

Distributed File Systems with



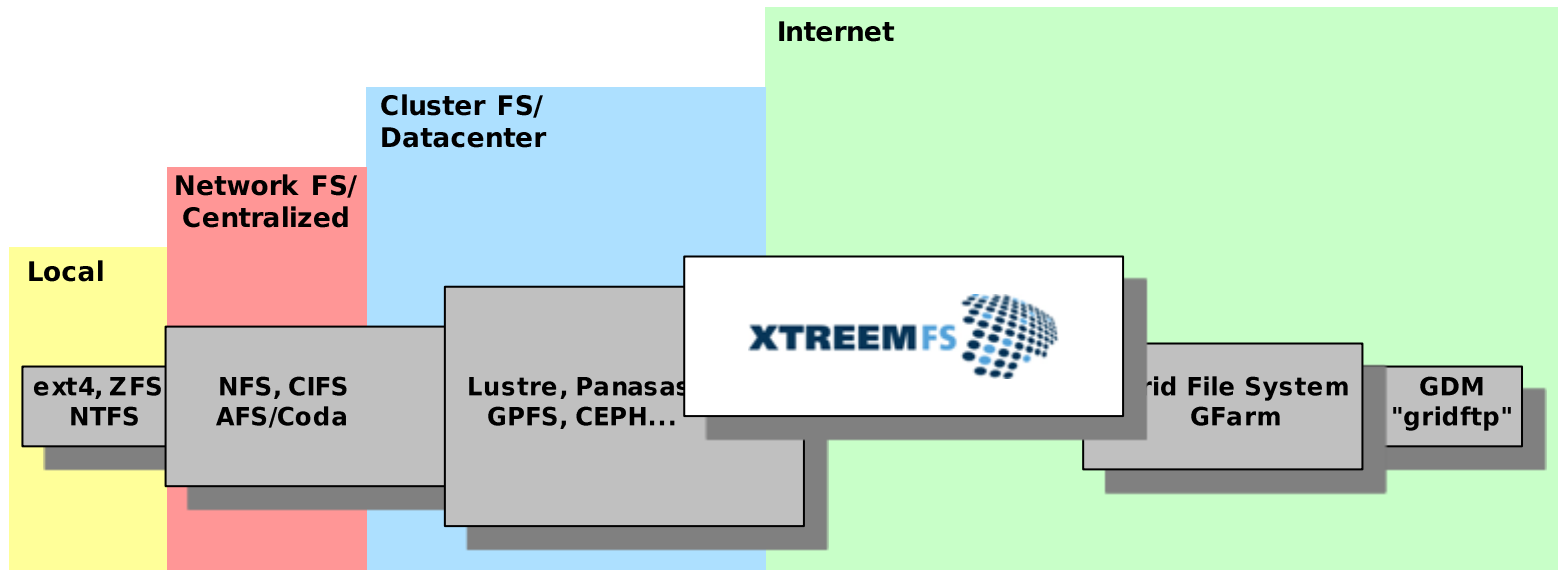
Christoph Kleineweber
Zuse Institute Berlin
kleineweber@zib.de

HARNESS Software Carpentry Workshop
July 17th, 2015

Outline

- **(Distributed) file systems**
- Introduction to XtremFS
- Use-cases for XtremFS
- Setting up your own XtremFS cluster

File System Landscape



How to Use a File System in Python?

Writing to a file:

```
f = open("myfile.txt", "w")  
f.write("Hello World\n")  
f.close()
```

Reading directory entries:

```
import os  
for f in os.listdir("/"):   
    print(f)
```

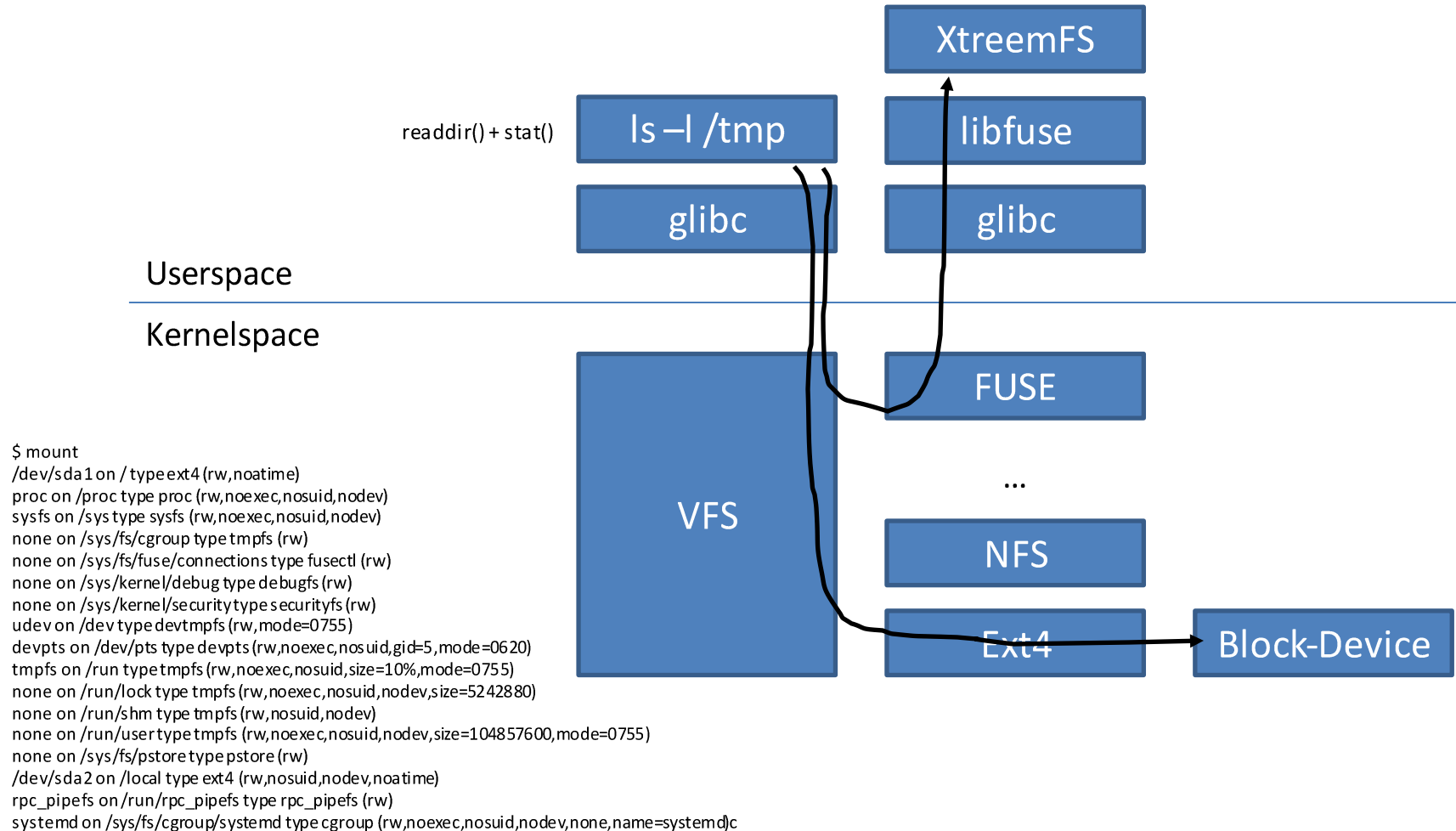
Reading a file:

```
f = open("myfile.txt", "r")  
for line in f:  
    print(line)  
f.close()
```

