# A parallel file system – made in Germany

March 7th 2012
Dr. Franz-Josef Pfreundt        Competence Center for HPC
Sven Breuner

**Fraunhofer**
ITWM

# Fraunhofer

## Non Profit  Applied Research

18,000 employees

62 institutes

1.8 billion € budget

### 7 alliances

- Microelectronics
- Production
- **Information and Communication**
- Materials and Components
- Life Sciences
- Surface Technology and Photonics
- Defense and Security Research

Itzehoe
Rostock
Lübeck
Bremerhaven
Bremen
Hannover
Berlin
Braunschweig
Potsdam
Teltow
Paderborn
Magdeburg
Cottbus
Oberhausen Dortmund
Halle Leipzig
Duisburg Schmallenberg
Schkopau
Dresden
Sankt Augustin
Erfurt
Aachen
Jena Chemnitz
Euskirchen
Ilmenau
Darmstadt Würzburg
Erlangen
St. Ingbert **Kaiserslautern**
Fürth
Saarbrücken
Nürnberg
Karlsruhe Pfinztal
Stuttgart
Freising
Freiburg
München
Oberpfaffenhofen
Efringen-Kirchen
Holzkirchen

ITWM-Overview: 3

# Fraunhofer Institut for Industrial Mathematics



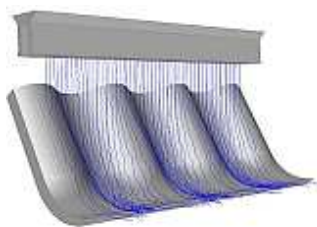Mathematical models

Algorithms

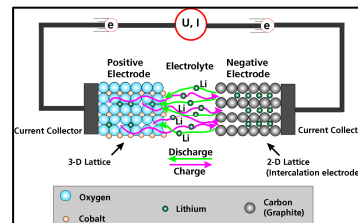Simulations

Software

Visualization

Data mining

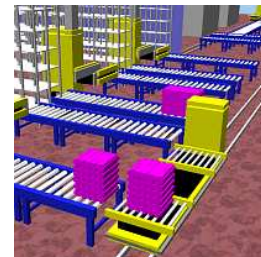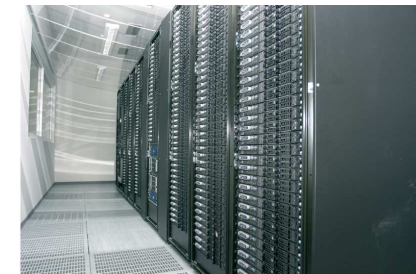Fluid Dynamics          LI-ION Battery Sim          Optimization          HPC

Fraunhofer
ITWM

# Fraunhofer Competence Center for HPC

## Business Fields



**FS**

**GPI-Space**

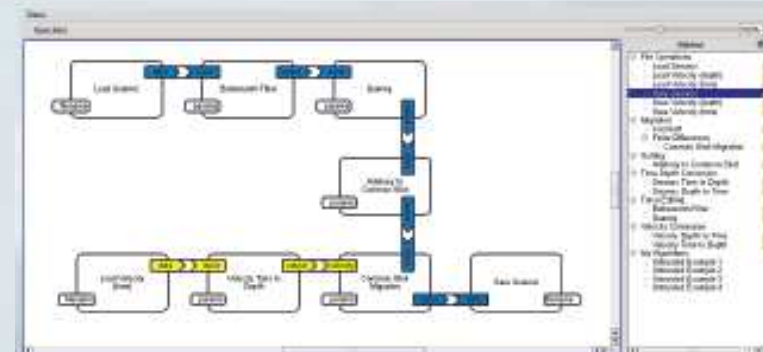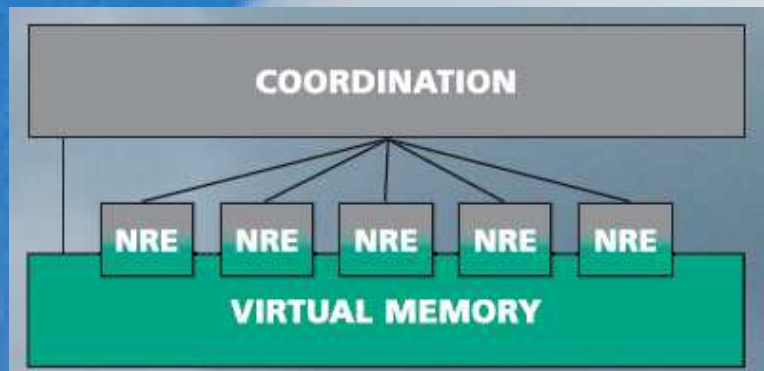| HPC –Tools | Visualization | Green IT | HPC Apps Seismic |

