# *Project 'Carme'*

## *An Open Source Framework for Multi-User, Interactive Machine Learning and Data Analytics on Distributed (GPU) Systems*



### *www.open-carme.org*

Competence Center High Performance Computing Fraunhofer ITWM, Kaiserslautern

Fraunhofer

# Motivation
## Data Analytics and Machine Learning

- Large investments in (multi)-GPU hardware for data analytics and machine (deep) learning

- Main problems AFTER buying the hardware:

  - How to manage the resources?
  - How to scale applications to more than one GPU?
  - How to manage data I/O and storage?

Fraunhofer

# Motivation
## Data Analytics and Machine Learning

- Large investments in (multi)-GPU hardware for data analytics and machine (deep) learning

- Main problems AFTER buying the hardware:

  - How to manage the resources?
  - How to scale applications to more than one GPU?
  - How to manage data I/O and storage?

**NOT a solution*:**



**+**



* we are seeing this being done

**Part I: Low Cost Hardware Setup**

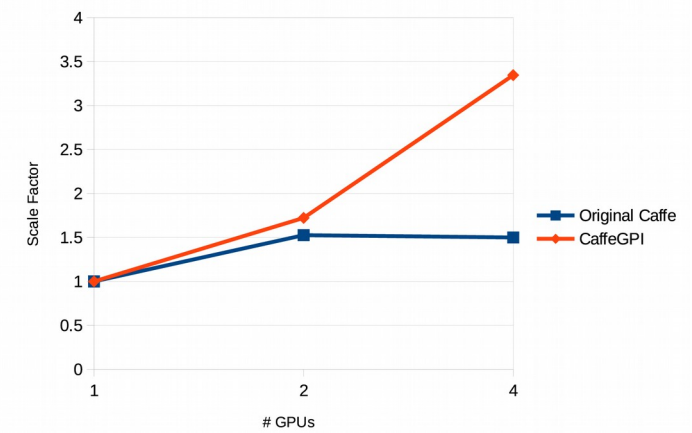**ITWM's high performance multi-GPU cluster build on gaming hardware**

**Part II: '*Carme*'**

**ITWM's open source software stack for multi-user GPU clusters**

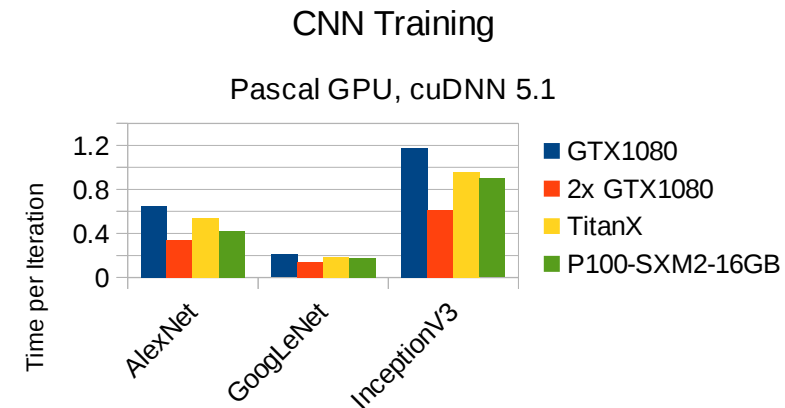Fraunhofer

# Deep Learning Setup
## Multi-GPU Hardware Overview

- **High Performance (NVLINK) Multi-GPU system are expensive!**

  - NVIDIA DGX-1 (8x P100/V100) ~ 130k EUR*

  - IBM Witherspoon (6x V100) ~ 60k EUR*

- **PCIe hosting of GPUs is to slow (!)**

  → **no multi-GPU training possible**

- **Gaming GPUs like GTX1080ti are a cheap alternative**

  - e.g. price GTX1080ti (11 GB) ~600 EUR*,

    P100 (16 GB)~ 6000 EUR*
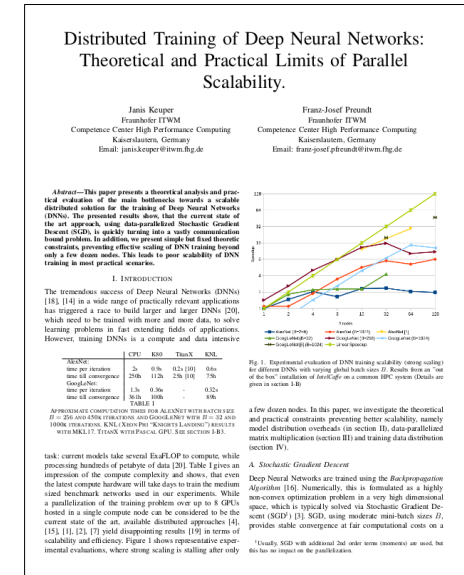
  - Some Benchmarks:

