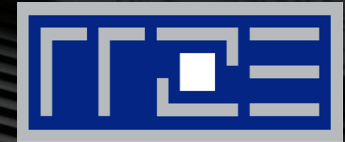


REGIONALES RECHENZENTRUM ERLANGEN [RRZE]



High Performance Computing

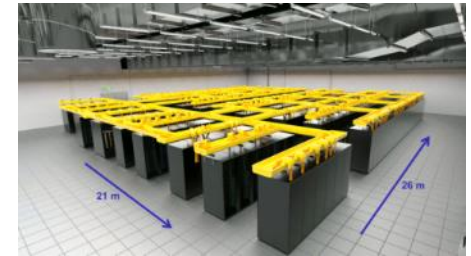
Systemausbildung – Grundlagen und Aspekte von
Betriebssystemen und systemnahen Diensten, 21.06.2017
HPC-Gruppe, RRZE

Agenda

- Was bedeutet HPC?
- Aus welchen Komponenten besteht ein typisches HPC-System?
- Mit dem HPC-System läuft
\$GTA5/\$Photoshop/\$Videoschnittprogramm doch
superschnell, oder?
- Shared-Memory-Parallelisierung vs. Message Passing
- Was ist dieses Queuing-System und wozu ist es gut?

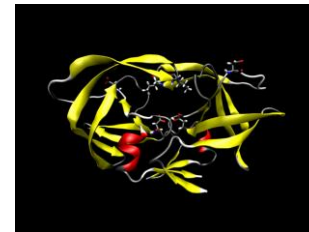
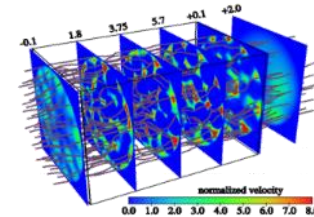
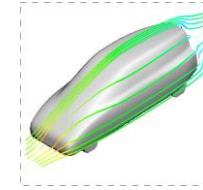
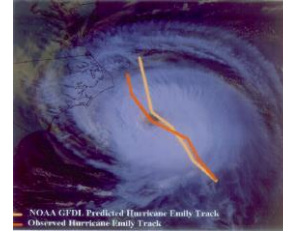
Was bedeutet HPC?

- Englisch: High Performance Computing
- Deutsch: **Hochleistungsrechnen**
- “Hochleistungsrechnen ist ein Bereich des computergestützten Rechnens. Er umfasst alle Rechenarbeiten, deren Bearbeitung einer hohen Rechenleistung oder Speicherkapazität bedarf.” (Wikipedia)
- genaue Abgrenzung ist schwierig



Einsatzbereiche von HPC-Systemen (Beispiele)

- Ingenieurwesen (Automotive)
 - Crashsimulation
 - Aerodynamik
 - Akustik
- Meteorologie
 - Wettervorhersage
 - Katastrophenvorhersage
- Biologie/Medizin
 - “Drug Design”
 - Aufklärung von Reaktionsprozessen



Einsatzbereiche von HPC-Systemen (Beispiele)

- Werkstoffwissenschaften
- Physik
 - Aufklärung von fundamentalen Wechselwirkungen
 - Struktur der Materie
- Filmindustrie
 - Rendering
 - CGI
- Zunehmend auch außerhalb der klassischen Natur-/Ingenieurwissenschaften, z.B. in Wirtschafts- oder Sprachwissenschaften



Ein Wort zur Warnung



Fricke / Lehmann: Der Rechenschieber
VEB Fachbuchverlag Leipzig 1961
7. Auflage

Man kann sich aber auch zu sehr an den Gebrauch des Rechenschiebers gewöhnen. Dann werden auch die einfachsten Berechnungen, die sonst spielend, fast ohne nachzudenken, im Kopf erledigt werden, unbedingt mit dem „Schieber“ erledigt. Bei diesen Leuten ist $3 \text{ mal } 3$ auf dem Schieber etwa 8,99. Es erübrigt sich wohl, erst zu sagen, daß diese Methode nicht ganz die richtige ist. Grundsätzlich soll der Rechenschieber das Kopfrechnen *nicht ersetzen*; denn dieses erhält den Geist beweglich und ist außerdem beim Gebrauch des Rechenschiebers unbedingt erforderlich.

