

Extreme Scale Computing – Challenges and Costs

Dieter Kranzlmüller

Munich Network Management Team
Ludwig-Maximilians-Universität München (LMU) &
Leibniz Supercomputing Centre (LRZ)
of the Bavarian Academy of Sciences and Humanities



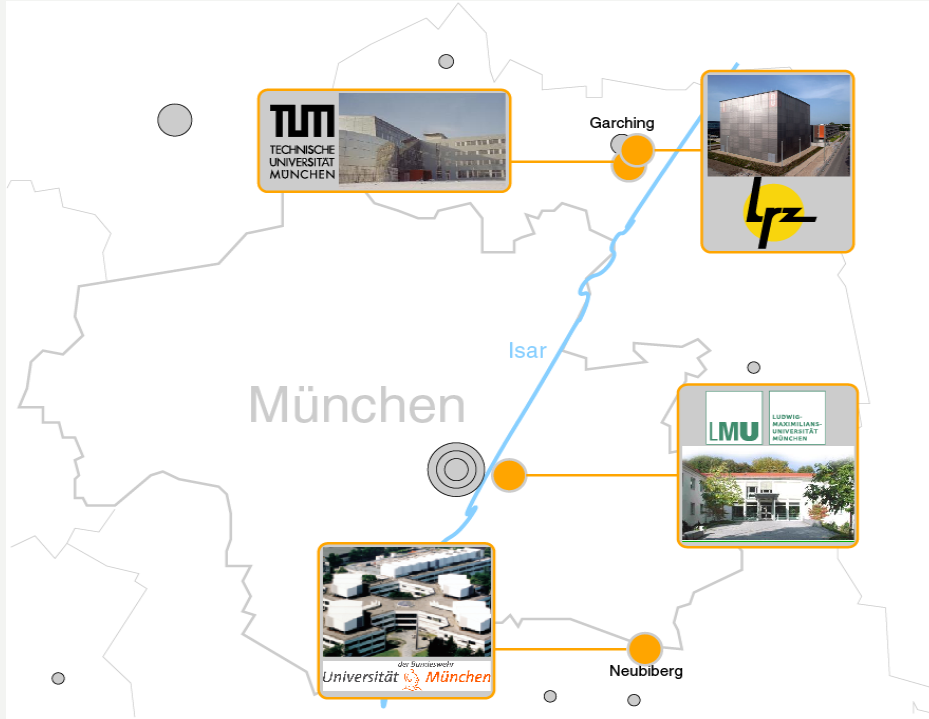
- Research at MNM-Team
- Introduction of the Leibniz Supercomputing Centre (LRZ)
- Roles of LRZ (Munich, Bavaria, Germany, Europe)
- Introduction of SuperMUC
- Challenges: Building, Scale, Power → Budget
- Discussion



LUDWIG-MAXIMILIANS-UNIVERSITÄT MÜNCHEN

MNM TEAM

MUNICH NETWORK MANAGEMENT TEAM



LUDWIG-MAXIMILIANS-UNIVERSITÄT MÜNCHEN

MNM TEAM

MUNICH NETWORK MANAGEMENT TEAM



Networks



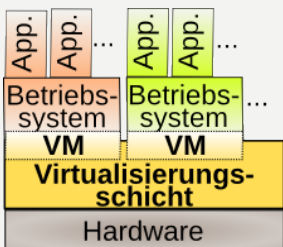
Grid computing



Cloud Computing

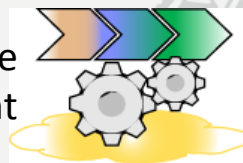


High Performance Computing



Virtualization

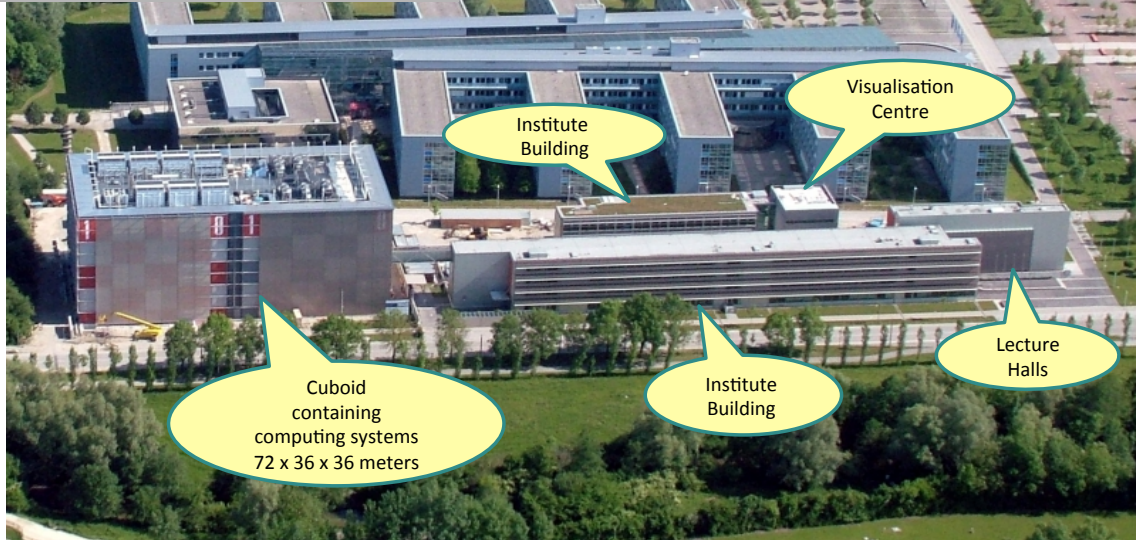
Service Management



IT Security



With 156 employees + 38 extra staff for more than 90.000 students and for more than 30.000 employees including 8.500 scientists



