# CALCUL PARALLELE
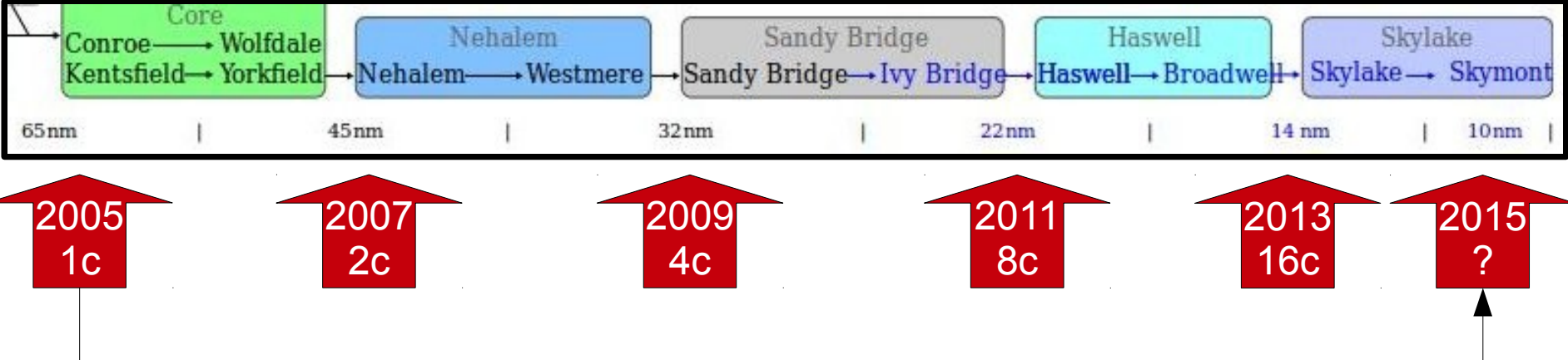# et
# APPLICATION AUX PLASMAS FROIDS

- Introduction

- Performance

- TPs
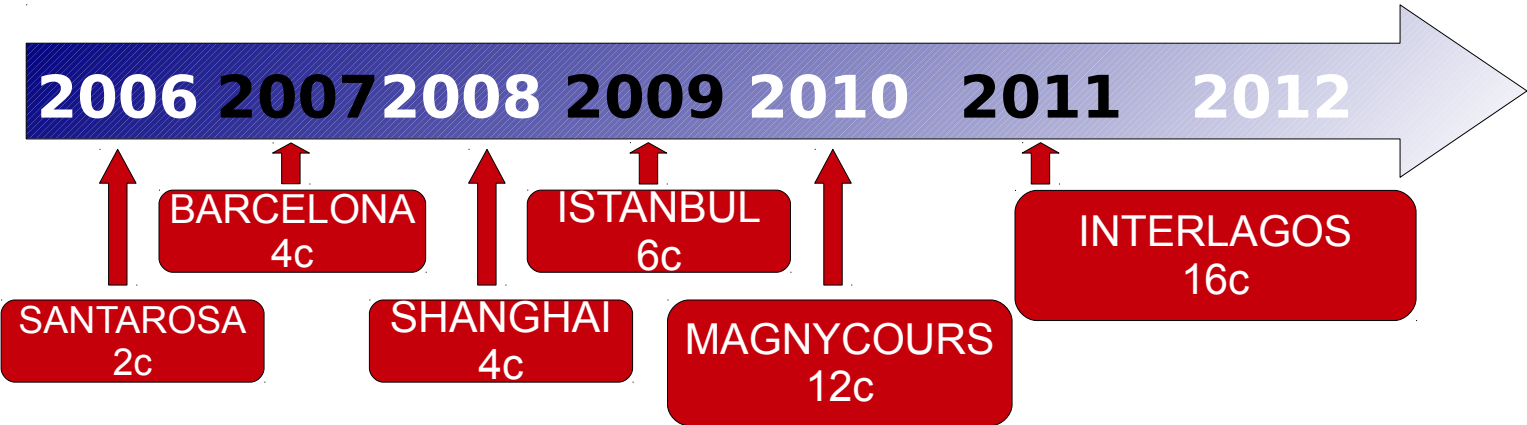
# Motivation



**INTEL ROADMAP**

| Core | | |
|---|---|---|
| Conroe → Wolfdale | Nehalem | Sandy Bridge | Haswell | Skylake |
| Kentsfield → Yorkfield | Nehalem → Westmere | Sandy Bridge → Ivy Bridge | Haswell → Broadwell | Skylake → Skymont |

65nm     |     45nm     |     32nm     |     22nm     |     14 nm     |   10nm    |

2005 1c    2007 2c    2009 4c    2011 8c    2013 16c    2015 ?

**AMD ROADMAP**

**2006 2007 2008 2009 2010 2011 2012**

SANTAROSA 2c

BARCELONA 4c

SHANGHAI 4c

ISTANBUL 6c

MAGNYCOURS 12c

INTERLAGOS 16c
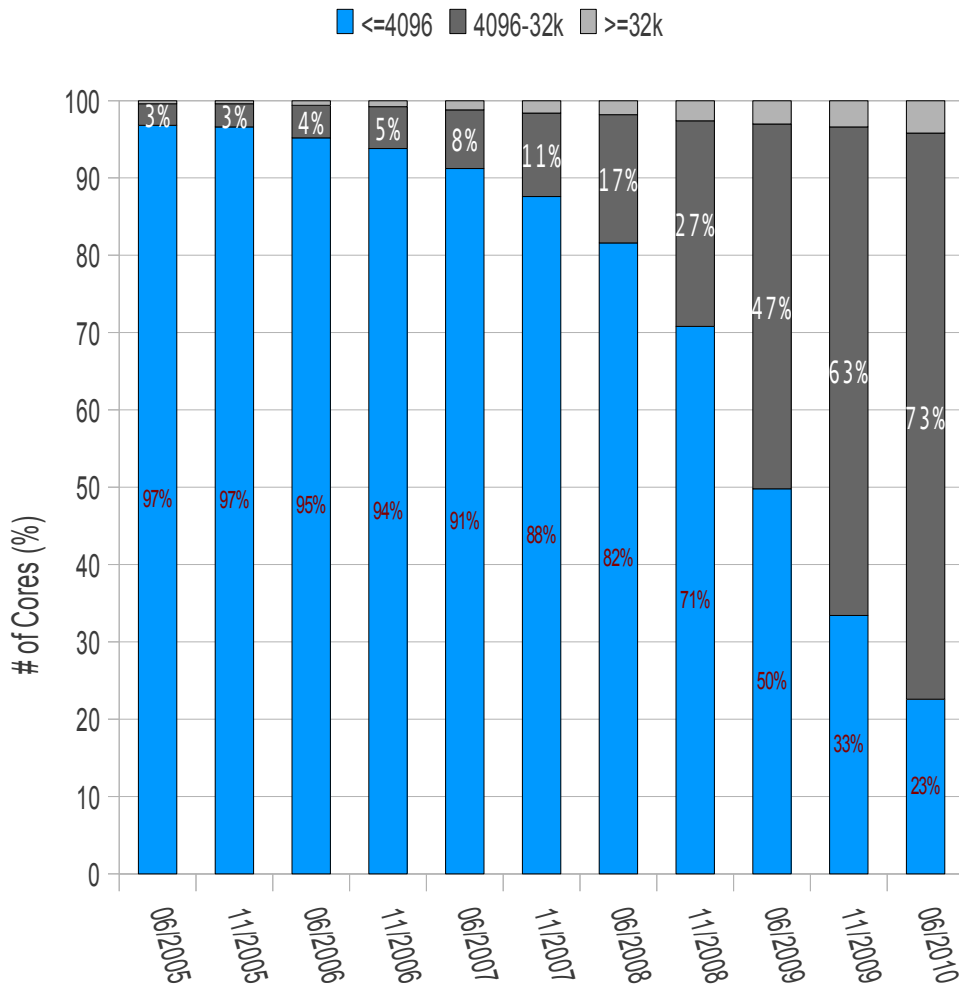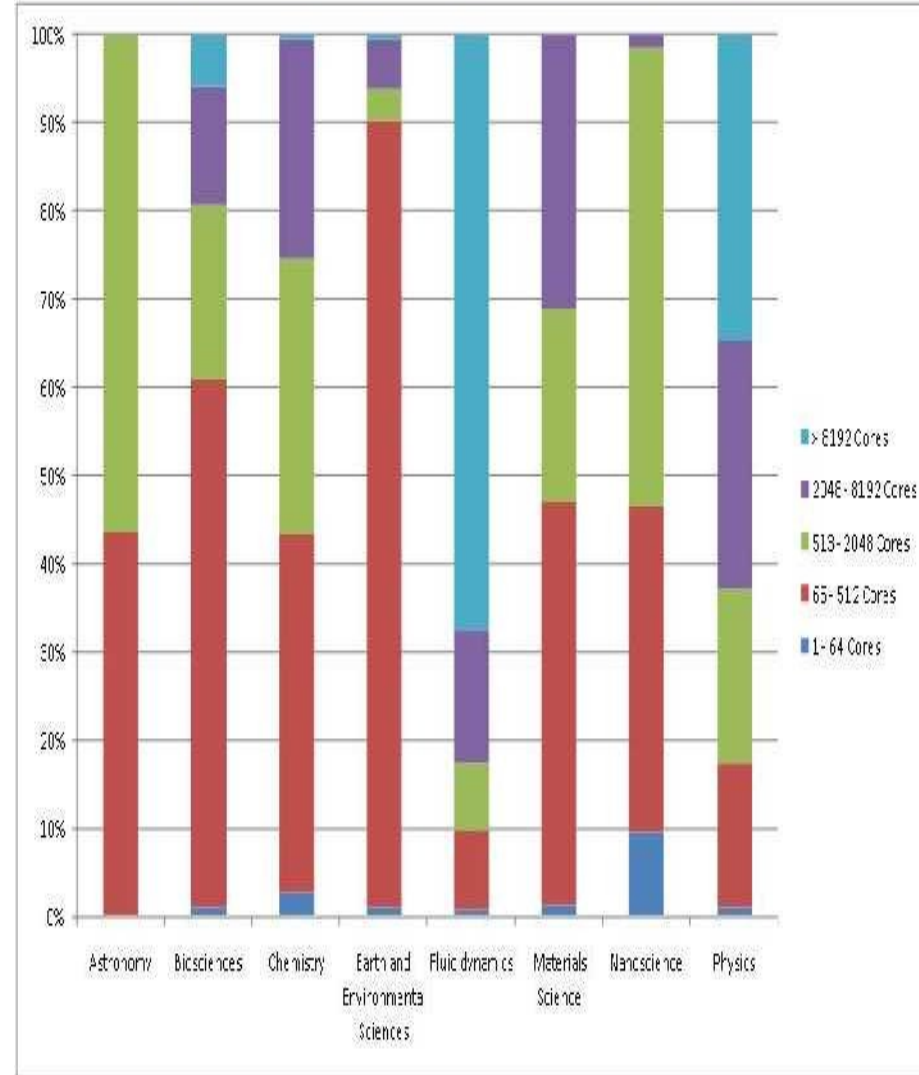
# Motivation



Evolution of systems in TOP500

Distribution of CSCS job sizes (2009)

# Example : Earth Weather forecasting (1 km resolution)

$$cells \approx 5.1E10 \ with$$
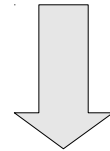$$\begin{cases} radius_{earth} &= 5.1 \, 10^8 \, km \\ levels &= 100 \\ cells &= 4\pi r^2 * levels \end{cases}$$

×

$$steps \approx 3.0E03 \ with$$
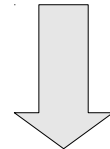$$\begin{cases} forecast &= 86400 \, sec \, (24h) \\ timeperstep &= 30 \, sec \\ steps &= forecast/timeperstep \end{cases}$$

×

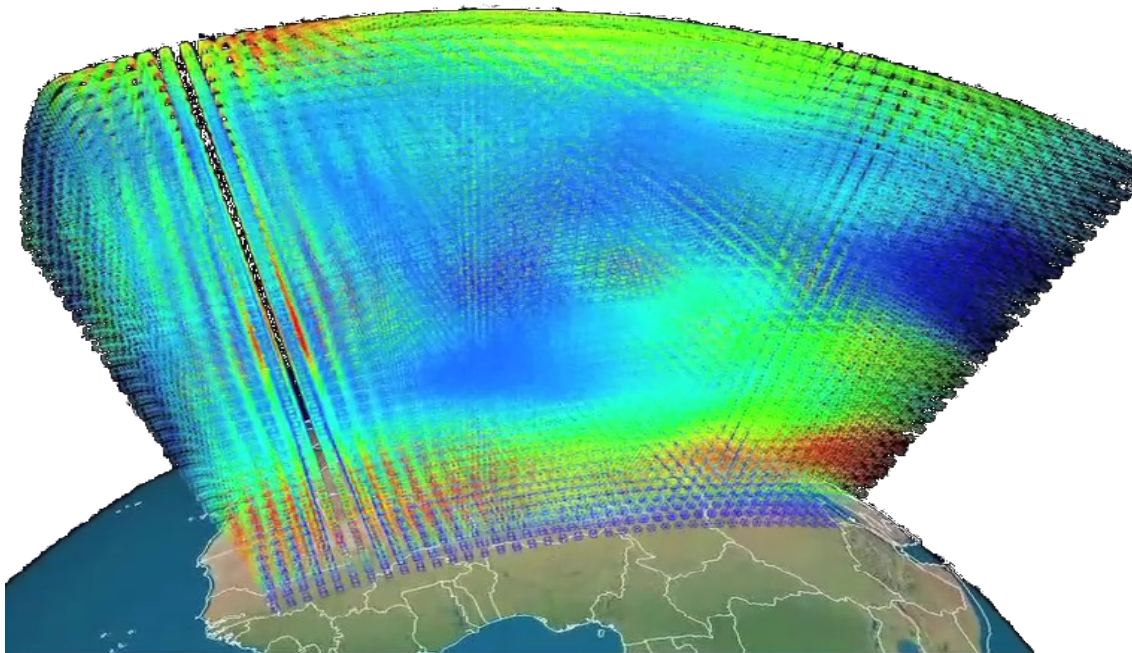$$op/cell/step \approx 1.0E03$$

1 simulated day = 1.5E17 flop

How much (real) time would it take to simulate 1 day ?

ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

CSCS
Swiss National Supercomputing Centre

# CSCS CRAY systems



| | ROSA CRAY XE6 (10/2011) | PALU CRAY XE6 (today) |
|---|---|---|
| Cores | **47872** cores<br>AMD Interlagos 2.1 Ghz<br>32 c/cnode<br>1496 cnodes | **4224** cores<br>AMD MagnyCours 2.1 Ghz<br>24 c/cnode<br>176 cnodes |
| Memory | **47** TB (DDR3)<br>32 GB/cnode<br>>= 1 GB/core | **5.6** TB (DDR3)<br>32 GB/cnode<br>>= 1.3 GB/core |
| RPeak Perf | **402** Tflops<br>8.4 Gflops/core<br>Water cooled | **35** Tflops<br>8.4 Gflops/core<br>Air cooled |

**ETH**
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

CSCS
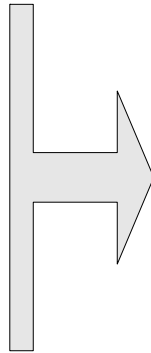Swiss National Supercomputing Centre

# Example : Earth Weather forecasting (1 km resolution)

1 simulated day = 1.5E17 flop

| Computational Performance | |
|---|---|
| 1 PetaFLOPS | 10^15 flop/sec |
| 1 TeraFLOPS | 10^12 flop/sec |
| 1 GigaFLOPS | 10^09 flop/sec |

| 1 simu. day | Rpeak performance | Time2solution |
|---|---|---|
| LAPTOP | 0.022 TFlops | 80 days |
| AUTRANS | 3.2 TFlops | 13 hours |
| ROSA | 212 TFlops | 12 min |
| JAGUAR | 2330 TFlops | 1.1 min |
| RIKKEN (#1) | 8000 TFlops | 19.1 sec |



Compared to my laptop :

Rosa is ~10 000 times faster,

Rikken is more than 360 000 times faster !

ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

CSCS
Swiss National Supercomputing Centre

# CRAY XT5

**Cray SeaStar2+ Interconnect**

9.6 GB/sec

9.6 GB/sec

9.6 GB/sec

9.6 GB/sec

9.6 GB/sec

9.6 GB/sec

**16 GB memory**

**25.6 GB/sec direct connect memory**

**6.4 GB/sec direct connect HyperTransport**

### CPU
- 12 cores per compute node (cnode)
  - *2.4 GHz AMD Opteron Istanbul*
- A total of 22128 compute cores
  - PEAK perf = 212 Tflop/s

### Memory
- *1,33 GB/core, 16 GB/cnode*
  - A total of 29.5TB
- 9.6 GB/s interconnect bandwidth

### I/O subsystem
- /scratch (287TB)
  - LUSTRE
- 12 GB/s sustained write bandwidth
- /project (400 TB) and /home (5.5T)

### Cray Linux Environment (CLE)
- Linux based operating system
  - Supports MPI and OpenMP
- Designed to run large scale apps

# Performance model

- Application name, version, language, parallelization, algo (?), etc...

- Machine name, configuration, etc...

- Programming env. name, version, etc... (compiler + opt flags), module list

- Scalability : Timings / # of cores => speedup, efficiency, etc...

- Profiling :

  - %MPI / %User code (CPU) / %IO

  - Memory usage

  - I/O : Outputs and inputs

    – File sizes / Number of files / File format, etc...

  - Hardware counters ==> gflops, % of peak, Cache hit/misses, TLB, etc...

  - Load imbalance,communication patterns, obstacles to scaling

  - More ? : OpenMP, tracing...

ETH
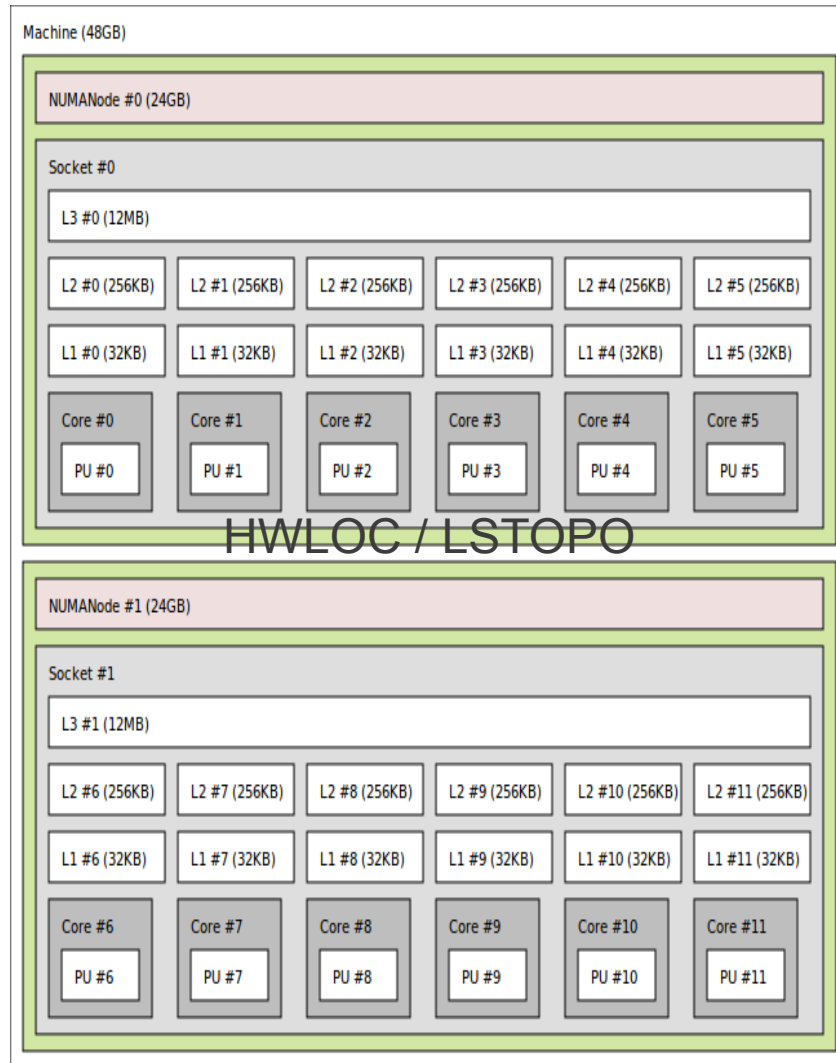Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

CSCS
Swiss National Supercomputing Centre

# Better know your machine : cpu/mpi



```
n1:/tmp $    more /proc/cpuinfo
processor       : 0
vendor_id       : GenuineIntel
cpu family      : 6
model           : 44
model name      : Intel(R) Xeon(R) CPU      X5660  @ 2.80GHz
stepping        : 2
cpu MHz         : 1600.000
cache size      : 12288 KB
```

```
n1:/tmp $    grep processor /proc/cpuinfo
processor       : 0
processor       : 1
processor       : 2
processor       : 3
processor       : 4
processor       : 5
processor       : 6
processor       : 7
processor       : 8
processor       : 9
processor       : 10
processor       : 11
```

Machine (48GB)

NUMANode #0 (24GB)

Socket #0

L3 #0 (12MB)

| L2 #0 (256KB) | L2 #1 (256KB) | L2 #2 (256KB) | L2 #3 (256KB) | L2 #4 (256KB) | L2 #5 (256KB) |

| L1 #0 (32KB) | L1 #1 (32KB) | L1 #2 (32KB) | L1 #3 (32KB) | L1 #4 (32KB) | L1 #5 (32KB) |

| Core #0 | Core #1 | Core #2 | Core #3 | Core #4 | Core #5 |
| PU #0 | PU #1 | PU #2 | PU #3 | PU #4 | PU #5 |

HWLOC / LSTOPO

NUMANode #1 (24GB)

Socket #1

L3 #1 (12MB)

| L2 #6 (256KB) | L2 #7 (256KB) | L2 #8 (256KB) | L2 #9 (256KB) | L2 #10 (256KB) | L2 #11 (256KB) |

| L1 #6 (32KB) | L1 #7 (32KB) | L1 #8 (32KB) | L1 #9 (32KB) | L1 #10 (32KB) | L1 #11 (32KB) |

| Core #6 | Core #7 | Core #8 | Core #9 | Core #10 | Core #11 |
| PU #6 | PU #7 | PU #8 | PU #9 | PU #10 | PU #11 |

ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

CSCS
Swiss National Supercomputing Centre

# Better know your machine : memory

```
n1:/tmp $     more /proc/meminfo
MemTotal:          49554172 kB
MemFree:           40151400 kB
Buffers:             253724 kB
```

```
valgrind --tool=cachegrind hostname
```

```
n1:/home/piccinali $    cg_annotate cachegrind.out.26322
I1 cache:           32768 B, 64 B, 4-way associative
D1 cache:           32768 B, 64 B, 8-way associative
LL cache:           12582912 B, 64 B, 24-way associative
```

```
top - 17:40:30 up  7:56,  5 users,  load average: 0.00, 0.01, 0.05
Tasks: 165 total,   1 running, 164 sleeping,   0 stopped,   0 zombie
Cpu0  :  6.7%us,  0.9%sy,  0.0%ni, 92.3%id,  0.1%wa,  0.0%hi,  0.0%si,  0.0%st
Cpu1  :  3.3%us,  0.4%sy,  0.0%ni, 96.3%id,  0.0%wa,  0.0%hi,  0.0%si,  0.0%st
Cpu2  :  0.9%us,  0.1%sy,  0.0%ni, 99.0%id,  0.0%wa,  0.0%hi,  0.0%si,  0.0%st
Cpu3  :  0.7%us,  0.1%sy,  0.0%ni, 99.2%id,  0.0%wa,  0.0%hi,  0.0%si,  0.0%st
Cpu4  :  0.7%us,  0.0%sy,  0.0%ni, 99.2%id,  0.0%wa,  0.0%hi,  0.0%si,  0.0%st
Cpu5  :  0.6%us,  0.1%sy,  0.0%ni, 99.4%id,  0.0%wa,  0.0%hi,  0.0%si,  0.0%st
Cpu6  :  6.6%us,  2.8%sy,  0.0%ni, 90.4%id,  0.1%wa,  0.0%hi,  0.0%si,  0.0%st
Cpu7  :  2.8%us,  0.4%sy,  0.0%ni, 96.7%id,  0.0%wa,  0.0%hi,  0.0%si,  0.0%st
Cpu8  :  0.7%us,  0.1%sy,  0.0%ni, 99.3%id,  0.0%wa,  0.0%hi,  0.0%si,  0.0%st
Cpu9  :  0.7%us,  0.1%sy,  0.0%ni, 99.2%id,  0.0%wa,  0.0%hi,  0.0%si,  0.0%st
Cpu10 :  0.7%us,  0.0%sy,  0.0%ni, 99.3%id,  0.0%wa,  0.0%hi,  0.0%si,  0.0%st
Cpu11 :  0.6%us,  0.0%sy,  0.0%ni, 99.4%id,  0.0%wa,  0.0%hi,  0.0%si,  0.0%st
Mem:  49554172k total,  9402852k used, 40151320k free,   253736k buffers
Swap: 50319356k total,        0k used, 50319356k free,  7337704k cached

  PID USER      PR  NI  VIRT  RES  SHR S %CPU %MEM    TIME+  COMMAND
    1 root      20   0 24140 2276 1324 S    0  0.0  0:02.65 init
    2 root      20   0     0    0    0 S    0  0.0  0:00.00 kthreadd
    3 root      20   0     0    0    0 S    0  0.0  0:02.21 ksoftirqd/0
```

HWLOC / LSTOPO

Socket #0

L3 #0 (12MB)

| L2 #0 (256KB) | L2 #1 (256KB) | L2 #2 (256KB) | L2 #3 (256KB) | L2 #4 (256KB) | L2 #5 (256KB) |

| L1 #0 (32KB) | L1 #1 (32KB) | L1 #2 (32KB) | L1 #3 (32KB) | L1 #4 (32KB) | L1 #5 (32KB) |

| Core #0 | Core #1 | Core #2 | Core #3 | Core #4 | Core #5 |
| PU #0 | PU #1 | PU #2 | PU #3 | PU #4 | PU #5 |

ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

ANGD 10/2011

CSCS
Swiss National Supercomputing Centre

# Better know your machine : memory caches (psinv / INTEL)

```
n1:/home/piccinali $  /softs/perfsuite/1.0.0/gnu/bin/psinv
System Information -
Node Name:           n1
OS Name:             Linux
OS Release:          2.6.38-11-server
OS Build/Version:    #48-Ubuntu SMP Fri Jul 29 19:20:32 UTC 2011
OS Machine:          x86_64
Processors:          12
Total Memory (MB):   48392.75
System Page Size (KB): 4.00


Processor Information -
Vendor:              Intel
Processor family:    Pentium Pro (P6)
Brand:               Intel(R) Xeon(R) CPU        X5660  @ 2.80GHz
Model (Type):        (unknown)
Revision:            2
Clock Speed:         1600.00 MHz
```

```
Cache and TLB Information -
Cache levels:              3

Cache Details -
Level 1:
        Type:              Instruction
        Size:              32 KB
        Line size:         64 bytes
        Associativity:     4-way set associative

        Type:              Data
        Size:              32 KB
        Line size:         64 bytes
        Associativity:     8-way set associative

Level 2:
        Type:              Unified
        Size:              256 KB
        Line size:         64 bytes
        Associativity:     8-way set associative

Level 3:
        Type:              Unified
        Size:              12.00 MB
        Line size:         64 bytes
        Associativity:     16-way set associative
```

ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

CSCS
Swiss National Supercomputing Centre

# Better know your machine : memory caches (psinv / AMD)

```
palul:~ $    aprun -n1 psinv |head -38
System Information -
Node Name:              nid00180
OS Name:                Linux
OS Release:             2.6.27.48-0.12.1_1.0301.5737-cray_gem_c
OS Build/Version:       #1 SMP Mon Mar 28 22:26:26 UTC 2011
OS Machine:             x86_64
Processors:             24
Total Memory (MB):      32311.14
System Page Size (KB):  4.00

Processor Information -
Vendor:                 AMD
Processor family:       K10
Brand:                  AMD Opteron(tm) Processor 6172
Model:                  (unknown)
Revision:               1
Clock Speed:            2100.00 MHz

Cache and TLB Information -
Cache levels:           3

Cache Details -
Level 1:
        Type:           Instruction
        Size:           64 KB
        Line size:      64 bytes
        Associativity:  2-way set associative

        Type:           Data
        Size:           64 KB
        Line size:      64 bytes
        Associativity:  2-way set associative

Level 2:
        Type:           Unified
        Size:           512 KB
        Line size:      64 bytes
        Associativity:  16-way set associative
```

```
Level 3:
        Type:           Unified
        Size:           10.00 MB
        Line size:      64 bytes
        Associativity:  96-way set associative

TLB Details -
Level 1:
        Type:           Instruction
        Entries:        32
        Pagesize (KB):  4
        Associativity:  Fully associative

        Type:           Data
        Entries:        48
        Pagesize (KB):  4
        Associativity:  Fully associative

Level 2:
        Type:           Data
        Entries:        512
        Pagesize (KB):  4
        Associativity:  4-way set associative

        Type:           Instruction
        Entries:        512
        Pagesize (KB):  4
        Associativity:  4-way set associative
```

ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

CSCS
Swiss National Supercomputing Centre

# Better know your machine : i/o

```
n1:/tmp $   df -h
Filesystem              Size  Used Avail Use% Mounted on
/dev/mapper/n1-root      87G  9.0G   73G  11% /
none                    24G  224K   24G   1% /dev
none                    24G     0   24G   0% /dev/shm
none                    24G  404K   24G   1% /var/run
none                    24G     0   24G   0% /var/lock
/dev/sda1              228M   45M  172M  21% /boot
```

```
rosa6:~ $    df -Ph | grep -E "project|scratch|home|File"
Filesystem                   Size  Used Avail Use% Mounted on
globalhome.cscs.ch:/apps     22T  1.5T   21T   7% /apps
globalhome.cscs.ch:/users    44T  3.7T   40T   9% /users
projects1.cscs.ch:/global   1.4P  767T  631T  55% /project
263@ptl:/scratch            287T  153T  119T  57% /scratch/rosa
```

ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

CSCS
Swiss National Supercomputing Centre

# Getting started (1)

mpif90 Laplace_mpi.F90

```
n1:/home/piccinali/trunk/debug/intro/f90 $   sbatch.sh
USAGE :
                    arg1=exe
                    arg2=mppwidth
                    arg3=mppnppn
                    arg4=mppdepth
                    arg5=exeargs
                    arg6=prempiexec
                    arg7=postmpiexec
```

```
n1:/home/piccinali/trunk/laplace/src $   sbatch.sh ./exe.03 2 2 1 "1920 1920 200 1.0d-5" "" -bind-to-core
+ export OMP_NUM_THREADS=1
+ OMP_NUM_THREADS=1
+ echo '/usr/bin/time -p  /softs/openmpi-1.4.3/bin/mpiexec -bind-to-core -n 2 -npernode 2 -x OMP_NUM_THREADS -hostfile /softs
/openmpi-1.4.3/h ./exe.03 1920 1920 200 1.0d-5'
```

```
n1:/home/piccinali/trunk/laplace/src $   cat o_exe.03.0002.2.2.1.1920-1920-200-1.0d-5--bind-to-core
/usr/bin/time -p  /softs/openmpi-1.4.3/bin/mpiexec -bind-to-core -n 2 -npernode 2 -x OMP_NUM_THREADS -hos
tfile /softs/openmpi-1.4.3/h ./exe.03 1920 1920 200 1.0d-5
 k =          100  err =   2.74535544735804438E-003
 k =          200  err =   1.31782906355876195E-003
 erreur_max finale :   1.31109160900838573E-003  it             201
Time (seconds) :      5.19
real 5.31
user 10.17
sys 0.39
```

```
#ifdef _ADIOS
#ifdef _HWMEM
#ifdef _VERBOSE
#ifdef _WITHO
#ifdef _XPAT
```

# Getting started (2)

module permet de modifier l'environnement utilisateur facilement,
en configurant les variables nécessaires à la compilation/exécution.