s/Rechenschieber/Supercomputer/g

s/Kopfrechnen/Common Sense/g

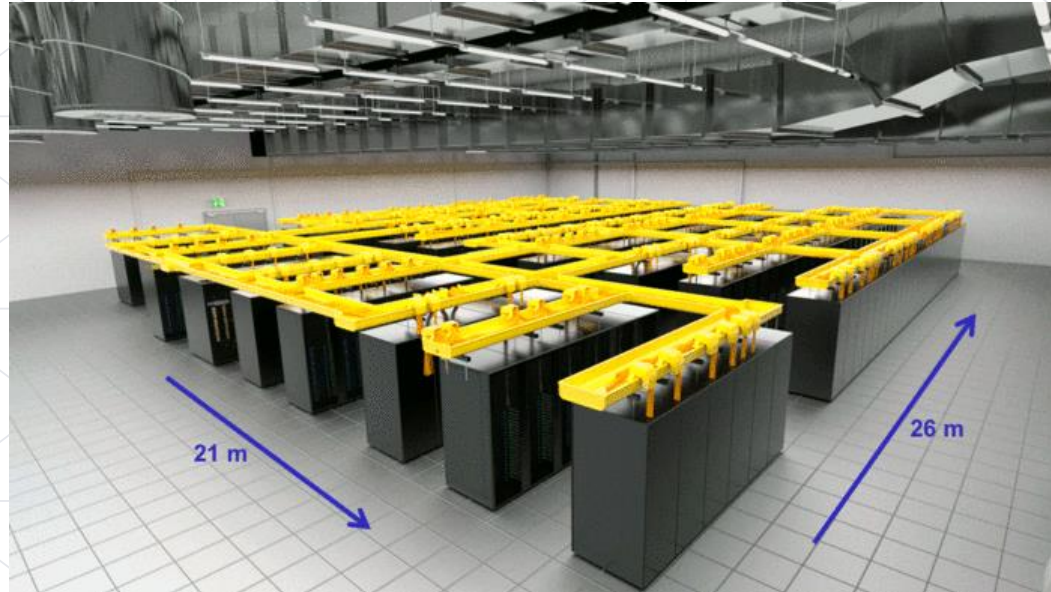
HPC-Systeme

- “Hochleistungsrechner sind Rechnersysteme, die geeignet sind, Aufgaben des Hochleistungsrechnens zu bearbeiten.”
- Früher: Spezialsysteme – Spezialhardware + -software
- **Heute: Fast ausschließlich Cluster aus Standard-Hardware,** meist mit schnellem Interconnect (Netzwerk)
 - Die einzelnen Clusterknoten sind nicht schneller als ein gut ausgestatteter Standardserver!
- Vektorrechner: Nischenprodukt

Beispiel: SuperMUC Phase 1 (2012) @ LRZ

- **9216 compute nodes**
 - 2x Intel Xeon E5-2680 CPU (8 Kerne, 2.7 GHz, „Sandy Bridge“)
 - 32 GB RAM
 - Ergibt in Summe 147456 Kerne und 288 TB RAM
- **Infiniband-Netz**
 - 18 Inseln a 512 Knoten
 - › Fully non-blocking innerhalb der Inseln
 - › ausgedünnt 4:1 zwischen den Inseln
 - FDR-10 (knapp 40 Gbit/s)
- **15 PB paralleles Filesystem @250 GByte/s**

“SuperMUC” Phase 1 (rendering)



