

Distributed Systems

Spring Semester 2020

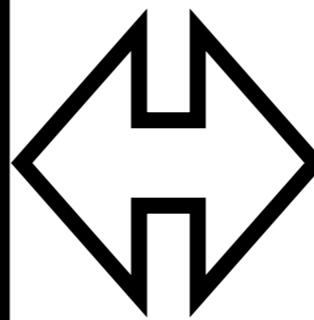
Lecture 17: Existential Consistency

John Liagouris
liagos@bu.edu

Why this paper

- Practical storage for giant read-heavy web sites — YET AGAIN ;-)
- Techniques for detecting anomalies
- Practical consistency concerns

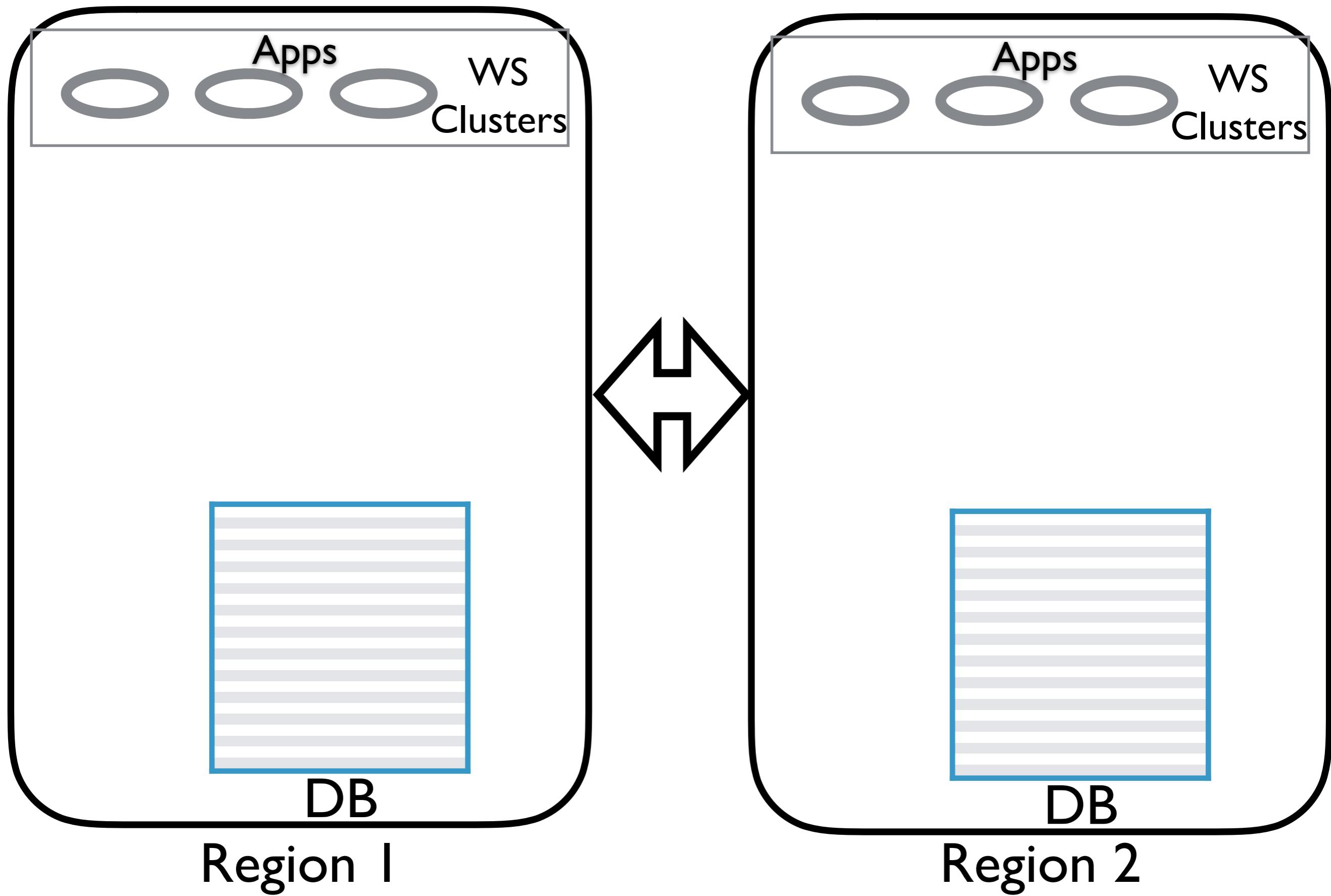
FB Infrastructure