Reading file attributes:

```
import os  
file_info = os.stat("myfile.txt")  
print(file_info)
```

File Systems on Linux with FUSE



Distributed File Systems

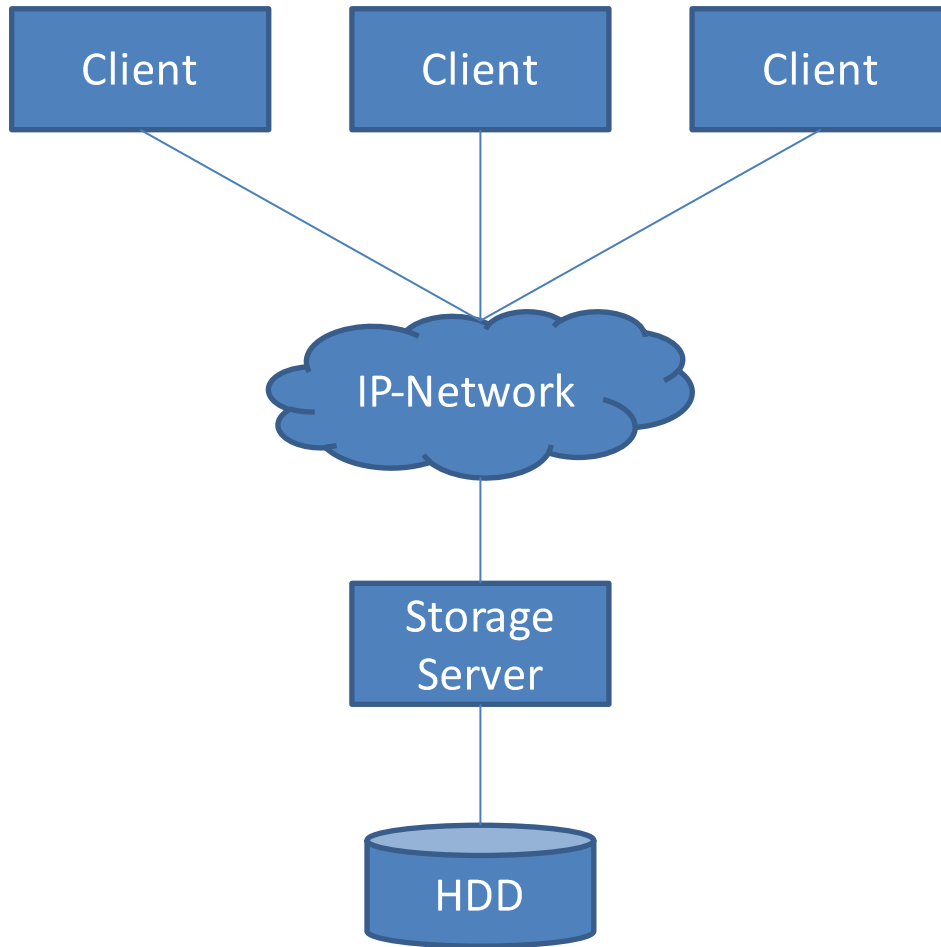
Goals:

- Transparency
 - Clients are unaware of file distribution and locality
 - Same behavior as a local file system
- Shared file access
 - Multiple users and clients
- Scalability, fault tolerance, consistency, performance, ...

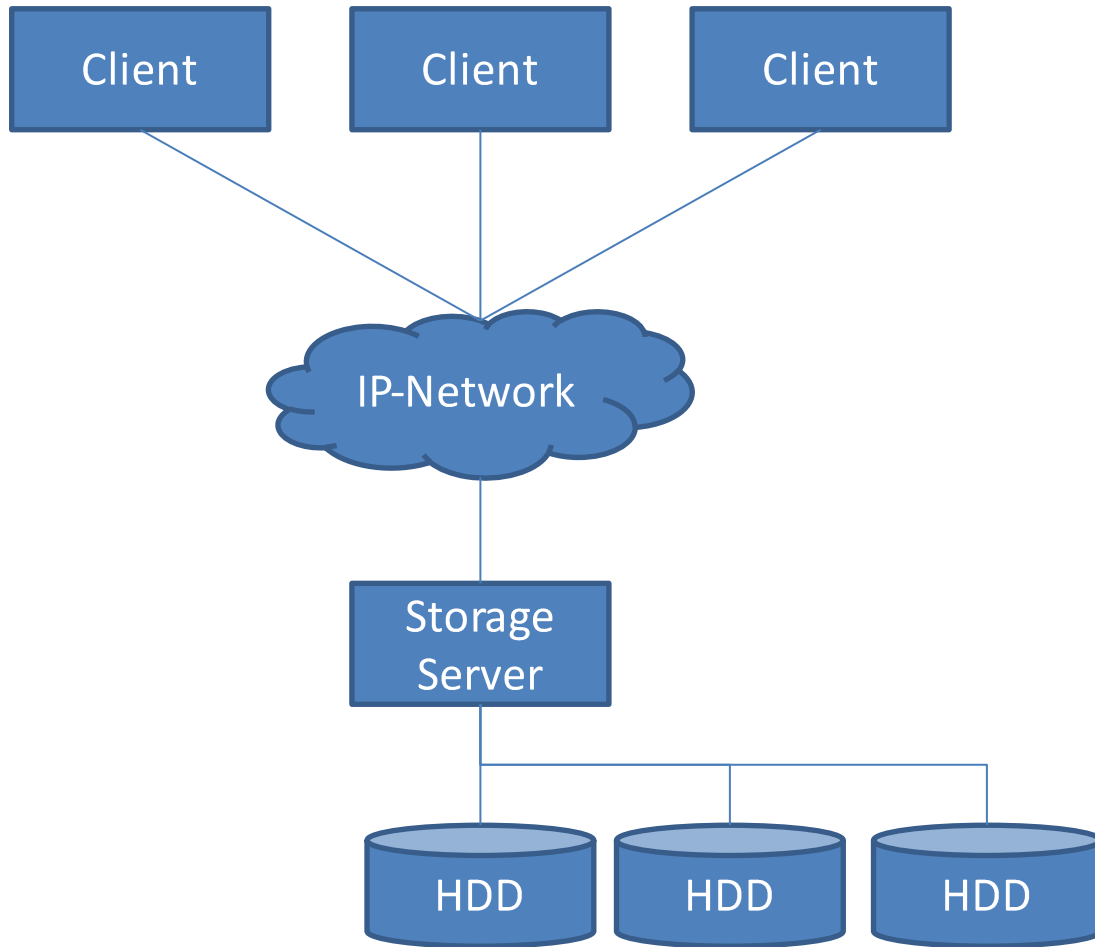
Use-cases:

- High performance computing
 - Multiple computers write checkpoints to a single file
 - Data analytics (e.g. MapReduce)
- Storing home directories
 - User wants the same environment on any PC
- Virtualization
 - Storing VM images

