

The logo for CEA (Commissariat à l'énergie atomique et aux énergies alternatives) consists of the lowercase letters 'cea' in a white, sans-serif font. A horizontal green line is positioned below the letters. The logo is enclosed within a white L-shaped border.The Phobos logo features a stylized planet with a cyan circular body and a ring that transitions from yellow at the top to red at the bottom. To the right of this graphic, the word 'Phobos' is written in a large, black, sans-serif font. The entire logo is set against a light pink rectangular background with a subtle pattern of small, darker pink circles.

Phobos

Phobos: an open-source object store implementing tape library support

Patrice LUCAS, patrice.lucas@cea.fr

The logo for CEA (Commissariat à l'énergie atomique et aux énergies alternatives) consists of the lowercase letters 'cea' in a white, sans-serif font. A horizontal green line is positioned below the letters. The logo is enclosed within a white L-shaped border that forms a partial square.The Phobos logo features a stylized planet with a cyan circular body and a ring system. The ring is a thick, curved line that transitions in color from yellow at the top to red at the bottom. To the right of this graphic, the word 'Phobos' is written in a large, bold, black, sans-serif font. The entire logo is set against a light pink rectangular background with a subtle pattern of small, darker pink circles.

Phobos

Phobos: an open-source object store implementing tape library support

Patrice LUCAS, patrice.lucas@cea.fr

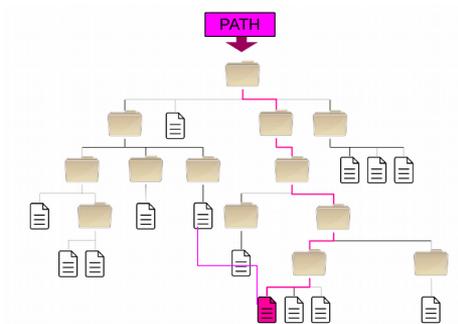
Next scale of mass storage

- Exaflop supercomputers in the 2020's
- Huge amounts of data to ingest: petabytes per day
- Huge amounts of data to store: exabytes

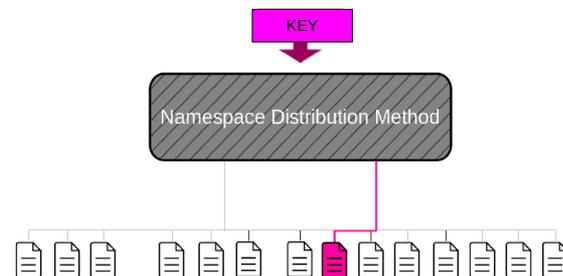


	Today	Tomorrow
Daily production	Hundreds of TB	Petabytes
Storage system capacity	Hundreds of PB	Exabytes

Suppressing POSIX filesystem's bottlenecks



*Addressing entries
in a POSIX file system*



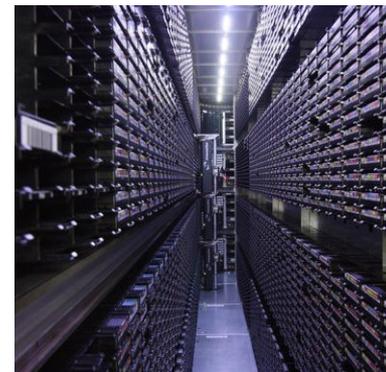
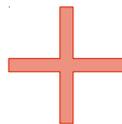
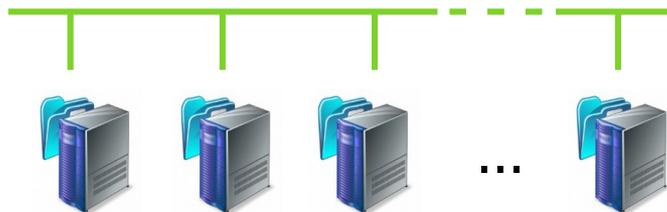
*Addressing entries
in an object-store*

- Object stores have proved their scalability
- Widely adopted for Internet services, Cloud computing, social networks...



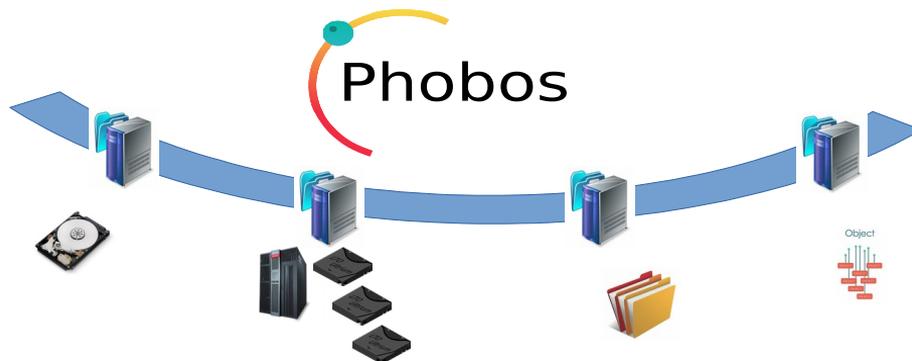
Needs for extremely scalable storage systems at a reasonable price

- Object store: horizontal scalability
- Tape library: safe long term storage at low cost



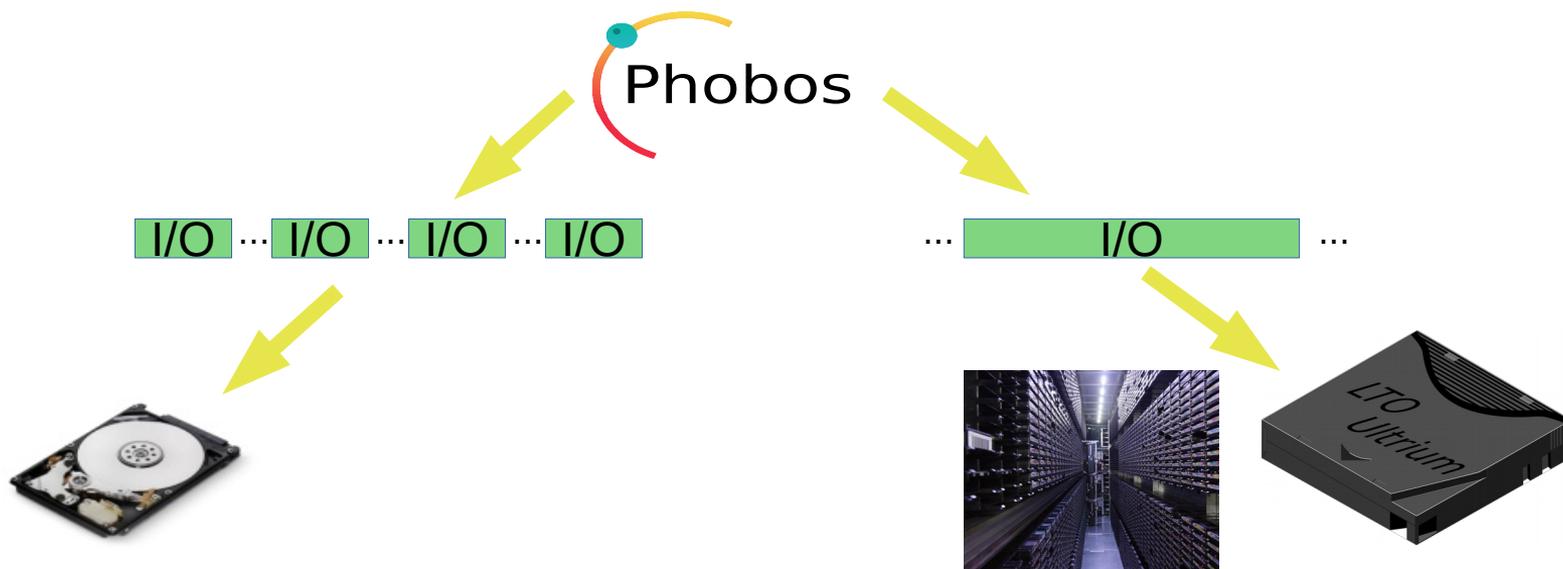
Phobos: Parallel Heterogeneous Object Store

