



**ceph** – a unified distributed storage system

sage weil  
cloudopen – august 29, 2012

# outline

- why you should care
- what is it, what it does
- how it works
  - architecture
- how you can use it
  - librados
  - radosgw
  - RBD
  - file system
- who we are, why we do this

why should you care about another  
storage system?

# requirements

- diverse storage needs
  - object storage
  - block devices (for VMs) with snapshots, cloning
  - shared file system with POSIX, coherent caches
  - structured data... files, block devices, or objects?
- **scale**
  - terabytes, petabytes, exabytes
  - heterogeneous hardware
  - reliability and fault tolerance

# time

- ease of administration
- no manual data migration, load balancing
- painless scaling
  - expansion **and** contraction
  - seamless migration

# cost

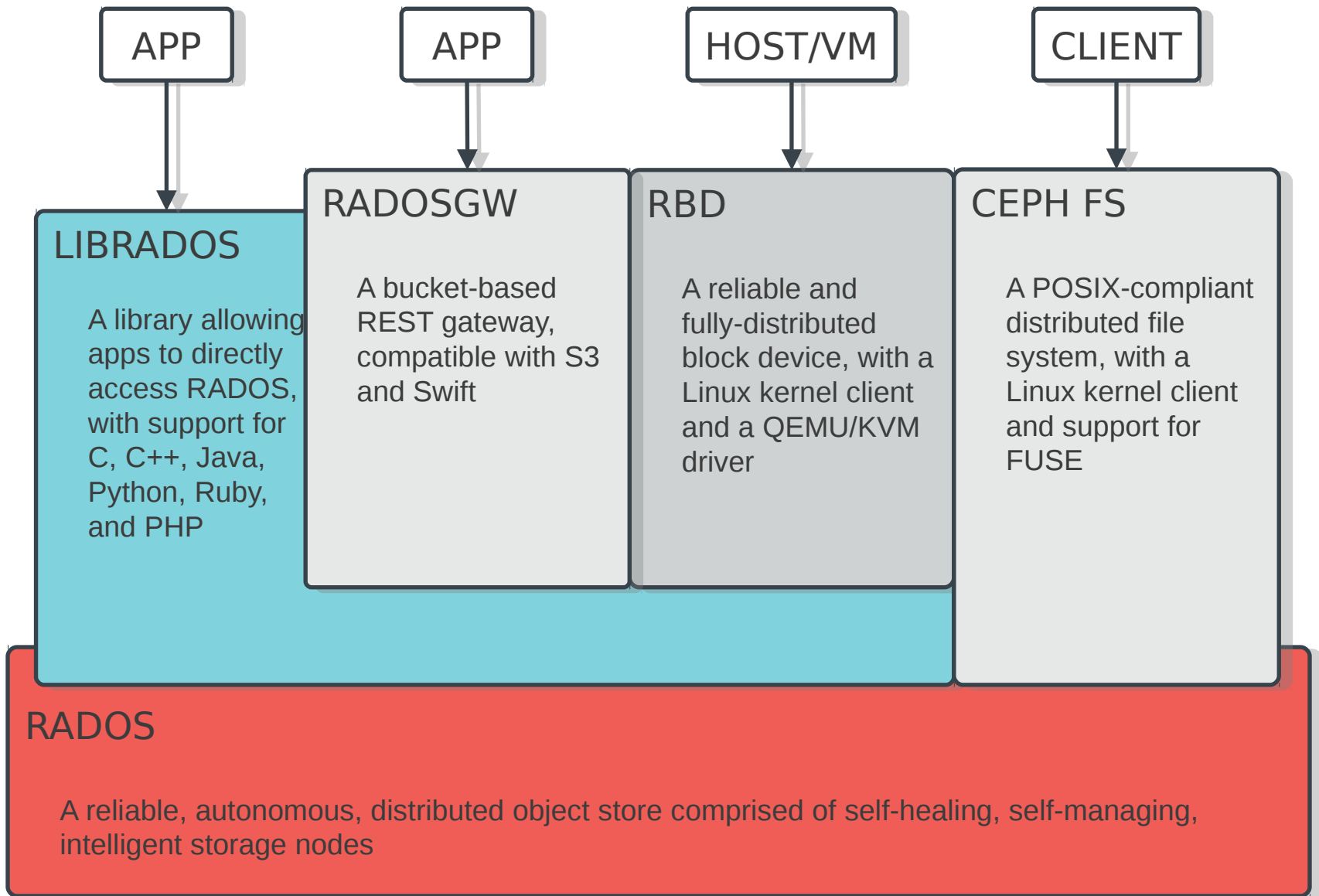
- linear function of size or performance
- incremental expansion
  - no fork-lift upgrades
- no vendor lock-in
  - choice of hardware
  - choice of software
- open

what is **ceph**?

# unified storage system

- objects
  - native
  - RESTful
- block
  - thin provisioning, snapshots, cloning
- file
  - strong consistency, snapshots





# open source

- LGPLv2
  - copyleft
  - ok to link to proprietary code
- no copyright assignment
  - no dual licensing
  - no “enterprise-only” feature set
- active community
- commercial support

# distributed storage system

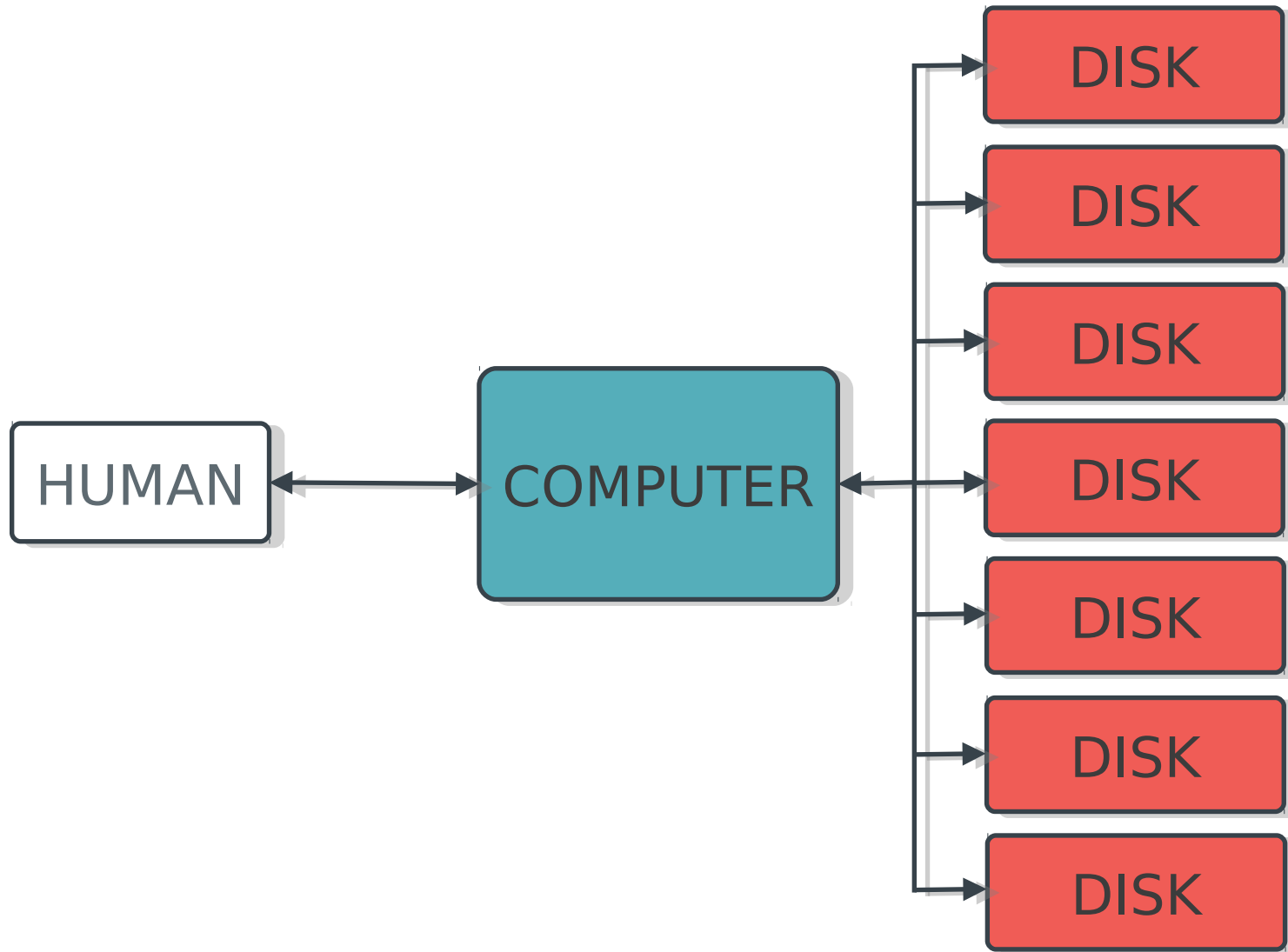
- data center scale
  - 10s to 10,000s of machines
  - terabytes to exabytes
- fault tolerant
  - no single point of failure
  - commodity hardware
- self-managing, self-healing

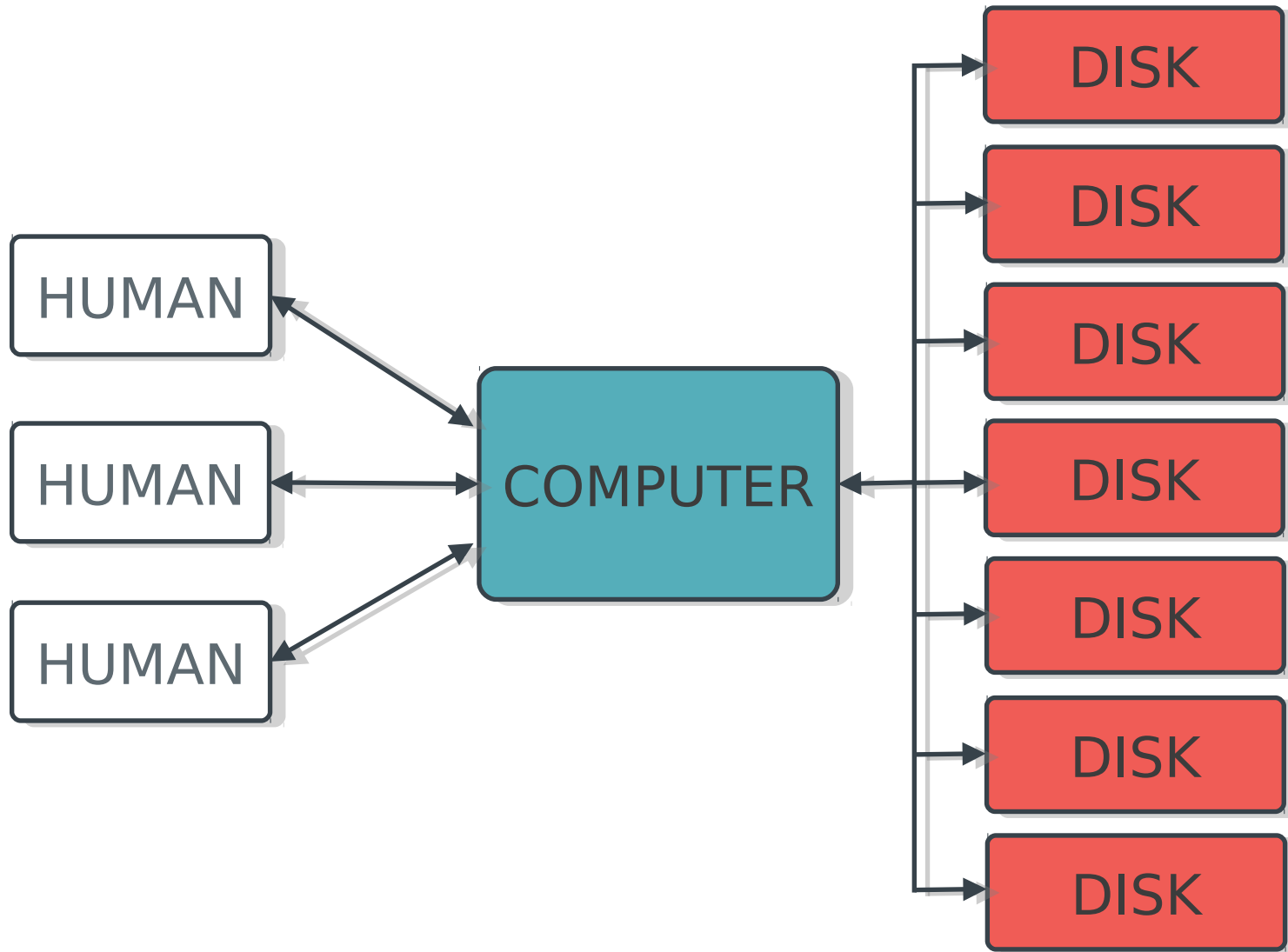
# ceph object model

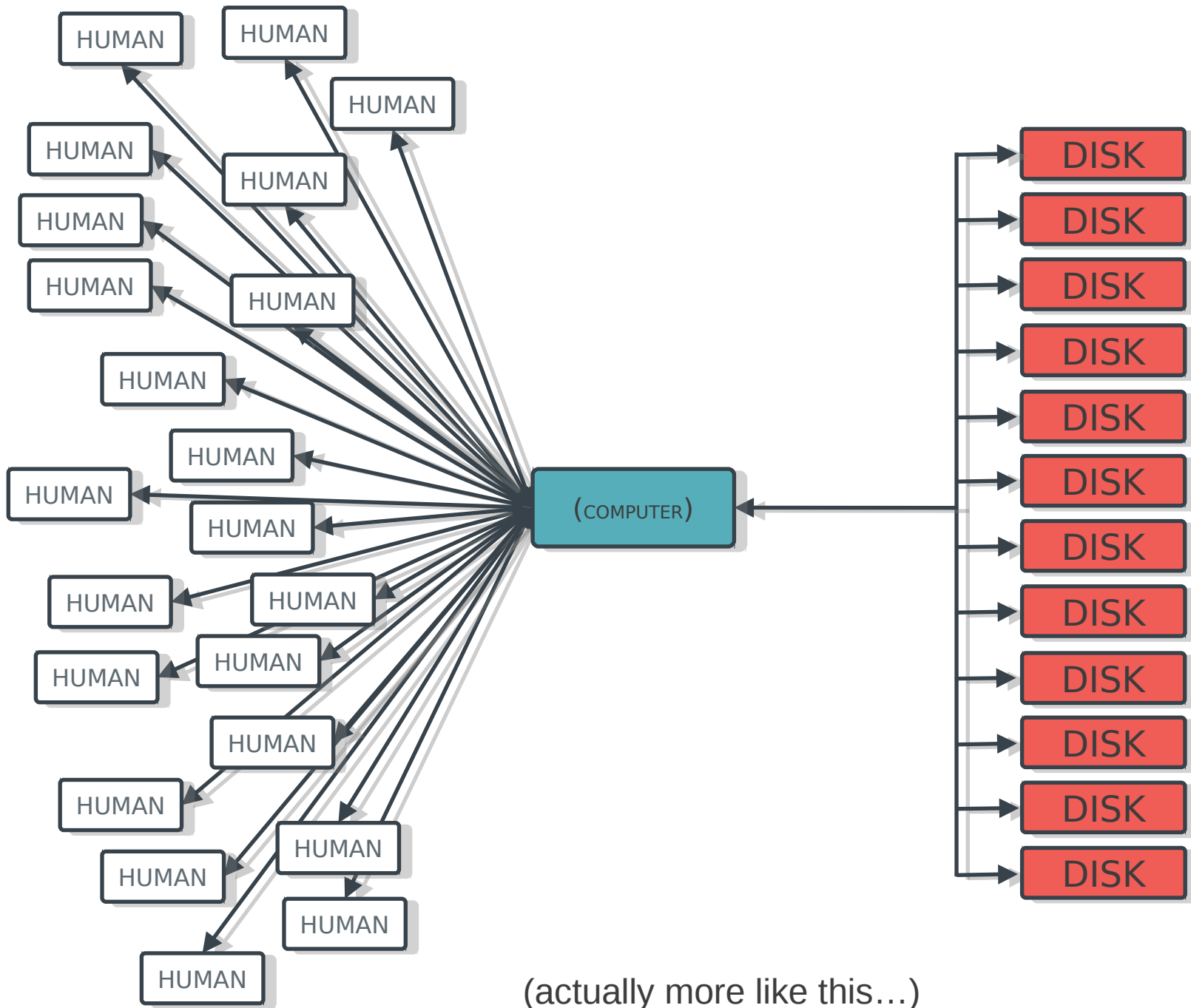
- pools
  - 1s to 100s
  - independent namespaces or object collections
  - replication level, placement policy
- objects
  - bazillions
  - blob of data (bytes to gigabytes)
  - attributes (e.g., “version=12”; bytes to kilobytes)
  - key/value bundle (bytes to gigabytes)

# why start with objects?

- more useful than (disk) blocks
  - names in a single flat namespace
  - variable size
  - simple API with rich semantics
- more scalable than files
  - no hard-to-distribute hierarchy
  - update semantics do not span objects
  - workload is trivially parallel

