

# Ceph Intro & Architectural Overview



Abbas Bangash  
Intercloud Systems

# About Me

Abbas Bangash  
Systems Team Lead,  
Intercloud Systems



[abangash@intercloudsys.com](mailto:abangash@intercloudsys.com)

[intercloudsys.com](http://intercloudsys.com)

# CLOUD SERVICES

COMPUTE

NETWORK

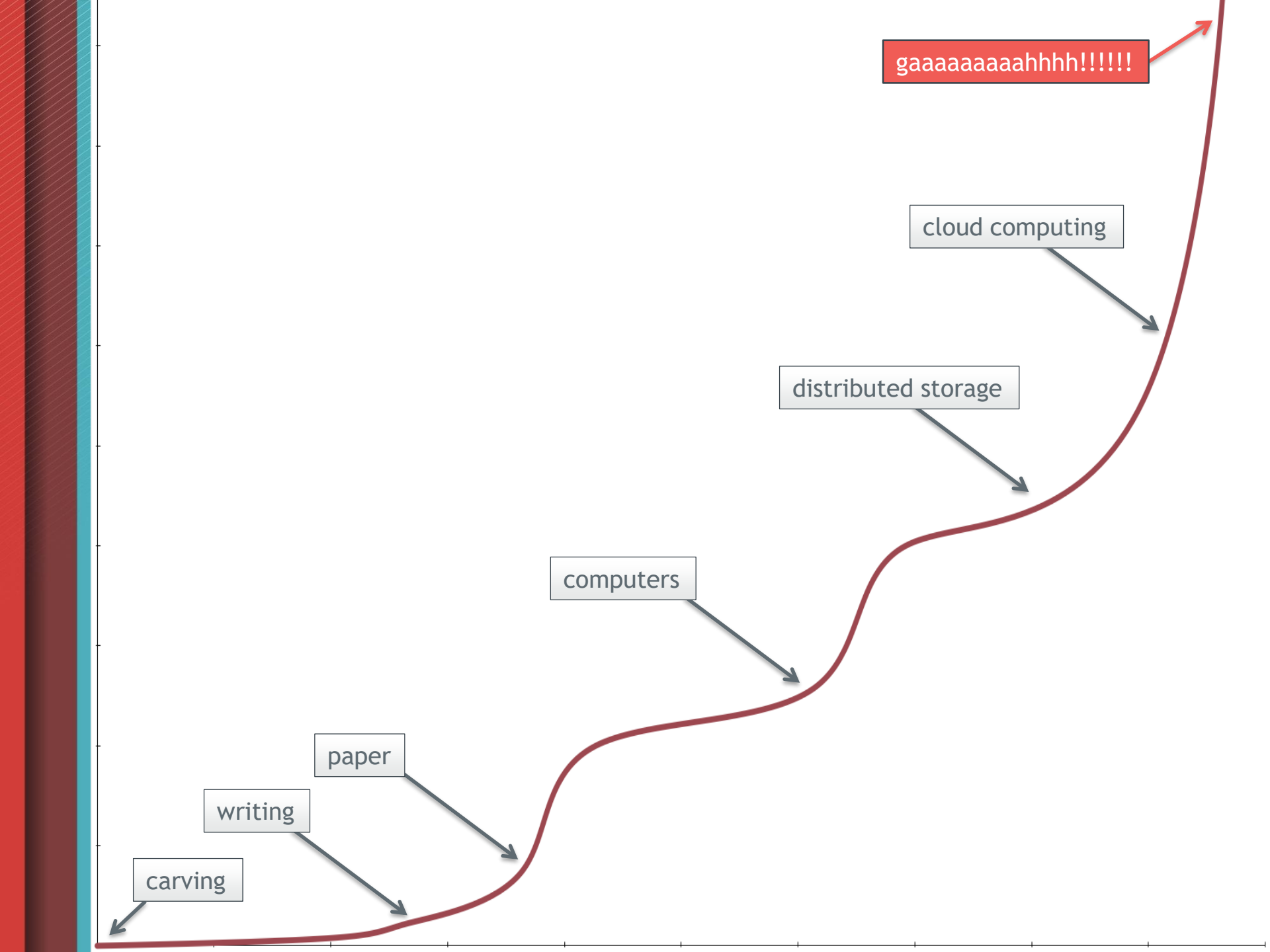
STORAGE

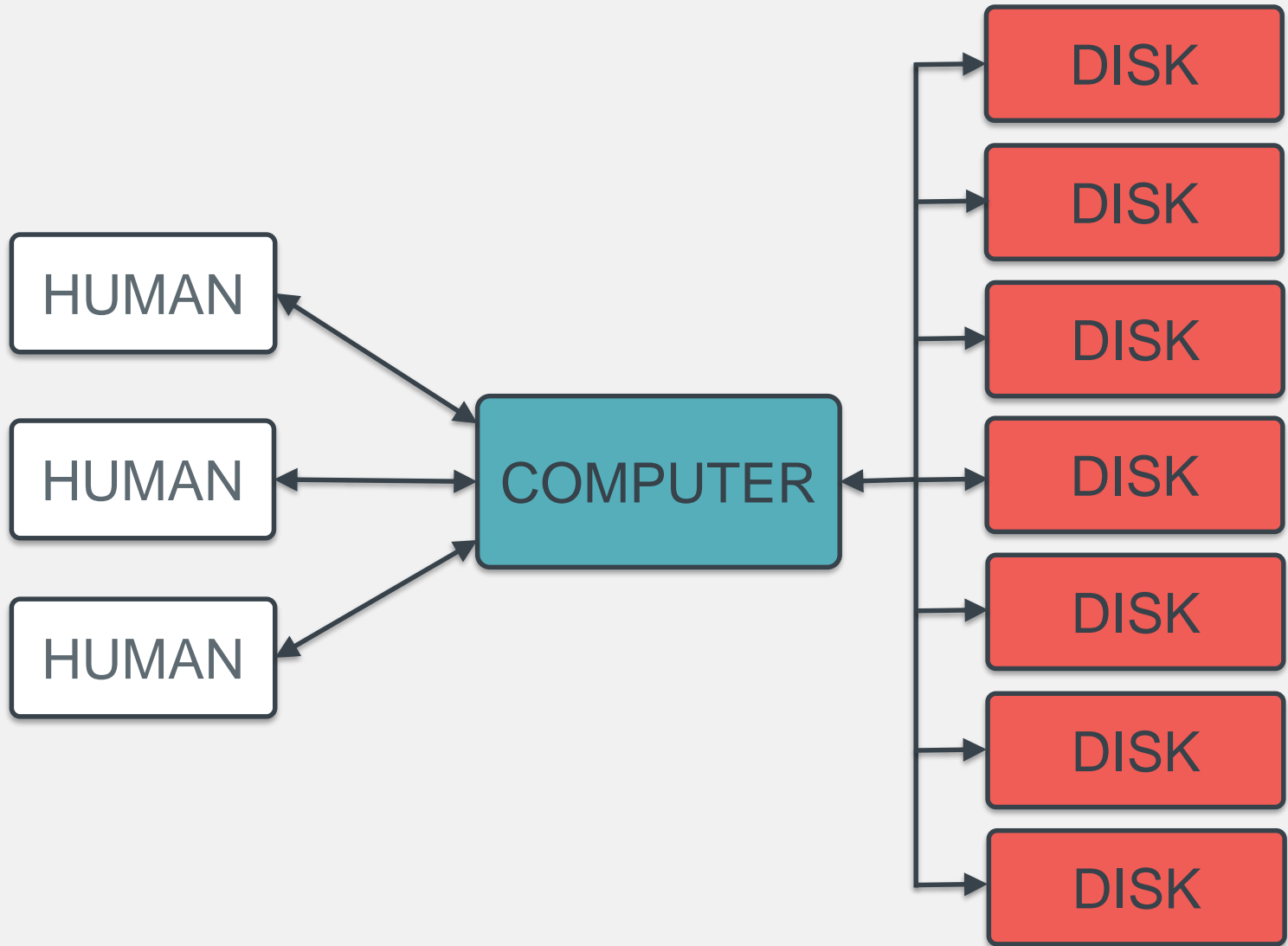


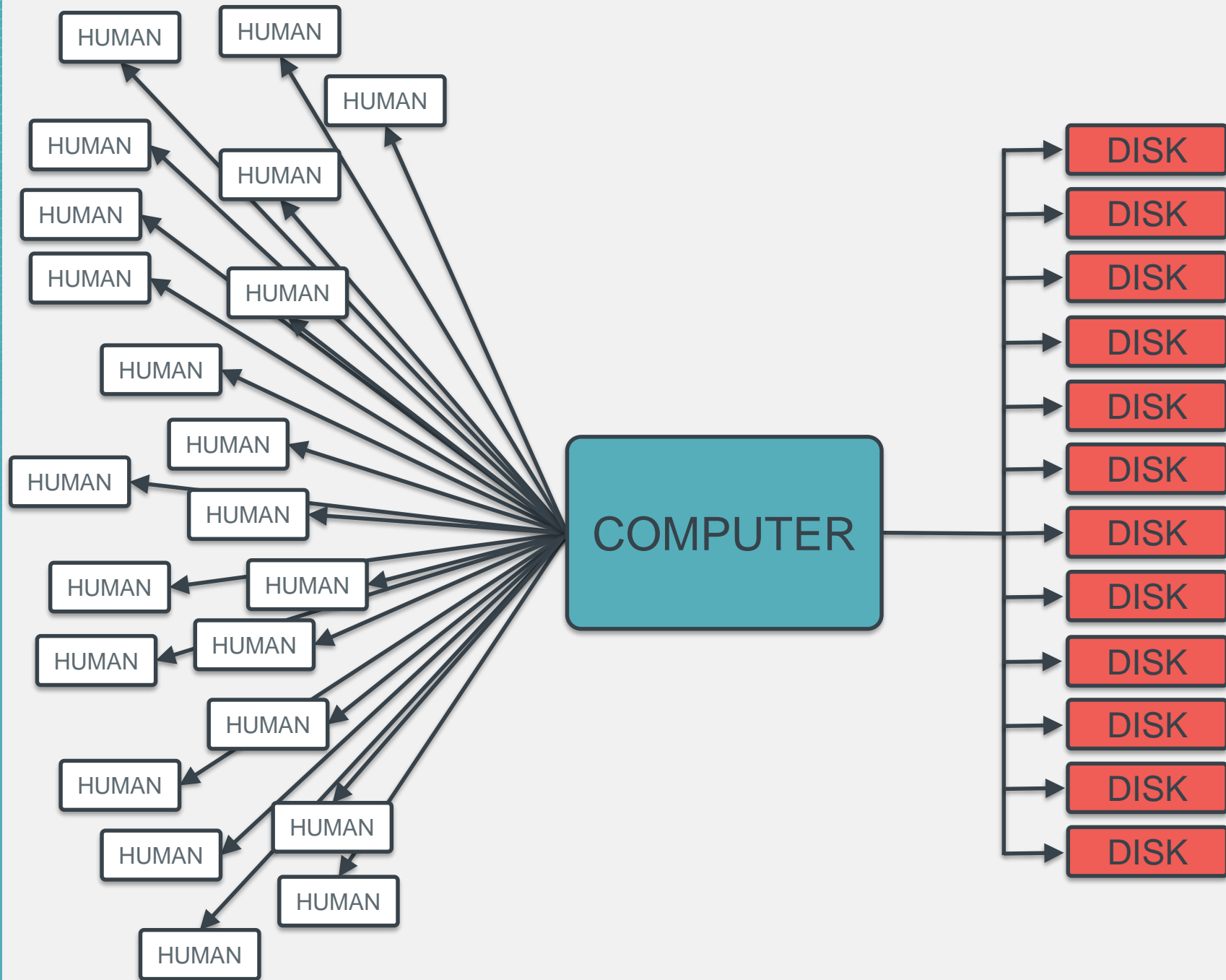
the future of storage™