- Manages a distributed set of file systems on different storage media technologies
 - hard drives or magnetic tapes



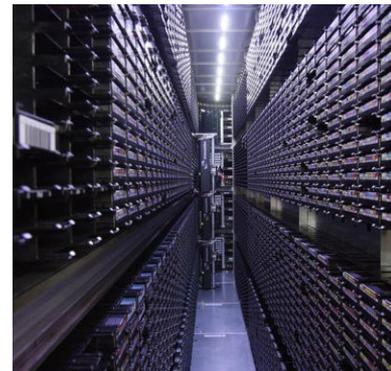
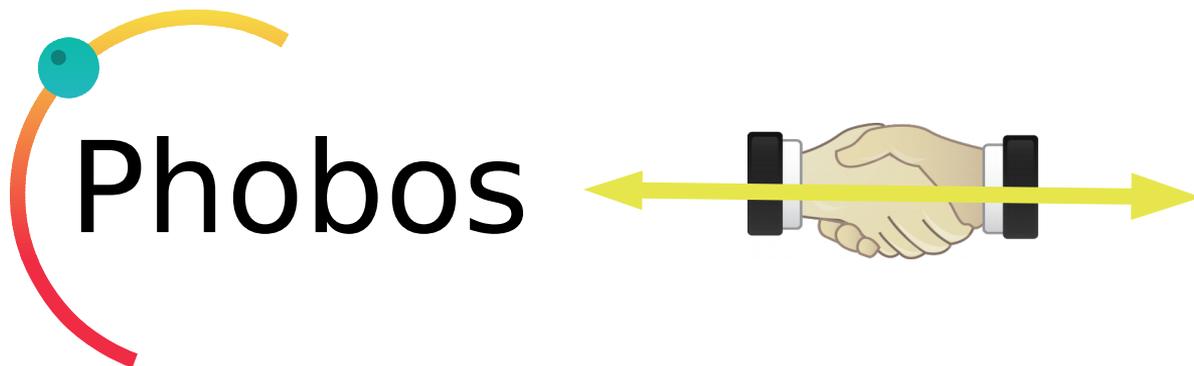
Phobos: optimizes I/O access depending on the storage technology

- minimizing data syncs for magnetic tapes



Phobos: directly manages a tape library (SCSI)

- Supports all common tape drives model



Design guidelines

- Scalability and fault-tolerance
- Based on open formats, open protocols, interoperable
 - E.g. LTFS as tape filesystem (ISO/IEC 20919:2016)
- Simple and common interfaces (REST, object stores API)
- Simple administration (intuitive, admin-friendly CLI)
- Light, easy to deploy, easy to maintain
 - As of today: 14k lines of C, 2.5k lines of Python

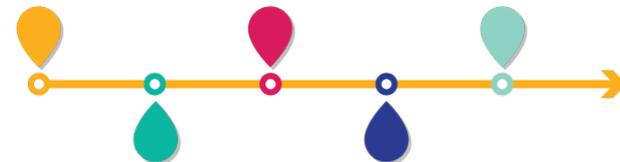


History of the project

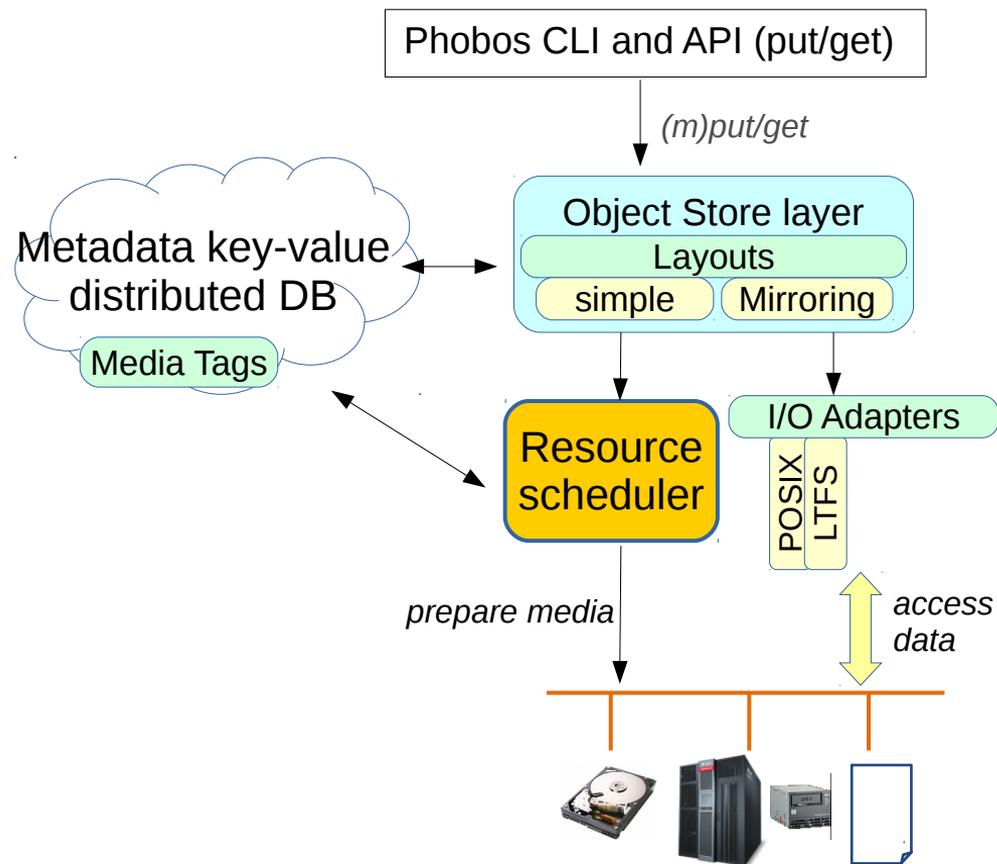
- 2013: high level design
- 2014-2015: development of the initial version