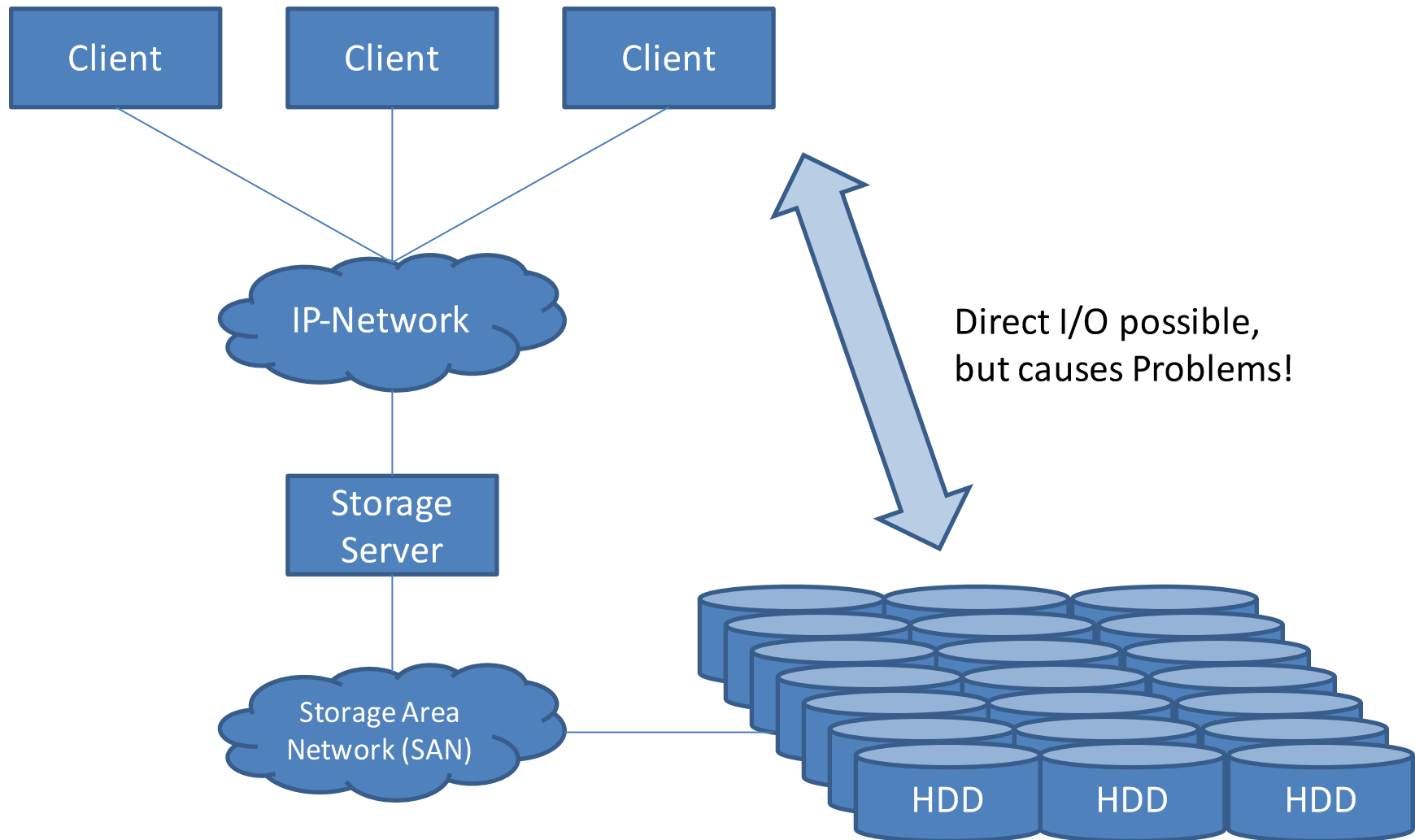
Scalable Storage Systems



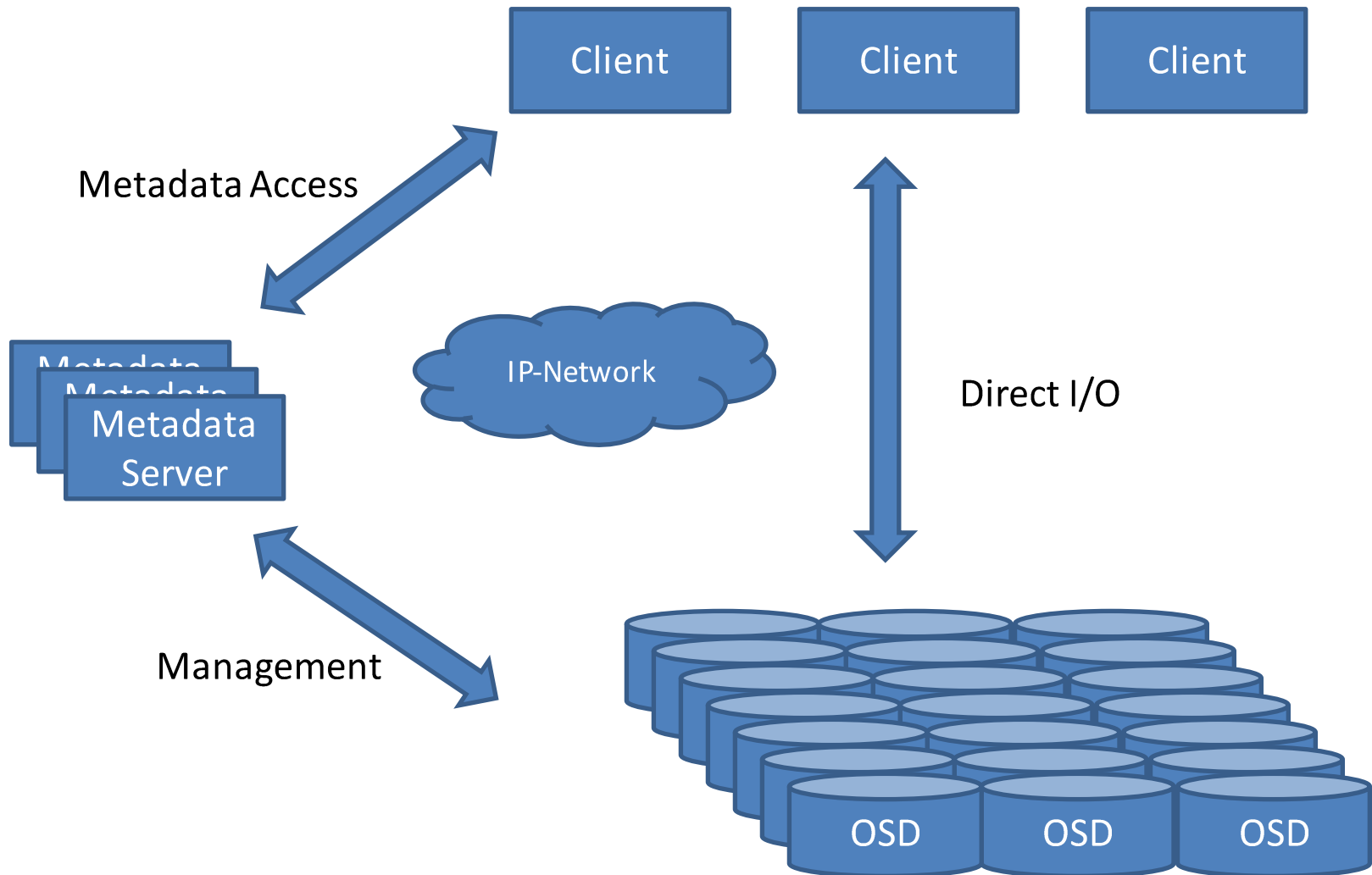
Scalable Storage Systems



Scalable Storage Systems



Object-based Storage



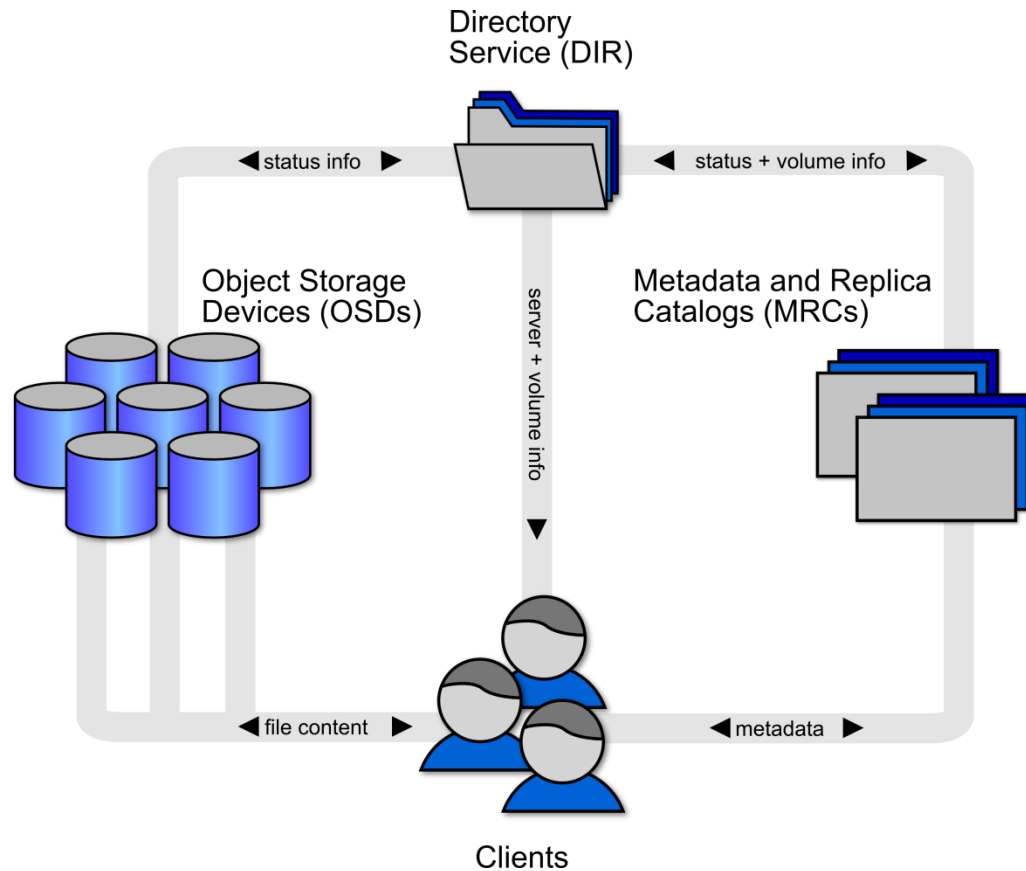
Outline

- (Distributed) file systems
- **Introduction to XtreemFS**
- Use-cases for XtreemFS
- Setting up your own XtreemFS cluster

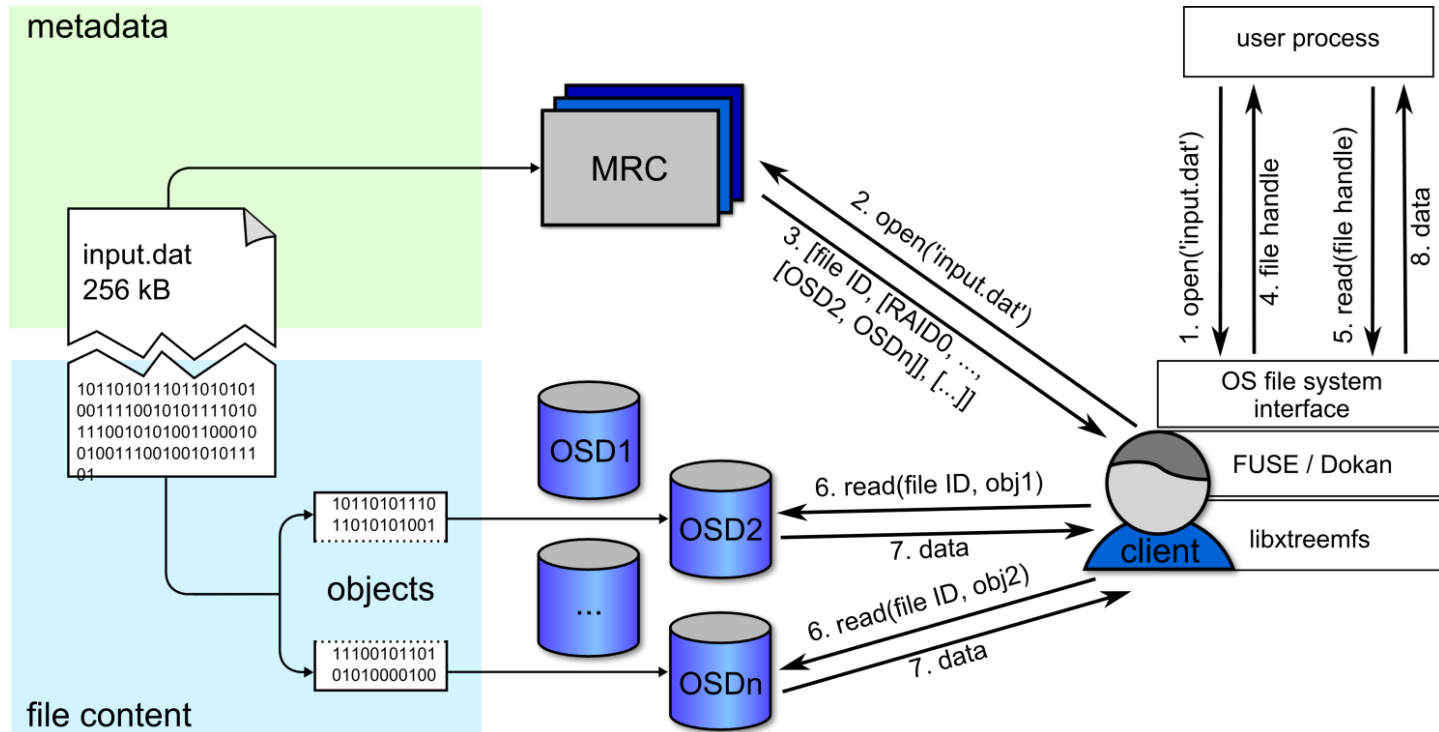
XtreemFS

- History
 - 2006 initial development in XtreamOS project
 - Since 2010 further development in different research projects
 - 2015 latest release 1.5.1 in March
- Project
 - Open source <http://www.xtreemfs.org> (BSD License)
 - POSIX compatible distributed file system
- Software
 - Client software (C++) runs on Linux, Mac OS X and Windows
 - Server software platform independent (Java)
 - Client library libxtreemfs for C++ and Java
 - Hadoop driver based on Java libxtreemfs