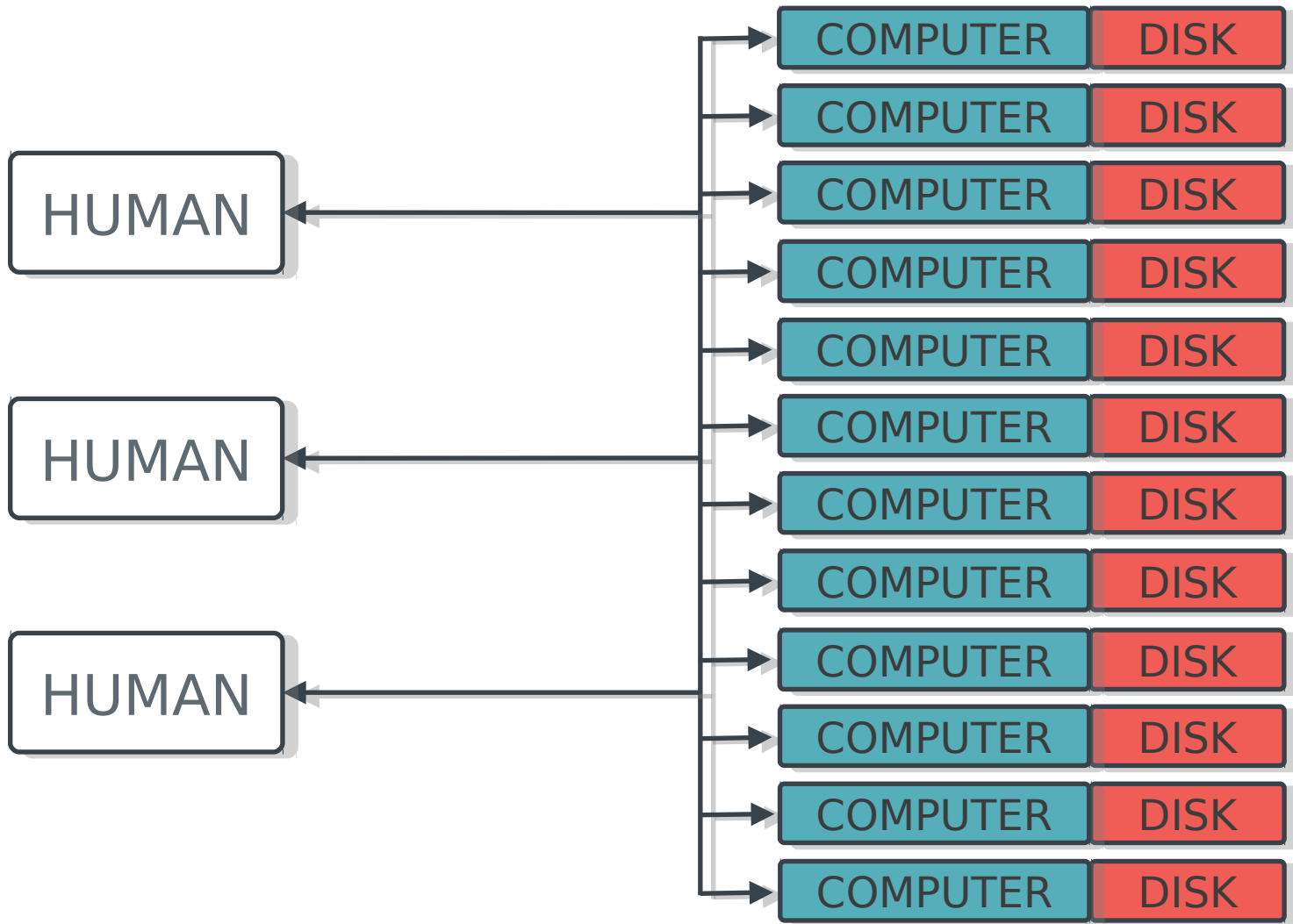


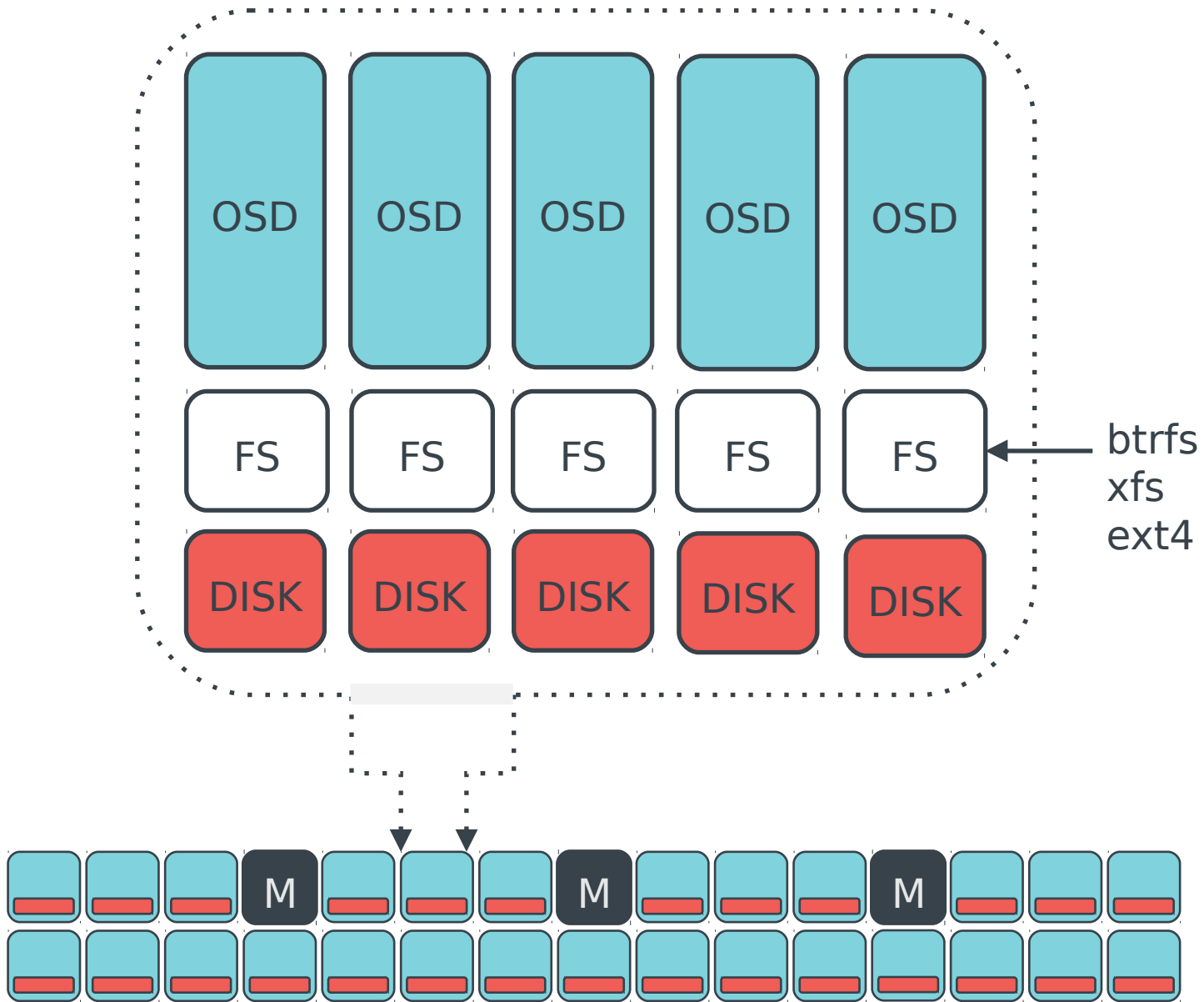




(actually more like this...)



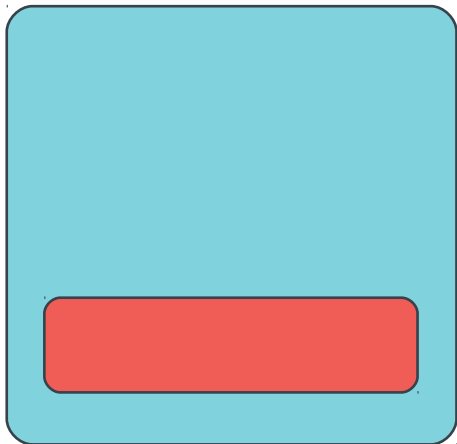






## Monitors:

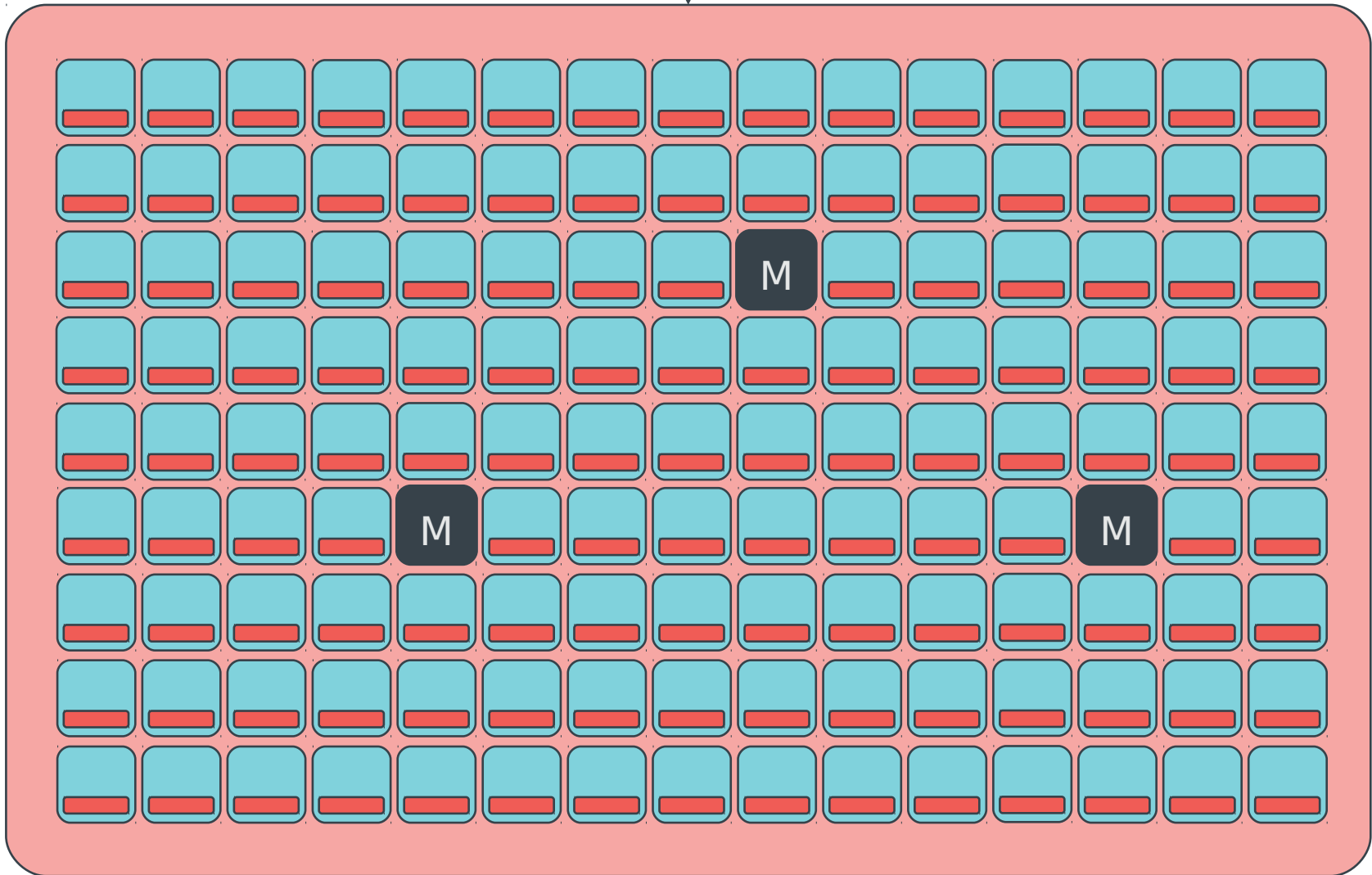
- Maintain cluster membership and state
- Provide consensus for distributed decision-making
- Small, odd number
- These do **not** serve stored objects to clients



## Object Storage Daemons (OSDs):

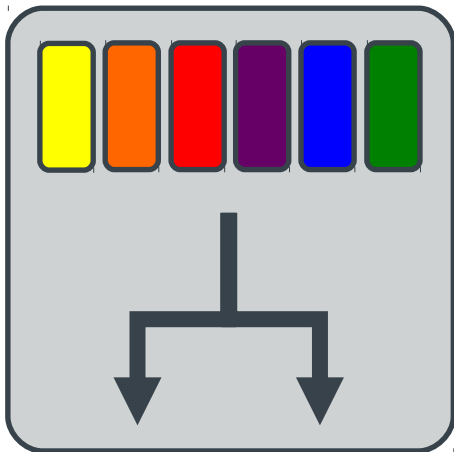
- At least three in a cluster
- One per disk or RAID group
- Serve stored objects to clients
- Intelligently peer to perform replication tasks

HUMAN



# data distribution

- all objects are replicated N times
- objects are automatically placed, balanced, migrated in a **dynamic** cluster
- must consider physical infrastructure
  - ceph-osds on hosts in racks in rows in data centers
- three approaches
  - pick a spot; remember where you put it
  - pick a spot; write down where you put it
  - calculate where to put it, where to find it



## CRUSH

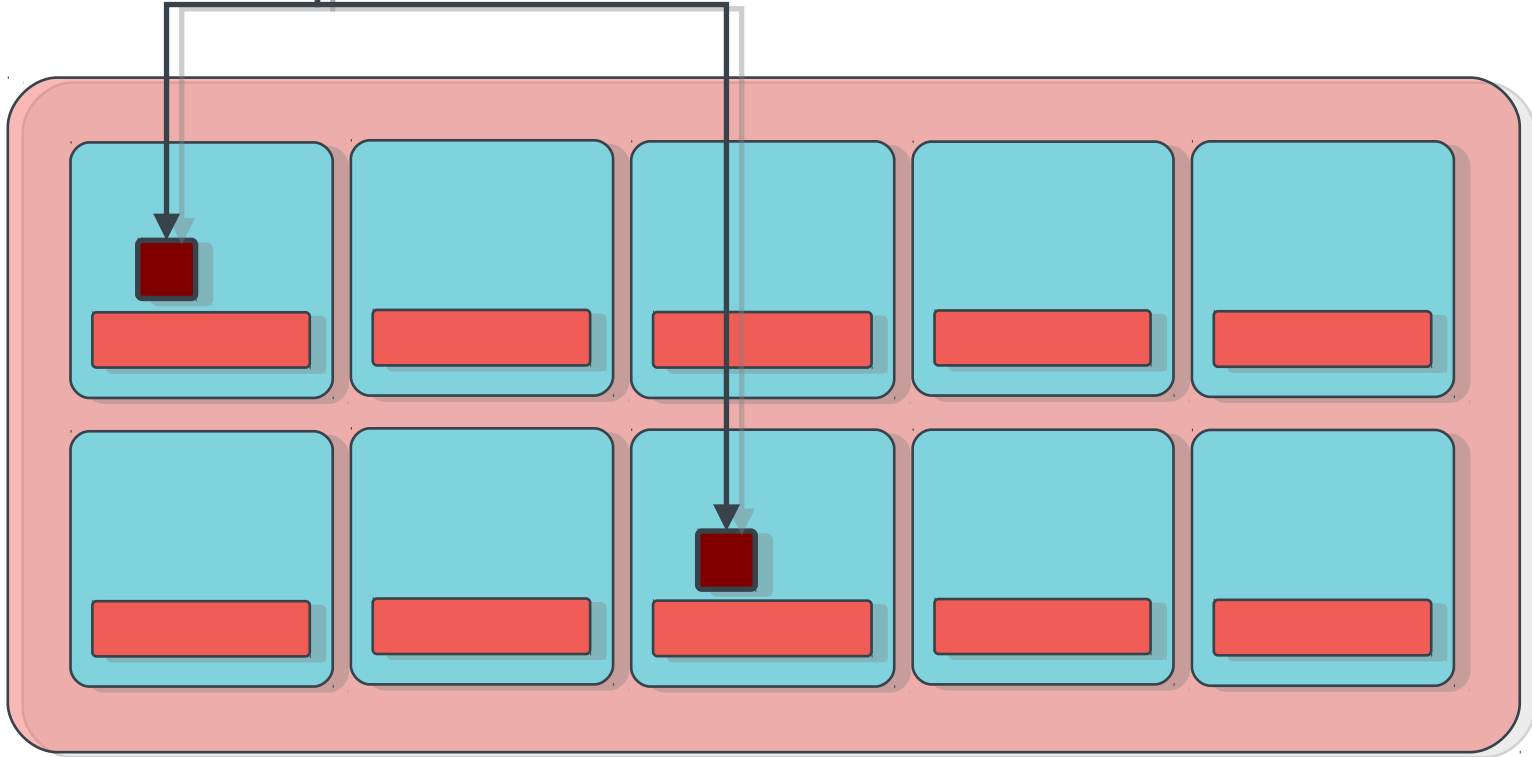
- Pseudo-random placement algorithm
- Fast calculation, **no lookup**
- Repeatable, deterministic
- Ensures even distribution
- Stable mapping
  - Limited data migration
- Rule-based configuration
  - specifiable replication
  - infrastructure topology aware
  - allows weighting

10 10 01 01 10 10 01 11 01 10

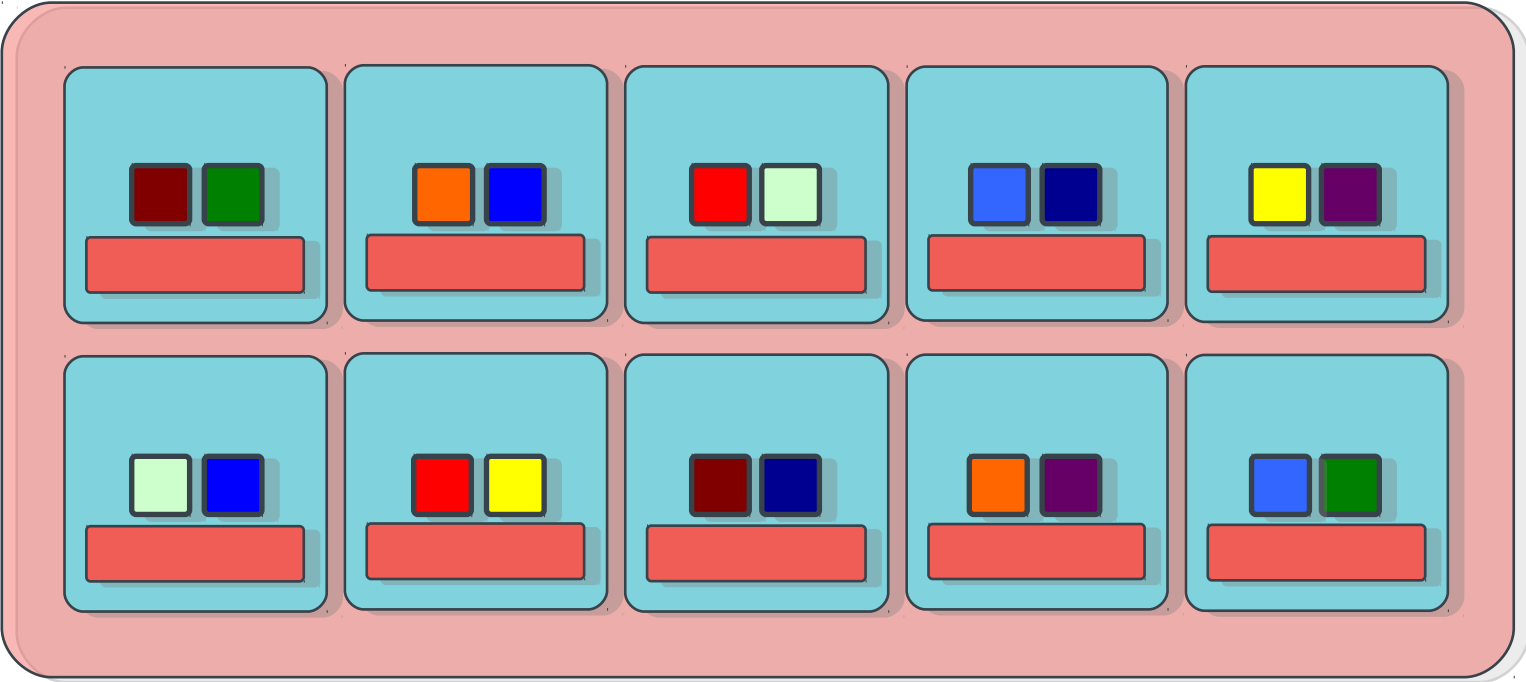
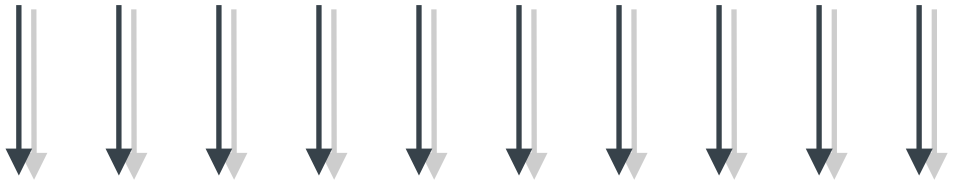
hash(object name) % num pg



