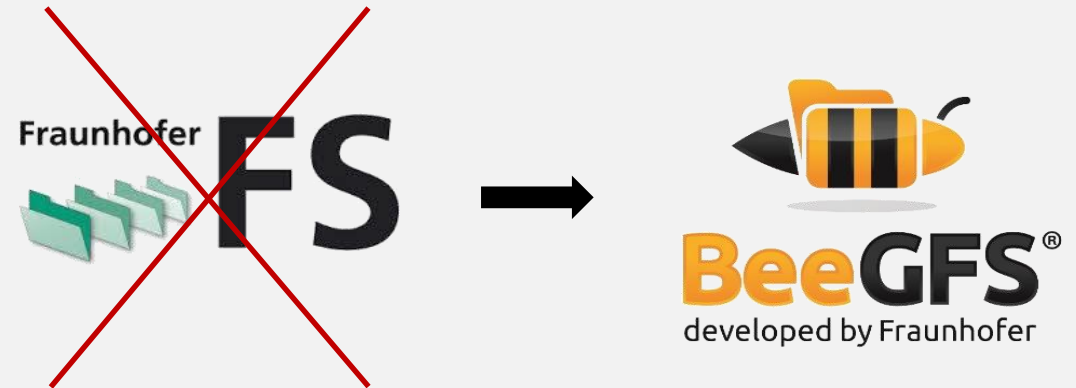# BeeGFS – not only for HPC

# What is BeeGFS?

- A hardware independent parallel filesystem

- Designed for high performance and high throughput environments

- Developed at Fraunhofer ITWM (original name: FhGFS)

- Productive installations since 2007

- First commercial installation in 2009

- Renamed to BeeGFS in 2014

- Free to use
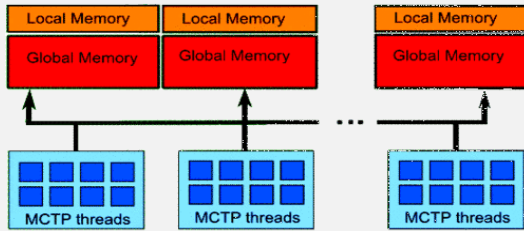
- Commercial support available

# Fraunhofer ITWM

- Institute for Industrial Mathematics

- Located in Kaiserslautern, Germany

- Staff: ~ 260 employees + ~ 60 PhD students in 8 departments
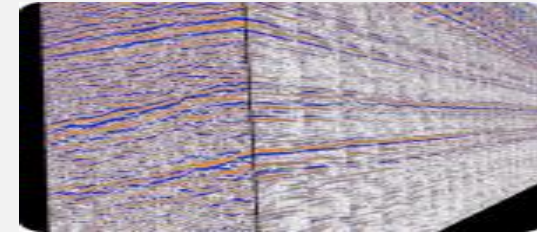


www.beegfs.com

# Fraunhofer ITWM – Competence Center HPC



Parallel Programming models & tools



Photo realistic real time ray tracing



Interactive seismic imaging



Parallel File Systems



Big Data



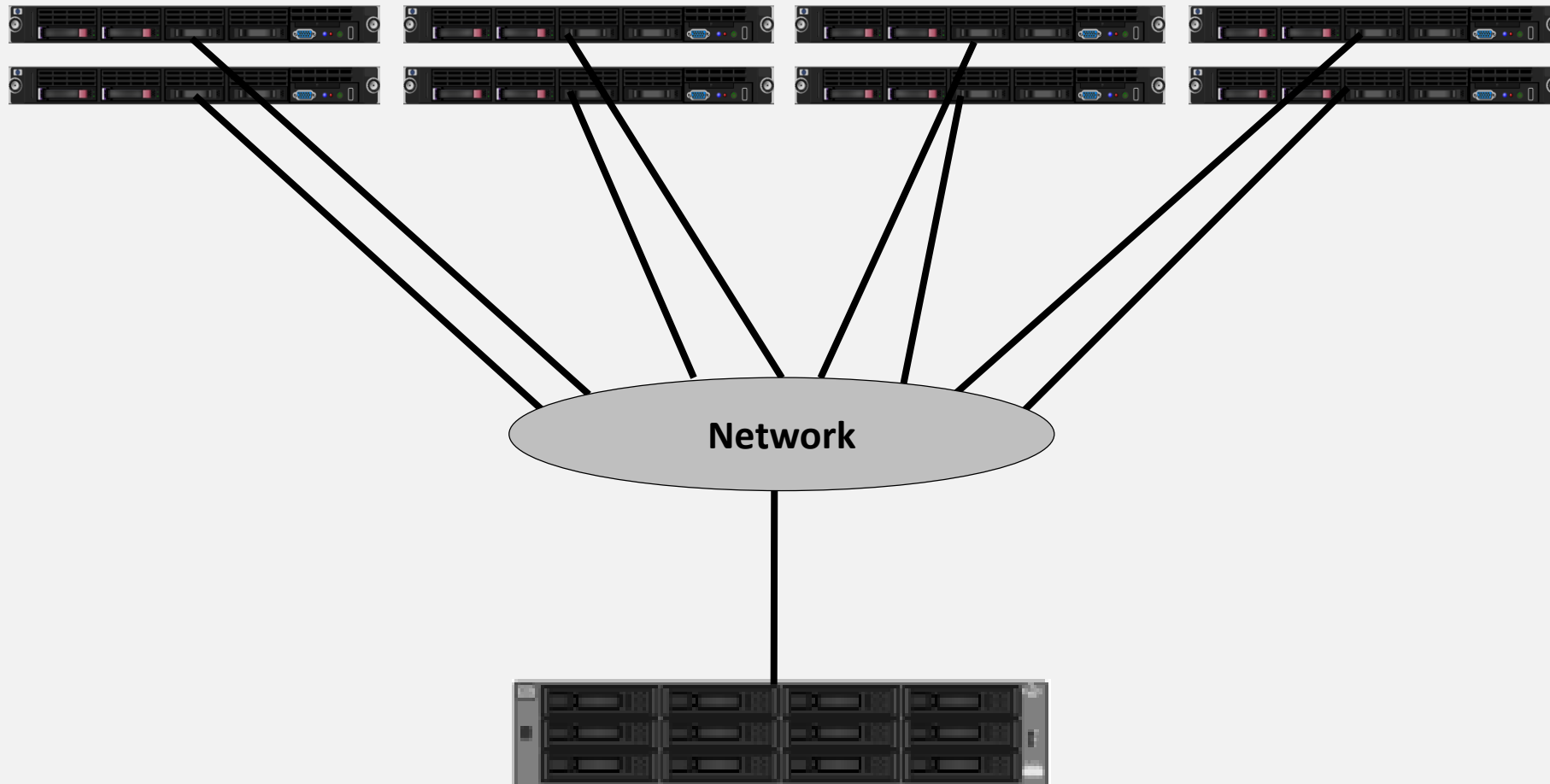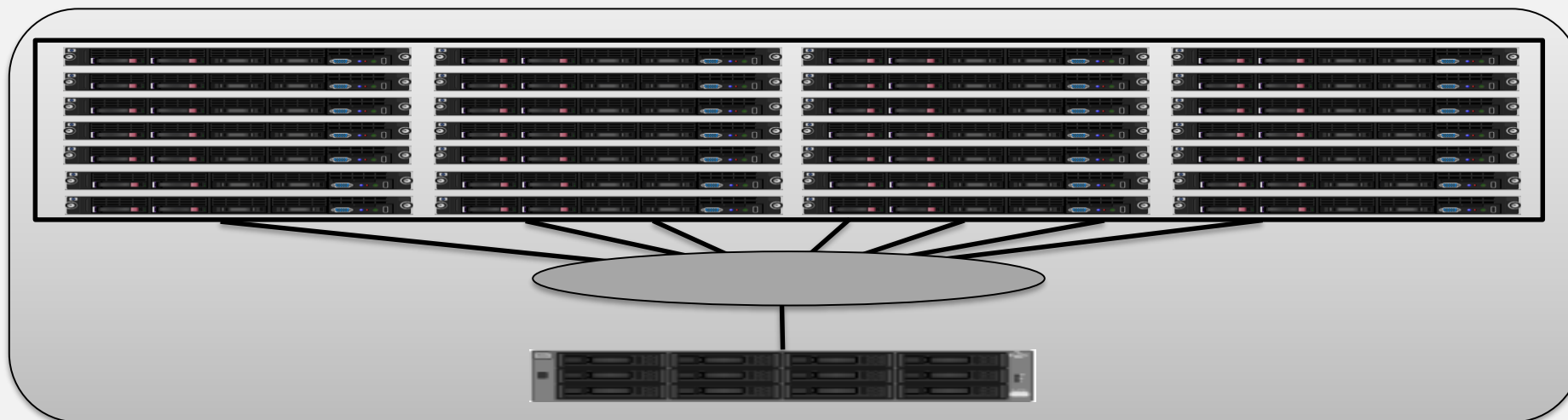Smart Energy / Green by IT