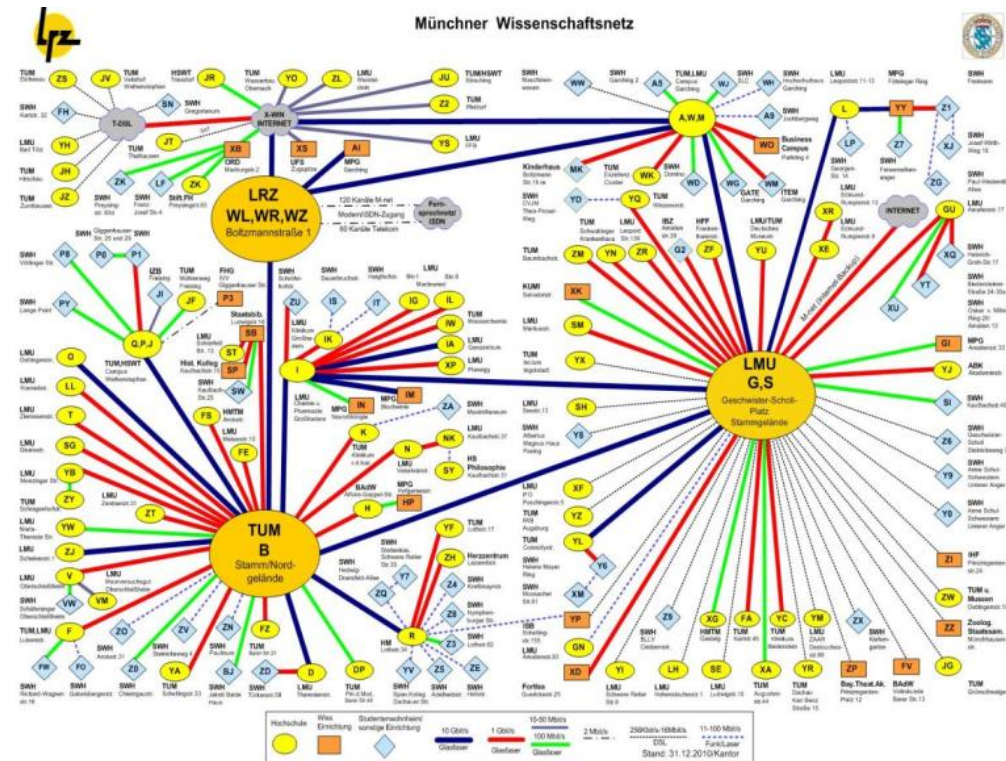
■ Computer Centre for all Munich Universities

IT Service Provider:

- Munich Scientific Network (MWN)
- Web servers
- e-Learning
- E-Mail
- Groupware
- Special equipment:
 - Virtual Reality Laboratory
 - Video Conference
 - Scanners for slides and large documents
 - Large scale plotters

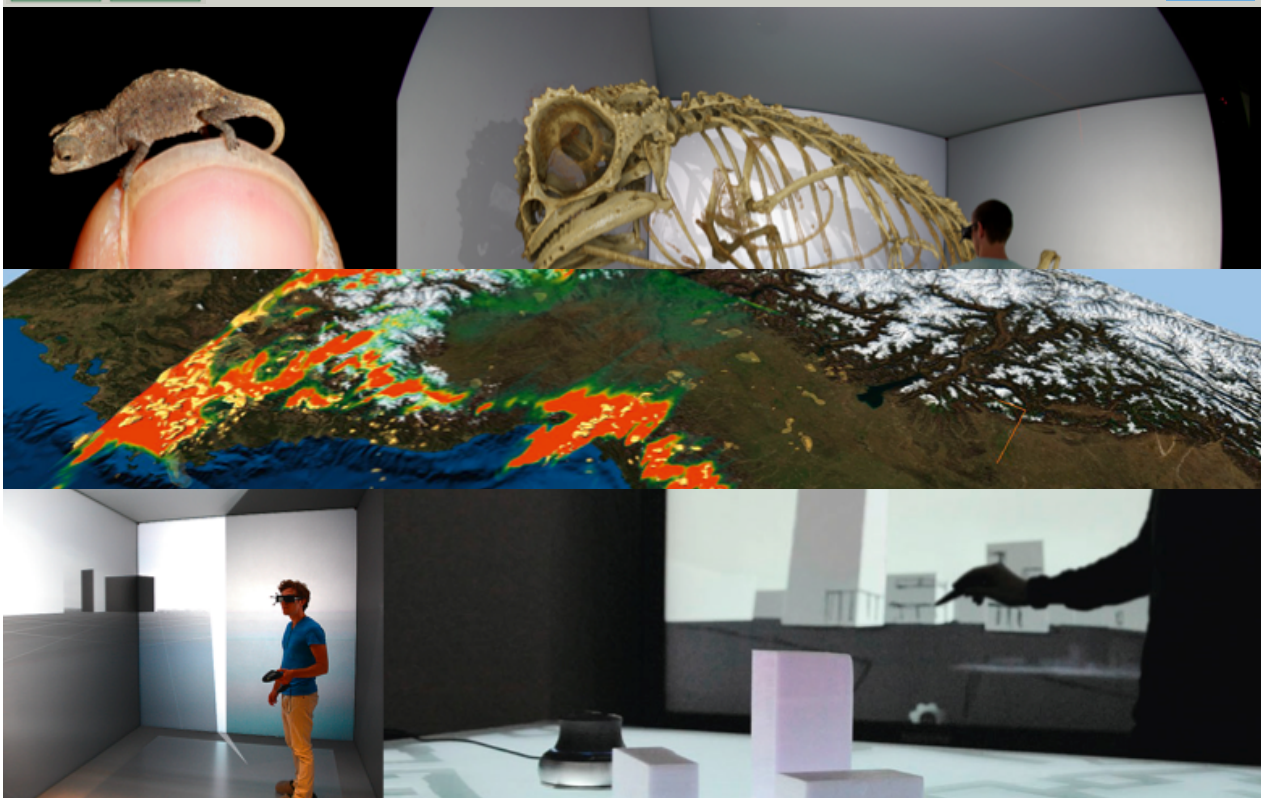
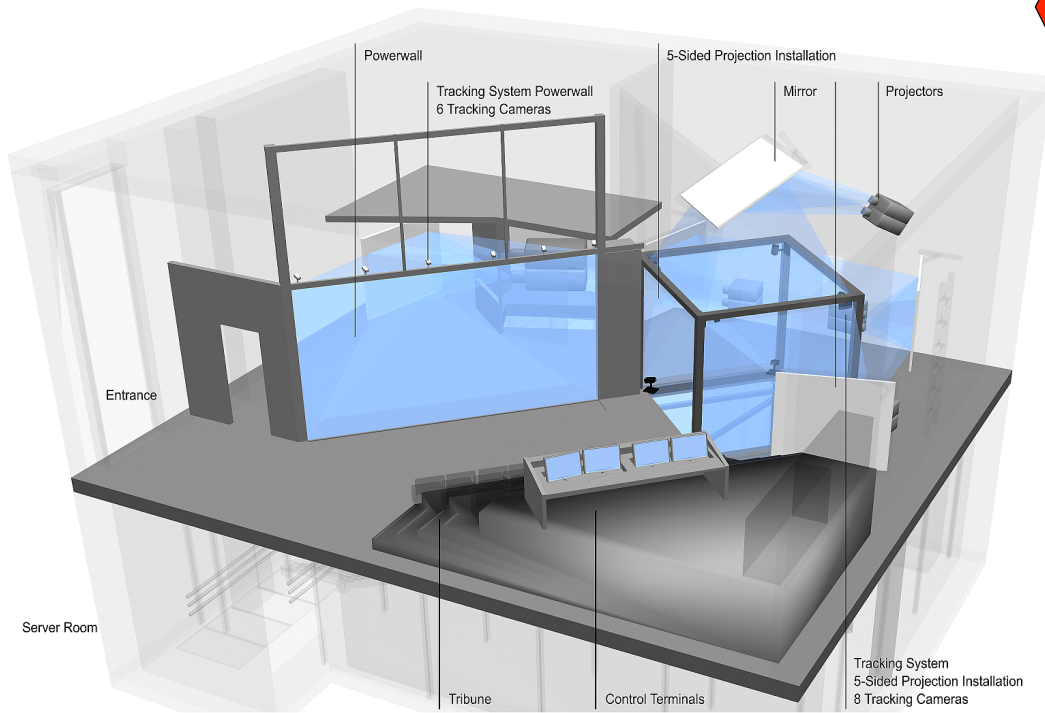
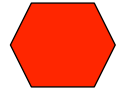
IT Competence Centre:

- Hotline and support
- Consulting (security, networking, scientific computing, ...)
- Courses (text editing, image processing, UNIX, Linux, HPC, ...)