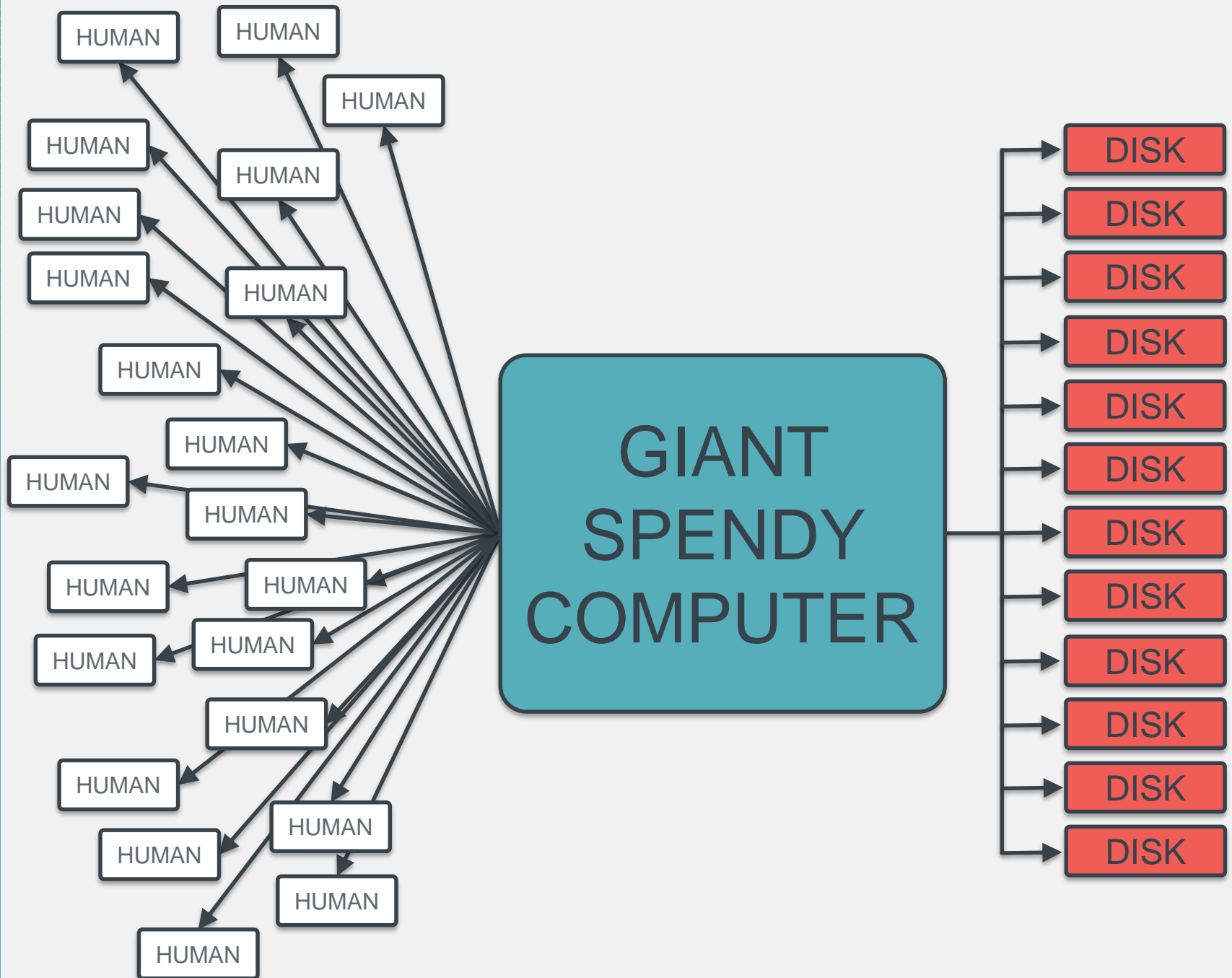


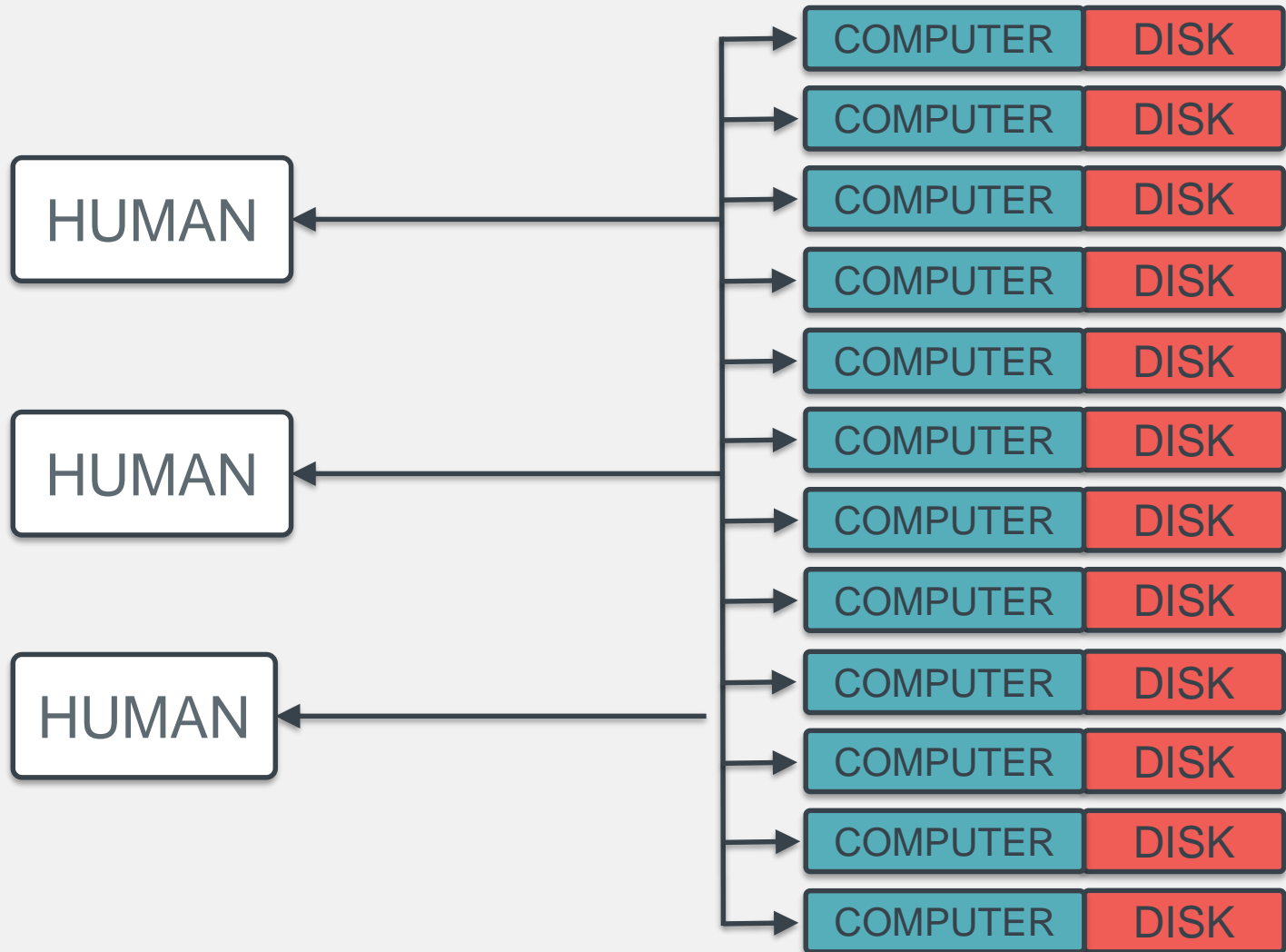


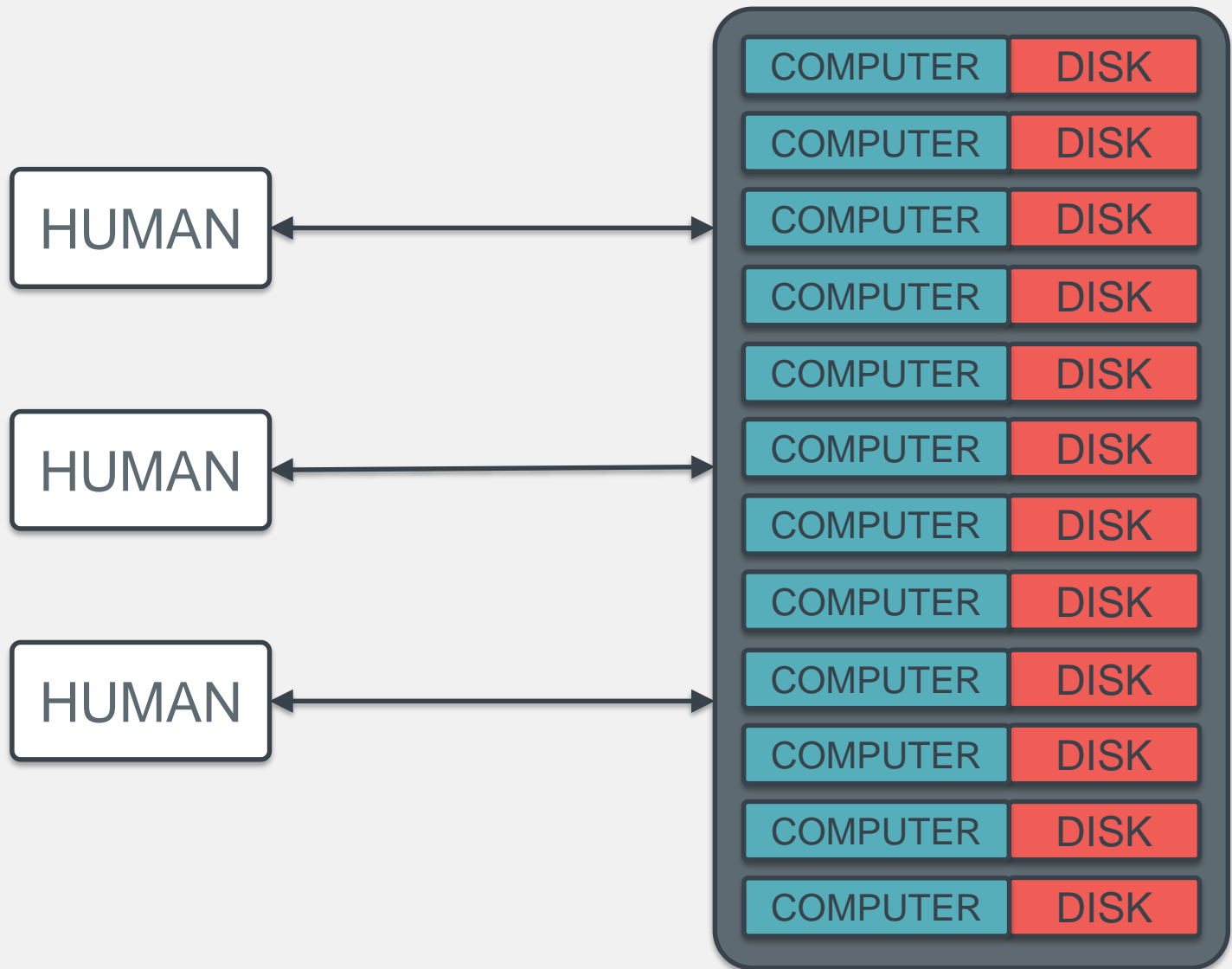




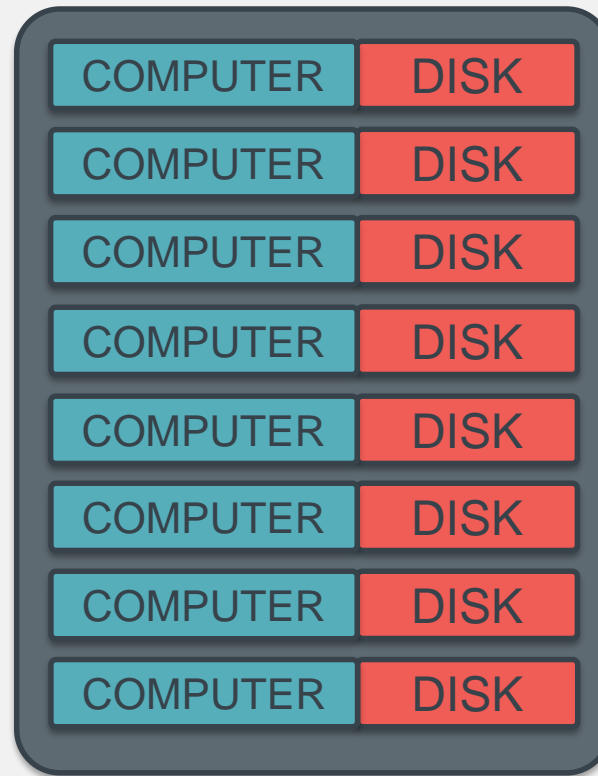






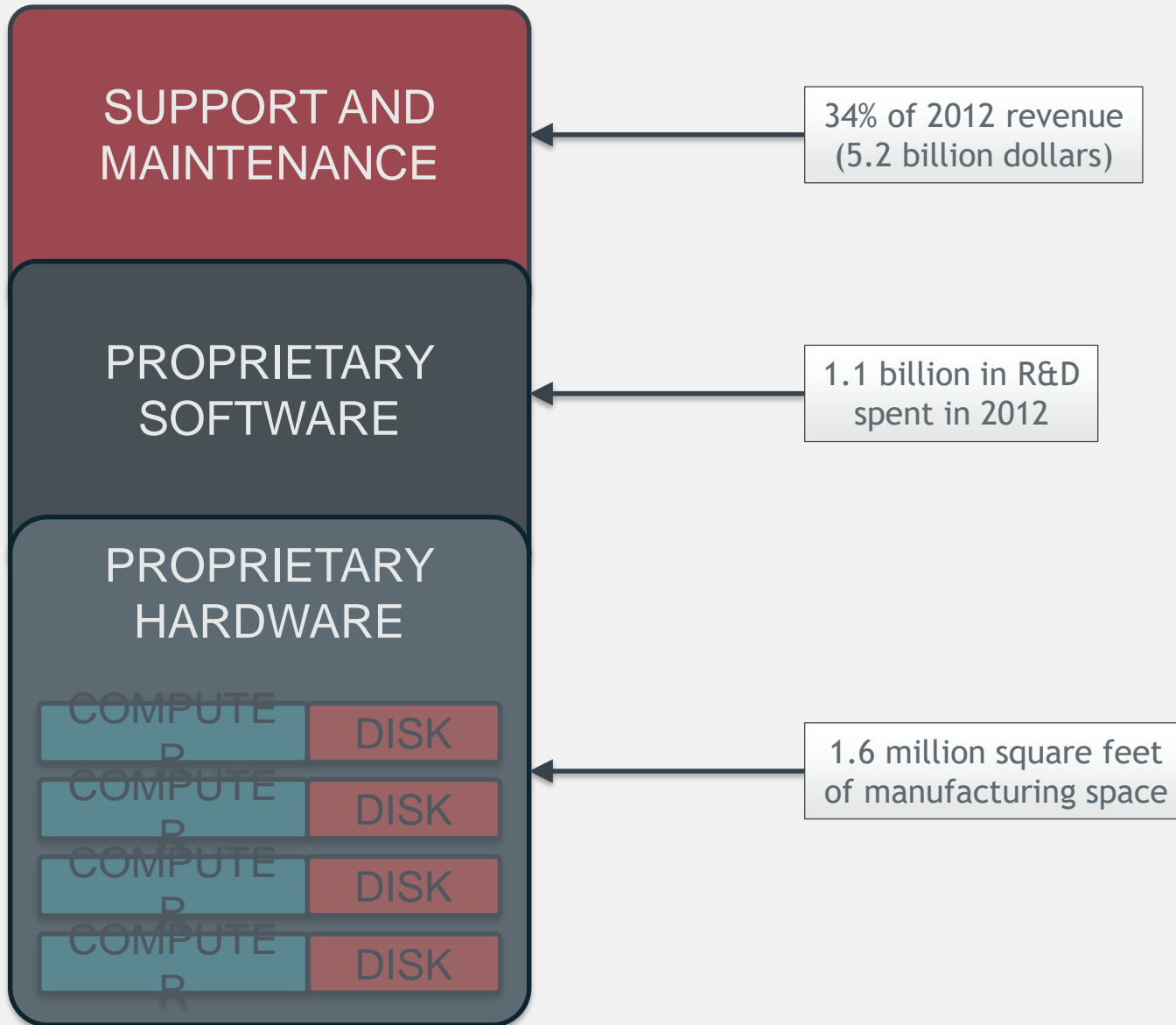


# “STORAGE APPLIANCE”



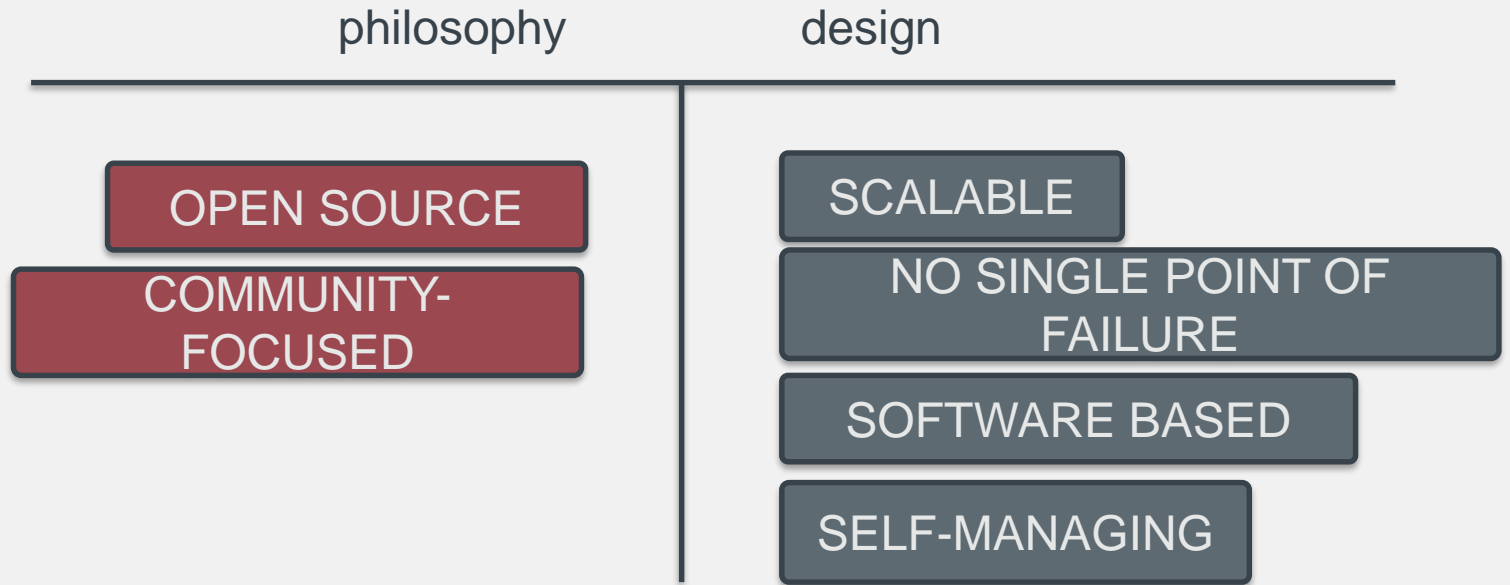


Storage Appliance

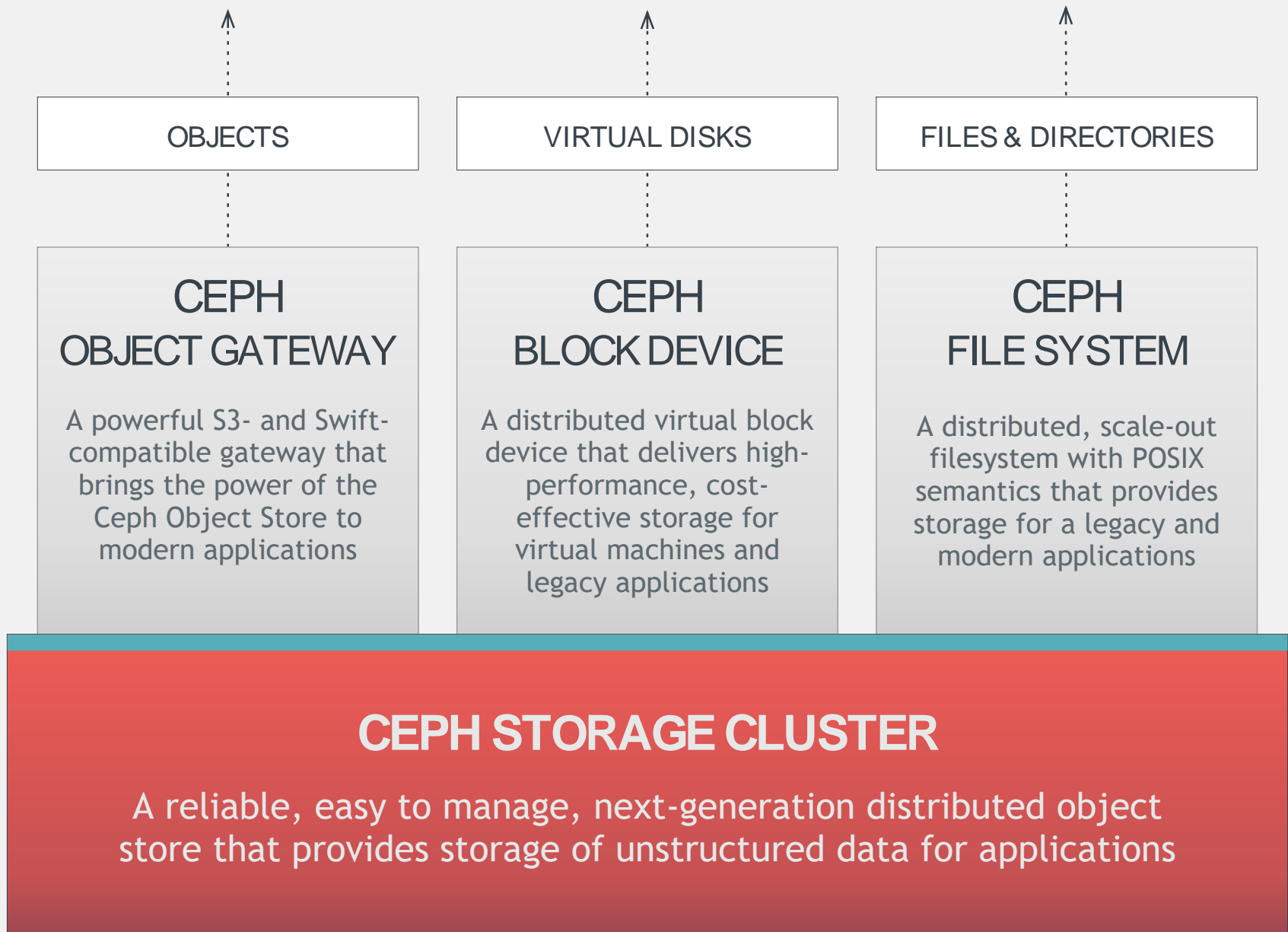


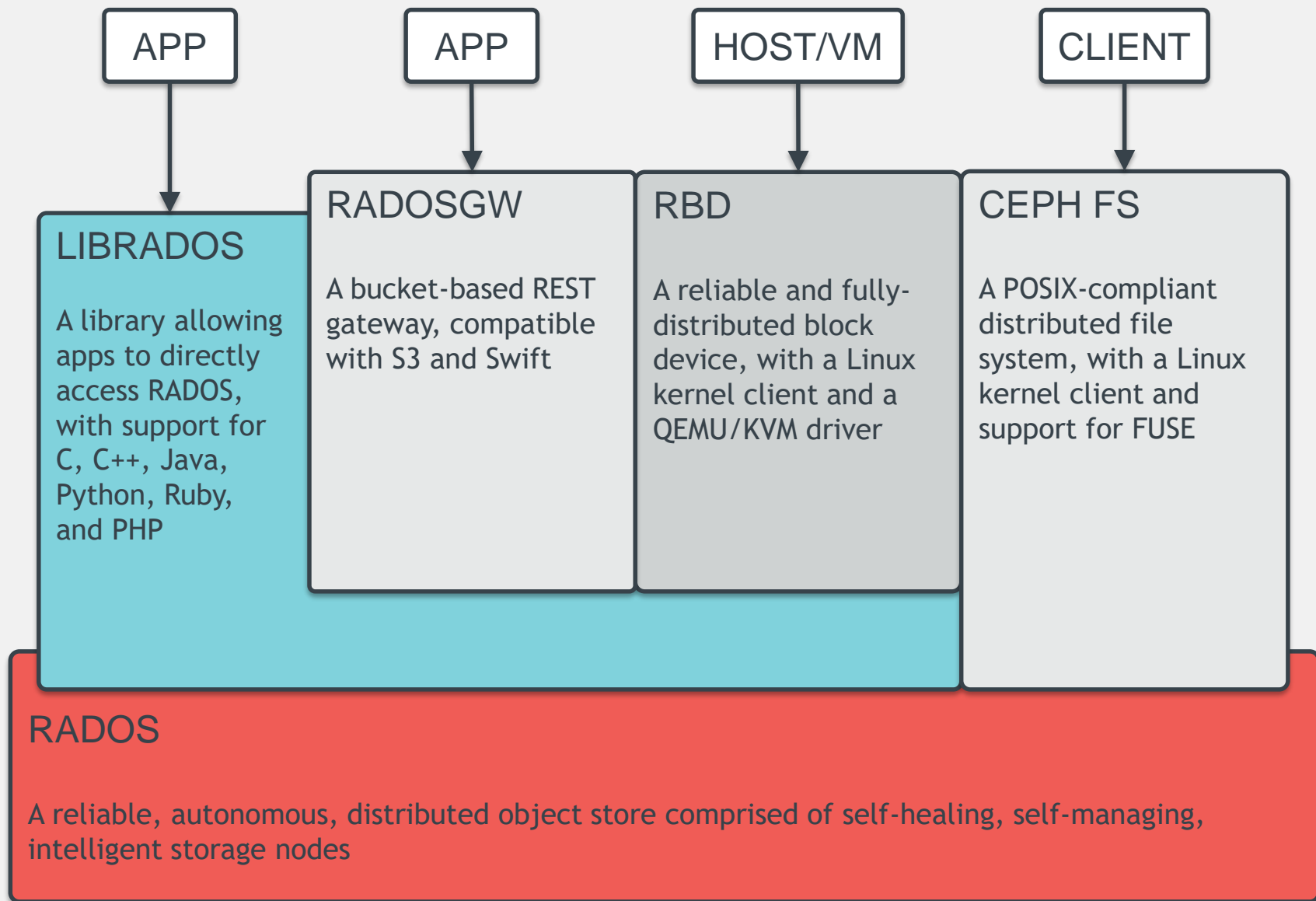


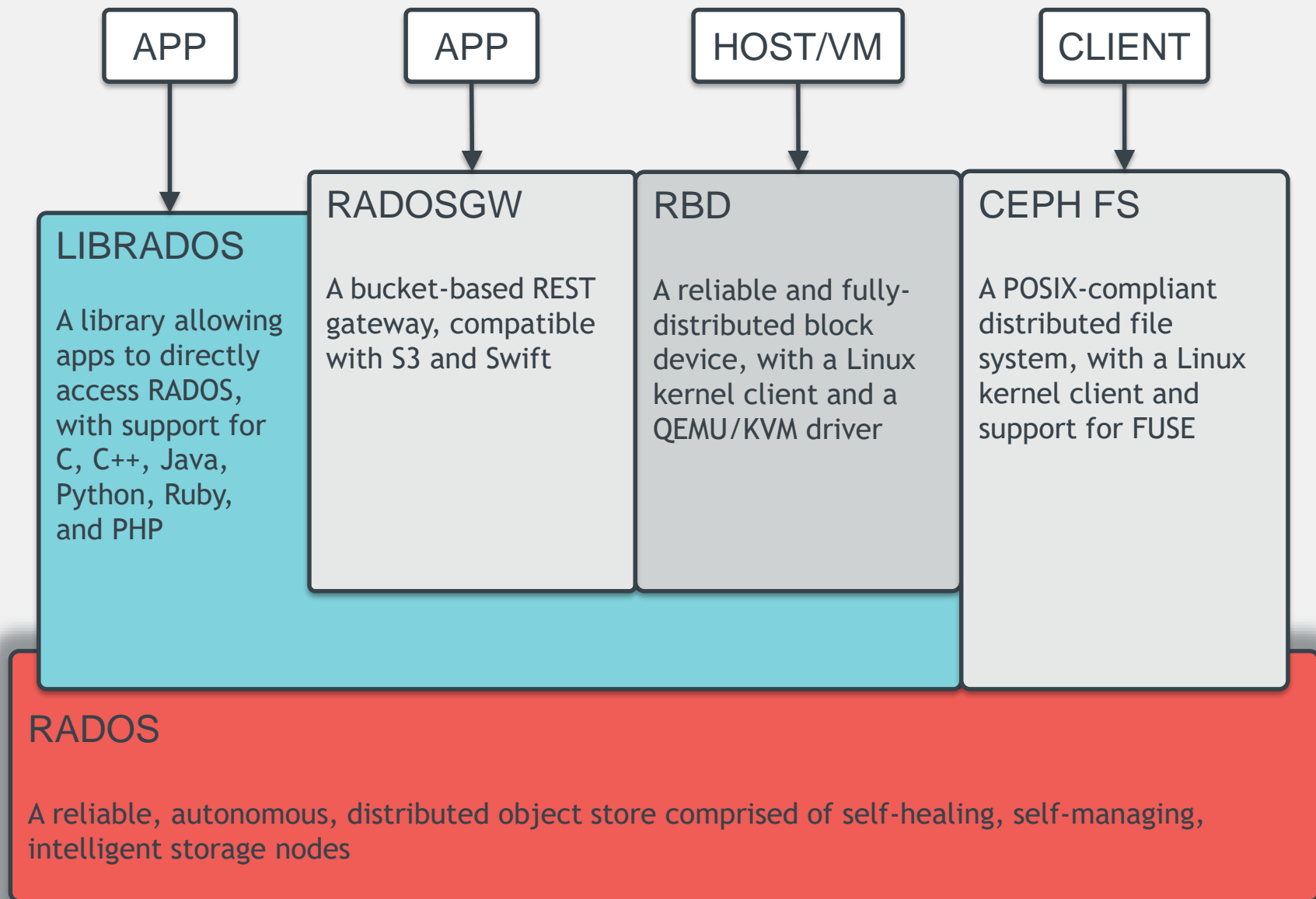
ceph

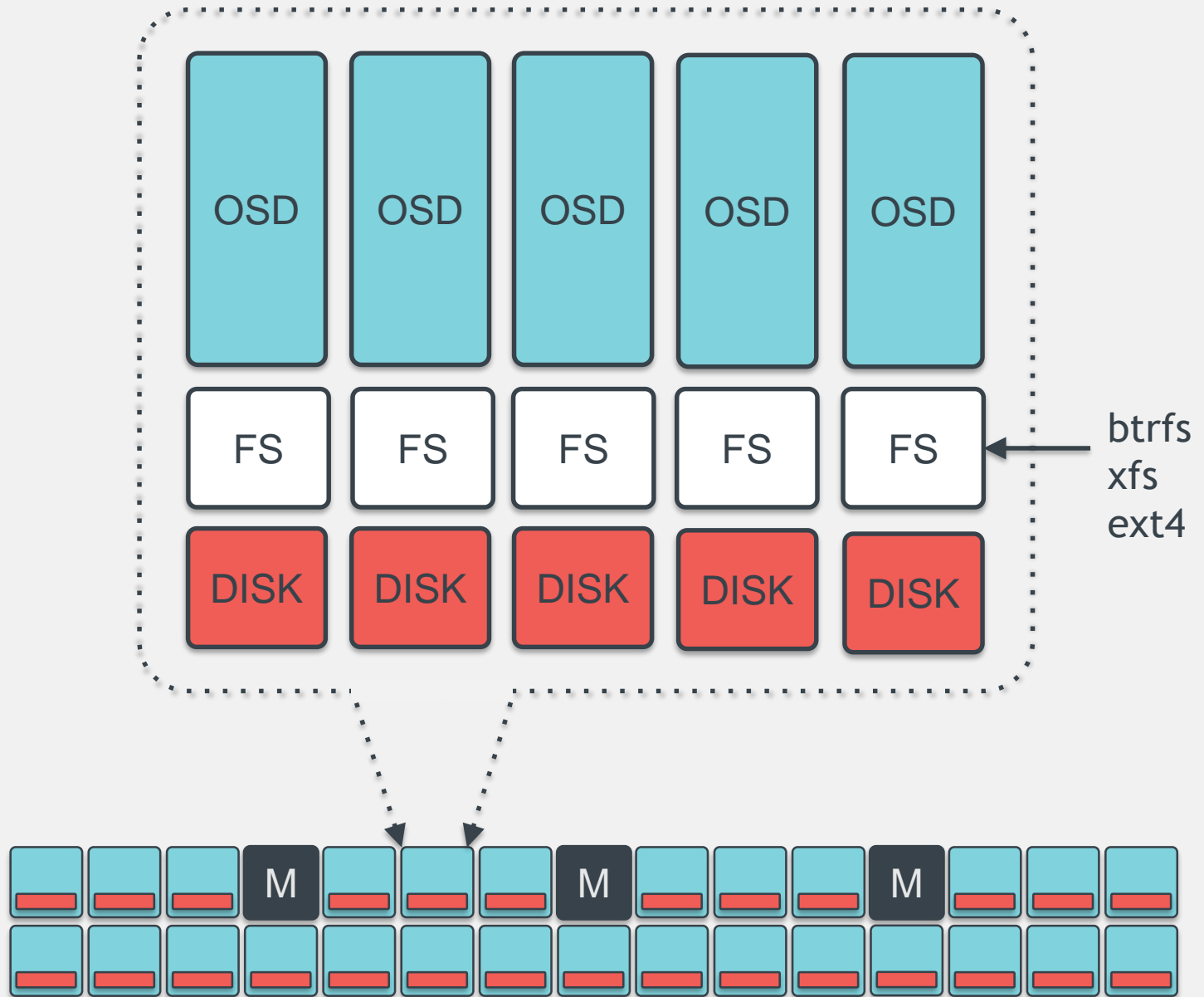