www.beegfs.com

# ThinkParQ?

- ThinkParQ...
  - ... is a spin-off from Fraunhofer ITWM to bring BeeGFS to the market
  - ... does consulting, services and support around BeeGFS
  - ... manages partner relationships
  - ... will develop own add-ons for BeeGFS

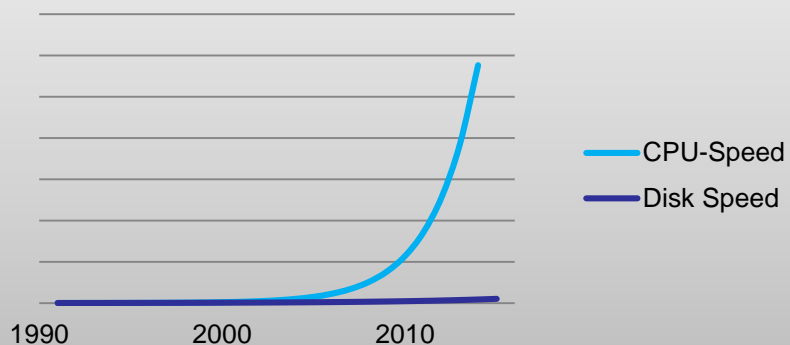# Most storage systems still look like this...

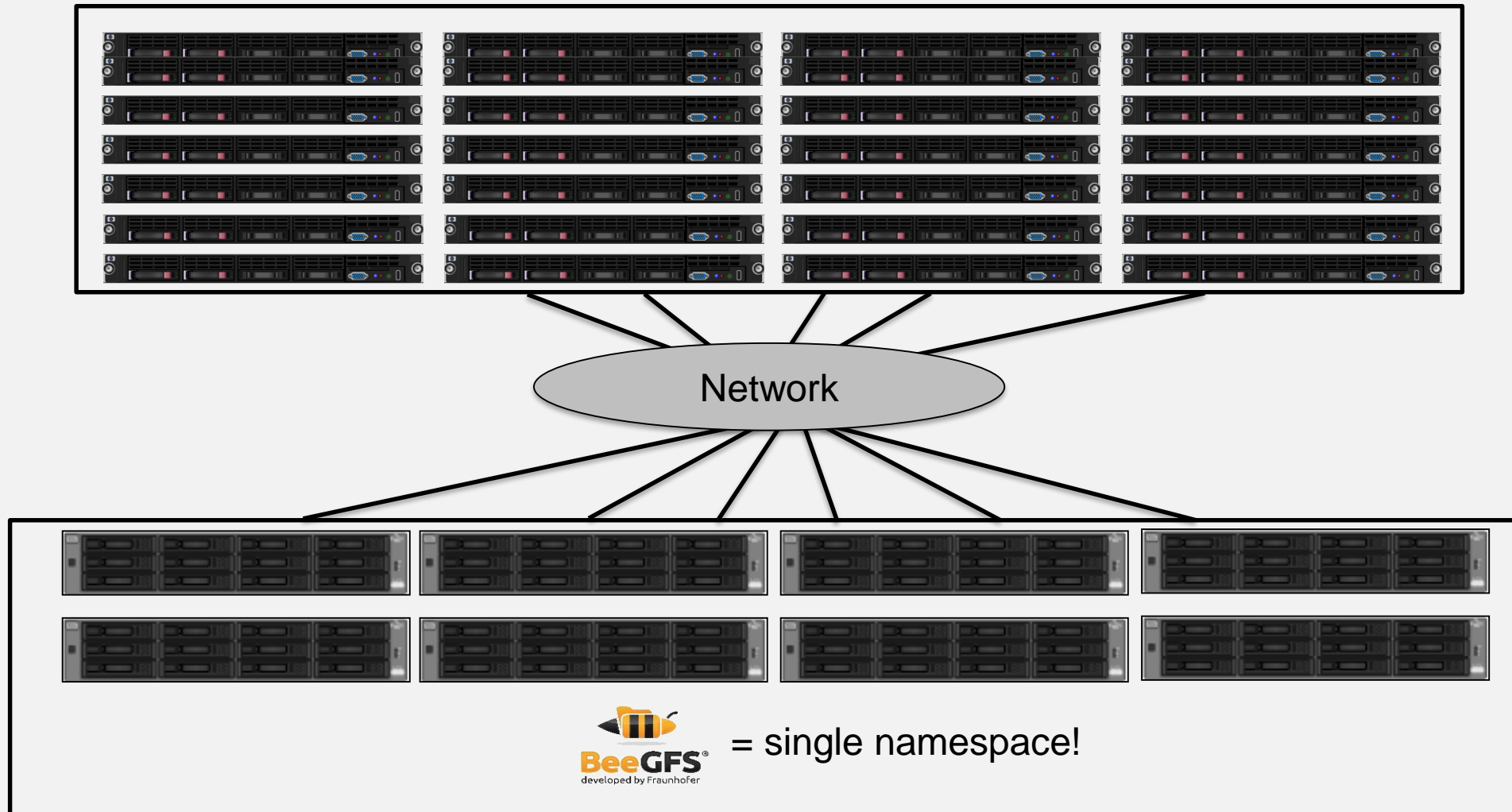

Network

www.beegfs.com

# But why bother?



### Disk vs. CPU-Speed



*„A supercomputer is a device for turning compute-bound problems into I/O-bound problems."*

*- Ken Batcher*
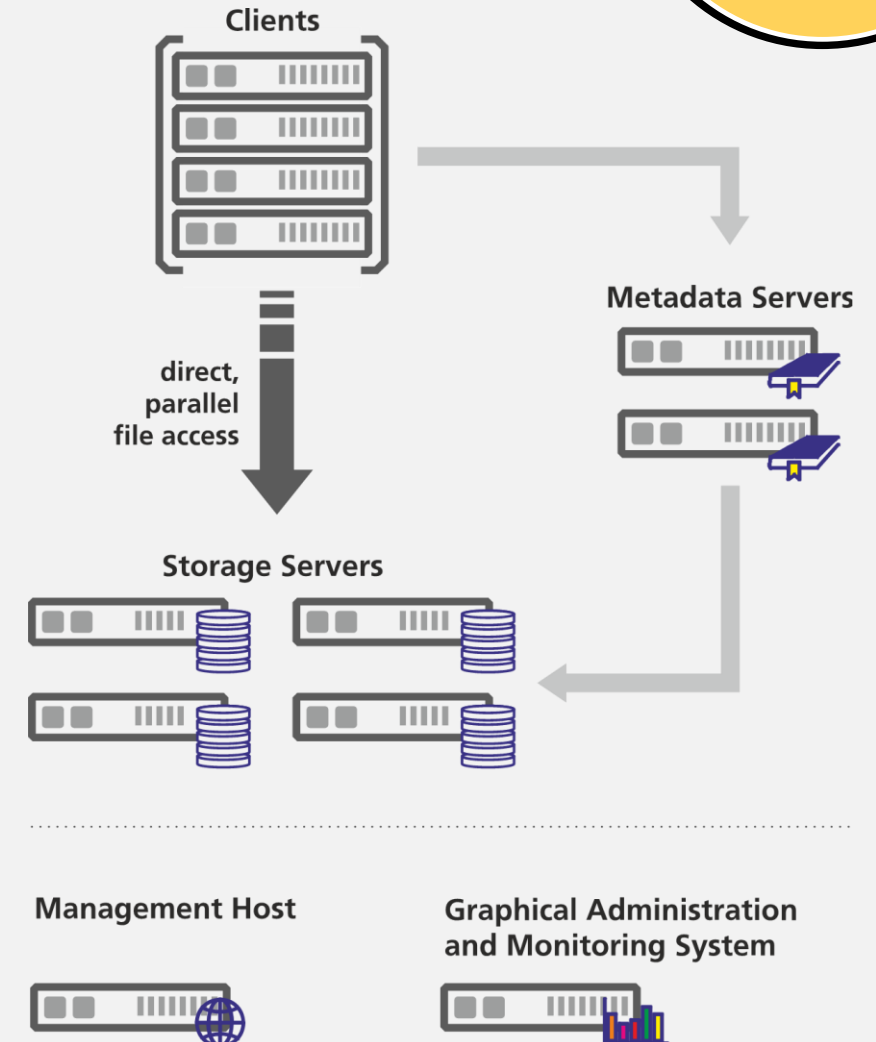
# Scale out storage



Network

BeeGFS® = single namespace!

# BeeGFS Architecture

- Management Host
  - Maintains a list of all components in the system
  - Provides all neccessary information to new components

- Storage Servers
  - Store the (distributed) file contents

- Metadata servers