## Research

- parallel programming models
- distributed computing
- parallel algorithms
- parallel file systems

- ray tracing in visualization
- seismic imaging
- distributed energy managment

**Staff : about 40 people**

**Fraunhofer**

**ITWM**

# GPI -SPACE

## Productive Development
## and Execution of Cluster&Cloud Applications

Fraunhofer
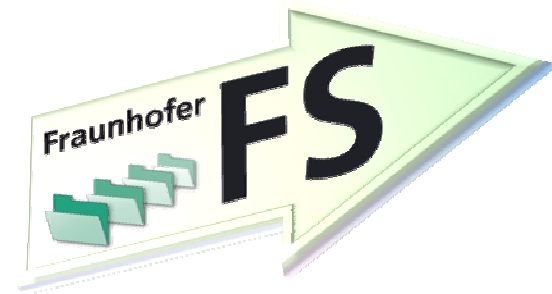ITWM

# GPI Space



| Seismic | Finance | Engineering | Life Sciences |

- ➢ One large distributed  virtual memory space

- ➢ Optimal throughput  - dynamic load balancing

- ➢ Failure tolerant execution

- ➢ Autoparallelization of complex workflows

Fraunhofer
ITWM

# FraunhoferFS      How it started

As part of a cooperation with Linux NetworX

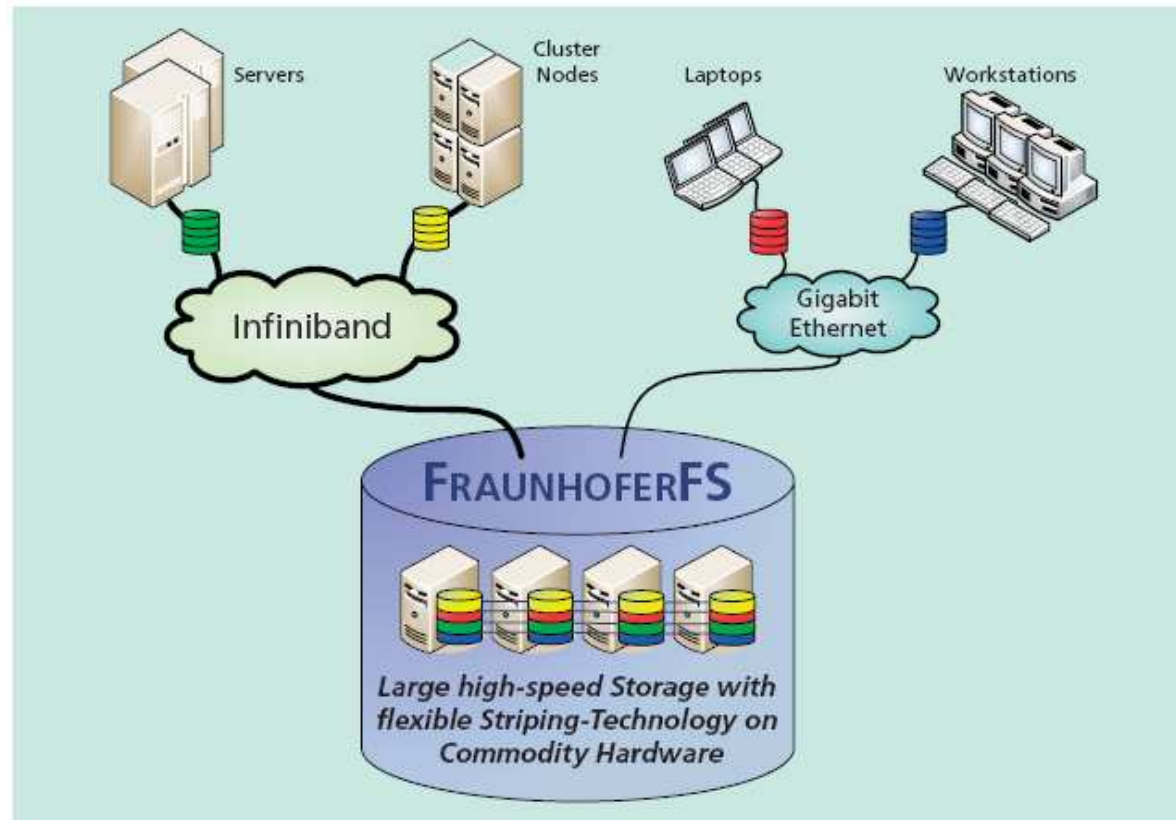2003  First Lustre Installation at Fraunhofer
2004 Port of the Blue Order Media Server System
        on top of Lustre

