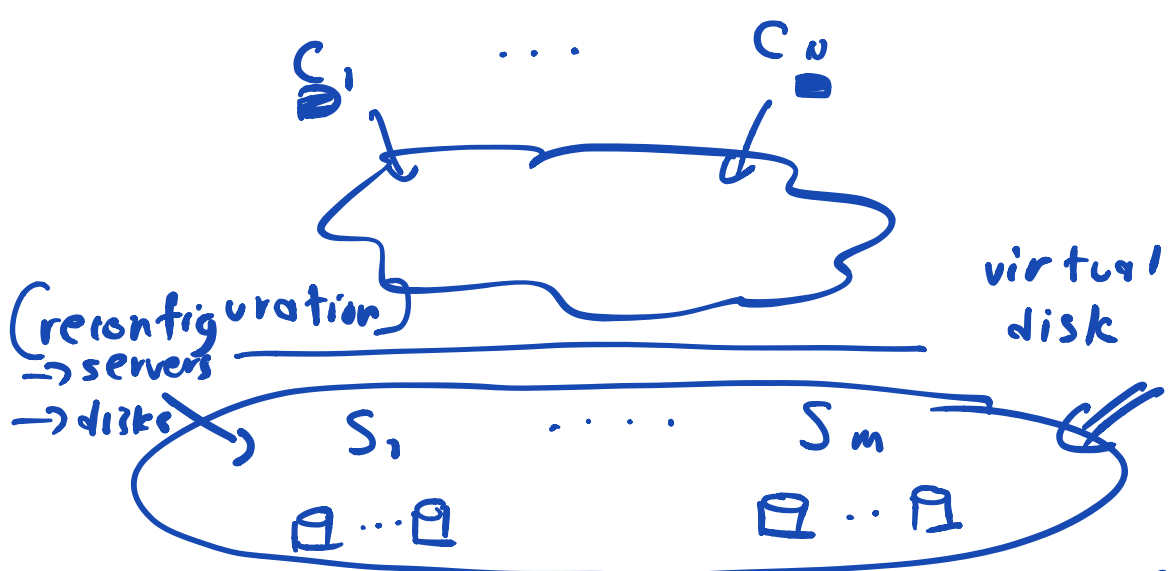


Petal / Fran gipani [Case Studies]

interface / abstraction :



foreach virtDisk: blocks (0 ... N-1)
read/write
(offset)
(delete)

[All servers cooperating to provide abstraction]

Liveness: heartbeats + majority consensus

Global state:

[membership
global mapping into] } widely replicated

=> PAXOS

=> common pattern { inner ring: paxos
outside: not }

(v → p translation)
=> < v disk, offset >

⇓ [key]