  - Manage the metadata of file system entries
  - Maintain striping information for files
  - Not involved in data access

- Client
  - Native client module to mount the file system

- Graphical Administration and Monitoring System
  - GUI to perform administrative tasks and monitor system information



Clients

direct, parallel file access

Metadata Servers

Storage Servers

Management Host

Graphical Administration and Monitoring System

# Key Concepts



**Performance Scalability**
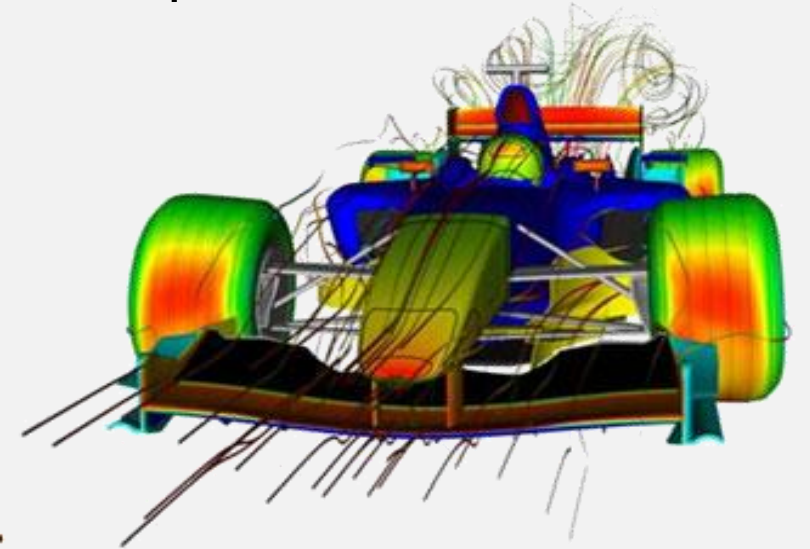
**Flexibility**

**Easy to use**

# Key Features

- Performance & scalability
  - Initially optimized for HPC
  - Completely multithreaded – lightweight design
  - Supports GigE/10GE/40GE (TCP/RoCE) and InfiniBand (TCP/RDMA)
  - Distributed file contents: aggregated throughput of multiple servers
  - Distributed metadata across multiple servers
  - Excellent single stream performance