* estimated retail prices



Scaling performance of a Multi-GPU
training on a PCIe System (blue line)

CNN Training

Pascal GPU, cuDNN 5.1

# Low Cost Deep Learning Setup
## ITWM Approach

- **Research results: GPU performance is hardly ever the bottleneck**
  - I/O is the dominating factor
    - GPU-2-GPU and GPU-2-CPU communication
    - Train data I/O

- **Main Ideas:**
  - Combine cheap gaming hardware (GPUs, CPU and main boards) with sophisticated HPC components (Network)
  - Use existing ITWM HPC Software-Stack
    - BeeGFS distributed File System
    - GPI Communication Model
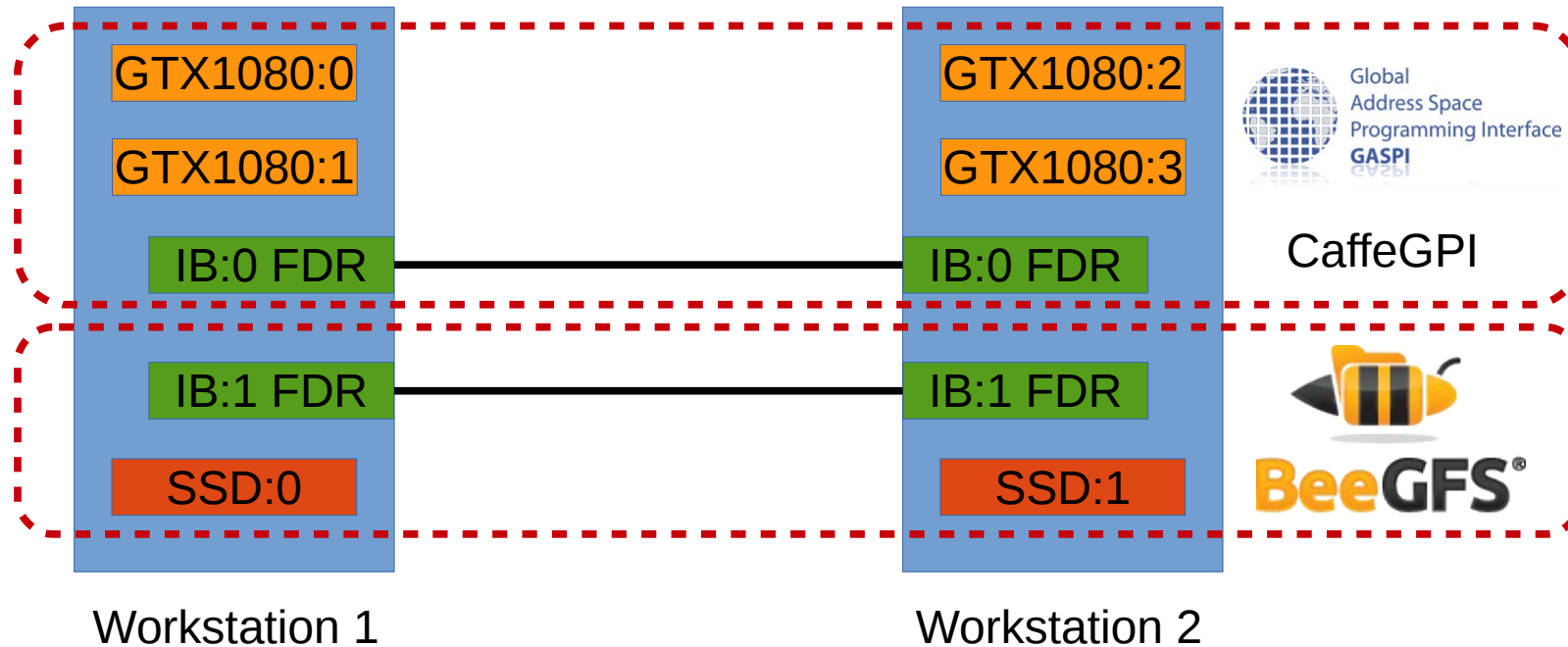    - CaffeGPI and HP-DLF Deep Learning Frameworks