* module avail : affiche l'environnement disponible
* module list : affiche l'environnement actuel
* module load <module> : ajoute application/bibliothèque/compilateur à l'environnement
* module unload <module> : enlève application/bibliothèque/compilateur de l'environnement
* module show : affiche le contenu du module
* module swap <module1> <module2> : remplace version1 avec version2
* module help <module> : help
* module whatis <module> : help

## http://modules.sf.net

ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

CSCS
Swiss National Supercomputing Centre

# Step0 : Compilation flags

```
gfortran -O3 Laplace_seq.F90 -o gnu.O3        ifort -O3 Laplace_seq.F90 -o int.O3
```

```
n1:/home/piccinali/trunk/laplace/src/seq $    echo 1920 1920 200 1.0d-5 | /usr/bin/time -p ./gnu.O3
 k =          100  erreur =   2.74535544735804438E-003
 k =          200  erreur =   1.31782906355876195E-003
real 7.08
user 6.99
sys 0.07
```

```
n1:~/trunk/laplace/src/seq > grep real o_gnu.O*
o_gnu.O0:real 29.92
o_gnu.O1:real 7.76
o_gnu.O2:real 7.16
o_gnu.O3:real 7.08
o gnu.O4:real 7.12
```

```
n1:~/trunk/laplace/src/seq > grep real o_int.O*
o_int.O0:real 45.79
o_int.O1:real 6.95
o_int.O2:real 6.59
o_int.O3:real 6.58
o_int.O4:real 6.59
```

ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

ANGD 10/2011

CSCS
Swiss National Supercomputing Centre

# Step0 : Compilers

```
gele2:~/GNU $ ~/ftn -O3 ~/Laplace_mpi.F90
gele2:~/GNU $ aprun -n 2 a.out 1920 1920 200 1.0d-5
 k =           100  err =   2.74535544735804438E-003
 k =           200  err =   1.31782906355876195E-003
 erreur_max finale :   1.31109160900838573E-003  it          201
Time (seconds) :      9.31
Application 6133 resources: utime 0, stime 0
gele2:~/GNU $
gele2:~/GNU $
```

```
⊞ All
gele2:~/INTEL $ ~/ftn -O3 ~/Laplace_mpi.F90
gele2:~/INTEL $ aprun -n 2 a.out 1920 1920 200 1.0d-5
 k =           100  err =   2.745355447358433E-003
 k =           200  err =   1.317829063558706E-003
 erreur_max finale :   1.311091609008386E-003  it          201
Time (seconds) :      9.08
Application 6132 resources: utime 0, stime 0
gele2:~/INTEL $
gele2:~/INTEL $
```

```
⊞ All
gele2:~/PGI $ ~/ftn -O3 ~/Laplace_mpi.F90
gele2:~/PGI $ aprun -n 2 a.out 1920 1920 200 1.0d-5
 k =           100  err =   2.7453554473582109E-003
 k =           200  err =   1.3178290635587619E-003
 erreur_max finale :   1.3110916090082192E-003  it          201
Time (seconds) :      9.42
Application 6134 resources: utime 0, stime 0
gele2:~/PGI $
gele2:~/PGI $
```

```
⊞ All
gele2:~/CRAY $ ~/ftn -O3 ~/Laplace_mpi.F90
gele2:~/CRAY $ aprun -n 2 a.out 1920 1920 200 1.0d-5
 k =   100  err =  2.74535544735809989E-3
 k =   200  err =  1.31782906355865093E-3
 erreur_max finale :  1.31109160900833022E-3  it  201
Time (seconds) :      9.18
Application 6135 resources: utime 0, stime 0
```

ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

CSCS
Swiss National Supercomputing Centre

# Where do I spend time (callgrind) ?



```
n1:/home/piccinali/trunk/laplace/src/seq $  echo 1920 1920 20 1.0d-5 | valgrind --tool=callgrind ./gnu.O3g
==24995== Callgrind, a call-graph generating
==24995== Copyright (C) 2002-2010, and GNU G
==24995== Using Valgrind-3.7.0.SVN and LibVE
==24995== Command: ./gnu.O3g
==24995==
==24995== For interactive control, run 'call
==24995==
```

callgrind.out.24995 [./gnu.O3g]

File   View   Go   Settings   Help

Open    Back    Forward    Up    » Cycle Estimation

**MAIN__**              Kcachegrind callgrind.out.24995
                              callgrind_annotate
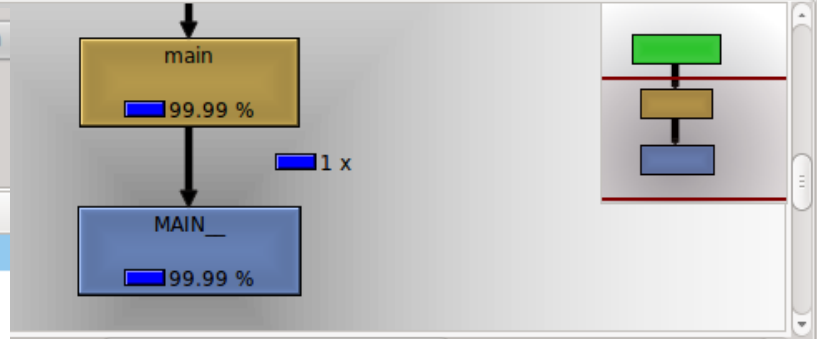
Types | Callers | All Callers | Callee Map | Source Code

| # | CEst | Source ('/home/piccinali/trunk/laplace/src/seq/Laplace seq.F90') |
|---|------|------------------------------------------------------------------|
| 192 | | ! Resolution de l'equation de Laplace |
| 193 | 0.00 | do j = 1, ny |
| 194 | 3.98 | do i = 1, nx |
| 195 | | V(i,j) = omega1 * V_old(i,j) & |
| 196 | | + ( sou(i,j) & |
| 197 | | + ve(i,j) * V(i+1,j ) + vw(i,j) * V(i-1,j ) & |
| 198 | 33.73 | + vn(i,j) * V(i ,j+1) + vs(i,j) * V(i ,j-1) ) * vc(i,j) |
| 199 | 29.75 | if ( abs(V(i,j) ) < 1d-20) V(i,j) = 0.0d0 |
| 200 | | end do |
| 201 | | end do |
| 202 | | |

**Flat Profile**

Search: _____ (No Grouping)

| Incl. | Self | Called | Function | Location |
|-------|------|--------|----------|----------|
| 99.99 | 0.00 | 1 | 0x00000000004008b0 | gnu.O3g |
| 99.99 | 0.00 | 1 | (below main) | libc-2.13.so: libc-start.c |
| 99.99 | 0.00 | 1 | main | gnu.O3g: Laplace_seq.F90 |
| 99.99 | 99.99 | 1 | MAIN__ | gnu.O3g: Laplace_seq.F90 |

Open    Back    Forward    Up    » Instruction Fetch

**MAIN__**

Types | Callers | All Callers | Callee Map | Source Code

| # | Ir | Source ('/home/piccinali/trunk/laplace/src/seq/Laplace seq.F90') |
|---|-----|------------------------------------------------------------------|
| 198 | 33.73 | + vn(i,j) * V(i ,j+1) + vs(i,j) * V(i ,j-1) ) * vc(i,j) |
| 199 | 29.75 | if ( abs(V(i,j) ) < 1d-20) V(i,j) = 0.0d0 |
| 208 | 11.90 | erreur_max = max (erreur_max, erreur) |

main
99.99 %

1 x

MAIN__
99.99 %

Parts | Callees | Call Graph | All Callees | Caller Map | Machine Code

callgrind.out.24995 [1] - Total Cycle Estimation Cost: 3 716 149 348

**ETH**
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

# What is my main memory usage (memcheck) ?

`echo 1920 1920 20 1.0d-5 | valgrind --tool=massif ./gnu.03g`

```
--------------------------------------------------------------------
Command:            ./gnu.03g          ms_print massif.out.25092
Massif arguments:   (none)
ms_print arguments: massif.out.25092
--------------------------------------------------------------------

    MB
253.4^@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@#
    |@                                                                              #
    |@                                                                              #
    |@                                                                              #
    |@                                                                              #
    |@                                                                              #
    |@                                                                              #
  0 +-------------------------------------------------------------------->Gi
    0                                                                      3.461

Number of snapshots: 71
 Detailed snapshots: [9, 19, 29, 39, 40 (peak), 50, 60, 70]
```

`echo 1920 1920 20 1.0d-5 | valgrind --tool=memcheck ./gnu.00g`

```
n1:~/trunk/laplace/src/seq > echo 1920 1920 20 1.0d-5 | valgrind --tool=memcheck ./gnu.00g
==25088== Memcheck, a memory error detector
==25088== Copyright (C) 2002-2010, and GNU GPL'd, by Julian Seward et al.
==25088== Using Valgrind-3.7.0.SVN and LibVEX; rerun with -h for copyright info
==25088== Command: ./gnu.00g
==25088==
==25088==
==25088== HEAP SUMMARY:
==25088==     in use at exit: 0 bytes in 0 blocks
==25088==   total heap usage: 34 allocs, 34 frees, 265,679,857 bytes allocated
==25088==
==25088== All heap blocks were freed -- no leaks are possible
==25088==
==25088== For counts of detected and suppressed errors, rerun with: -v
==25088== ERROR SUMMARY: 0 errors from 0 contexts (suppressed: 2 from 2)
```

| n | time(i) | total(B) | useful-heap(B) | extra-heap(B) | stacks(B) |
|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 | 0 |
| 1 | 206,404 | 680 | 672 | 8 | 0 |

| n | time(i) | total(B) | useful-heap(B) | extra-heap(B) | stacks(B) |
|---|---|---|---|---|---|
| 40 | 3,716,093,177 | 265,715,184 | 265,678,657 | 36,527 | 0 |
| 69 | 3,716,101,900 | 56 | 48 | 8 | 0 |
| 70 | 3,716,101,949 | 0 | 0 | 0 | 0 |

```
00.00% (0B) (heap allocation functions) malloc/new/new[], --alloc-fns, etc.
->00.00% (0B) in 1+ places, all below ms_print's threshold (01.00%)
```

# What is my L1/L3 cache usage (cachegrind) ?

```
n1:/home/piccinali/trunk/laplace/src/seq $    echo 1920 1920 20 1.0d-5 | valgrind --tool=cachegrind ./gnu.O3g
==25248== Cachegrind, a cache and branch-prediction profiler
==25248== Copyright (C) 2002-2010, and GNU GPL'd, by Nicholas Nethercote et al.
==25248== Using Valgrind-3.7.0.SVN and LibVEX; rerun with -h for copyright info
==25248== Command: ./gnu.O3g
==25248==
--25248-- warning: L3 cache found, using its data for the LL simulation.
--25248-- warning: pretending that LL cache has associativity 24 instead of actual 16
==25248==
==25248== I   refs:      3,716,149,724
==25248== I1  misses:            1,395
==25248== LLi misses:            1,369
==25248== I1  miss rate:          0.00%
==25248== LLi miss rate:          0.00%
==25248==
==25248== D   refs:      1,269,911,041  (1,048,976,325 rd  + 220,934,716 wr)
==25248== D1  misses:       135,328,344  (  121,463,449 rd  +  13,864,895 wr)
==25248== LLd misses:       115,735,040  (  101,907,212 rd  +  13,827,828 wr)
==25248== D1  miss rate:         10.6% (         11.5%  +        6.2%  )
==25248== LLd miss rate:          9.1% (          9.7%  +        6.2%  )
==25248==
==25248== LL refs:         135,329,739  (  121,464,844 rd  +  13,864,895 wr)
==25248== LL misses:       115,736,409  (  101,908,581 rd  +  13,827,828 wr)
==25248== LL miss rate:          2.3% (          2.1%  +        6.2%  )
```

```
n1:/home/piccinali/trunk/laplace/src/seq $    cg_annotate --auto=yes cachegrind.out.25248
--------------------------------------------------------------------
I1 cache:         32768 B, 64 B, 4-way associative
D1 cache:         32768 B, 64 B, 8-way associative
LL cache:         12582912 B, 64 B, 24-way associative
Command:          ./gnu.O3g
Data file:        cachegrind.out.25248
Events recorded:  Ir I1mr ILmr Dr D1mr DLmr Dw D1mw DLmw
Events shown:     Ir I1mr ILmr Dr D1mr DLmr Dw D1mw DLmw
Event sort order: Ir I1mr ILmr Dr D1mr DLmr Dw D1mw DLmw
Thresholds:       0.1 100 100 100 100 100 100 100 100
Include dirs:
User annotated:
Auto-annotation:  on
```

| | | Open | Back | Forward | Up | » | Data Read Access |

MAIN_                                                    kcachegrind

| Types | Callers | All Callers | Callee Map | Source Code |

| # | Dr | Source ('/home/piccinali/trunk/laplace/src/seq/Laplace seq.F90') |
|---|---|---|
| 198 | 77.31 | + vn(i,j) * V(i ,j+1) + vs(i,j) * V(i ,j-1) ) * vc(i,j) |
| 207 | 14.06 | erreur = abs (V(i,j) - V_old(i,j) ) |
| 188 | 7.03 | V_old(i,j) = V(i,j) |

```
-------------------------------------------------------------------------------------------
       Ir  I1mr  ILmr            Dr         D1mr        DLmr           Dw       D1mw         DLmw
-------------------------------------------------------------------------------------------
3,716,149,724 1,395 1,369 1,048,976,325 121,463,449 101,907,212 220,934,716 13,864,895 13,827,828  PROGRAM TOTALS
```

# Where is my L1/L3 cache usage (kcachegrind) ?

ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

CSCS
Swiss National Supercomputing Centre

Objectif :

* Montrer l'importance des indices de boucles sur la performance du code

Instructions :

* cp /home/piccinali/trunk/matmul/seq/matmult.F90  $HOME ; cd $HOME

* Compilez le programme : gfortran -D_A -g -O3 matmult.F90 -o A ;

* Executez le programme et notez le temps d'execution.