# Key Features

- Performance & scalability
- Flexibility
  - Multiple daemons (any combination) can run on the same machine
  - Flexible striping per file/per directory
  - Add servers without downtime
  - On demand filesystem „per job" possible
  - Client runs on any kernel >2.6.18
  - Client runs on Xeon PHI
  - ARM port available
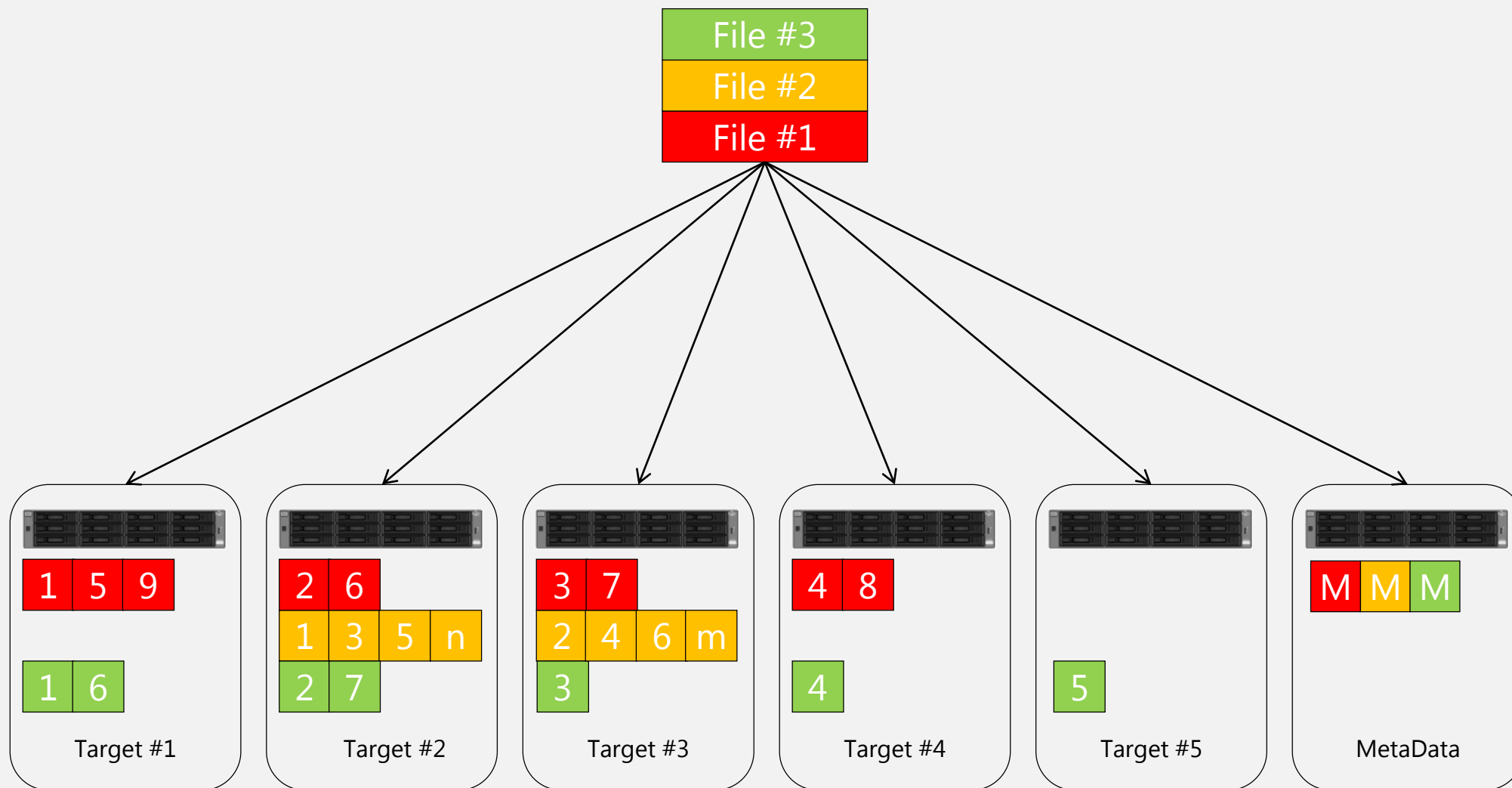  - NFS & SMB/CIFS re-export possible

www.beegfs.com

# Key Features

- Performance & scalability

- Flexibility

- Easy to use
  - Servers run in user space
  - No kernel patches
  - Servers use existing local filesystems (ext4, xfs, zfs, ...)
  - Packages for RHEL/SL/CentOS/SLES/Debian/Ubuntu
  - Hardware independent
  - Graphical monitoring tool

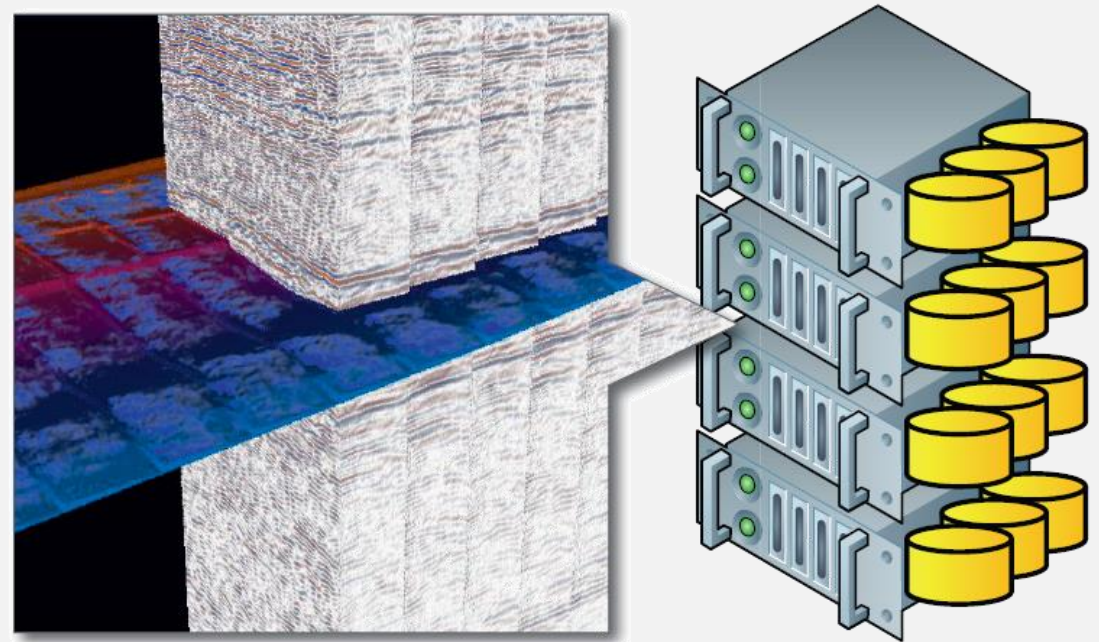# Striping

# BeeOND – BeeGFS On Demand

- Create a parallel file system instance with one command

- Use cases: cloud computing, test systems, cluster compute nodes, .....