2004 Decision to develop the Fraunhofer FS

Requirements:    Distributed Metadata
                 No Kernel patches, zero config clients
                 Scalable multithreaded architecture
                 Native IB and Ethernet
                 Easy to install and maintain
                 Use P2P technology

# Anouncement after 3 years of development (ISC Dresden)



Fraunhofer Parallel Filesystem

Available Q4/2007

Fraunhofer
ITWM

# SC 2007 Reno Introduction of the FhGFS

# SC 2010 New Orleans



**Customer Base**

**Oil&Gas**
**Universities**

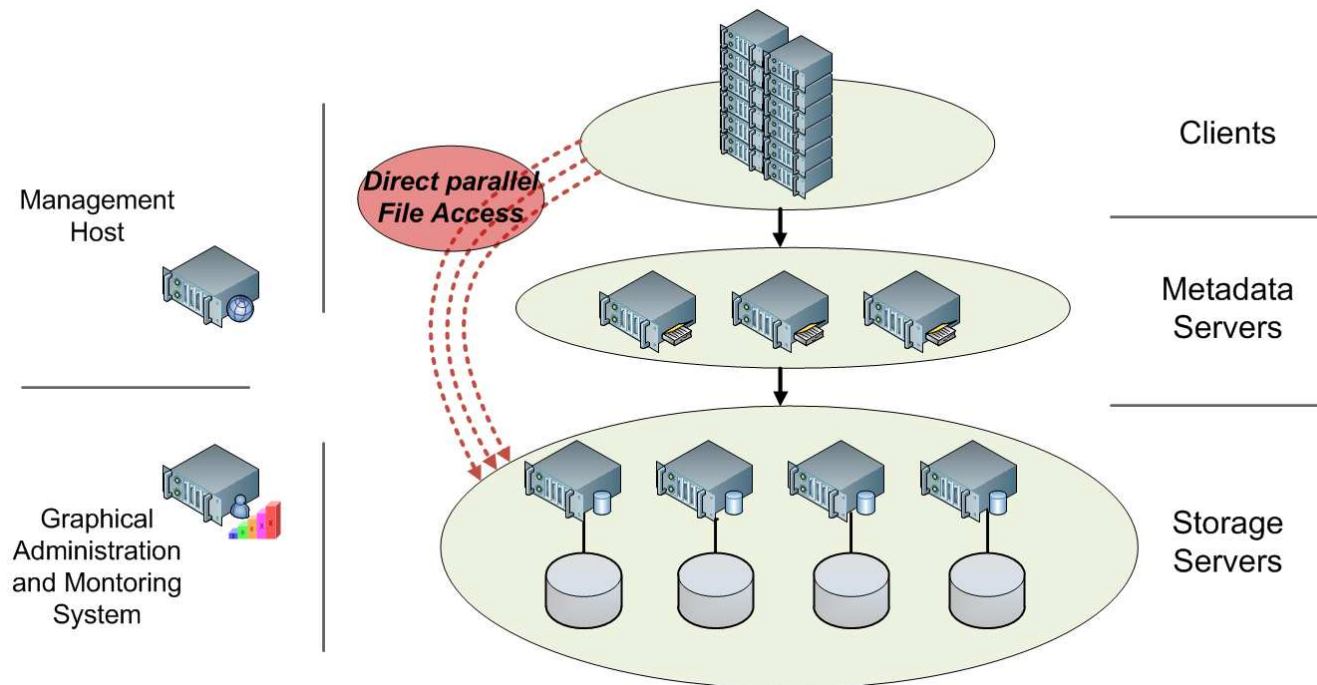**~ 30 supported**