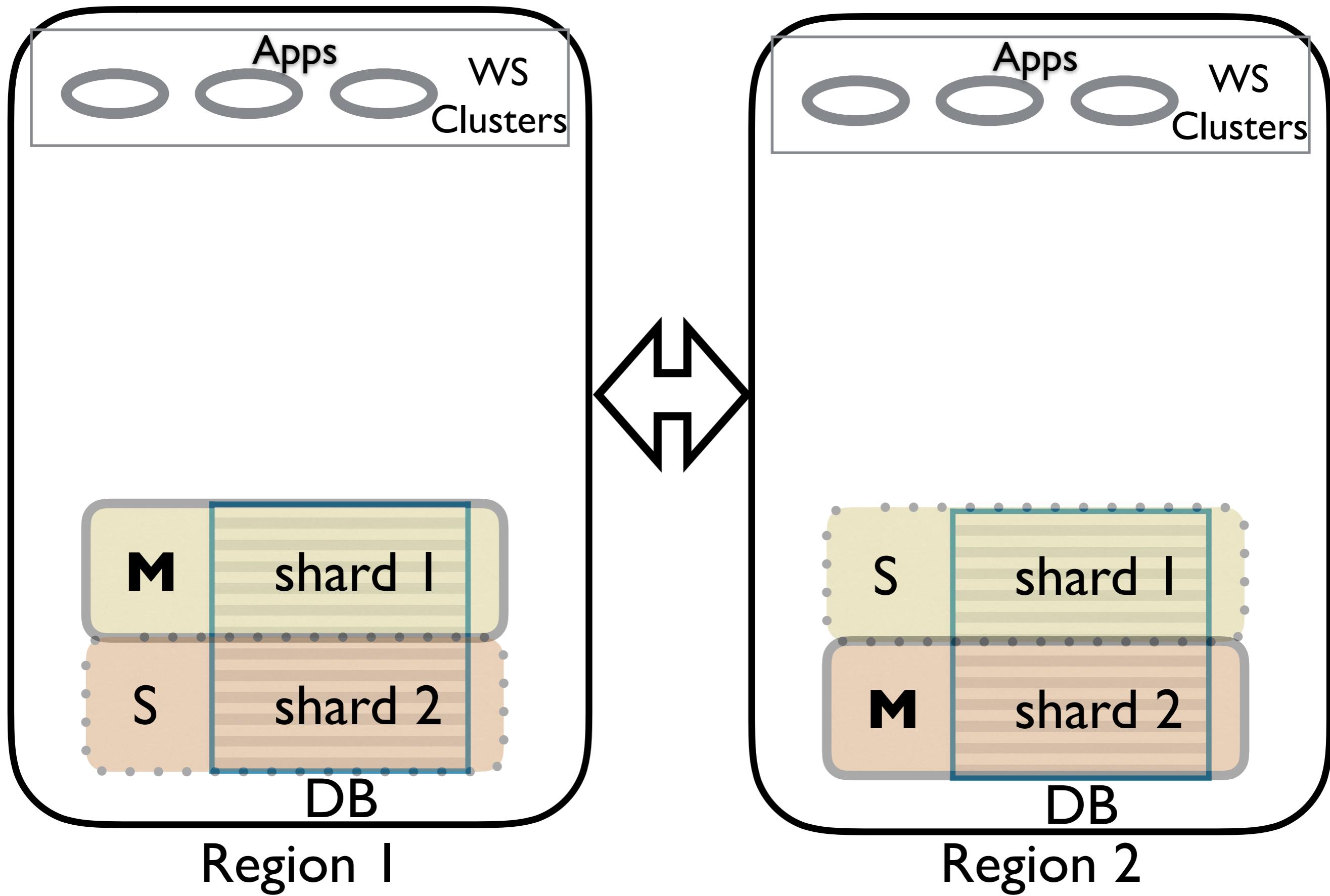
Region 1

Region 2

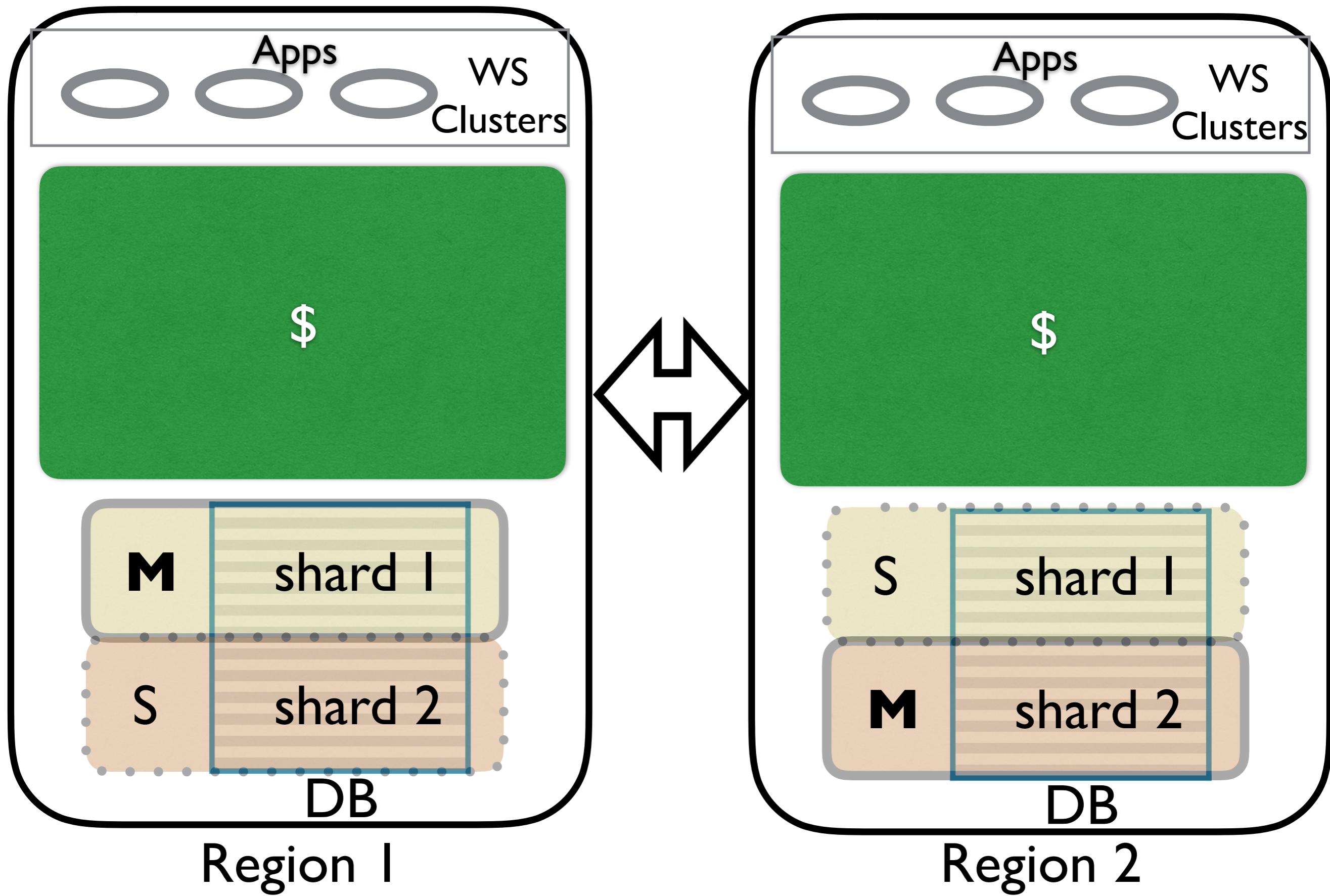
FB Infrastructure



FB Infrastructure



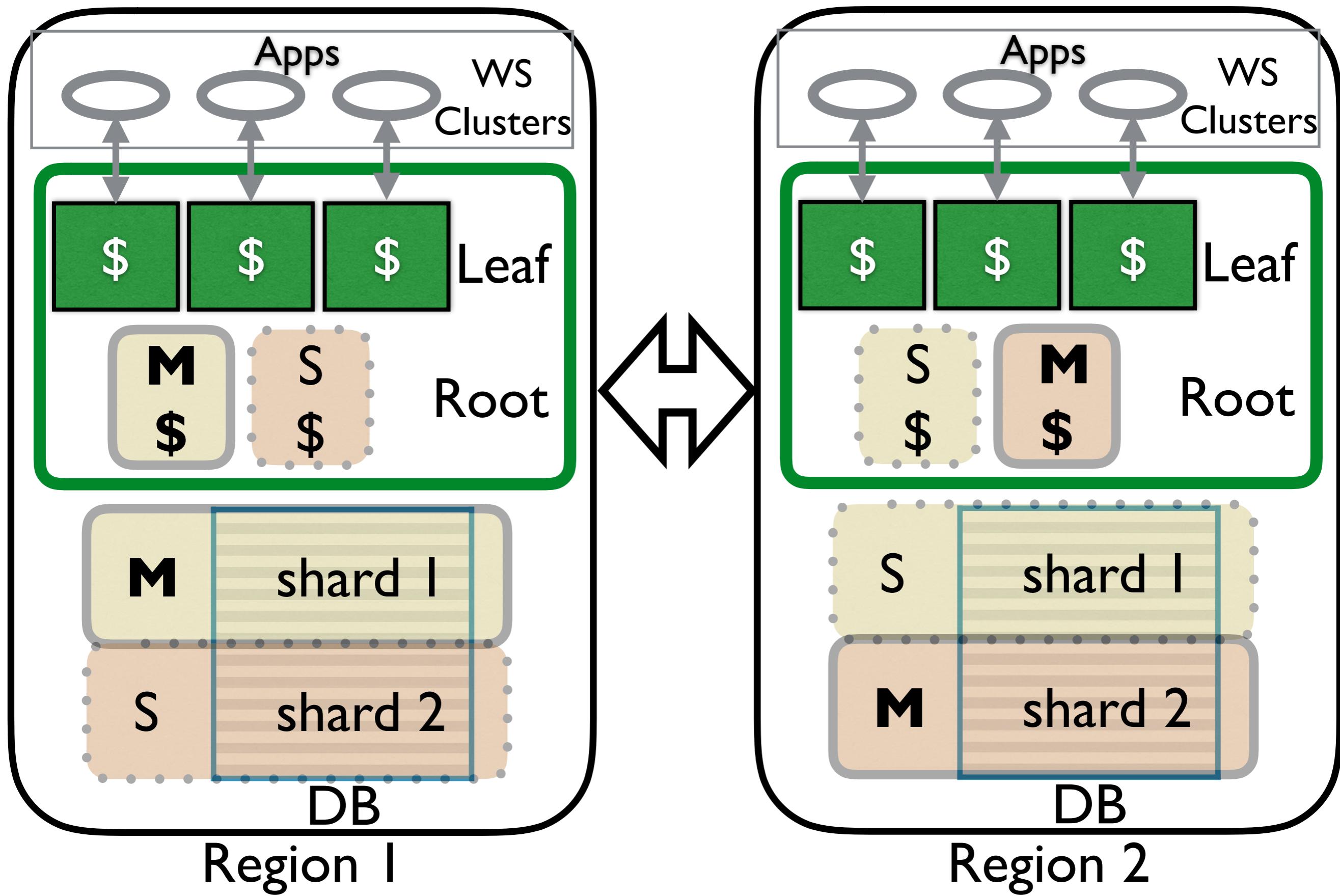
FB Infrastructure