Based on our paper from
ACM SuperComputing 16

# Low Cost Deep Learning Setup
## Prototype 2-Node System

**Price: <8000 EUR, standard components**
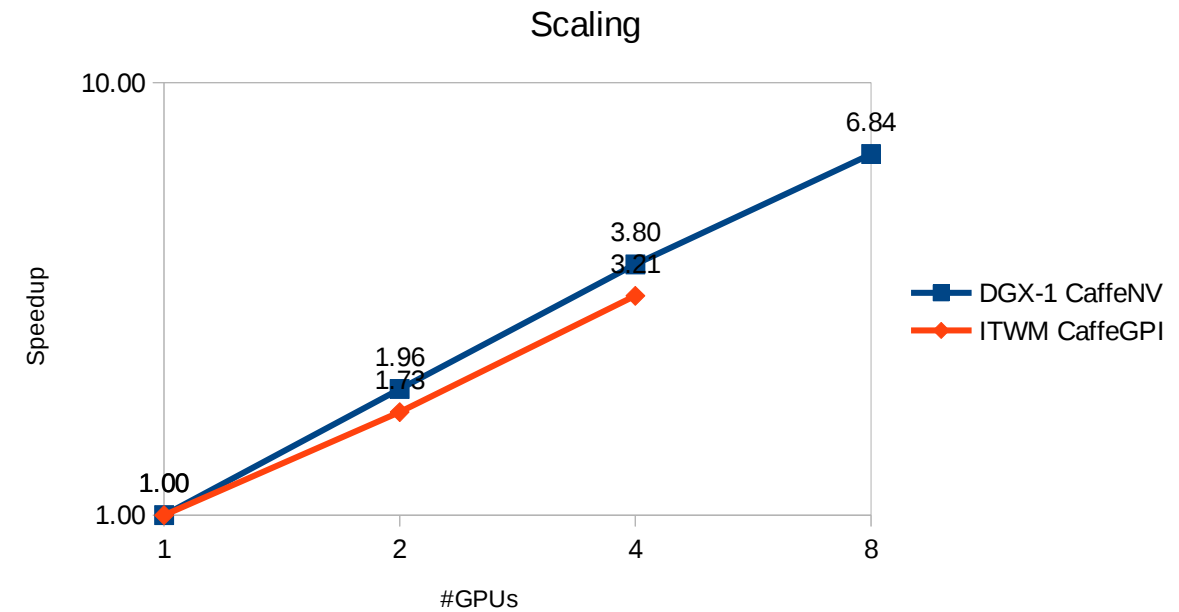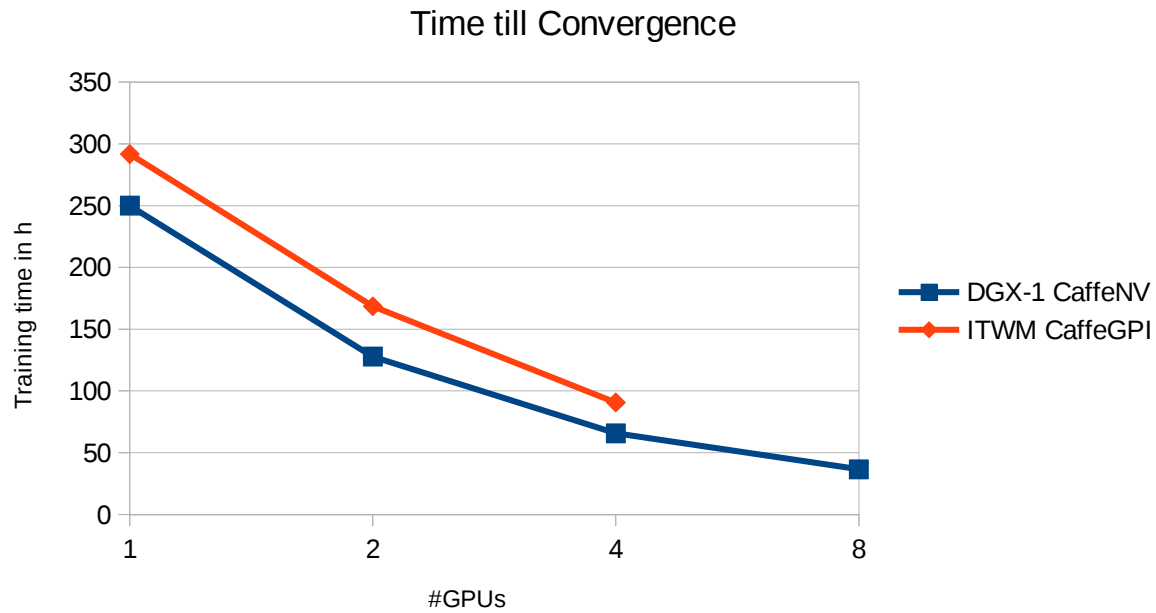