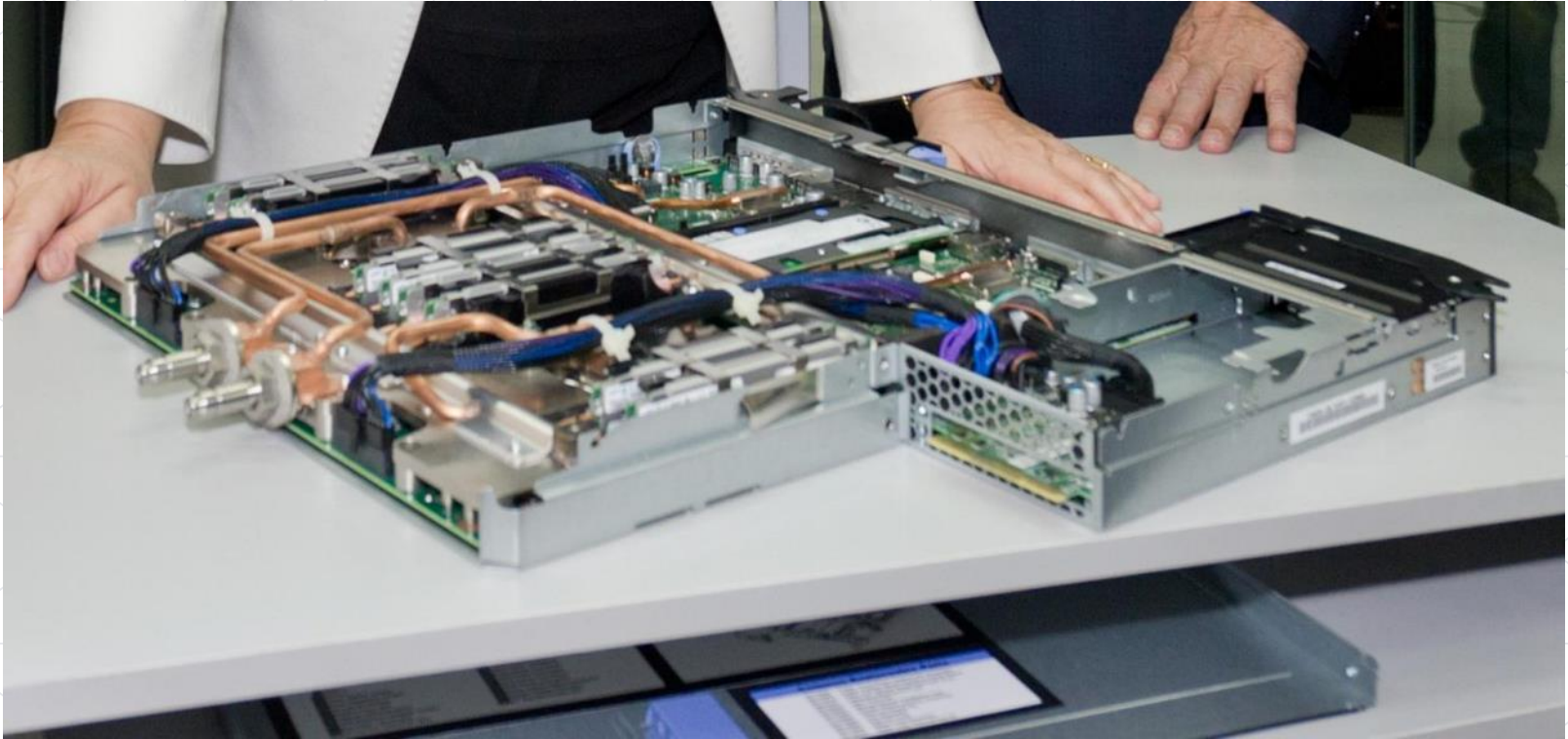
Quelle: www.lrz.de

SuperMUC 2012 (realer Kabelsalat)



Quelle: <https://en.wikipedia.org/wiki/File:SuperMUC.jpg>, Wikipedia-User Mdw77, CC-BY-SA 4.0

SuperMUC (Rechenknoten mit Wasserkühlung)