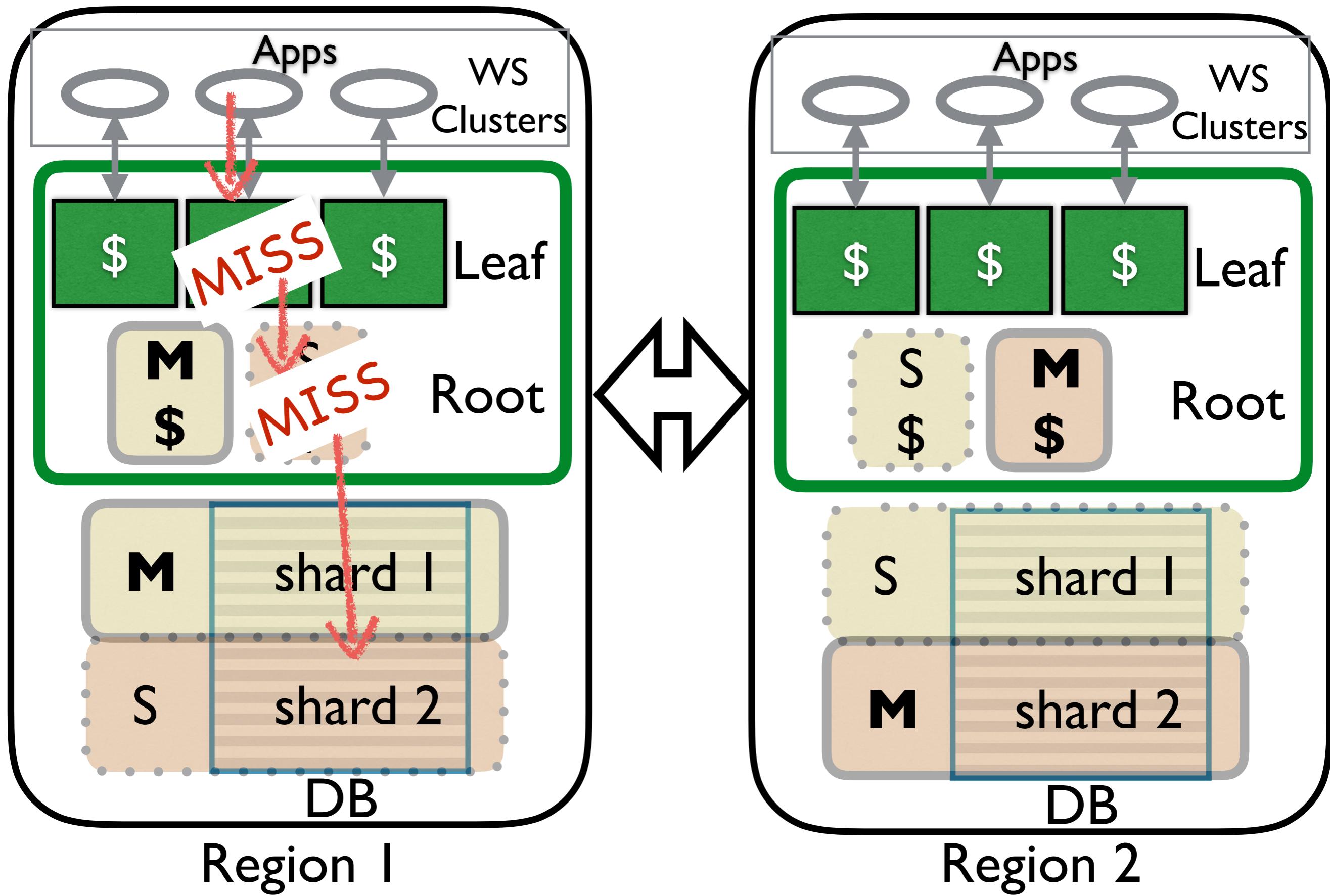
Why this overall arrangement?

- why multiple geographic sites (regions)?
- why complete replica of DB for each region?
- why DB?
- why shared DB servers?