■ Regional Computer Centre for all Bavarian Universities

■ Computer Centre for all Munich Universities



- National Supercomputing Centre
- Regional Computer Centre for all Bavarian Universities
- Computer Centre for all Munich Universities

- Combination of the 3 German national supercomputing centers:
 - John von Neumann Institute for Computing (NIC), Jülich
 - High Performance Computing Center Stuttgart (HLRS)
 - Leibniz Supercomputing Centre (LRZ), Garching n. Munich
- Founded on 13. April 2007
- Hosting member of PRACE
(Partnership for Advanced Computing in Europe)

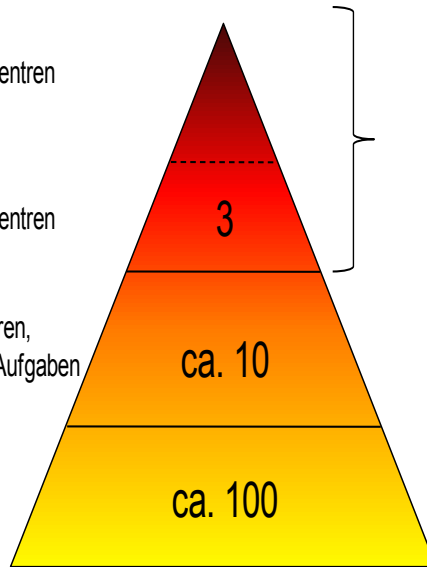


Europäische
Höchstleistungsrechenzentren
(Tier 0)

Nationale
Höchstleistungsrechenzentren
(Tier 1)

Thematische HPC-Zentren,
Zentren mit regionalen Aufgaben
(Tier 2)

HPC-Server
(Tier 3)



Gauss Centre for
Supercomputing (GCS)
(Garching, Stuttgart, Jülich)

Aachen, Berlin, DKRZ, Dresden,
DWD, Karlsruhe, Hannover,
MPG/RZG, udgl.

Hochschule/Institut



■ Establishment of the legal framework

- PRACE AISBL created with seat in Brussels in April (Association Internationale Sans But Lucratif)
- 20 members representing 20 European countries
- Inauguration in Barcelona on June 9



■ Funding secured for 2010 - 2015

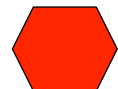
- 400 Million € from France, Germany, Italy, Spain Provided as Tier-0 services on TCO basis
- Funding decision for 100 Million € in The Netherlands expected soon
- 70+ Million € from EC FP7 for preparatory and implementation Grants INFSO-RI-211528 and 261557 Complemented by ~ 60 Million € from PRACE members



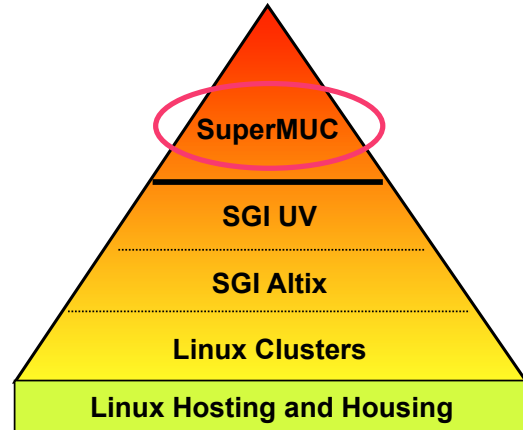
- **Curie @ GENCI:**
Bull Cluster, 1.7 PFlop/s
- **FERMI @ CINECA:**
IBM BG/Q, 2.1 PFlop/s
- **Hermit @ HLRS:**
Cray XE6, 1 Pflop/s
- **JUQUEEN @ FZJ:**
IBM Blue Gene/Q, 5.9 PFlop/s
- **MareNostrum @ BSC:**
IBM System X iDataPlex, 1 PFlop/s
- **SuperMUC @ LRZ:**
IBM System X iDataPlex, 3.2 PFlop/s



- Single pan-European Peer Review
- <http://www.prace-project.eu/Call-Announcements?lang=en>
- Early Access Call in May 2010
 - 68 proposals asked for 1870 Million Core hours
 - 10 projects granted with 328 Million Core hours
 - Principal Investigators from D (5), UK (2) NL (1), I (1), PT (1)
 - Involves researchers from 31 institutions in 12 countries
- Further calls being scheduled (every 6 months)
 - Call open February > Access September same year
 - Call open September > Access March next year
- Example from 8th Regular Call closed on 15 October 2013
 - Spatially adaptive radiation-hydrodynamical simulations of reionization
 - Project leader: Dr Andreas Pawlik, Max Planck Society, GERMANY
 - Research field: Universe Sciences
 - Resource Awarded: 33,800,000 core hours on SuperMUC, Germany;



- European Supercomputing Centre
- National Supercomputing Centre
- Regional Computer Centre for all Bavarian Universities
- Computer Centre for all Munich Universities