     **customers**

# FhGFS    Key Features

❑ **Maximum Scalability**
- ○ Distributed File Contents & Metadata
- ○ Low server load – efficient multithreading
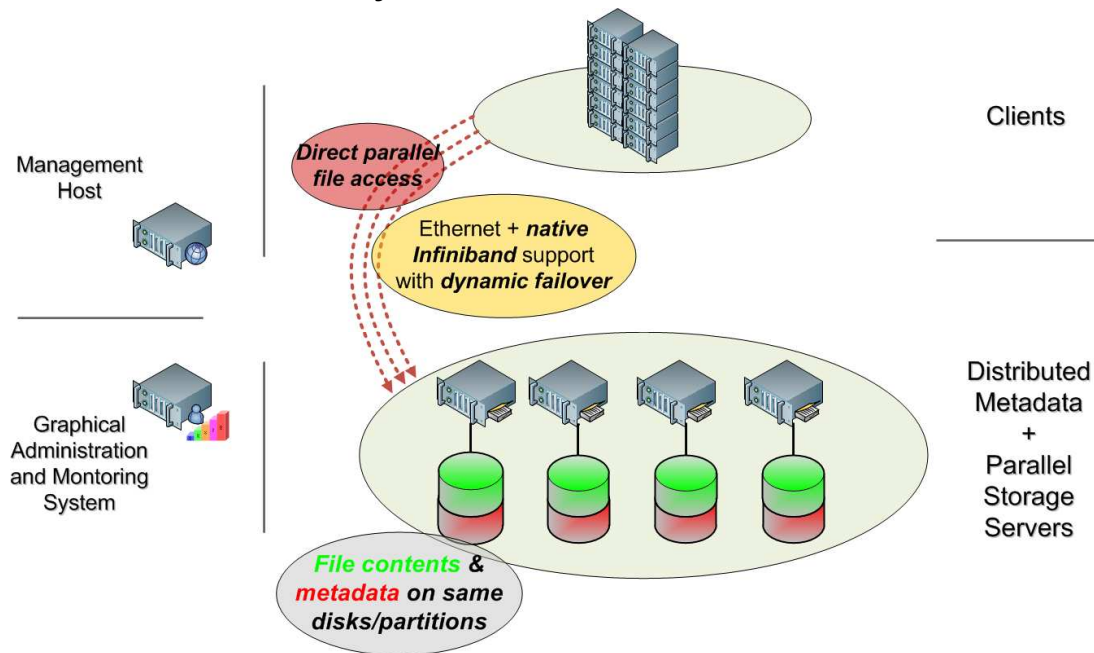
Installed system
More than 300 servers



Object storage , servers use a local file system (XFS, EXT,…)

# FhGFS    Key features

❑ **Flexibility**

- Add Clients and Servers without Downtime
- Client and Servers can run on same Machine
- On-the-fly storage init (mkfs)
- Multiple Networks with dynamic Failover