XtreemFS Architecture

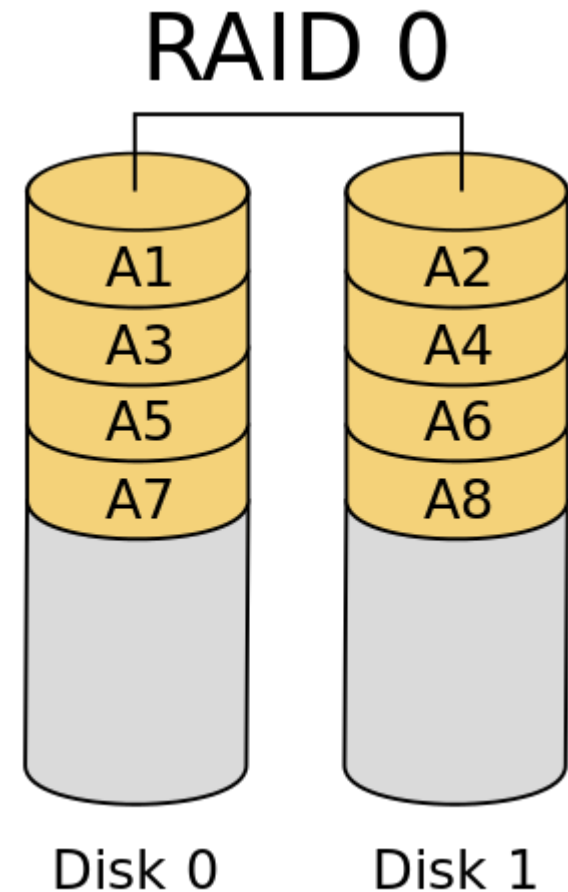


File Access



Striping Files over OSDs

- Store file segments on different OSDs
 - Similar to RAID0
- File size can exceed single server capacity
- Possible performance gain by parallel access
- Caution: Device failure corrupts the whole data!



https://en.wikipedia.org/wiki/Standard_RAID_levels

Replication

Why?

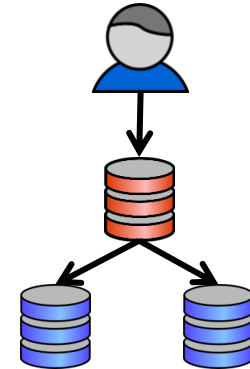
- Reliability
 - Access your data even in the case of failing servers
- Performance
 - Spread the load to multiple servers

Problems:

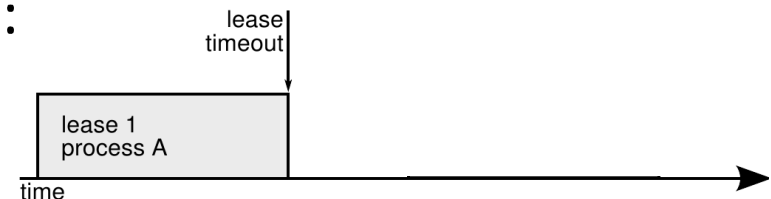
- Which replica holds the latest data?
 - Strong consistency desired!
- How to avoid “split-brain” situations?

Replication with Strong Consistency

- Primary/backup scheme
 - POSIX requires total order of update operations
 - Primary/backup
 - Primary fail-over?

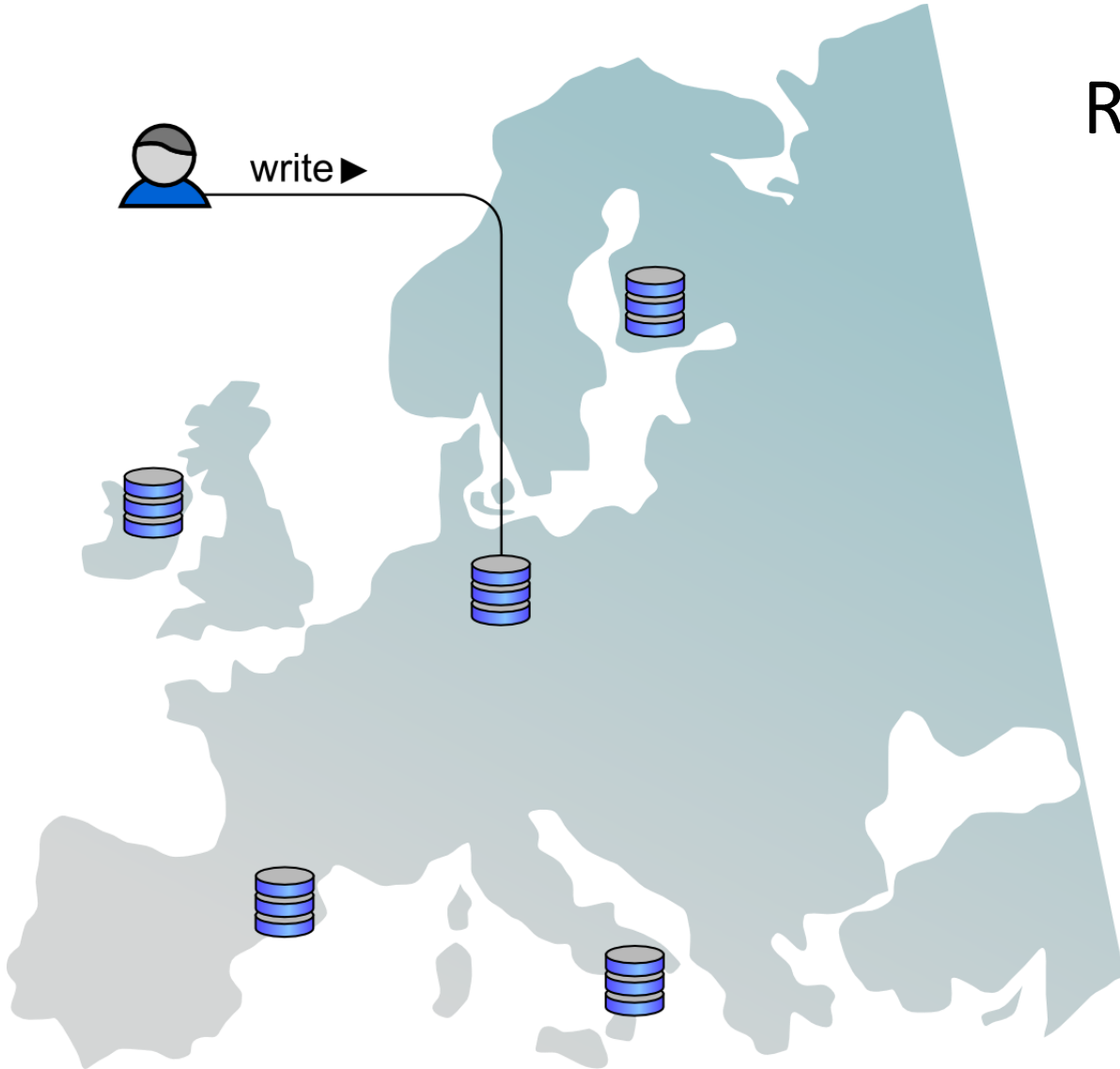


- Leases
 - Grants access to a resource (here: primary role) for a predefined period of time
 - Failover after timeout possible
 - Assumption: loosely synchronized clocks
 - Max drift ϵ

