

# Policy Based Storage System in Heterogeneous Environment

Dai Qin, Ashvin Goel, Angela Demke Brown

University of Toronto

{mike,ashvin}@eecg.toronto.edu demke@cs.toronto.edu

## Motivation

Storage environments are highly heterogeneous

### Workloads and Applications

- Databases
- File systems
- NFS servers
- Scientific clusters

### Storage and Hardware

- SSDs (PCIe/SATA)
- Disks (SAS/SATA)
- Remote storage
- Tapes

### Management Challenges

#### Current solutions are:

Applications manage underlying hardware

Examples:

File systems: ZFS/btrfs

- ZFS as volume manager
- Integrated: metadata replication, RAID-Z,...

NFS Servers: NetApp Appliance

- Even integrate special hardware

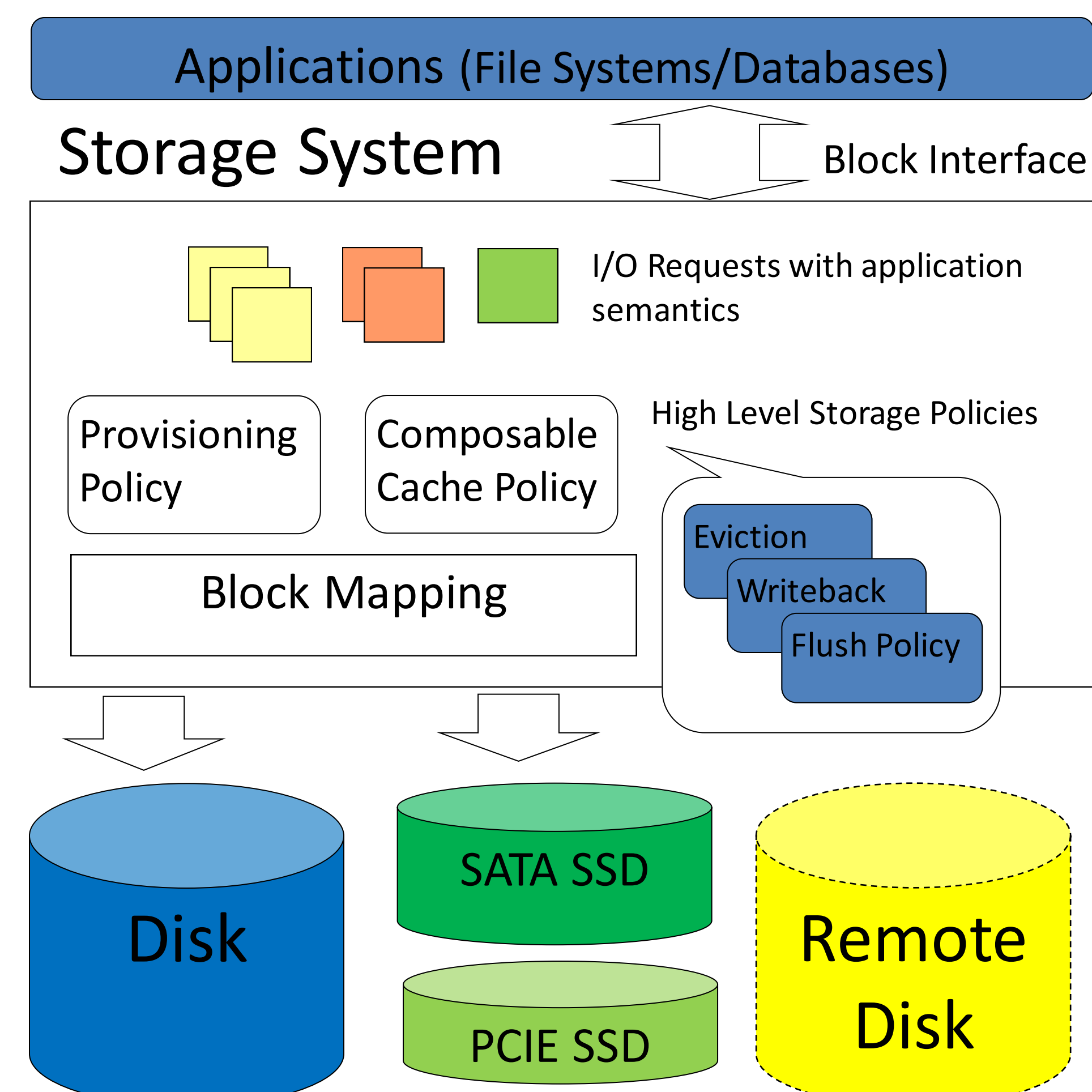
#### Problems with current solutions: What if ...?