➡ Storage Cluster
Compute  + Storage



○ Flexible Striping: individual Settings on a per-File /per-Directory Basis

Fraunhofer
ITWM

# Fraunhofer Seislab   - interactive seismic imaging
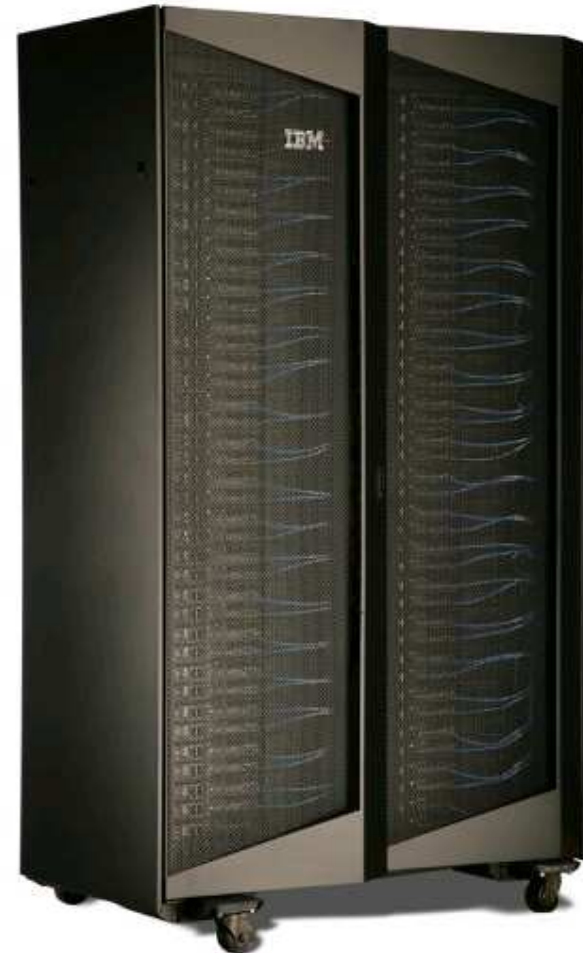
## Compute & Storage

20 Compute Nodes
   48 -96 GB RAM
   4 x 256 GB SSD striped
   QDR Infiniband

 5 Compute&Storage Nodes
   20 TB SATA , RAID5 (Archive)
   QDR Infiniband

On demand SSD based FhGFS
per job up to 20 TB
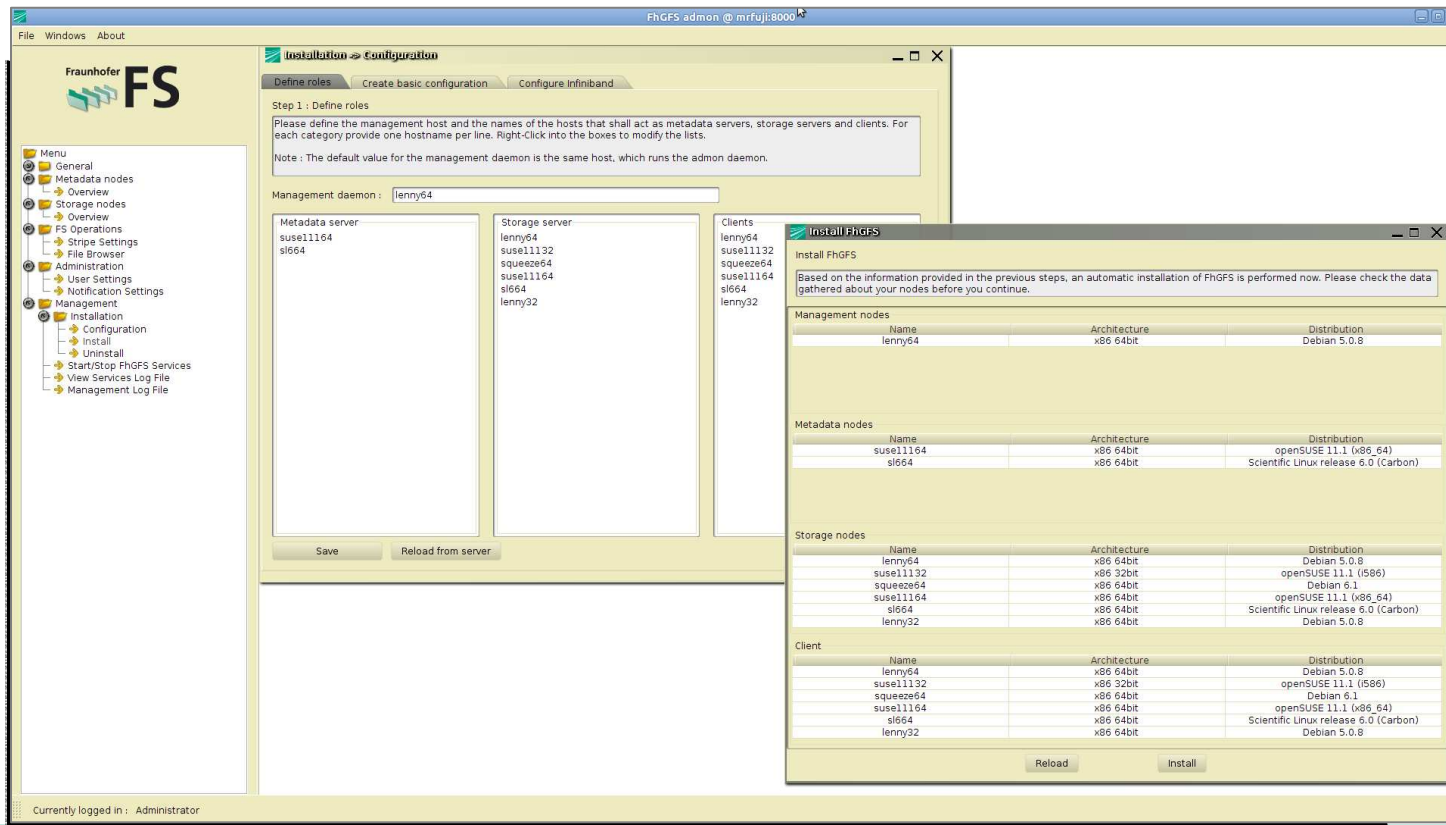
Read: 30 GB/sec  Write: 20GB/sec