Quelle: www.lrz.de

LiMa Rechenknoten

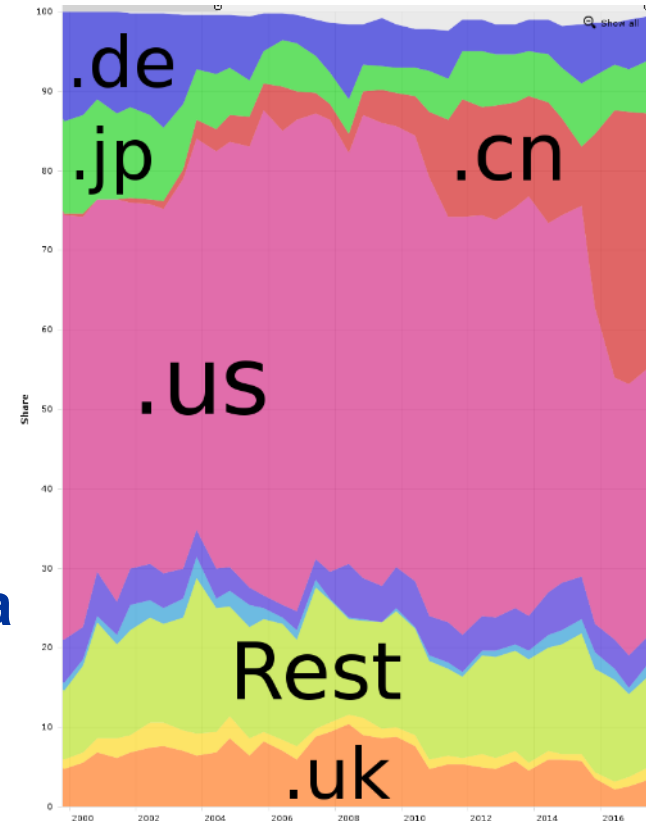


Mit dem HPC-System läuft \$GTA5 doch super, oder?

- HPC-Systeme führen nicht auf magische Weise beliebige Software schneller aus
- Software muss explizit dafür programmiert werden, dass sie das HPC-System nutzen kann → **Parallelisierung!**
- → Selbst wenn das HPC-Cluster unter Windows laufen würde, könnte \$GTA5 nicht mehr als einen Knoten eines Clusters benutzen.

Top500 <http://www.top500.org>

- Liste der **500 schnellsten HPC-Systeme**
 - ...die öffentlich bekannt sein sollen
 - Gemessen mit “**Linpack**”-Benchmark (lösen eines großen linearen Systems)
 - Erscheint 2x pro Jahr (Nov. + Juni)
- **#1 Juni 2017: Sunway TaihuLight, China**
 - 10 649 600 cores, chinesische CPU
- **#2: Tianhe-2, China**
 - 3 120 000 Intel Xeon Phi cores



Trends Top500-Liste Juni 2017

- 43% mit Infiniband- oder OmniPath-Vernetzung
- **>99% mit Linux** (Windows: zuletzt 07/2015 1 System)
- 33% der Systeme in den USA, 32% in China (.de: 5%);
 - Top 10: 2 .cn, 5 USA, 2 .jp, 1 .ch
 - Zum Vergleich 11/2014: 46% USA, 12% .cn.
Top10: 1 .cn, 6 USA, 1 .jp, 1 .ch, 1 .de
- **90% mit Intel Xeon CPUs** (11/2014: 85%)

Infiniband / Omni-Path

- **De-facto Standard** für HPC-Cluster
 - Außer Ethernet praktisch keine Konkurrenz mehr
- **Schnell** und mit sehr geringer Latenz
 - **Switched** fabric, oft fully non-blocking (bis 648 Ports)
 - QDR: 40 Gbit roh, 32 Gbit Nutz, 1.3 μ s Latenz (seit 2007)
 - FDR: 56 Gbit roh, 54 Gbit Nutz, 0.7 μ s Latenz (seit 2011)
 - Omni-Path: im Prinzip Infiniband @100 GBit, nur Intel-proprietär und mit Infiniband inkompatibel :-/

Infiniband - Probleme

- Grauenhafte (Nicht-)API
 - Vom Endanwender unbemerkt
- Geringe Zuverlässigkeit
- Besorgniserregende Marktkonsolidierung
 - Nur noch zwei Hersteller für QDR, nur einer für FDR
 - Stark verlangsamte Entwicklung, steigende Preise
 - Wird durch Omni-Path nicht besser, im Gegenteil...

Ein Teil des Infiniband-Netzes von LiMa...



Acceleratoren

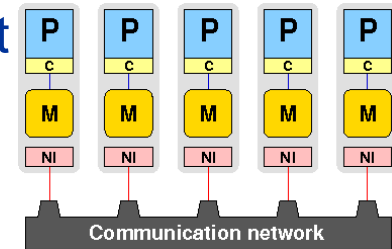
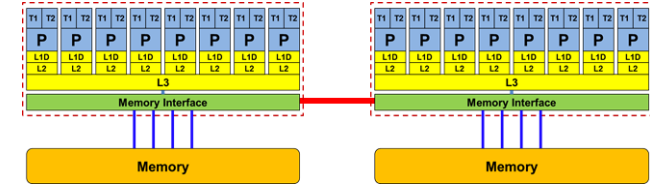
- Einsteckkarten, auf die Berechnungen ausgelagert werden können
 - GPUs
 - Intel Xeon Phi
- Noch schwieriger effizient zu programmieren als normale shared-memory-Systeme oder Cluster.