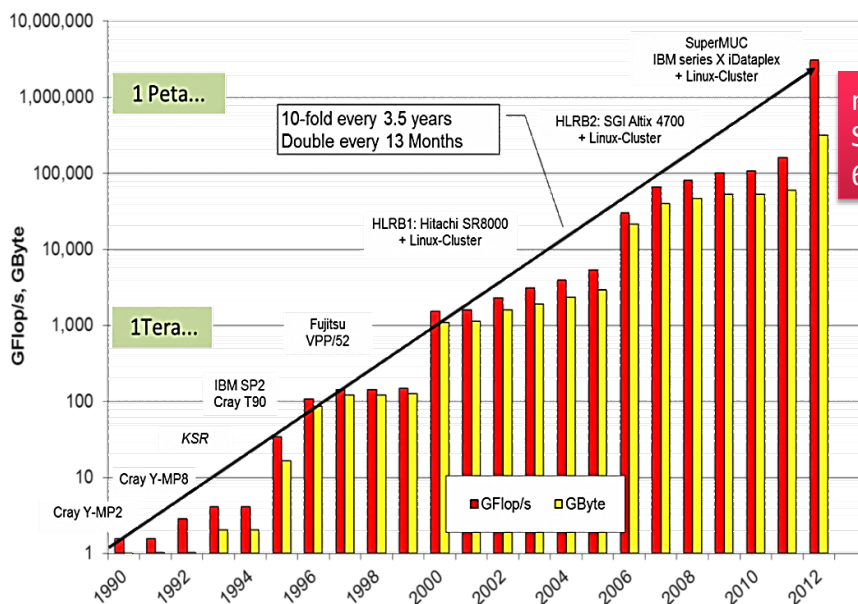


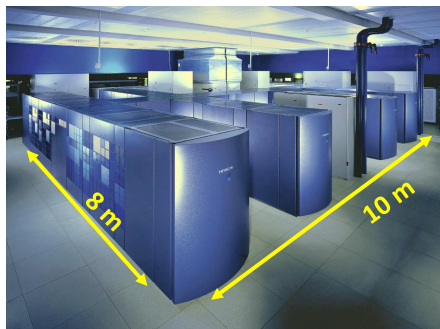
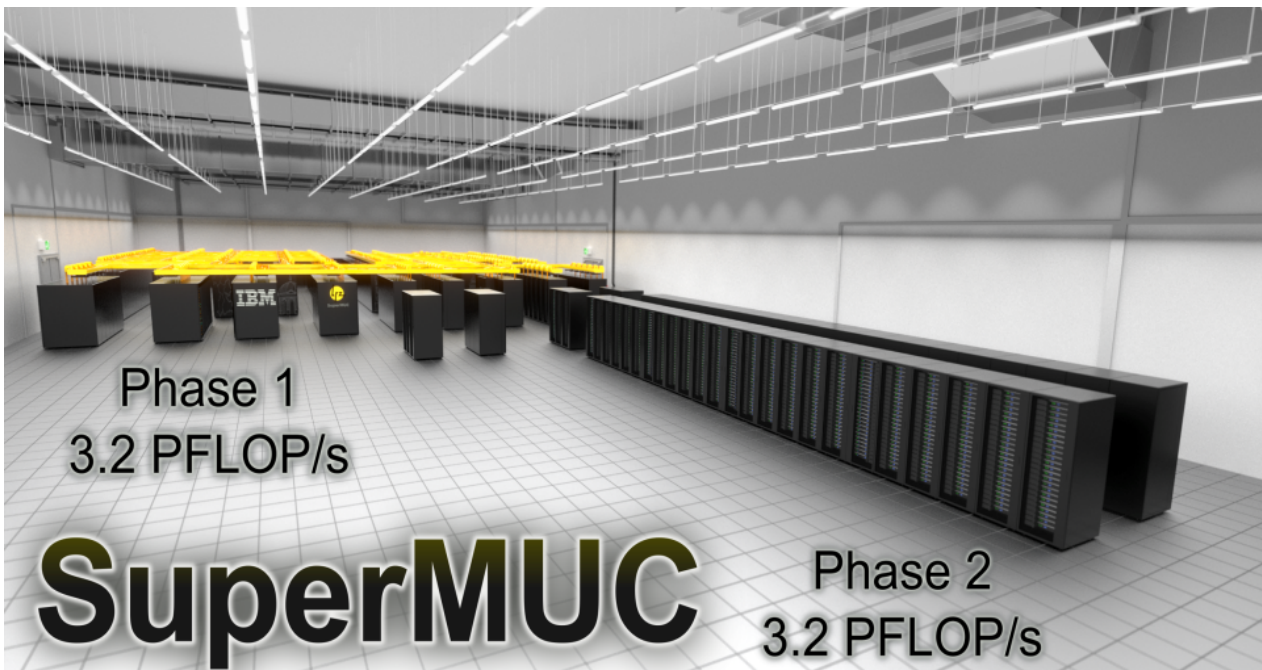
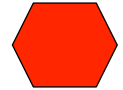
Video: SuperMUC rendered on SuperMUC by LRZ