Scope:

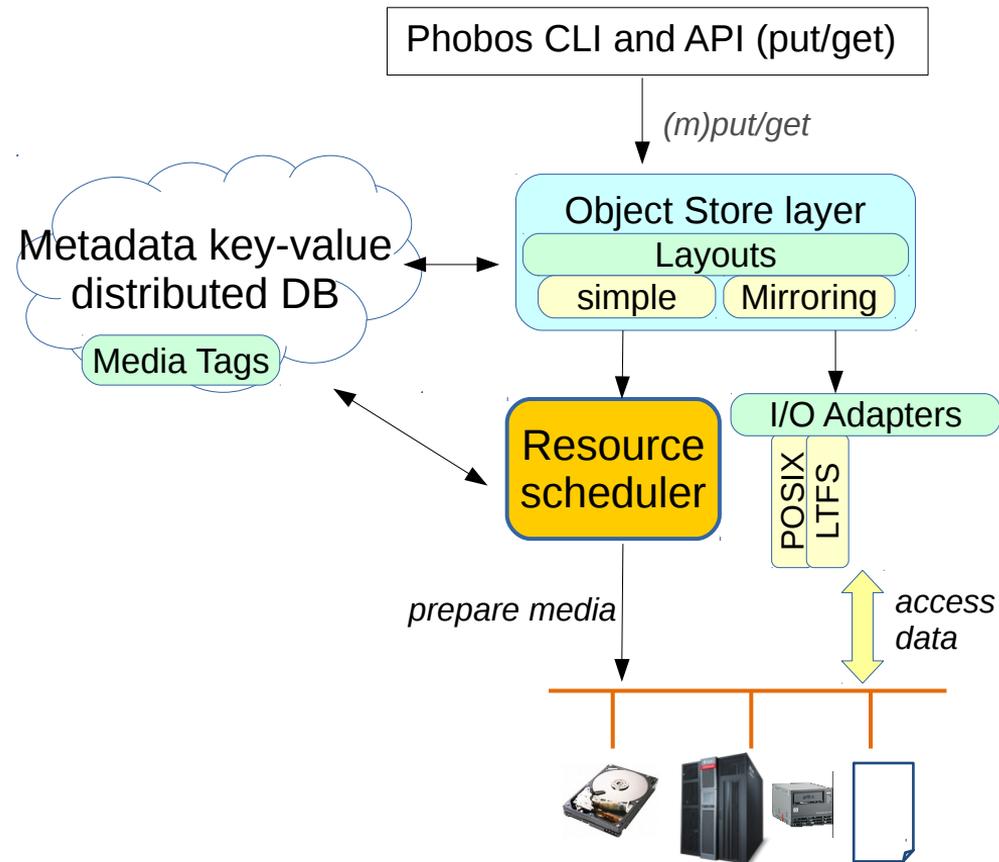
- Storage on tape, or in a filesystem
- SCSI-controlled tape library and LTO drives
- Single server
- 2016: Phobos v1.0 in production
 - Multi-Petabyte storage of genomics data
 - IBM TS3500 library, LTO5/6 drives
- 2019: Phobos made available on github as open-source (LGPL v2.1)
- 2020: Implementation of an S3 front-end (collab. with ICHEC and DDN) + required features in the Phobos core
- Next steps:
 - Advanced IO scheduling
 - Parallelization across multiple servers



