# Parallel computing, data and storage

## "Grids, Clouds and distributed filesystems"

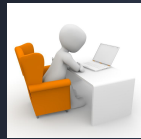Mário David - LIP Lisbon

david@lip.pt

# Overview

➢ Batch clusters

➢ Grid computing

➢ Cloud computing

    ○ Infrastructure as a Service (IaaS)

    ○ Platform as a Service (PaaS)

    ○ Software as a Service (SaaS)

➢ Distributed data and storage

    ○ Objects, Blocks and Filesystems (POSIX)

    ○ Parallel filesystems (Lustre filesystem case)

    ○ Object storage (Ceph case)

# Batch clusters: Introduction

➢ Set of compute nodes connected through a LAN
➢ Execute computational tasks
➢ Orchestrated by a master server:
   ○ Scheduler
   ○ Batch (queue) system
➢ Compute nodes in general: homogeneous hardware and operating system (OS):

   ○ Different hardware and OS can be grouped into different partitions (batch queues)
➢ Input and Output data for the computational tasks are served through a shared/distributed filesystem

# Batch clusters: usage (simplified view)

**Task/job submission**

**Task/job scheduling for execution**

**Compute nodes**

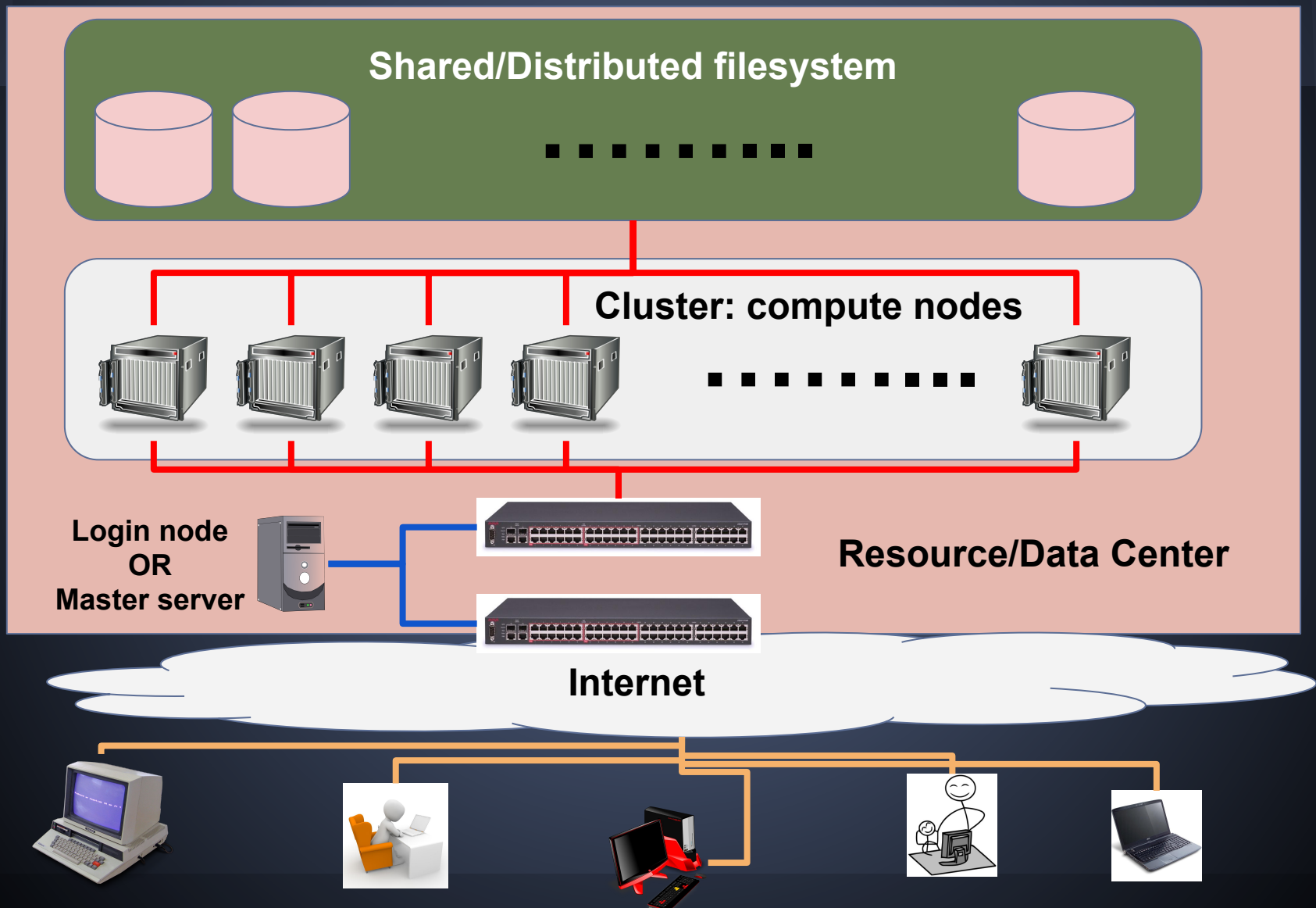**User submits a computational task to the master server**