- Management features are not built in?
- Storage virtualization/deeply stacked storage?
- Users run other applications?

Monolithic solutions are inflexible

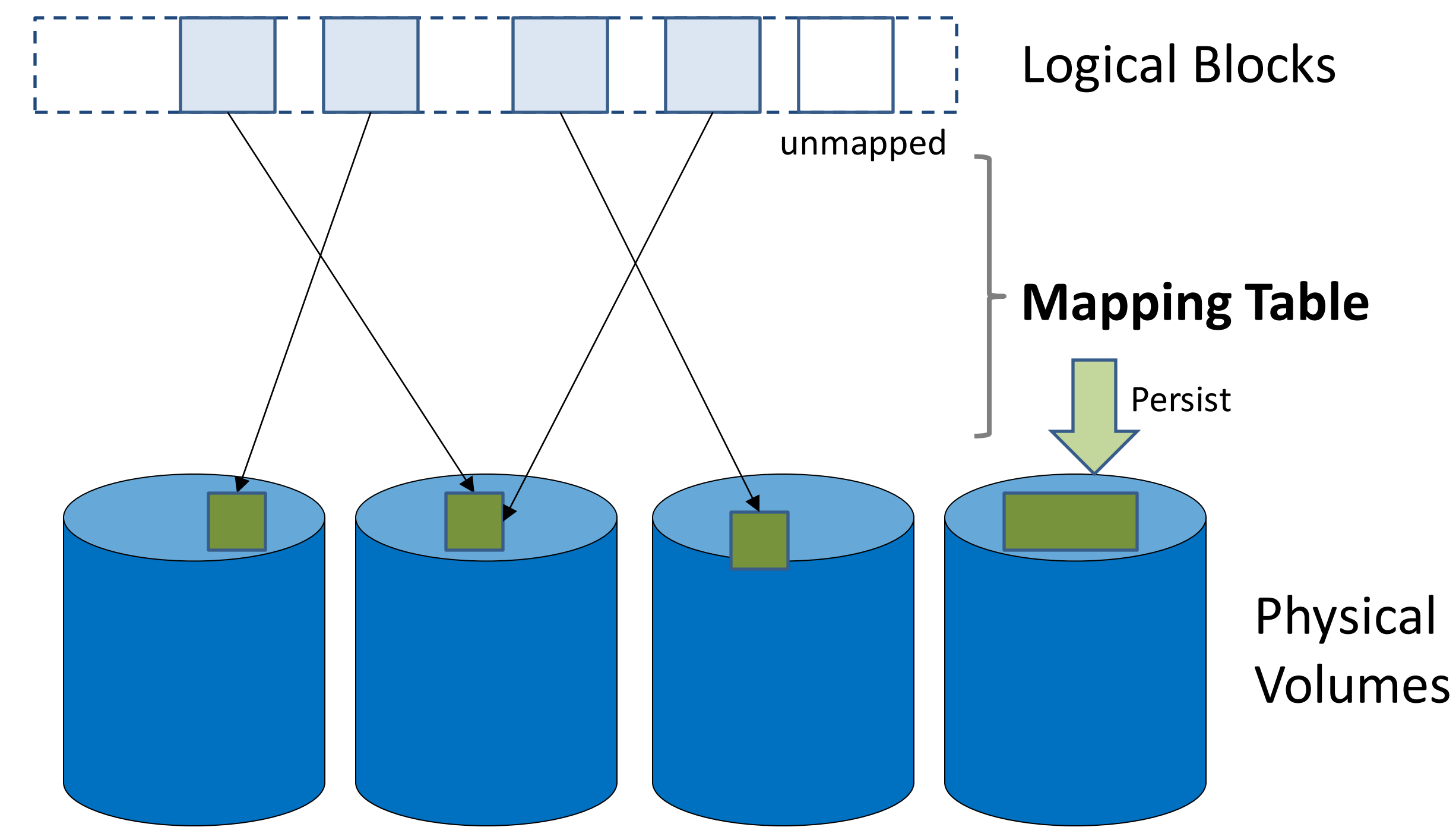
## Policy Based Management

- Storage hardware is associated with **properties** (e.g. performance, reliability)
- Storage system is aware of application semantics
  - Metadata vs data
  - Allocated vs unallocated
- Flexible policies map application requests to hardware capability
- Composable policies make it easier to create new solutions.