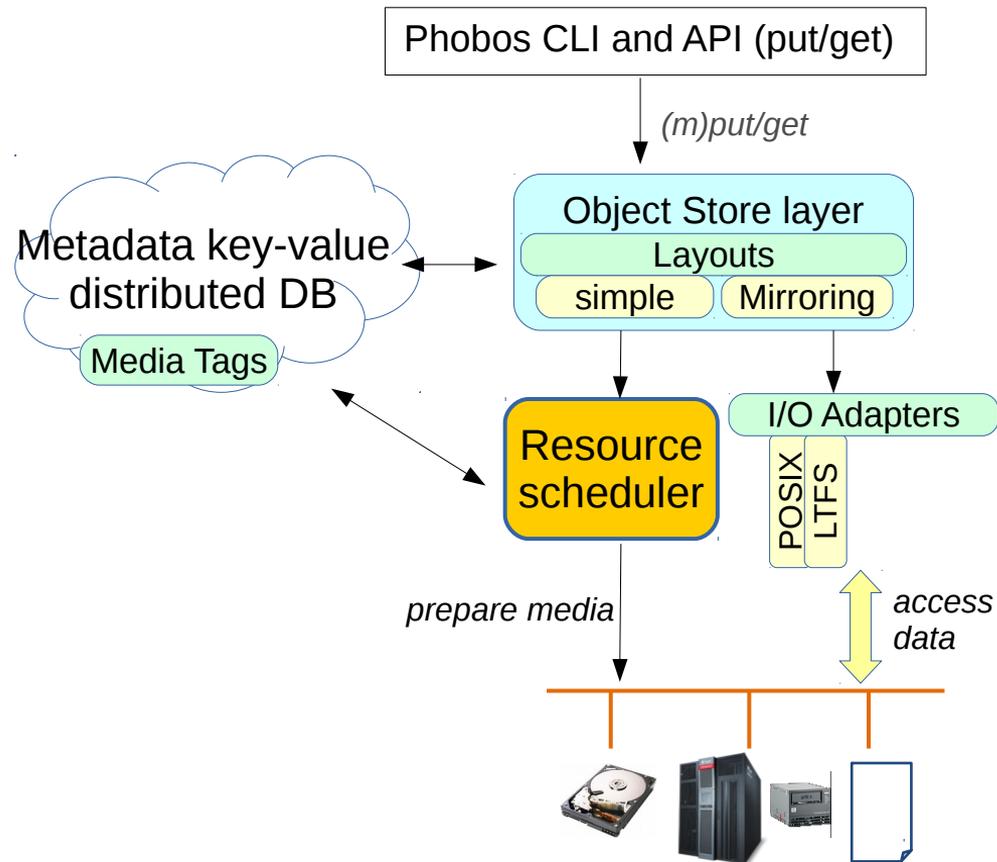
FRANCE GÉNOMIQUE



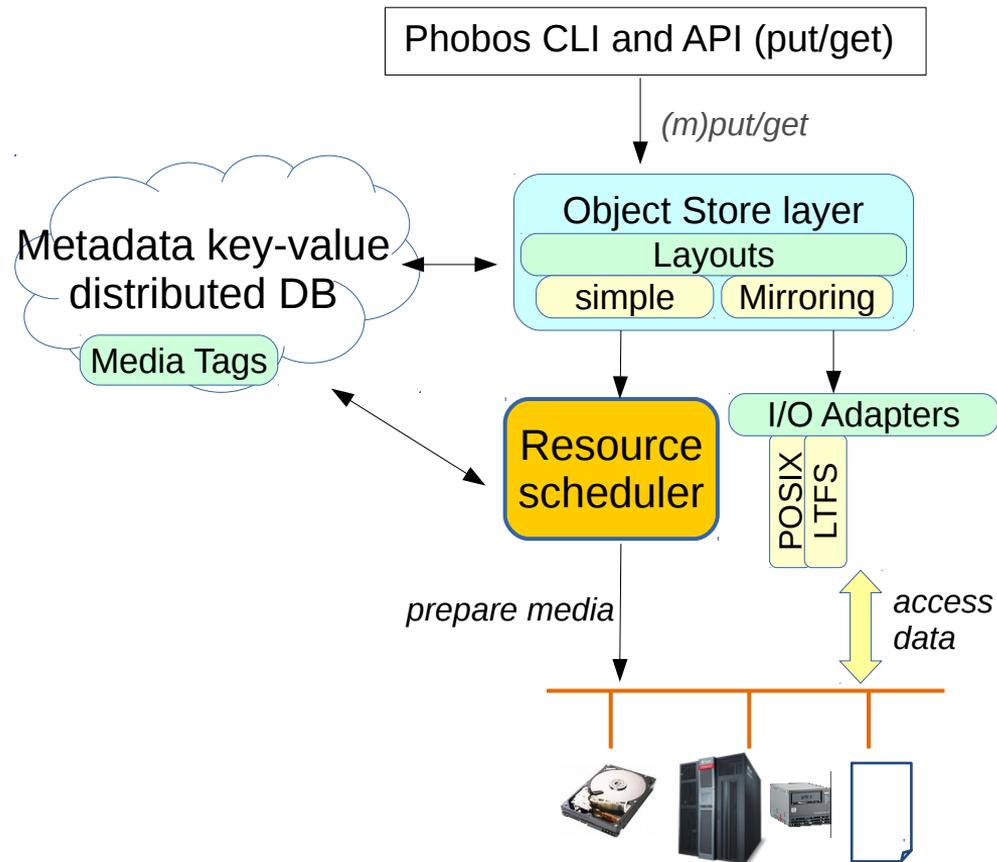
- IO adapters: multiple storage technologies (Posix, LTFS)



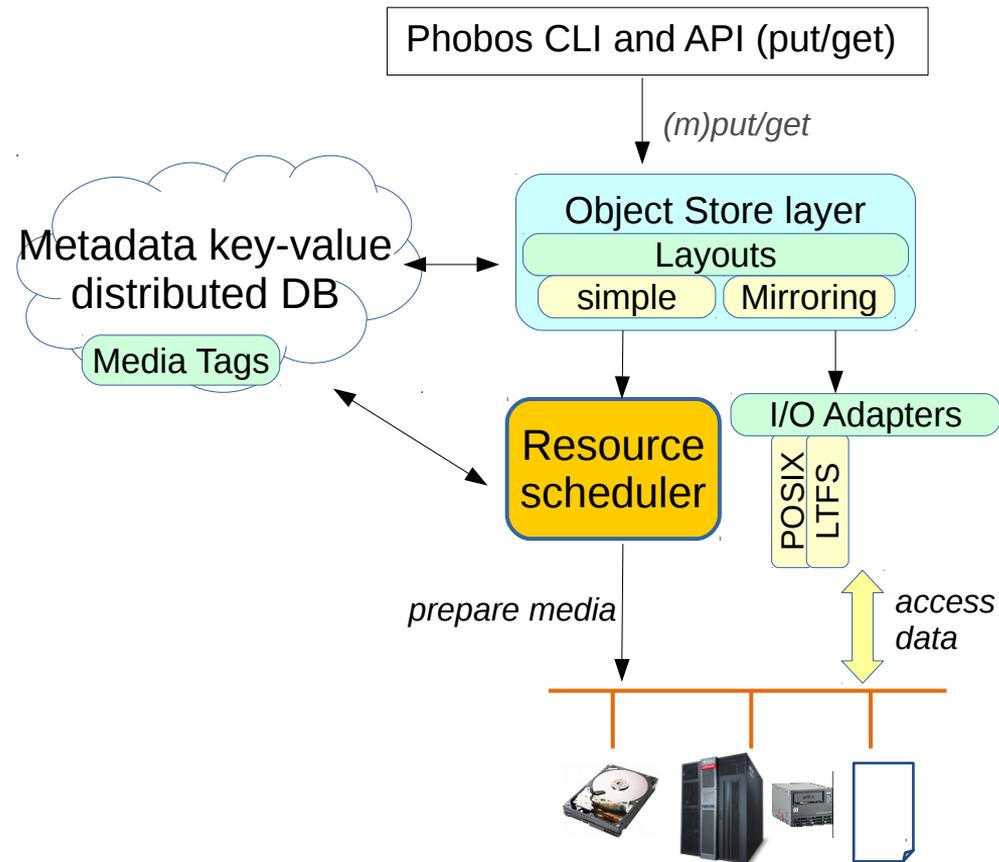
- IO adapters: multiple storage technologies (Posix, LTFS)
- Layout plugins: performance and fault-tolerance



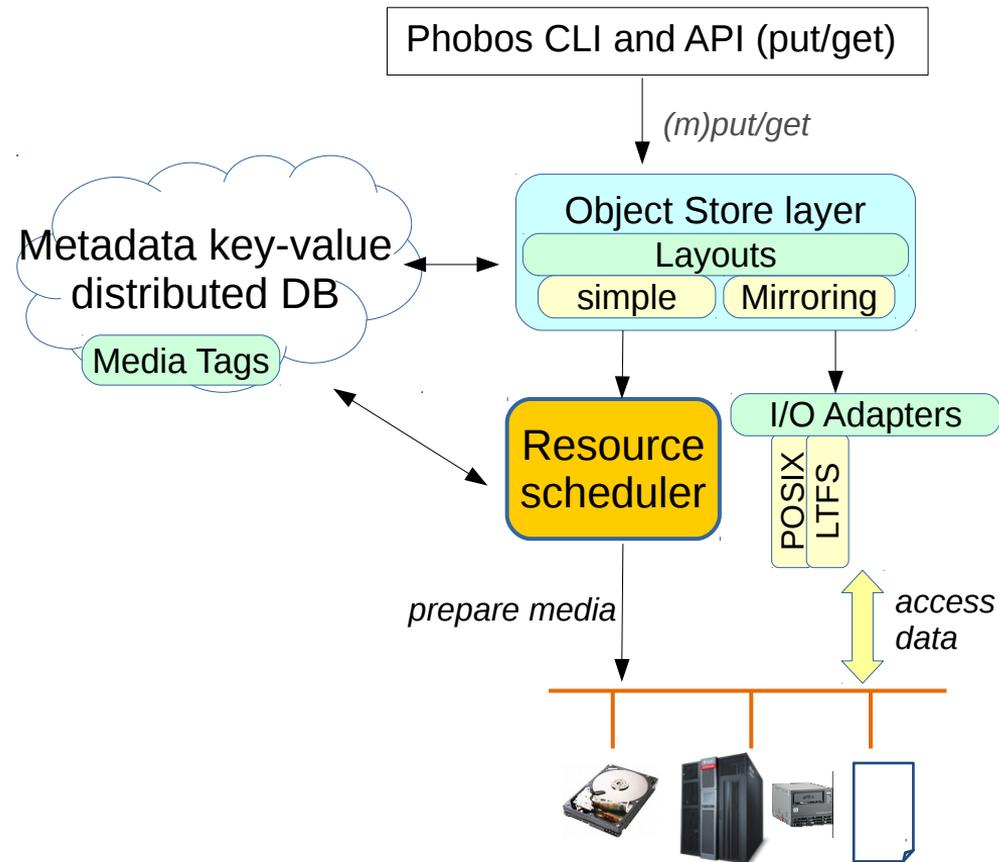
- IO adapters: multiple storage technologies (Posix, LTFS)
- Layout plugins: performance and fault-tolerance
- Tags: storage partitioning