## Network bisection BW  ~  I/O performance

Fraunhofer
ITWM

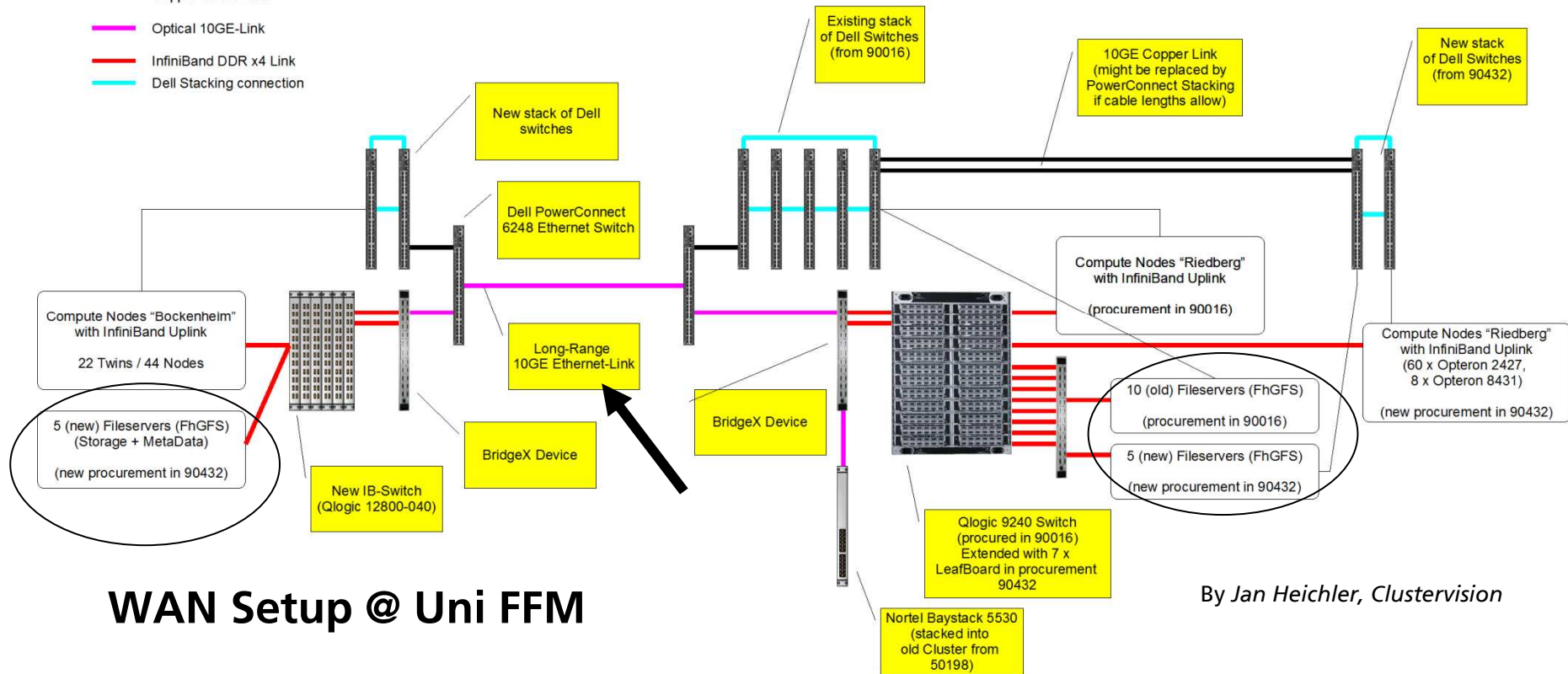# FhGFS   Key Features

## Easy to use
- Automated Cluster Installation
- Kernel Module
- Graphical System Administration & Monitoring

- No specific Linux Distribution,
- no special Hardware required



© Fraunhofer ITWM 2012

# FhGFS  Key Features

➢ Light-weight Client Kernel Module

❑ High Single-Stream Throughput (>2.7GB/s on QDR IB)

➢ Server Preference

❑ Clients can prefer a Subset of Servers  => Support for multiple Data Centers
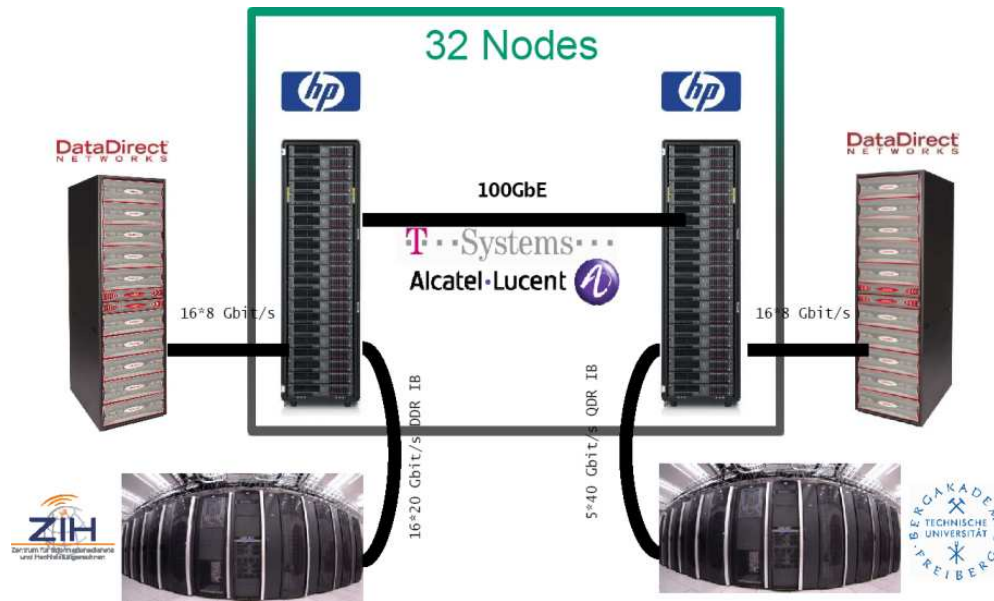


**WAN Setup @ Uni FFM**

By *Jan Heichler, Clustervision*

# Typical Size and performance of current installations

Frankfurt University :   12 servers ,1 PByte  , 20GB/sec , 900 clients