- Can be integrated in cluster batch system (e.g. PBS)
  - Suitable for a private „per-job parallel file system"

- Used in Fraunhofer Seislab
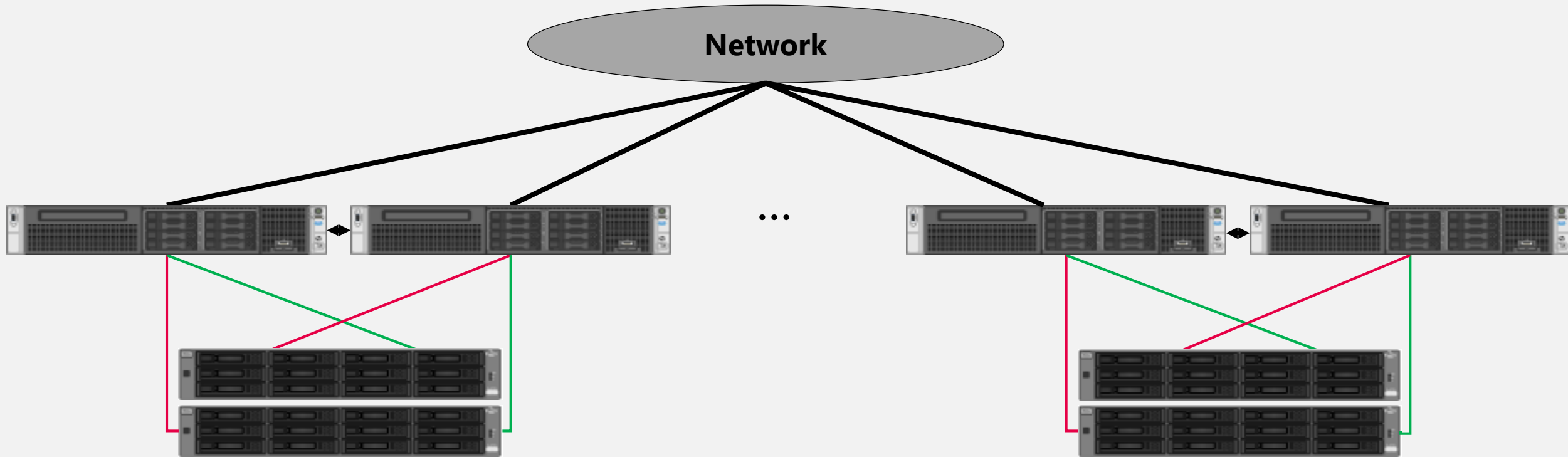  - Take load from global storage
  - Speed-up certain I/O patterns

# BeeOND – Use in Fraunhofer Seislab

- Fraunhofer Seislab
  - In-house cluster of CC-HPC at Fraunhofer ITWM
  - 92 compute nodes with 1 TB of SSDs each
  - Global BeeGFS storage on 3,5" SATA drives
- Create BeeOND on SSDs on job startup
- Stage-in input data, work on BeeOND, stage-out results