* Recommencez avec -D_B,-D_C,-D_D,-D_E,-D_F et comparez les resultats.

* Utilisez cachegrind pour trouver une explication

ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

ANGD 10/2011

CSCS
Swiss National Supercomputing Centre

# TP : matmul (solution)

```
gfortran -D_A -g -O3 matmult.F90 -o A

/usr/bin/time -p ./A
```

```
A/oo.A:real 80.51
B/oo.B:real 155.29
C/oo.C:real 78.16
D/oo.D:real 3.21
E/oo.E:real 155.73
F/oo.F:real 4.07
```



```
                                    cachegrind.out.26268 [./A]
File  View  Go  Settings  Help

  Open      Back        Forward        Up         »  L1 Data Read Miss

MAIN__

Types  Callers  All Callers  Callee Map  Source Code

#   D1mr      Source ('/home/piccinali/trunk/matmul/seq/A/../matmult.F90')
17
18            ! Code to illustrate example given in talk
19            ! (c) Manchester Computing  Spring 1999
...   ...
41
42            #ifdef _A
43                print *,"ijk"
44                do i=1,n
45                  do j=1,n
46                    do k=1,n
47                      a(i,j) = a(i,j) + b(i,k)*c(k,j)
48    99.61         end do
```

17 825 792 misses

```
                                    ../D/cachegrind.out.26095 [./D]
File  View  Go  Settings  Help

  Open      Back        Forward        Up      %  Relative   »  L1 Data Read Miss

MAIN__

Types  Callers  All Callers  Callee Map  Source Code

#   D1mr      Source ('/home/piccinali/trunk/matmul/seq/D/../matmult.F90')
13            ! ms_print massif.out.29212
14            ! valgrind --time-unit=B => FORTRAN ?
15            ! It is worth emphasising that by default Massif measures only heap memory, i.e.
16    0.39      program matrix_multiplication
17
18            ! Code to illustrate example given in talk
19            ! (c) Manchester Computing  Spring 1999
...   ...
78                  do k=1,n
79                do i=1,n
80                    a(i,j) = a(i,j) + b(i,k)*c(k,j)
81    99.26         end do
82              end do
83            end do
84            #endif
```

1 052 672 misses

ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

CSCS
Swiss National Supercomputing Centre

```
/softs/openmpi-1.4.3/bin/mpiexec -n 2 -bind-to-core -npernode 2 -hostfile /softs/openmpi-1.4.3/h ./exe.03 1920 1920 200 1.0d-5
```

**-bind-to-core**

```
n1:/home/piccinali/trunk/laplace/src $   grep Time o.00 o.01 o.02 o.03 o.04
o.00:Time (seconds) :     14.76
o.01:Time (seconds) :      5.54
o.02:Time (seconds) :      5.23
o.03:Time (seconds) :      5.21
o.04:Time (seconds) :      5.24
```

**bind-to-none**

```
n1:/home/piccinali/trunk/laplace/src $   grep Time o.g00 o.g01 o.g02 o.g03 o.g04
o.g00:Time (seconds) :     15.85
o.g01:Time (seconds) :      5.57
o.g02:Time (seconds) :      3.93
o.g03:Time (seconds) :      5.20
o.g04:Time (seconds) :      5.50
```

**bind-to-socket**

```
n1:/home/piccinali/trunk/laplace/src $   grep Time o.s00 o.s01 o.s02 o.s03 o.s04
o.s00:Time (seconds) :     14.68
o.s01:Time (seconds) :      5.45
o.s02:Time (seconds) :      5.20
o.s03:Time (seconds) :      5.16
o.s04:Time (seconds) :      5.18
```

**ETH**
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

CSCS
Swiss National Supercomputing Centre

# Step1 : Scalability study

## Speedup = t1 / tn

t1 : temps de l'agorithme sequentiel,

tn : temps de l'agorithme parallele sur n processeurs,

Linear speedup : acceleration lineaire avec le nombre de cpus (rare mais possible)

Efficiency : Speedup / n

Weak scaling : la taille du probleme augmente avec le nombre de cpus

Strong scaling : la taille du probleme ne varie pas en fonction du nombre de cpus (harder)

```
0002cores  t=5.19
0006cores  t=5.21
0012cores  t=3.10
0024cores  t=2.06
0036cores  t=1.72
0048cores  t=1.47
```

http://www.sc2000.org/bell/twelve-ways.txt

Twelve Ways to Fool the Masses When Giving Performance Results on Parallel Computers

David H. Bailey

ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

CSCS
Swiss National Supercomputing Centre

# Step1 : Loi Amdahl (1967)

$$Speedup = \frac{1}{(1-P) + \dfrac{P}{N}}$$

P = fraction // du code ; 1-P = fraction non // du code

N = nombre de cpus

**Corollaire 1 :**

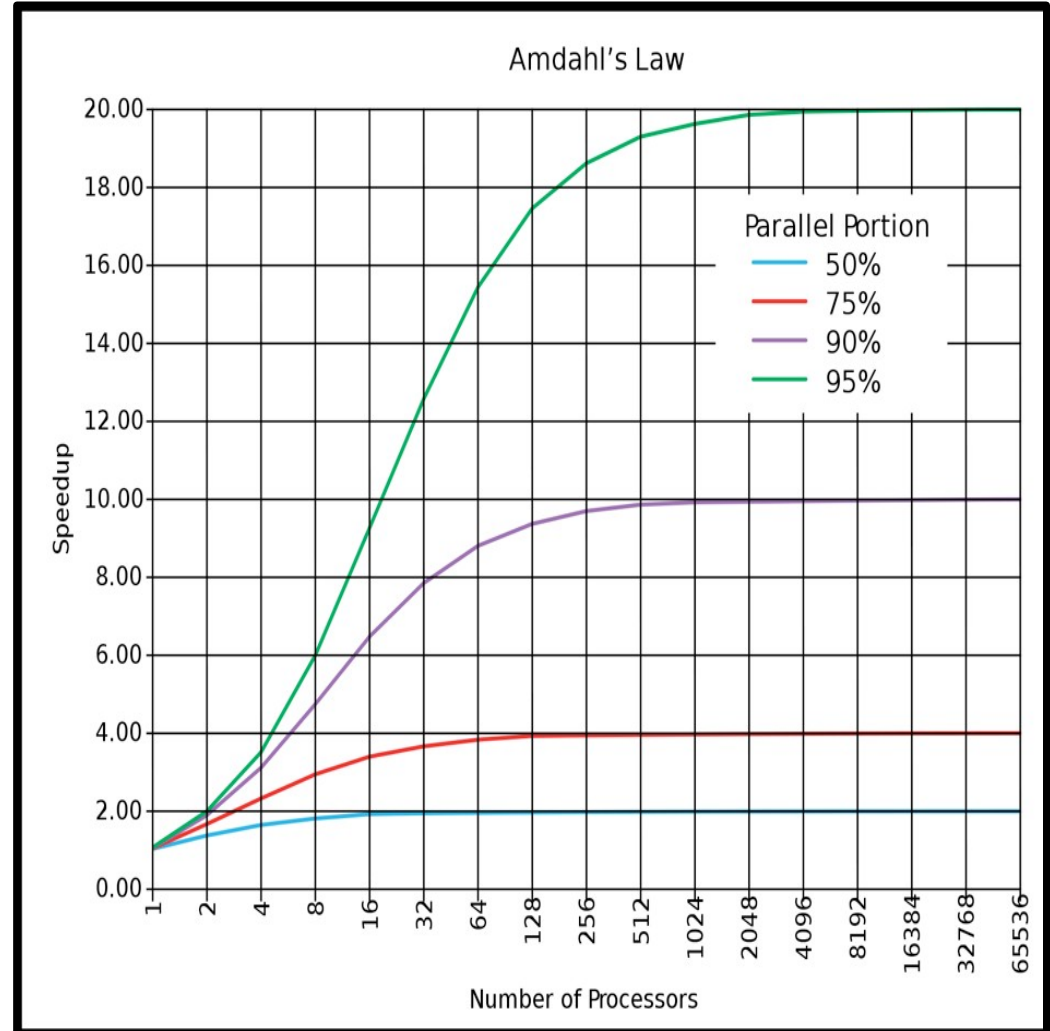=> si $N \rightarrow +\infty$

=> alors Speedup <= 1/(1-P)

=> performance majoree par la partie non parallele du code (1-P), independemment du nombre de cpus :-(

=> ex : si P=90% alors Speedup <= 10

**Corollaire 2 :**

=> estimation de P a partir du Speedup

=> P ~= (1-Sp)*N / (Sp * (1-N))

(avec N=cores)



Amdahl's Law

Parallel Portion
— 50%
— 75%
— 90%
— 95%

Speedup

Number of Processors

ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

CSCS
Swiss National Supercomputing Centre

# How to time MPI programs ?

Use the linux time command :

/usr/bin/time -p aprun -n12 exe

When job finishes, time writes timing statistics :

- User CPU time
  - It's the time spent executing the timed command,
- System CPU time
  - It's the time spent executing system calls on behalf of your program,
- Real "wallclock" time = total time
  - It's the total elapsed time taken by the job.

ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

CSCS
Swiss National Supercomputing Centre

# Profiling : Inclusive vs. Exclusive

```
int main( )
{ /* takes 100 secs */
  f1(); /* takes 20 secs */
 /* other work */
  f2(); /* takes 50 secs */
  f1(); /* takes 20 secs */

 /* other work */
}

/* similar for other metrics, such
as hardware performance counters,
etc. */
```
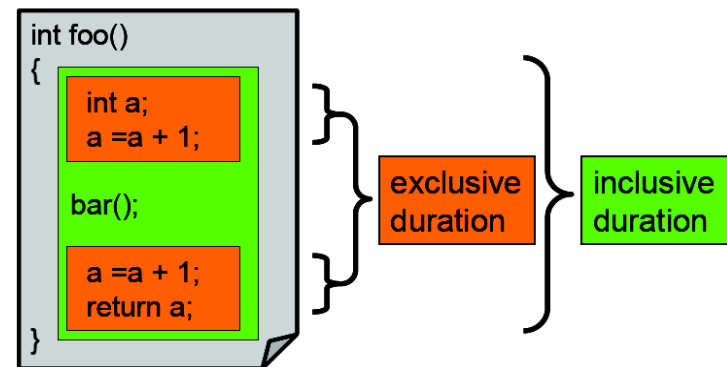
- **Inclusive** time for **main**

  - 100 secs

- **Exclusive** time for **main**

  - 100-20-50-20=10 secs

- Exclusive time sometimes called "self"

ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

CSCS
Swiss National Supercomputing Centre

# How to time MPI regions ?

The elapsed (wall-clock) time between two points in an MPI program can be computed using **MPI_Wtime**:

```
double precision :: t1, t2
t1 = MPI_Wtime()
...
t2 = MPI_Wtime()
print *, 'time is ', t2-t1
```



```
n1:/tmp $  ompi_info |grep MPI_WTIME
      MPI_WTIME support: gettimeofday
```

**MPI_Wtick** returns the resolution of MPI_Wtime in seconds, i.e the number of seconds between successive hardware clockticks.

Performance tools allow you to measure time without modifying your src code...

ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

CSCS
Swiss National Supercomputing Centre

# Step2 : Profiling/Sampling vs Tracing

Level of performance information depends on 2 types of experiments :

- Sampling experiments :
  - ➔ captures values at specified time intervals or when a specified counter overflows,
  - ➔ provides a summary of performance events and timings for the execution as a whole,

- Tracing experiments :
  - ➔ count some event such as the number of times a specific system call is executed,
  - ➔ records the chronology of events,
  - ➔ amount of data increases with the runtime, and can lead to rather large tracing files,

  A profile is sufficient to pinpoint load imbalance due to problem decomposition and/or identify the origin of excessive communication time.

  A trace is useful for detailed examination of timing issues occurring within a code.

ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

CSCS
Swiss National Supercomputing Centre

# Overview of some tools

* opensrc :
  > scalasca
  > tau
* vampir
* vendors :
  > cray, ibm, sgi, etc...

# XT-CRAYPAT : Features

- Craypat basics

  - Developped by CRAY (Luiz DeRose)

  - Multiple functionalities

    - Most consuming routines

    - Load balance across computing resources

    - Communication overhead and Cache utilization

    - FLOPS and HW counters

    - SSE instructions (Vectorization)

    - Ratio of computation vs communication

  - Integrated and easy to use

    - pat_build : utility to instrument the application

    - pat_report : utility to create performance report

    - pat_help : provides craypat infos

    - apprentice2 : graphical performance analysis tool

ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

CSCS
Swiss National Supercomputing Centre

# Step2 : pat_report (imbalance)

```
gele1:/scratch/gele/piccinal/laplace $  pat_report -O profile exe+apa+14148-16t.xf
pat_report:  Using existing file:   exe+apa+14148-16t.ap2
Processing table 1 of 1
Table 1:   Profile by Function Group and Function

 Time%  |    Time   |    Imb.   |   Imb.    |   Calls  |Group
        |           |    Time   |   Time%   |          | Function
        |           |           |           |          |   PE=HIDE
        |           |           |           |          |     Thread=HIDE

 100.0% | 5.863119  |     --    |     --    | 386624.0 | Total
|---------------------------------------------------------------------
|  39.4% | 2.309052 |     --    |     --    |      2.0 |USER           <---
||--------------------------------------------------------------------
|  39.4% | 2.308966 | 0.123713  |    5.5%   |      1.0 | laplace_
||====================================================================
|  31.1% | 1.820974 |     --    |     --    | 386420.0 |MPI            <---
||--------------------------------------------------------------------
||  29.3% | 1.716654 | 1.905820  |   57.4%   | 128200.0 |mpi_recv
||   1.1% | 0.063398 | 0.017869  |   24.0%   | 128600.0 |MPI_ISEND
||====================================================================
|  29.6% | 1.733093 |     --    |     --    |    202.0 |MPI_SYNC       <---
||--------------------------------------------------------------------
|  29.6% | 1.732687 | 1.728773  |   99.8%   |    200.0 | mpi_allreduce_(sync)
|=====================================================================
```

ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

CSCS
Swiss National Supercomputing Centre

# Detecting load imbalance



### Motivation

- Increase of systems complexity

- Increase applications scaling

### Imbalance time

- Metric based on execution time

- User functions : $t_{Maximum} - t_{Average}$

- MPI sync time : $t_{Average} - t_{Minimum}$

- Identifies computational code regions and synchronization calls that could benefit most from load balance optimization

- Estimates how much overall program time could be saved if corresponding section of code had a perfect balance

- Represents upper bound on "potential savings"

- Assumes other processes are waiting, not doing useful work while slowest member finishes

- Minimize computing resources 'waste'

ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

CSCS
Swiss National Supercomputing Centre

```
gele1:/scratch/gele/piccinal/laplace $    pat_report -O load_balance_m exe+apa+14148-16t.xf
pat_report: Using existing file:   exe+apa+14148-16t.ap2
Processing table 1 of 1
Table 1:  Load Balance with MPI Message Stats
```

| Time% | Time | MPI Msg Count | MPI Msg Bytes | Avg MPI Msg Size | Group PE=[mmm] |
|---|---|---|---|---|---|
| 100.0% | 6.171748 | 96618.7 | 2561620.0 | 26.51 | Total |
| 37.4% | 2.309054 | 0.0 | 0.0 | -- | USER |
| 39.4% | 2.432767 | 0.0 | 0.0 | -- | pe.10 |
| 37.4% | 2.310209 | 0.0 | 0.0 | -- | pe.6 |
| 33.1% | 2.039806 | 0.0 | 0.0 | -- | pe.5 |
| 34.5% | 2.129440 | 96618.7 | 2561620.0 | 26.51 | MPI |
| 65.7% | 4.055387 | 128602.0 | 2817620.0 | 21.91 | pe.5 |
| 34.1% | 2.104354 | 802.0 | 2561620.0 | 3194.04 | pe.10 |
| 6.4% | 0.392414 | 128402.0 | 1793620.0 | 13.97 | pe.0 |
| 28.1% | 1.733254 | 0.0 | 0.0 | -- | MPI_SYNC |
| 56.3% | 3.477124 | 0.0 | 0.0 | -- | pe.0 |
| 26.7% | 1.649781 | 0.0 | 0.0 | -- | pe.7 |
| 0.1% | 0.004096 | 0.0 | 0.0 | -- | pe.11 |

max, average, min

max, average, min

max, average, min

# Apprentice2 : Load imbalance view

# Apprentice2 : Call tree view

# Step2 : pat_report (src line numbers)

```
gele1:/scratch/gele/piccinal/laplace $   pat_report -O ct+src exe+apa+14148-16t.xf
pat_report:  Using existing file:   exe+apa+14148-16t.ap2
Processing table 1 of 1
Table 1:  Calltree View with Callsite Line Numbers

 Time%  |    Time   |   Calls  |Calltree
        |           |          | PE=HIDE

 100.0% | 5.864769 | 386624.0 | Total
|--------------------------------------------------------------
|   39.4% | 2.308966 |      1.0 |laplace_:Laplace_mpi.F90:line.1
|   29.6% | 1.737061 |    400.0 |laplace_:Laplace_mpi.F90:line.371
|   29.5% | 1.732687 |    200.0 | mpi_allreduce_(sync)
|   28.9% | 1.695851 |    200.0 |laplace_:Laplace_mpi.F90:line.318
|        |           |          | mpi_recv
|    1.6% | 0.092724 | 256000.0 |laplace_:Laplace_mpi.F90:line.326
|    1.0% | 0.058471 | 128000.0 | MPI_ISEND
|==============================================================
```

ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

CSCS
Swiss National Supercomputing Centre

```
gele1:/scratch/gele/piccinal/laplace $   pat_report -O mpi_callers exe+apa+14148-16t.xf
pat_report:  Using existing file:    exe+apa+14148-16t.ap2
Processing table 1 of 1
Table 1:  MPI Message Stats by Caller
```

| MPI Msg Bytes | MPI Msg Count | MsgSz <16B Count | 256B<= MsgSz <4KB Count | 4KB<= MsgSz <64KB Count | Function Caller PE=[mmm] |
|---|---|---|---|---|---|
| 2561620.0 | 96618.7 | 96202.0 | 266.7 | 150.0 | Total |
| 2560000.0 | 96416.7 | 96000.0 | 266.7 | 150.0 | MPI_ISEND |
| 2560000.0 | 96416.7 | 96000.0 | 266.7 | 150.0 | laplace_ |
| 3\|\| 3584000.0 | 128600.0 | 128000.0 | 400.0 | 200.0 | pe.4 |
| 3\|\| 2560000.0 | 128400.0 | 128000.0 | 400.0 | 0.0 | pe.1 |
| 3\|\| 1792000.0 | 400.0 | 0.0 | 200.0 | 200.0 | pe.11 |
| 1600.0 | 200.0 | 200.0 | 0.0 | 0.0 | MPI_ALLREDUCE |
| 1600.0 | 200.0 | 200.0 | 0.0 | 0.0 | laplace_ |
| 3\|\| 1600.0 | 200.0 | 200.0 | 0.0 | 0.0 | pe.0 |
| 3\|\| 1600.0 | 200.0 | 200.0 | 0.0 | 0.0 | pe.6 |
| 3\|\| 1600.0 | 200.0 | 200.0 | 0.0 | 0.0 | pe.11 |

Average of 12

```
gele1:/scratch/gele/piccinal/laplace $   pat_report -O mpi_dest_bytes exe+apa+14148-16t.)
pat_report:  Using existing file:    exe+apa+14148-16t.ap2
Processing table 1 of 1
Table 1:  MPI Sent Message Stats by Distance
```

| Sent Msg Total Bytes% | Sent Msg Total Bytes | Sent Msg Count | MsgSz <16B Count | 256B<= MsgSz <4KB Count | 4KB<= MsgSz <64KB Count | Sent Distance |
|---|---|---|---|---|---|---|
| 100.0% | 30720000 | 1157000 | 1152000.0 | 3200.0 | 1800.0 | Total |
| 60.0% | 18432000 | 1153800 | 1152000.0 | 0.0 | 1800.0 | 3 |
| 40.0% | 12288000 | 3200 | 0.0 | 3200.0 | 0.0 | 1 |

Sum of 12

ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

CSCS
Swiss National Supercomputing Centre

# Step3 : Memory usage (INTEL)

```
make clean ; make FFLAGS="-O3 -w -D_HWMEM"
```

```
t=5.3(sec) 140(MB) o_exe.mem.0002.2.2.1.1920-1920-200-1.0d-5--bind-to-core
t=6.3(sec) 53(MB) o_exe.mem.0006.6.6.1.1920-1920-200-1.0d-5--bind-to-core
t=4.2(sec) 32(MB) o_exe.mem.0012.12.12.1.1920-1920-200-1.0d-5--bind-to-core
t=4.8(sec) 51(MB) o_exe.mem.0012.12.6.1.1920-1920-200-1.0d-5--bind-to-core
t=3.6(sec) 41(MB) o_exe.mem.0024.24.12.1.1920-1920-200-1.0d-5--bind-to-core
```

Domain decomposition =>
Less memory / process

```
o_exe.mem.0004.4.1.1.1920-1920-200-1.0d-5--bind-to-core t=4.87sec m=91MB
o_exe.mem.0004.4.2.1.1920-1920-200-1.0d-5--bind-to-core t=5.50sec m=90MB
o_exe.mem.0004.4.3.1.1920-1920-200-1.0d-5--bind-to-core t=7.30sec m=73MB
o_exe.mem.0004.4.4.1.1920-1920-200-1.0d-5--bind-to-core t=6.69sec m=73MB
```

<=============== FAST

<== SLOW but only 1 cnode
and less memory / process

ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

CSCS
Swiss National Supercomputing Centre

# Step3 : Memory usage (CRAY)

```
gele2:/scratch/gele/piccinal/laplace $   grep 'Table 4:' -A10 xf.txt
Table 4:  Wall Clock Time, Memory High Water Mark

  Process  |  Process  |PE=[mmm]
     Time  |    HiMem  |
           |  (MBytes) |

  5.621835 |    45.109 |Total
  |----------------------------
  | 5.621915 |    45.094 |pe.9
  | 5.621880 |    45.074 |pe.4
  | 5.621253 |    45.199 |pe.0
```

```
0002cores, 152.160 MBytes
0004cores, 88.953 MBytes
0006cores, 67.980 MBytes
0012cores, 47.240 MBytes
0024cores, 102.260 MBytes
0072cores, 95.305 MBytes
0144cores, 93.491 MBytes
0192cores, 93.059 MBytes
```

```
make clean ; make FFLAGS="-O3 -w -D_HWMEM"
```

ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

CSCS
Swiss National Supercomputing Centre

# How to print memory usage ?

Every Linux based system supports a /proc/self/status file :

- 1 file for each active process id (mpi rank)

  ➜ VmPeak: Peak virtual memory usage
  ➜ VmSize:  Current virtual memory usage
  ➜ VmLck:   Current mlocked memory
  ➜ VmHWM:  Peak resident set size       ### high water mark ###
  ➜ VmRSS:  Resident set size            ### current memory usage ###
  ➜ VmData: Size of "data" segment
  ➜ VmStk:   Size of stack
  ➜ VmExe:   Size of "text" segment
  ➜ VmLib:    Shared library usage
  ➜ VmPTE:   Pagetable entries size
  ➜ VmSwap: Swap space used

# PAPI Features and usage

## PAPI characteristics

- University of Tennessee and vendors

    http://icl.cs.utk.edu/papi

- PAPI (Performance Application Programming Interface) is a portable interface to hardware performance counters

- Collecting low level performance metrics (e.g. clock cycles and instruction counts, memory cache misses, functional units, etc)

## PAPI features

- 2 types of performance events :

- Platform dependent native events

- Platform independent preset events

- Preset events are derived from multiple native events

- PAPI support on most HPC platforms

    - IBM, CRAY, INTEL, etc...

## Information on PAPI and AMD hardware counters

- Load the xt-papi (and xt-craypat) modulefile(s)

    - module load xt-papi ; module load xt-craypat

    - man intro_papi ; man papi_counters ; man hwpc

    - pat_help counters

## PAPI Utilities (to be executed on compute node)

- papi_avail shows which predefined events are available on the system

- papi_native_avail lists all AMD native events available on the system

- derived metrics are made from native counter event names

- papi_event_chooser reports information about the current PAPI installation and supported preset events

- papi_command_line adds named events from the command line to a PAPI EventSet and does some work with that EventSet. This serves as a handy way to see if events can be counted together, and if they give reasonable results for known work.

- papi_mem_info provides information on the memory architecture of the processor

ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

CSCS
Swiss National Supercomputing Centre

# Better know your machine : papi_avail

```
n1:/home/piccinali/trunk/laplace/src $  papi_avail
Available events and hardware information.

-------------------------------------------------------------------

PAPI Version          : 4.1.4.0
Vendor string and code   : GenuineIntel (1)
Model string and code    : Intel(R) Xeon(R) CPU        X5660  @ 2.80GHz (44)
CPU Revision          : 2.000000
CPUID Info            : Family: 6  Model: 44  Stepping: 2
CPU Megahertz         : 1600.000000
CPU Clock Megahertz   : 1600
Hdw Threads per core  : 1
Cores per Socket      : 6
NUMA Nodes            : 2
CPU's per Node        : 6
Total CPU's           : 12
Number Hardware Counters : 7
Max Multiplex Counters   : 64
-------------------------------------------------------------------
```

```
    Name        Code      Avail Deriv Description (Note)
PAPI_L1_DCM  0x80000000  Yes   No    Level 1 data cache misses
PAPI_L1_ICM  0x80000001  Yes   No    Level 1 instruction cache misses
PAPI_L1_TCM  0x80000006  Yes   Yes   Level 1 cache misses
PAPI_L1_LDM  0x80000017  Yes   No    Level 1 load misses
PAPI_L1_STM  0x80000018  Yes   No    Level 1 store misses
PAPI_L1_DCH  0x8000003e  No    No    Level 1 data cache hits
PAPI_L1_DCA  0x80000040  No    No    Level 1 data cache accesses
PAPI_L1_DCR  0x80000043  No    No    Level 1 data cache reads
PAPI_L1_DCW  0x80000046  No    No    Level 1 data cache writes
PAPI_L1_ICH  0x80000049  Yes   No    Level 1 instruction cache hits
PAPI_L1_ICA  0x8000004c  Yes   No    Level 1 instruction cache accesses
PAPI_L1_ICR  0x8000004f  Yes   No    Level 1 instruction cache reads
PAPI_L1_ICW  0x80000052  No    No    Level 1 instruction cache writes
PAPI_L1_TCH  0x80000055  No    No    Level 1 total cache hits
PAPI_L1_TCA  0x80000058  No    No    Level 1 total cache accesses
PAPI_L1_TCR  0x8000005b  No    No    Level 1 total cache reads
PAPI_L1_TCW  0x8000005e  No    No    Level 1 total cache writes
Of 107 possible events, 57 are available, of which 14 are derived.
```

ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

CSCS
Swiss National Supercomputing Centre

# How to use PAPI API ?

```fortran
#include fpapi.h
integer :: num_events, event(2), values(2)
call PAPIf_num_counters( num_events )
```

- Initialise PAPI, also shows how many hardware events are supported (num_events),

```fortran
call PAPIf_query_event(PAPI_FP_INS, retval)
event(1) = PAPI_FP_INS    # Total floating point operations,
event(2) = PAPI_TOT_CYC   # Time used
```

- check if the PAPI Preset event can be counted on the architecture,
- If the event CAN be counted, the function returns PAPI_OK, if the event CANNOT be counted, the function returns an error code,

```fortran
call PAPIf_start_counters( event, num_events, retval)
call PAPIf_read_counters(values, num_events,retval)
```

- start counting the events named in the "event" array,
- and reset the counters before entering the region of code to be measured,
- size of the event array should be no longer than the value returned by PAPIf_num_counters,

```fortran
call PAPIf_stop_counters(values,num_events,retval)
```

- stop the counters and copy the counts into the values array after the src code region to be measured.

ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

CSCS
Swiss National Supercomputing Centre

**Objectif :** Calcul de la performance crete

**Instructions :**

* Calculez le nombre de Flops theorique (rpeak) pour 1, 2, 12, 24, 36, 48 cores.

Indices :          * Ghz = Nombre de Cycles / 1 seconde

                        * Flop / Cycle = 24 (intel xeon X5660)

                        * Flops = floating point operations / 1 seconde

* Compilez et executez ~piccinali/trunk/matmul/seq/papi/matmult.F90 :

        * gfortran -O3 -I/softs/papi/4.1.4/gnu/include -D_A matmul.F90 -L/softs/papi/4.1.4/gnu/lib -lpapi -o A ; ./A

        * gfortran -O3 -I/softs/papi/4.1.4/gnu/include -D_D matmul.F90 -L/softs/papi/4.1.4/gnu/lib -lpapi -o D ; ./D

        * Quel est le nombre de flop ? Quel est la performance (Flops) par rapport a la crete ?

* Changez la taille de la matrice (PARAMETER ligne 12) : y a-t-il une difference ?

ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

CSCS
Swiss National Supercomputing Centre

flops = [ cycles / sec ] * [ flop / cycle ]
1 core  => 2.8 10^9 * 24 *  1 = 67.3 GFLOPS
2 cores => 2.8 10^9 * 24 *  2 =  134 GFLOPS
12 cores => 2.8 10^9 * 24 * 12 =  806 GFLOPS
24 cores => 2.8 10^9 * 24 * 24 =  1.6 TFLOPS
36 cores => 2.8 10^9 * 24 * 36 =  2.4 GFLOPS
48 cores => 2.8 10^9 * 24 * 48 =  3.225 TFLOPS

jki (D)
FP Instructions: 67110331
Cycles        : 107688026
=> D = 1.75 GFLOPS
        2.6% peak (2.8*67110331/107688026)/67.3*100)

ijk (A)
FP Instructions: 279047242
Cycles        : 898272636
=> A = 0.86 GFLOPS
        1.3% peak (2.8*279047242/898272636)/67.3*100)

# Step5 : I/O

n1:~/trunk/laplace/src/io
make clean
Comparez les temps d'exe de :
* make FFLAGS="-D_WITHO -O3 -w"
* make FFLAGS="-O3 -w"

/usr/bin/time -p /softs/openmpi-1.4.3/bin/mpiexec
-machinefile /softs/openmpi/h -n 4 ./exe 1920 1920 25 1.0d-5
Recommencez en variant le nombre de taches mpi...

```
CORES / TOTALTIME / CPU / IO
  2cores t=13 : 10.7 +  2.8 seconds
  4cores t=11 :  9.9 +  1.8 seconds
  6cores t=10 :  8.7 +  1.5 seconds
 12cores t= 6 :  5.5 +  1.4 seconds
 24cores t= 7 :  5.0 +  3.0 seconds
 72cores t=12 :  3.7 +  8.7 seconds
144cores t=19 :  4.2 + 15.7 seconds
192cores t=21 :  4.0 + 17.3 seconds
```

ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

CSCS
Swiss National Supercomputing Centre

The Total value for each data item is the sum of the File Name values.
The File Name value for Write B/Call is the avg of the PE values.
The File Name value for each of Write MB, Writes, Write Time, Write Rate
  MB/sec is the sum of the PE values.
  (To specify different aggregations, see:  pat_help report options s1)

This table shows only lines with Writes > 0.