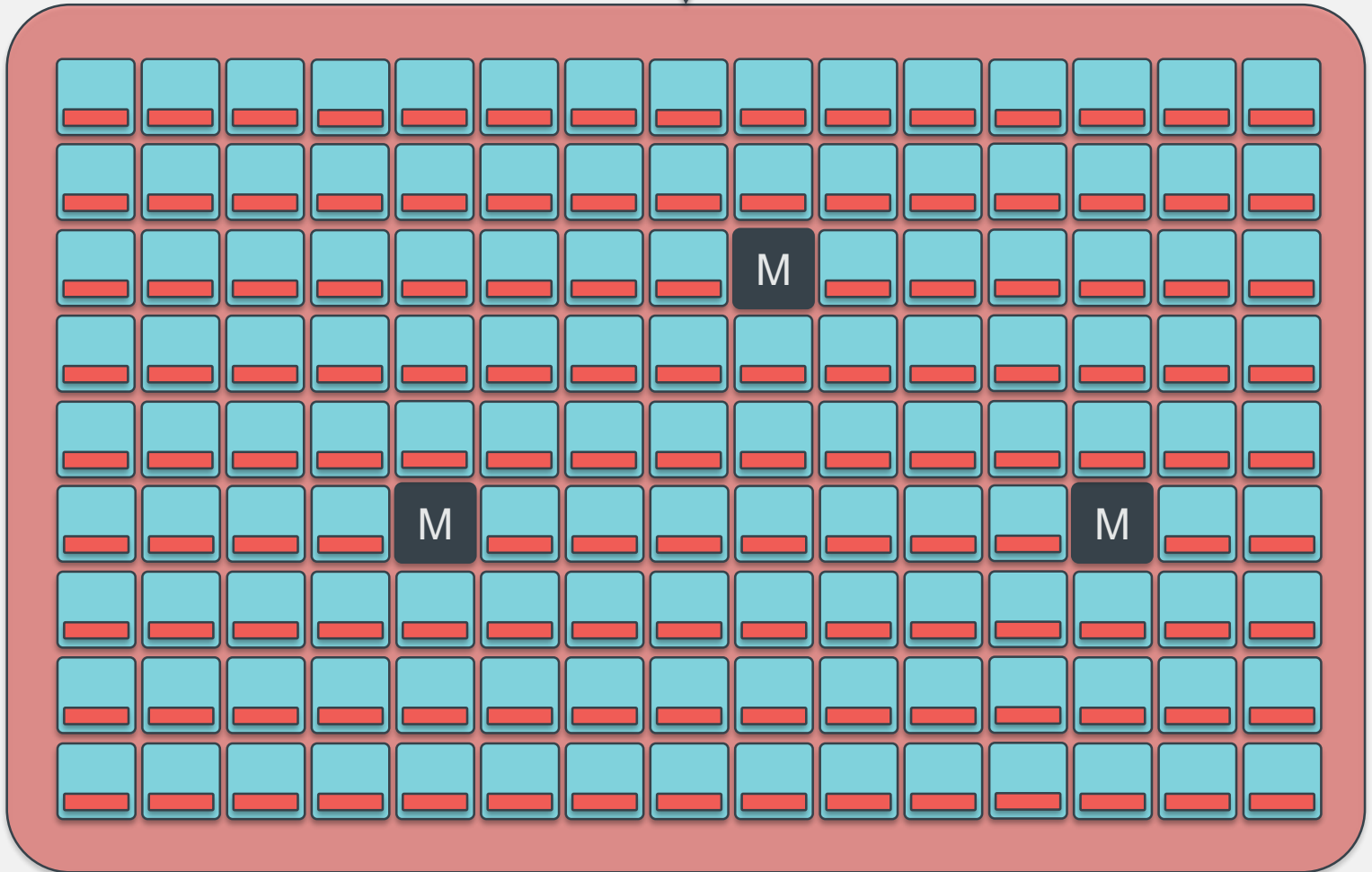
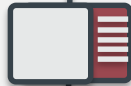


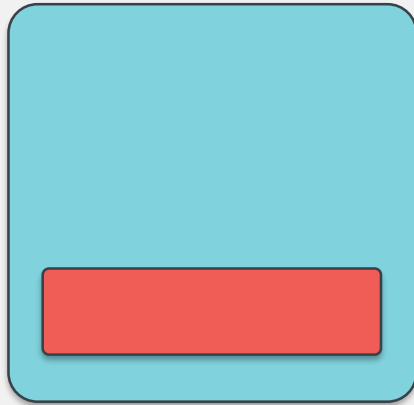






HUMAN





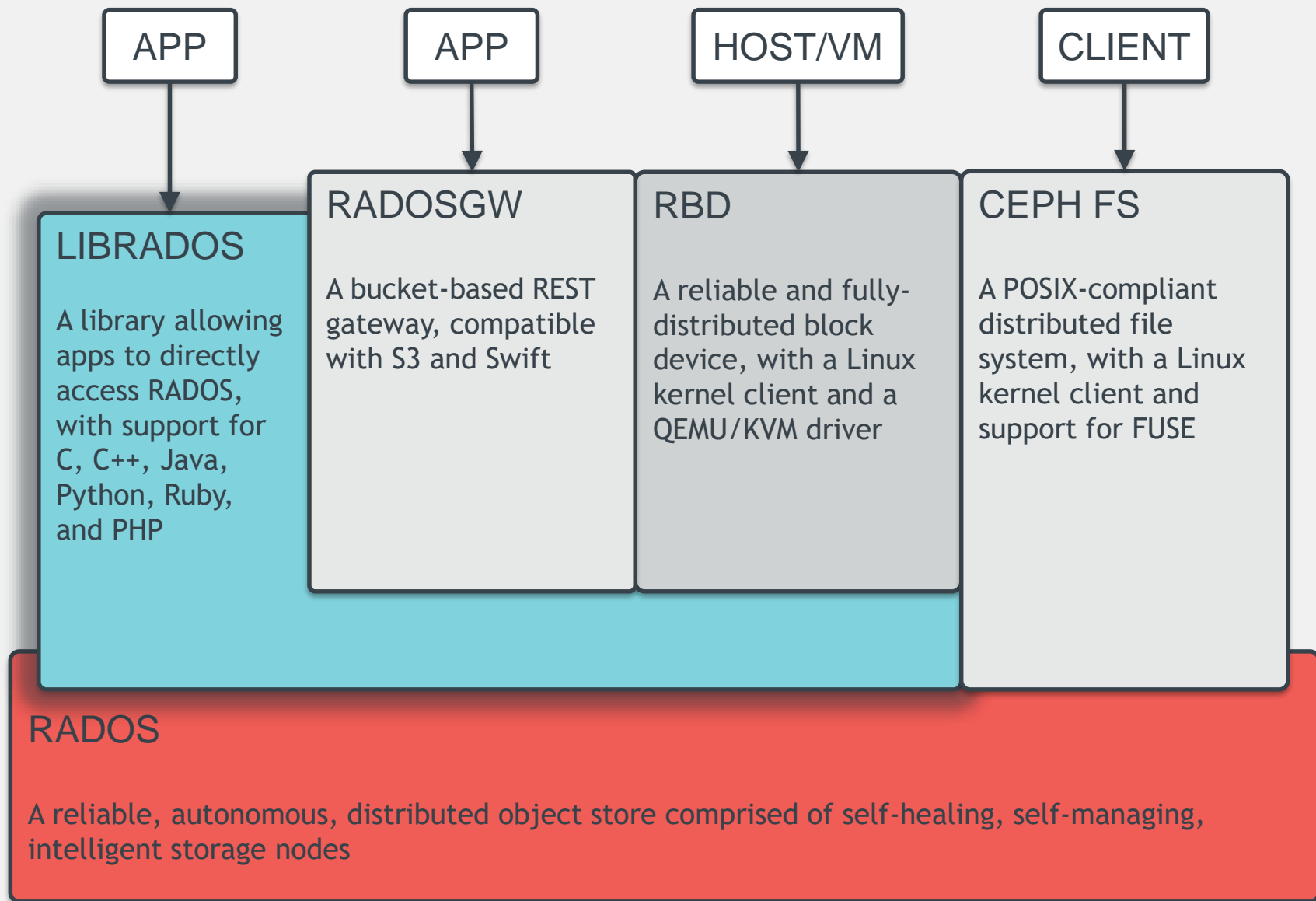
### OSDs:

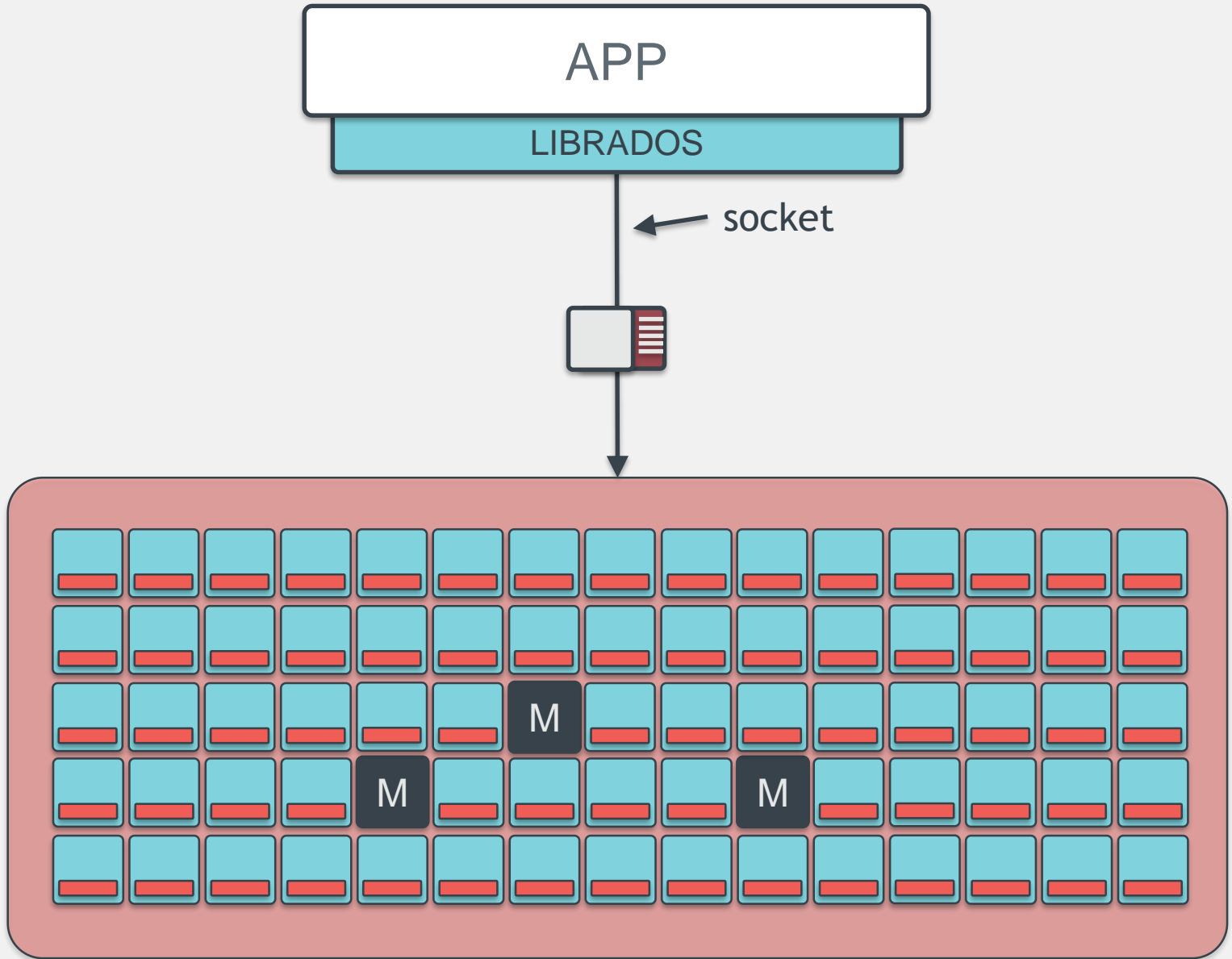
- 10s to 10000s in a cluster
- One per disk
  - (or one per SSD, RAID group...)
- Serve stored objects to clients
- Intelligently peer to perform replication and recovery tasks



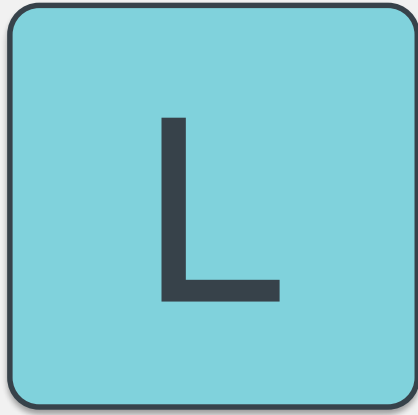
### Monitors:

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



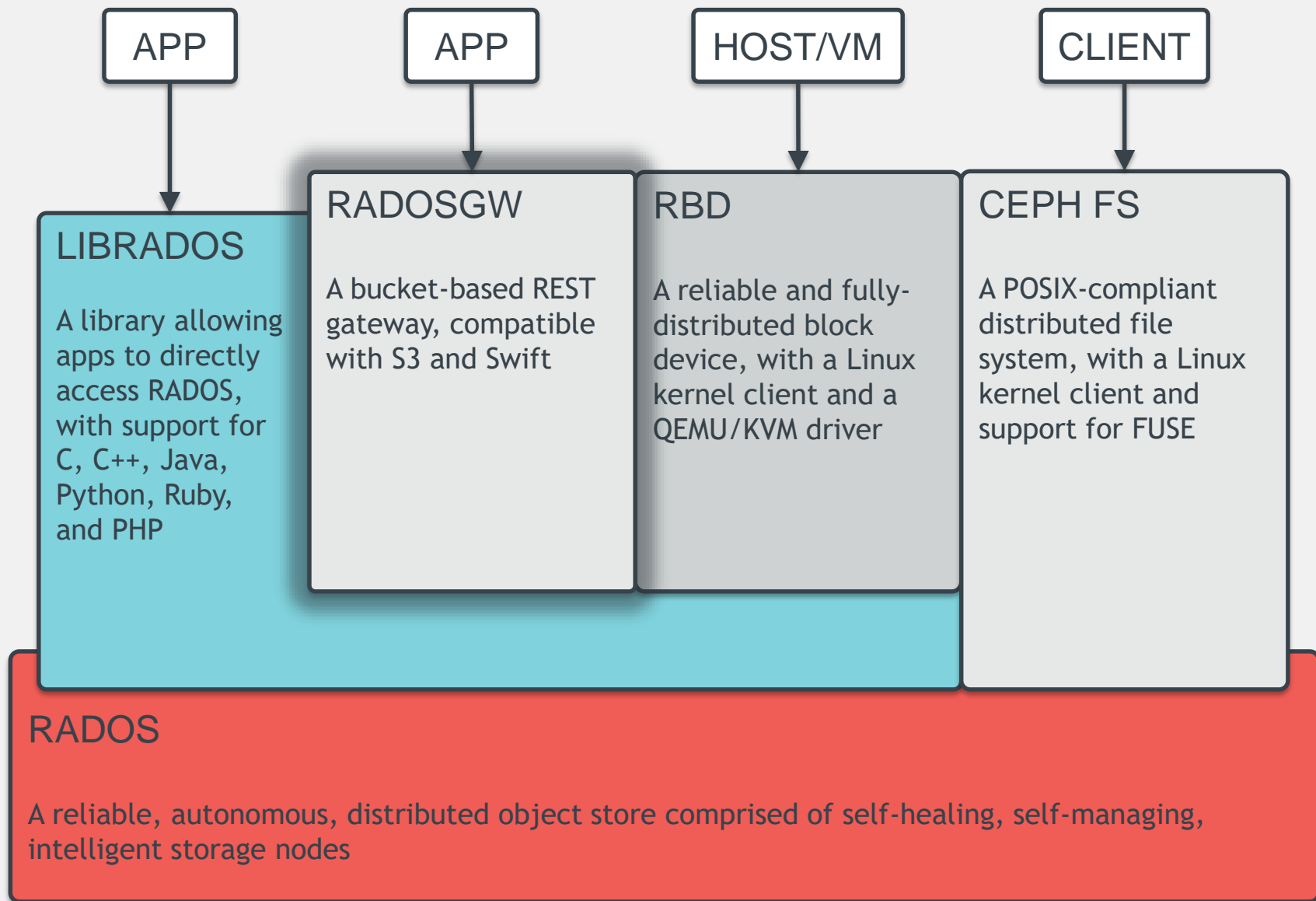


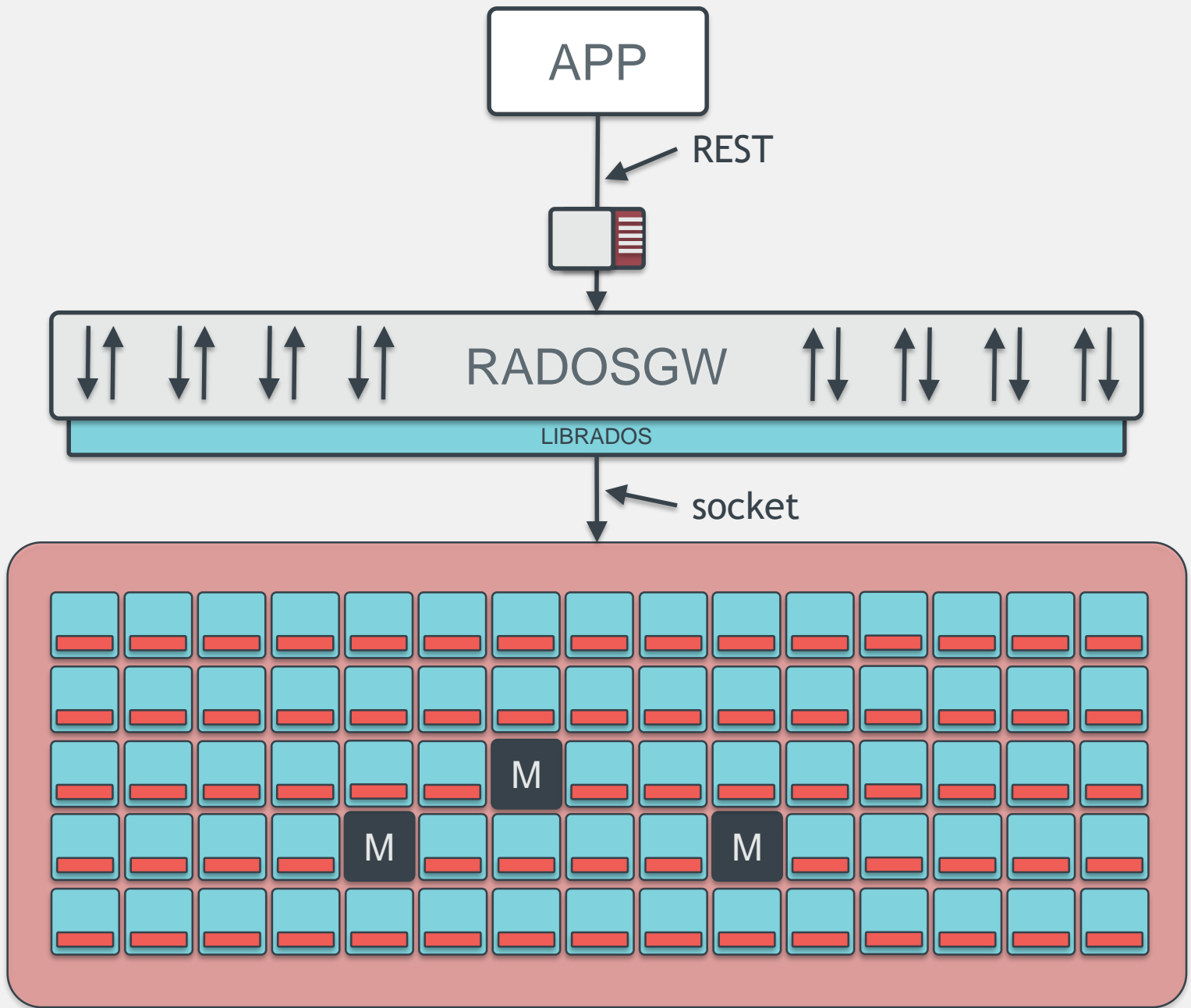


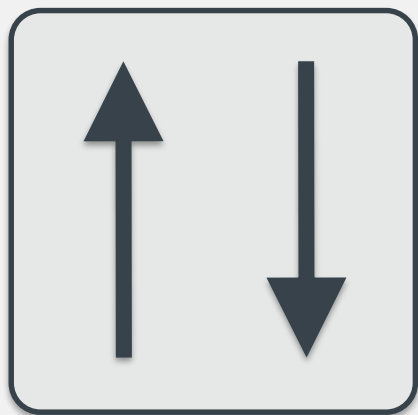


## LIBRADOS

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

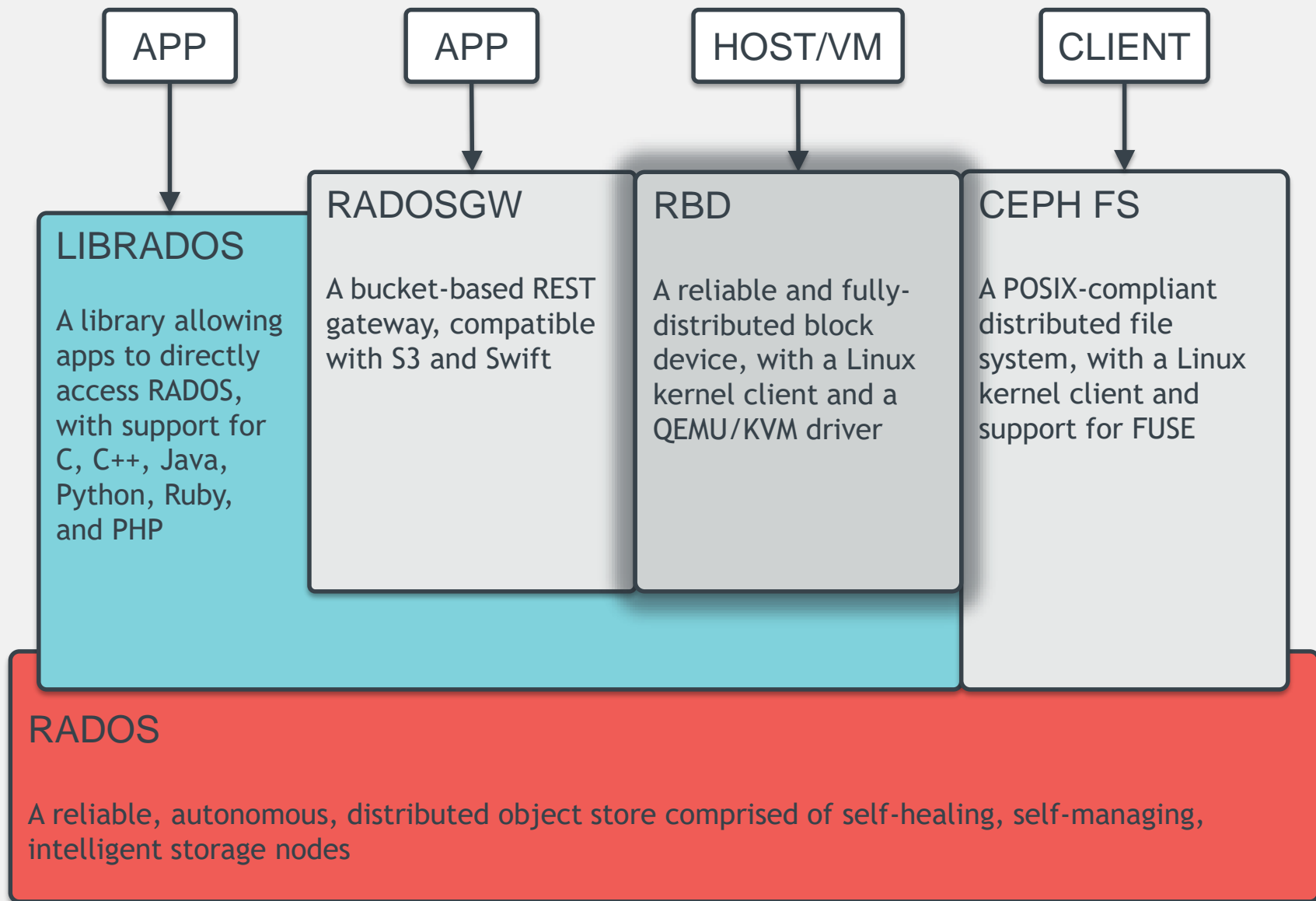


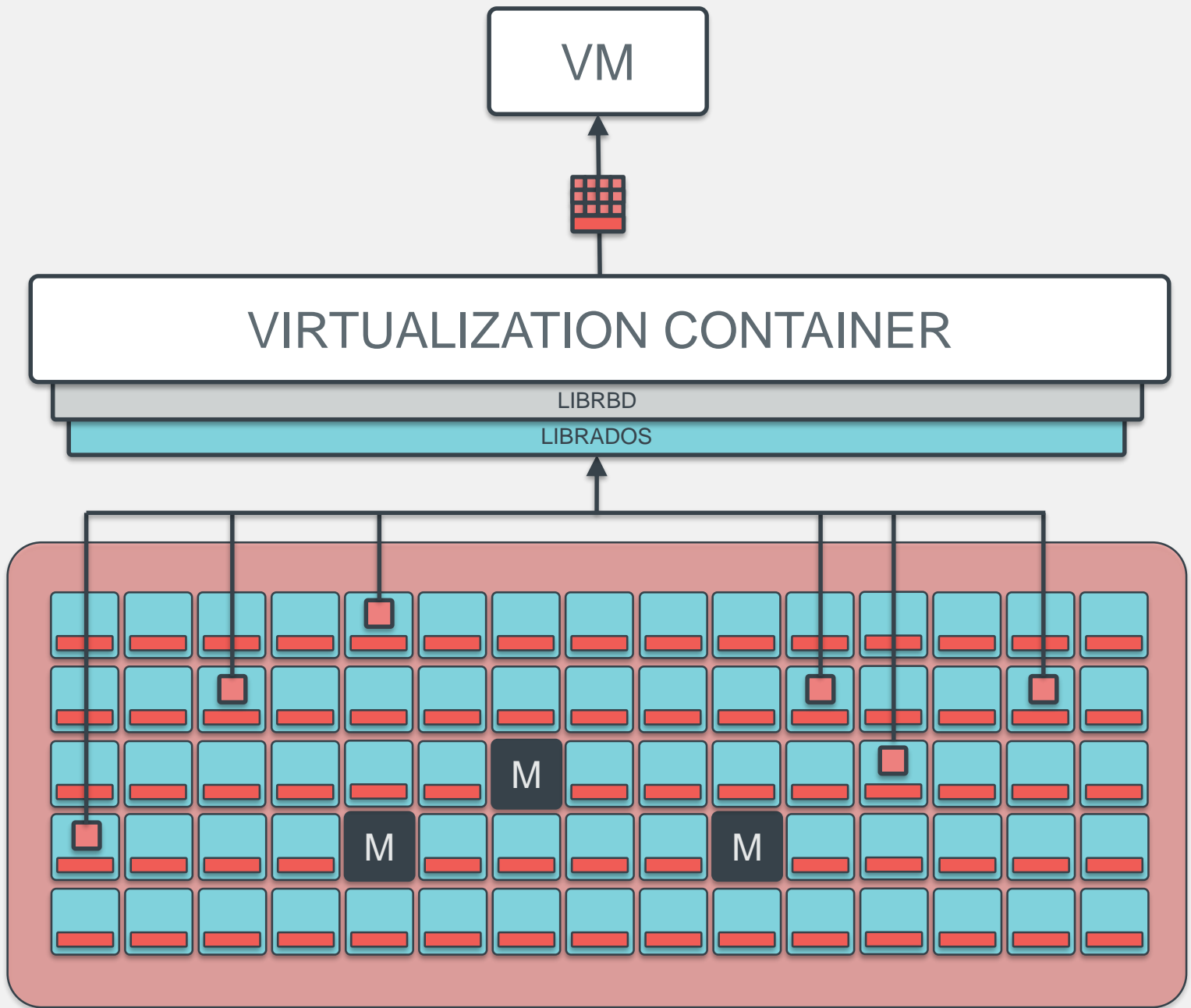


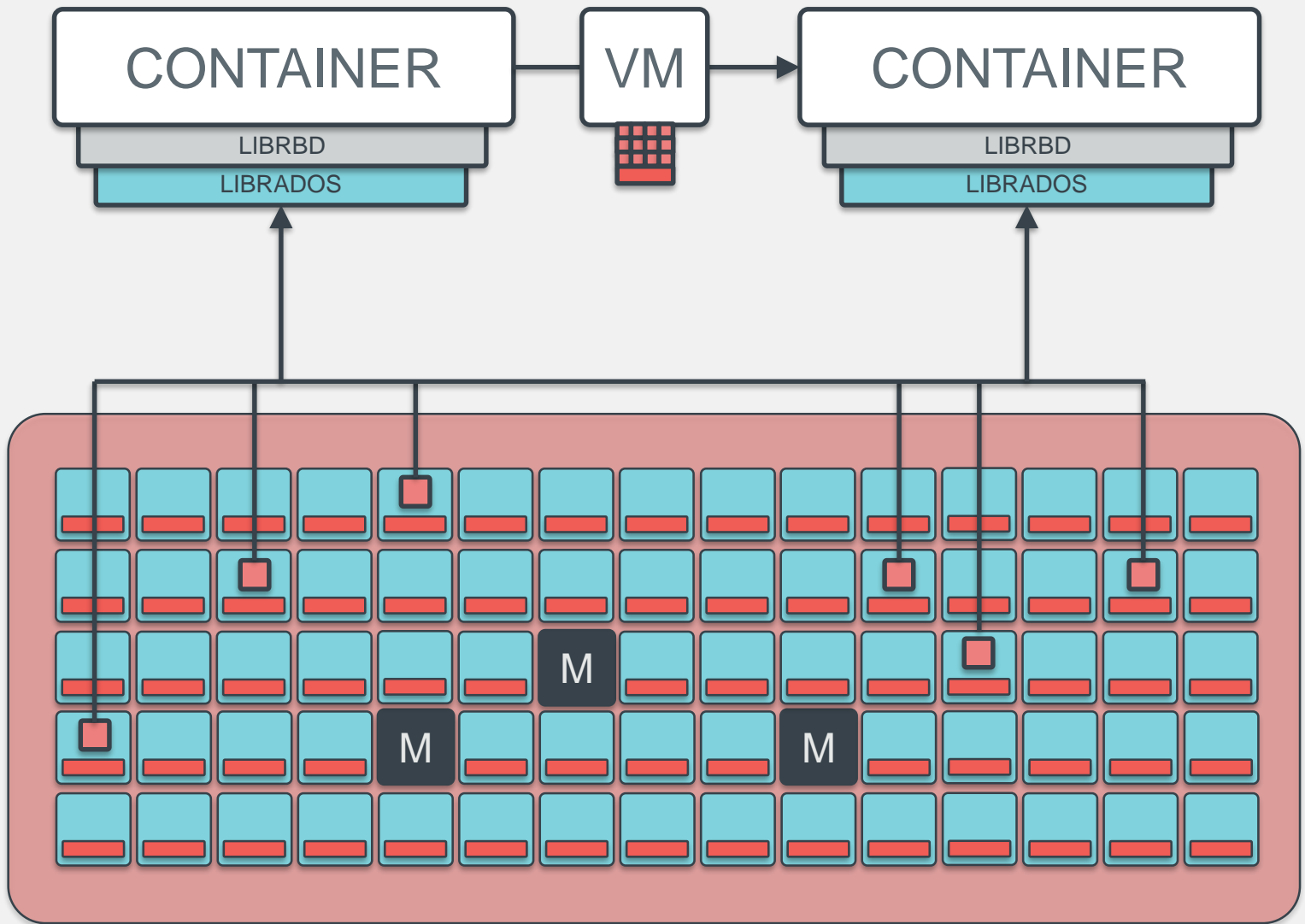


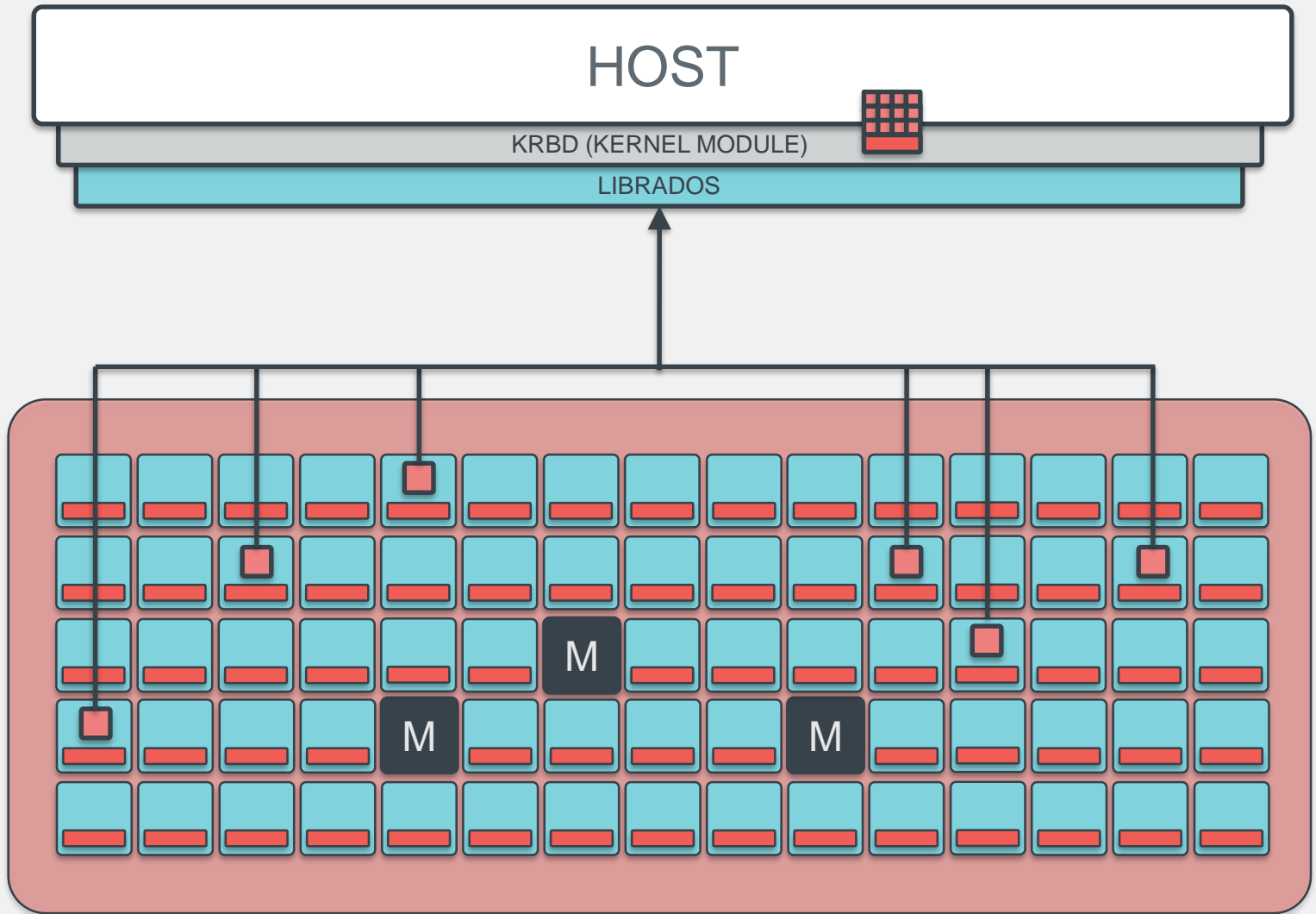
## RADOS Gateway:

- REST-based object storage proxy
- Uses RADOS to store objects
- API supports buckets, accounts
- Usage accounting for billing
- Compatible with S3 and Swift applications