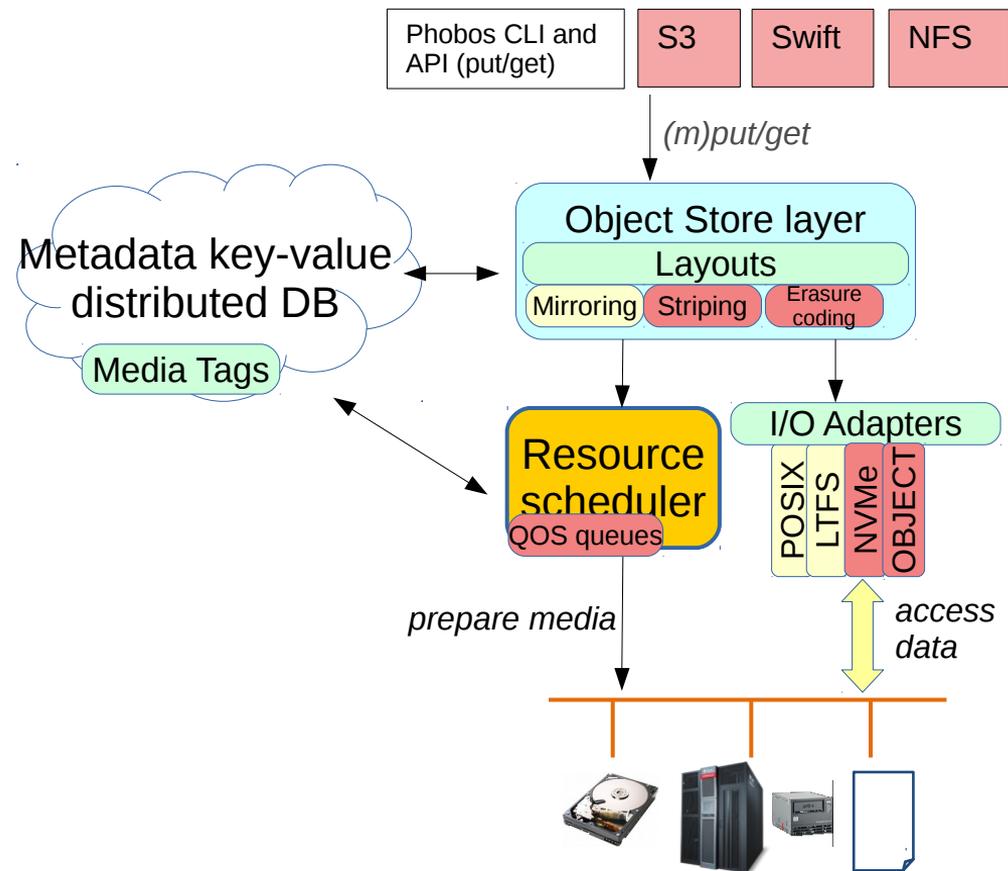


- IO adapters: multiple storage technologies (Posix, LTFS)
- Layout plugins: performance and fault-tolerance
 - simple
 - Mirroring
- Tags: storage partitioning
- Resource scheduling: optimizes tape fill rate, minimizes tapes mounts

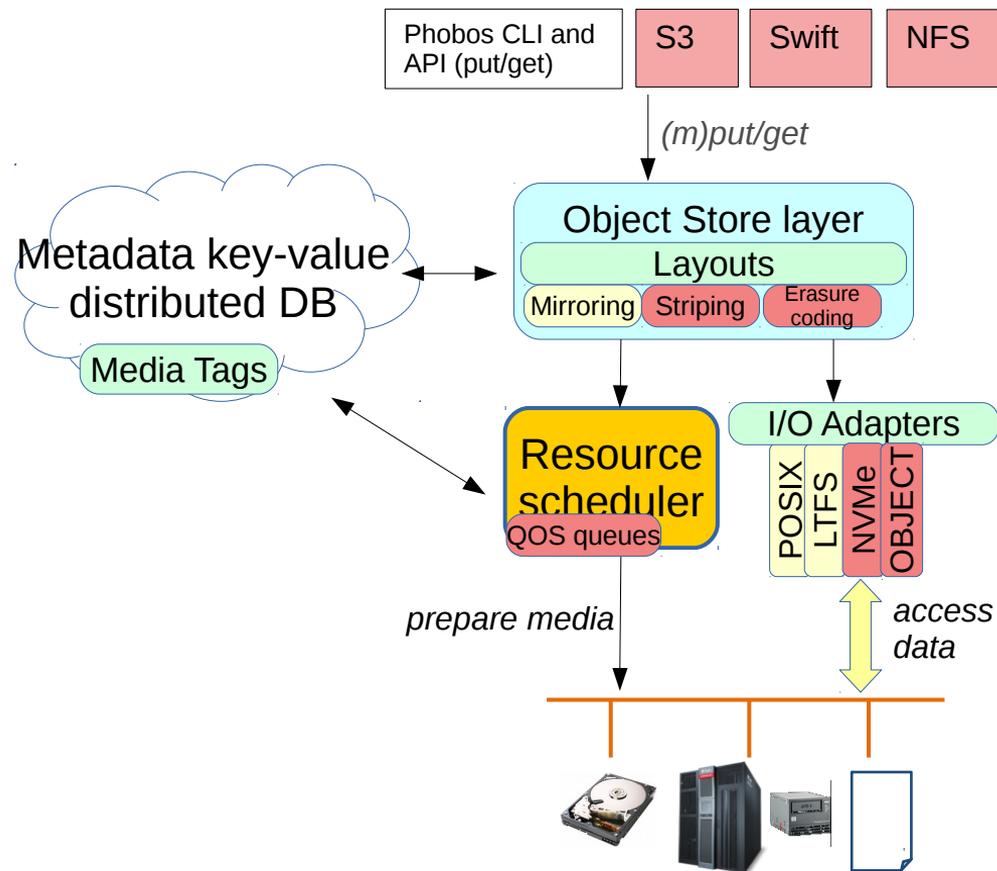


- IO adapters: multiple storage technologies (Posix, LTFS)
- Layout plugins: performance and fault-tolerance
 - simple
 - Mirroring
- Metadata key-value distributed DB
 - Media Tags
- Tags: storage partitioning
- Resource scheduling: optimizes tape fill rate, minimizes tapes mounts
- Key-value metadata schema:
 - Distributed NoSQL Database
 - Saved within objects on media (recovery, tape import)