**Master server:**
- **Task is inserted into a batch queue**
- **The scheduler, schedules the task to run in one (or more) of the compute nodes that are "free"**
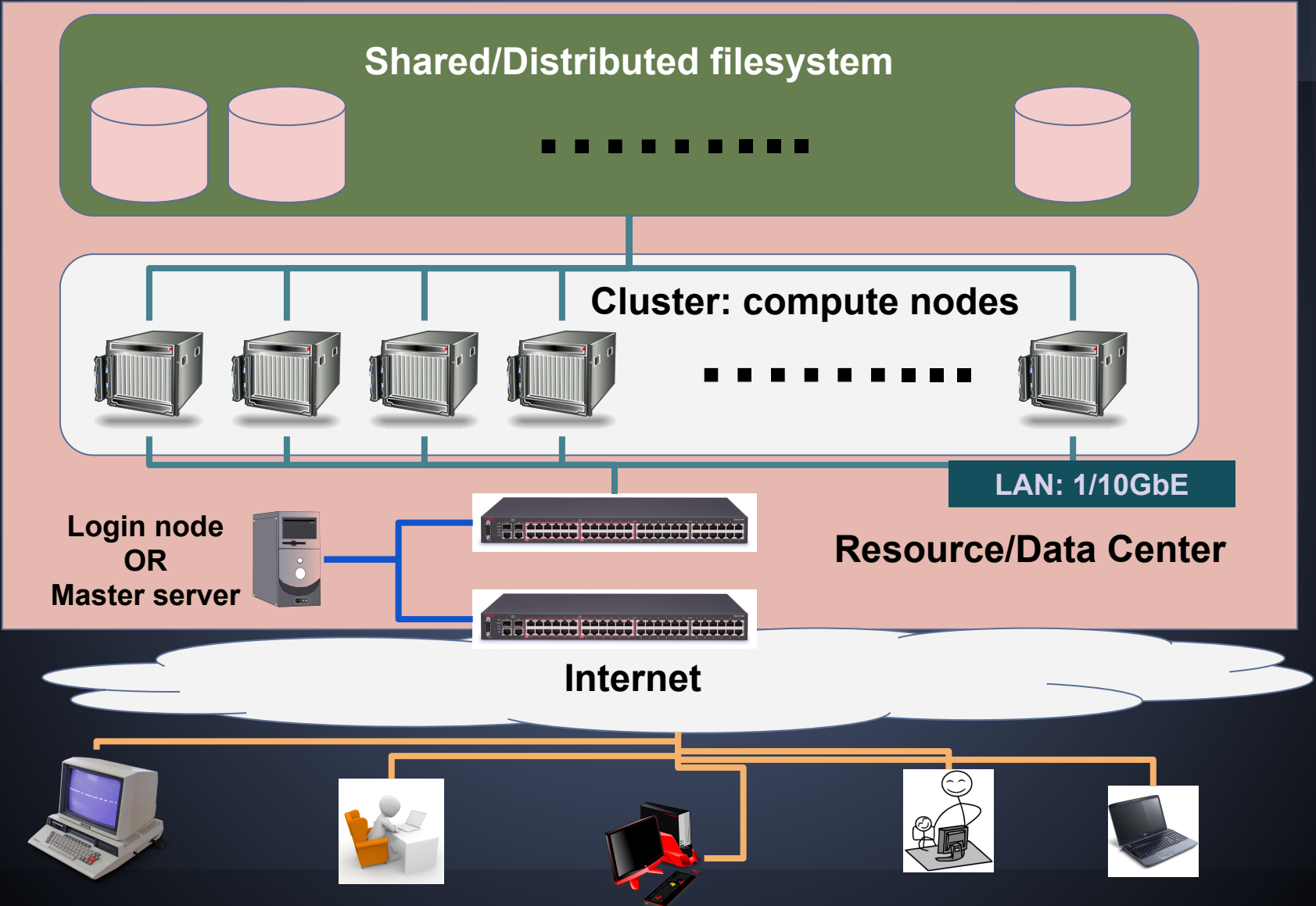- **If there are no free compute nodes the will stay in "wait" for free node(s)**

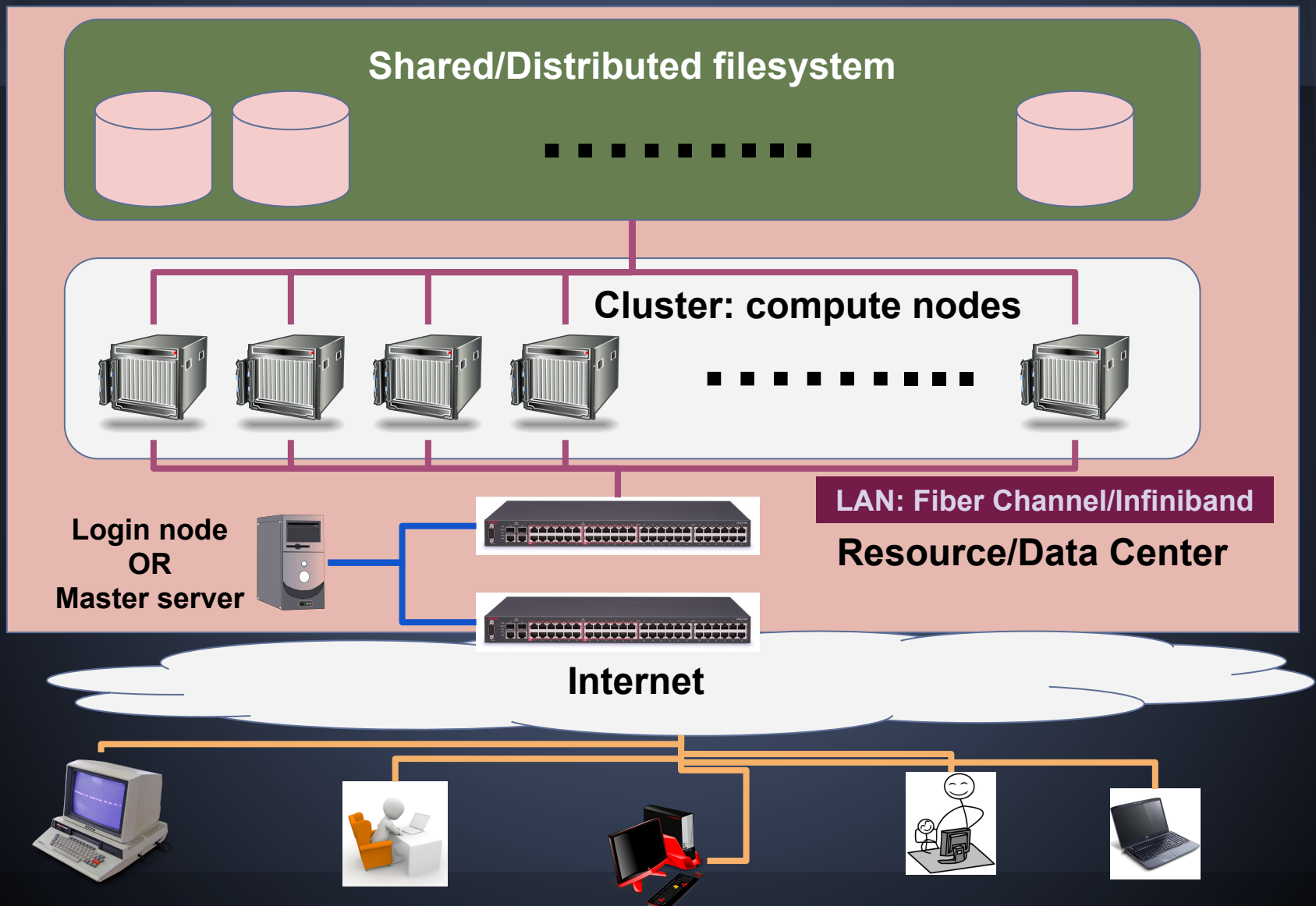# Batch clusters: Typical architecture



**Shared/Distributed filesystem**

**Cluster: compute nodes**

**Login node OR Master server**

**Resource/Data Center**

**Internet**

# Batch clusters: High Throughput Computing



**Shared/Distributed filesystem**

**Cluster: compute nodes**

**LAN: 1/10GbE**

**Login node OR Master server**

**Resource/Data Center**

**Internet**

# Batch clusters: High Performance Computing

**Shared/Distributed filesystem**

▪ ▪ ▪ ▪ ▪ ▪ ▪ ▪ ▪ ▪

**Cluster: compute nodes**

▪ ▪ ▪ ▪ ▪ ▪ ▪ ▪ ▪

Login node
OR
Master server

**LAN: Fiber Channel/Infiniband**

**Resource/Data Center**

**Internet**

# Batch clusters: Types - HTC

➢ High Throughput Computing (HTC):

- Computing paradigm that focuses on the efficient execution of a large number of loosely-coupled tasks.

- Adequate both for data intensive (I/O bound) and compute intensive (CPU bound) applications
- Adequate for serial applications.
- Embarrassingly parallel applications.