Programmiermodelle

- NVIDIA CUDA, OpenCL, OpenACC (GPUs)
- OpenMP (Intel Xeon Phi)
- **Wunder darf man auch im Idealfall nicht erwarten!**

Parallelisierungsarten

- Shared Memory
 - **OpenMP**: Compiler-Erweiterung
 - Pthreads: POSIX library
 - TBB, OmpSs, Cilk+,...
 - Automatische Parallelisierung: nicht praxisrelevant
- Distributed Memory
 - **MPI** (Message Passing Interface)
 - › Version 1.0: Juni 1994; Aktuelle Version: MPI 3.1 (06/2015)
 - Sprachbasierte Parallelität
 - X10, Chapel, ...



OpenMP

- Spracherweiterungen, die von einem OpenMP-fähigen Compiler ausgewertet und umgesetzt werden
- Trivial-Beispiel aus Wikipedia (in C):

```
int i, a[100000];
```

```
#pragma omp parallel for  
for (i = 0; i < N; i++)  
    a[i] = 2 * i;
```

- Jeder Thread bekommt einen Teil der Schleife zugewiesen und arbeitet ihn ab (Shared Memory!)

MPI

- Alt aber konstant weiterentwickelt
- **Library-basiert** (Funktionsaufrufe)
 - Programmiersprache ist weiterhin seriell
- Praktisch jede Software die parallel auf mehr als einem Rechner laufen kann setzt darauf auf
- Alle hochskalierbaren Produktionscodes auf Supercomputern sind MPI-basiert.

MPI

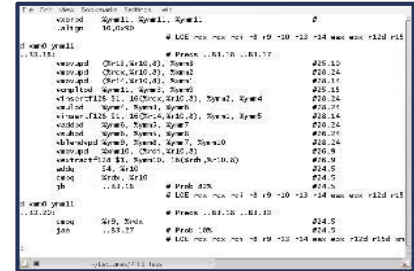
- MPI-Standard definiert ein **Set von Funktionen** für die Programmierung von Parallelrechnern
 - Senden / Empfangen von Nachrichten
 - Parallel I/O
 - Synchronisation
 - › warten bis alle Prozesse eine Barriere erreicht haben

“Hello World” in MPI (C)

```
int main(int argc, char ** argv) {  
    int numprocs, myid;  
    MPI_Init(&argc, &argv);  
    MPI_Comm_size(MPI_COMM_WORLD, &numprocs);  
    MPI_Comm_rank(MPI_COMM_WORLD, &myid);  
    printf("Hi, ich bin Prozess %d von %d\n",  
           myid, numprocs);  
    MPI_Finalize();  
    return 0;  
}
```