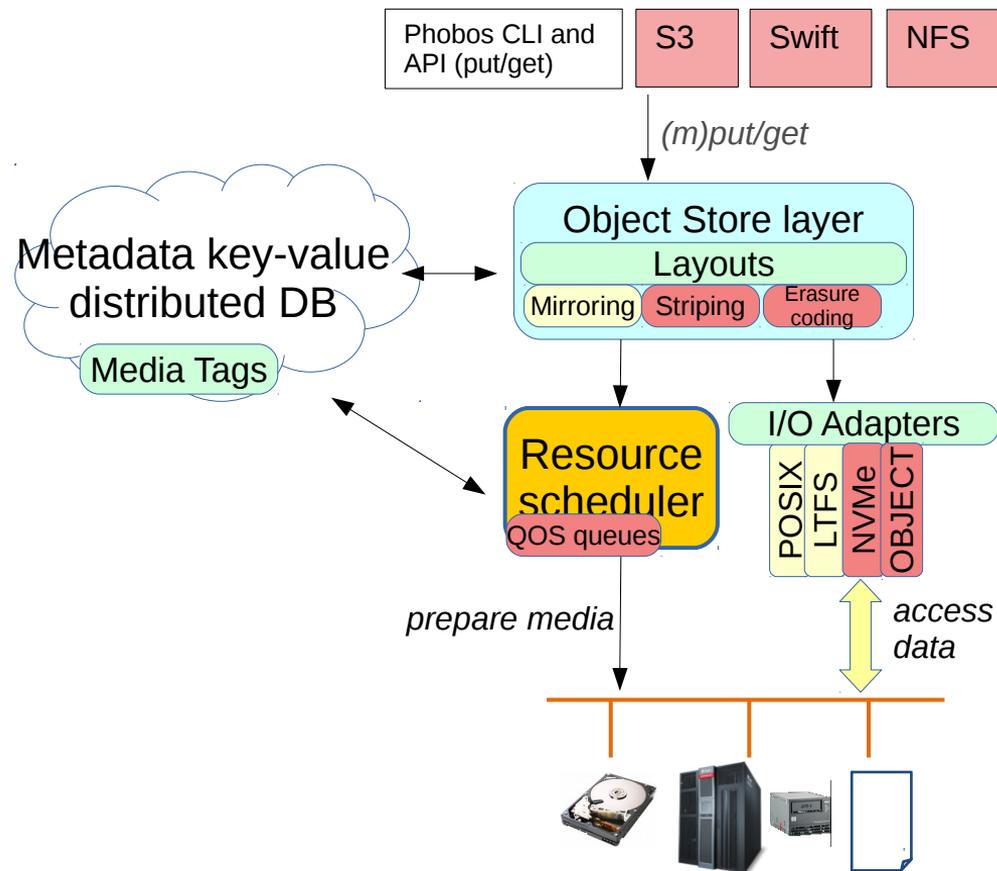




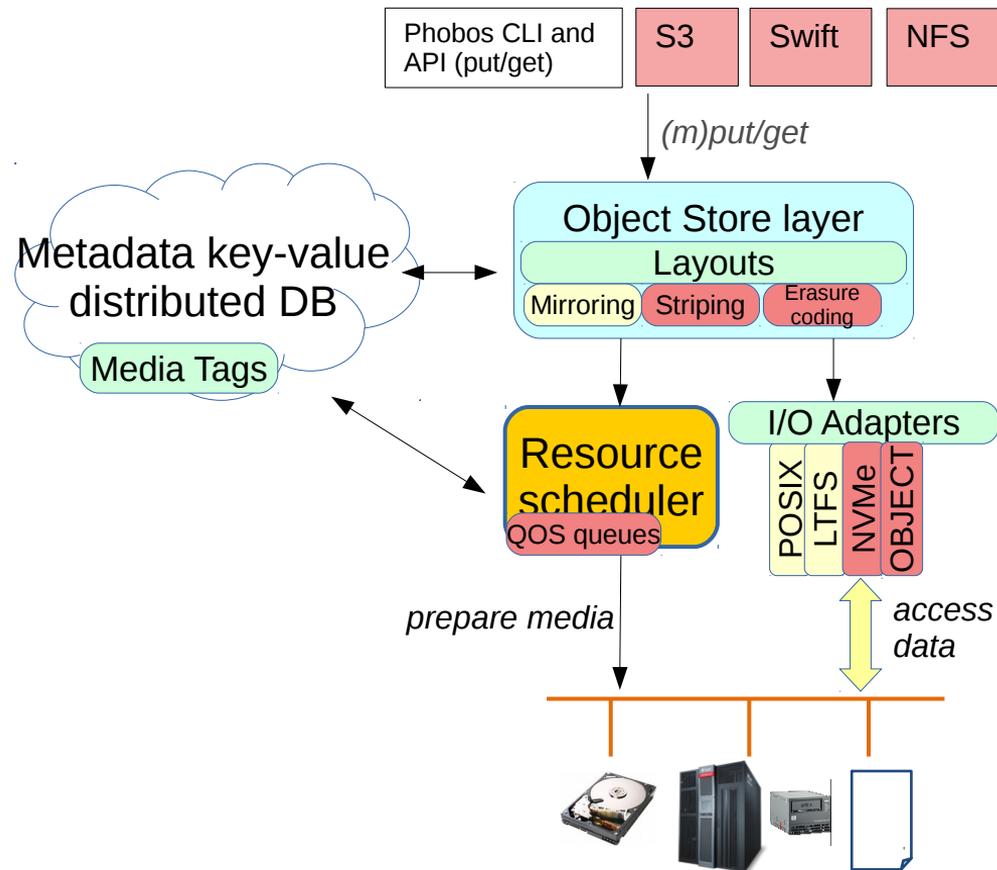
- Providing S3, Swift and NFS connectors