< server, disk,
local offset >

Data Access

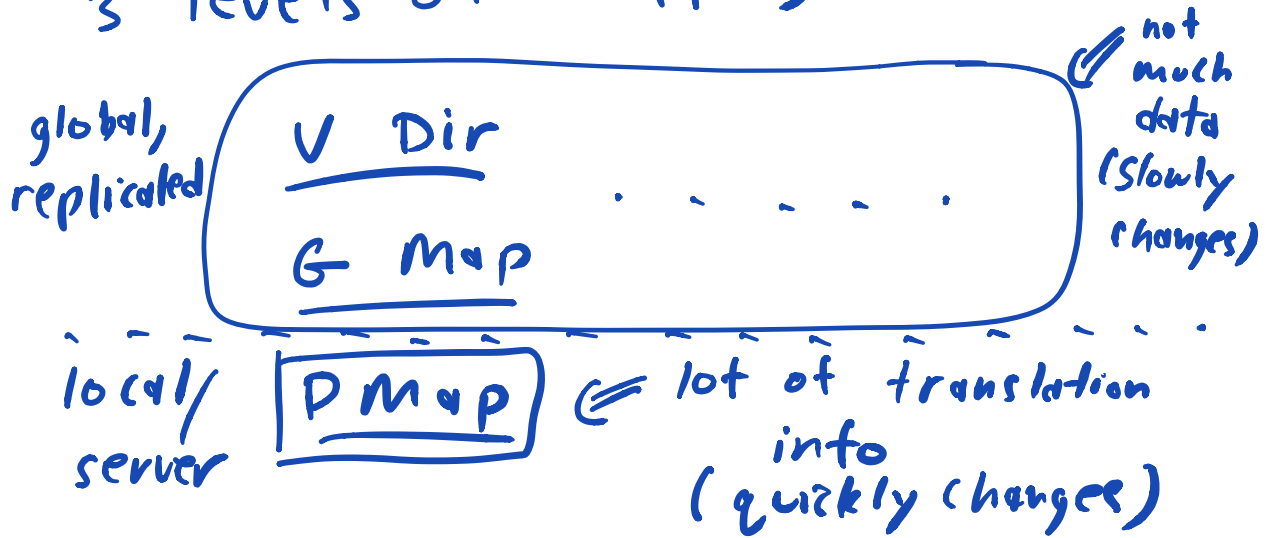
=> mirroring,
striping

Recovery

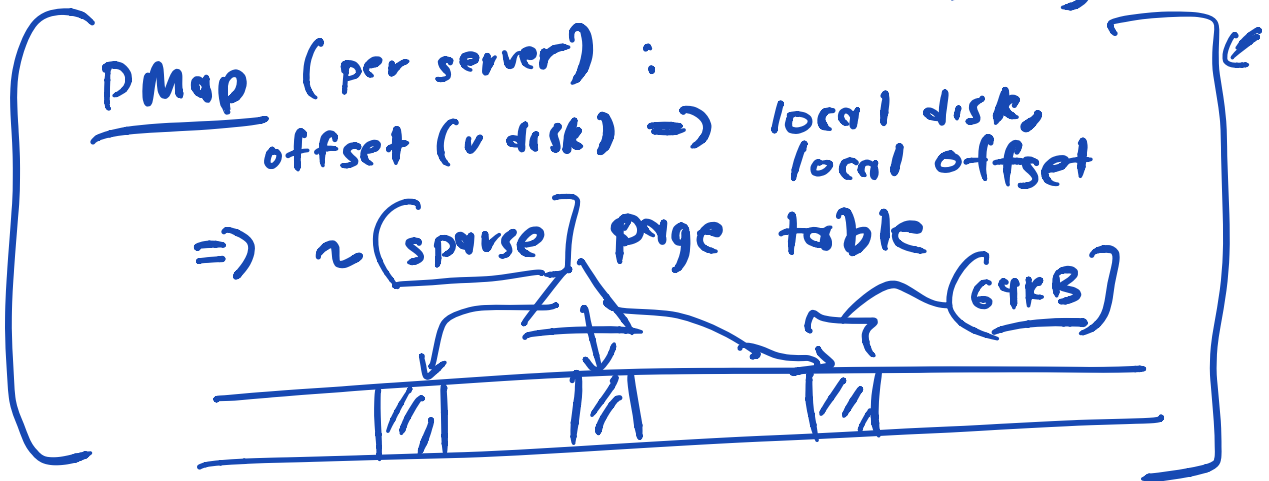
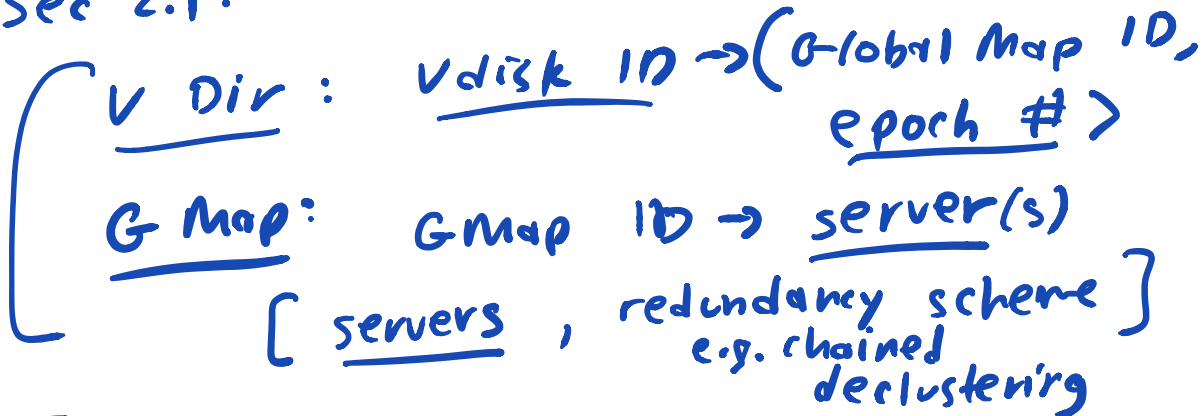
=>

$V \rightarrow P$ (virtual \rightarrow physical)

3 levels of mapping



Sec 2.1:



Backup: snapshots

(GMap ID, epoch)

GMap:
immutable

=> are not "consistent"
(from app perspective)

Incremental Reconfiguration :

1) local



actions:

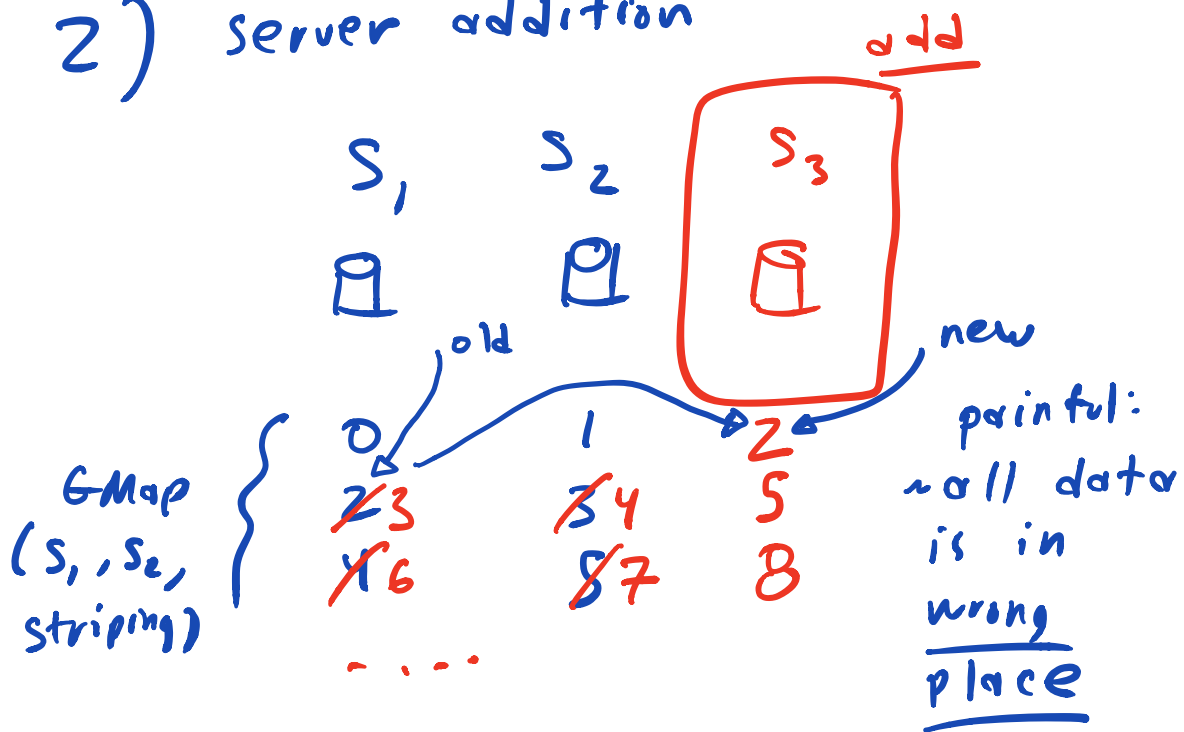
1) writes → easy

2) reads → hard

offline / background rebalance

All local : "easy" (no Paxos,
consensus)

2) server addition



Process: reconfig

create new gmap'
change vdir \rightarrow gmap'
move all data
continue normal operation

"it depends"

"Basic" algorithm:

gmapnew ↵

gmapold

writes: new

reads: use new map
but, it fails,
consult old map

(bg: moving data old map
→ new)
@ some rate

Problem? too slow

Solution: relocate parts of
map @ time

v Disk:

regions: old, new, fenced ↵

↗
use
old

↗
use
new

↗
in flight

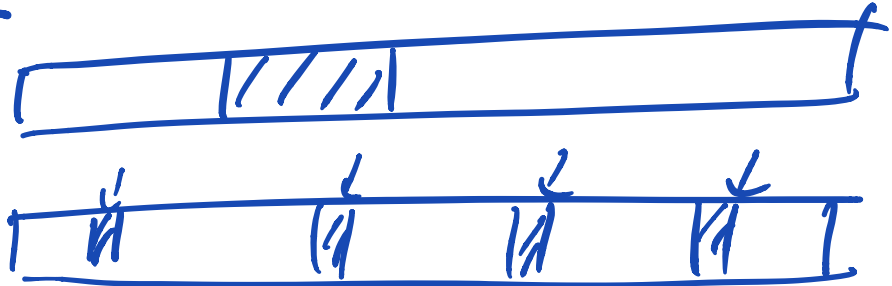
(new, old?)

use "basic"

Perf problem? Hot spots

Locality in w/L
(workload)

Solve? => guessing that



Data Access / Recovery

1) Redundancy scheme

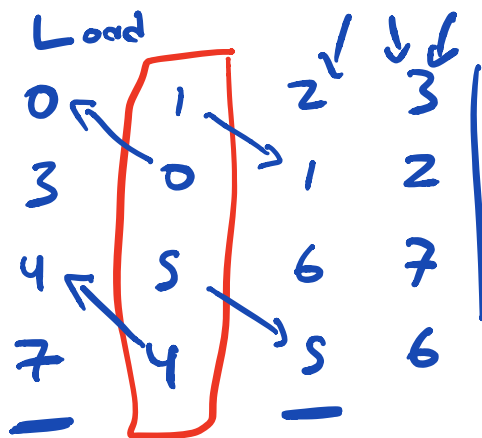
=> striping

=> chained declustering '89

mirror: failure

0	0	1	1
2	2	3	3
4	4	5	5
6	6	7	7
⋮			

load problem
balance

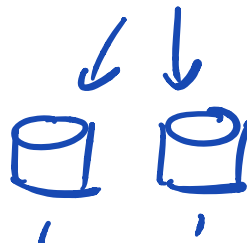


better: ✓

dynamic load balancing:
 clients track outstanding #
 of reqs / disk
 => balance

Petal : replicas

approach: logging: busy
Primary / Backup
 writes



(=> writes: slower
logging)

next optimization:
 => lazy in marking "not busy"

Petal :

virtual Disk

target for FS's, DBMS's

Limits: scale

reconfig \Rightarrow expensive