**Specs**:

- 32 GB GPU Mem
- 64 GB PGAS Mem
- 2TB BeeGFS for Train Data
- GPU interconnect: PCIe

# Low Cost Deep Learning Setup
## Prototype 2-Node System Benchmarks
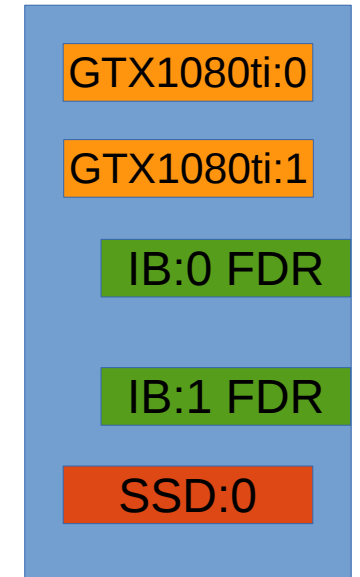
**Time till Convergence**



**Scaling**



Specs: Topology: GoogLeNet, Cuda 8, cuDNN 5.1, CaffeNV 16.4, Batch Size/Node: 64

# Low Cost Deep Learning Setup
## Currently building: Prototype 16-Node System

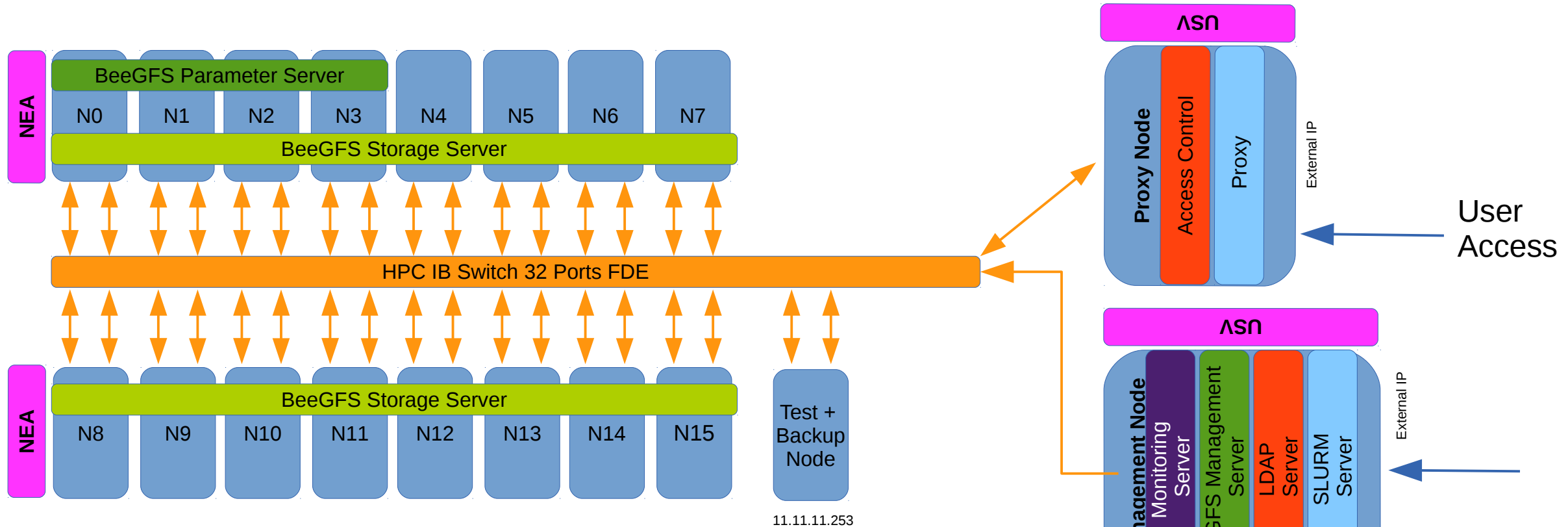**for BMBF founded Fraunhofer Consortium @ITWM**