CRUSH(pg, cluster state, policy)



10 10 01 01 10 10 01 11 01 10

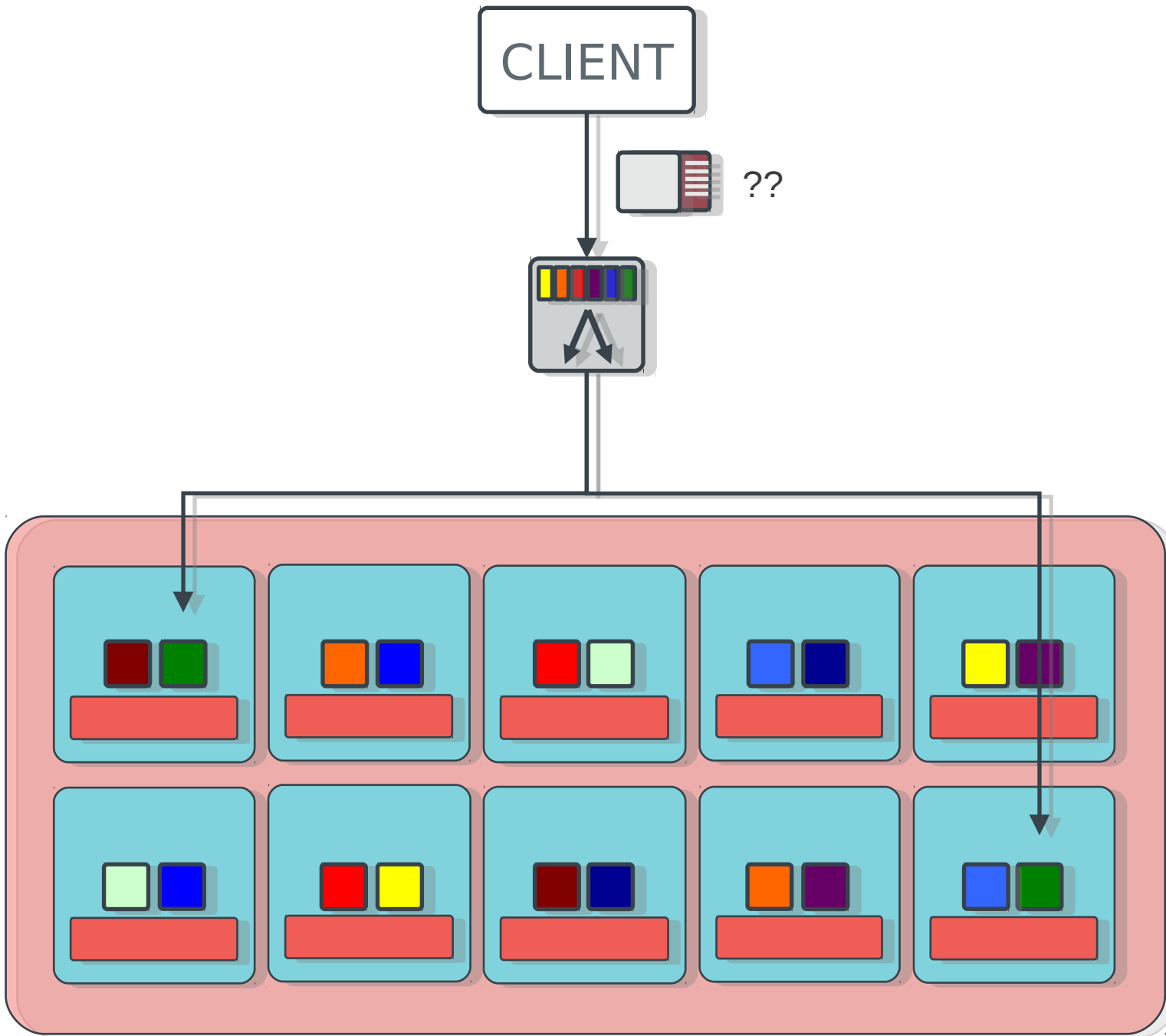


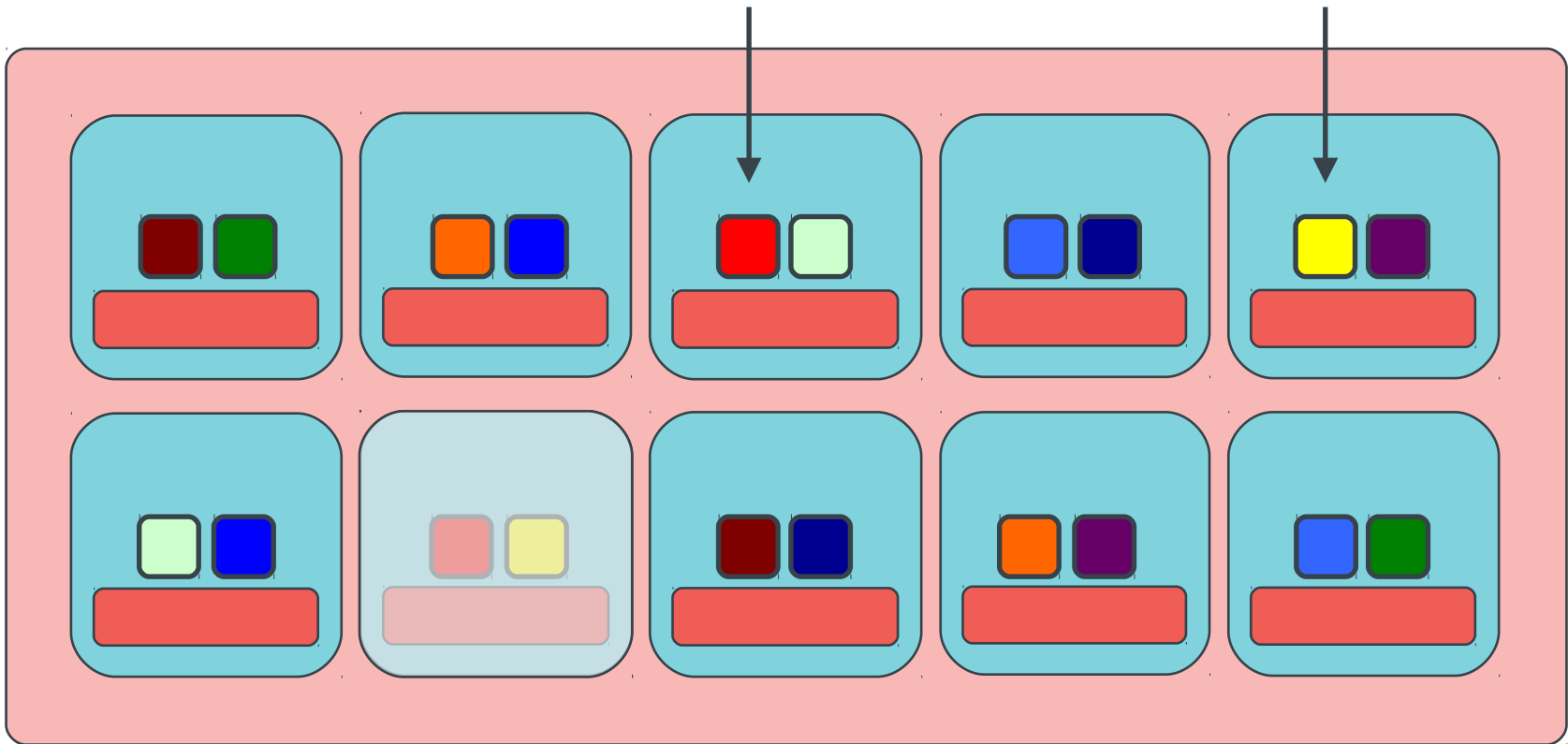


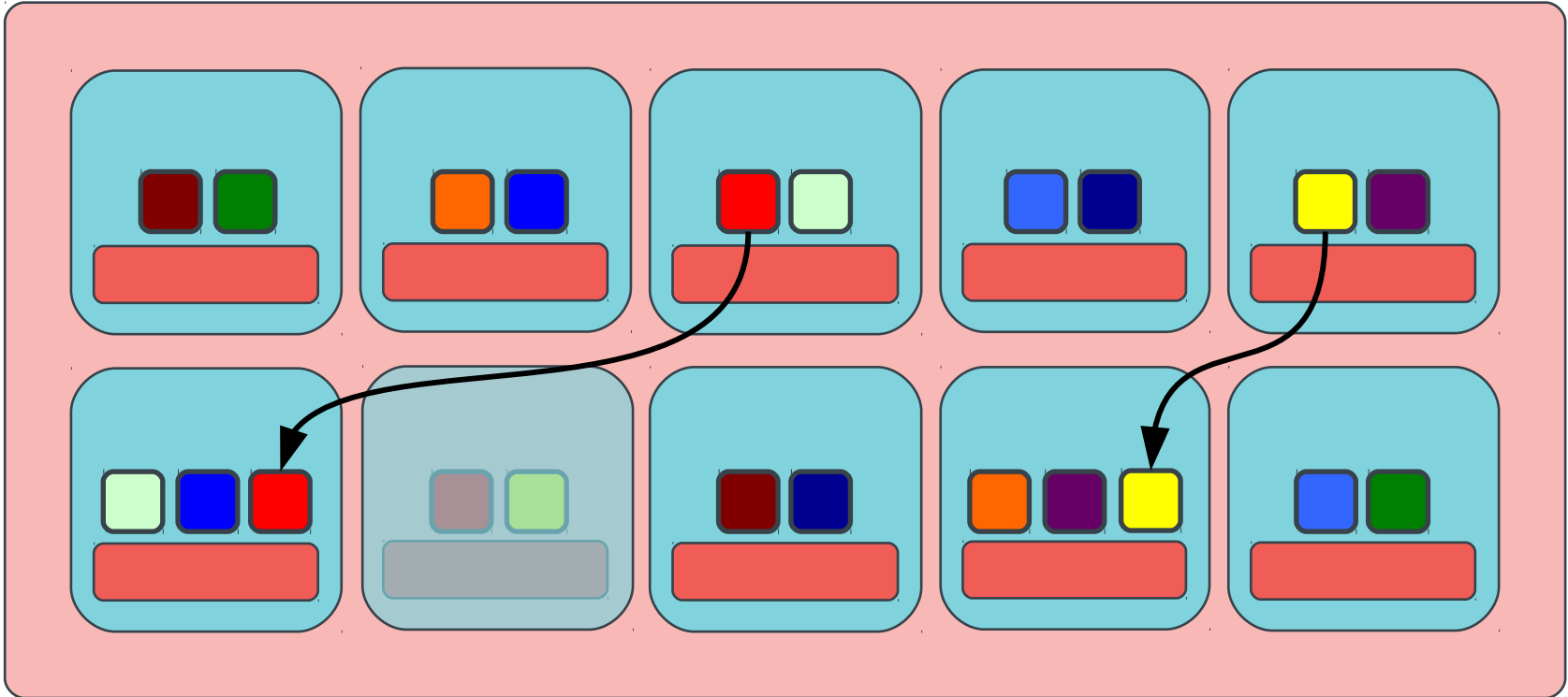
# RADOS

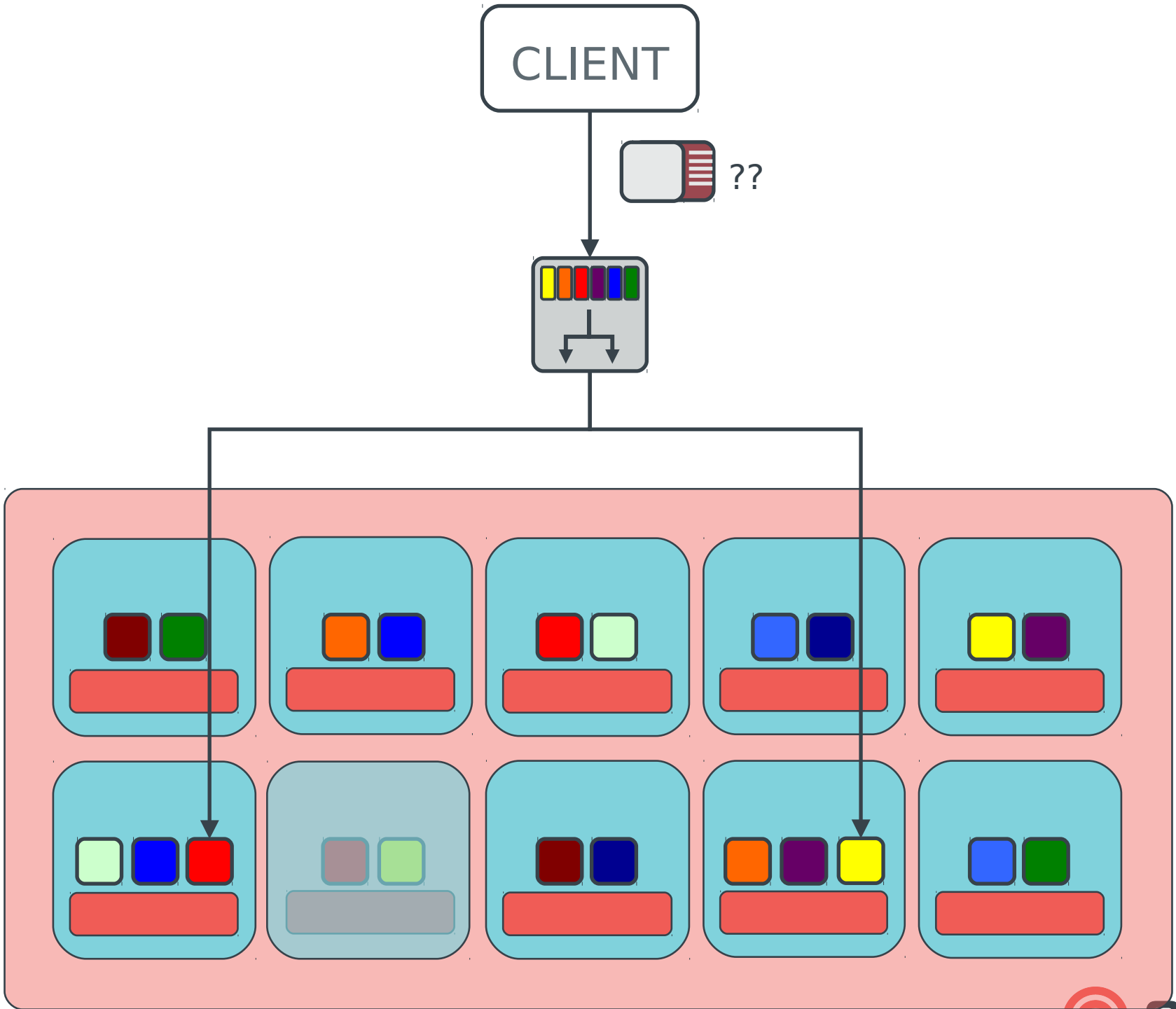
- monitors publish **osd map** that describes cluster state
  - ceph-osd node status (up/down, weight, IP)
  - CRUSH function specifying desired data distribution
- object storage daemons (OSDs)
  - safely replicate and store object
  - migrate data as the cluster changes over time
  - coordinate based on shared view of reality
- decentralized, distributed approach allows
  - massive scales (10,000s of servers or more)
  - the illusion of a single copy with consistent behavior