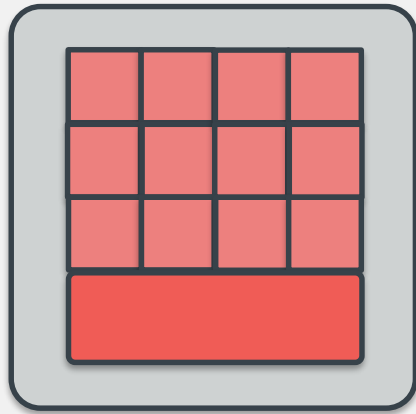






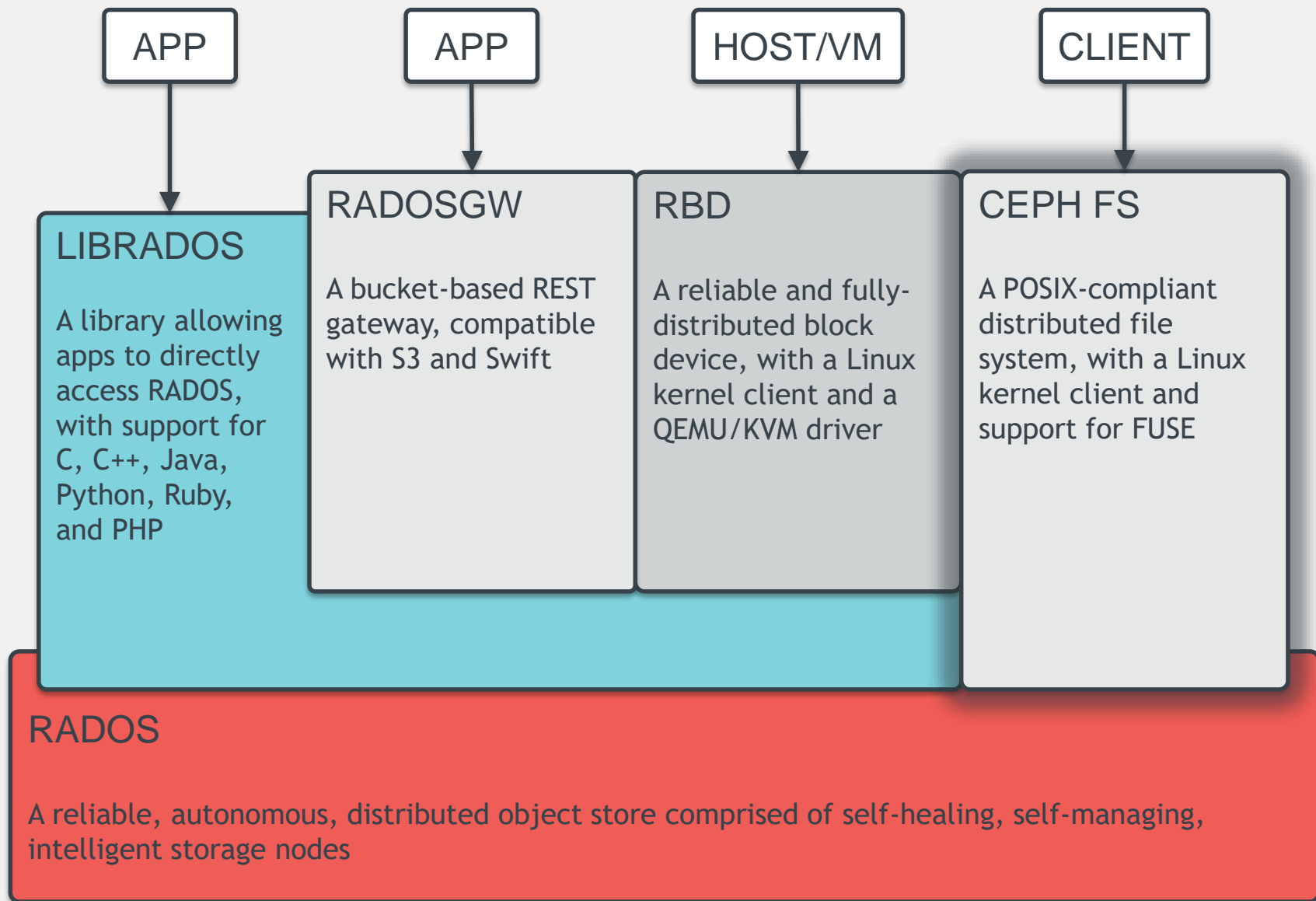


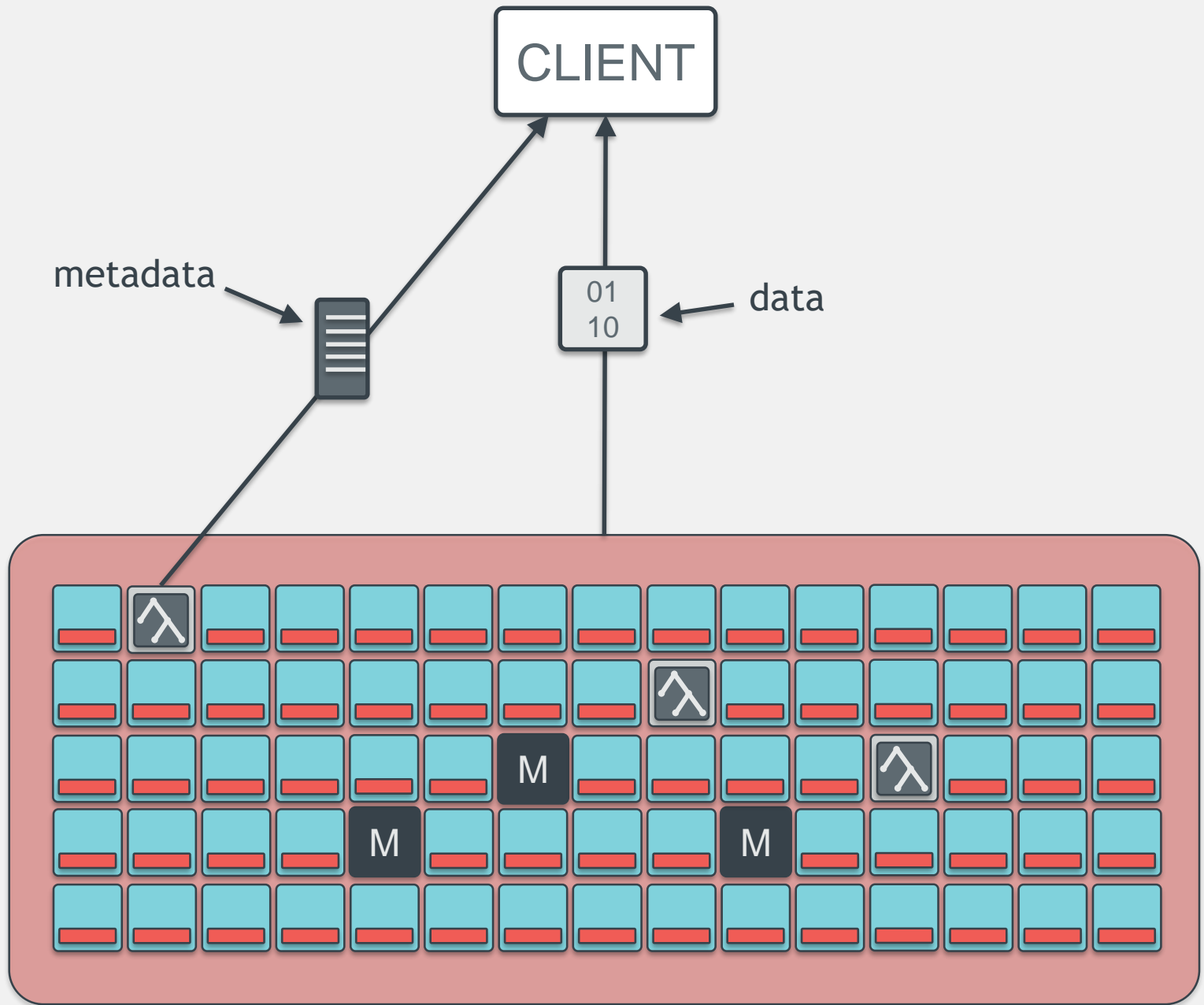


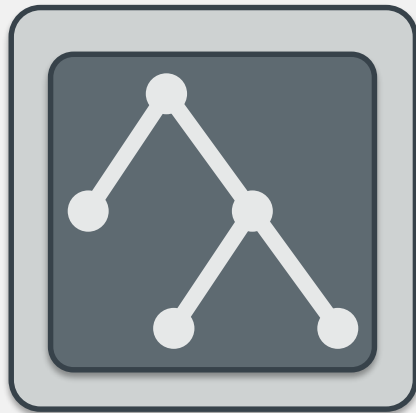


## RADOS Block Device:

- Storage of disk images in RADOS
- Decouples VMs from host
- Images are striped across the cluster (pool)
- Snapshots
- Copy-on-write clones
- Support in:
  - Mainline Linux Kernel (2.6.39+)
  - Qemu/KVM, native Xen coming soon
  - OpenStack, CloudStack, Nebula, Proxmox







## 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

# What Makes Ceph Unique?

Part one: CRUSH

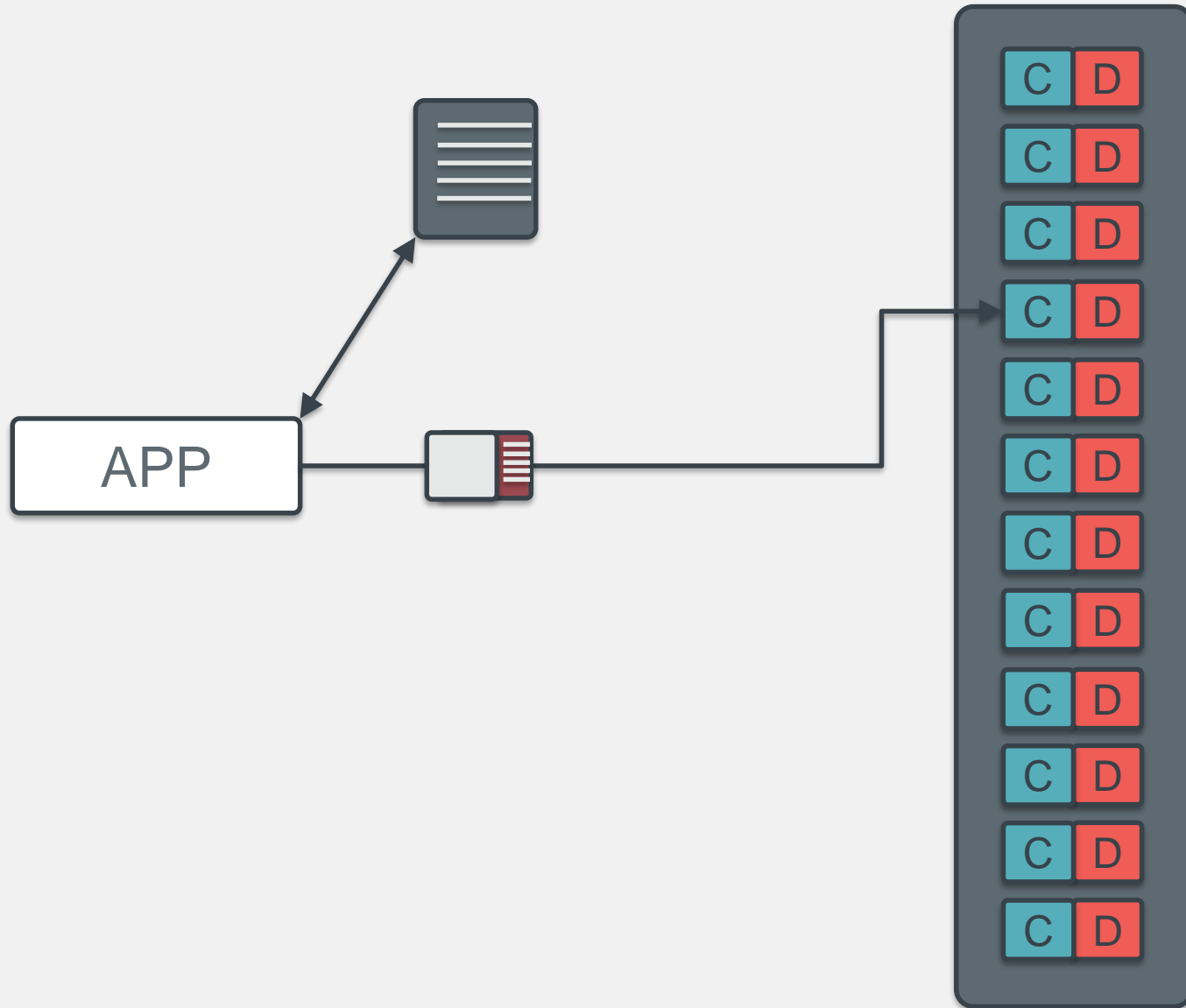






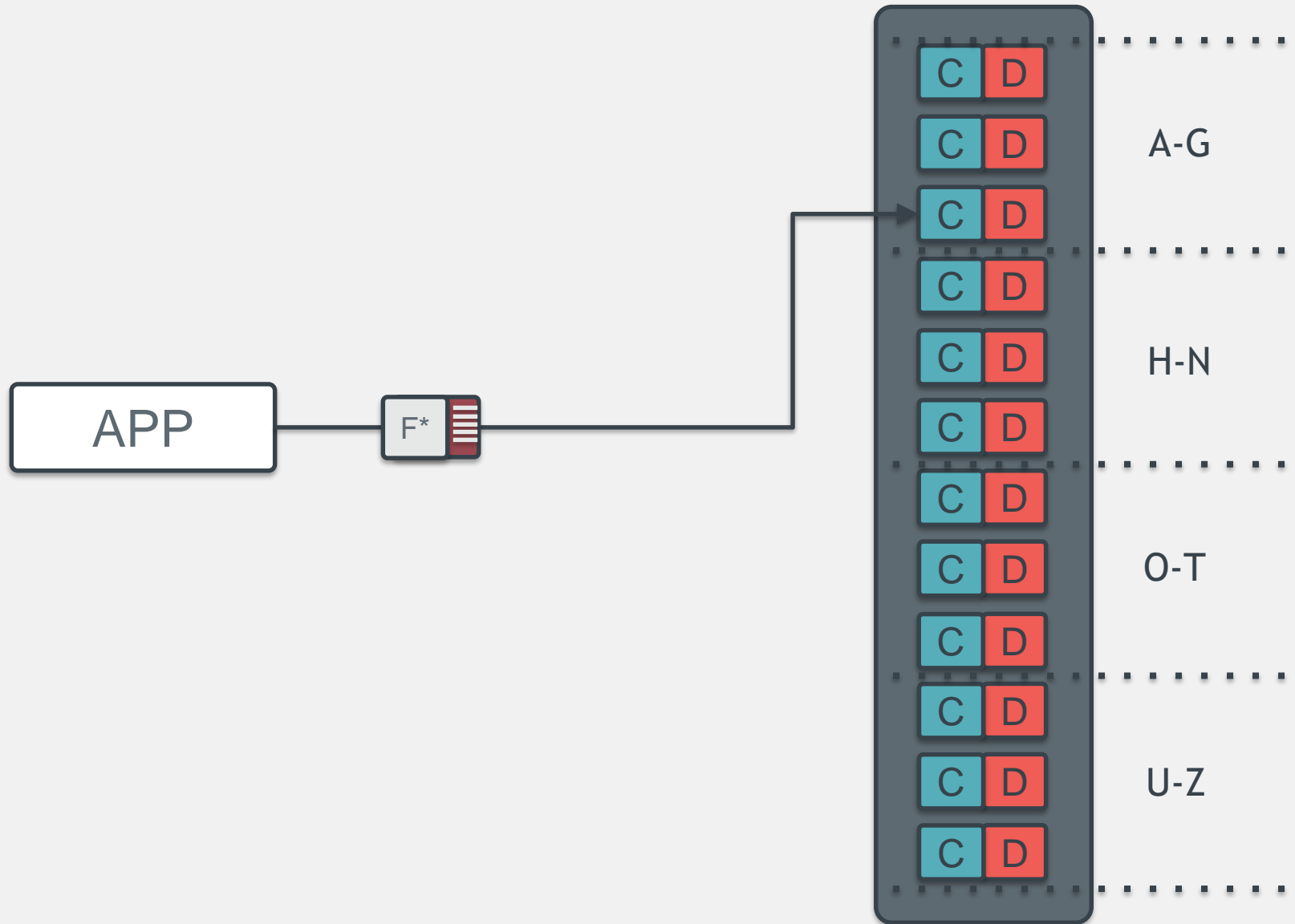
How Long Did It Take You To Find Your Keys This Morning?

azmeen, Flickr / CC BY 2.0











HOW DO YOU  
FIND YOUR KEYS  
WHEN YOUR HOUSE  
IS  
**INFINITELY BIG**  
AND  
**ALWAYS CHANGING?**





**The Answer: CRUSH!!!!**

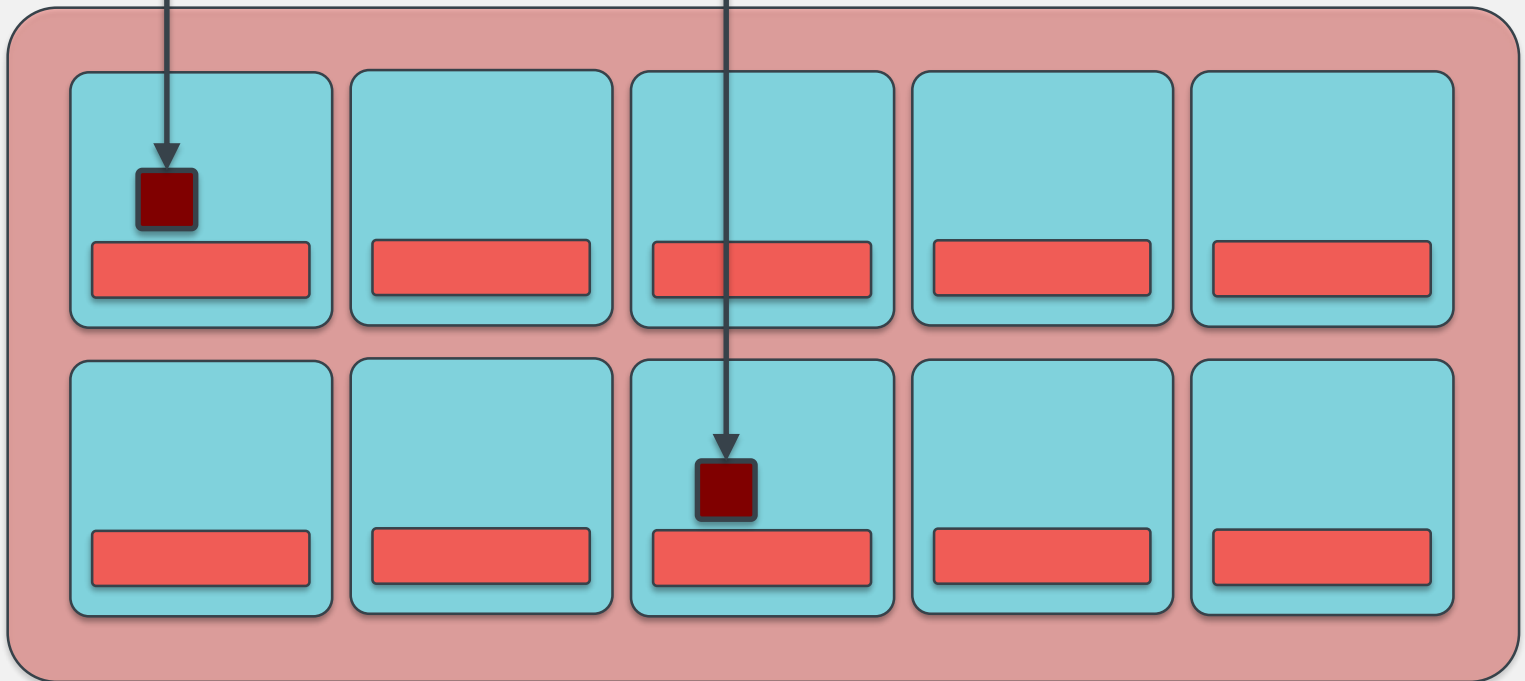
pasukaru76, Flickr / CC SA 2.0

10 10 01 01 10 10 01 11 01 10

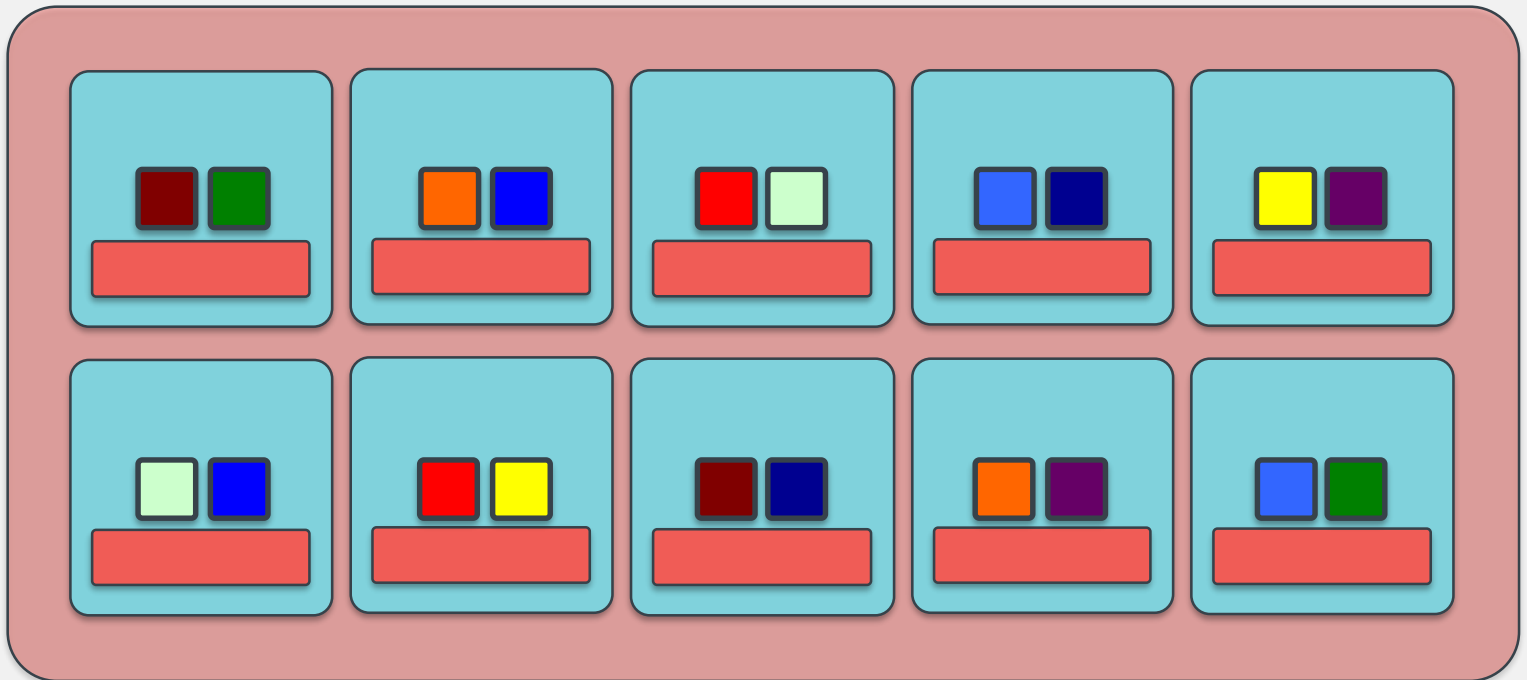
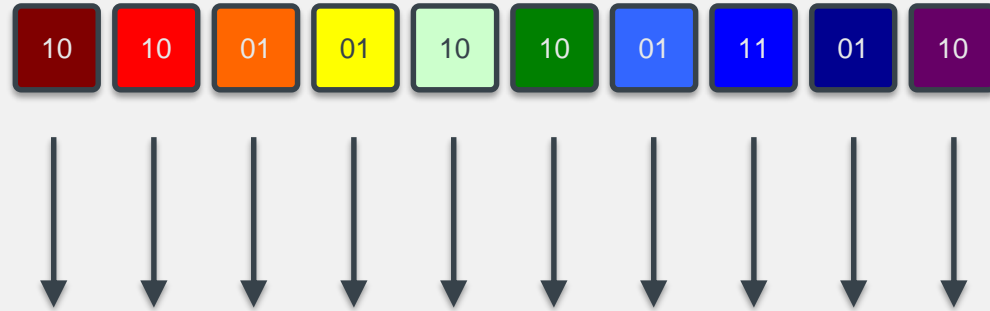
$\text{hash}(\text{object name}) \% \text{num pg}$