# High Availability - Shared Storage
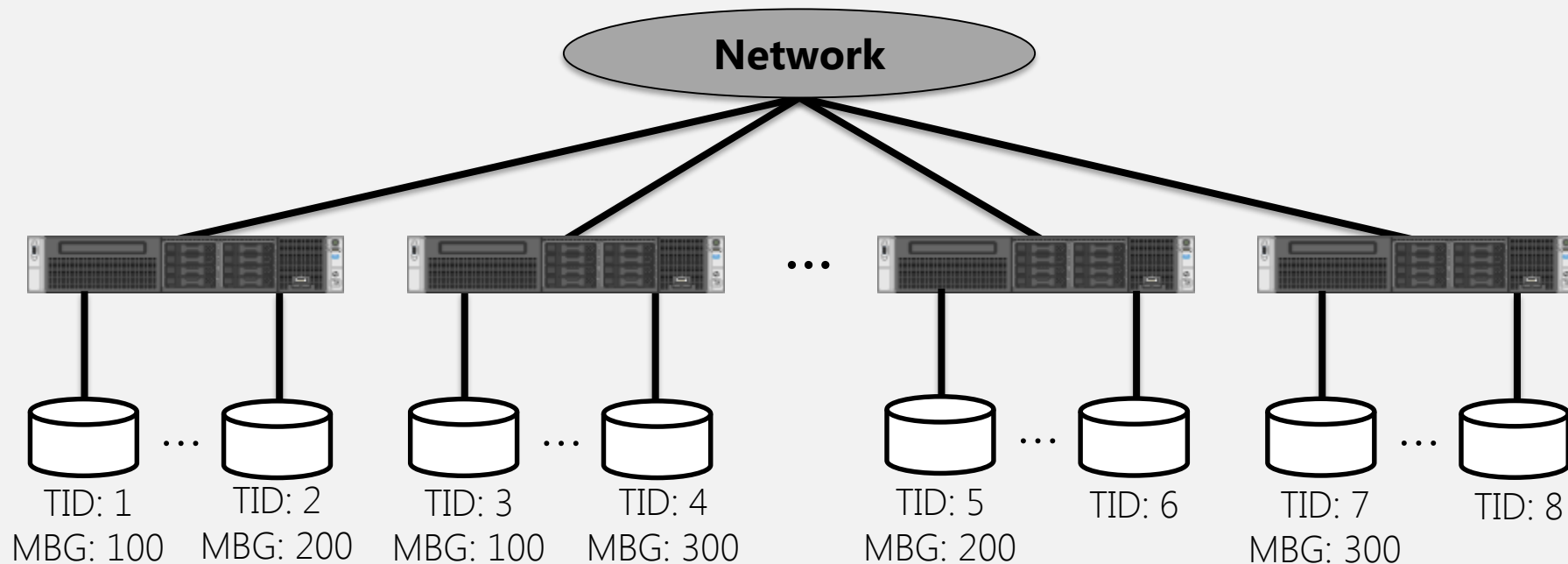
# High Availability - Shared Storage

- No system downtime in case of server failure
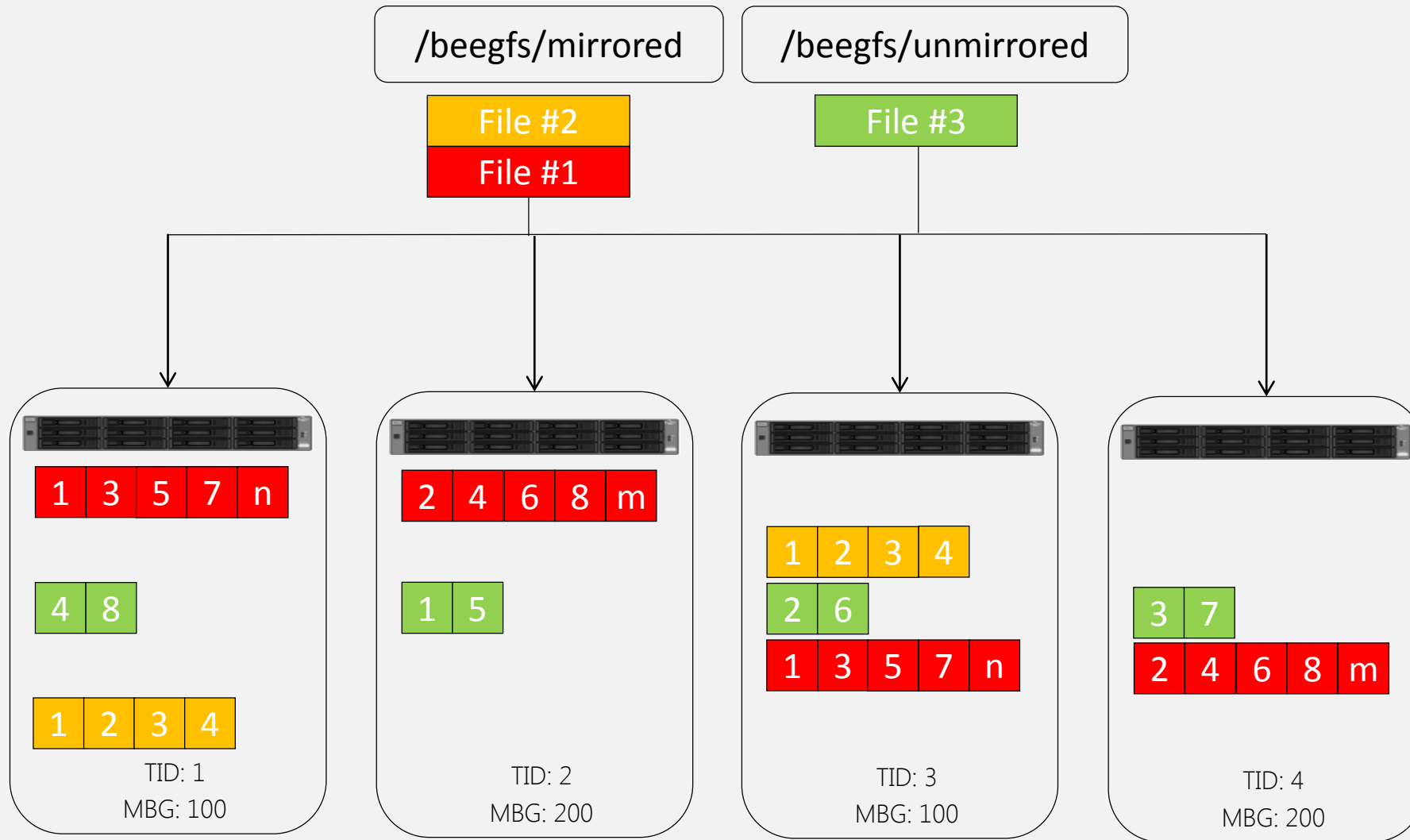- No additional storage capacity needed

- Expensive storage components needed
- 3rd party software components needed
- Complex to set up and maintain
- Failover Risk
- No increased data safety

www.beegfs.com

# High Availability – Built-in Replication

- Assign targets to „mirror buddy groups"

- MBGs replicate chunks (but can also store non-replicated data)