For example, processing/analysis of independent events ↦

IF:

you have 1000 events, and 100 CPUs

THEN:

distributing the processing of 10 events/CPU would yield a gain of 100 over a serial processing of all events in a single CPU

# Batch clusters: Types - HPC

➢ High Performance Computing (HPC):
  ○ Focus on tightly coupled parallel jobs and fast job execution.
  ○ Main difference in HW with respect to HTC, LAN is "**low latency**" such as Infiniband.
  ○ Adequate for compute intensive applications (CPU bound)
  ○ Adequate for parallel applications:

➢ Very common making use of **parallel programing**, such as using some implementation of MPI (Message Passing Interface) standard.

➢ Processes/Tasks need to communicate (send/receive messages) from other Processes/Tasks.

**1/10 GbE over TCP/IP    Latency ~ 10-100μs**
**Infiniband 10 - 100 Gb/s   Latency ≲ 1μs**

# Grid computing: Introduction

➢ Federation of clusters that are **geographically** distributed:
  ○ Each cluster is independent from the others: increase in heterogeneity with respect to a single cluster.
  ○ It has different administrative domains and policies
  ○ BUT, the users/researchers want a single way of "interaction" with all resources/clusters of the Grid:
    ■ Common APIs, CLIs
    ■ Common/single Authentication and Authorization system

# Grid computing: Architecture

# Grid computing: Grid middleware I

**Grid Middleware**

- **Common Authentication mechanism**:
  - X.509 certificates
  - Certification Authorities
- **Common Authorization mechanism**:
  - Users grouped by Virtual Organizations (VOs)
  - Resource providers authorize VOs to access and use their resources (computing and storage)
- **Compute Element (CE)**:
  - Frontend service exposing the local computing cluster to users through a common API/CLI
- **Storage Element (SE)**:
  - Frontend service exposing the local storage system to users through a common API/CLI

# Grid computing: Grid middleware II

**Grid Midleware**

- **Information services (IS)**:
  - Gather and publish information about the resources
- **Global data catalogs**:
  - Global view of data/files that are spread through several Storage Elements
  - Provide information about the physical location of the files
- **Orchestrator service/Resource Broker**:
  - Schedules compute tasks to Compute Elements based on the Information service and authorization policies (supported VOs)
- **File Transfer Service**:
  - Management of data movement/transfer between resource providers

13

# Grid computing: Grid middleware III

# Cloud Computing I

"(…) a model for enabling ubiquitous, convenient, on-demand network access to a shared pool of configurable computing resources (e.g., networks, servers, storage, applications, and services) that can be rapidly provisioned and released with minimal management effort or service provider interaction."

➢ **On-demand Self Service**
- ○ The users are able to provision and manage their own computing environment according to their needs, without further intervention from the provider.

➢ **Elastic Provisioning and Scalability**
- ○ The cloud model tries to deliver easily and rapidly the resources to the users, in a short-deadline basis.
- ○ Users are able to scale in and out their infrastructure so as to satisfy the real demand, not only by increasing their capacity, but also by shrinking it whenever it is not needed.

15

# Cloud Computing II

➢ **Metered Usage and Billing**

  ○ Resources accounted by their usage, rather than following a subscription mode.

➢ **Multi-tenancy and Dynamic Resource Pooling**

  ○ Ability for a software or provider to deliver a service to several parties simultaneously.

  ○ Services owned by several users are being co-located in the same resources.

  ○ Tenants resources are isolated from each other.

  ○ Each tenant manages creates and manages it's own compute, storage and local network.

  ○ Important for organizations supporting multiple projects/groups, and service providers supporting multiple users.

# Cloud Computing: <u>X</u> <u>a</u>s <u>a</u> <u>S</u>ervice models

**Increase of abstraction**

**SaaS** — Cloud applications

**PaaS** — Cloud services

**IaaS** —
Orchestration layer

Provisioning layer

Compute Resources | Network Resources | Storage Resources

# Cloud Computing: Classification I

➢ **Infrastructure as a Service (IaaS):**
  ○ Lowest level of abstraction
  ○ Considered as the foundation of the cloud model.
  ○ Offers its infrastructure resources: computing, networking and storage.
  ○ Users can deploy its own OS, software, network configuration, etc.
  ○ Abstracts the underlying fabric (physical resources) into a uniform resource layer:
    ■ **Virtualization** or encapsulation the raw resources.
    ■ Users get **transparent access** to this layer as if they were using the bare metal resources.
    ■ Able to deploy any infrastructure on top of it without the extra burden of directly managing the different physical resources.