                         measured single stream performance : 2,0 GB/sec

TU Vienna             :   12 servers,  300 TByte, 6GB/sec, 1200 clients
                         12  metadata server(SSD), x00 000 I/O Ops/sec

RSI (Houston)         :   12 server, 300 Tbyte, 6 GB/sec, 28 clients
                         client and server on same machine

Fraunhofer Seislab  :   20 servers, 20 TB SSD,120TB SATA, 30 GB/sec
                         server and clients

DTU Kopenhagen     :   5 servers, 200 TByte, 5GB/sec, 100 clients
                         port to BSD UNIX

Fraunhofer
ITWM

# 100GBit Testbed (Dresden <-> Freiberg)



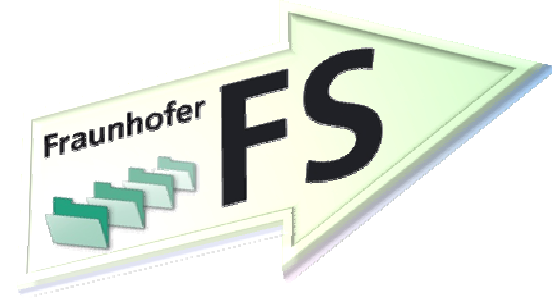By *Michael Kluge, TU Dresden*

## Uni-directional

- GPFS  - 10,1 GB/s (60 km)
- Lustre - 11,8 GB/s (60 km)
- FhGFS - 12,4 GB/s (400 km)