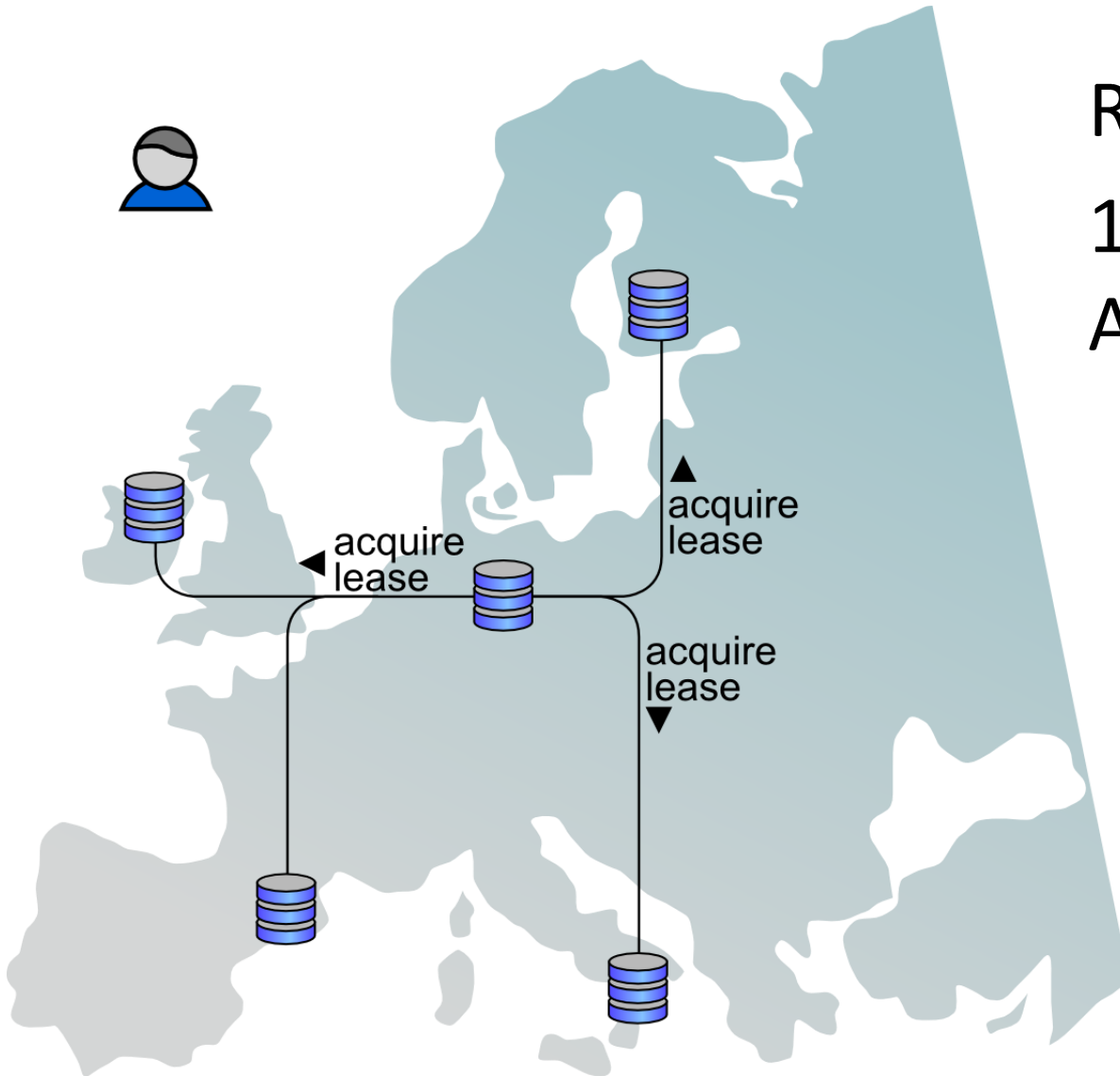


Example: Replicated Write

Replicated write():



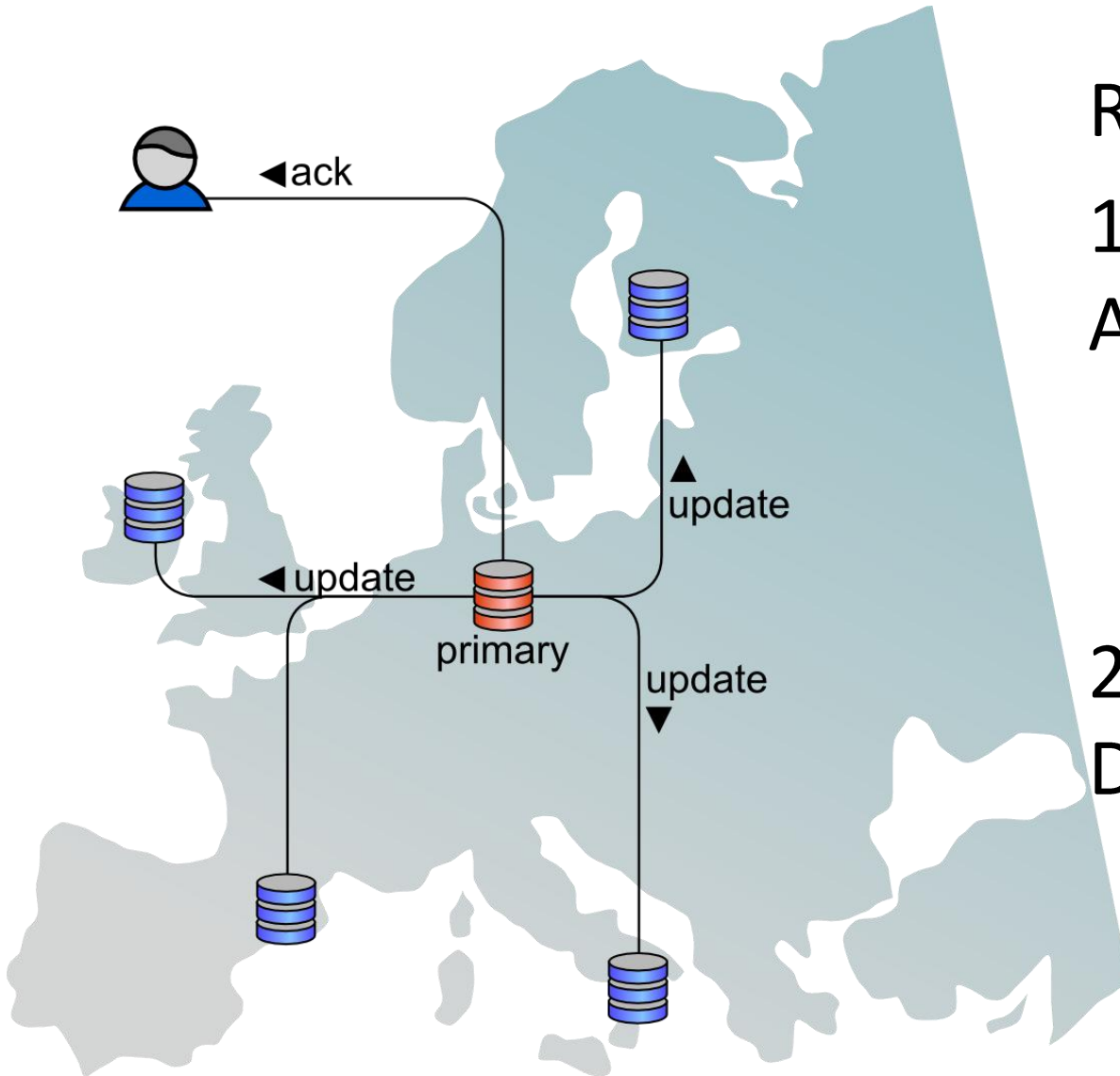
Example: Replicated Write



Replicated write():