# IaaS: Openstack CFM

**Instantiate a VM machine**



**Openstack (Cloud Management Framework)**

Web - GUI/Dashboard

CLI

```
$ nova boot ...
```

# IaaS: Openstack CFM

**Instantiate a VM machine**



**Openstack (Cloud Management Framework)**

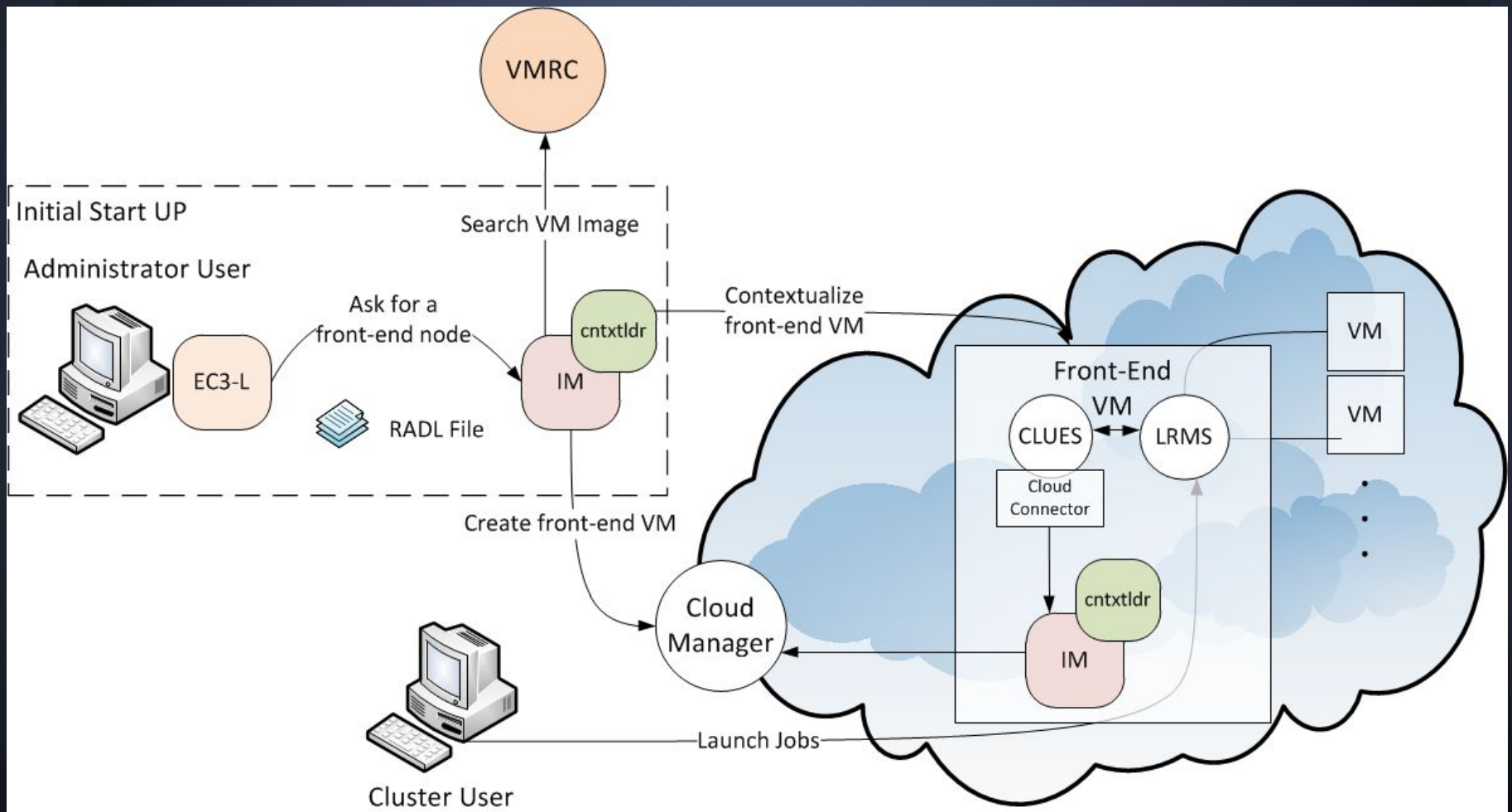Web - GUI/Dashboard

CLI

# IaaS: Openstack CFM

## Login into the VM

# Cloud Computing: Classification II

➢ **Platform as a Service (PaaS)**:
  ○ Second step in the abstraction level
  ○ Resources coming from an IaaS are **composed** so that they can be consumed by the users without requiring the management or the knowledge of the underlying infrastructure.
  ○ Offer an environment where a user **can deploy and manage its applications** using the libraries, software, tools, APIs, etc. supported by the provider
  ○ Makes possible to deliver **complex applications and services** involving different components to end-users:
    ■ No need of direct managing of the machine configurations and deployment
    ■ Allows to define the requirements of those applications, so that the platform layer is **able to orchestrate the resources**.

# PaaS: Infrastructure Manager (IM)

**From Univ. Valencia: http://www.grycap.upv.es/im/index.php**

**EC3 (Elastic Cloud Computing Cluster)**

# Cloud Computing: Classification III

➢ **Software as a Service (SaaS):**
- ○ The highest level of abstraction.
- ○ Comprises the applications that are running on top of a cloud infrastructure.
- ○ Access to SaaS applications are normally addressed using ad-hoc thin clients executed inside web browsers or applications that are executed on tablets or smartphones, directly addressing the end user.
- ○ Change in paradigm: FROM buy software TO buy service:
  - ■ Delegate the software management (to service provider) and focus on the use of the service/software
  - ■ One example: the Primavera ERP http://www.famcorp.pt/publico/Solu%C3%A7%C3%B5es-Solu%C3%A7%C3%B5es%20Online-Primavera%20SaaS.aspx

# Cloud Computing: Classification III

# SaaS: Galaxy science portal

# SaaS: Public provider

# Grids/Clusters versus Clouds

| Grids/Clusters | Clouds |
|---|---|
| Fixed environments: OS, applications | Flexible environments: users choose OS, applications, through virtualization (1) |
| Amount of resources are fixed apriori | Amount of resources are elastic: can increase or decrease according to users needs |
| Applications/tasks are executed during some fixed amount of time | Can run applications/tasks during a fixed amount of time, but can also run long term services such as web and scientific portals, databases, etc. |
| Applications are scheduled to batch queues | On demand "almost" real time provisioning of resources |
| Grid is a federation of clusters (resource providers) | Federation of clouds is still quite difficult and a strong hot topic. One such case is the EGI FedCloud infrastructure (2) |

(1)    **It's also possible to deploy flexible environment in bare metal (physical) nodes**
(2)    **EGI - European Grid Infrastructure: Provides a Grid <u>and</u> a Federated Cloud infrastructure at European level**

# EGI Federated cloud: Architecture

**EGI Cloud Infrastructure Platform**

- Instance Mgmt — OCCI
- Storage Mgmt — CDMI

OVF

**Providers Cloud Management Framework**
(Based on standards, cloud middleware agnostic)

GLUE2 · GSI · SAM · UR

- Service Registry
- Information Discovery
- Federated AAI
- Monitoring
- Accounting

- Helpdesk Support
- Security Coordination
- Training Outreach
- Sustainable Business Models

**EGI Core Platform**

**EFI Collaboration tools**

**EGI Cloud Service Marketplace**

# Distributed Data and Storage

# Distributed data and storage

➢ Data + Metadata
  - Object storage
  - Block storage
  - File storage and filesystems (POSIX)

**MetaData**

**Data**

# Distributed data and storage

➢ Data + Metadata
➢ **Object storage**
➢ Block storage
➢ File storage and filesystems (POSIX)



**MetaData**

**[Key=Value]** → ObjectID →

**Data**

**Binary block**

# Distributed data and storage

➢ Data + Metadata
➢ Object storage
➢ Block storage
➢ **File storage and filesystems (POSIX)**

**MetaData**

**INodes**

**Pointers to** →

**Data**

**Binary block**

# Storage: POSIX vs Objects

## POSIX filesystem

- **Directories**
  - Only metadata
  - List of filenames and corresponding inode number, etc.
- **Files**
  - Metadata (inodes)
  - Data

- Hierarchical - Tree structure

## Object storage

- **Containers**
  - Only metadata
  - Information about the objects contained in _this_ container
- **Objects**
  - Metadata AND data AND ObjID

- Flat structure
  - Horizontal scalability

# Storage: POSIX vs Objects

## POSIX (partial list)

- open
- read
- write
- close
- lseek
- llseek
- _llseek
- lseek64
- stat
- fstat
- stat64
- chmod
- fchmod
- access
- rename
- mkdir
- getdents
- fcntl
- unlink
- fseek
- rewind
- ftell
- fgetpos

- fsetpos
- fclose
- fsync
- creat
- readdir
- opendir
- fopendir
- rewinddir
- scandir
- seekdir
- telldir
- flock
- lockf
- lseekm
- lstat
- fstatat
- fopen
- fdopen
- freopen
- remove
- chown
- fchown
- fchmodat

- fchownat
- faccessat
- utime
- futimes
- lutimes
- futimesat
- link
- linkat
- unlinkat
- symlink
- symlinkat
- rmdir
- mkdirat
- getxattr
- lgetxattr
- fgetxattr
- xetxattr
- lsetxattr
- fsetxattr
- listxattr
- llistxattr
- flistxattr
- removexattr

## Objects (RESTful API)

**PUT** ("write"): PUT the **object** into the storage

**GET** ("read"): GET the **object** from the storage

**DELETE:** delete the **object** which is the file

**POST**: create, update, delete **metadata**

**HEAD**: returns an object's **metadata**

# Object Storage

➢ Object storage is a storage architecture that manages data as objects
➢ Each object typically includes
  ○ the data itself
  ○ a variable amount of metadata
  ○ a globally unique identifier: Object ID.
➢ Access through RESTful API


  ○ Example: Ceph object storage offers access through S3 and SWIFT APIs

# Block Storage

A Device:

- ➢ Harddisk (and/or disk partition)
- ➢ CD
- ➢ DVD
- ➢ Disk array
- ➢ ...