## Challenges

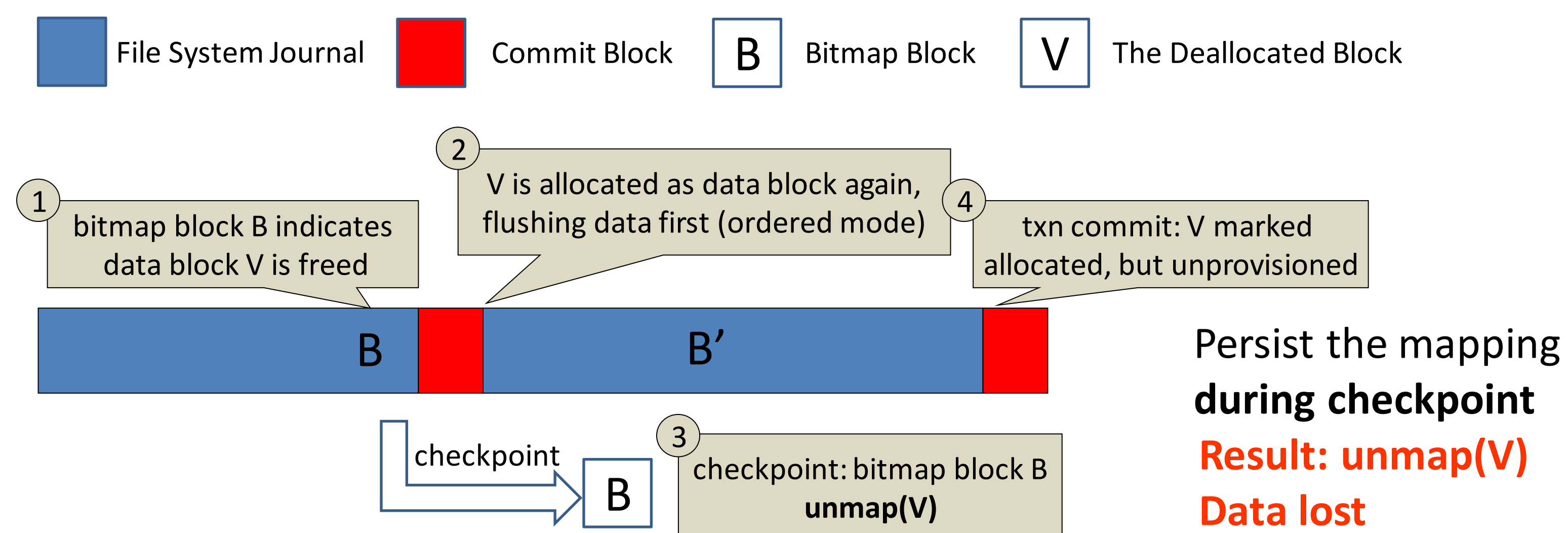
- Policies may need to persist their own metadata
  - Key-Value pair like mapping



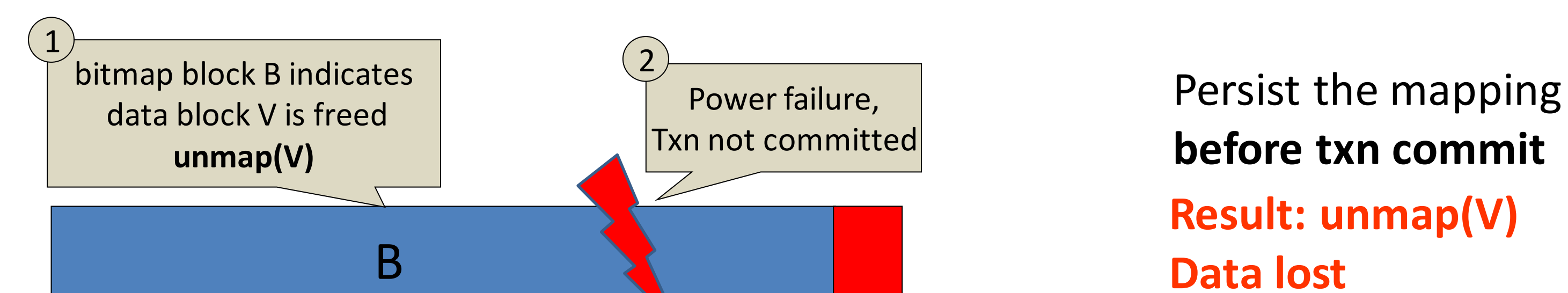
- Mapping needs to be **consistent** with application state
- Mapping shouldn't violate **durability expectations**