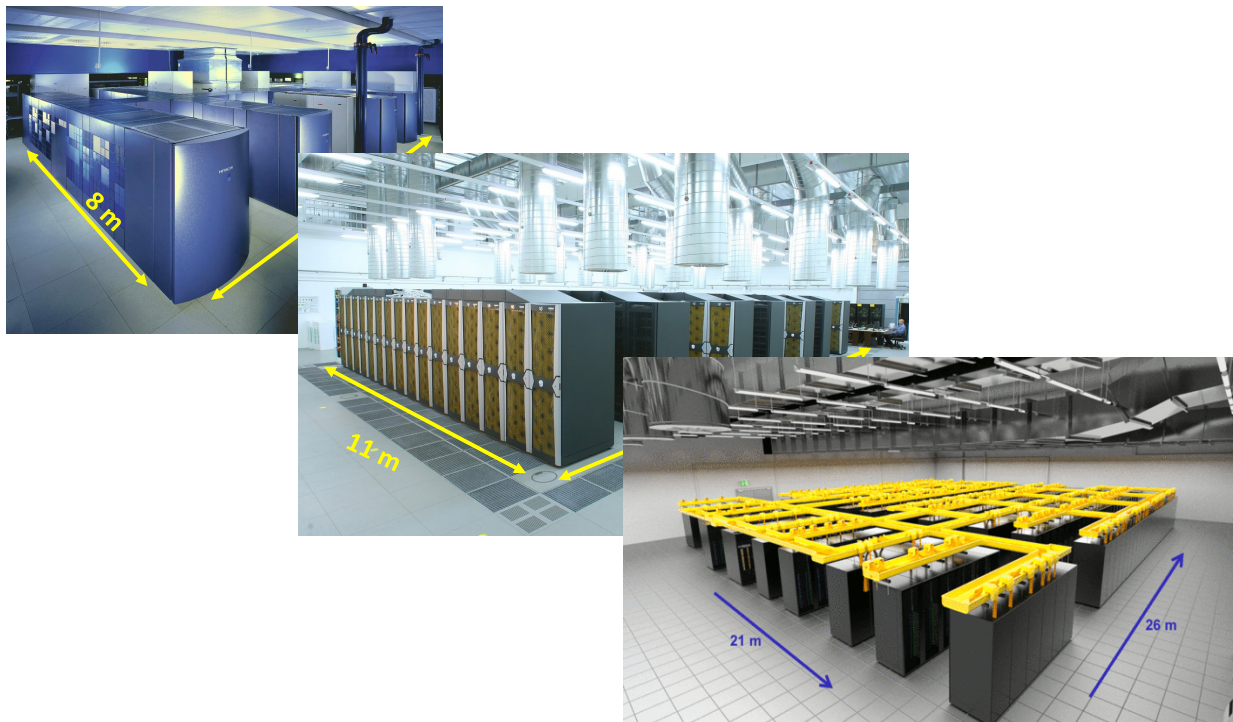
<http://youtu.be/OIAS6iiqWrQ>

| Rank | Site | Computer/Year Vendor | Cores | R _{max} | R _{peak} | Power |
|------|---|--|---------|------------------|-------------------|---------|
| 1 | DOE/NNSA/LLNL United States | Sequoia - BlueGene/Q, Power BQC 16C 1.60 GHz, Custom / 2011 IBM | 1572864 | 16324.75 | 20132.66 | 7890.0 |
| 2 | RIKEN Advanced Institute for Computational Science (AICS) Japan | K computer, SPARC64 VIIIfx 2.0GHz, Tofu interconnect / 2011 Fujitsu | 705024 | 10510.00 | 11280.38 | 12659.9 |
| 3 | DOE/SC/Argonne National Laboratory United States | Mira - BlueGene/Q, Power BQC 16C 1.60GHz, Custom / 2012 IBM | 786432 | 8162.38 | 10066.33 | 3945.0 |
| 4 | Leibniz Rechenzentrum Germany | SuperMUC - iDataPlex DX360M4, Xeon E5-2680 8C 2.70GHz, Infiniband FDR / 2012 IBM | 147456 | 2897.00 | 3185.05 | 3422.7 |
| 5 | National Supercomputing Center in Tianjin China | Tianhe-1A - NUDT YH MPP, Xeon X5670 6C 2.93 GHz, NVIDIA 2050 / 2010 NUDT | 186368 | 2566.00 | 4701.00 | 4040.0 |
| 6 | DOE/SC/Oak Ridge National Laboratory United States | Jaguar - Cray XK6, Opteron 6274 16C 2.200GHz, Cray Gemini interconnect, NVIDIA 2090 / 2009 Cray Inc. | 298592 | 1941.00 | 2627.61 | 5142.0 |
| 7 | CINECA Italy | Fermi - BlueGene/Q, Power BQC 16C 1.60GHz, Custom / 2012 IBM | 163840 | 1725.49 | 2097.15 | 821.9 |
| 8 | Forschungszentrum Juelich (FZJ) Germany | JuQUEEN - BlueGene/Q, Power BQC 16C 1.60GHz, Custom / 2012 IBM | 131072 | 1380.39 | 1677.72 | 657.5 |
| 9 | CEA/TGCC-GENCI France | Curie thin nodes - Bullx B510, Xeon E5-2680 8C 2.700GHz, Infiniband QDR / 2012 Bull | 77184 | 1359.00 | 1667.17 | 2251.0 |
| 10 | National Supercomputing Centre in Shenzhen (NSCS) China | Nebulae - Dawning TC3600 Blade System, Xeon X5650 6C 2.66GHz, Infiniband QDR, NVIDIA 2050 / 2010 Dawning | 120640 | 1271.00 | 2984.30 | 2580.0 |

www.top500.org







Picture: Horst-Dieter Steinhöfer

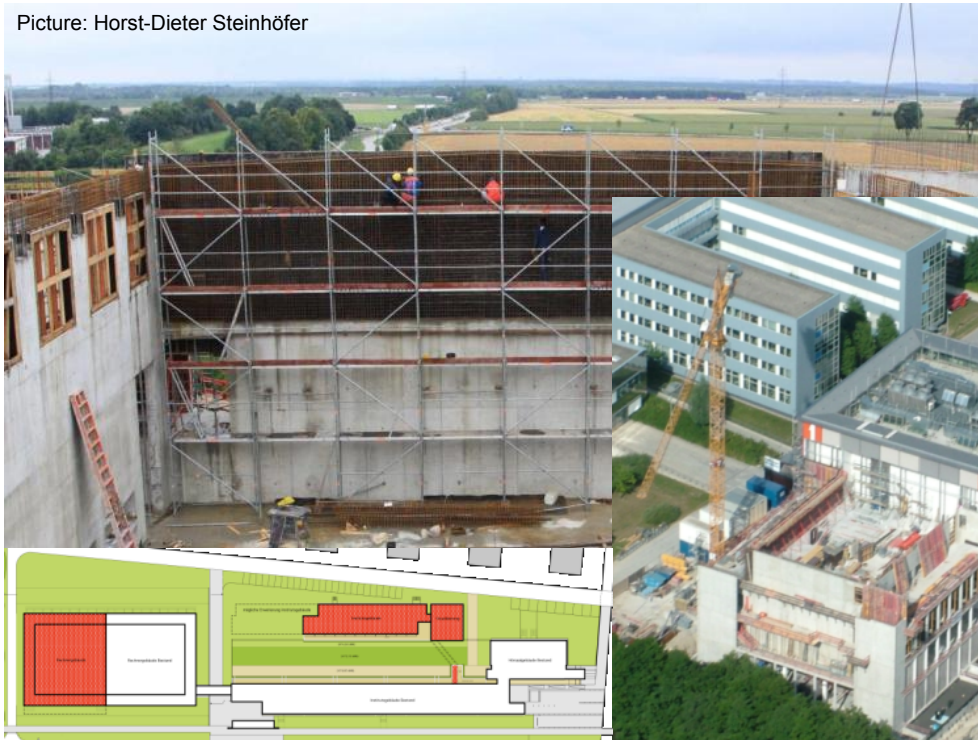
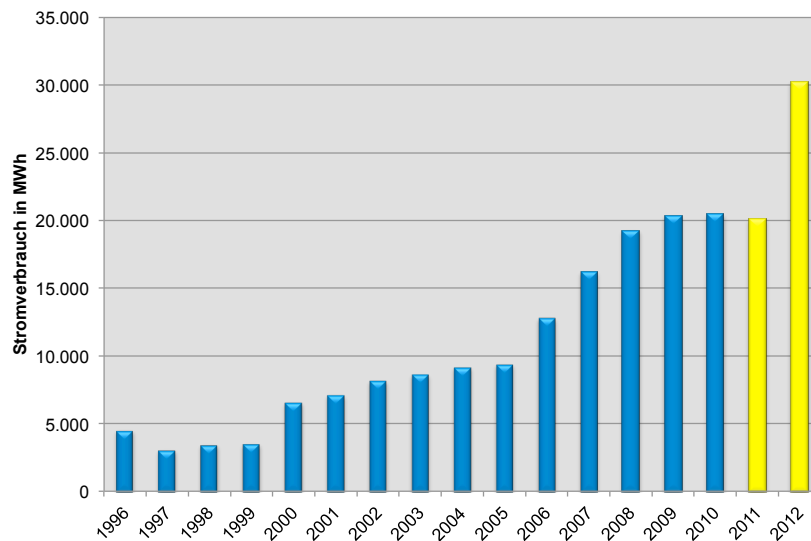
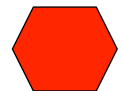
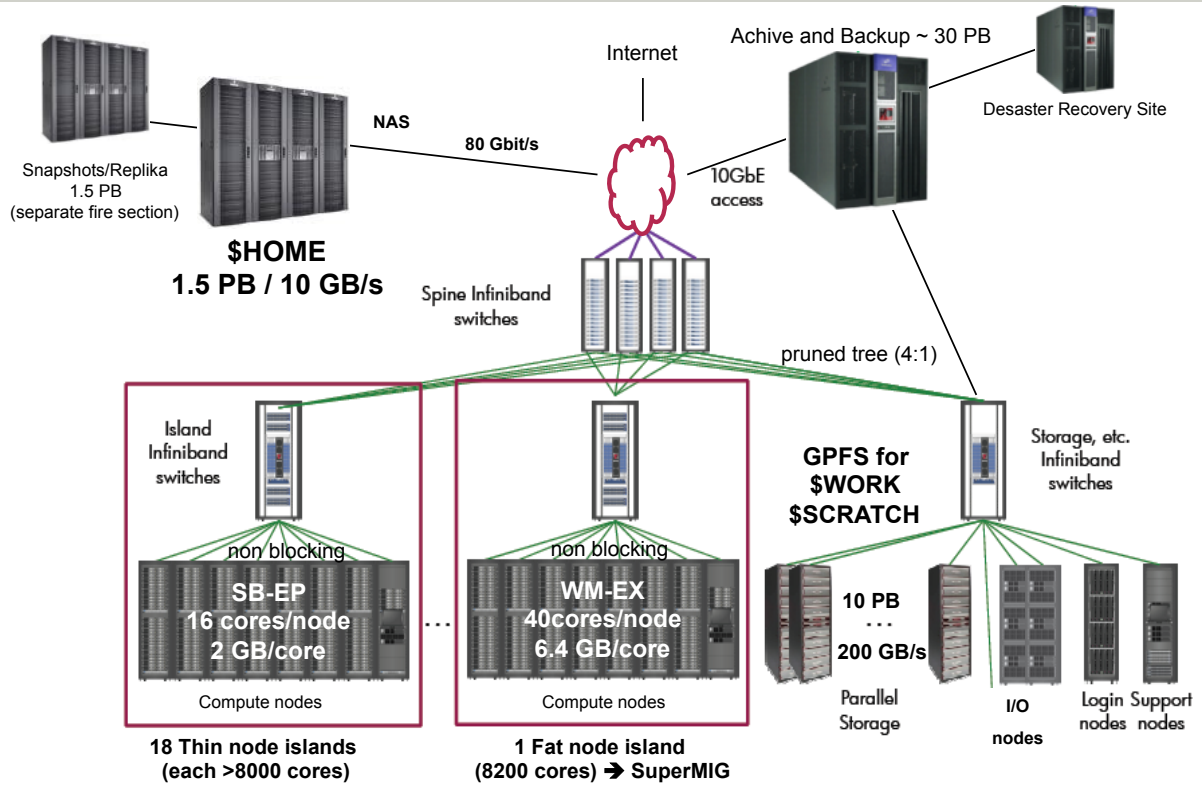
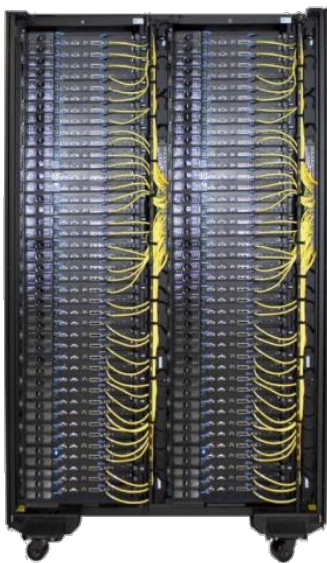