- **Low cost Hardware**
  - Consumer GPUs (GTX1080ti)
  - Novel AMD architecture provides enough PCIe-lanes
  - Hosting Cost per GPU ~ 1.25 k EUR. Compared to DGX-1 ~ 10k

- **Fast HPC Interconnect and Data I/O**
  - Infiniband Network
  - Parallel HPC file system with local SSD

- **Multi-GPU performance**
  - Scalable multi-GPU training

GTX1080ti:0

GTX1080ti:1

IB:0 FDR

IB:1 FDR

SSD:0

Node configuration

# Low Cost Deep Learning Setup
## Currently building: Prototype 16-Node System

# Project Carme

## An open source software stack for multi-user GPU clusters

*Carme* (/ˈkɑːrmiː/ KAR-mee; Greek: Κάρμη) is a **Jupiter** moon, also giving the name for a **Cluster** of Jupiter moons (the carme group).

Or in our case:

an open source frame work to mange resources for  multiple users running **Jupyter** notebooks on a **Cluster** of compute nodes.

# Project Carme

## An open source software stack for multi-user GPU clusters

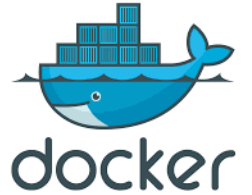**Common problems in GPU-Cluster operation:**

- **Interactive, secure multi user environment**
  - ML and Data Science users want interactive GUI access to compute resources

- **Resource Management**
  - How to assign (GPU) resources to competing users?
    - User management
    - Accounting
    - Job scheduling
    - Resource reservation

- **Data I/O**
  - Get user data to compute nodes (I/O Bottleneck)

- **Maintenance**
  - Meet (fast changing and diverse) software demands of users

≡ Fraunhofer

# Project Carme

## An open source software stack for multi-user GPU clusters

***Carme* core idea:**

- **Combine established open source ML and DS tools with HPC back-ends**
  - Use containers
    - (for now) Docker
  - Use Jupyter Notebooks as main web based GUI-Frontend
    - All web front-end (OS independent, no installation on user side needed)
  - Use HPC job management and scheduler
    - SLURM
  - Use HPC data I/O technology
    - ITWM's BeeGFS
  - Use HPC maintenance and monitoring tools

# Project Carme
## An open source software stack for multi-user GPU clusters

**Carme features:**

- **Open source**
  - Carme uses only opensource components that allow commercial usage
  - Carme is open source, allowing commercial usage
- **User Management**
  - User quotas (GPU time, priority, GPUs per job, jobs per time, Disk quota)
  - Different User Roles (Quotas, right to add containers)
- **Container Management**
  - Container store (user selects from predefined containers)
  - Adding of user defined containers

Fraunhofer

# Project Carme
## An open source software stack for multi-user GPU clusters

***Carme* features:**

- **Scheduler**
  - Resource reservation (calender)
  - Job queues for large jobs and instant interactive access for small jobs
- **Data Management and I/O**
  - Redundant, global file system (BeeGFS), mounts into container
  - Temporary job FS on local SSDs for max performance (BeeOND)
- **Web-Interface**
  - HTTPS and SSH (if allowed) access via proxy
  - Web front-end (management and IDE)