## Challenge: Consistency

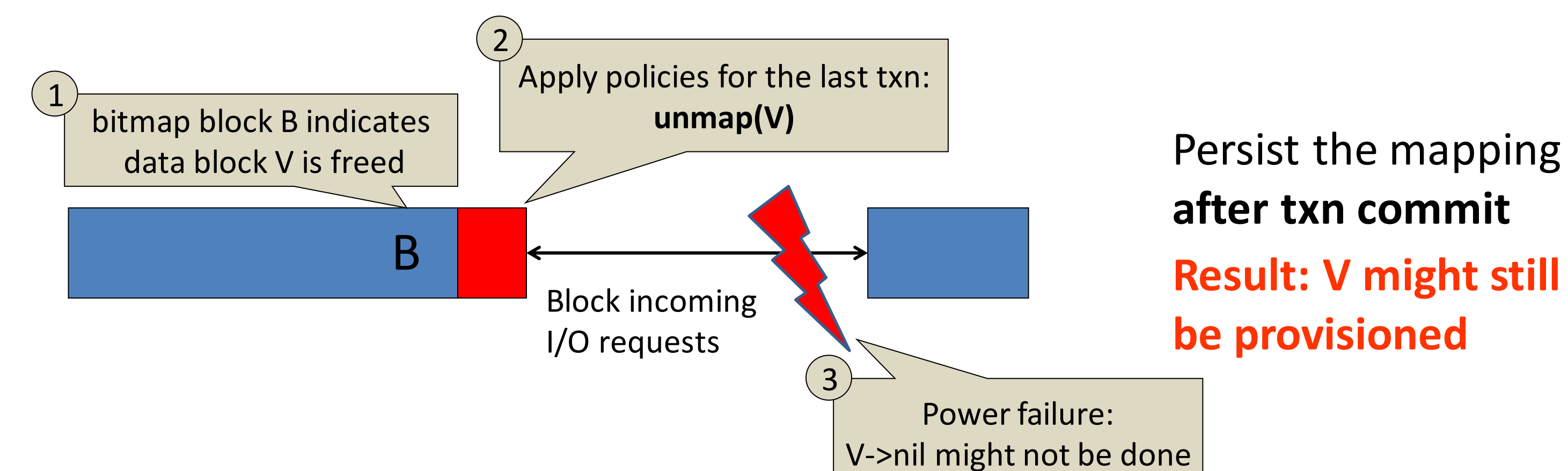
- Example Policy: Thin Provisioning in a journaled file system
  - When bitmap deallocate a block, reclaim it