Figure: Herzog+Partner für StBAM2 (staatl. Hochbauamt München 2)

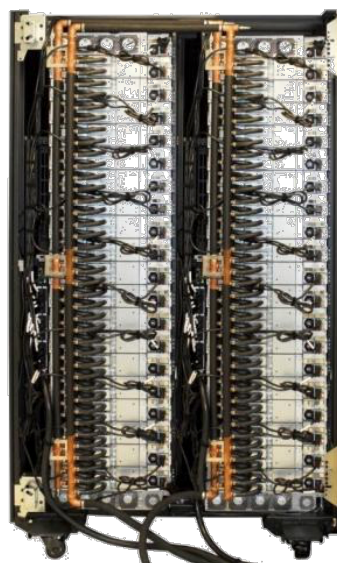
Picture: Ernst A. Graf

| Date | System | Flop/s | Cores |
|------|-------------------|-------------------|--------|
| 2000 | HLRB-I | 2 Tflop/s | 1512 |
| 2006 | HLRB-II | 62 Tflop/s | 9728 |
| 2012 | SuperMUC | 3200 Tflop/s | 155656 |
| 2014 | SuperMUC Phase II | 3.2 + 3.2 Pflop/s | 229960 |



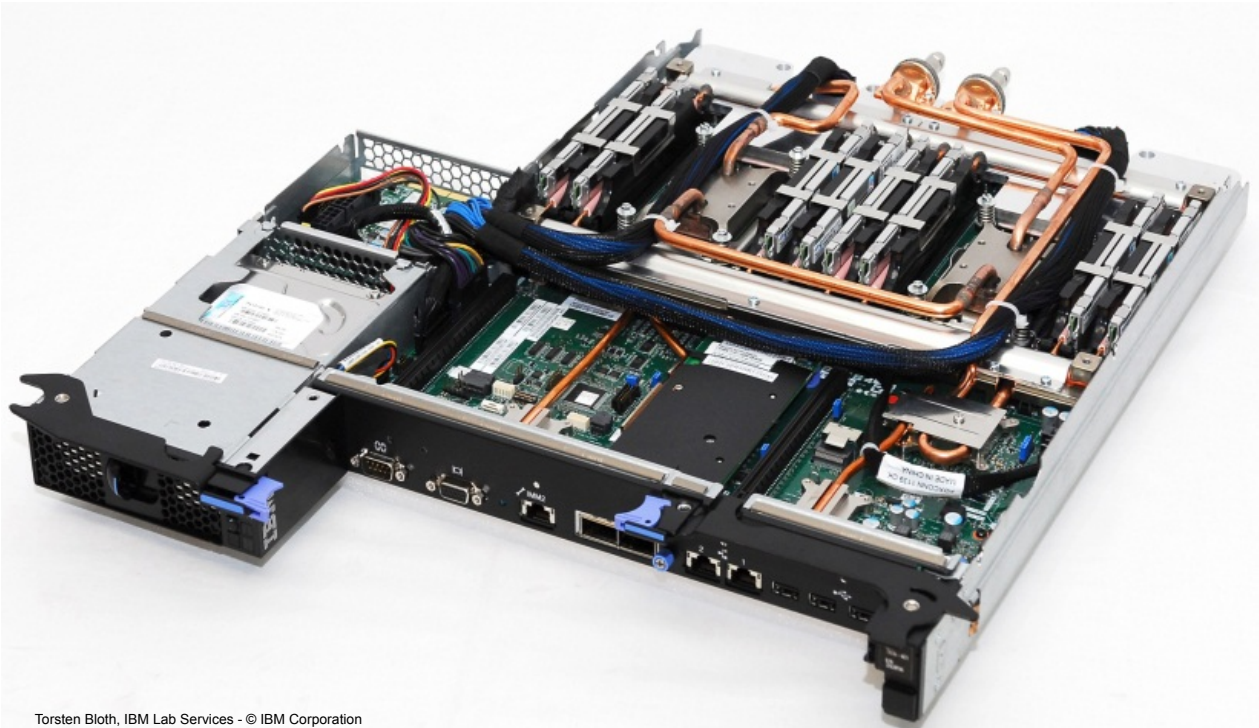


iDataPlex DWC Rack
w/ water cooled nodes
(front view)



iDataPlex DWC Rack
w/ water cooled nodes
(rear view of water manifolds)