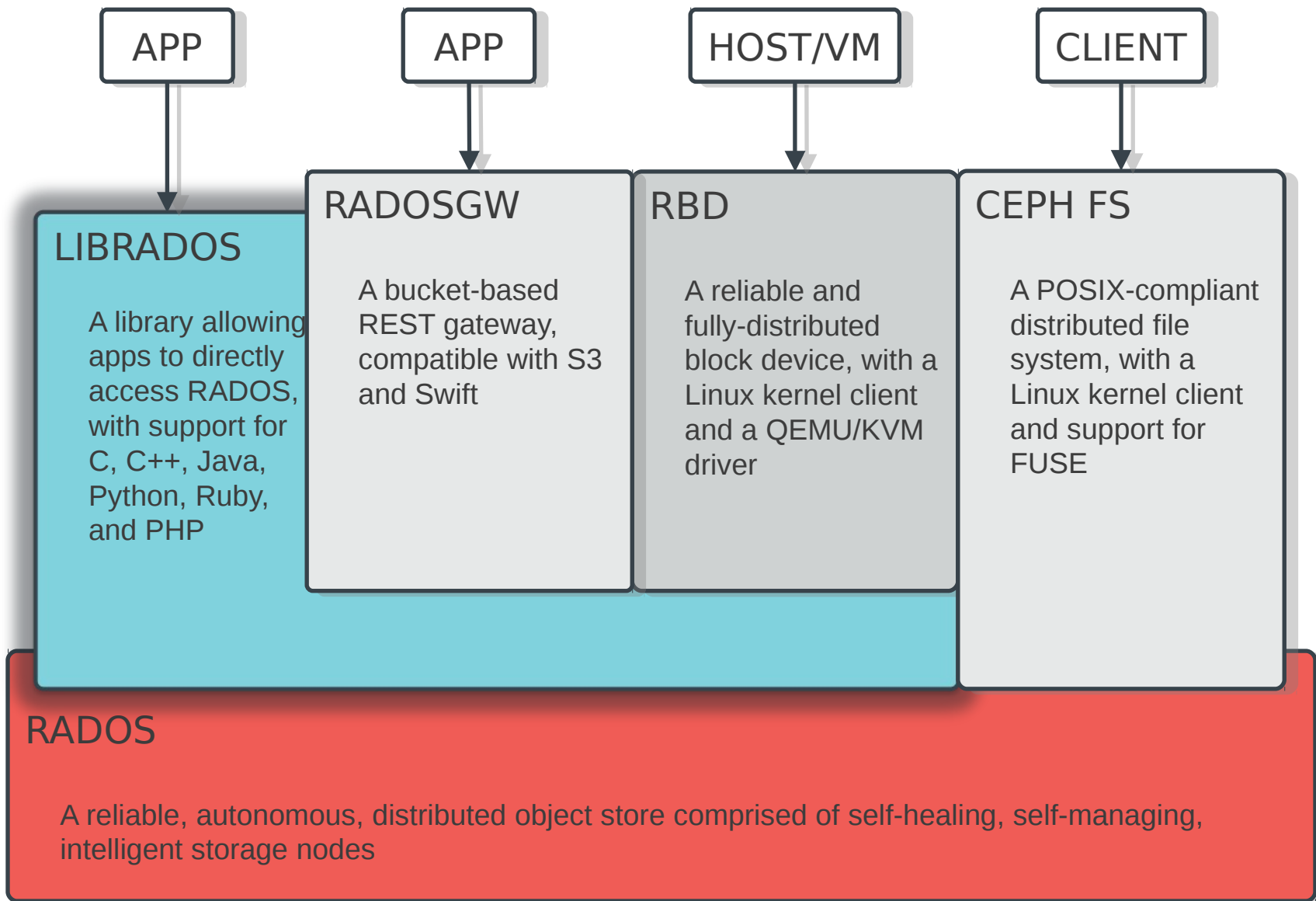


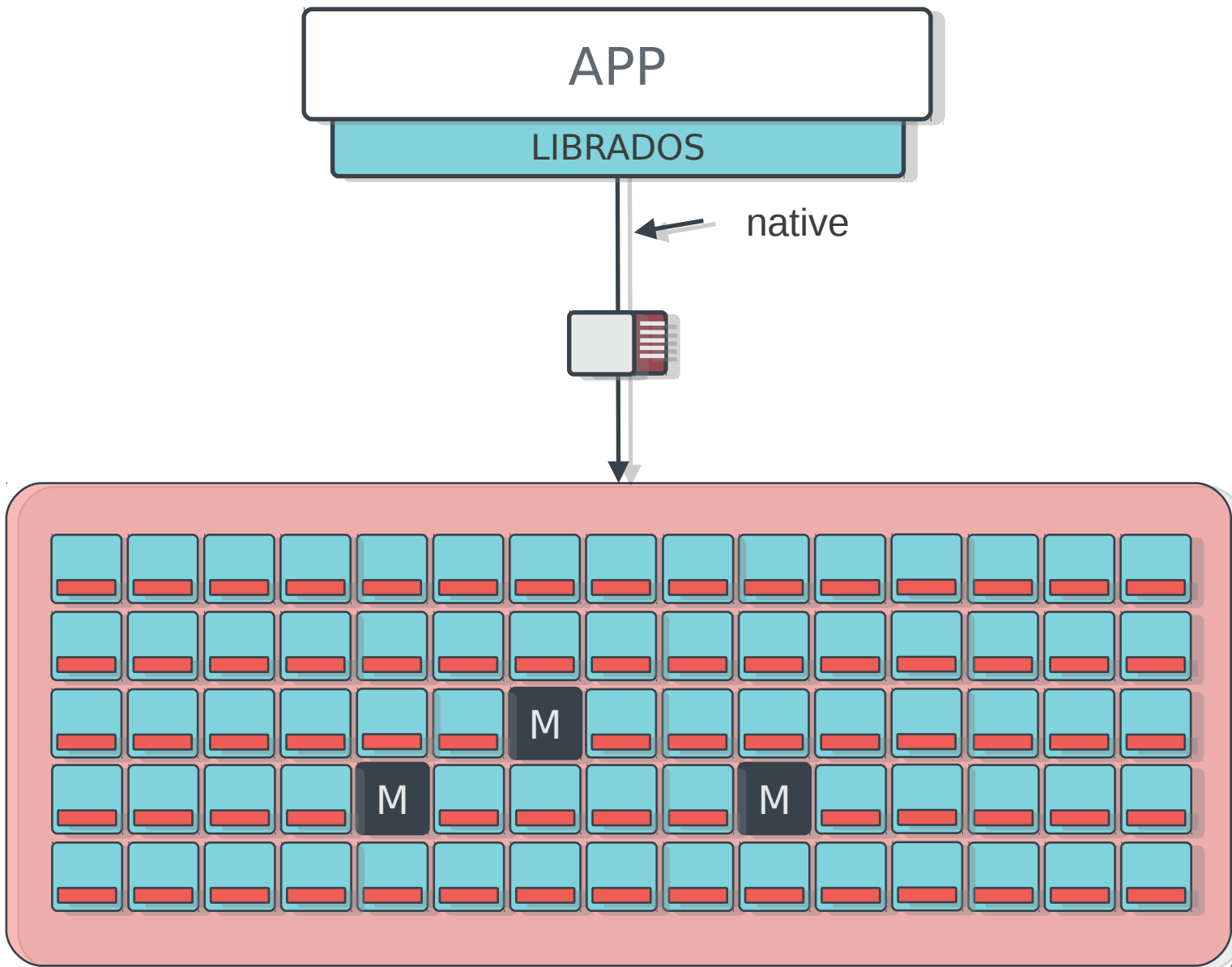










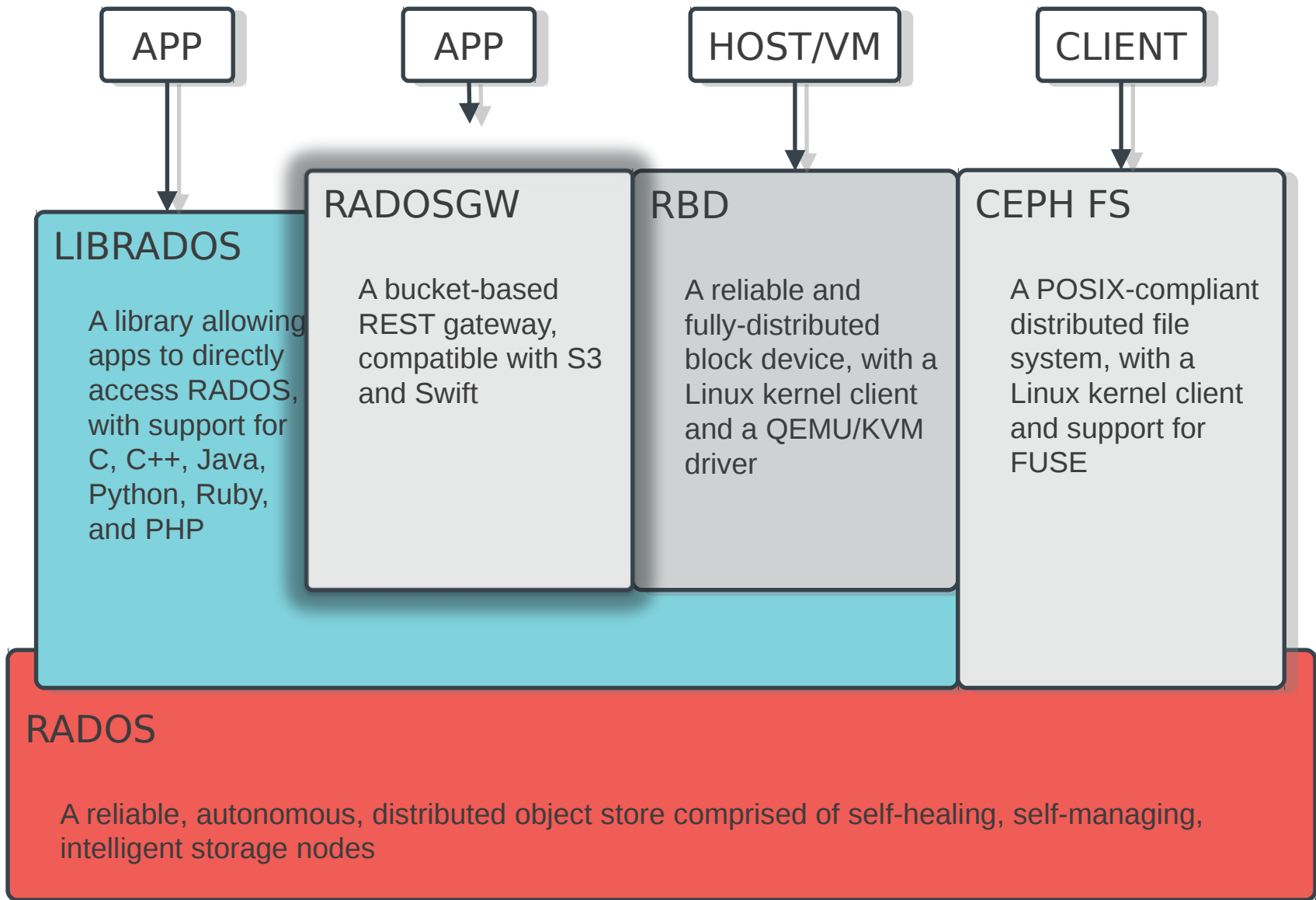


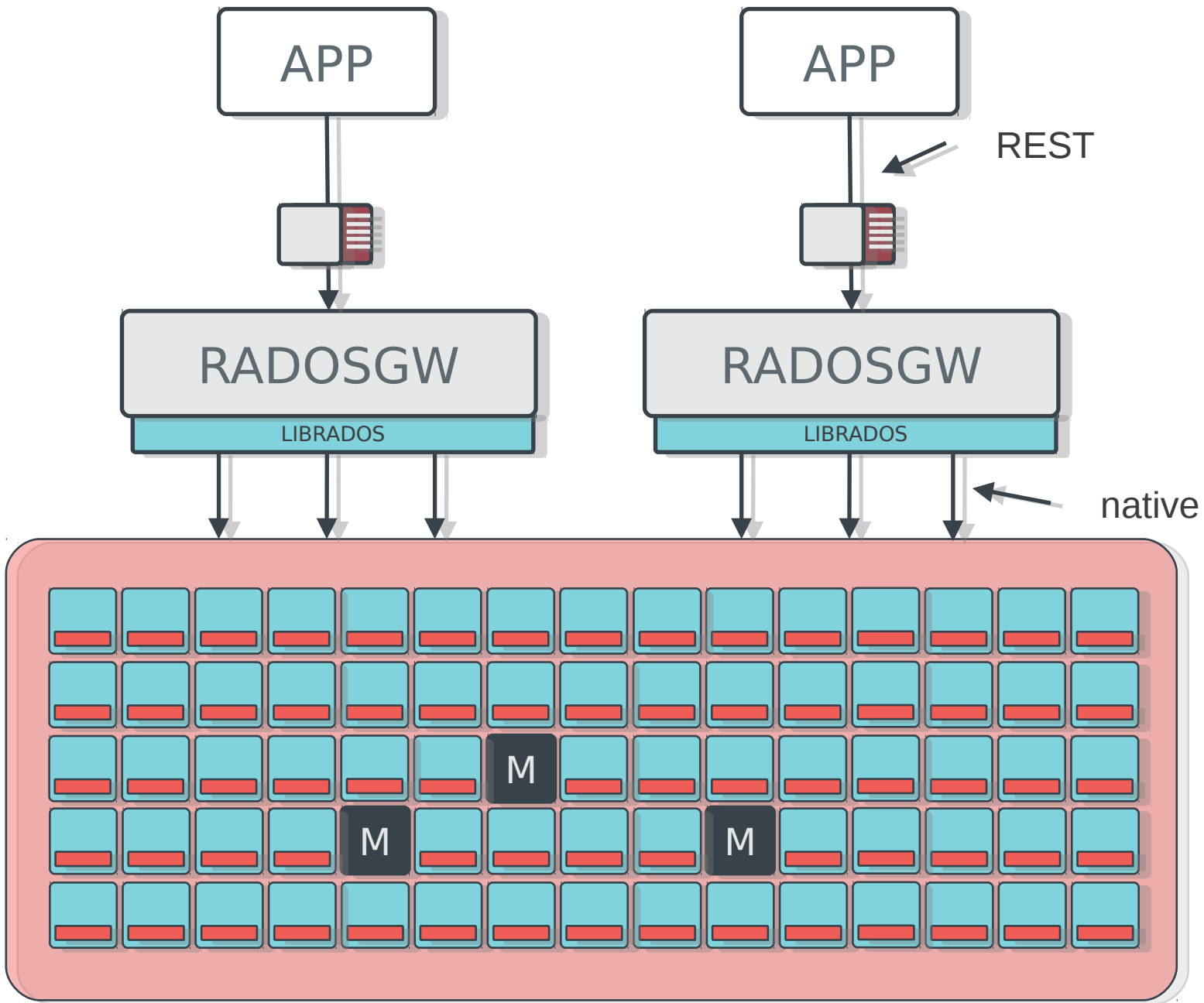


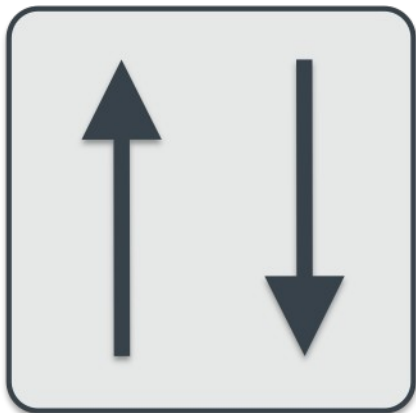
## LIBRADOS

- Provides direct access to RADOS for applications
- C, C++, Python, PHP, Java
- No HTTP overhead