Nutzung von HPC-Clustern in der Praxis: Batchbetrieb

- Fast wie in der guten alten Zeit: Lochkarten mit Programmen darauf abgeben, auf Verarbeitung warten, Ergebnis abholen
 - Nur heutzutage natürlich ohne Lochkarten...
- Man kann sich nicht einfach direkt auf beliebigen Knoten des Clusters einloggen
- Stattdessen: Einloggen auf **Frontend**, Jobs submittieren
 - Frontends sind interaktiv zugänglich (meist SSH)
 - Dort können Jobs vorbereitet und abgeschickt werden
 - Jobs müssen komplett ohne Benutzereingaben auskommen!



```
slurm@frontend:~$ squeue -o '%10i,%10j,%10k,%10l,%10m,%10n,%10o,%10p,%10q,%10r,%10s,%10t,%10u,%10v,%10w,%10x,%10y,%10z,%10aa,%10ab,%10ac,%10ad,%10ae,%10af,%10ag,%10ah,%10ai,%10aj,%10ak,%10al,%10am,%10an,%10ao,%10ap,%10aq,%10ar,%10as,%10at,%10au,%10av,%10aw,%10ax,%10ay,%10az,%10ba,%10bb,%10bc,%10bd,%10be,%10bf,%10bg,%10bh,%10bi,%10bj,%10bk,%10bl,%10bm,%10bn,%10bo,%10bp,%10bq,%10br,%10bs,%10bt,%10bu,%10bv,%10bw,%10bx,%10by,%10bz,%10ca,%10cb,%10cc,%10cd,%10ce,%10cf,%10cg,%10ch,%10ci,%10cj,%10ck,%10cl,%10cm,%10cn,%10co,%10cp,%10cq,%10cr,%10cs,%10ct,%10cu,%10cv,%10cw,%10cx,%10cy,%10cz,%10da,%10db,%10dc,%10dd,%10de,%10df,%10dg,%10dh,%10di,%10dj,%10dk,%10dl,%10dm,%10dn,%10do,%10dp,%10dq,%10dr,%10ds,%10dt,%10du,%10dv,%10dw,%10dx,%10dy,%10dz,%10ea,%10eb,%10ec,%10ed,%10ee,%10ef,%10eg,%10eh,%10ei,%10ej,%10ek,%10el,%10em,%10en,%10eo,%10ep,%10eq,%10er,%10es,%10et,%10eu,%10ev,%10ew,%10ex,%10ey,%10ez,%10fa,%10fb,%10fc,%10fd,%10fe,%10ff,%10fg,%10fh,%10fi,%10fj,%10fk,%10fl,%10fm,%10fn,%10fo,%10fp,%10fq,%10fr,%10fs,%10ft,%10fu,%10fv,%10fw,%10fx,%10fy,%10fz,%10ga,%10gb,%10gc,%10gd,%10ge,%10gf,%10gg,%10gh,%10gi,%10gj,%10gk,%10gl,%10gm,%10gn,%10go,%10gp,%10gq,%10gr,%10gs,%10gt,%10gu,%10gv,%10gw,%10gx,%10gy,%10gz,%10ha,%10hb,%10hc,%10hd,%10he,%10hf,%10hg,%10hh,%10hi,%10hj,%10hk,%10hl,%10hm,%10hn,%10ho,%10hp,%10hq,%10hr,%10hs,%10ht,%10hu,%10hv,%10hw,%10hx,%10hy,%10hz,%10ia,%10ib,%10ic,%10id,%10ie,%10if,%10ig,%10ih,%10ii,%10ij,%10ik,%10il,%10im,%10in,%10io,%10ip,%10iq,%10ir,%10is,%10it,%10iu,%10iv,%10iw,%10ix,%10iy,%10iz,%10ja,%10jb,%10jc,%10jd,%10je,%10jf,%10jg,%10jh,%10ji,%10jj,%10jk,%10jl,%10jm,%10jn,%10jo,%10jp,%10jq,%10jr,%10js,%10jt,%10ju,%10jv,%10jw,%10jx,%10jy,%10jz,%10ka,%10kb,%10kc,%10kd,%10ke,%10kf,%10kg,%10kh,%10ki,%10kj,%10kk,%10kl,%10km,%10kn,%10ko,%10kp,%10kq,%10kr,%10ks,%10kt,%10ku,%10kv,%10kw,%10kx,%10ky,%10kz,%10la,%10lb,%10lc,%10ld,%10le,%10lf,%10lg,%10lh,%10li,%10lj,%10lk,%10