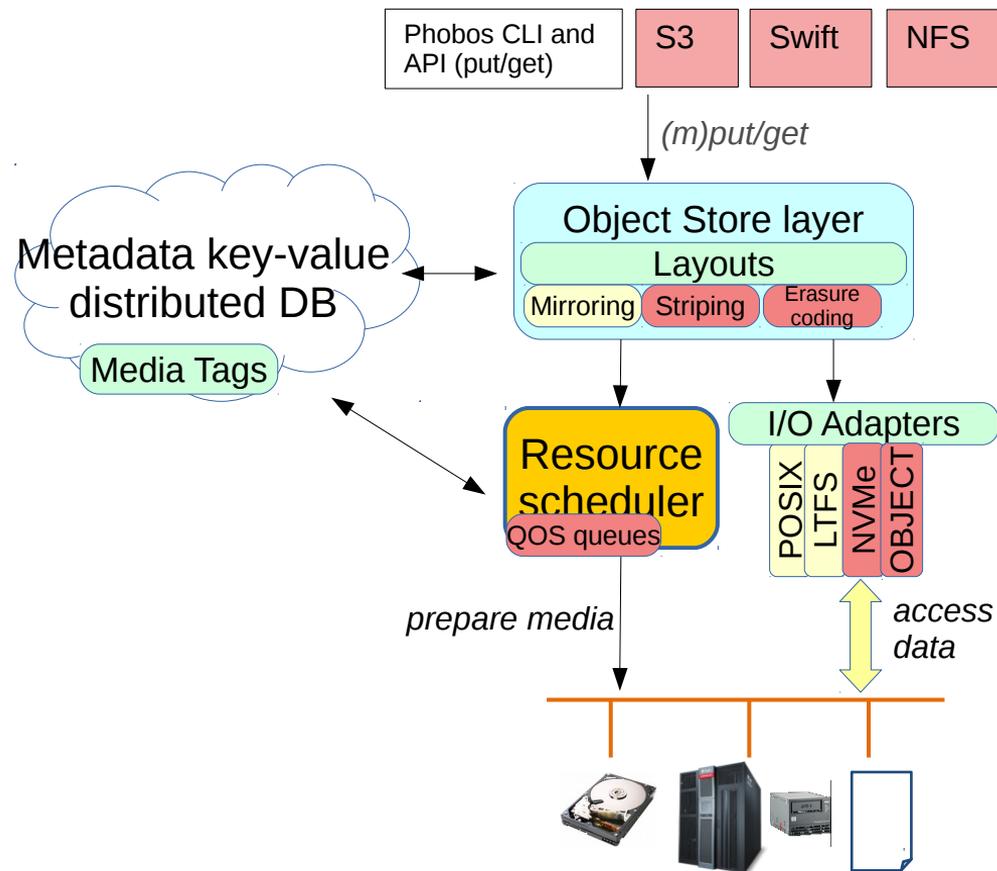
- Providing S3, Swift and NFS connectors
- Adding new layouts: striping, erasure-coding



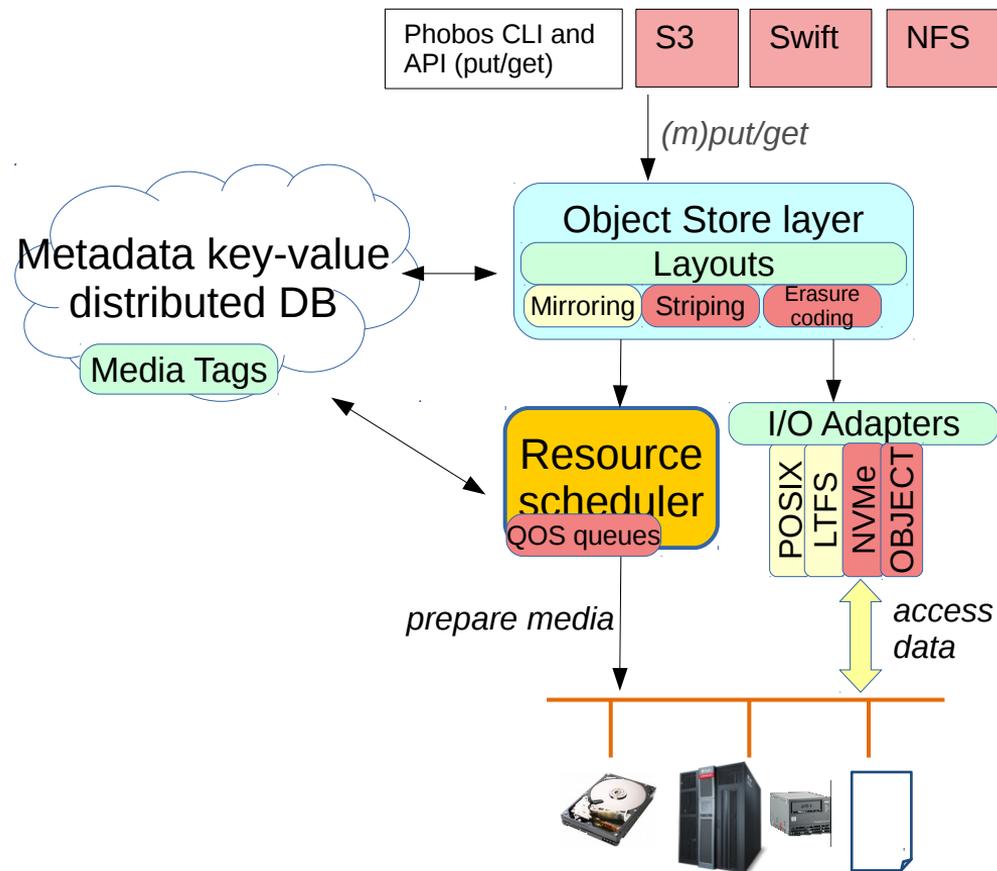
- Providing S3, Swift and NFS connectors
- Adding new layouts: striping, erasure-coding
- New IO adapters: NVMe, Object

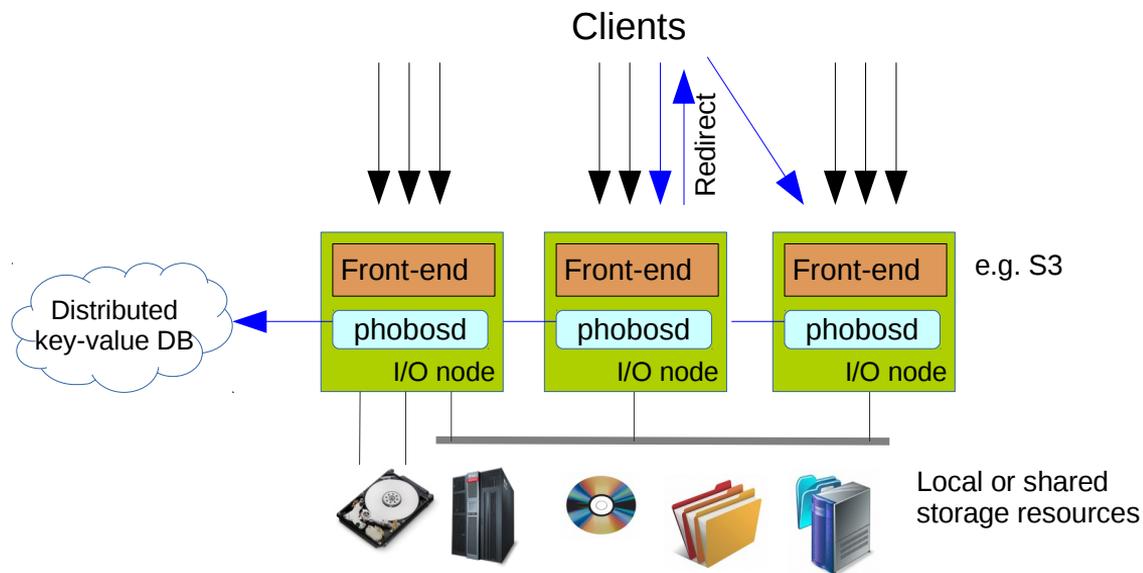


- Providing S3, Swift and NFS connectors
- Adding new layouts: striping, erasure-coding
- New IO adapters: NVMe, Object
- Media life cycle management: automatic migrations between storage technologies