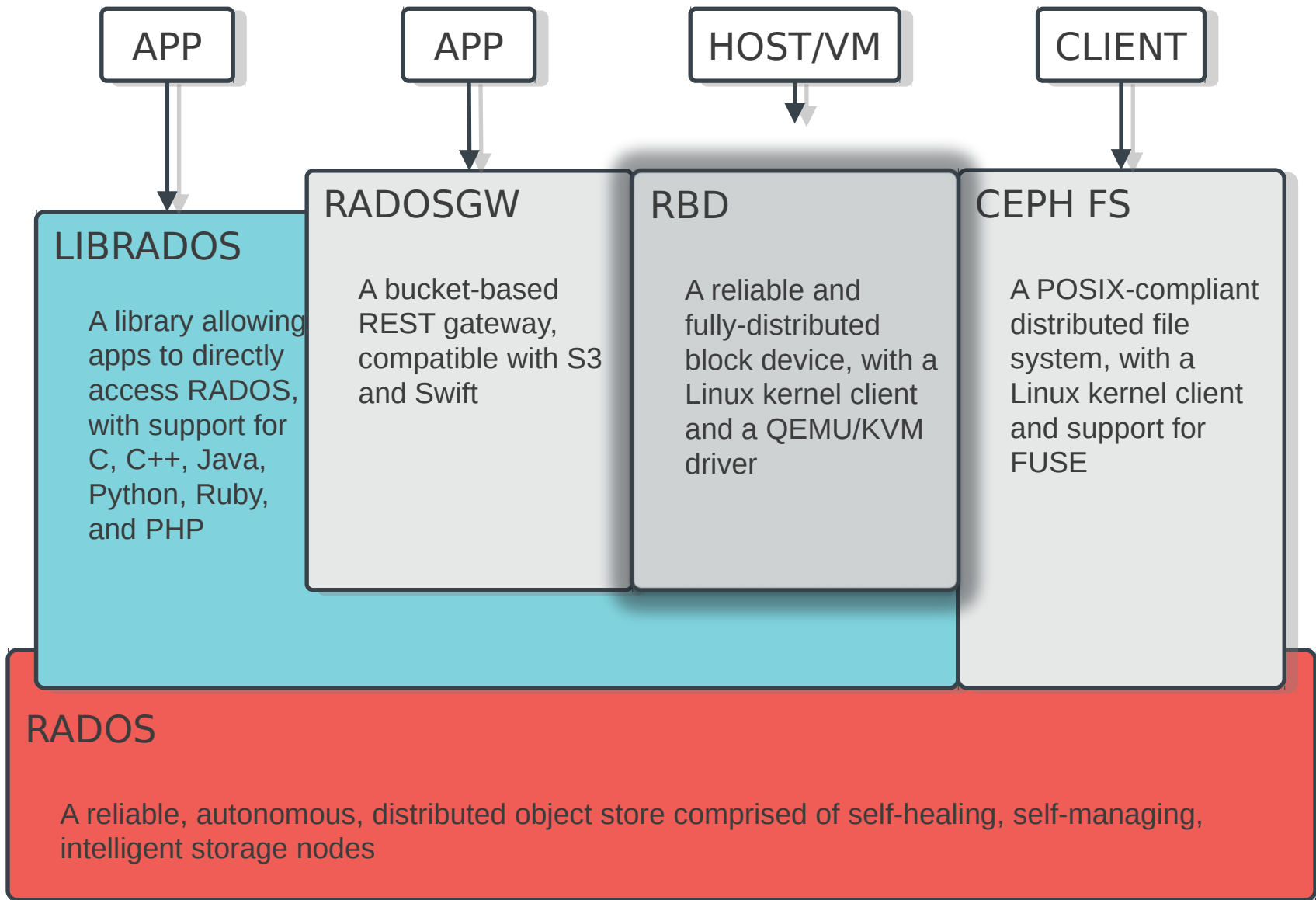


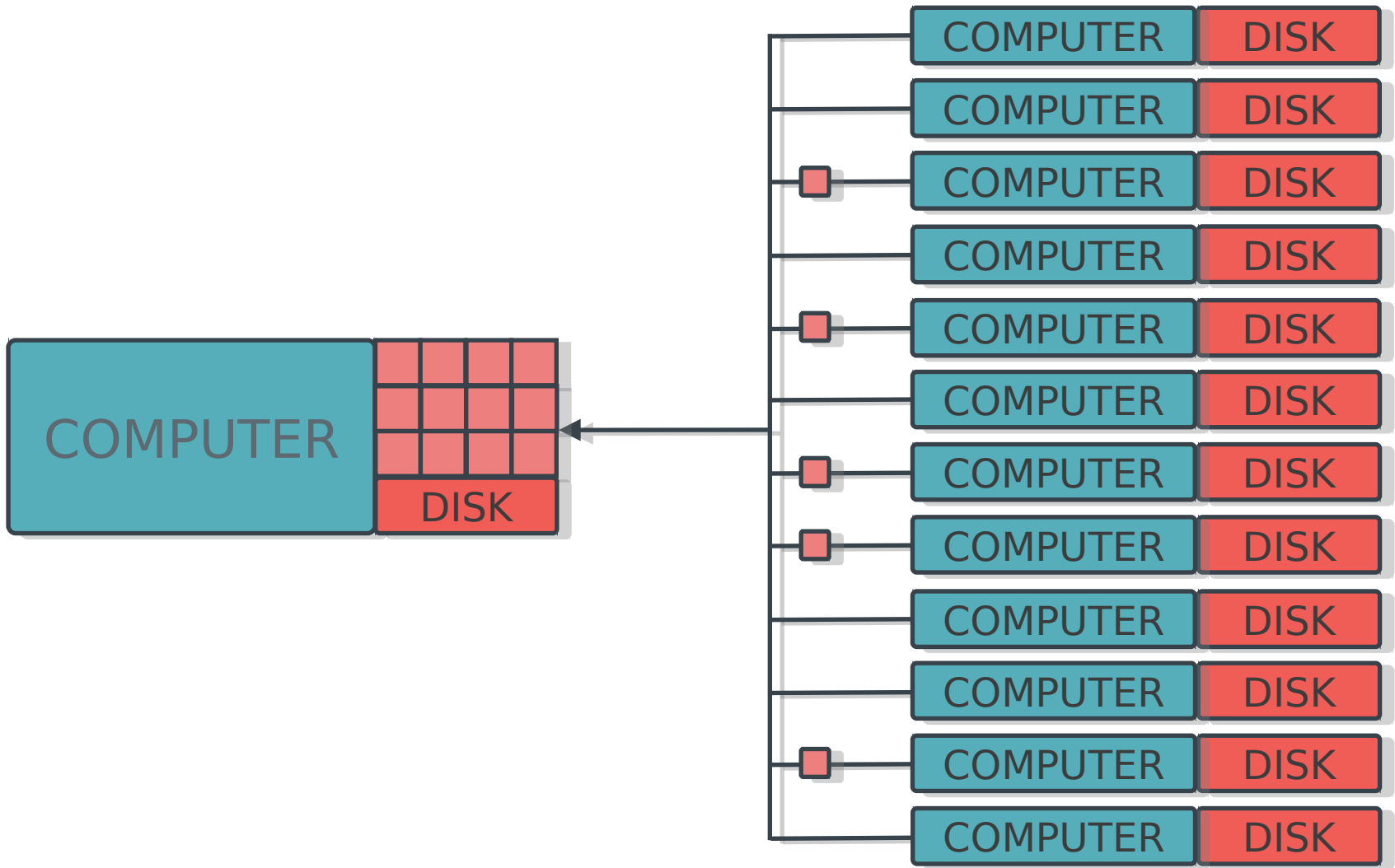


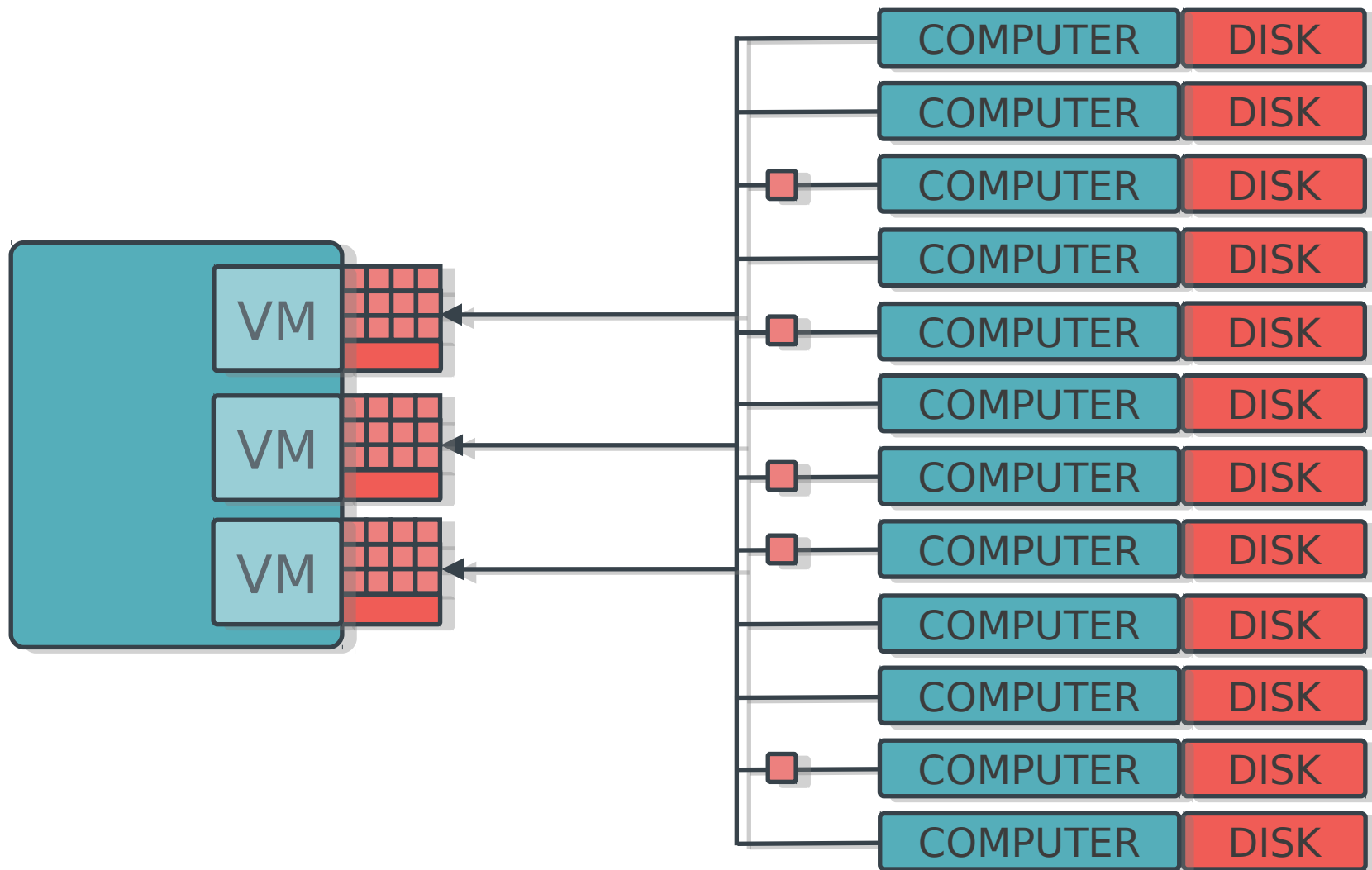


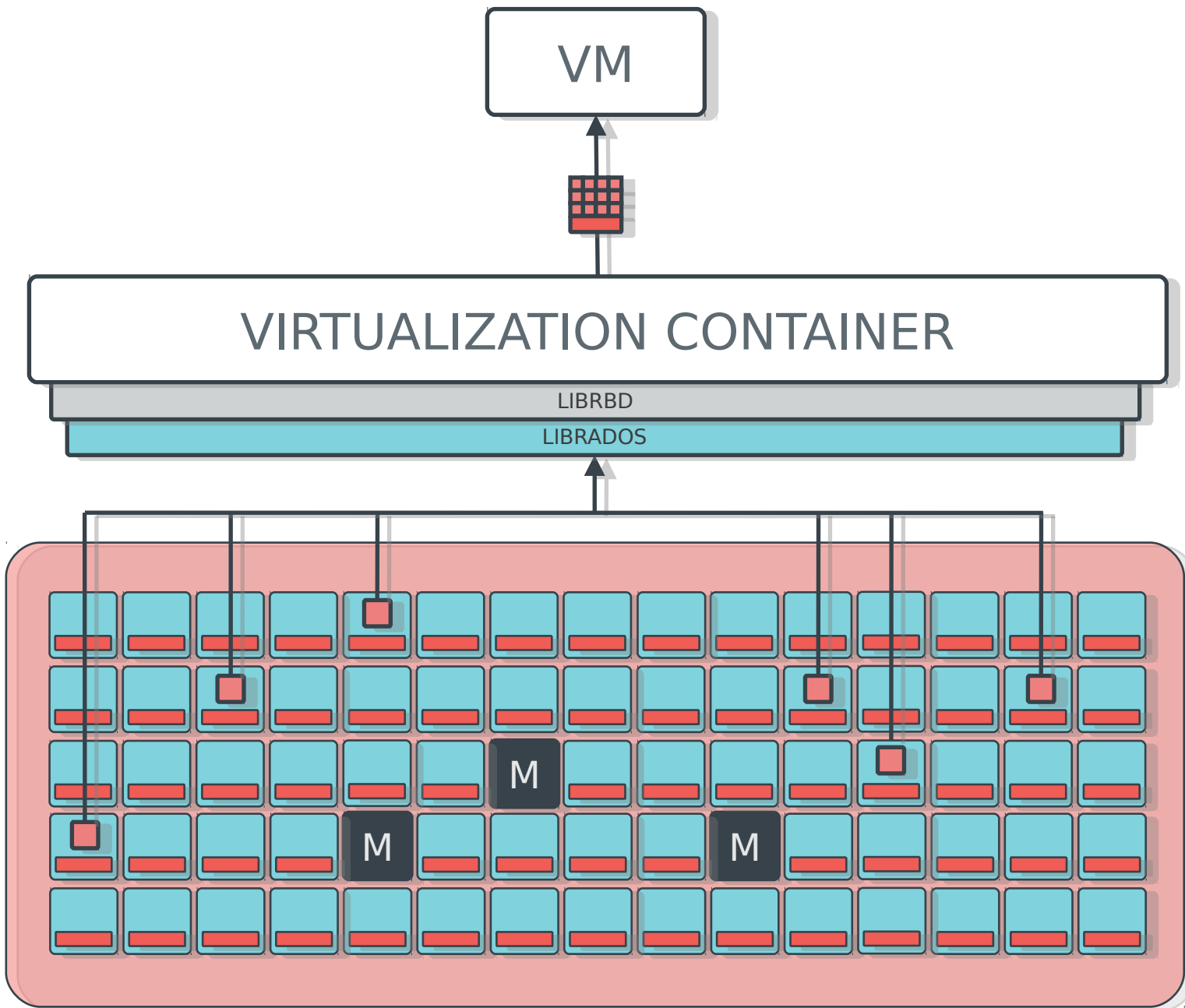
## RADOS Gateway:

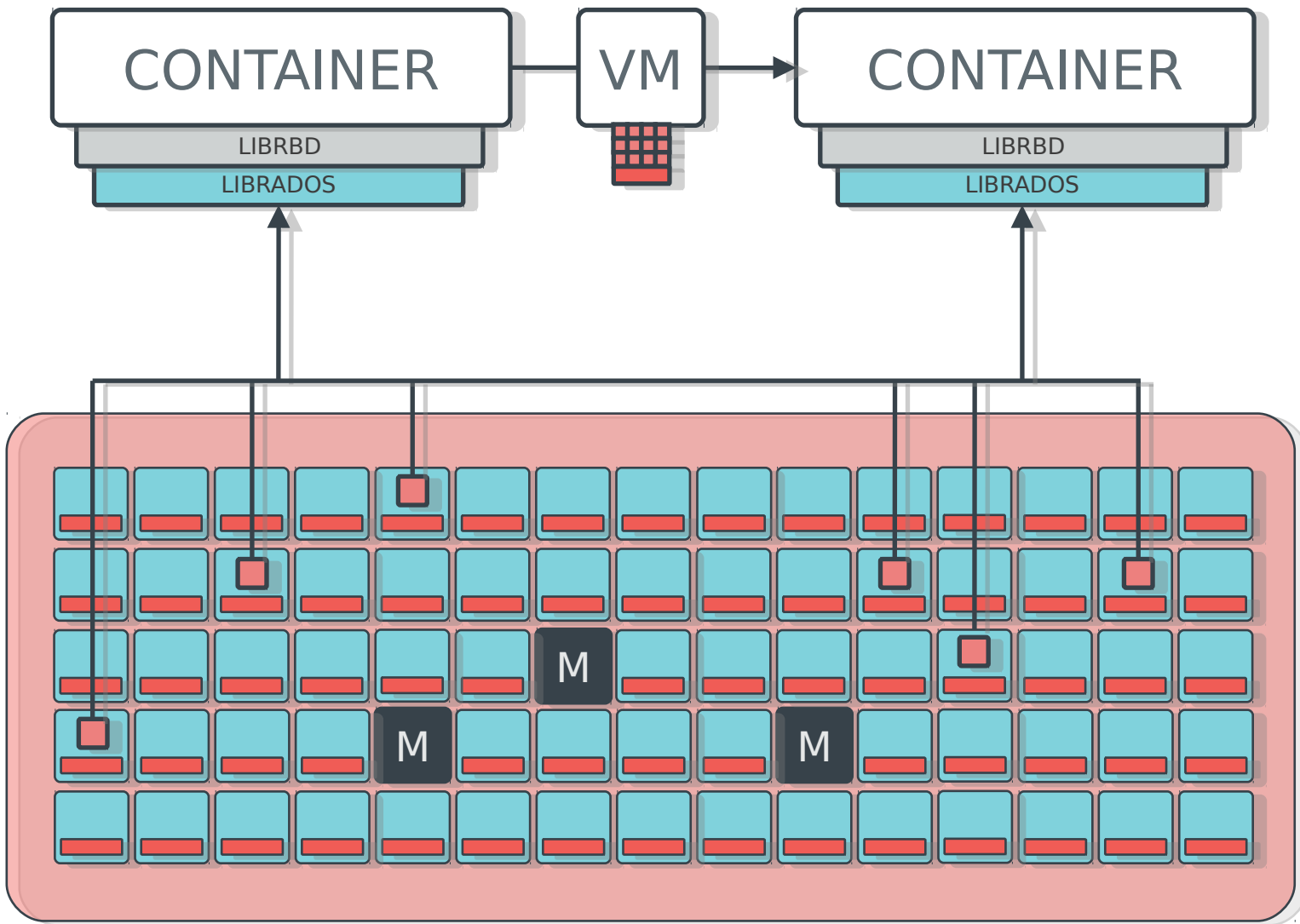
- REST-based interface to RADOS
- Supports buckets, accounting
- Compatible with S3 and Swift applications