Fraunhofer

# Project Carme
## An open source software stack for multi-user GPU clusters

**_Carme_ features:**

- **Scalable Framework**
  - Single GPU to distributed multi-GPU scaling
  - Add GPUs to running job
  - Strong and weak scaling of DL training
  - Works alongside existing _Slurm_ systems on existing (HPC) clusters
- **Cluster Maintenance + Monitoring**
  - auto Worker updates
  - Easy hardware scale up (adding more compute hardware later)

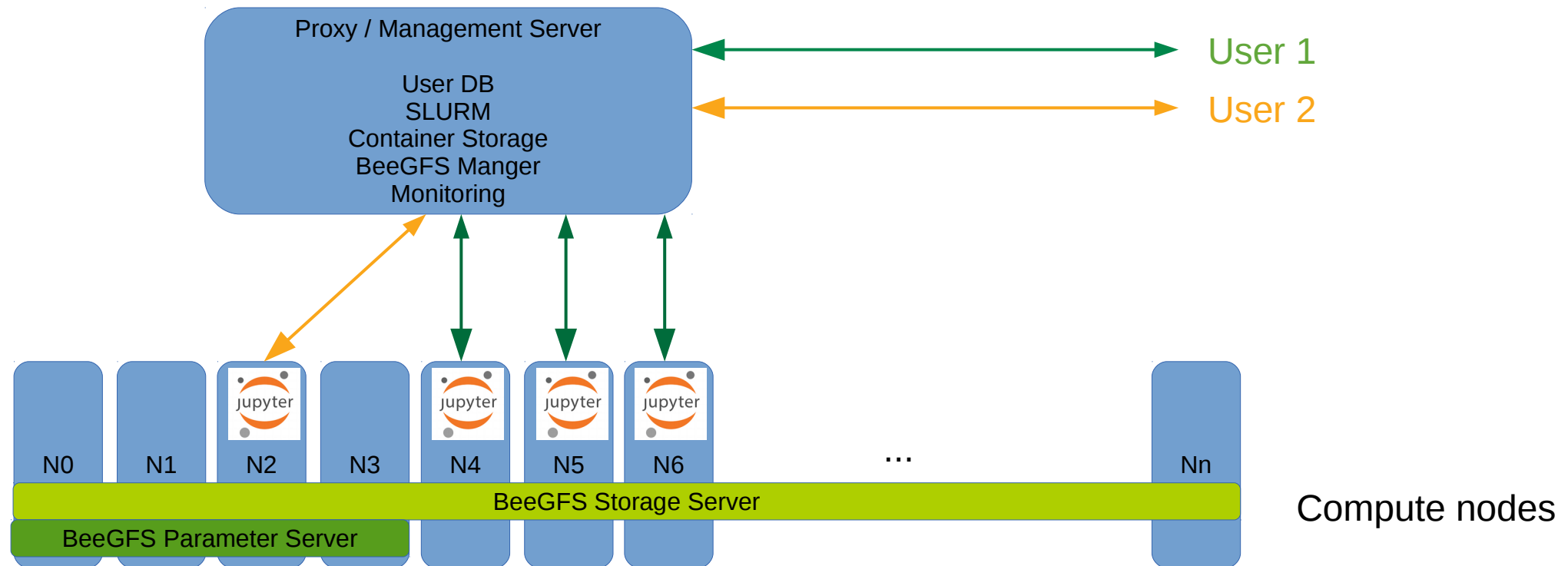Fraunhofer

# Project Carme
## An open source software stack for multi-user GPU clusters

***Carme* features:**

- **Web-Interface**
  - HTTPS and SSH (if allowed) access via proxy
  - Web front-end (management and IDE)
  - Jupyter Notebooks + plugins
  - Tensorboard Server  (other web application possible)

Fraunhofer

# Project Carme
An open source software stack for multi-user GPU clusters

# Project Carme
## An open source software stack for multi-user GPU clusters

**Carme road map:**

Running beta Version by 04/18

Public 0.9 beta release by 06/18 (ISC High Performance Conference)

→ AWS live demo

Version 1.0 with deployment tools by 12/18 (NIPS)

Contact:       janis.keuper@itwm.fhg.de

Fraunhofer