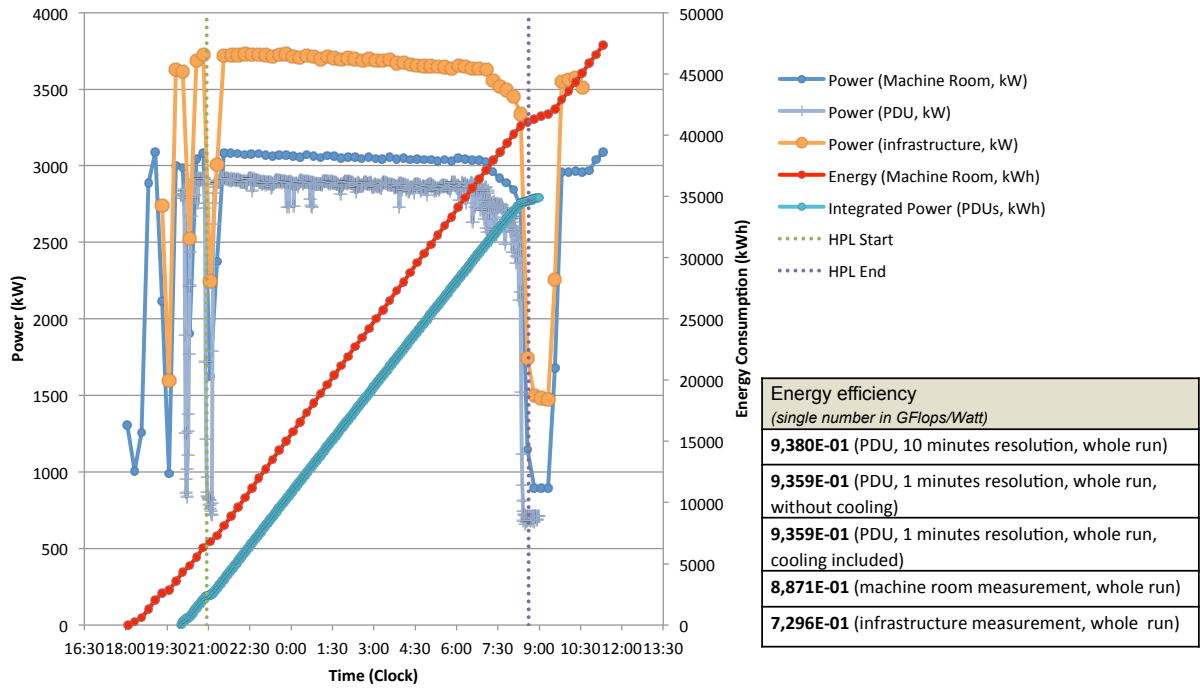
Torsten Bloth, IBM Lab Services - © IBM Corporation



Torsten Bloth, IBM Lab Services - © IBM Corporation



Photos: StBAM2 (staatl. Hochbauamt München 2)





- Computational Fluid Dynamics: Optimisation of turbines and wings, noise reduction, air conditioning in trains**
- Fusion: Plasma in a future fusion reactor (ITER)**
- Astrophysics: Origin and evolution of stars and galaxies**
- Solid State Physics: Superconductivity, surface properties**
- Geophysics: Earth quake scenarios**
- Material Science: Semiconductors**
- Chemistry: Catalytic reactions**
- Medicine and Medical Engineering: Blood flow, aneurysms, air conditioning of operating theatres**
- Biophysics: Properties of viruses, genome analysis**
- Climate research: Currents in oceans**

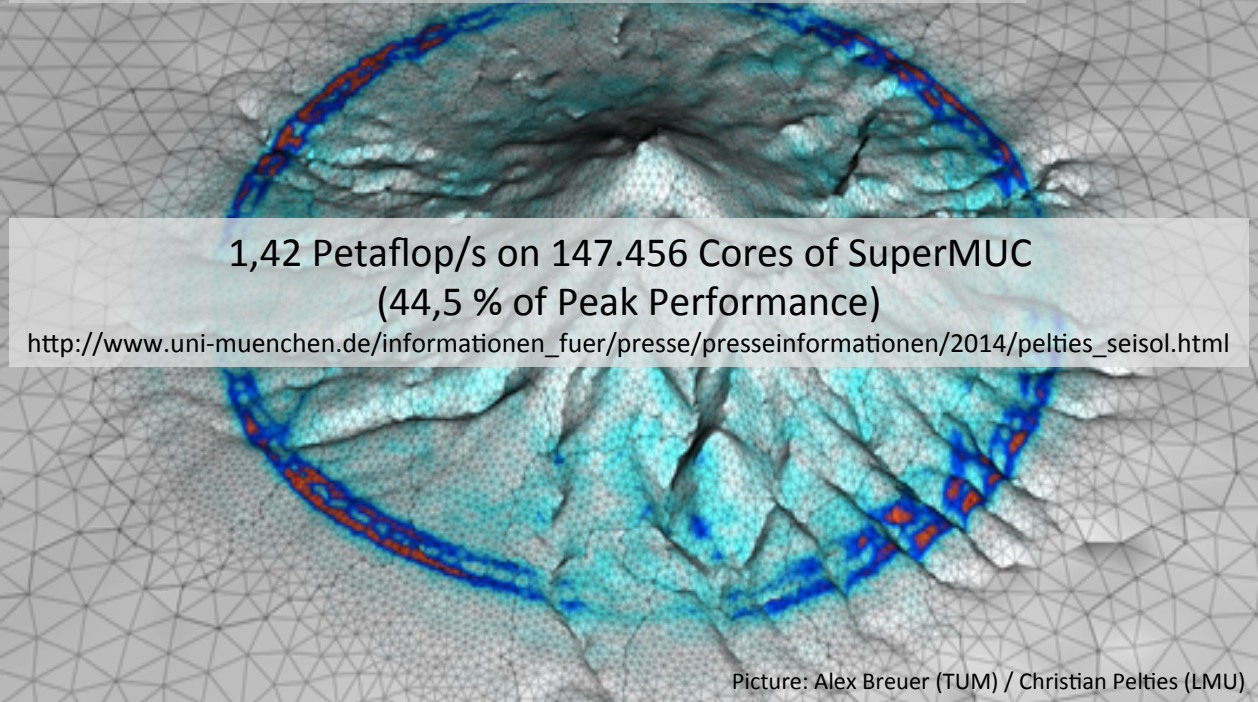



Results (Sustained TFlop/s on 128000 cores)



| Name | MPI | # cores | Description | TFlop/s/island | TFlop/s max |
|----------|------------|----------|---------------------|----------------|-------------|
| Linpac | IBM | ★ 128000 | TOP500 | 161 | 2560 |
| Vertex | IBM | ★ 128000 | Plasma Physics | 15 | 245 |
| GROMACS | IBM, Intel | ★ 64000 | Molecular Modelling | 40 | 110 |
| Seissol | IBM | ★ 64000 | Geophysics | 31 | 95 |
| waLBerla | IBM | ★ 128000 | Lattice Boltzmann | 5.6 | 90 |
| LAMMPS | IBM | ★ 128000 | Molecular Modelling | 5.6 | 90 |
| APES | IBM | ★ 64000 | CFD | 6 | 47 |
| BQCD | Intel | ★ 128000 | Quantum Physics | 10 | 27 |