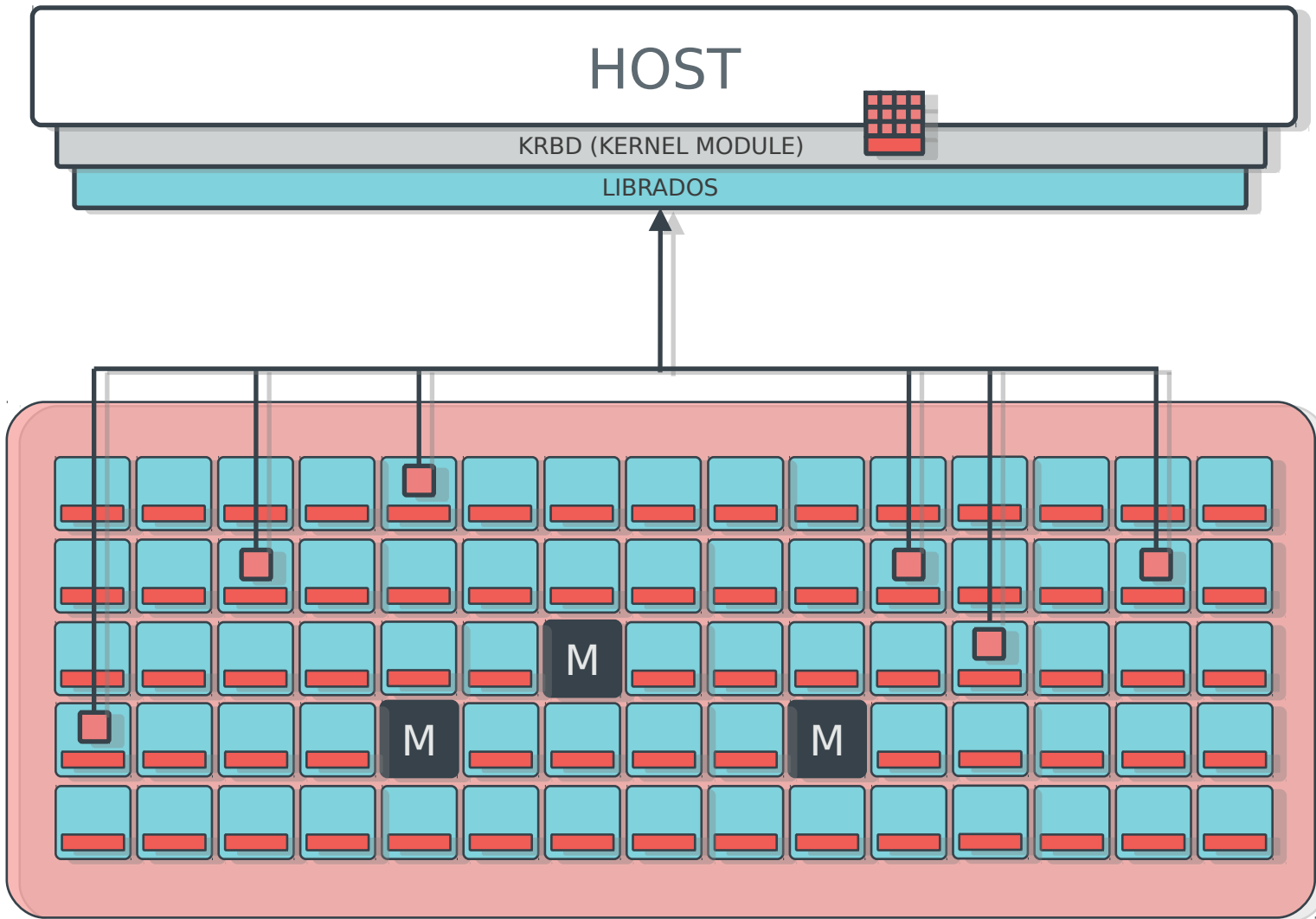


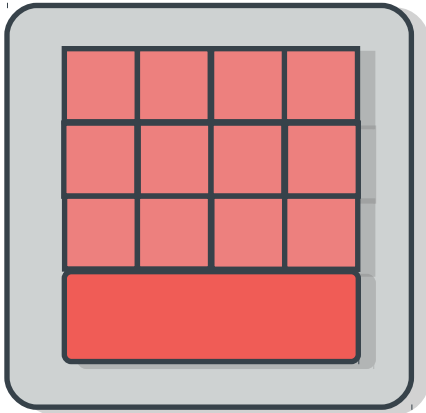








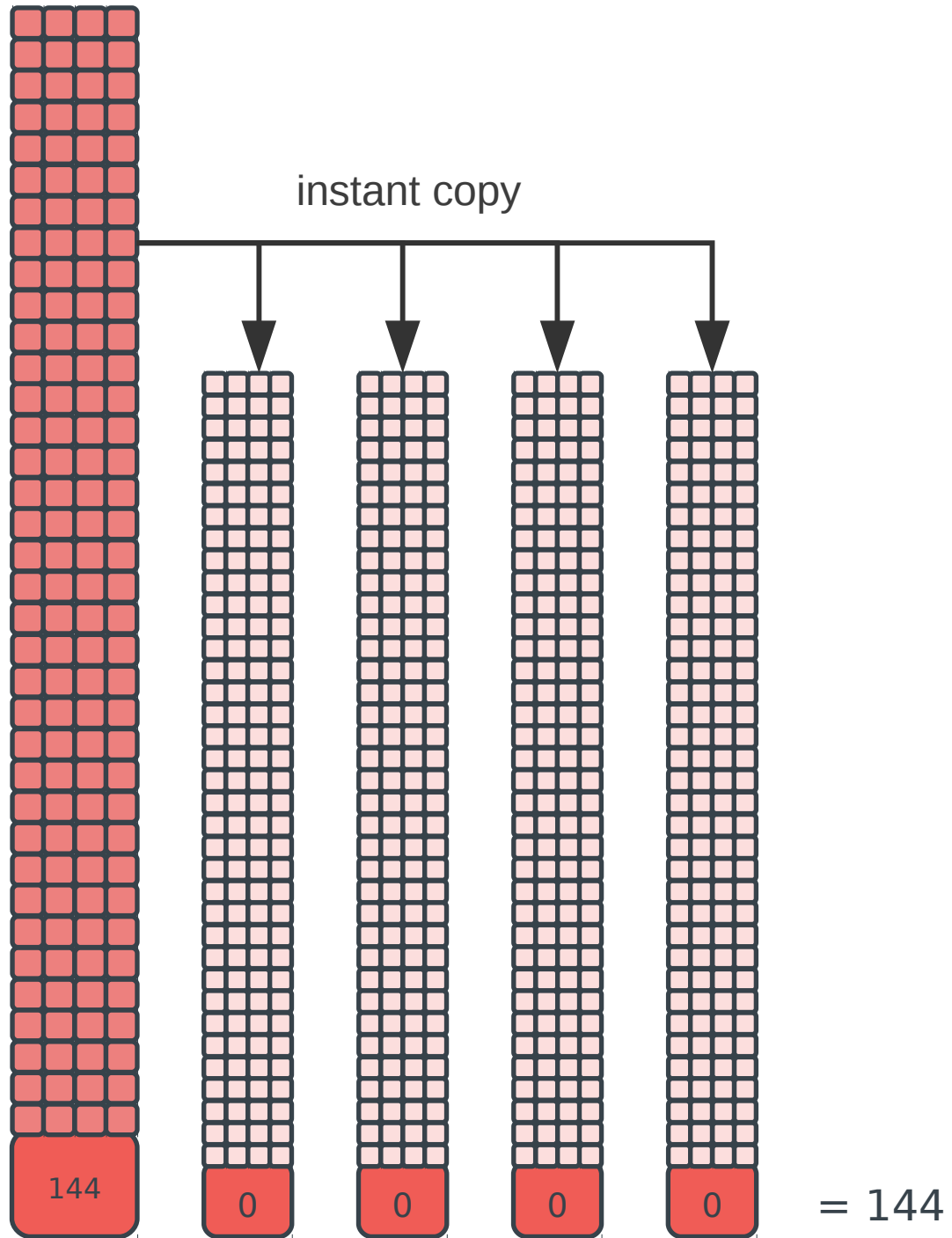


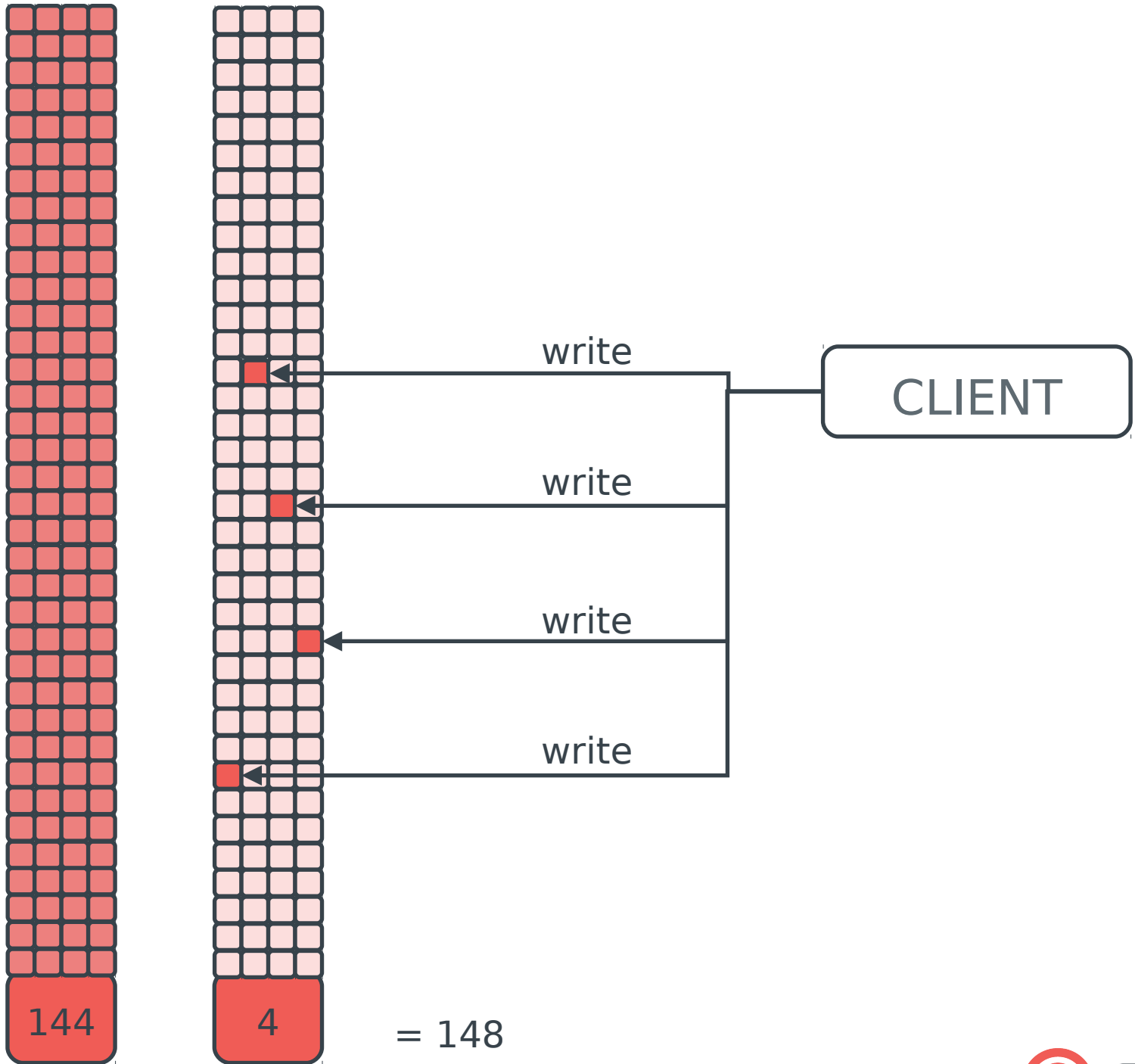


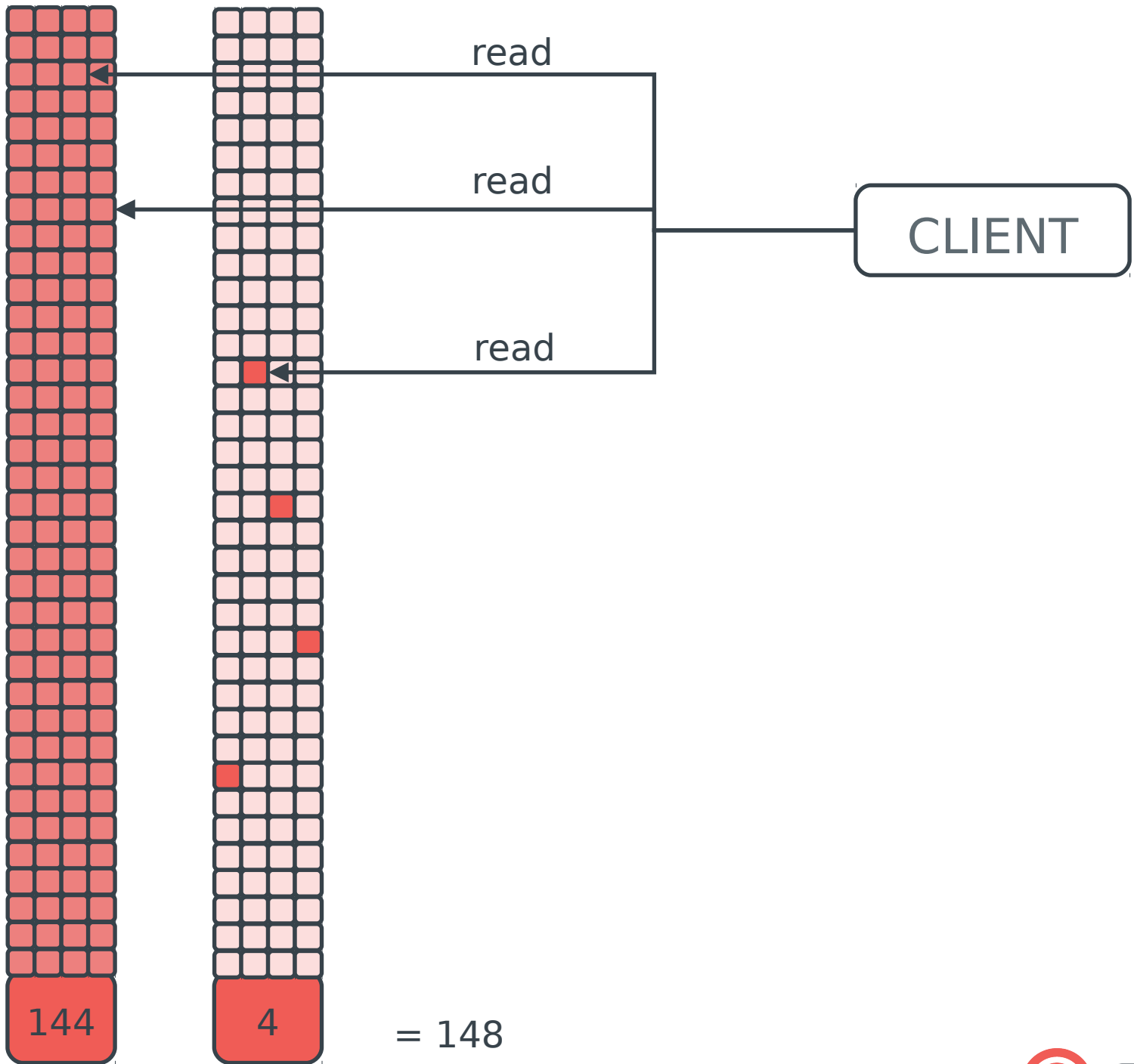
### RADOS Block Device:

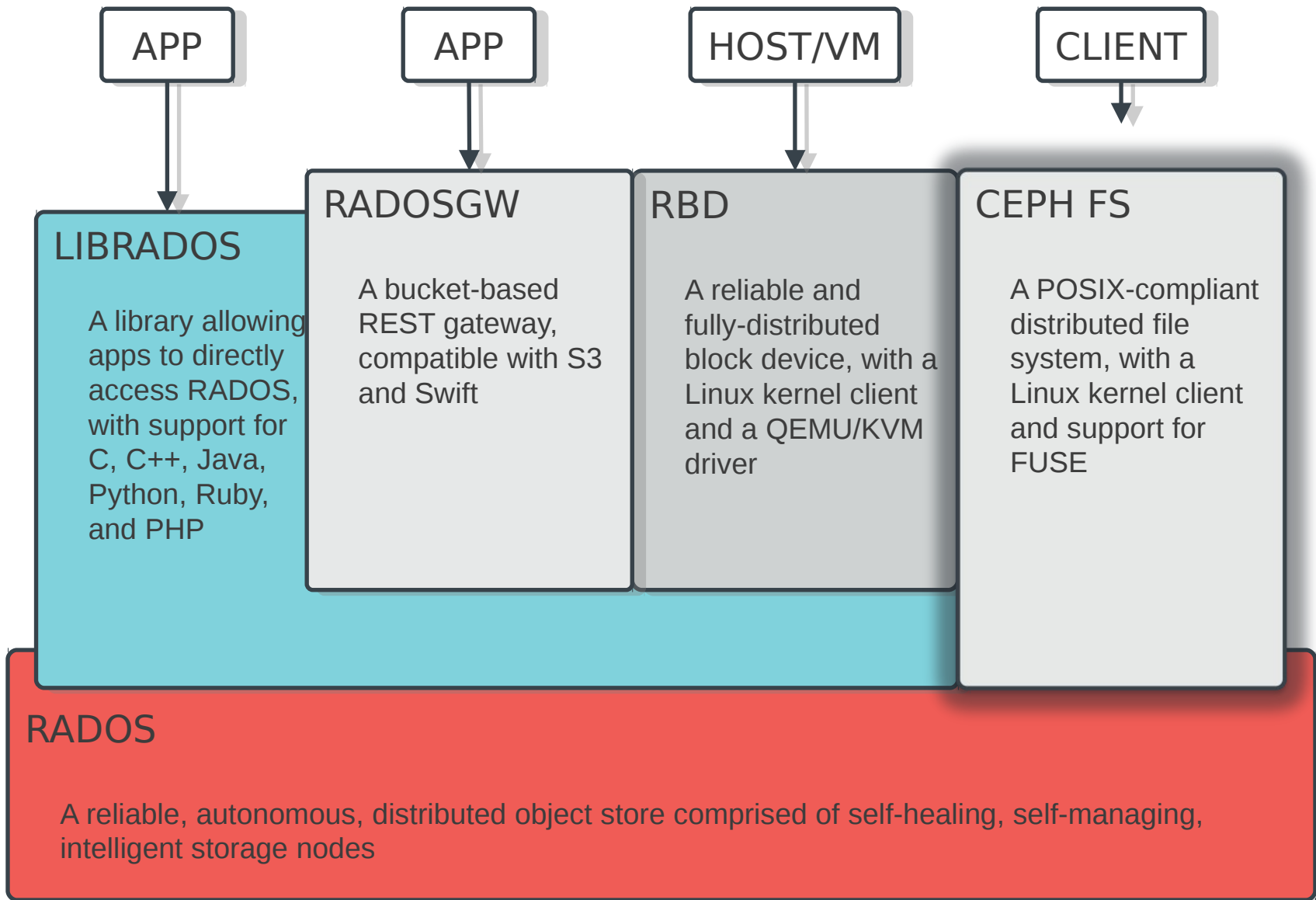
- Storage of virtual disks in RADOS
- Decouples VMs and containers
  - Live migration!
- Images are striped across the cluster
- Snapshots!
- Support in
  - Qemu/KVM
  - OpenStack, CloudStack
  - Mainline Linux kernel

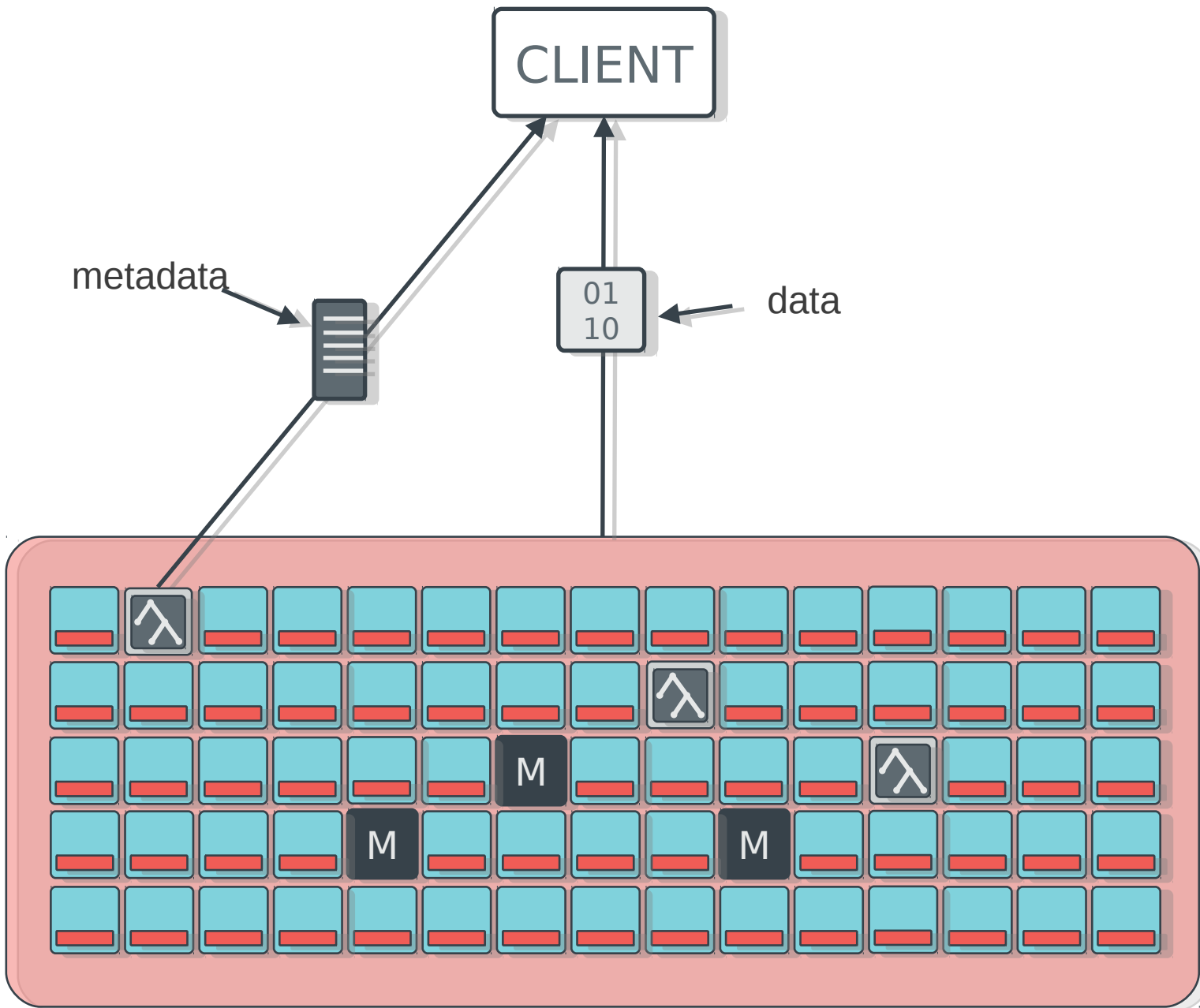
HOW DO YOU  
SPIN UP  
THOUSANDS OF VMs  
**INSTANTLY**  
AND  
**EFFICIENTLY?**