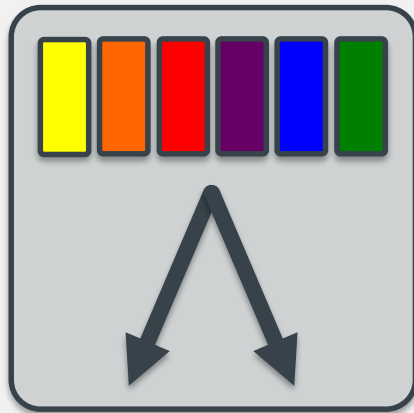


$\text{CRUSH}(\text{pg}, \text{cluster state}, \text{rule set})$



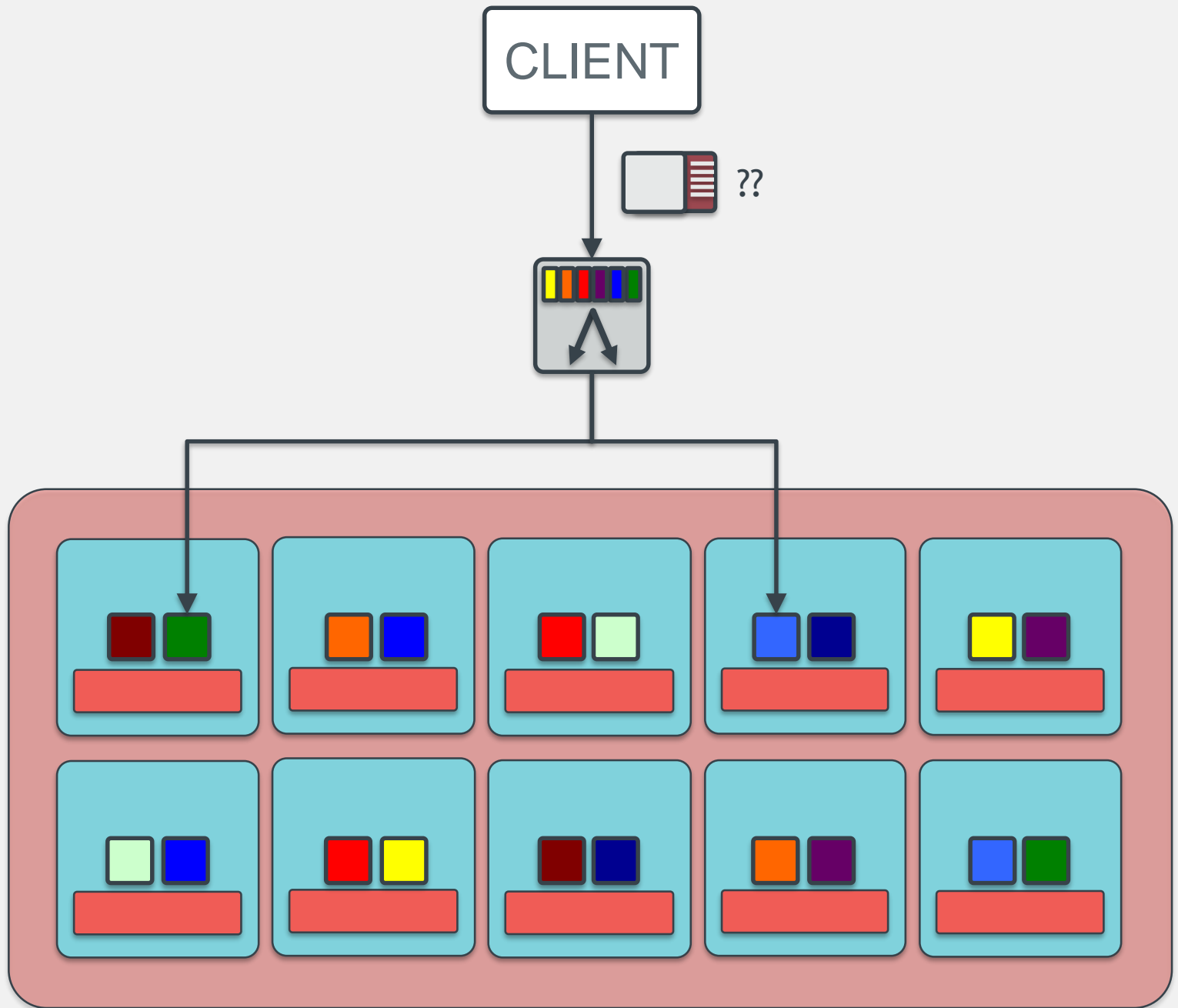
10 10 01 01 10 10 01 11 01 10



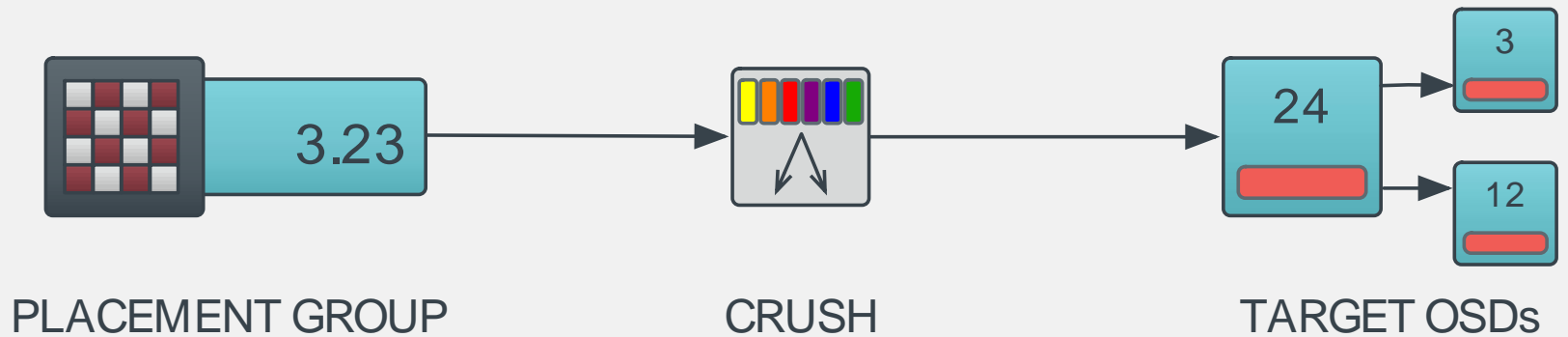
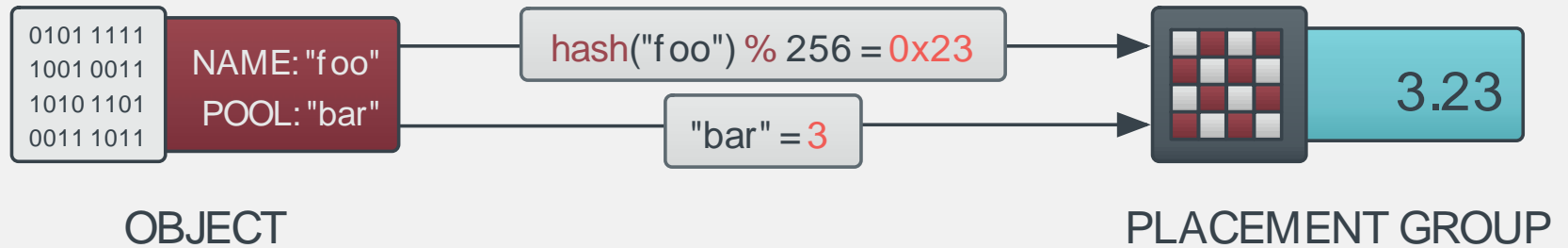


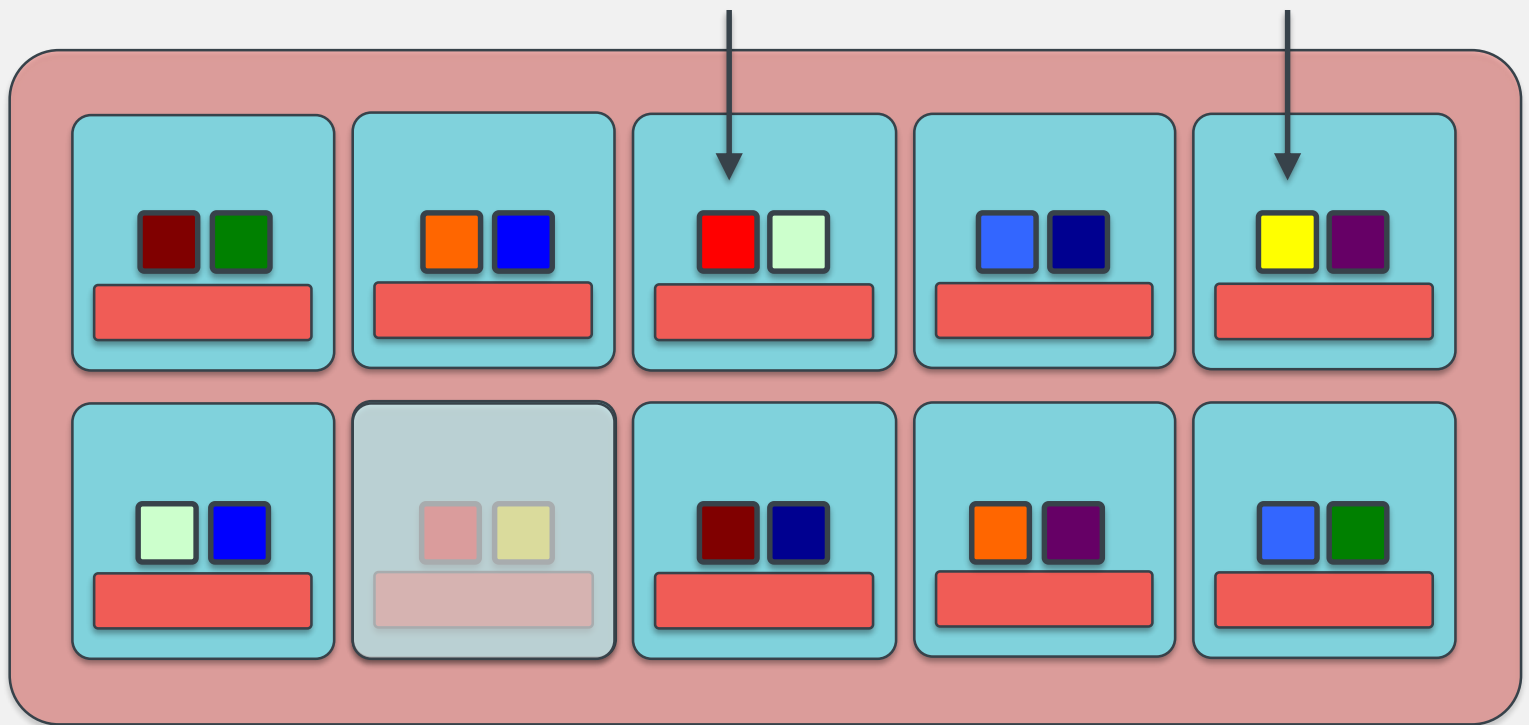
## CRUSH

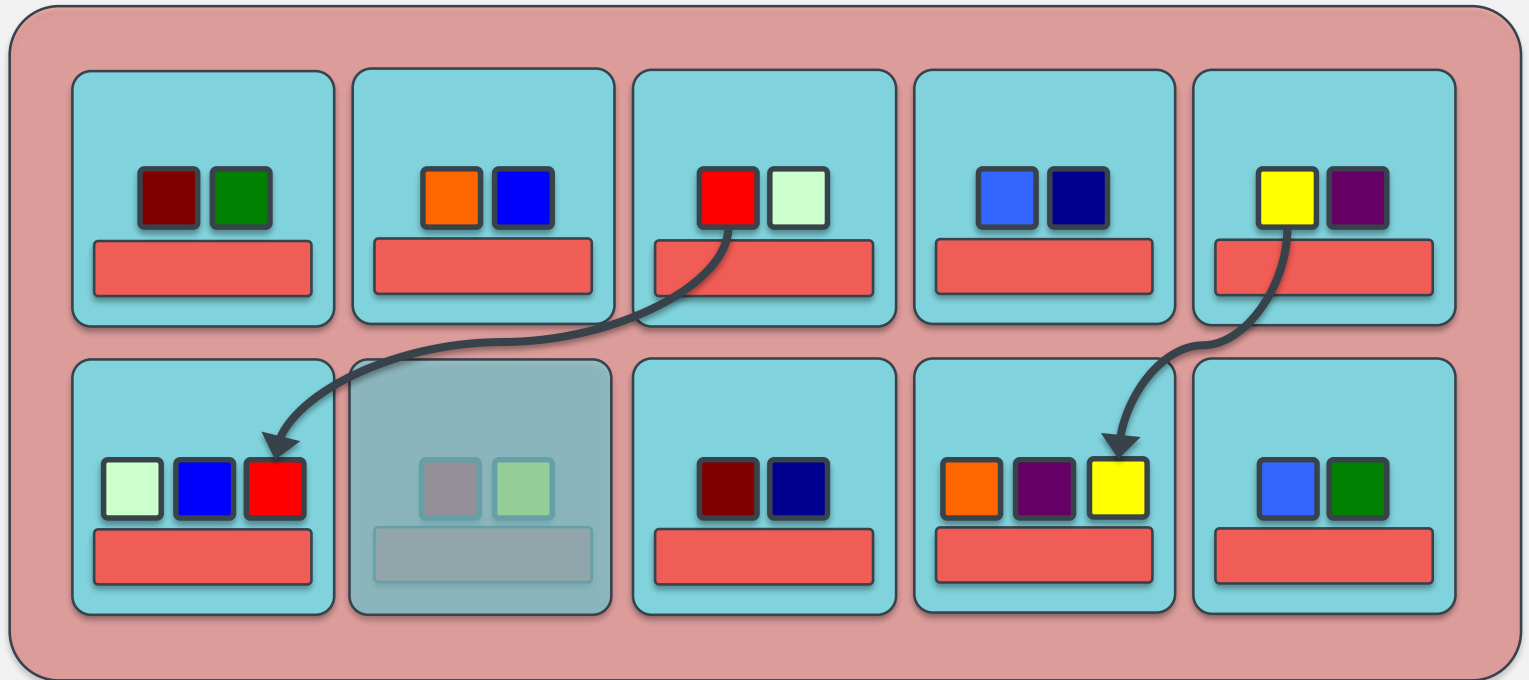
- Pseudo-random placement algorithm
  - Fast calculation, **no lookup**
  - Repeatable, deterministic
- Statistically uniform distribution
- Stable mapping
  - Limited data migration on change
- Rule-based configuration
  - Infrastructure topology aware
  - Adjustable replication
  - Weighting