ll,%10lm,%10ln,%10lo,%10lp,%10lq,%10lr,%10ls,%10lt,%10lu,%10lv,%10lw,%10lx,%10ly,%10lz,%10ma,%10mb,%10mc,%10md,%10me,%10mf,%10mg,%10mh,%10mi,%10mj,%10mk,%10ml,%10mm,%10mn,%10mo,%10mp,%10mq,%10mr,%10ms,%10mt,%10mu,%10mv,%10mw,%10mx,%10my,%10mz,%10na,%10nb,%10nc,%10nd,%10ne,%10nf,%10ng,%10nh,%10ni,%10nj,%10nk,%10nl,%10nm,%10nn,%10no,%10np,%10np,%10nq,%10nr,%10ns,%10nt,%10nu,%10nv,%10nw,%10nx,%10ny,%10nz,%10oa,%10ob,%10oc,%10od,%10oe,%10of,%10of,%10og,%10oh,%10oi,%10oj,%10ok,%10ol,%10om,%10on,%10oo,%10op,%10op,%10oq,%10or,%10os,%10ot,%10ou,%10ov,%10ow,%10ox,%10oy,%10oz,%10pa,%10pb,%10pc,%10pd,%10pe,%10pf,%10pf,%10pg,%10ph,%10pi,%10pj,%10pk,%10pl,%10pm,%10pn,%10po,%10pp,%10pp,%10pq,%10pr,%10ps,%10pt,%10pu,%10pv,%10pw,%10px,%10py,%10pz,%10qa,%10qb,%10qc,%10qd,%10qe,%10qe,%10qf,%10qh,%10qi,%10qj,%10qk,%10ql,%10qm,%10qn,%10qo,%10qp,%10qp,%10qq,%10qr,%10qs,%10qt,%10qu,%10qu,%10qv,%10qw,%10qx,%10qy,%10qz,%10ra,%10rb,%10rc,%10rd,%10re,%10re,%10rf,%10rh,%10ri,%10rj,%10rk,%10rl,%10rm,%10rn,%10ro,%10rp,%10rp,%10rq,%10rr,%10rs,%10rt,%10ru,%10ru,%10rv,%10rw,%10rx,%10ry,%10rz,%10sa,%10sb,%10sc,%10sd,%10se,%10se,%10sf,%10sh,%10si,%10sj,%10sk,%10sl,%10sm,%10sn,%10so,%10sp,%10sp,%10sq,%10sr,%10ss,%10st,%10su,%10su,%10sv,%10sw,%10sx,%10sy,%10sz,%10ta,%10tb,%10tc,%10td,%10te,%10te,%10tf,%10th,%10ti,%10tj,%10tk,%10tl,%10tm,%10tn,%10to,%10tp,%10tp,%10tq,%10tr,%10ts,%10tt,%10tu,%10tu,%10tv,%10tw,%10tx,%10ty,%10tz,%10ua,%10ub,%10uc,%10ud,%10ue,%10ue,%10uf,%10uh,%10ui,%10uj,%10uk,%10ul,%10um,%10un,%10uo,%10up,%10up,%10uq,%10ur,%10us,%10ut,%10uu,%10uu,%10uv,%10uw,%10ux,%10uy,%10uz,%10va,%10vb,%10vc,%10vd,%10ve,%10ve,%10vf,%10vh,%10vi,%10vj,%10vk,%10vl,%10vm,%10vn,%10vo,%10vp,%10vp,%10vq,%10vr,%10vs,%10vt,%10vu,%10vu,%10vv,%10vw,%10vx,%10vy,%10vz,%10wa,%10wb,%10wc,%10wd,%10we,%10we,%10wf,%10wh,%10wi,%10wj,%10wk,%10wl,%10wm,%10wn,%10wo,%10wp,%10wp,%10wq,%10wr,%10ws,%10wt,%10wu,%10wu,%10wv,%10ww,%10wx,%10wy,%10wz,%10xa,%10xb,%10xc,%10xd,%10xe,%10xe,%10xf,%10xh,%10xi,%10xj,%10xk,%10xl,%10xm,%10xn,%10xo,%10xp,%10xp,%10xq,%10xr,%10xs,%10xt,%10xu,%10xu,%10xv,%10xw,%10xx,%10xy,%10xz,%10ya,%10yb,%10yc,%10yd,%10ye,%10ye,%10yf,%10yh,%10yi,%10yj,%10yk,%10yl,%10ym,%10yn,%10yo,%10yp,%10yp,%10yq,%10yr,%10ys,%10yt,%10yu,%10yu,%10yv,%10yw,%10yx,%10yz,%10za,%10zb,%10zc,%10zd,%10ze,%10ze,%10zf,%10zh,%10zi,%10zj,%10zk,%10zl,%10zm,%10zn,%10zo,%10zp,%10zp,%10zq,%10zr,%10zs,%10zt,%10zu,%10zu,%10zv,%10zw,%10zx,%10zy,%10zz'
slurm@frontend:~$
```