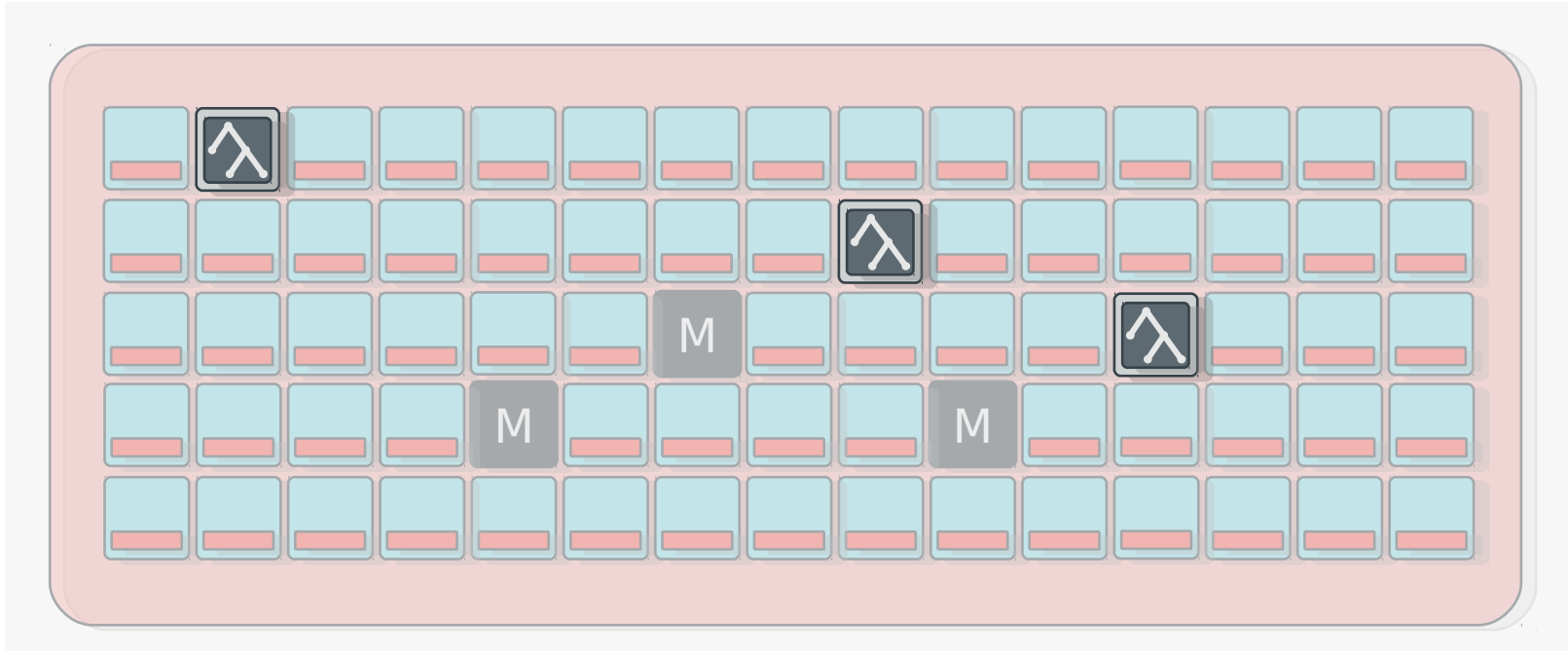


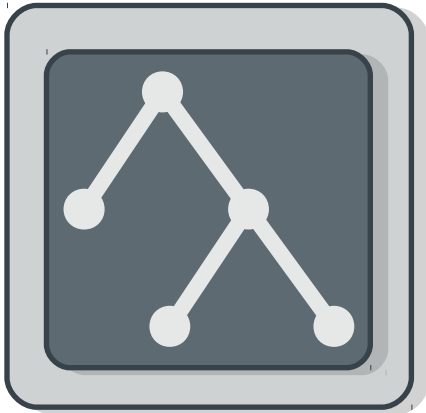






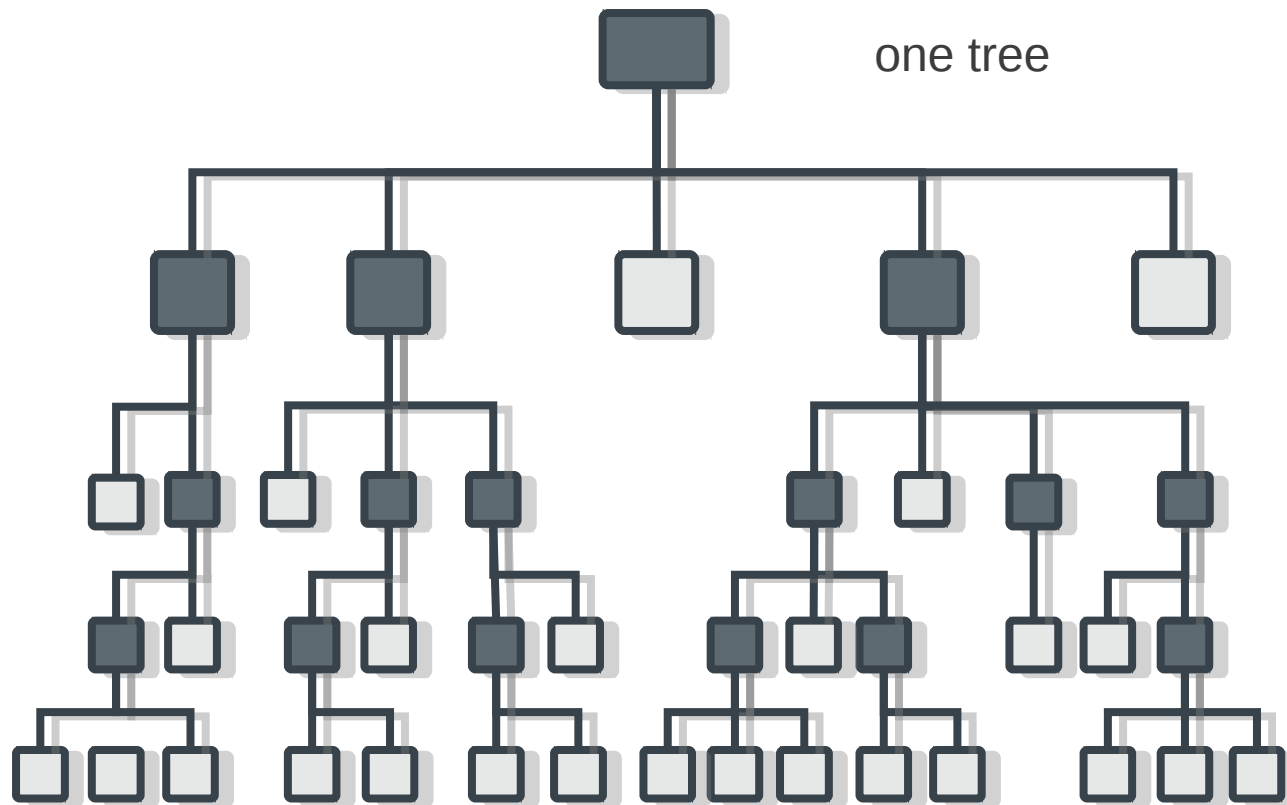






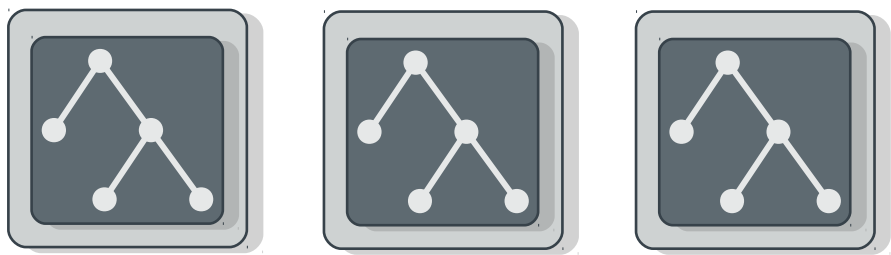
## Metadata Server

- Manages metadata for a POSIX-compliant shared filesystem
  - Directory hierarchy
  - File metadata (owner, timestamps, mode, etc.)
- Stores metadata in RADOS
- Does **not** serve file data to clients
- Only required for shared filesystem

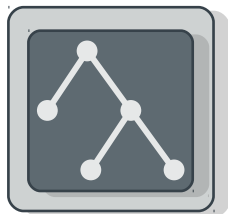
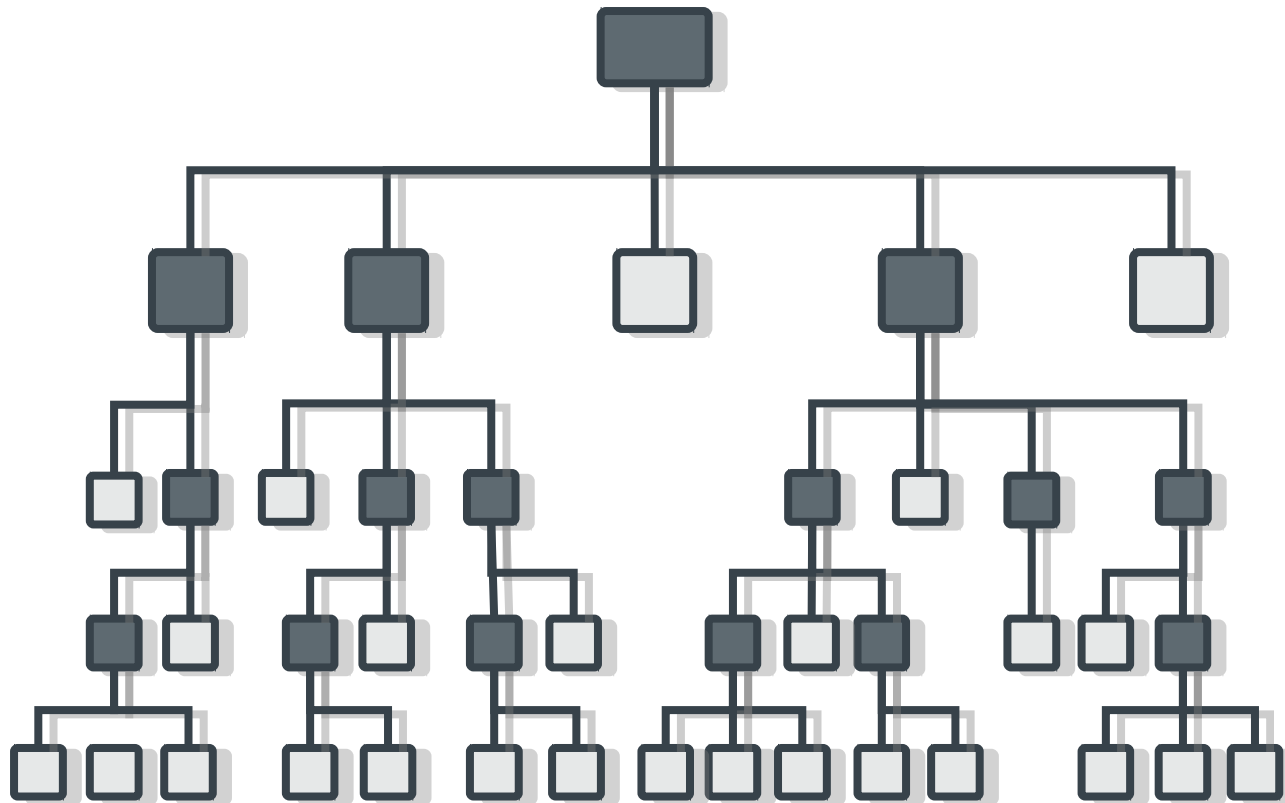


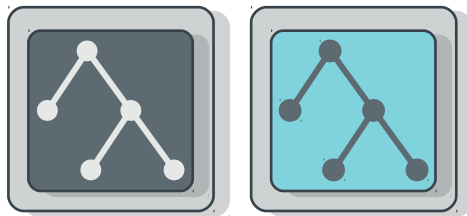
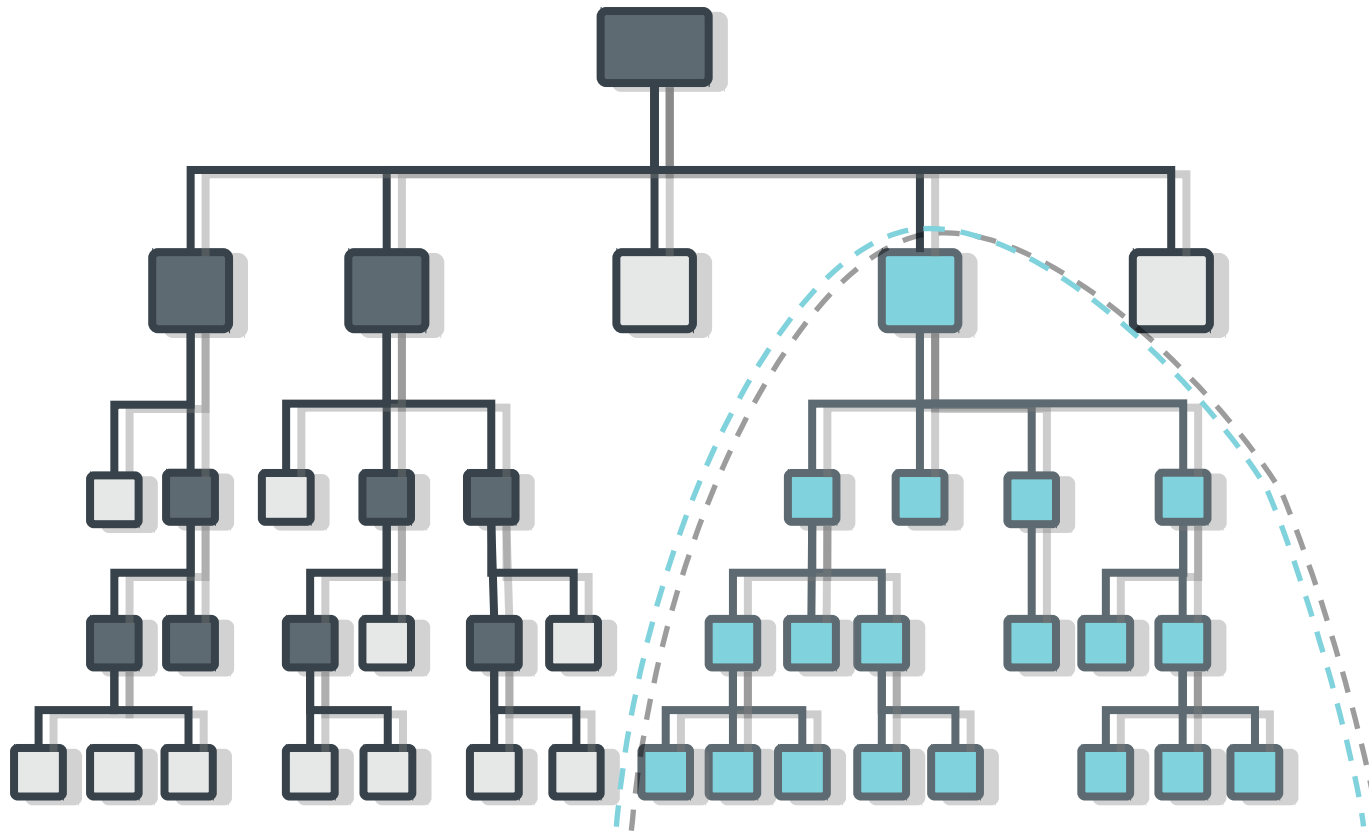
one tree

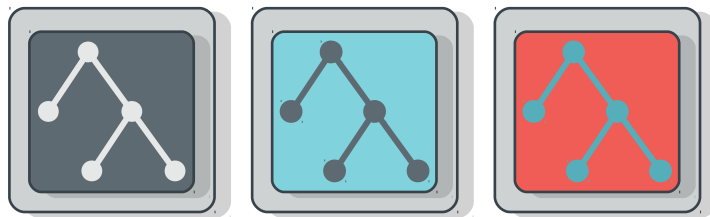
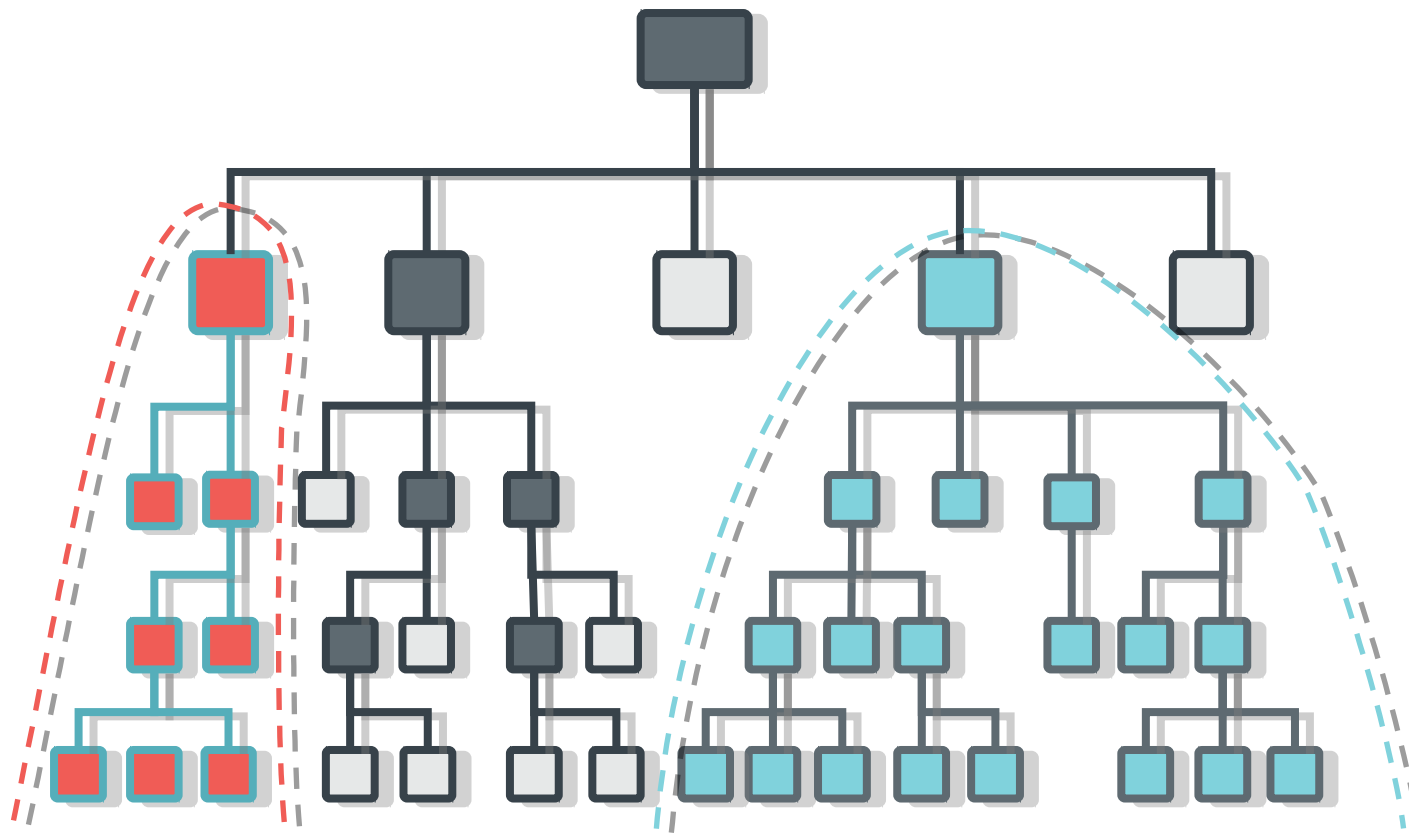
three metadata servers