Lessons: ordering must be preserved



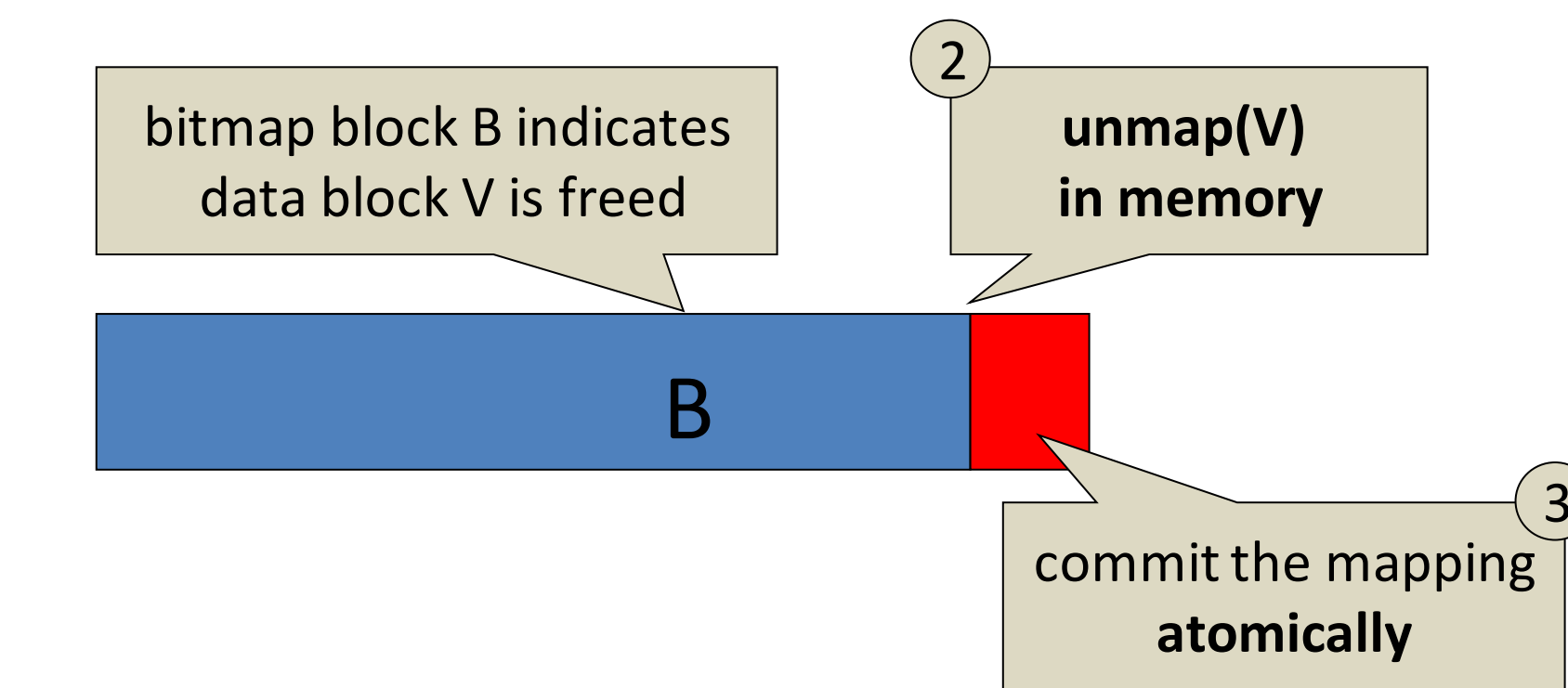
Lessons: before txn commit would fail



Lessons: after txn commit would fail

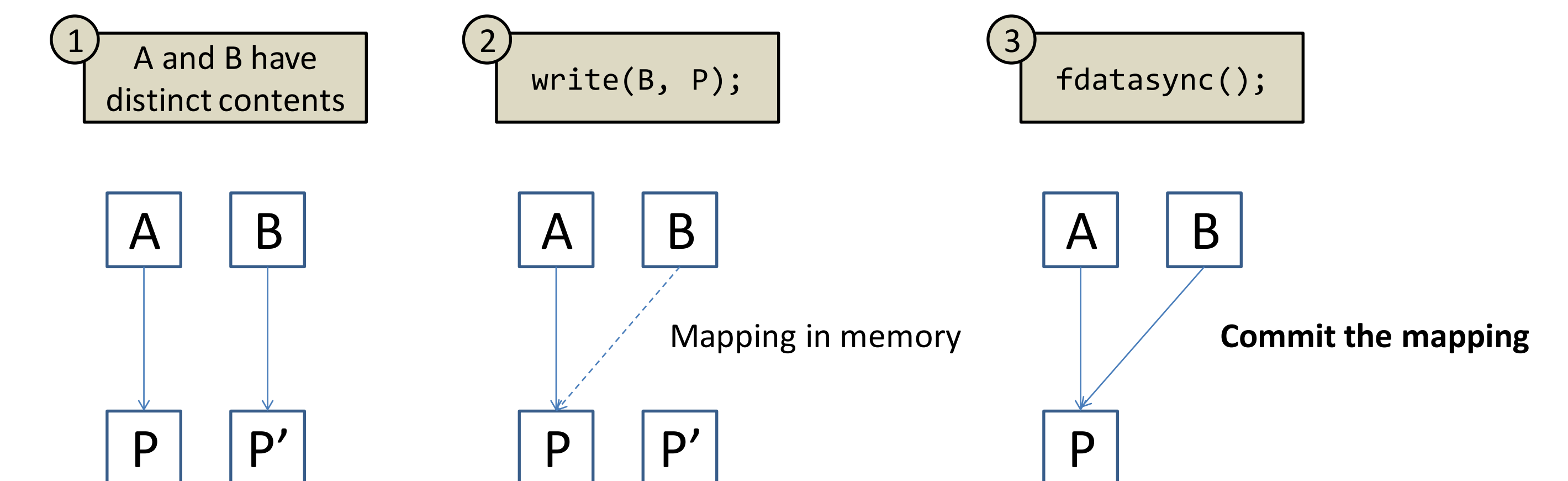
**Solution:** commit the mapping **atomically** with file

system commit record



## Challenge: Durability

- Application have durability expectations
  - After disk barrier applications expect that file data are durable
- Example Policy: Inline Deduplication
  - Blocks with same content mapped to same location
  - Need to persist the mapping on `fdatasync()`



## Summary of Solutions

Guarantee	Durability	Consistency
When to create mapping	On data update	After txn finish Before commit
When to commit mapping	On disk barrier	Commit atomically wit txn