On the cloud:

- ➢ On demand request for a disk volume
  - ○ Attach to a VM instance, as local storage to increase the storage capacity available to the instance.
  - ○ Can be formated with whatever filesystem the user wants.

# Network File System: NFS

➢ Distributed filesystem:
   ○ Protocol originally developed by Sun Microsystems in 1984
   ○ Client/Server architecture.
   ○ Based on RPCs (Remote Procedure Calls)
   ○ Reads and writes on the client are **mapped** to read and writes on the server.
   ○ Only the server accesses the filesystem, managing all calls from the multiple clients.

# Parallel filesystems



Metadata

Clients

"queries"/"responses"

Data I/O

Data

# Parallel filesystems



Multiple
"queries"/"responses"

Multiple streams
Data I/O

# Lustre filesystem: Architecture

**Lustre file system components**

# Lustre filesystem I

**Layout Extended Attribute (EA) on MDT pointing to file data on OSTs**

# Lustre filesystem II

**Lustre client requesting file data**

# Lustre filesystem III
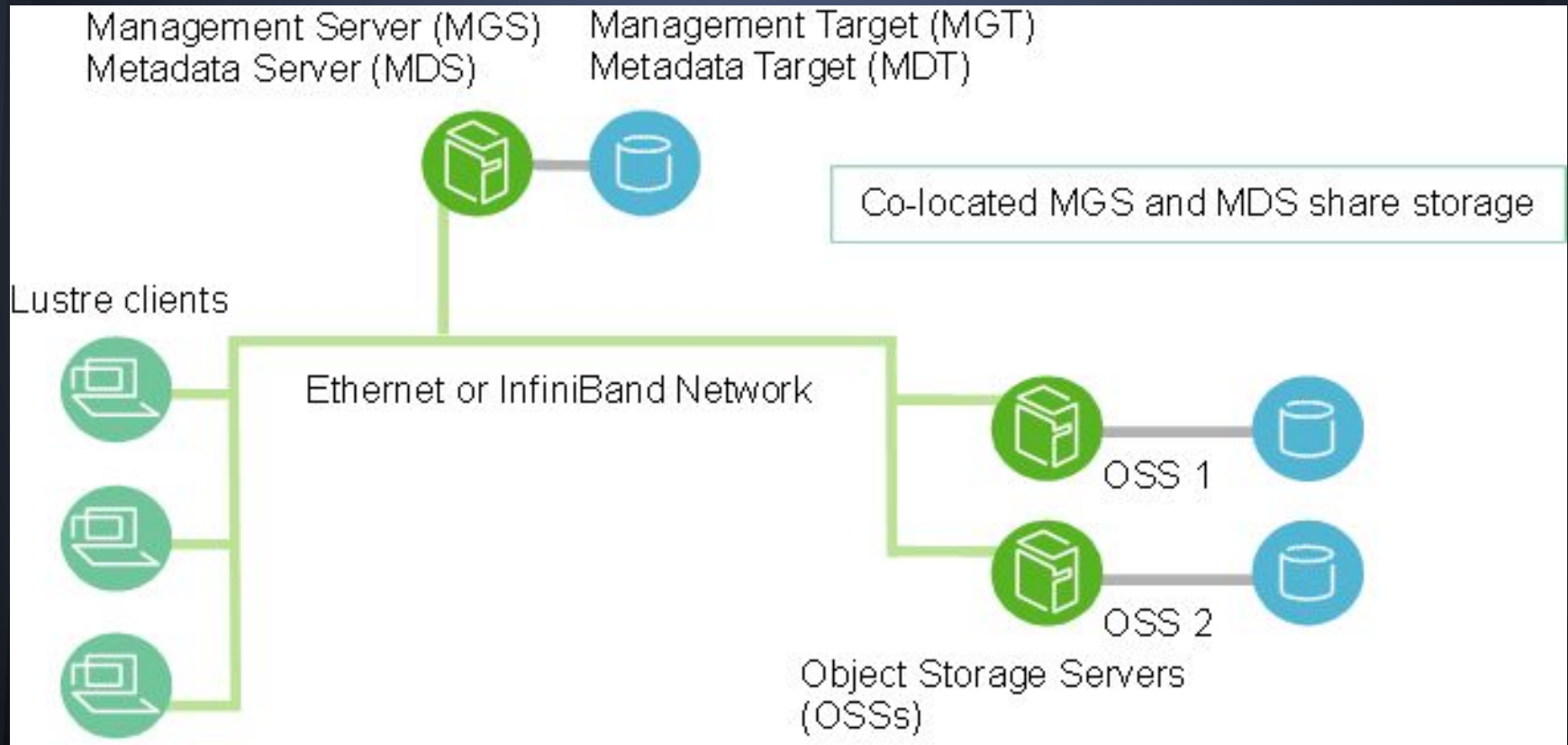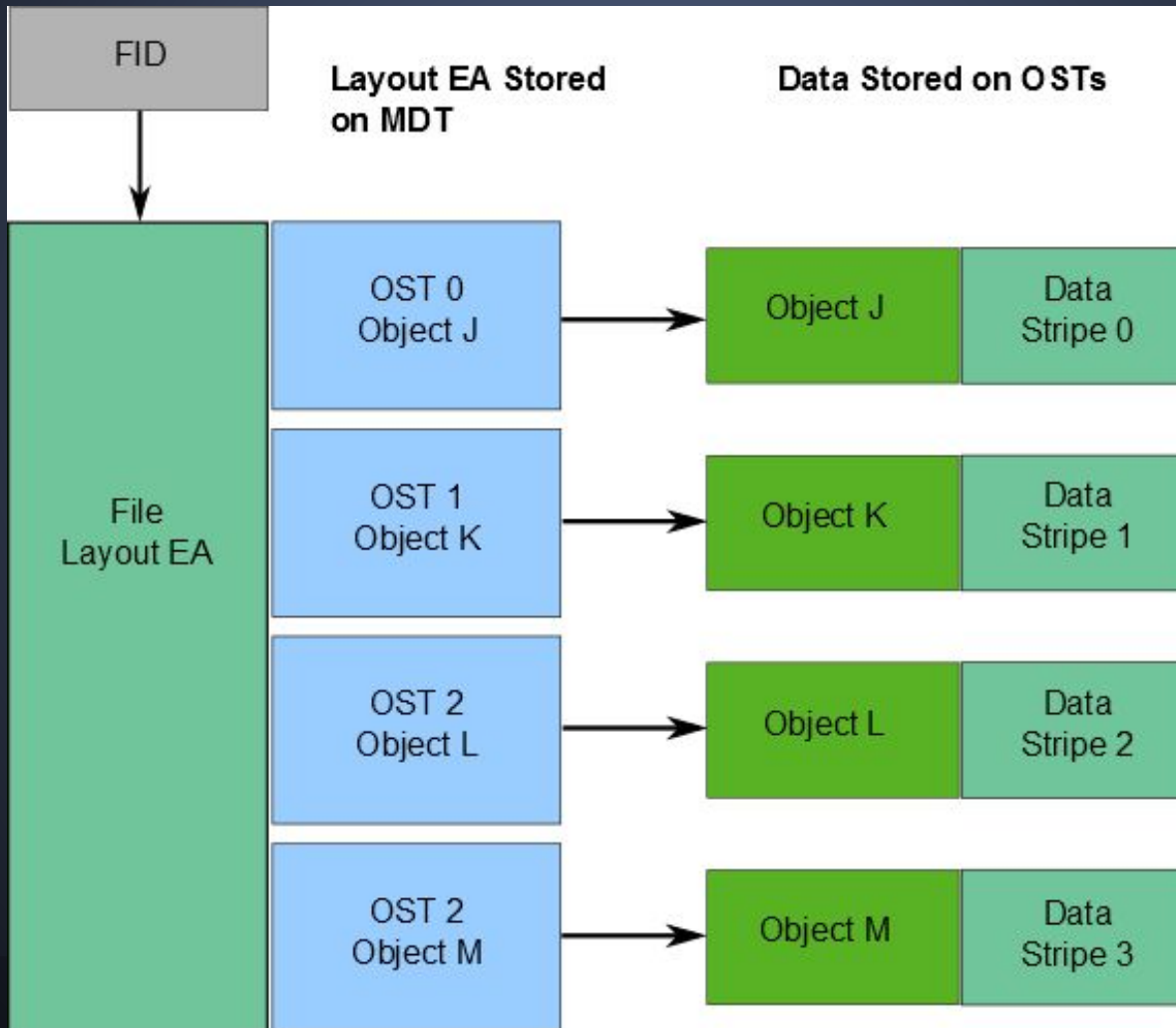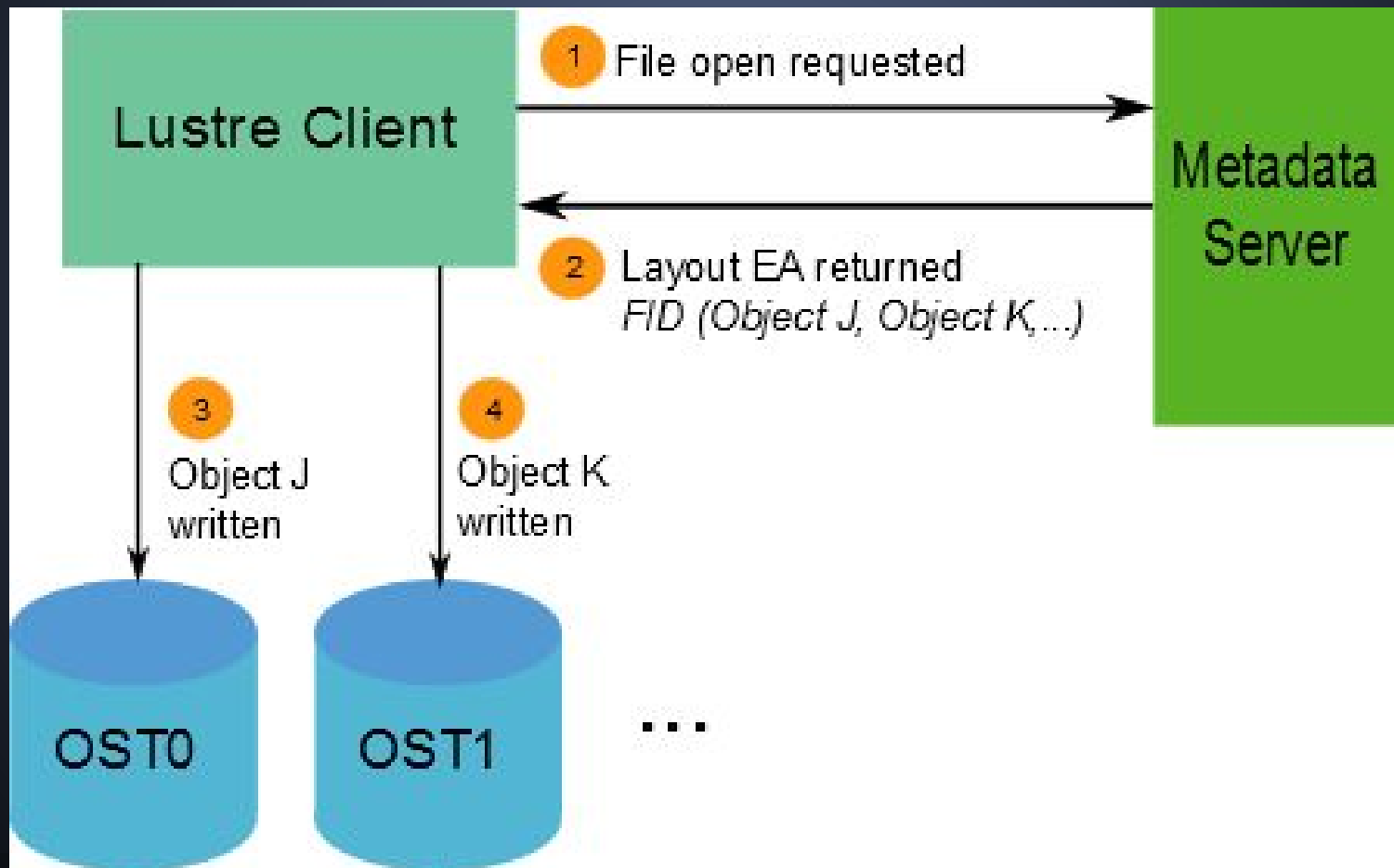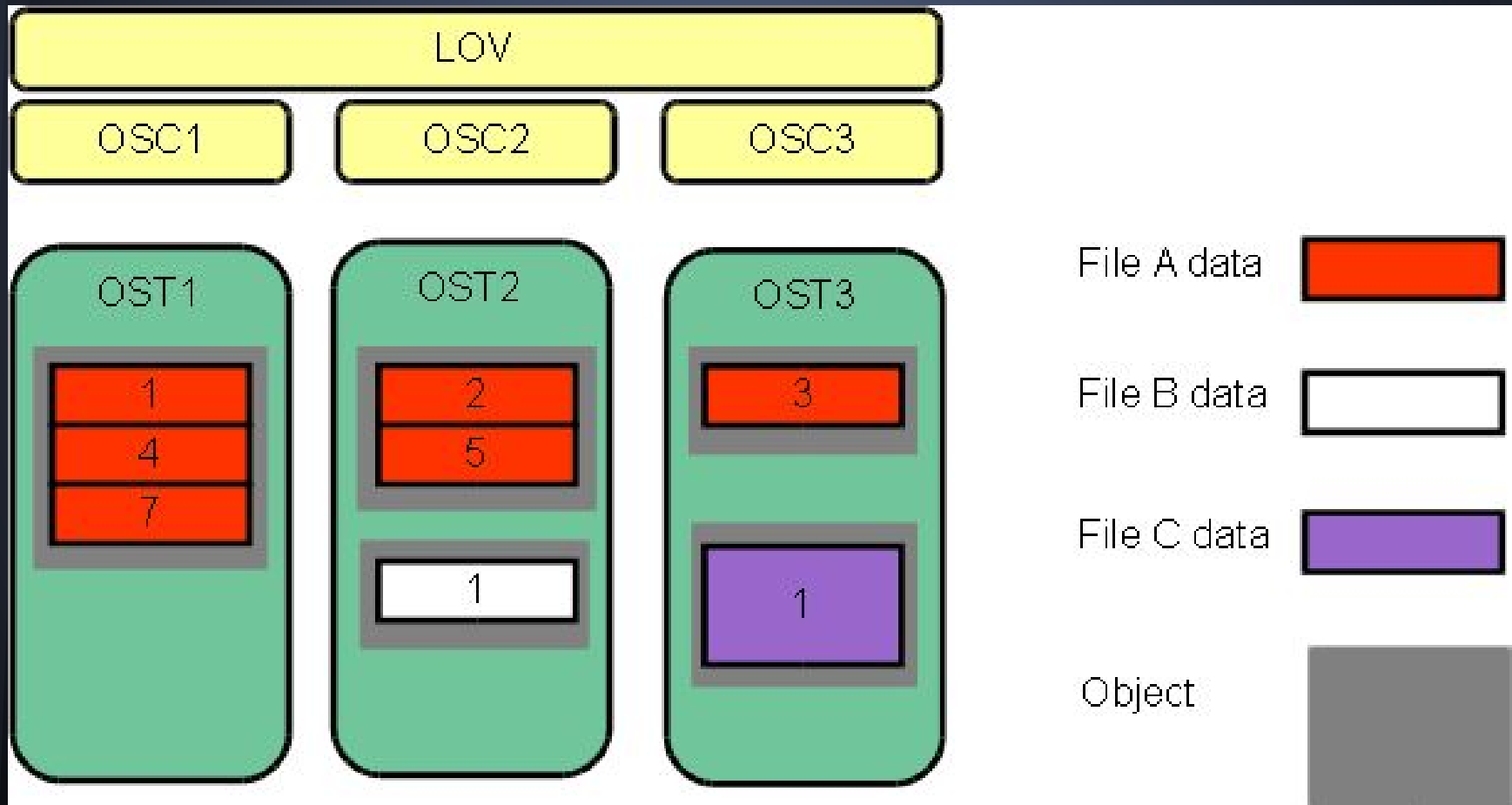
## File striping on a Lustre file system