- Providing S3, Swift and NFS connectors
- Adding new layouts: striping, erasure-coding
- New IO adapters: NVMe, Object
- Media life cycle management: automatic migrations between storage technologies
- Optimizing resource scheduler policies: prioritizing and grouping I/O





- Synchronization in the distributed mode:
 - Through the distributed key-value DB (object location, resource reservation...)
 - Redirection of client requests to the preferred I/O node (max 1 hop)

Setting up a tape storage in a couple of commands

```
phobos drive add --unlock /dev/st1
```

```
phobos tape add -t lto6 [073200-073222]L6
```

```
phobos tape format --unlock [073200-073222]L6
```

That's done! Your system is ready for I/Os.

Example of use-cases



- Multi-petabyte genomics datasets
- In production since 2016

DNA sequencers



Phobos

- IBM TS3500 tape library (SCSI)
- LTO6 and LTO8 drives

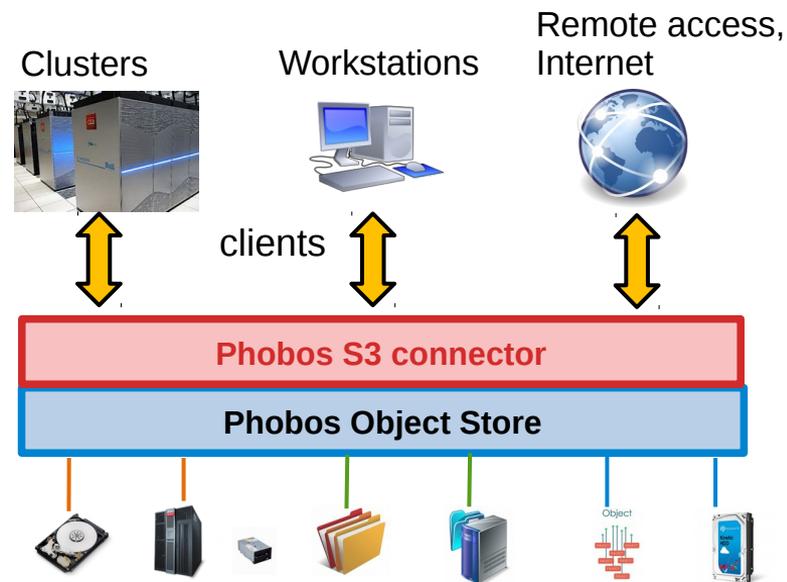


HPC data clusters



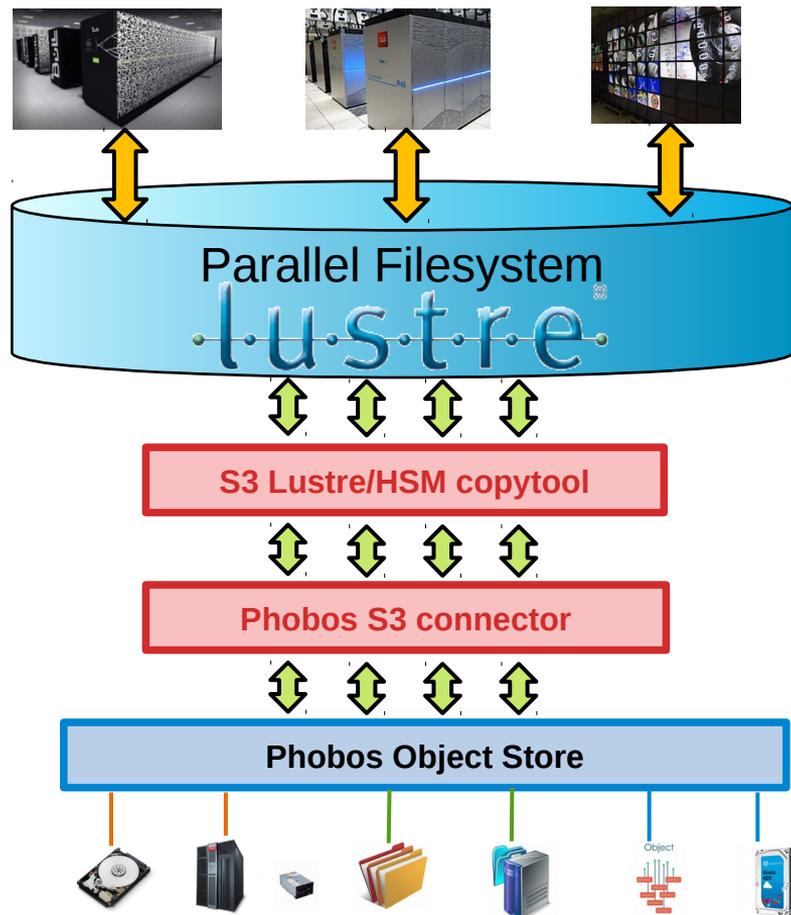
Object store with an S3 interface

- S3 interface exposed to end-users
- Phobos: high-performance, scalable storage
 - Can manage a wide variety of capacitive storage, including tape libraries
 - Provides an easy/uniform management of these storages



Phobos as a Lustre/HSM backend

- Lustre: filesystem user front-end
- Phobos as capacitive and scalable backend (hierarchical storage)



Summary

- Tape object storage at scale (and more)
- Phobos is open-source, available on github:
<https://github.com/cea-hpc/phobos>
- Contributions are welcome, as well as testers!



The CEA logo consists of the lowercase letters 'cea' in a white, rounded, sans-serif font. A thin green horizontal line is positioned directly beneath the letters. The logo is set against a dark red background that features a faint, repeating pattern of small white circles.

DE LA RECHERCHE À L'INDUSTRIE

The Phobos logo is contained within a semi-transparent light red square. It features a stylized orbital path represented by a curved line that transitions from red at the bottom to yellow at the top. A small cyan circle with a white dot in the center is positioned at the top of the curve, representing the planet Phobos. To the right of this graphic, the word 'Phobos' is written in a large, black, sans-serif font.

<https://github.com/cea-hpc/phobos>

Thank you for your attention!

Patrice LUCAS, patrice.lucas@cea.fr