Table 6:  File Output Stats by Filename

| Write Time | Write MB | Write Rate MB/sec | Writes | Write B/Call | File Name PE='HIDE' |
|---|---|---|---|---|---|
| 0.515503 | 43.457216 | 84.300636 | 2467865.000000 | 18.46 | Total |
| 0.044181 | 3.613281 | 81.783283 | 204800.000000 | 18.50 | Potentiel2D_008.dat |
| 0.043875 | 3.613281 | 82.354246 | 204800.000000 | 18.50 | Potentiel2D_006.dat |
| 0.043224 | 3.613281 | 83.593739 | 204800.000000 | 18.50 | Potentiel2D_001.dat |
| 0.043046 | 3.613281 | 83.939147 | 204800.000000 | 18.50 | Potentiel2D_011.dat |
| 0.042989 | 3.613281 | 84.051665 | 204800.000000 | 18.50 | Potentiel2D_010.dat |
| 0.042737 | 3.613281 | 84.547593 | 204800.000000 | 18.50 | Potentiel2D_000.dat |
| 0.042590 | 3.613281 | 84.838995 | 204800.000000 | 18.50 | Potentiel2D_007.dat |
| 0.042543 | 3.613281 | 84.932644 | 204800.000000 | 18.50 | Potentiel2D_005.dat |
| 0.042341 | 3.613281 | 85.337258 | 204800.000000 | 18.50 | Potentiel2D_004.dat |
| 0.042326 | 3.613281 | 85.368827 | 204800.000000 | 18.50 | Potentiel2D_009.dat |
| 0.042148 | 3.613281 | 85.729068 | 204800.000000 | 18.50 | Potentiel2D_003.dat |
| 0.041697 | 3.613281 | 86.656569 | 204800.000000 | 18.50 | Potentiel2D_002.dat |
| 0.000423 | 0.024413 | 57.667058 | 2559.000000 | 10.00 | Potentiel1D_004.dat |
| 0.000423 | 0.024413 | 57.669839 | 2559.000000 | 10.00 | Potentiel1D_007.dat |
| 0.000422 | 0.024413 | 57.903052 | 2559.000000 | 10.00 | Potentiel1D_001.dat |
| 0.000417 | 0.024413 | 58.489748 | 2559.000000 | 10.00 | Potentiel1D_010.dat |
| 0.000121 | 0.000189 | 1.558470 | 29.000000 | 6.83 | stdout |

ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

CSCS
Swiss National Supercomputing Centre

# Scalasca :   Features

- Scalasca is an open-source toolset that can be used to analyze the performance behavior of parallel applications :

  - Developed by the Jülich Supercomputing Centre, Germany

  - Designed for use on large-scale systems (Cray XT, IBM) and also for medium scale,

  - Supports message passing and threads (MPI, OpenMP, SHMEM, CAF)

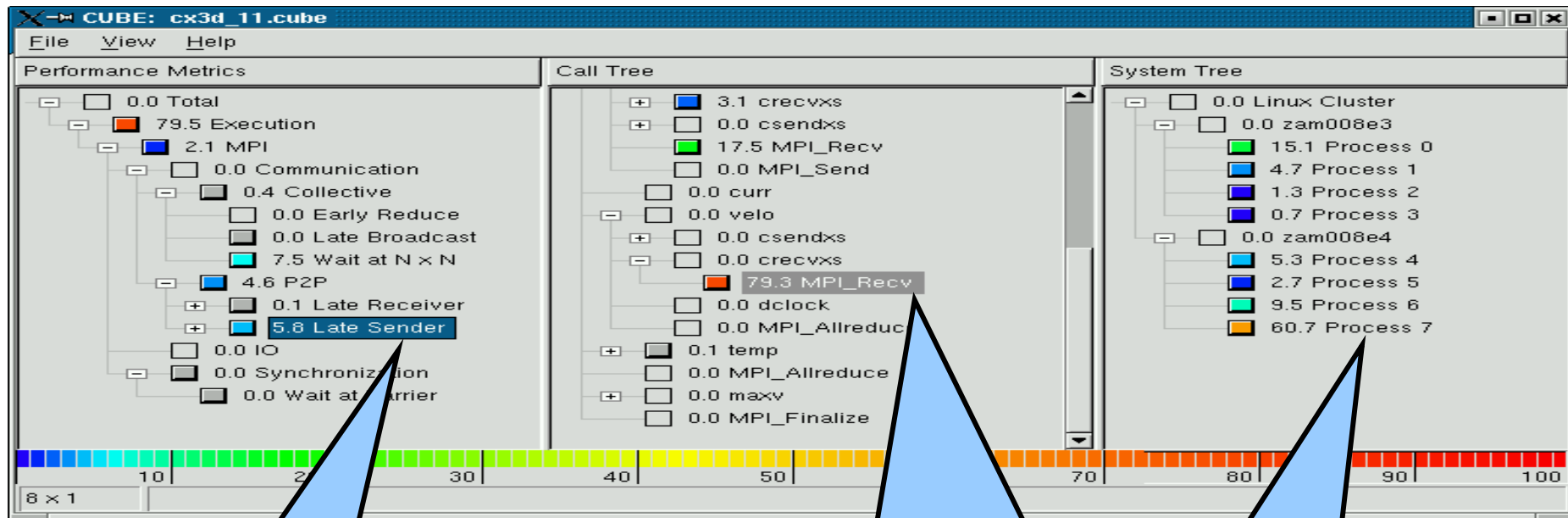  - Provides a low-overhead performance summary of CPU, parallel and memory usage.

SOURCE CODE

MULTI-LEVEL INSTRUMENTATION

EXECUTABLE

EXECUTION ON PARALLEL MACHINE

EVENT TRACE

PARALLEL PATTERN ANALYSIS

RUNTIME SUMMARY

WAIT STATE REPORT

RESULT PRESENTATION

ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

CSCS
Swiss National Supercomputing Centre

# Using SCALASCA

- Profiling your code with SCALASCA involves the following steps :

  - Load the SCALASCA module

    - ➔ module load *scalasca*

    - – Using the CRAY wrappers (ftn, cc, CC) will automatically link your MPI code with the ipm library : no need to modify your compilation line, the wrapper will do it !

  - Instrument your code (cannot handle compiling and linking in one step)

    - ➔ scalasca -instrument *mpif90 -c test.f90*

    - – scalasca -instrument *mpif90 -o test.exe test.o*

  - Execute the resultant executable by submitting your PBS job

    - ➔ scalasca -analyse *mpirun -np 128 test.exe*

  - Visualize the results

    - ➔ scalasca -examine epik_a

    - ➔ cube3 epik_a/epitome.cube        (java GUI)

      - can be done on your workstation !

ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

CSCS
Swiss National Supercomputing Centre

# Using CUBE

- CUBE is the tool to interactively examine the parallel application execution analysis reports

  → Easy to use and portable

  → Uses a relatively simple XML input file structure, simple operations can be performed om multiple inputs (diff, merge)

  → Displays tree-based views of collected infos (calls, walltime, etc...),

  → Is compatible with IPM/TAU log files



**Which type of problem?**

**Where in the source code? Which call path?**

**Which process / thread ?**

ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

CSCS
Swiss National Supercomputing Centre

# CUBE example output : Laplace (scalasca)



- **Topology display :**
- Shows distribution of pattern over HW topology

# References

- SCALASCA website (M. Geimer, B. Wylie, F. Wolf)

  - http://www.scalasca.org

- CSCS website

  - **User Entry Point >**

    - **Software and Programming Environment >**

      - **Debugging and Performance Analysis > Performance > Scalasca**

# TAU Features and usage

- TAU characteristics

  → Flat MPI and callgraph profiling

  → Hardware counter data collection

  → OpenMP & pthread profiling

  → MPI tracing

  → Memory profiling

- TAU features

  → Auto-instrumentation utility (PDT)

  → Custom configurations

  → Interoperability with other tools

Profiling your code with TAU involves the following steps

- Load the TAU modulefile

  → module load tau ; module help tau

- Recompile your code using the TAU compiler scripts

  → tau_f90.sh -c file.f90

  → tau_f90.sh -o exe file.o

  → A simple relink will skip the autoinstrumentation step

- Execute the resultant executable by submitting your PBS job

  → mpirun -n 12 exe

- Visualize the results

  → pprof profile.*

  → paraprof &

  → can be done on your workstation !

ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

CSCS
Swiss National Supercomputing Centre

# TAU

```
n1:/home/piccinali/trunk/laplace/src $    make clean ; make FCOMPILER=tau_f90.sh FFLAGS="-O3 -w" PROGRAM=exe.tau
```

```
module help tau

----------- Module Specific Help for 'tau/2.20.3' -----------------


            modulefile : tau/2.20.3 help
            tau Version 2.20.3 :
            CSCS Users Documentation : http://www.cscs.ch
            http://www.cs.uoregon.edu/research/tau

            This version was compiled with gnu
            See : /softs/tau/2.20.3/gnu

            Usage :
              module load tau/2.20.3   will set $TAU_MAKEFILE to Makefile.tau-papi-mpi-pdt,
             currently   $TAU_MAKEFILE = Makefile.tau-papi-mpi-pdt
             You can also choose to set TAU_MAKEFILE to other Makefile configurations
             available in /softs/tau/2.20.3/gnu/x86_64/lib/ :
             export TAU_MAKEFILE=/softs/tau/2.20.3/gnu/x86_64/lib/Makefile.tau-mpi-pdt
             export TAU_MAKEFILE=/softs/tau/2.20.3/gnu/x86_64/lib/Makefile.tau-phase-papi-mpi-pdt
             export TAU_MAKEFILE=/softs/tau/2.20.3/gnu/x86_64/lib/Makefile.tau-param-mpi-pdt
             export TAU_MAKEFILE=/softs/tau/2.20.3/gnu/x86_64/lib/Makefile.tau-depthlimit-mpi-pdt
             export TAU_MAKEFILE=/softs/tau/2.20.3/gnu/x86_64/lib/Makefile.tau-papi-pdt
             export TAU_MAKEFILE=/softs/tau/2.20.3/gnu/x86_64/lib/Makefile.tau-papi-mpi-pdt
             export TAU_MAKEFILE=/softs/tau/2.20.3/gnu/x86_64/lib/Makefile.tau-papi-pthread-pdt
             export TAU_MAKEFILE=/softs/tau/2.20.3/gnu/x86_64/lib/Makefile.tau-pthread-pdt
             export TAU_MAKEFILE=/softs/tau/2.20.3/gnu/x86_64/lib/Makefile.tau-pdt

            Use tau_cc.sh or tau_cxx.sh or  tau_f90.sh to compile your code,
            Run your batch job as usual,
            Use   pprof profile.0.0.0 OR  paraprof (= GUI) to visualize your results.
```

```
sbatch.sh ../exe.tau 6 6 1 "1920 1920 200 1.0d-5" "" "-bind-to-core"
```

```
n1:/home/piccinali/trunk/laplace/src/tau $    pprof -s
Reading Profile files in profile.*

FUNCTION SUMMARY (total):
---------------------------------------------------------------------------
%Time    Exclusive     Inclusive      #Call      #Subrs  Inclusive Name
             msec    total msec                          usec/call
---------------------------------------------------------------------------
100.0       12,459        38,416          6 2.68034E+06    6402682 LAPLACE
 18.9        7,244         7,244       1200           0       6037 MPI_Allreduce()
 18.6        7,128         7,128     576600           0         12 MPI_Recv()
 16.1        6,170         6,170          6           0    1028374 MPI_Init()
 12.5        4,789         4,789     600006           0          8 SOLVER [THROTTLED]
  0.6          219           219     300003           0          1 MPI_Recv() [THROTTLED]
```

```
FUNCTION SUMMARY (mean):
---------------------------------------------------------------------------
%Time    Exclusive     Inclusive      #Call      #Subrs  Inclusive Name
             msec    total msec                          usec/call
---------------------------------------------------------------------------
100.0        2,076         6,402          1      446722    6402682 LAPLACE
 18.9        1,207         1,207        200           0       6037 MPI_Allreduce()
 18.6        1,188         1,188      96100           0         12 MPI_Recv()
 16.1        1,028         1,028          1           0    1028374 MPI_Init()
 12.5          798           798     100001           0          8 SOLVER [THROTTLED]
```