Dr. Christian Pelties, Department of Earth and Environmental Sciences (LMU)
Prof. Michael Bader, Department of Informatics (TUM)



1,42 Petaflop/s on 147.456 Cores of SuperMUC
(44,5 % of Peak Performance)

http://www.uni-muenchen.de/informationen_fuer/presse/presseinformationen/2014/pelties_seisol.html

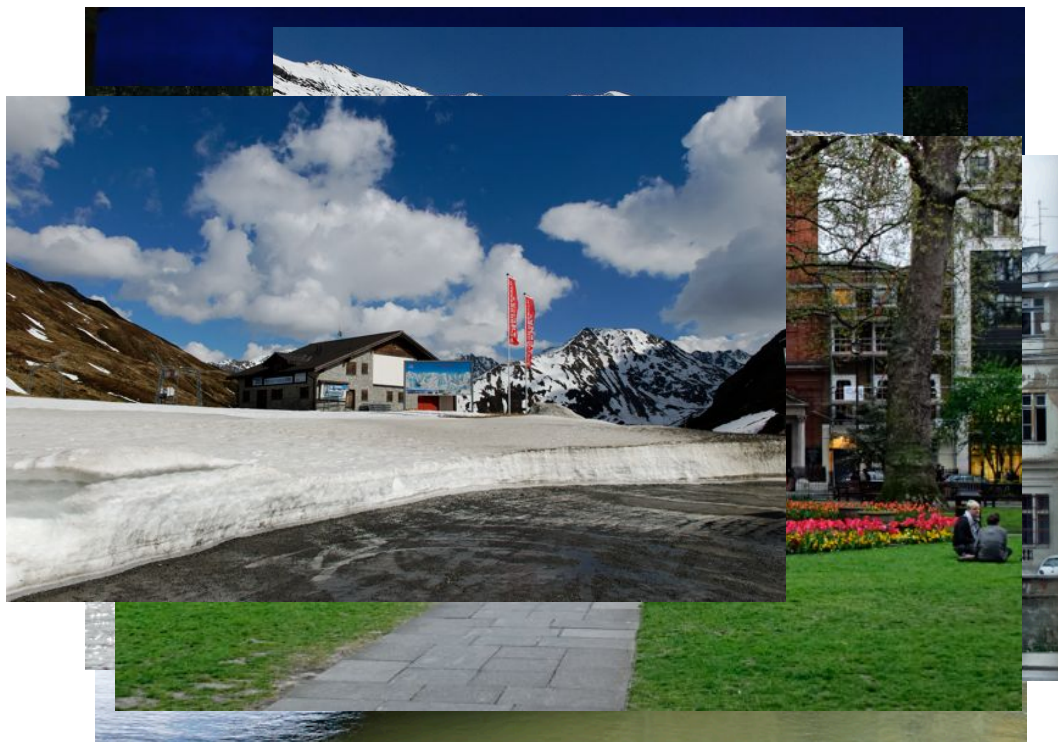
Picture: Alex Breuer (TUM) / Christian Pelties (LMU)

What are the Costs of such an e-Infrastructure?

- Contract n°: RI-283449 (CSA-SA)
- Start date: 01/08/2011
- Duration: 18 months (end 31/1/2013)
- Total budget: 392.523 €
- Total funded effort in PMs: 33.75
- Partners: AUEB-RC, EGI.eu, NUI Galway, ETL
- Web site: www.efiscal.eu

- Observations/Assumptions:
 - Cloud, commodisation catalysts for cost assessment
 - Inevitable with growing scale – also without catalysts
 - Cloud/dedicated cost ratios (literature)
 - Measuring costs is easy (Track spending)
 - Comparison is easy (Death of the distance/location, virtualization)
 - No major surprises tomorrow (Mature technologies/business models)

Slide Courtesy Matti Heikkurinen



Slide © Matti Heikkurinen

- Tracking costs is not easy
 - Absence of tracking technologies
 - Full cost model often not used in academia/research
 - Past or future
- Comparison is not trivial
 - HPC, HTC, HPC Cloud, HTC Cloud
 - Different kinds of services
- Future is uncertain
 - Plummeting prices as a rule, but
 - Flood in Thailand -> HD prices double
 - New upcoming technologies
 - Energy costs, green regulations
 - Changes in user behavior

Slide Courtesy Matti Heikkurinen

- | | |
|---|--|
| <ul style="list-style-type: none"> ■ Categories <ul style="list-style-type: none"> – Infrastructure <ul style="list-style-type: none"> • Building/Housing • Cooling – Machines <ul style="list-style-type: none"> • SuperMUC • Other Systems (HTC, ...) • Federation of Systems – Operating <ul style="list-style-type: none"> • Power/Electricity • Others – Personnel <ul style="list-style-type: none"> • Fixed • Project-based | <ul style="list-style-type: none"> ■ Sources <ul style="list-style-type: none"> – Users – Local/Universities – Regional/State – Federal/National – European/International – Industry (Prototypes) ■ Intervals <ul style="list-style-type: none"> – Pay-per-use – Regular (annual budget) – Irregular (project budget) |
|---|--|

- You can know the budget
 - New areas, hype
 - Strategic investment from governments, EC

- You can know the resources needed
 - Commodity solutions, commonly agreed “pain points”
 - Constant budget pressure (per user)

- You can't know both at the same time
 - New capability: unknown volume of use
 - Commodity: market disruptions
 - Commoditisation process tends to be painful: technology vs. governance, developing new metrics, political engagement...

Slide © Matti Heikkurinen

- What is the Total-Cost-of-Ownership (TCO)?

- Difference between IT-Services for commodity and research

- Inherent complexity of cost categories and funding sources

- Cost only one of the factor influencing the selection of system
 - Convenience, flexibility
 - Intangible factors
 - Non-standard requirements
 - Risk management

- Professionalizing service delivery: FedSM/FitSM (<http://fedsm.eu>)

- **Attract, develop, and retain talent**

Extreme Scale Computing: Challenges and Costs

Dieter Kranzlmüller
kranzlmueLLer@lrz.de