1. Lease Acquisition

Example: Replicated Write

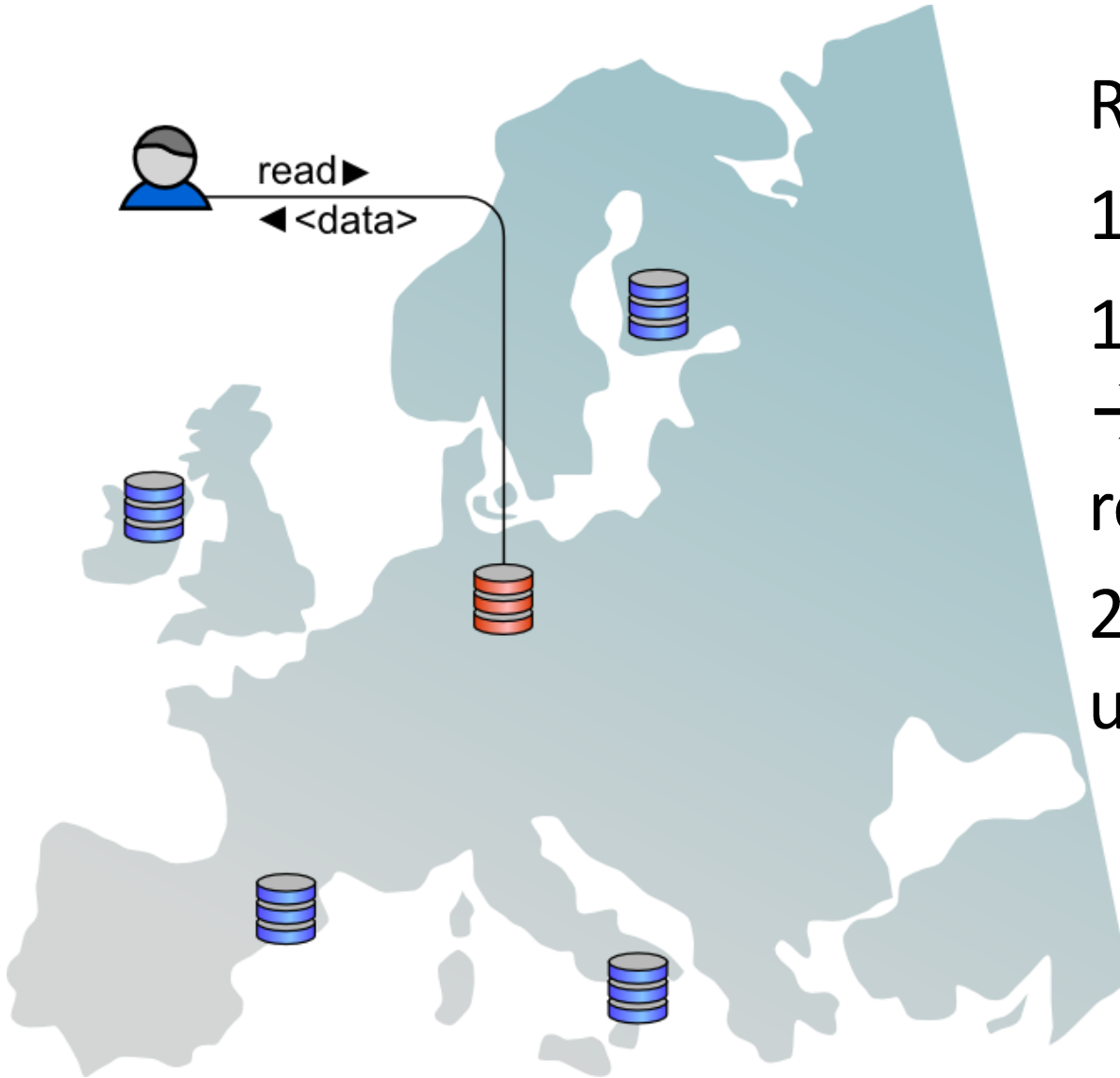


Replicated write():

1. Lease
Acquisition

2. Data
Dissemination

Example: Replicated Read



Replicated read():

1. Lease Acquisition

1b. "Replica Reset"
→ update primary's
replica

2. Respond to read()
using local replica

Replication of Immutable Files

- Only for “write-once” files
 - File must be marked as “read-only”
 - done automatically after close()
 - Use Case: CDN
 - client can read from and replica
- Replica Types:
 1. Full replicas
 - complete copy, fills itself as fast as possible
 2. Partial replicas
 - Initially empty
 - on-demand fetching of missing objects
 - P2P-like efficient transfer between all replicas

File System Snapshots

- Persistent file system content
- Use-cases:
 - File versioning
 - Get consistent state for backups
- Challenge:
 - Files are spread across multiple servers

Snapshots in XtremFS

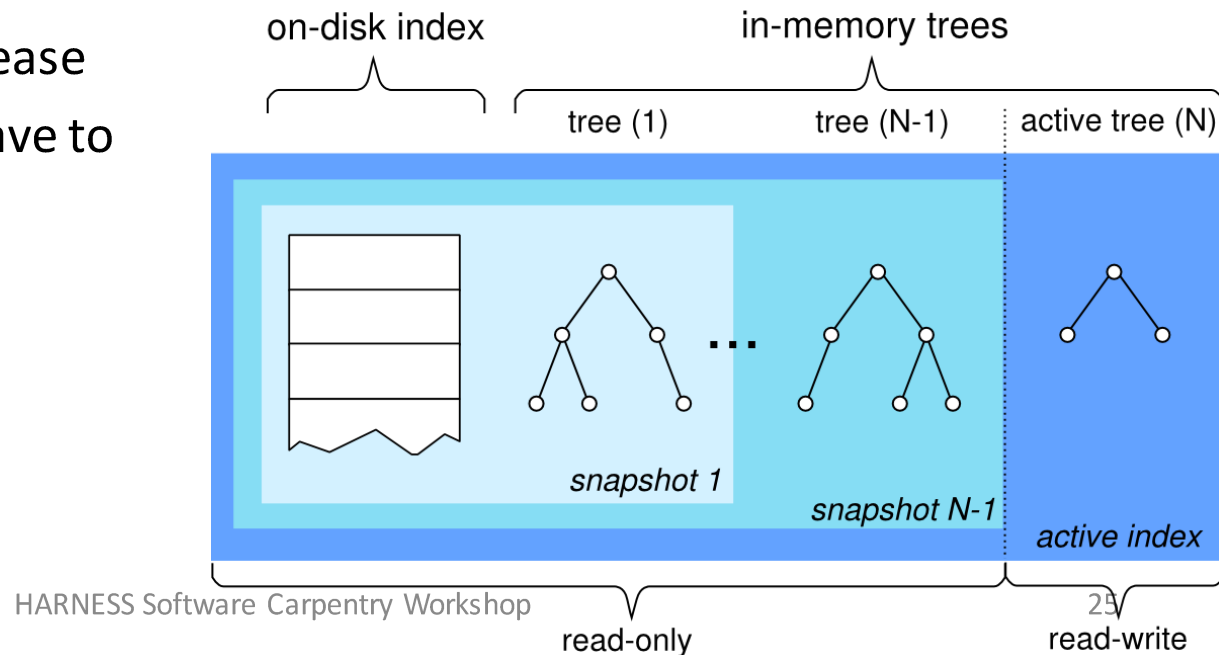
- Use copy-on-write on OSDs
 - Objects are not overwritten
 - New version for each write
- Take metadata Snapshot on MRC
- Periodic cleanup process deletes object versions that are not assigned to snapshots

Metadata Storage

- DIR and MRC use key-value store BabuDB for persistent storage
- Uses log-structured merge-tree internally
 - Known from Google's BigTable
 - Well suited for snapshots
- DB replication for fault tolerance
 - Master replica holds lease
 - Majority of replicas have to be updated