- Internal HA/failover mechanisms

# Built-in Replication - Striping

# High Availability – Built-in Replication

- Flexible (replication configurable per-directory)
- Easy to scale/extend
- No 3rd party tools for monitoring and failover functionality
- Any storage backend can be used
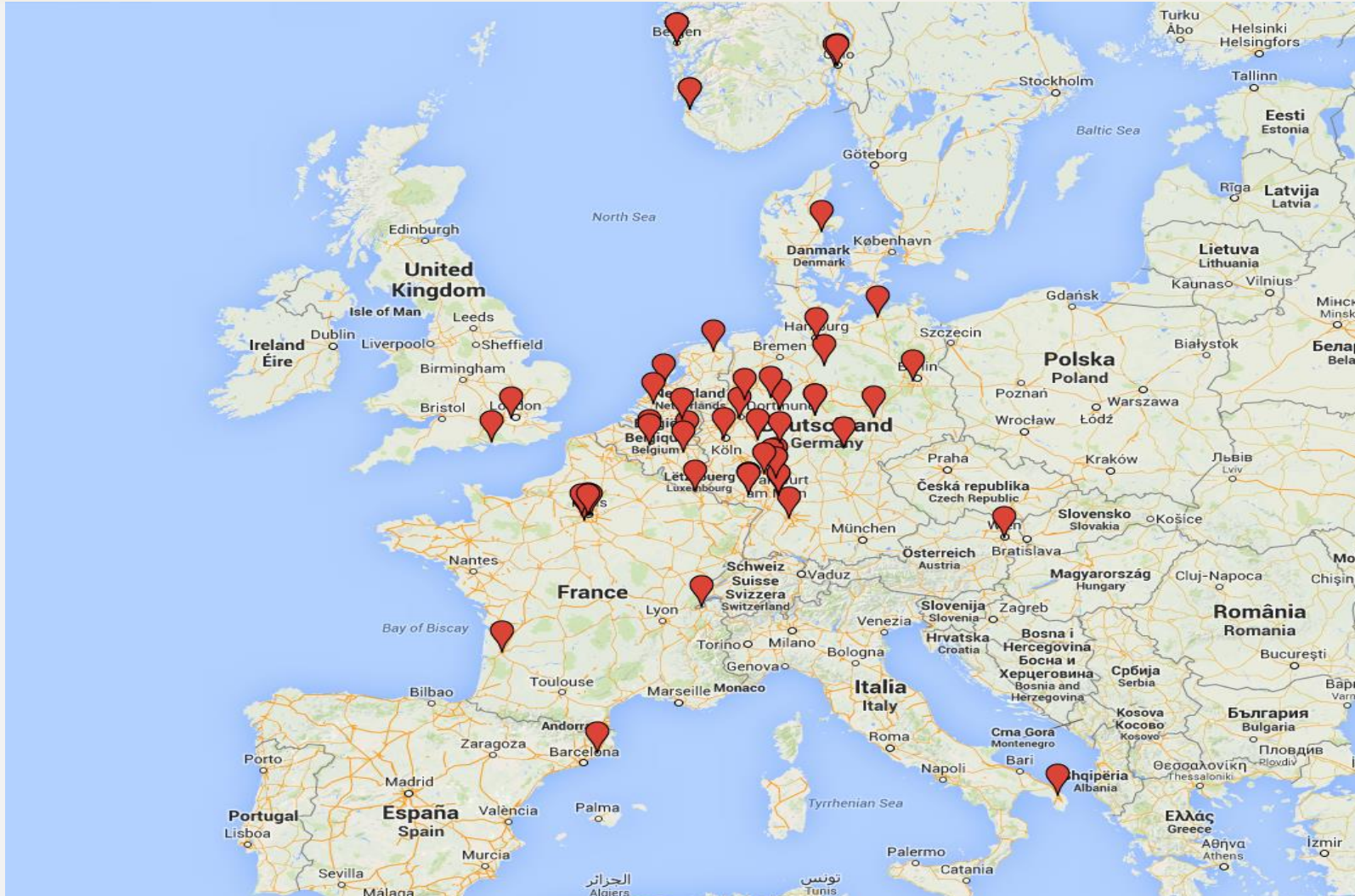- Additional data safety

- Overhead in storage capacity
- Write penalty for replicated data
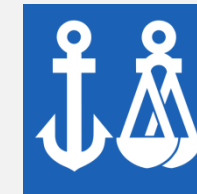
# More than 100 happy customers

# More than 100 happy customers

www.beegfs.com

# Customer Examples

www.beegfs.com

# OcULUS – A Typical HPC Installation

Campus Network

2 HeadNodes

4 LoginNodes

Ethernet

40 GPU Nodes

BeeGFS
developed by Fraunhofer

2 MDS

7 OSS

Capacity: 560 TB
Sustained sequential read/write: 21GB/s

InfiniBand Fabric

572 CPU Nodes
(9600 Cores)

UNIVERSITÄT PADERBORN
Die Universität der Informationsgesellschaft
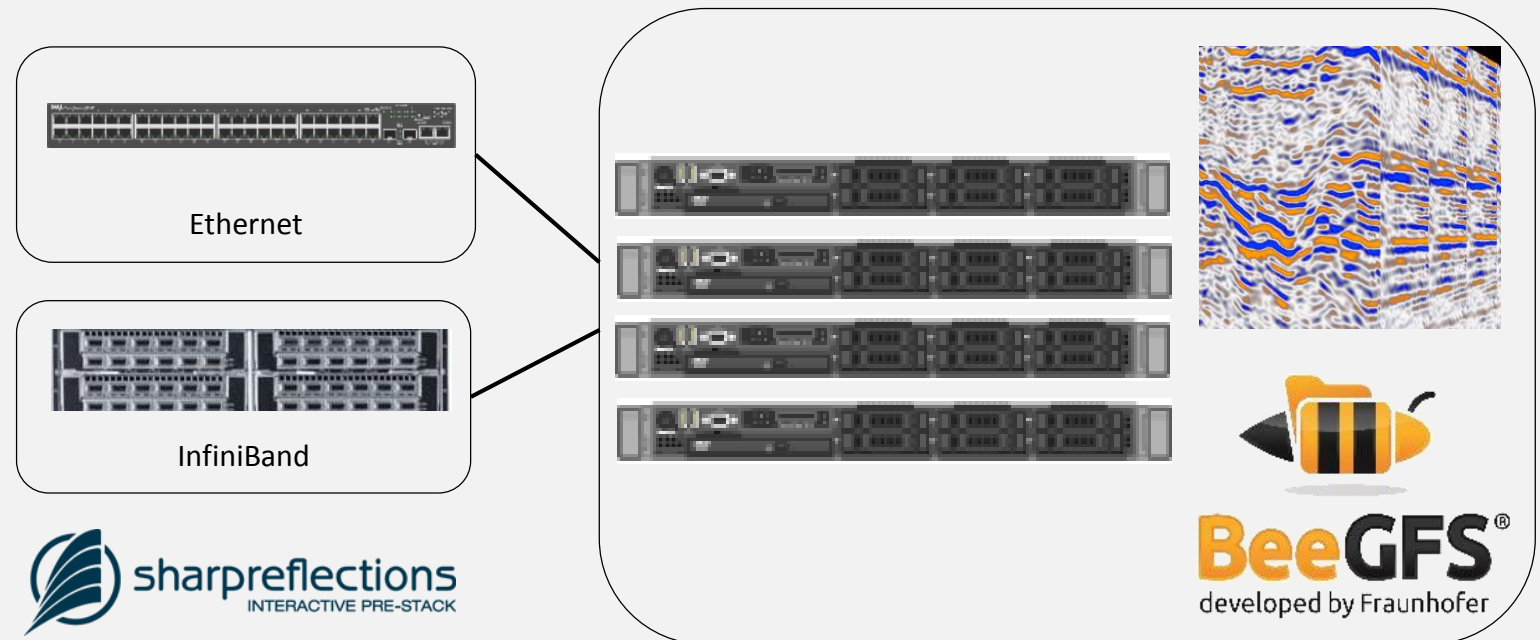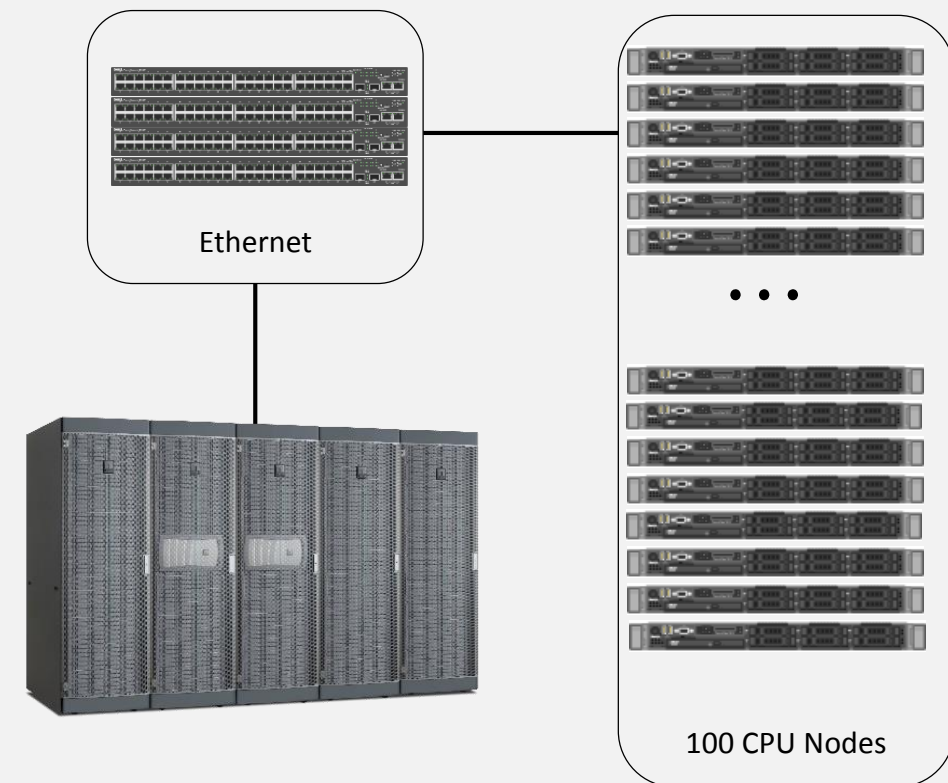
# Sharp Reflections – A Lightweight Solution