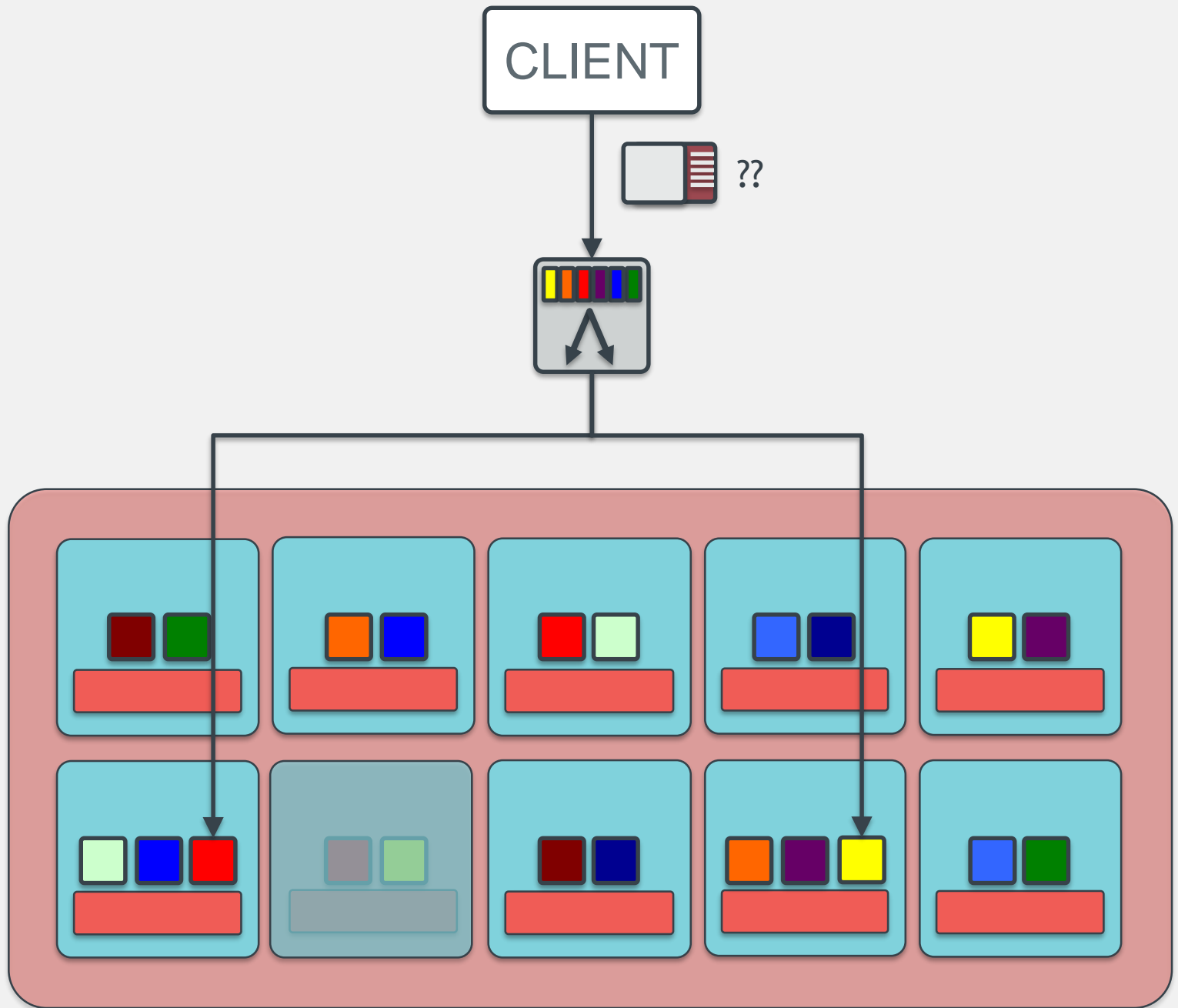








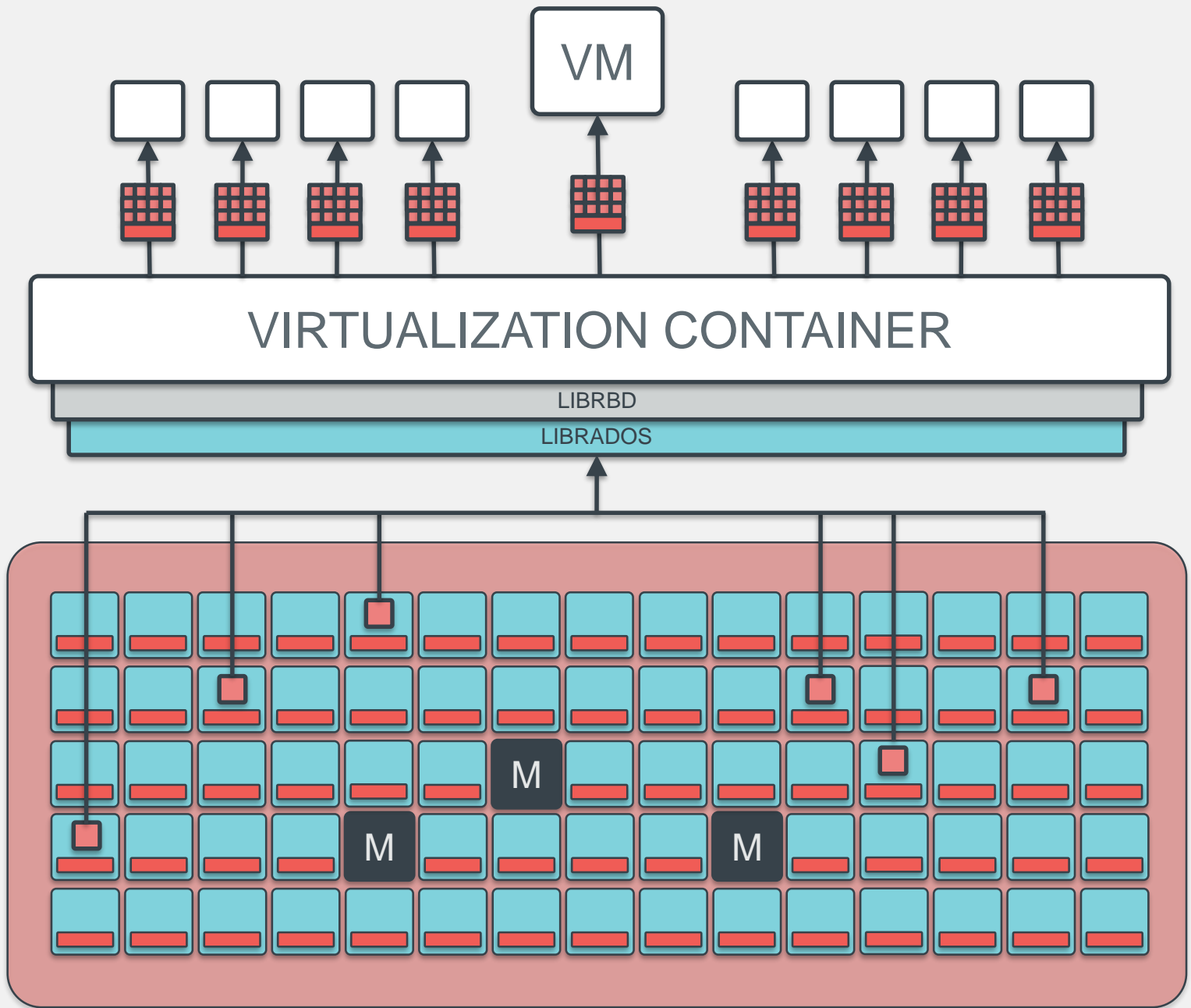





# What Makes Ceph Unique

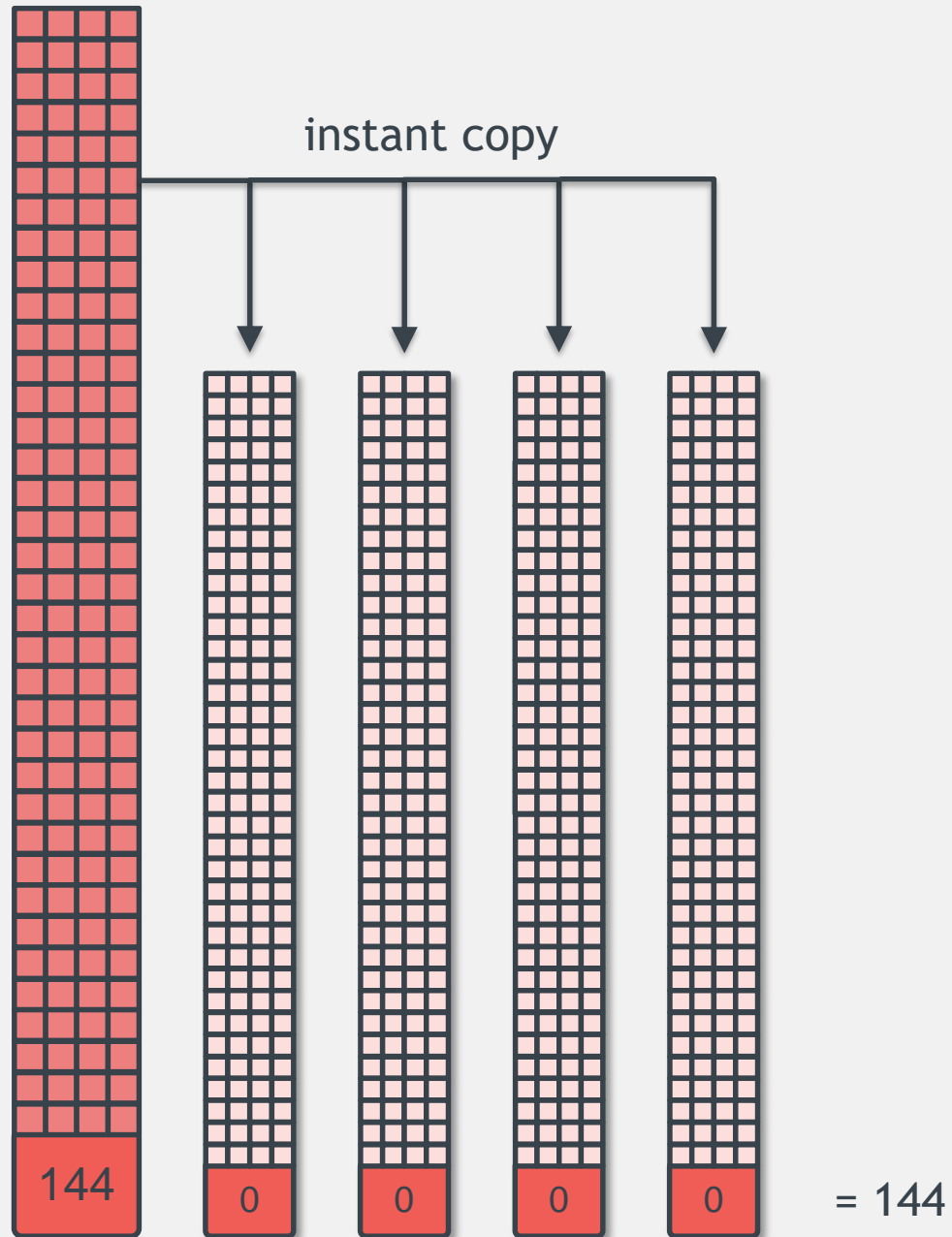
Part two: thin provisioning



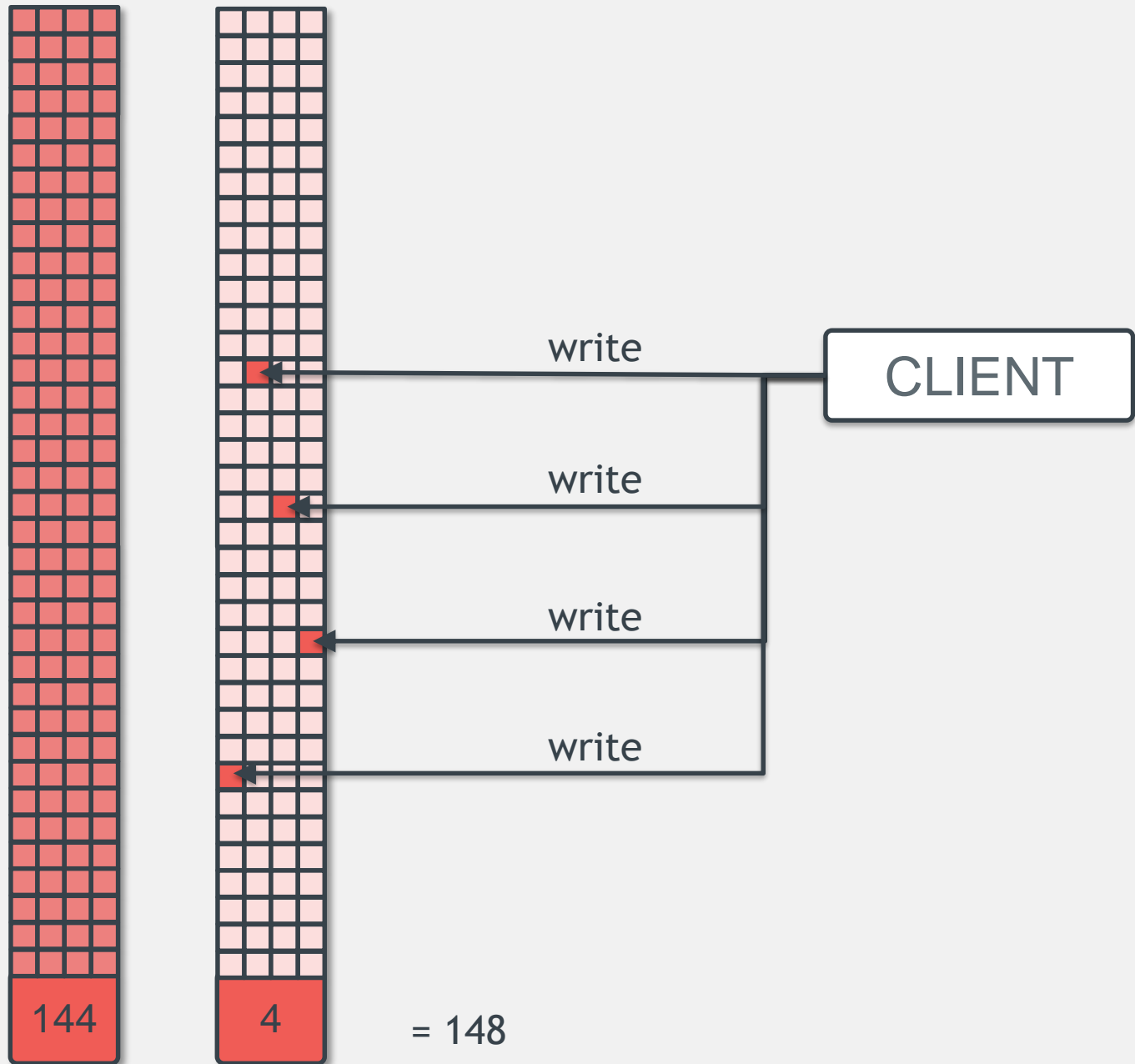


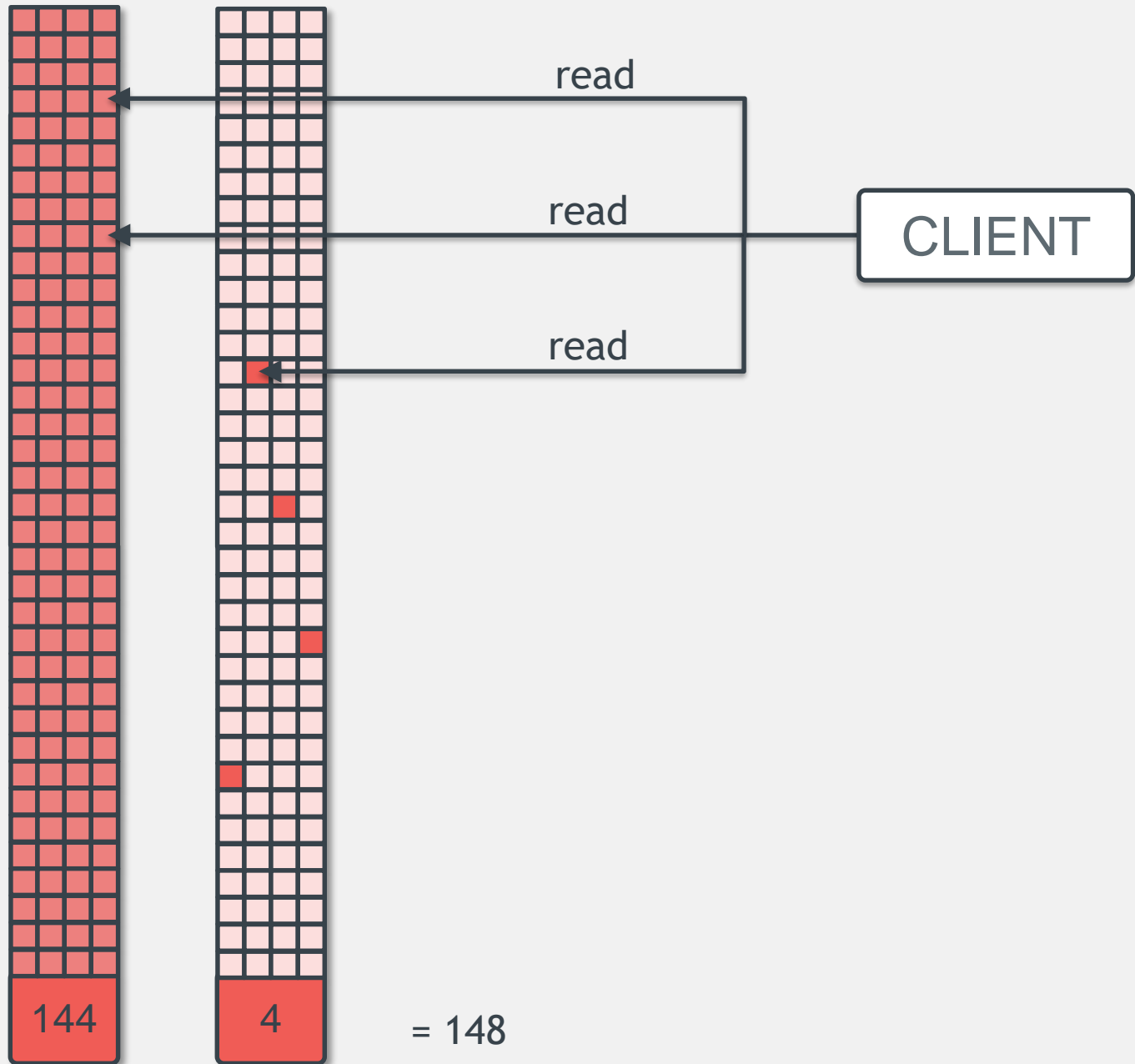


HOW DO YOU  
SPIN UP  
THOUSANDS OF VMs  
**INSTANTLY**  
AND  
**EFFICIENTLY?**









# What Makes Ceph Unique?

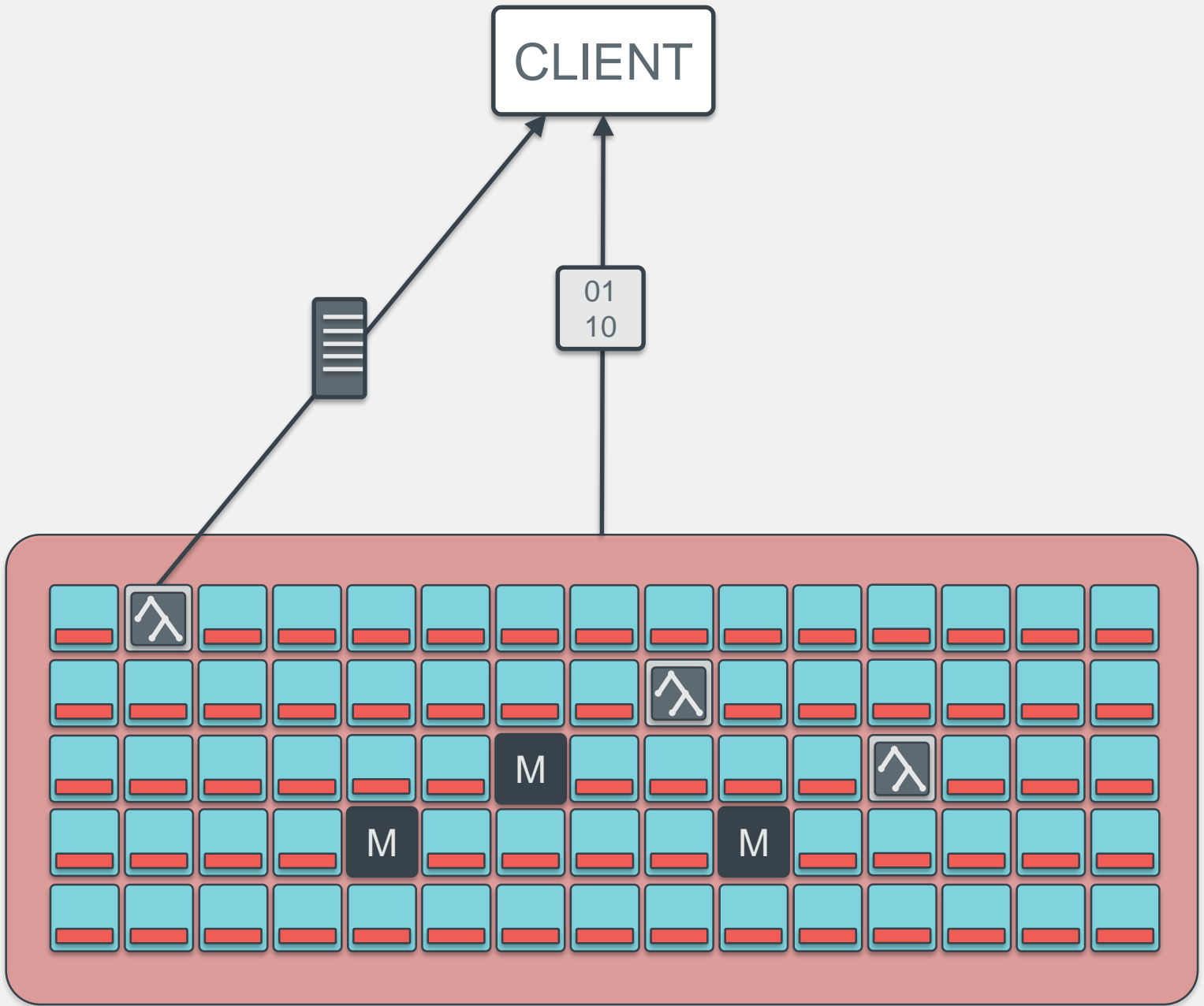
Part three: clustered metadata

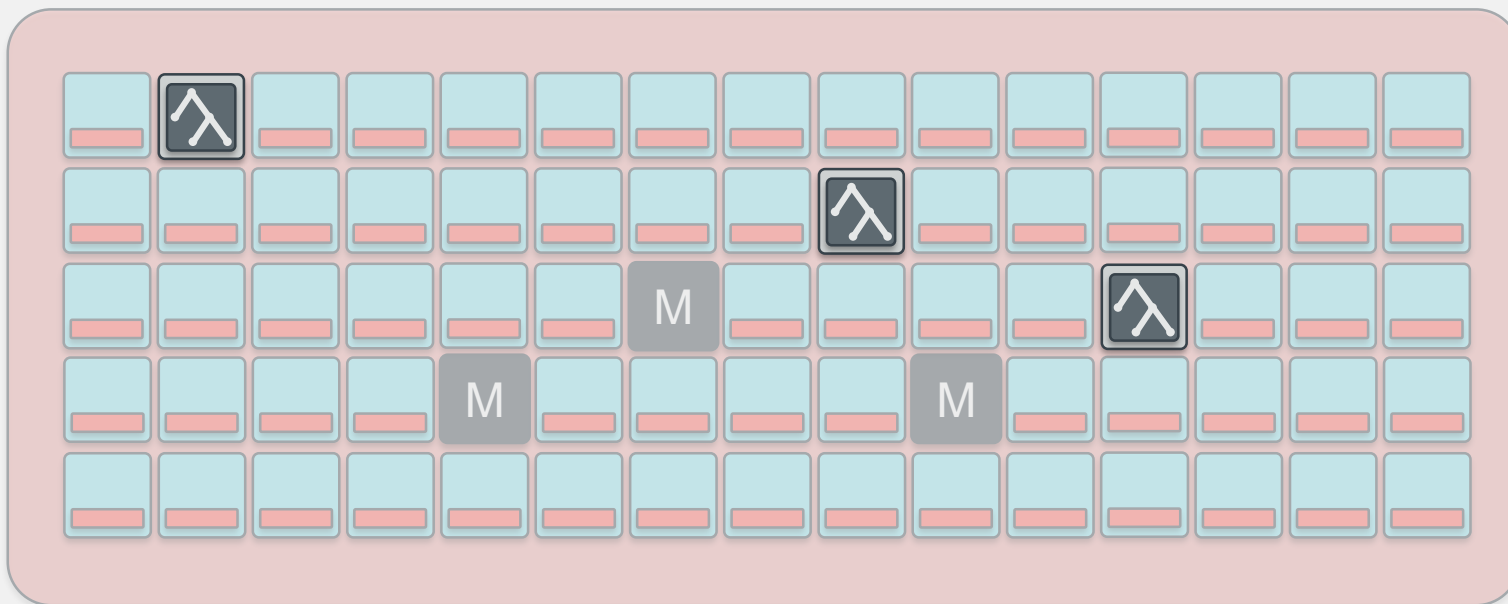


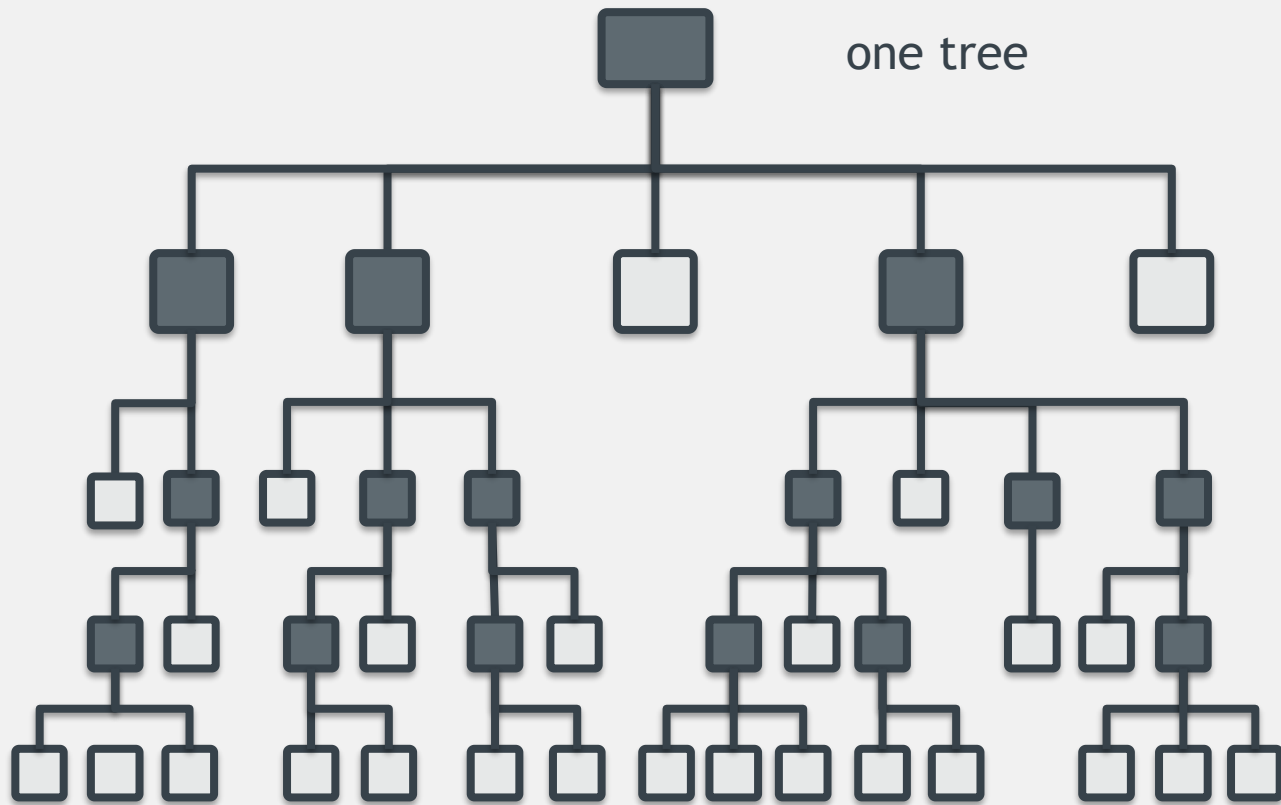
lrwxrwxrwx	1	root	root	26	Apr	26	17:54	libgssapi_krb5.so -> mit-krb5/libgssapi_krb5.so
lrwxrwxrwx	1	root	root	21	Apr	26	17:55	libgssapi_krb5.so.2 -> libgssapi_krb5.so.2.2
-rw-r--r--	50	root	root	216824	Jul	31	2012	libgssapi_krb5.so.2.2
lrwxrwxrwx	1	root	root	13	Apr	26	17:54	libgs.so.8 -> libgs.so.8.71
-rw-r--r--	17	root	root	9478048	Jan	25	2011	libgs.so.8.71
lrwxrwxrwx	1	root	root	21	Apr	26	17:55	libgssrpc.so -> mit-krb5/libgssrpc.so
lrwxrwxrwx	1	root	root	16	Apr	26	17:55	libgssrpc.so.4 -> libgssrpc.so.4.1
-rw-r--r--	50	root	root	115352	Jul	31	2012	libgssrpc.so.4.1
-rw-r--r--	50	root	root	21832	Sep	8	2010	libgthread-2.0.a
-rw-r--r--	50	root	root	972	Sep	8	2010	libgthread-2.0.la
lrwxrwxrwx	1	root	root	26	Apr	26	17:55	libgthread-2.0.so -> libgthread-2.0.so.0.2400.2
lrwxrwxrwx	1	root	root	26	Apr	26	17:55	libgthread-2.0.so.0 -> libgthread-2.0.so.0.2400.2
-rw-r--r--	50	root	root	17704	Sep	8	2010	libgthread-2.0.so.0.2400.2
drwxr-xr-x	2	root	root	4096	Apr	26	18:00	libgtk2.0-0
-rw-r--r--	49	root	root	9275282	Oct	14	2010	libgtk-x11-2.0.a
-rw-r--r--	49	root	root	981	Oct	14	2010	libgtk-x11-2.0.la
lrwxrwxrwx	1	root	root	26	Apr	26	17:55	libgtk-x11-2.0.so -> libgtk-x11-2.0.so.0.2000.1
lrwxrwxrwx	1	root	root	26	Apr	26	17:55	libgtk-x11-2.0.so.0 -> libgtk-x11-2.0.so.0.2000.1
-rw-r--r--	49	root	root	4319784	Oct	14	2010	libgtk-x11-2.0.so.0.2000.1
lrwxrwxrwx	1	root	root	15	Apr	26	17:55	libgvc.so -> libgvc.so.5.0.0
lrwxrwxrwx	1	root	root	15	Apr	26	17:55	libgvc.so.5 -> libgvc.so.5.0.0
-rw-r--r--	49	root	root	504424	Jul	5	2010	libgvc.so.5.0.0
lrwxrwxrwx	1	root	root	16	Apr	26	17:55	libgvpr.so -> libgvpr.so.1.0.0
lrwxrwxrwx	1	root	root	16	Apr	26	17:55	libgvpr.so.1 -> libgvpr.so.1.0.0
-rw-r--r--	50	root	root	482856	Jul	5	2010	libgvpr.so.1.0.0
-rw-r--r--	50	root	root	267948	Apr	13	2009	libHalf.a
lrwxrwxrwx	1	root	root	16	Apr	26	17:55	libHalf.so -> libHalf.so.6.0.0
lrwxrwxrwx	1	root	root	16	Apr	26	17:55	libHalf.so.6 -> libHalf.so.6.0.0
-rw-r--r--	50	root	root	269992	Apr	13	2009	libHalf.so.6.0.0
-rw-r--r--	50	root	root	52850	Nov	1	2009	libhistory.a