# Ceph: Architecture

Ceph uniquely delivers object, block, and file storage in one unified system

# Ceph: Components

- ➢ **OSDs:**
  - ○ Stores data, handles data replication, recovery, etc.
  - ○ Provides some monitoring information to Ceph Monitors
- ➢ **Monitors:**
  - ○ Maintains maps of the cluster state, including the monitor map, the OSD map, etc.
  - ○ Maintains a history (called an "epoch") of each state change in the Ceph Monitors, Ceph OSD Daemons, etc.
- ➢ **MDSs:**
  - ○ Stores metadata on behalf of the Ceph Filesystem (POSIX).
  - ○ Ceph Block Devices and Ceph Object Storage do not use MDS.
  - ○ Make it feasible for POSIX file system users to execute basic commands like ls, find, etc. without placing an enormous burden on the Ceph Storage Cluster.

# Ceph: Object storage I
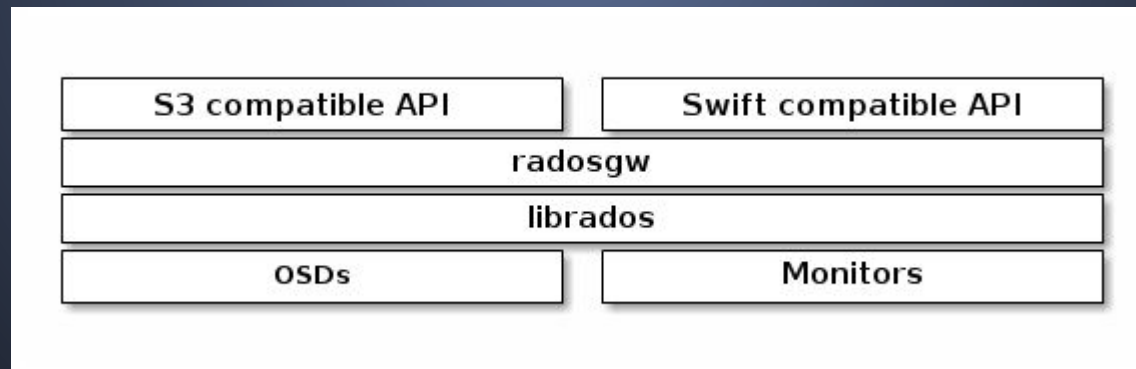
➢ Ceph Object Gateway - **radosgw**:
- ○ Object storage interface built on top of librados to provide applications with a RESTful gateway to Ceph Storage Clusters.

➢ Ceph Object Storage supports two interfaces:
- ○ **S3-compatible**: Provides object storage functionality with an interface that is compatible with a large subset of the Amazon S3 RESTful API.
- ○ **Swift-compatible**: Provides object storage functionality with an interface that is compatible with a large subset of the OpenStack Swift API.

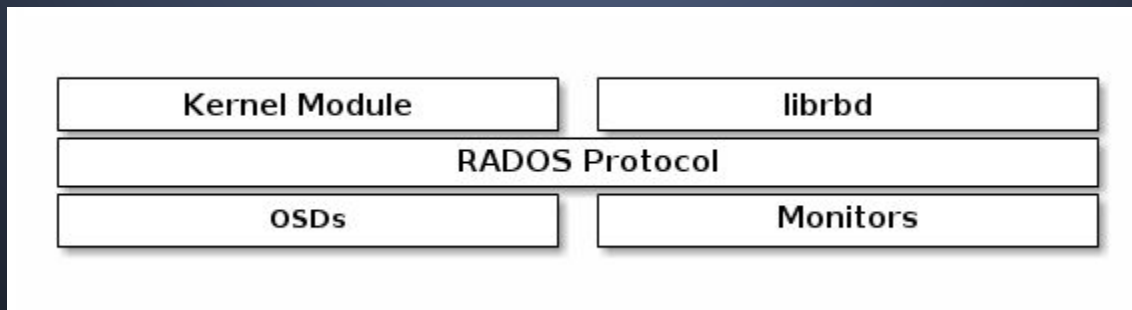| ID | Binary Data | Metadata | |
|------|-------------------------------------------|--------|--------|
| 1234 | 01010101010101001101010010 | name1 | value1 |
| | 01011000010101001101010010 | name2 | value2 |
| | 01011000010101001101010010 | nameN | valueN |