- Four compute nodes for seismic data interpretation

- 12 3.5" SATA drives per compute node

- BeeGFS running on compute nodes



Ethernet

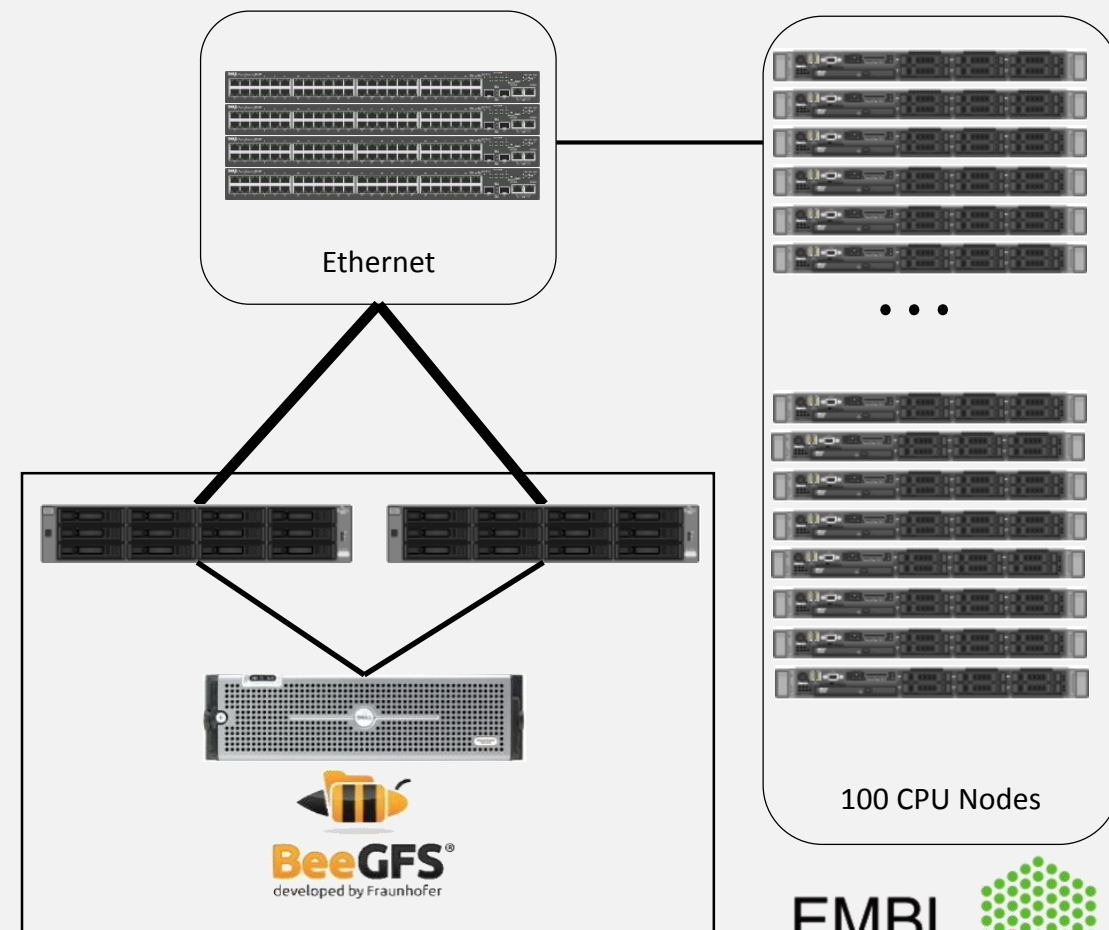InfiniBand

# EMBL – A Workload Optimized Solution

- Situation
  - 100 compute nodes (connected with GigE)
  - „Traditional" NAS storage
  - Only single core compute jobs (life science)
  - Jobs working on the same data
  - Random data access
  - 100k+ IOPS needed
  - ~ 40 jobs saturate NAS



Ethernet

...

100 CPU Nodes

www.beegfs.com

# EMBL – A Workload Optimized Solution

- BeeGFS solution
  - 2 storage servers
  - 30 disk drives each
  - 0.5 TB RAM each
  - 40GbE uplinks for storage
  - Up to 600 jobs at peak performance

Ethernet

···

100 CPU Nodes

BeeGFS®
developed by Fraunhofer

EMBL

# Questions?



- Wiki       wiki.beegfs.com

- Twitter       www.twitter.com/BeeGFS

- MailingList       fhgfs-user@googlegroups.com

- NewsList       beegfs-news@googlegroups.com

- Mail       sales@thinkparq.com

      support@beegfs.com