July 17th, 2015



HARNESS Software Carpentry Workshop

Outline

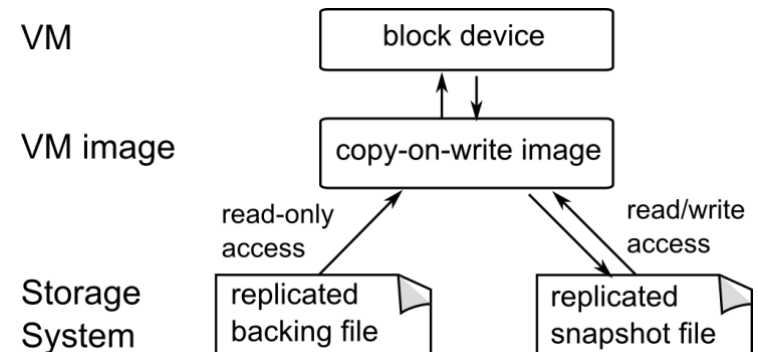
- (Distributed) file systems
- Introduction to XtreemFS
- **Use-cases for XtreemFS**
- Setting up your own XtreemFS cluster

Use-Case: XtreamFS for VM Images

- VM deployment
 - Create copy (clone) of original VM image
 - Run cloned VM image on hypervisor
 - (Discard cloned image after VM shutdown)
- Problems
 - Cloning is time-consuming
 - Waste of space (redundancy when running multiple instances of same image)

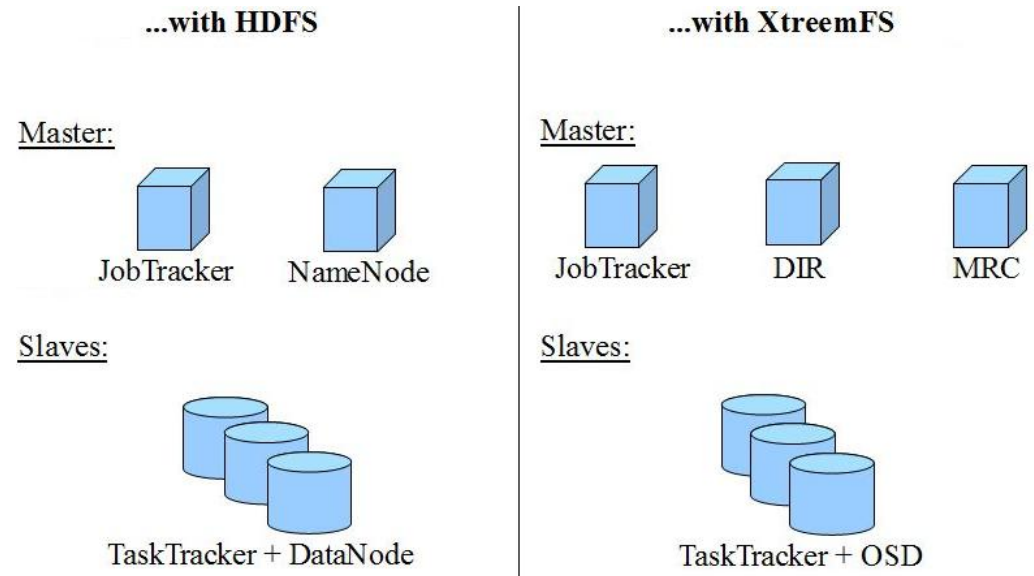
XtreemFS Replication for VM Storage

- qcow2 VM image format
 - Allows snapshots
 1. Immutable backing file
 2. Mutable, initially empty snapshot file
 - Instead of cloning, snapshot original VM image (< 1 second)
 - Use Read/Write replication for snapshot file
- Problem left: run multiple VMs simultaneously
 - Snapshot file: R/W replication scales with # OSDs and # files
 - Backing file: bottle neck
 - Use *immutable file* replication



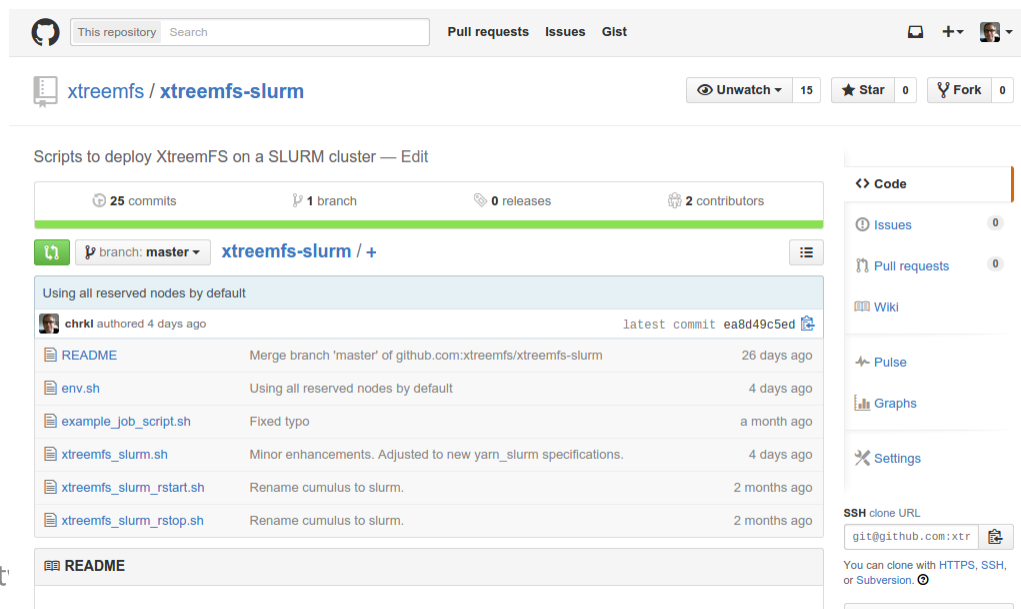
Use-Case: XtreamFS for Data Analytics

- Replace HDFS in Hadoop cluster
 - Avoid limitations of HDFS
 - Hadoop file system driver based on libxtreemfs



XtreemFS on Demand

- XtreemFS runs on commodity hardware
 - Often with large SATA disks
 - Slightly utilized disks on compute clusters
- Scripts to deploy XtreemFS on SLURM are available!



The screenshot shows the GitHub repository page for `xtreemfs / xtreemfs-slurm`. The repository has 25 commits, 1 branch, 0 releases, and 2 contributors. The main branch is `master`. The repository description is "Scripts to deploy XtreemFS on a SLURM cluster — Edit".

The commit history table is as follows:

File	Commit Message	Time Ago
README	Merge branch 'master' of github.com:xtreemfs/xtreemfs-slurm	26 days ago
env.sh	Using all reserved nodes by default	4 days ago
example_job_script.sh	Fixed typo	a month ago
xtreemfs_slurm.sh	Minor enhancements. Adjusted to new yarn_slurm specifications.	4 days ago
xtreemfs_slurm_rstart.sh	Rename cumulus to slurm.	2 months ago
xtreemfs_slurm_rstop.sh	Rename cumulus to slurm.	2 months ago

The right sidebar contains links to `Code`, `Issues`, `Pull requests`, `Wiki`, `Pulse`, `Graphs`, and `Settings`. The SSH clone URL is `git@github.com:xtreemfs/xtreemfs-slurm`.

Questions?

Website: <http://www.xtreemfs.org>
Mailinglist: xtreemfs@googlegroups.com
Repository: <https://github.com/xtreemfs/xtreemfs>
Twitter: [@xtreemfs](https://twitter.com/xtreemfs)

Outline

- (Distributed) file systems
- Introduction to XtremFS
- Use-cases for XtremFS
- **Setting up your own XtremFS cluster**

Hands-on Session

- Getting familiar with Vagrant
- Starting XtremFS services
- Creating and mounting volumes
- Configure striping
- Configure replications
- Creating snapshots