# Outputs (TAU)

# Functions profile (tau)

## What routines account for the most time ?
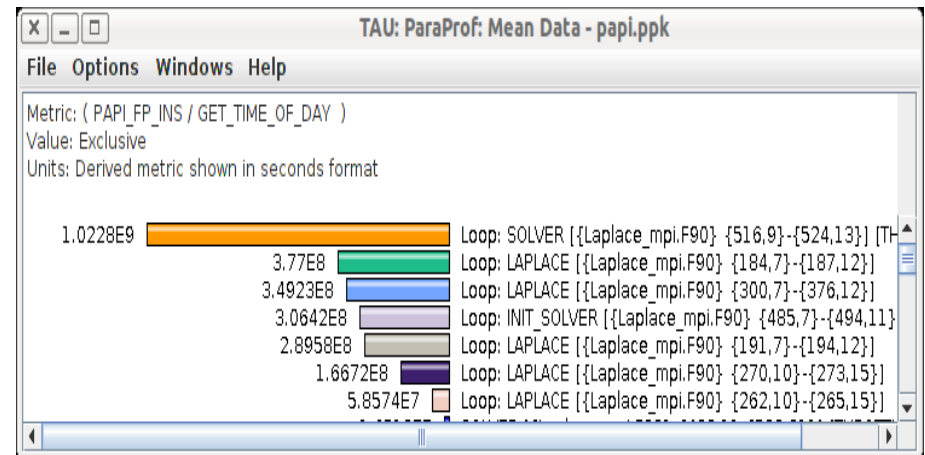
# PAPI (tau)



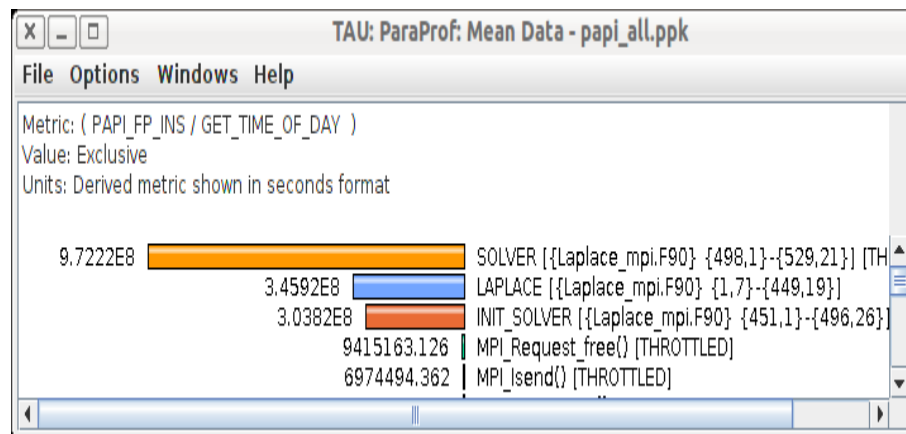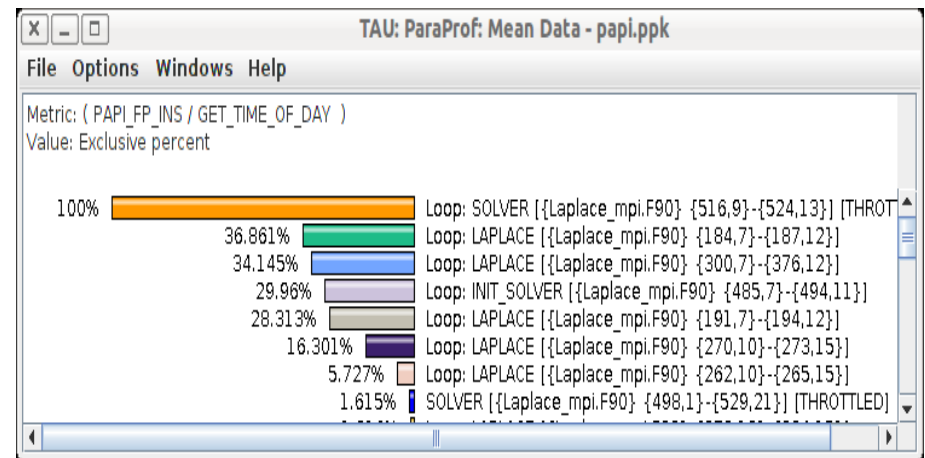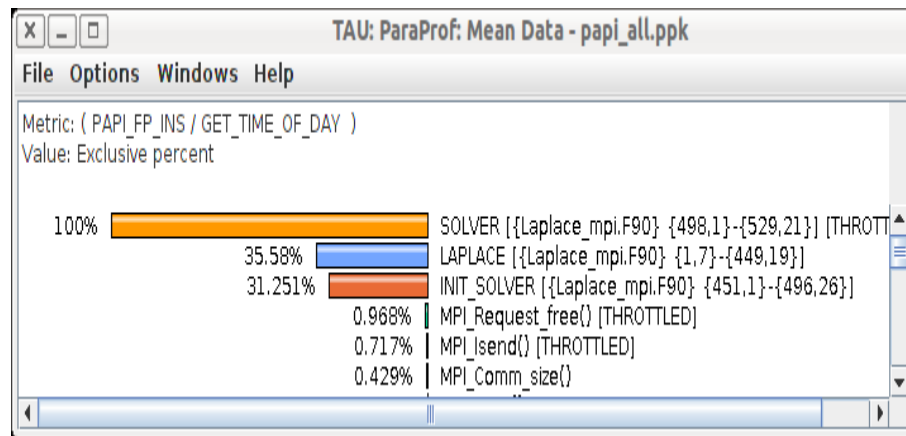ANGD 10/2011

# Loops profiling (tau)

```
export TAU_OPTIONS="-optTauSelectFile=select.tau"
```

```
export TAU_METRICS="GET_TIME_OF_DAY:PAPI_FP_INS"
```

```
paraprof --pack papi.ppk
```
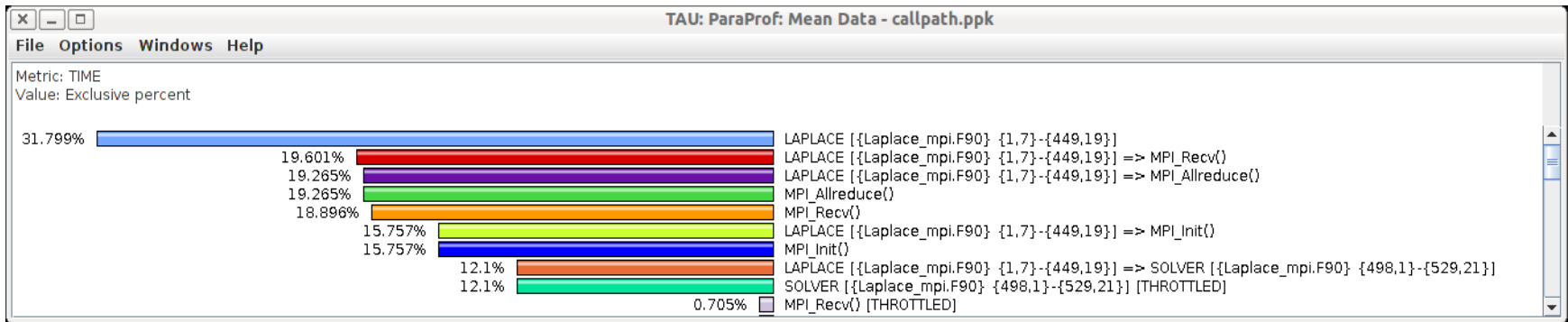
```
paraprof papi.ppk
```

```
BEGIN_INSTRUMENT_SECTION
loops routine="#"
END_INSTRUMENT_SECTION
```
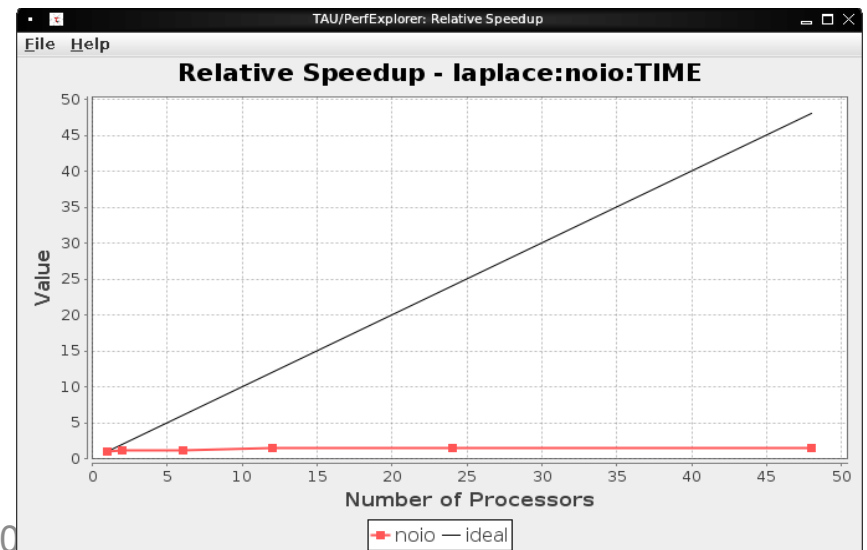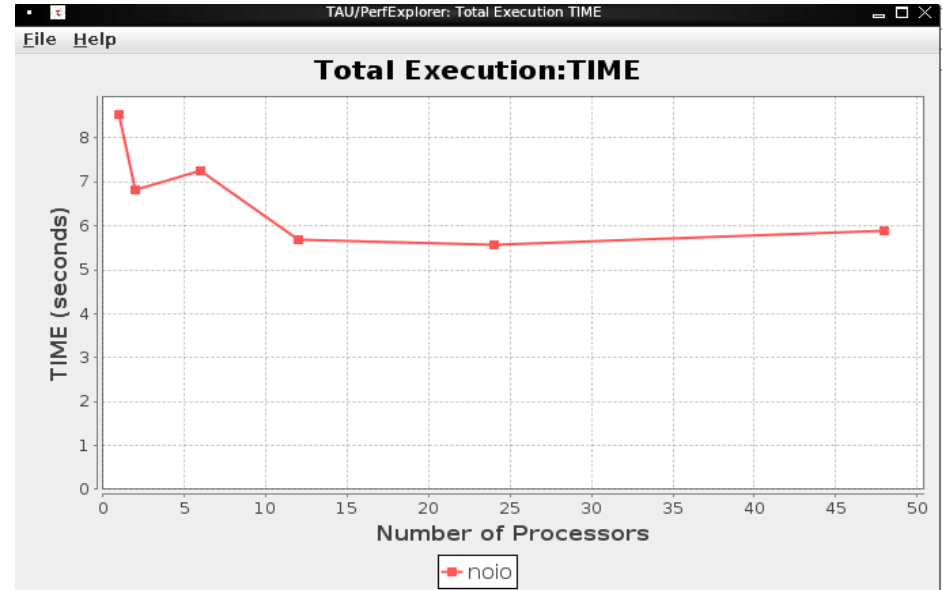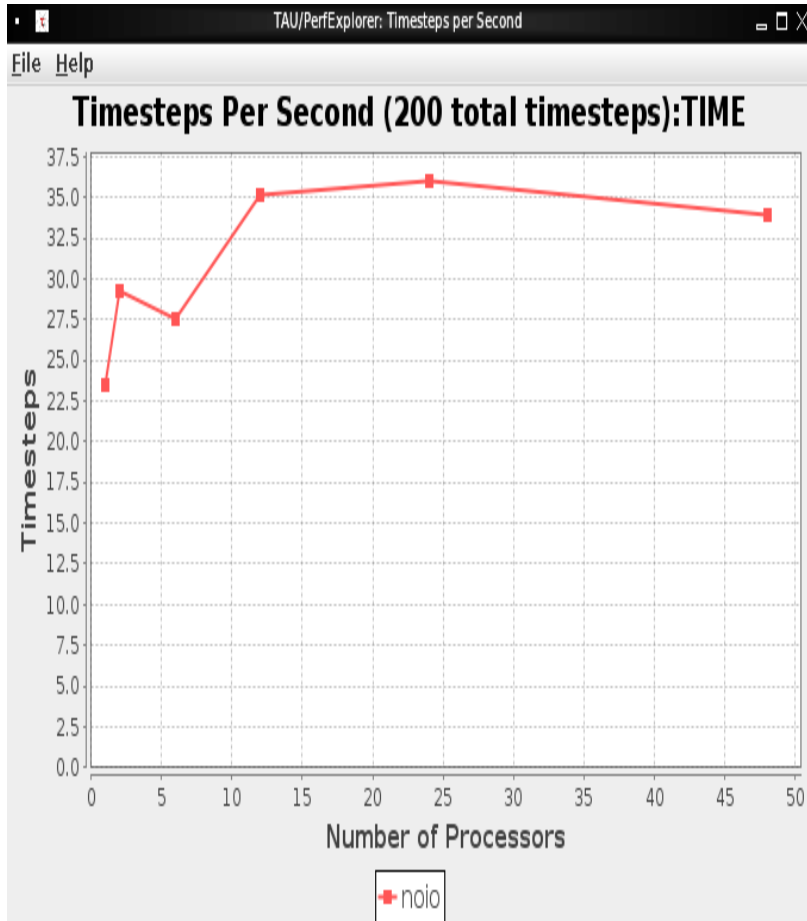
```
export TAU_CALLPATH=1
```

```
paraprof --pack callpath.ppk
```



```
n1:~/trunk/laplace/src/tau > pprof -s
Reading Profile files in profile.*

FUNCTION SUMMARY (total):
---------------------------------------------------------------------
%Time    Exclusive    Inclusive      #Call      #Subrs  Inclusive Name
             msec   total msec                          usec/call
---------------------------------------------------------------------
100.0      12,451       39,157          6 2.68034E+06   6526308 LAPLACE
 19.6       7,675        7,675     876603            0         9 LAPLACE => MPI_Recv()
 19.3       7,543        7,543       1200            0      6287 LAPLACE => MPI_Allreduce()
 19.3       7,543        7,543       1200            0      6287 MPI_Allreduce()
 18.9       7,399        7,399     576600            0        13 MPI_Recv()
 15.8       6,170        6,170          6            0   1028347 LAPLACE => MPI_Init()
 15.8       6,170        6,170          6            0   1028347 MPI_Init()
 12.1       4,738        4,738     600006            0         8 LAPLACE => SOLVER
 12.1       4,738        4,738     600006            0         8 SOLVER [THROTTLED]
  0.7         276          276     300003            0         1 MPI_Recv() [THROTTLED]
```

ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

CSCS
Swiss National Supercomputing Centre

# PARAPROF 3D (tau)

ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

CSCS
Swiss National Supercomputing Centre

# PerfExplorer (tau)

# Selective instrumentation/profiling (tau)

```
export TAU_OPTIONS="-optTauSelectFile=select.tau"
```

```
BEGIN_INSTRUMENT_SECTION
loops routine="#"
END_INSTRUMENT_SECTION
```

```
n1:/home/piccinali/trunk/laplace/src $    make clean ; make FCOMPILER=tau_f90.sh FFLAGS="-O3 -w"
```

```
n1:/home/piccinali/trunk/laplace/src/tau $    pprof -s
Reading Profile files in profile.*

FUNCTION SUMMARY (total):
---------------------------------------------------------------------
%Time    Exclusive    Inclusive      #Call      #Subrs  Inclusive Name
             msec   total msec                          usec/call
---------------------------------------------------------------------
100.0          25       38,646          6         156   6441096 LAPLACE
 83.4      12,414       32,213          6   2.68022E+06  5368984 Loop: LAPLACE [{Laplace_mpi.F90} {300,7}-{376,12}]
 18.9       7,304        7,304       1200           0      6087 MPI_Allreduce()
 18.5       7,165        7,165     576600           0        12 MPI_Recv()
 16.0       6,172        6,172          6           0   1028826 MPI_Init()
 12.6          88        4,872     600006      600006         8 SOLVER [THROTTLED]
 12.4       4,784        4,784     600006           0         8 Loop: SOLVER [{Laplace_mpi.F90} {516,9}-{524,13}]
```
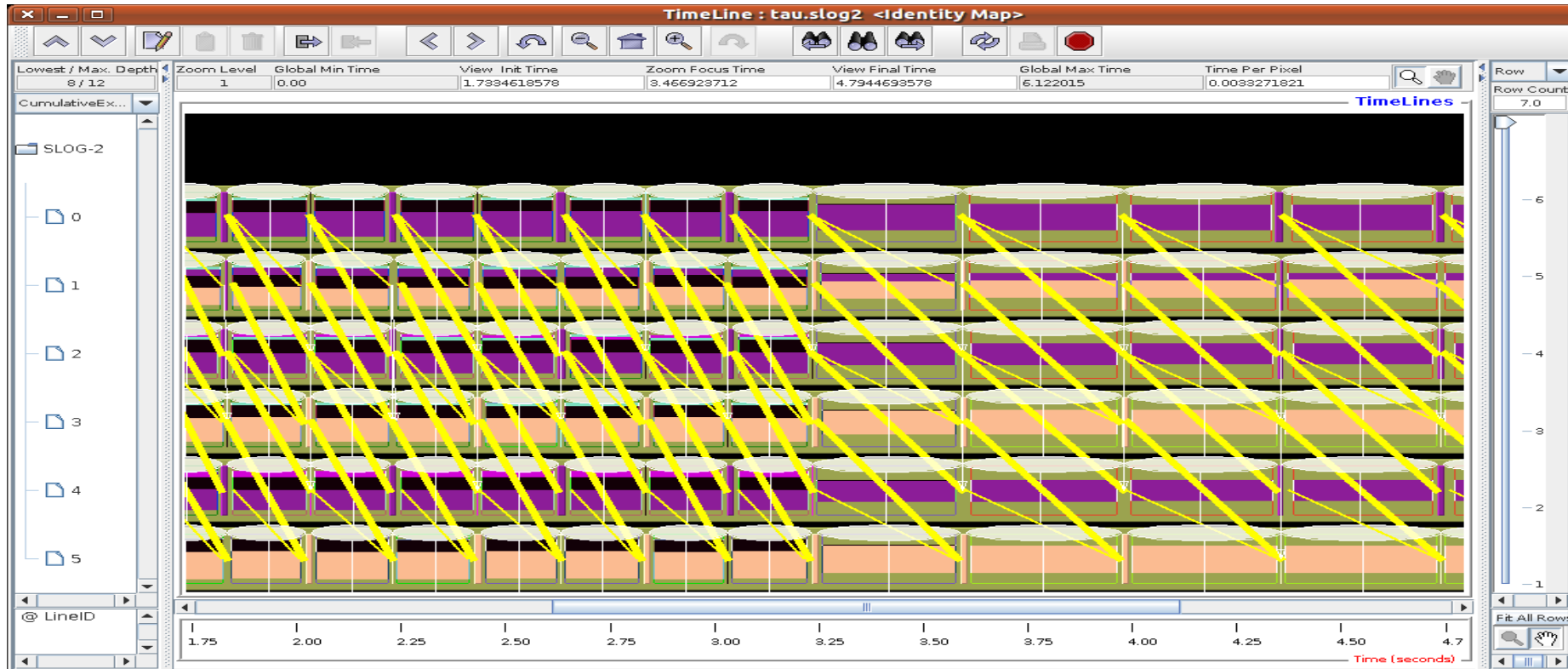
ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