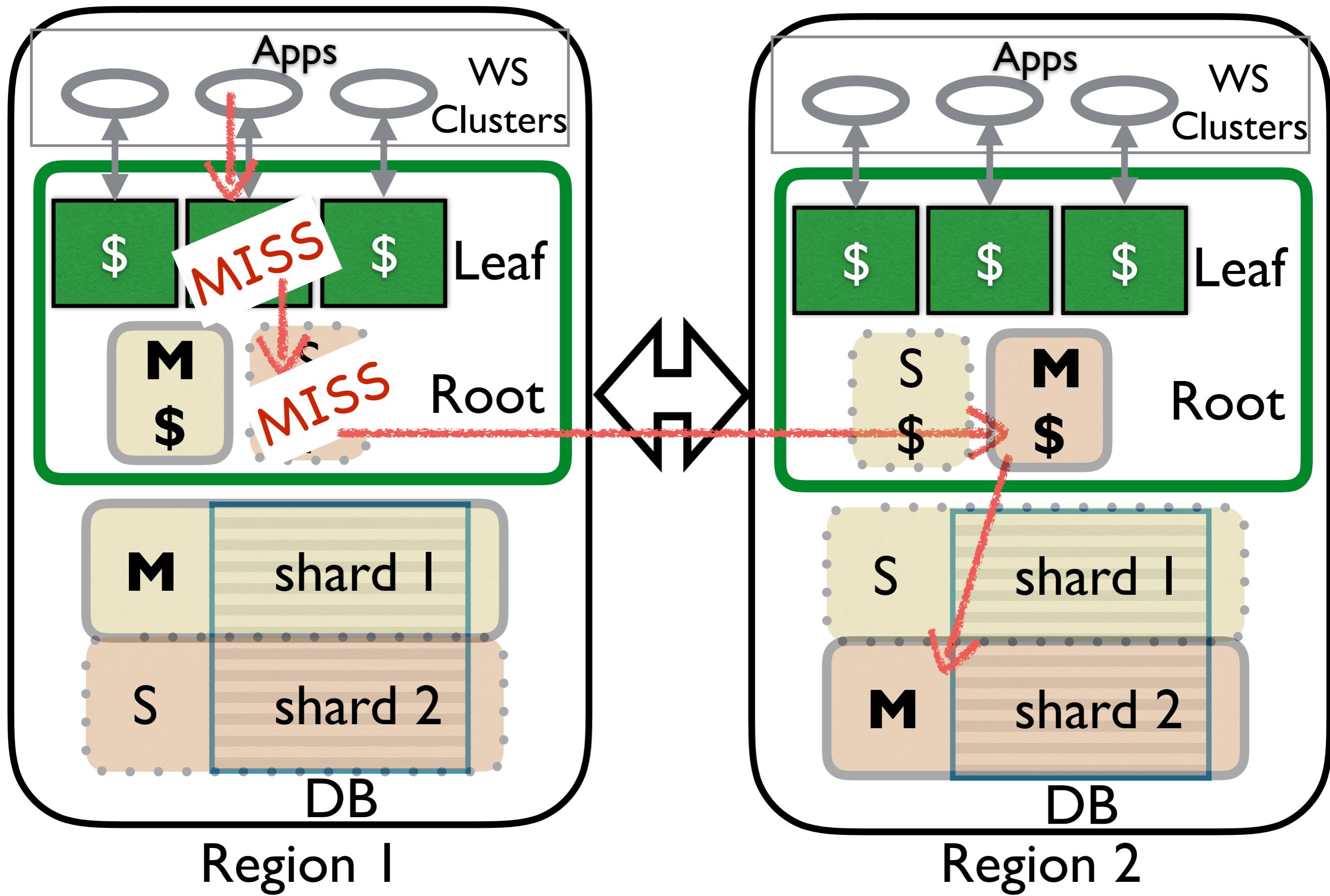
FB TAO



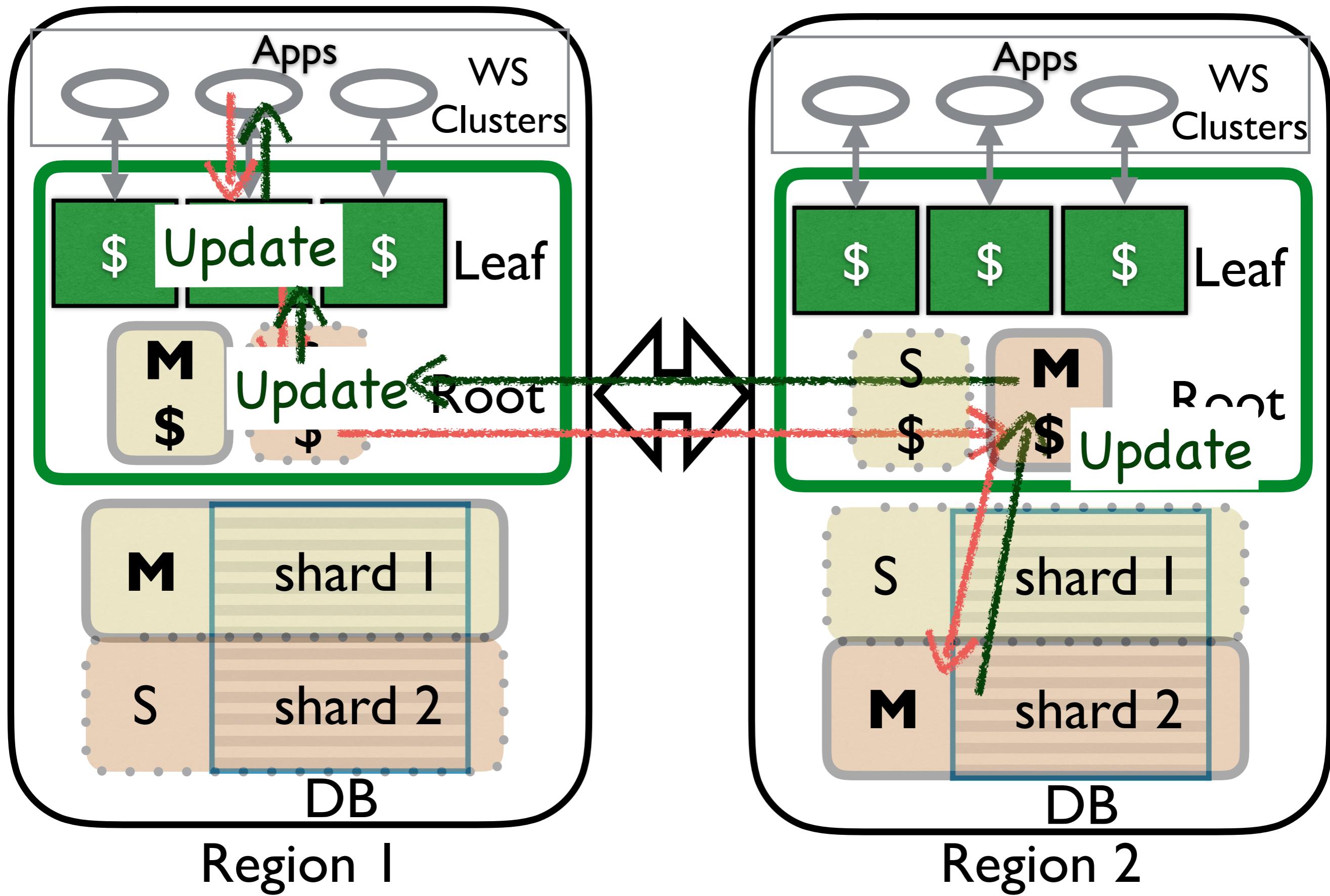
READS



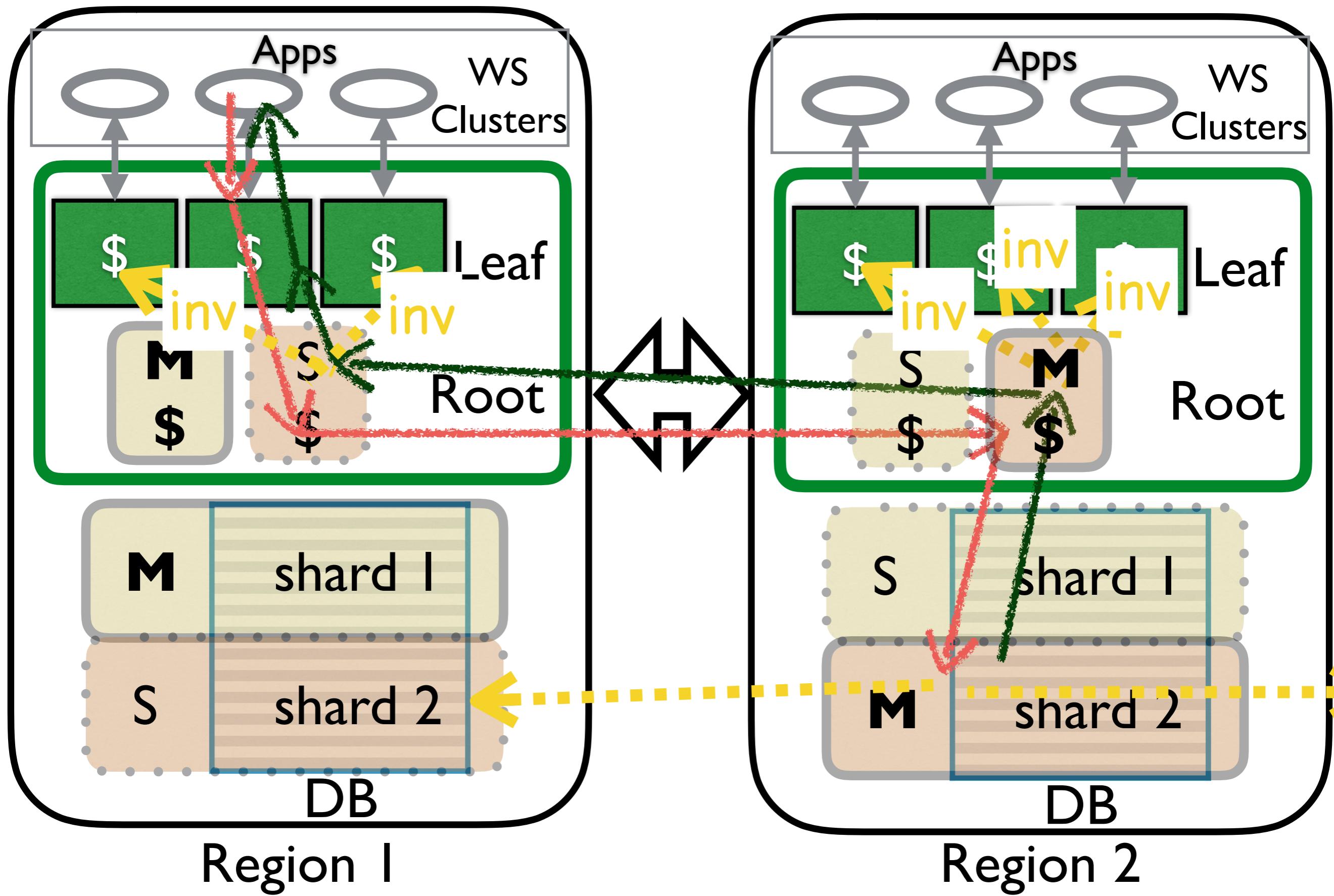
Writes



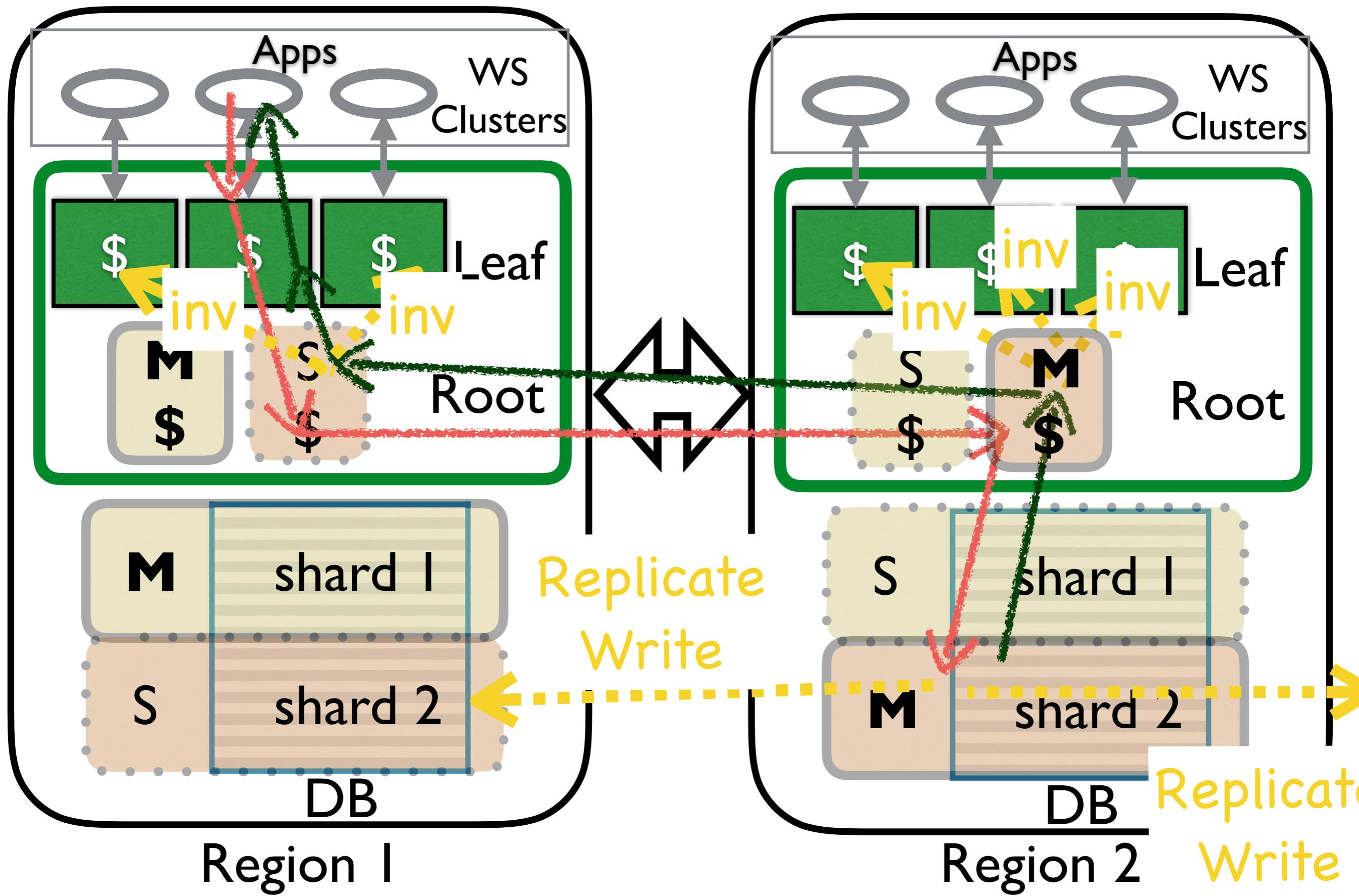
Writes



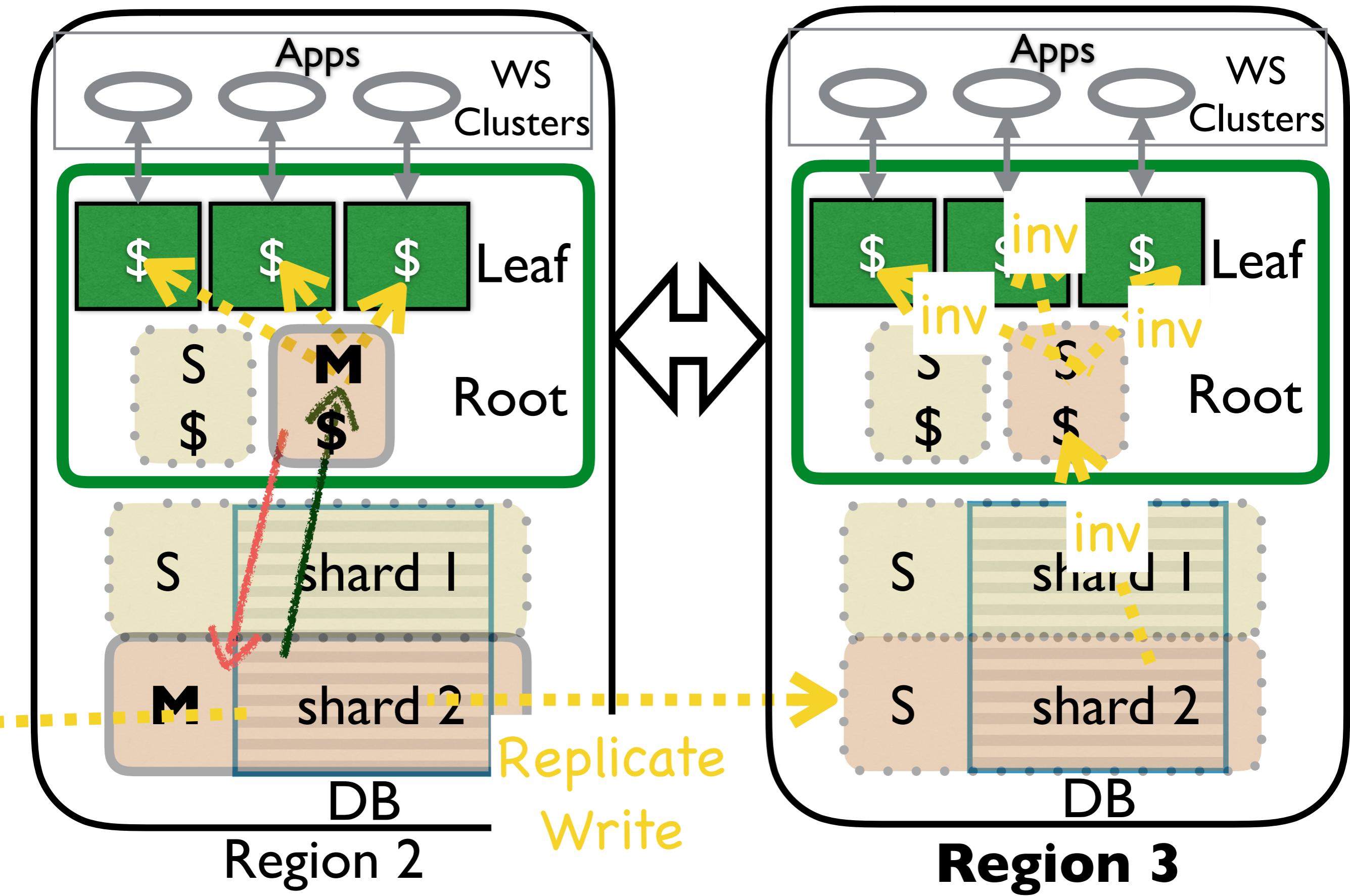
Writes



Writes



Writes



TAO Consistency Model

- 2.2 : “per-object sequential consistency and read-after-write consistency within a cache [and thus shard] and eventual consistency across caches”
- per-object sequential consistency: like PNUTS timeline consistency — cache’s may lag but never go backwards — master coordinates and orders all updates — eventual consistency never reorders — users never see data go backwards
- Read-after-write consistency: see your own write’s — important for FB user experience

What is "within a cache"?

Guarantees only hold if
app uses a single cluster
if cluster fails and app
switches, no guarantees

Weak Consistency Model

- Allows stale reads from caches (until invalidate arrives) — what would a linearizable design look like?
- Single copy of each key eg. writes+reads all go to master DB — like PNUTS read-latest
- caching, but write must **wait** for all invalidates to complete?
- What would the cost be? — Either much slower reads, or much slower writes

- KEY OBSERVATION:
 - Linearizability costs you on every operation — But may only make a difference in smallish window after each write

Why Consistency Checker?

- WANT — Linearizability but don't want to pay for it ;-)
 - So want to know how TAO makes out in practice
 - So check against linearizability as well as their actual guarantees
- They want to systematically detect the kind of bugs they had with look-aside caching