## POSIX Filesystem Metadata

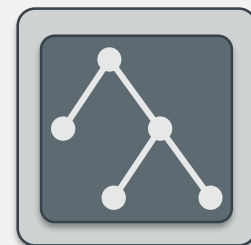
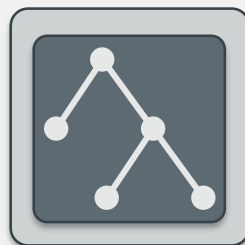
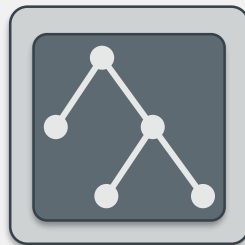
Barnaby, Flickr / CC BY 2.0



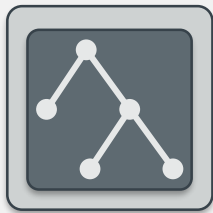




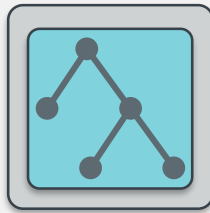
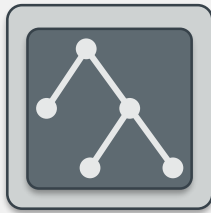
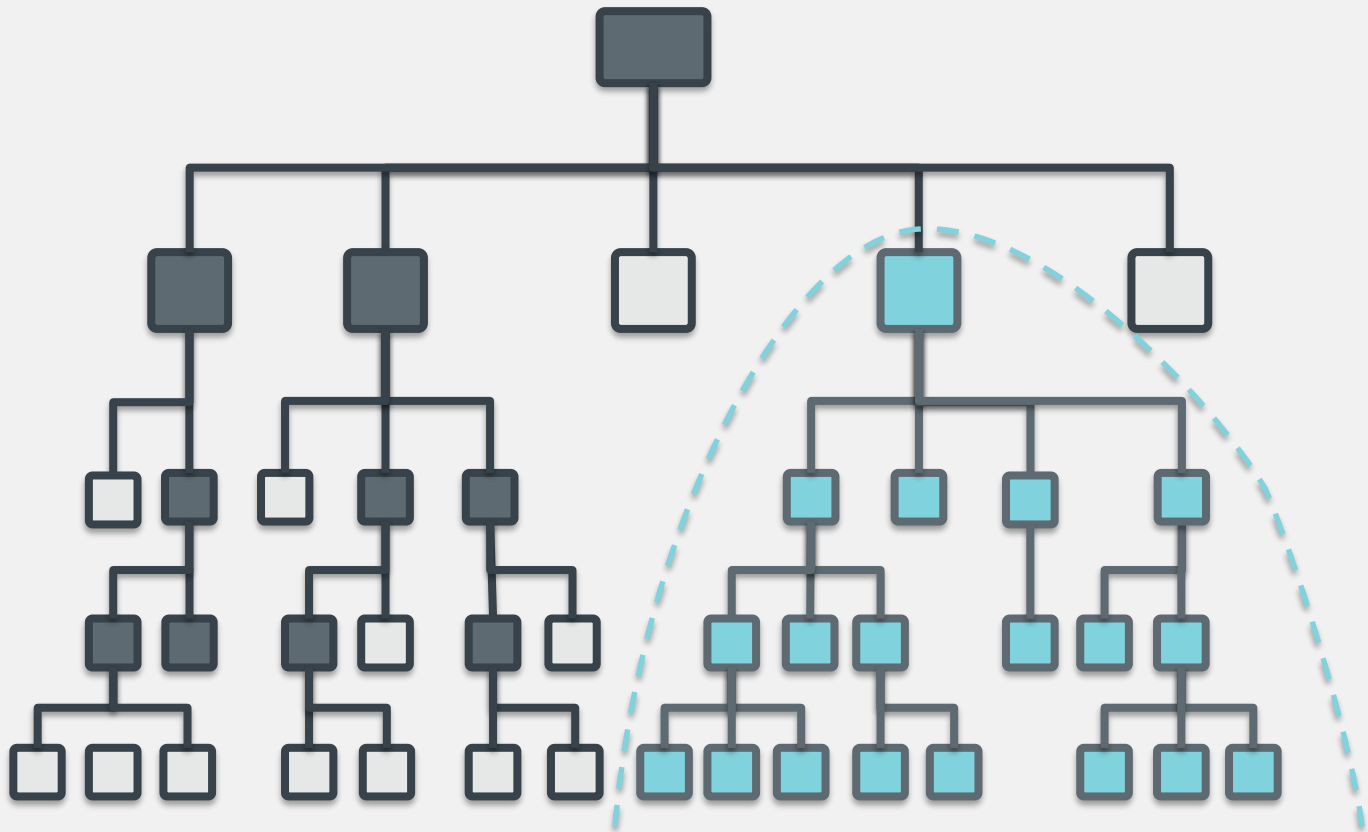
three metadata servers

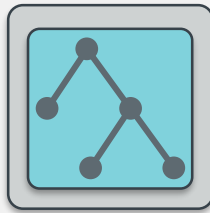
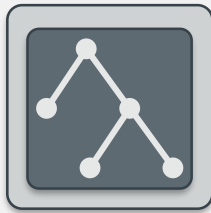
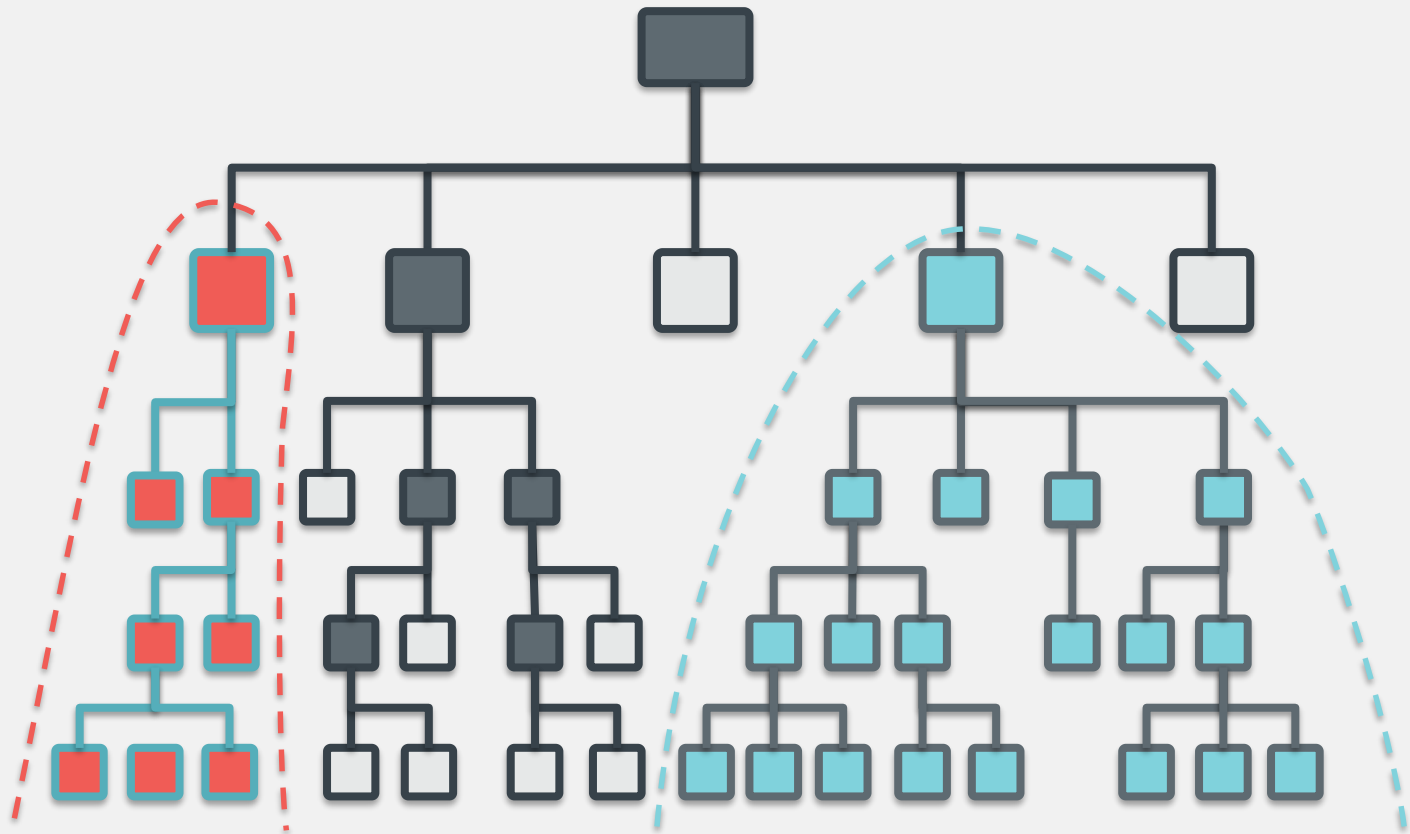


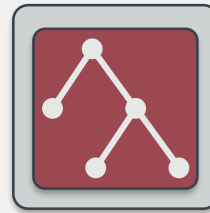
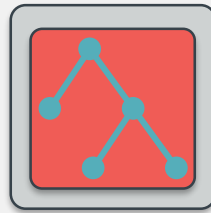
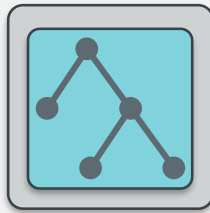
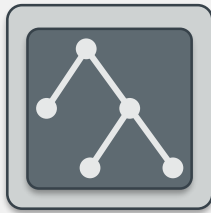
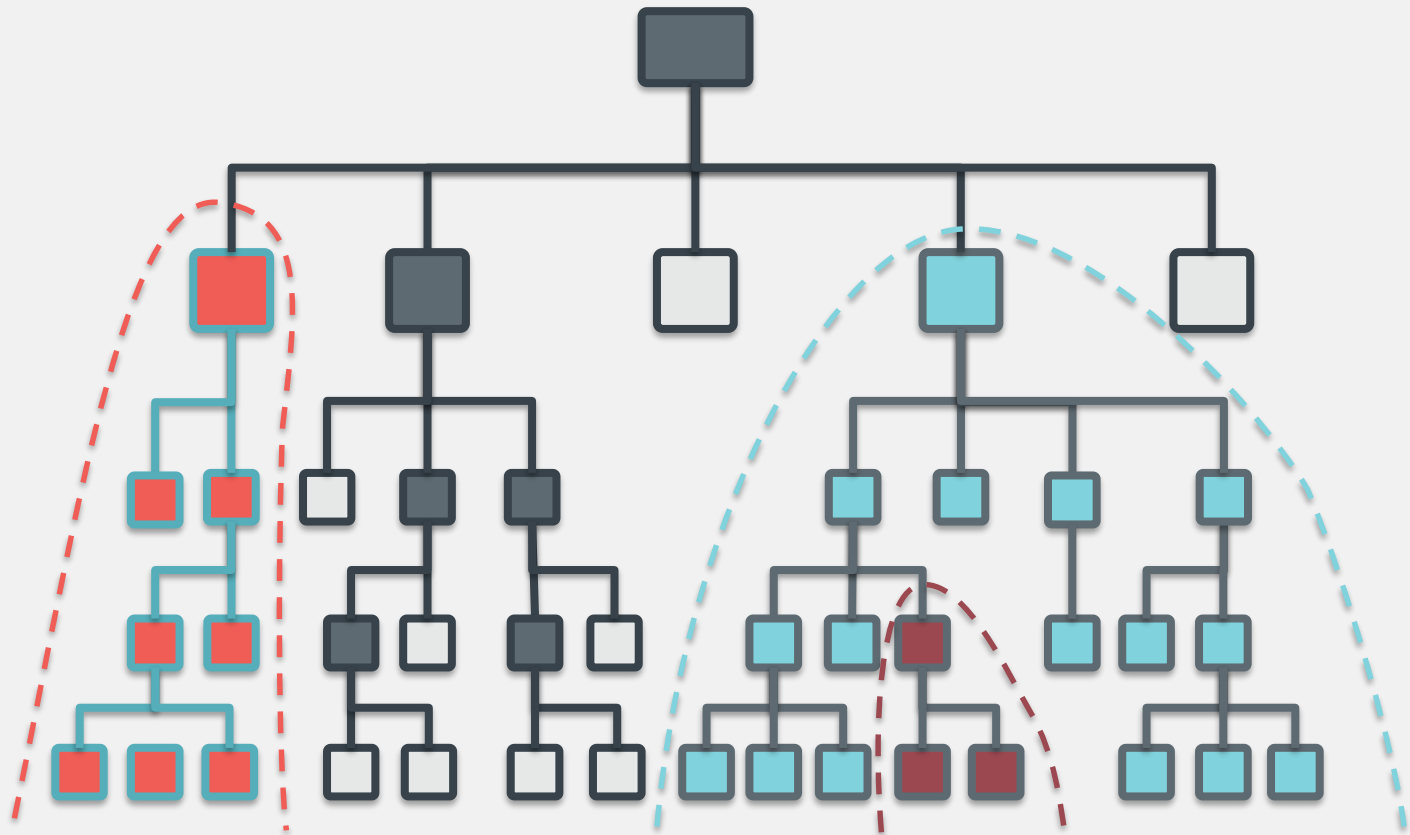
??

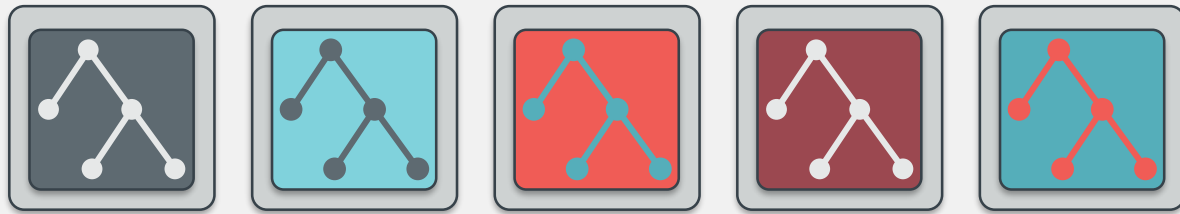
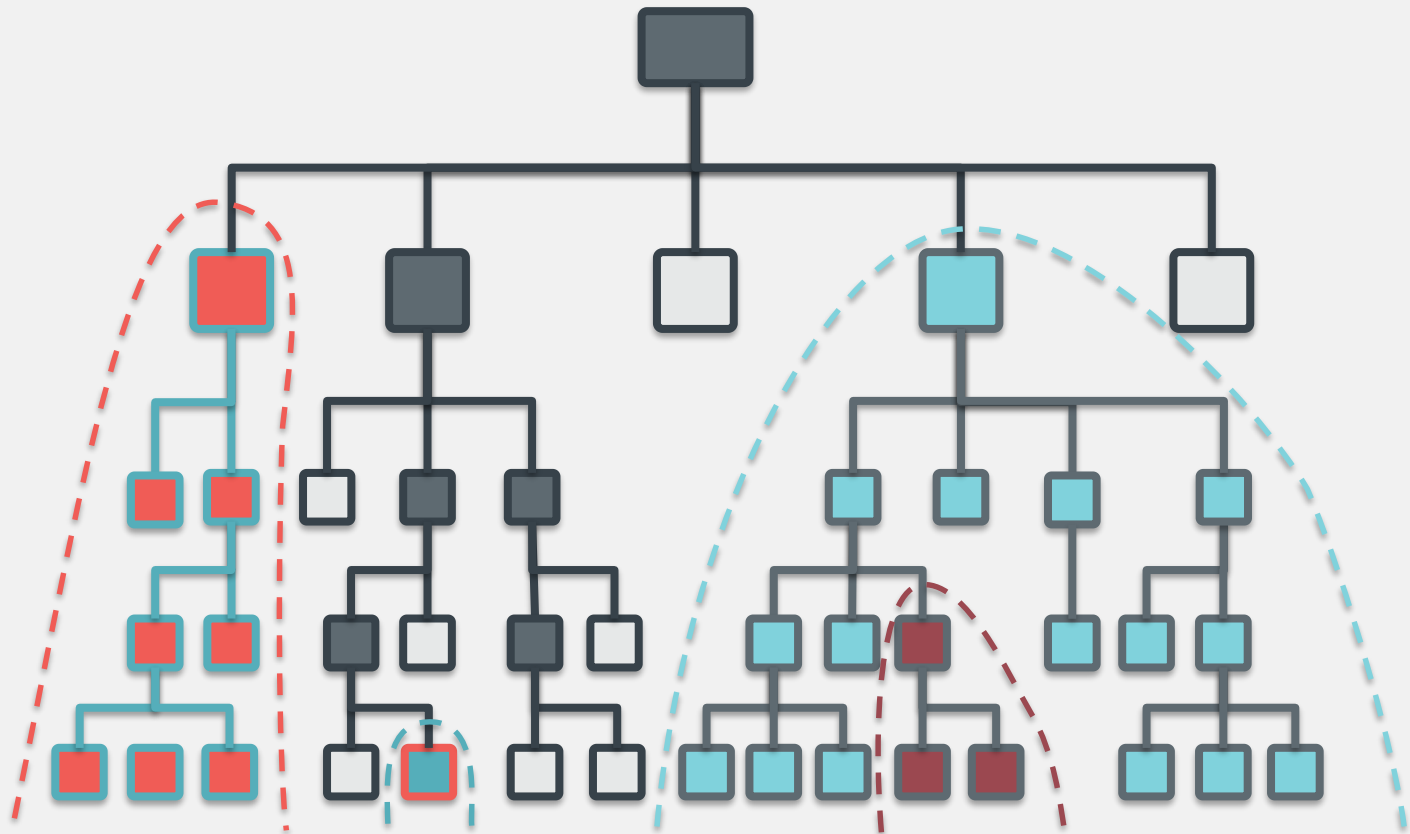












DYNAMIC SUBTREE PARTITIONING

# OpenStack + Ceph

Used for *Glance* Images, *Cinder* Volumes and *Nova* ephemeral disk (coming soon)

Ceph + OpenStack offers compelling features:

- CoW clones, layered volumes, snapshots, boot from volume, live migration
- Cost effective with Thin Provisioning
  - ~110TB “used”, ~45TB \* replicas on disk

Ceph is the most popular network block storage backend for OpenStack

# Deployment

Automated deployment using Cephdeploy

Automated machine commissioning and maintenance

- Add a server to the hostgroup (osd, mon, radosgw)
- OSD disks are detected, formatted, prepared, auth'd
  - Also after disk replacement
- Auto-generated ceph.conf
- Last step is manual/controlled: service ceph start

Cephdeploy for bulk operations on the servers

- Ceph rpm upgrades
- daemon restarts

# Getting Started With Ceph

Have a working cluster up quickly.

Read about the latest version of Ceph.

- The latest stuff is always at <http://ceph.com/get>

Deploy a test cluster using ceph-deploy.

- Read the quick-start guide at <http://ceph.com/qsg>

Deploy a test cluster on the AWS free-tier using Juju.

- Read the guide at <http://ceph.com/juju>

Read the rest of the docs!

- Find docs for the latest release at <http://ceph.com/docs>

# Getting Involved With Ceph

Help build the best storage system around!

Most project discussion happens on the mailing list.

- Join or view archives at <http://ceph.com/list>

IRC is a great place to get help (or help others!)

- Find details and historical logs at <http://ceph.com/irc>

The tracker manages our bugs and feature requests.

- Register and start looking around at <http://ceph.com/tracker>

Doc updates and suggestions are always welcome.

- Learn how to contribute docs at <http://ceph.com/docwriting>