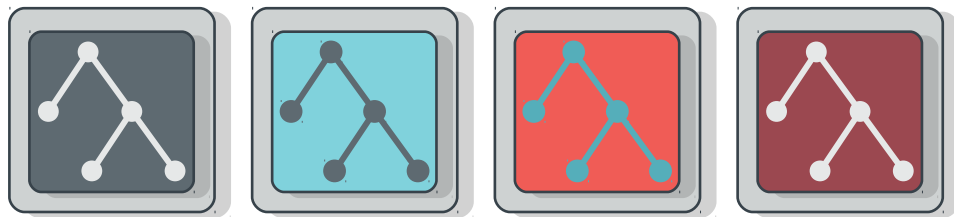
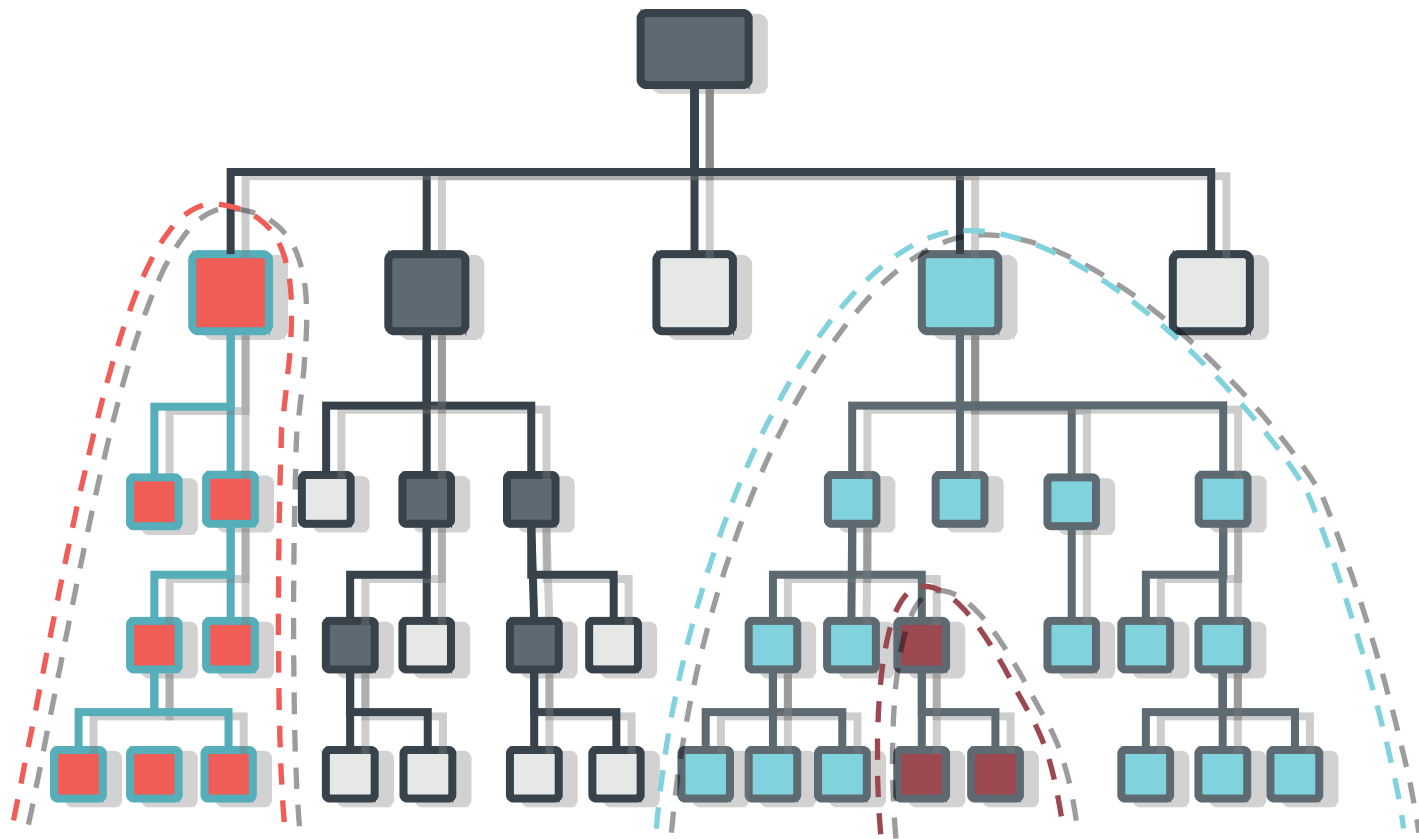


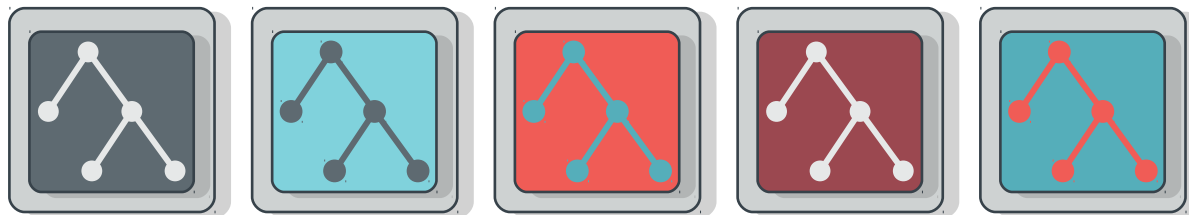
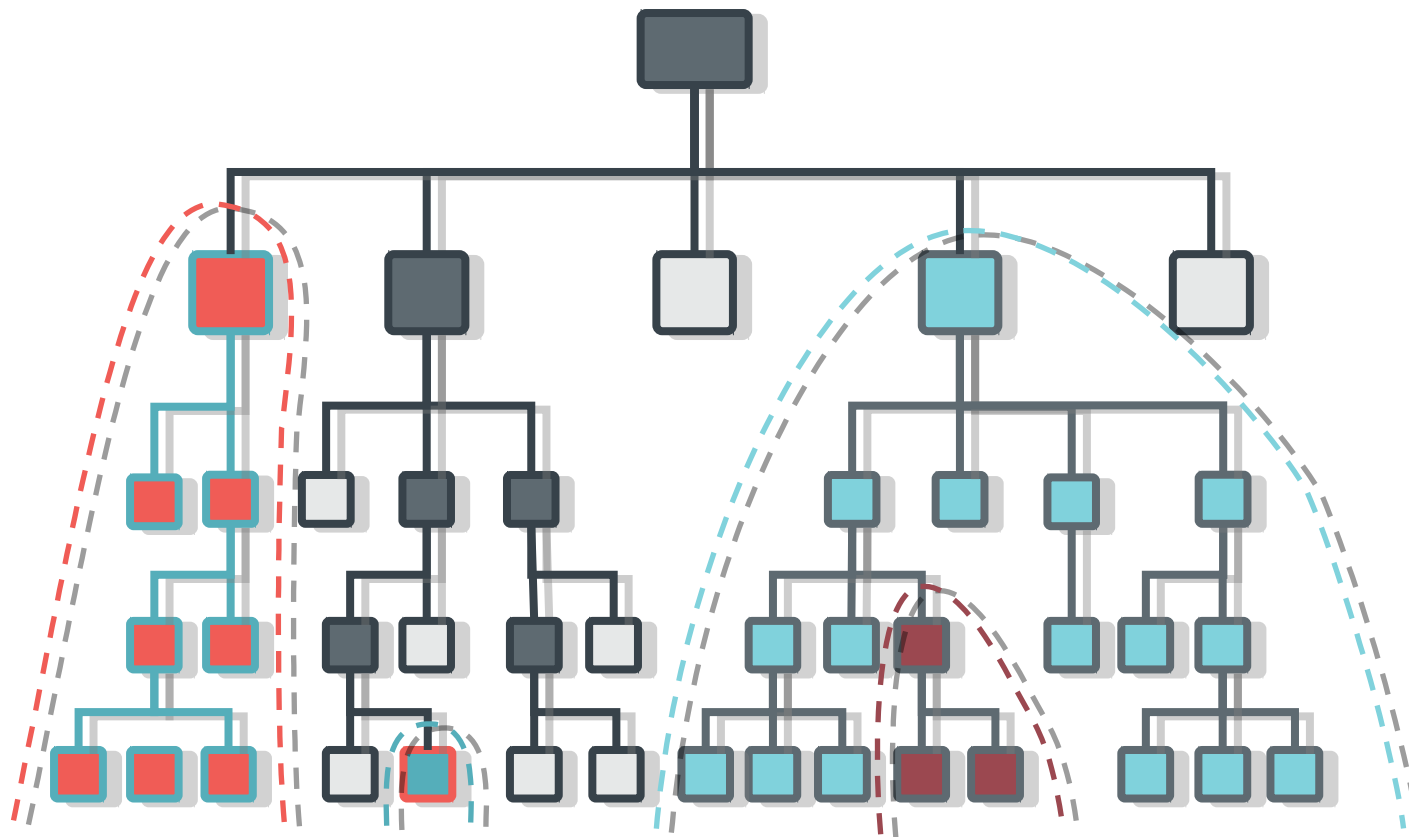
??











DYNAMIC SUBTREE PARTITIONING



# recursive accounting

- ceph-mds tracks recursive directory stats
  - file sizes
  - file and directory counts
  - modification time
- virtual xattrs present full stats
- efficient

```
$ ls -alSh | head
total 0
drwxr-xr-x 1 root          root          9.7T 2011-02-04 15:51 .
drwxr-xr-x 1 root          root          9.7T 2010-12-16 15:06 ..
drwxr-xr-x 1 pomceph      pg4194980    9.6T 2011-02-24 08:25 pomceph
drwxr-xr-x 1 mcg_test1    pg2419992    23G  2011-02-02 08:57 mcg_test1
drwx--x--- 1 luko        adm          19G  2011-01-21 12:17 luko
drwx--x--- 1 eest        adm          14G  2011-02-04 16:29 eest
drwxr-xr-x 1 mcg_test2    pg2419992    3.0G 2011-02-02 09:34 mcg_test2
drwx--x--- 1 fuzyceph     adm          1.5G 2011-01-18 10:46 fuzyceph
drwxr-xr-x 1 dallasceph   pg275        596M 2011-01-14 10:06 dallasceph
```

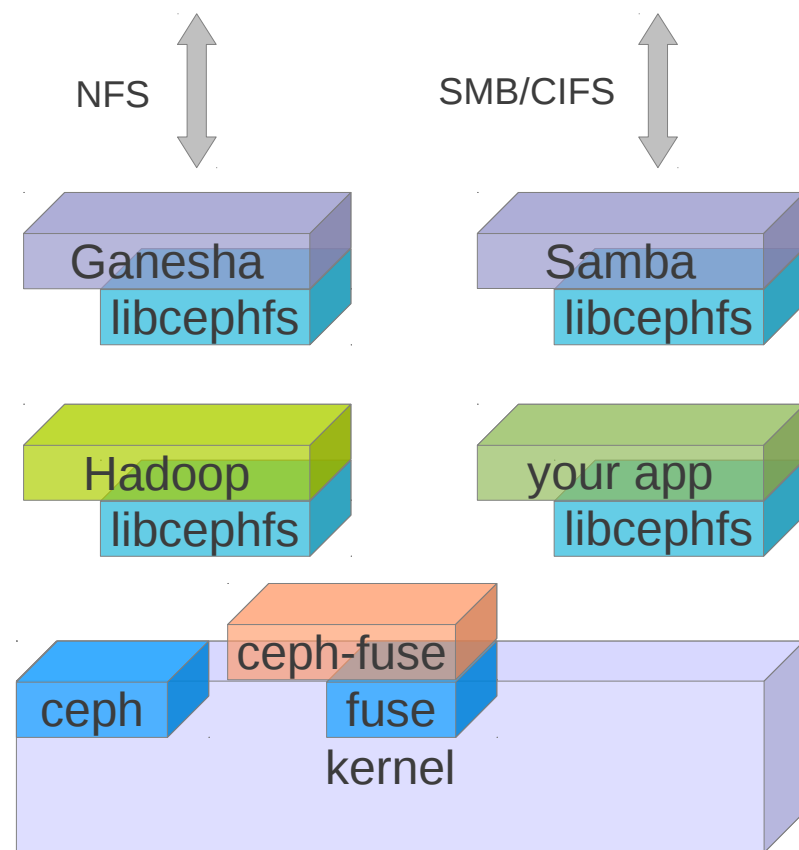
# snapshots

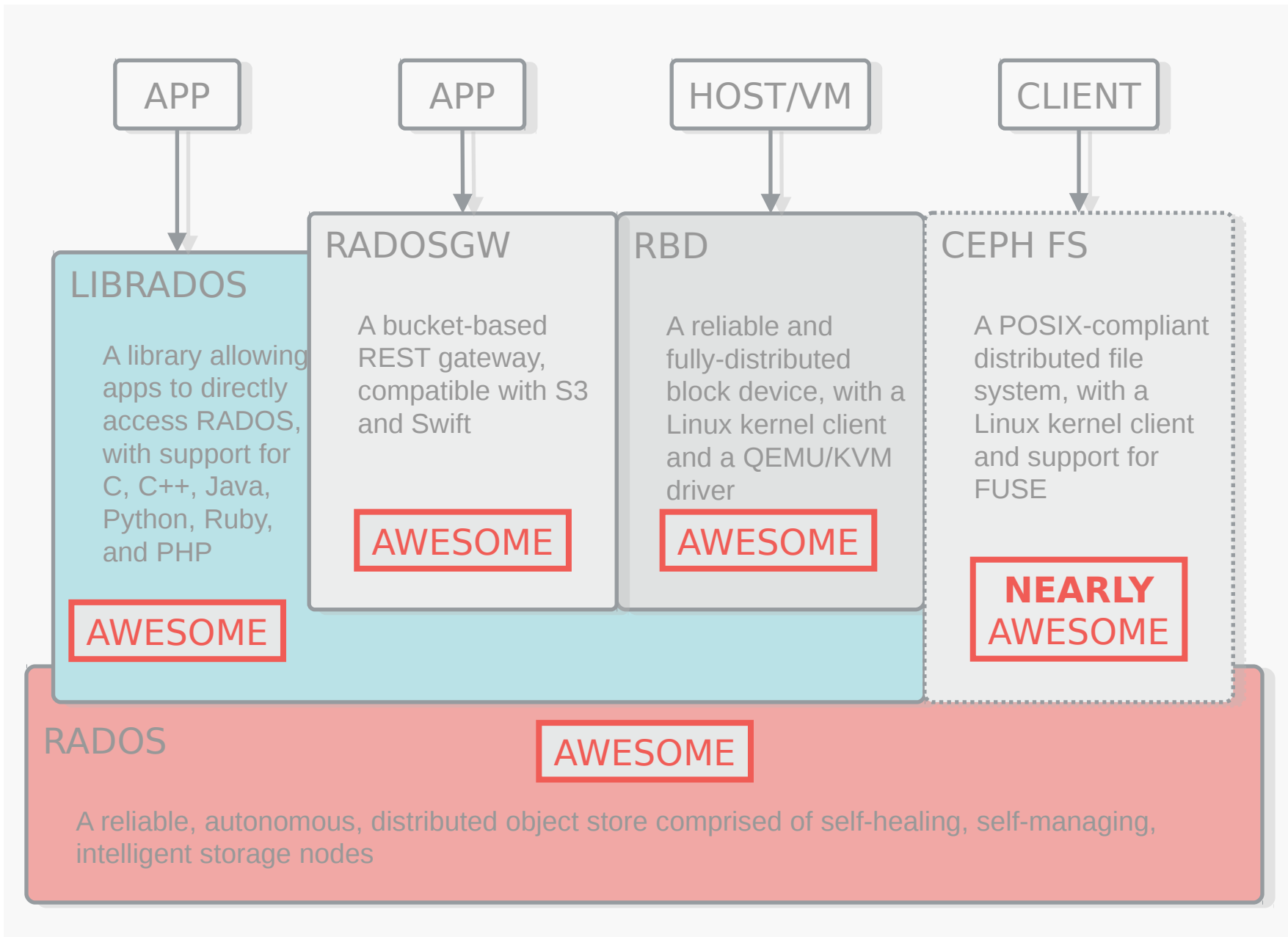
- volume or subvolume snapshots unusable at petabyte scale
  - snapshot arbitrary subdirectories
- simple interface
  - hidden '.snap' directory
  - no special tools

```
$ mkdir foo/.snap/one      # create snapshot
$ ls foo/.snap
one
$ ls foo/bar/.snap
_one_1099511627776      # parent's snap name is mangled
$ rm foo/myfile
$ ls -F foo
bar/
$ ls -F foo/.snap/one
myfile bar/
$ rmdir foo/.snap/one    # remove snapshot
```

# multiple protocols, implementations

- Linux kernel client
  - `mount -t ceph 1.2.3.4:/ /mnt`
  - export (NFS), Samba (CIFS)
- ceph-fuse
- libcephfs.so
  - your app
  - Samba (CIFS)
  - Ganesha (NFS)
  - Hadoop (map/reduce)





# why we do this

- limited options for scalable open source storage
- proprietary solutions
  - expensive
  - don't scale (well or out)
  - marry hardware and software
- industry needs to change

# who we are

- Ceph created at UC Santa Cruz (2007)
- supported by DreamHost (2008-2011)
- Inktank (2012)
  - Los Angeles, Sunnyvale, San Francisco, remote
- growing user and developer community
  - Linux distros, users, cloud stacks, SIs, OEMs

<http://ceph.com/>

# thanks

BoF tonight @ 5:15

sage weil

[sage@inktank.com](mailto:sage@inktank.com)

@liewegas

<http://github.com/ceph>

<http://ceph.com/>





# why we like btrfs

- pervasive checksumming
- snapshots, **copy-on-write**
- efficient metadata (xattrs)
- inline data for small files
- transparent compression
- integrated volume management
  - software RAID, mirroring, error recovery
  - SSD-aware
- online fsck
- active development community