Beispiel-Batchscript

```
#!/bin/bash -l
#
#PBS -M michael.meier@fau.de
#PBS -m abe
#PBS -N LINPACK
#PBS -l walltime=01:45:00
#PBS -l nodes=16:ppn=24
#
```

Spezifikation von
Ressourcen und
anderen
Parametern

```
cd $PBS_O_WORKDIR
module load intelmpi
mpirun_rrze -npernode 12 -pinexpr "S0:0-5@S1:0-5" ./xhpl
```

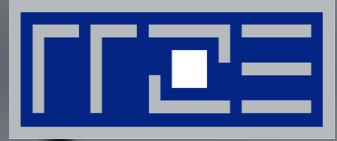
Arbeit!

Batchbetrieb

- **Scheduler** legt fest, wann der Job ausgeführt wird und berücksichtigt dabei Dinge wie...
 - **Priorität des Nutzers** und seiner Gruppe
 - wieviel Rechenzeit hat der Nutzer und seine Gruppe in der jüngeren Vergangenheit schon bekommen („Fair Share“)
 - Andere laufende Jobs
 - Reservierungen (z.B. für Nutzer oder Wartungsarbeiten)
 - Verfügbarkeit von Lizenzen die der Job braucht



ORGANISATORISCHES



- Die Vorträge im Überblick
- Andere Vortragsreihen des RRZE
- Ablageorte Vortragsfolien
- RRZE-Veranstaltungskalender / Mailingliste abonnieren
- Themenvorschläge & Anregungen

Weitere Vorträge zur „Systemausbildung“

26.04.2017 – Geschichte der Betriebssysteme

03.05.2017 – Unixoiden Betriebssysteme (Unix, Linux, OS X)

10.05.2017 – Systemüberwachung / Monitoring

17.05.2017 – Storage & Filesysteme

31.05.2017 – Windows-Betriebssysteme

21.06.2017 – High Performance Computing

28.06.2017 – Benutzerverwaltung: MS Active Directory

05.07.2017 – Virtualisierung

12.07.2017 – Backup / Archiv

19.07.2017 – Kerberos

26.07.2017 – IT-Sicherheit

- Immer mittwochs (ab 14 c.t.),
- Raum 2.049 im RRZE

Andere Vortragsreihen des RRZE

Campustreffen

- immer donnerstags ab 15 Uhr c.t.
- vermittelt Informationen zu den Dienstleistungen des RRZE
- befasst sich mit neuer Hard- & Software, Update-Verfahren sowie Lizenzfragen
- ermöglicht den Erfahrungsaustausch mit Spezialisten

Netzwerkausbildung „Praxis der Datenkommunikation“

- immer mittwochs in den Wintersemestern, ab 14 Uhr c.t.
- Vorlesungsreihe, die in die Grundlagen der Netztechnik einführt
- stellt die zahlreichen aktuellen Entwicklungen auf dem Gebiet der (universitären) Kommunikationssysteme dar

Vortragsfolien

Die Vortragsfolien werden nach der Veranstaltung auf der Webseite des RRZE abgelegt:

<http://www.rrze.fau.de/news/systemausbildung.shtml>

RRZE-Veranstaltungskalender & Mailinglisten

- Kalender abonnieren oder bookmarken
 - Alle Infos hierzu stehen auf der Webseite des RRZE unter:
<http://www.rrze.fau.de/news/kalender.shtml>
- Mailingliste abonnieren
 - Wöchentliche Terminhinweise werden zusätzlich an die Mailingliste [RRZE-Aktuelles](#) gesendet.
 - Auch diese Liste kann man abonnieren:
<https://lists.fau.de/mailman/listinfo/rrze-aktuelles>

Themenvorschläge & Anregungen

Themenvorschläge und Anregungen nehmen wir gerne entgegen!

Bitte schreiben Sie uns einfach eine E-Mail an:
rrze-zentrale@fau.de (Betreff: Systemausbildung)

REGIONALES RECHENZENTRUM ERLANGEN [RRZE]



Vielen Dank für Ihre Aufmerksamkeit!

Regionales RechenZentrum Erlangen [RRZE]

Martensstraße 1, 91058 Erlangen

<http://www.rrze.fau.de>