## Bi-directional

- GPFS  - *n/a*
- Lustre - 21,9 GB/s (60 km)
- FhGFS - 22,5 GB/s (400 km)

# Last major release August 2011

- Client operation counters

- All file attributes stored on metadata server

- **Distributed POSIX file locking**

- Simplified automatic updates via software repositories

- Multiple storage targets per server

- Re-designed metadata request handling to scale to high numbers of CPU cores

- Parallel online file system check/repair

*Faster,*

*more flexible,*

*easier to use*

**Business Model**

# No license fees

# Pay for support and maintenance

Open Source  - on a individual basis
So far not a community request

Fraunhofer
**ITWM**

# Our supported customers  ( ~ 50)

**HPC Centers**

GOETHE UNIVERSITÄT FRANKFURT AM MAIN

universität wien

**University Oslo**

Universität Stuttgart

DESY

University of Zagreb

UNIVERSITÉ BORDEAUX 1 Sciences Technologies

university of groningen

**Oil&Gas**

MARATHON

Statoil

NORECO

RSI

REPSOL

DETNORSKE

... And  more

Cloud Computing

Social Media

No system Halt for Software resons

Happy users

Fraunhofer ITWM

# About the FhGFS Roadmap

Some FhGFS roadmap pillars are
fixed, e.g.:

- HA
- HSM

We leave some room to implement
interesting user ideas, e.g.:

- Server affinity
- Client operation counters

We learned that we need to leave
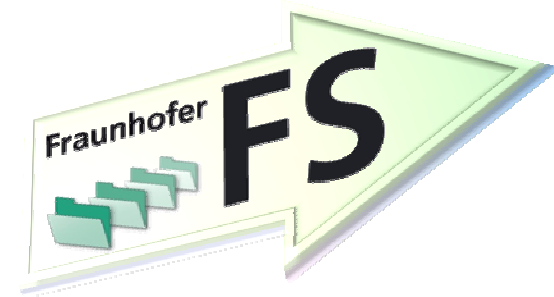some room to improve Linux
kernel / tools, e.g.:

- tail, ls -l, Linux RDMA

And we have enough people in the
institute that develop HPC
appplications with disruptive new
ideas

- Fraunhofer Seislab

  - 20 compute nodes with SSDs

  - Runs FhGFS on-demand

  - Jobs store temporary data on SSDs and
    move it to dedicated servers afterwards

  - 20GB/s write, 25GB/s read sustained

# Next major release 2012

- Data/metadata mirroring over multiple FhGFS servers

- Configurable on a per-file (per-directory) basis

- Server groups for remote mirroring

- Quota/ACL support

- MAC support（Q2 2012)

→ Next major release  Q2 2012

Fraunhofer
ITWM

# Hierarchical Storage Management



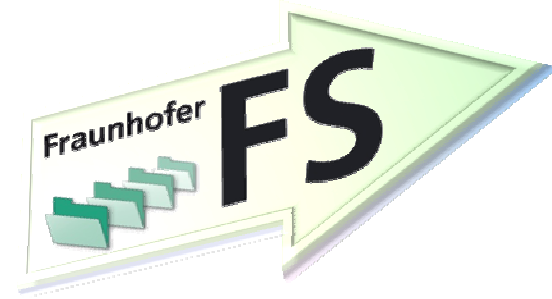GrauData provides Grau ArchiveManager (GAM) as a solid single-server HSM solution

Fraunhofer and Grau teamed up to integrate GAM and FhGFS

The combined solution will support…

- Parallel data migration
  (e.g. recall all file chunks at once)

- Collocation IDs

- Asynchronous recalls

First prototype will run at HLRS Q3 2012

# Questions ?



## http://www.fhgfs.com

Franz-Josef Pfreundt , Sven Breuner
pfreundt@itwm.fhg.de ,breuner@itwm.fhg.de

Fraunhofer
ITWM