# Ceph: Object storage II

➢ Ceph Object Gateway - **radosgw**:
  ○ Has its own user management
  ○ It can store data in the same Ceph Storage Cluster used to store data from Ceph Filesystem clients or Ceph Block Device clients.
  ○ The S3 and Swift APIs share a common namespace, so you may write data with one API and retrieve it with the other.
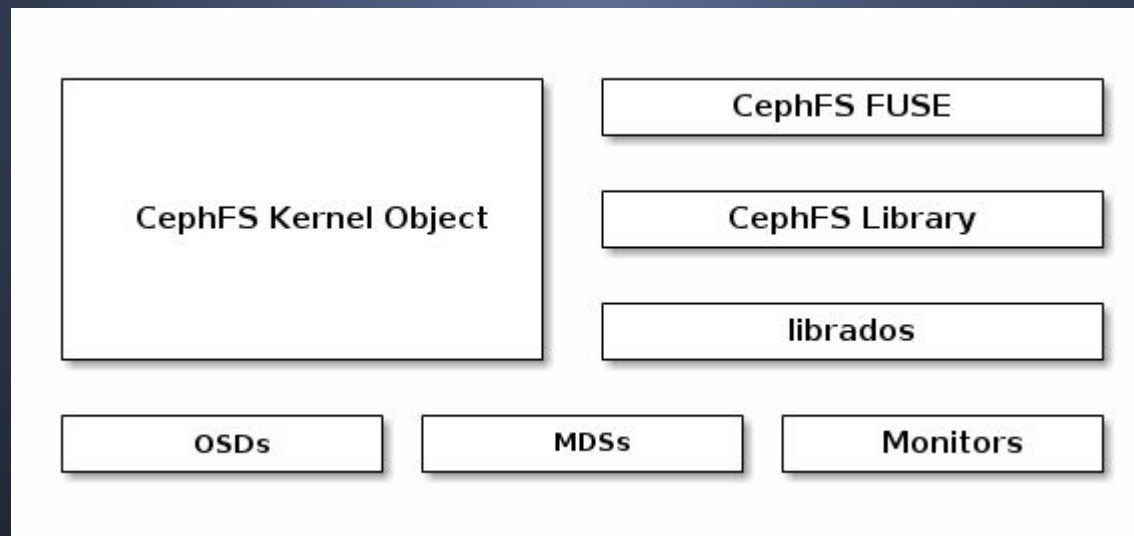
| S3 compatible API | Swift compatible API |
|---|---|
| radosgw | |
| librados | |
| OSDs | Monitors |

# Ceph: Block storage

➢ Ceph block devices are:

   ○ Thin-provisioned, resizable and store data striped over multiple OSDs in a Ceph cluster.

   ○ Leverage RADOS capabilities such as snapshotting, replication and consistency.

➢ Ceph's RADOS Block Devices (RBD) interact with OSDs using kernel modules or the librbd library.
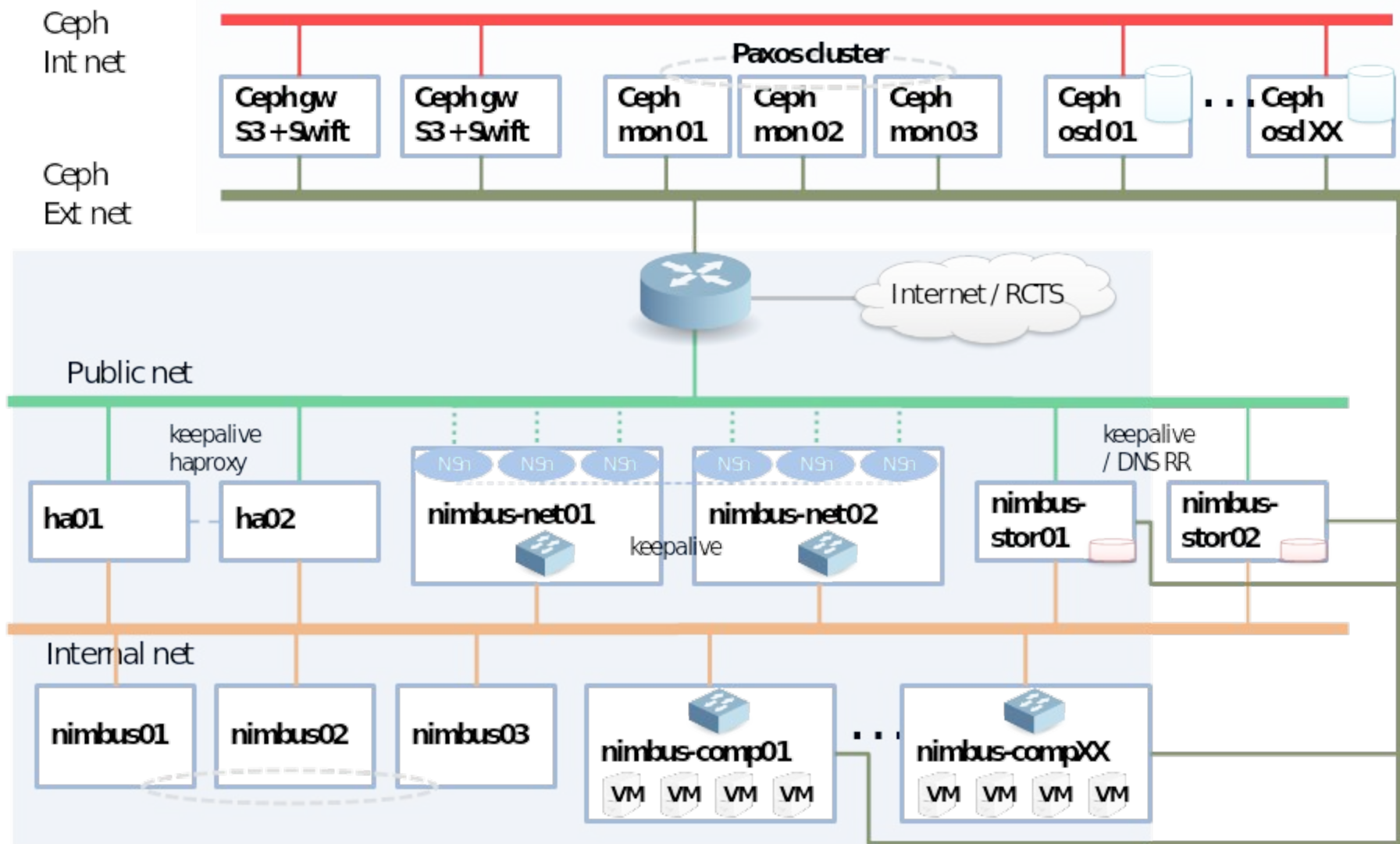
| Kernel Module | librbd |
|---|---|
| RADOS Protocol | |
| OSDs | Monitors |

# Ceph: Filesystem

➢ The Ceph Filesystem (Ceph FS):

○ POSIX-compliant filesystem that uses a Ceph Storage Cluster to store its data.

○ Uses the same Ceph Storage Cluster system as Ceph Block Devices, Ceph Object Storage with its S3 and Swift APIs, or native bindings (librados).

# Putting it all together
# The IaaS Openstack infrastructure

# Thanks!!

# Questions??