CSCS
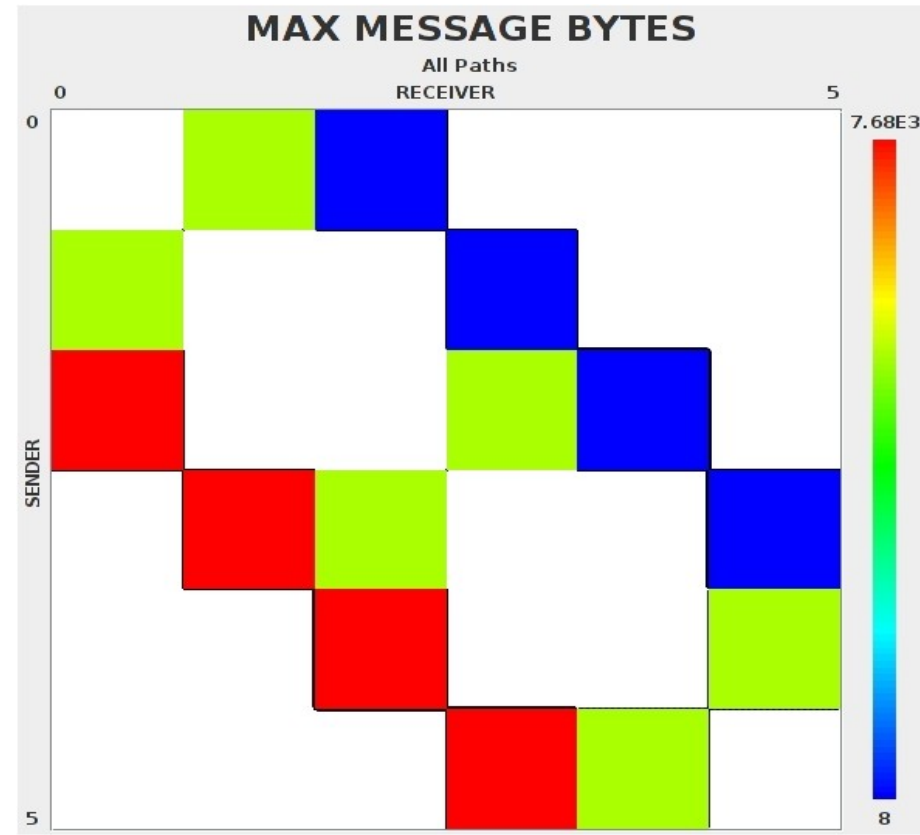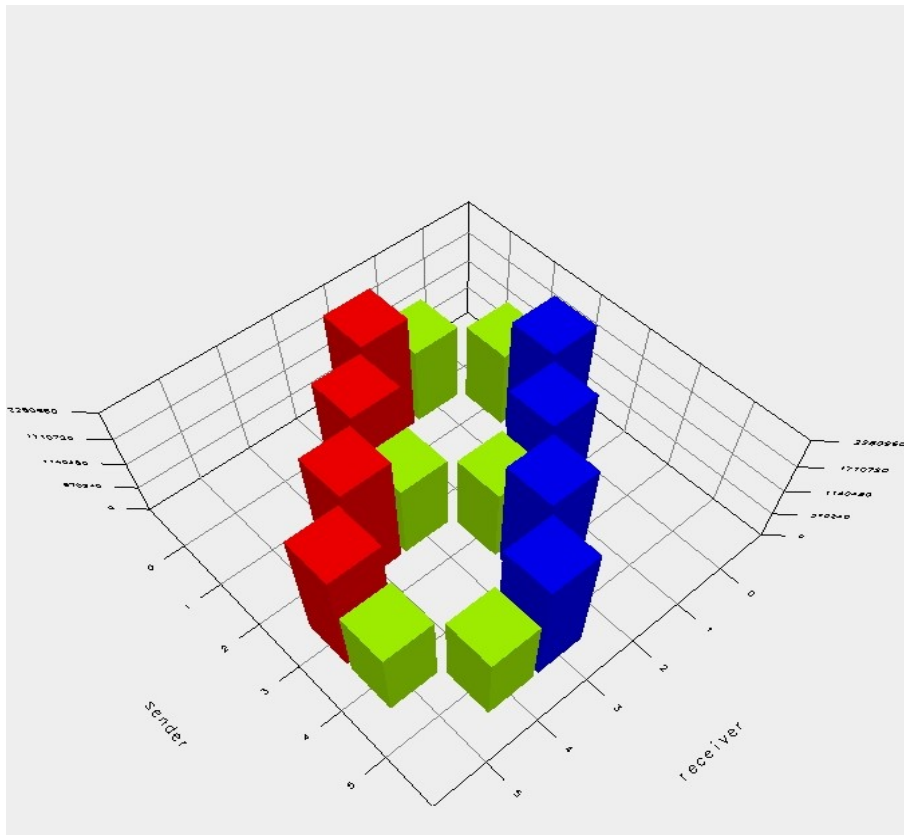Swiss National Supercomputing Centre

# Tracing (tau)

```
export TAU_TRACE=1

sbatch.sh ../exe.tau 6 6 1 "1920 1920 200 1.0d-5" "" "-bind-to-core"
tau2slog2 tau.trc tau.edf -o tau.slog2
jumpshot tau.slog2
```

ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich
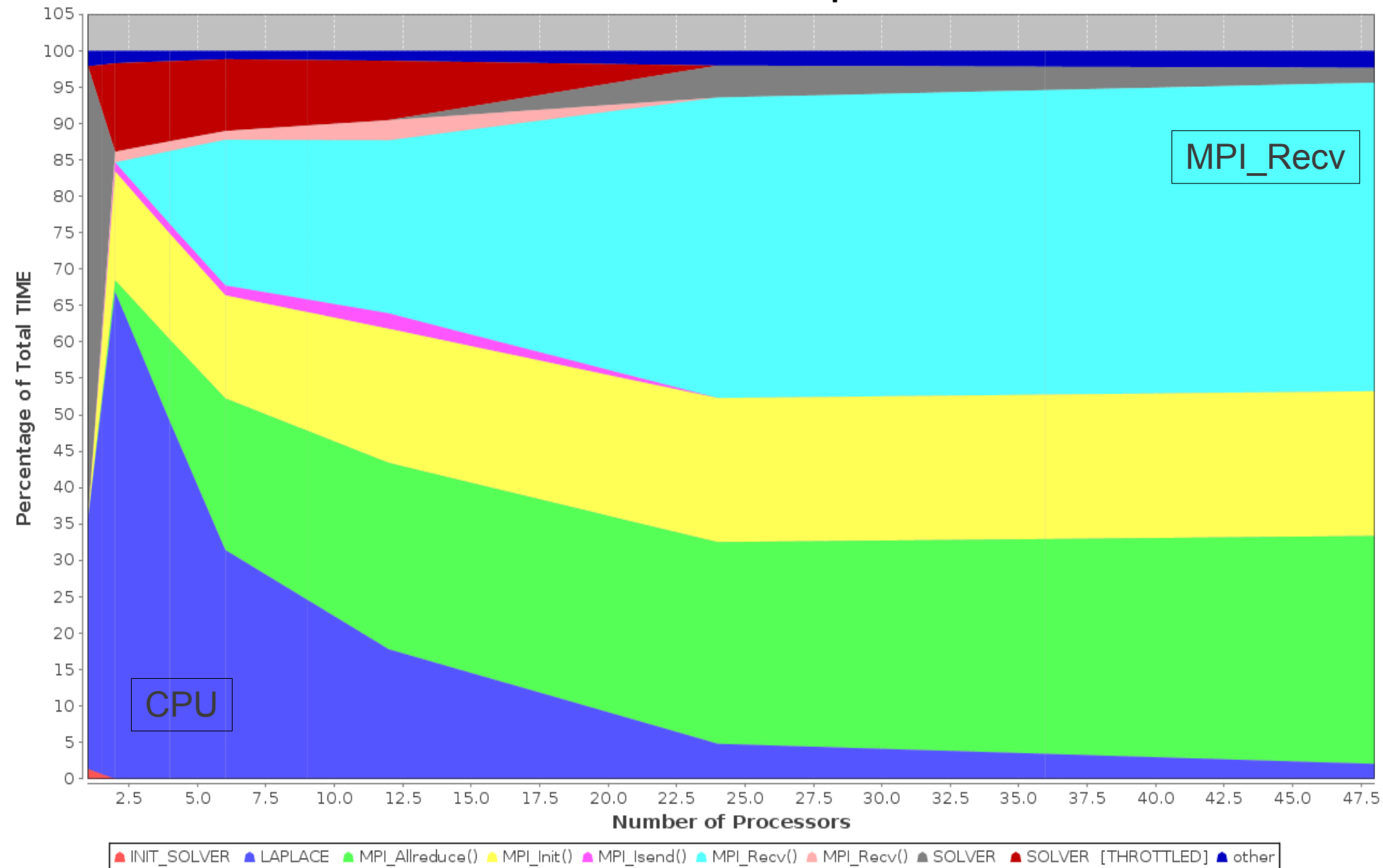
CSCS
Swiss National Supercomputing Centre

# Communication matrix (tau)

```
export TAU_TRACE=0 ; export TAU_COMM_MATRIX=1

paraprof
```

Total TIME Breakdown for laplace:noio

# References

- TAU website
  - ✓ http://www.cs.uoregon.edu/research/tau/home.php
  - ✓ Documentation
  - ✓ Tutorials
  - ✓ Downloads (offline analysis utilities)

- CSCS Rosa website
  - ✓ http://www.cscs.ch > User Entry Point
    - ✓ > Software and Programming Environment
      - ✓ > Debugging and Performance Analysis
        - ✓ > Performance
          - ✓ > TAU
  - ✓ Module setup
  - ✓ Usage information

# TAU : Perfdmf / Perfexplorer

```
n1:~ > perfdmf_configure
Configuration file NOT found...
a new configuration file will be created.

Welcome to the configuration program for PerfDMF.
This program will prompt you for some information necessary to ensure
the desired behavior for the PerfDMF tools.


You will now be prompted for new values, if desired.  The current or default
values for each prompt are shown in parenthesis.
To accept the current/default value, just press Enter/Return.

Please enter the name of this configuration.
():autrans
Please enter the database vendor (oracle, postgresql, mysql, db2 or derby).
(derby):
Please enter the JDBC jar file.
(/softs/tau/2.20.3/gnu/x86_64/lib/derby.jar):
Please enter the JDBC Driver name.
(org.apache.derby.jdbc.EmbeddedDriver):
Please enter the path to the database directory.
(/home/piccinali/.ParaProf/perfdmf):
Please enter the database username.
():piccinali
Store the database password in CLEAR TEXT in your configuration file? (y/n):y
Please enter the database password:
Please enter the database password:Please enter the PerfDMF schema file.
(/softs/tau/2.20.3/gnu/etc/dbschema.derby.txt):

Writing configuration file: /home/piccinali/.ParaProf/perfdmf.cfg.autrans

Now testing your database connection.

Database created, command: jdbc:derby:/home/piccinali/.ParaProf/perfdmf;create=true

Uploading Schema: /softs/tau/2.20.3/gnu/etc/dbschema.derby.txt
Found /softs/tau/2.20.3/gnu/etc/dbschema.derby.txt  ... Loading
Successfully uploaded schema

Database connection successful.
Configuration complete.
n1:~ >
```

```
n1:~/workspacetau/laplace/bin > perfexplorer_configure
What is the name of your PerfDMF Configuration: autrans

Configuring scripts to use the following values:
----------------------------------------------
tauroot = /softs/tau/2.20.3/gnu
architecture = x86_64
taushell = sh
targetdir = /softs/tau/2.20.3/gnu
server = localhost
configfile = /home/piccinali/.ParaProf/perfdmf.cfg.autrans
tmpdir = /tmp

TAU: installing tools in /softs/tau/2.20.3/gnu

/home/piccinali/.ParaProf/weka-3-6-1.jar not found.

Would you like to attempt to automatically download the Weka jar file? (y/n) y
Getting weka-3-6-1.zip... please be patient...
18926k bytes... done.926k bytes

/home/piccinali/.ParaProf/drools-core-3.0.6.jar not found.

Would you like to attempt to automatically download the required jar files? (y/n) y
Getting PE2_jars.tgz... please be patient...
7008k bytes... done.08k bytes
mv: `PE2_jars.tgz' and `/home/piccinali/.ParaProf/./PE2_jars.tgz' are the same file
jbossrules/drools-compiler-3.0.6.jar
jbossrules/drools-core-3.0.6.jar
jbossrules/drools-decisiontables-3.0.6.jar
jbossrules/drools-jsr94-3.0.6.jar
jbossrules/lib/antlr-3.0ea8.jar
jbossrules/lib/commons-jci-core-1.0-406301.jar
jbossrules/lib/commons-jci-eclipse-3.2.0.666.jar
jbossrules/lib/core-3.2.0.666.jar
jbossrules/lib/commons-logging-api-1.0.4.jar
jbossrules/lib/commons-lang-2.1.jar
jbossrules/lib/stringtemplate-2.3b6.jar
jbossrules/lib/antlr-2.7.6.jar
jbossrules/lib/jsr94-1.1.jar
jbossrules/lib/jxl-2.4.2.jar
jbossrules/lib/junit-3.8.1.jar

Now testing your database connection.

Configuration file found...
Parsing config file...
Cannot connect to server.
Connection String: jdbc:derby:/home/piccinali/.ParaProf/perfdmf
Exception Message: Failed to start database '/home/piccinali/.ParaProf/perfdmf', see the next exception for details.

Please make sure that your DBMS is configured correctly, and the database /home/piccinali/.ParaProf/perfdmf has been

Configuration complete!
  If you haven't already done so,
  Please add  /softs/tau/2.20.3/gnu/x86_64/bin  to your path
n1:~/workspacetau/laplace/bin >
```

ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

CSCS
Swiss National Supercomputing Centre

# TAU → ECLIPSE



```
n1:/softs/tau/src/tau-2.20.3/tools/src/eclipse > ./install_plugins.sh /softs/eclipse/3.7
Installing to /softs/eclipse/3.7/dropins
...
Eclipse plugins installed!
n1:/softs/tau/src/tau-2.20.3/tools/src/eclipse >
```

ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

CSCS
Swiss National Supercomputing Centre
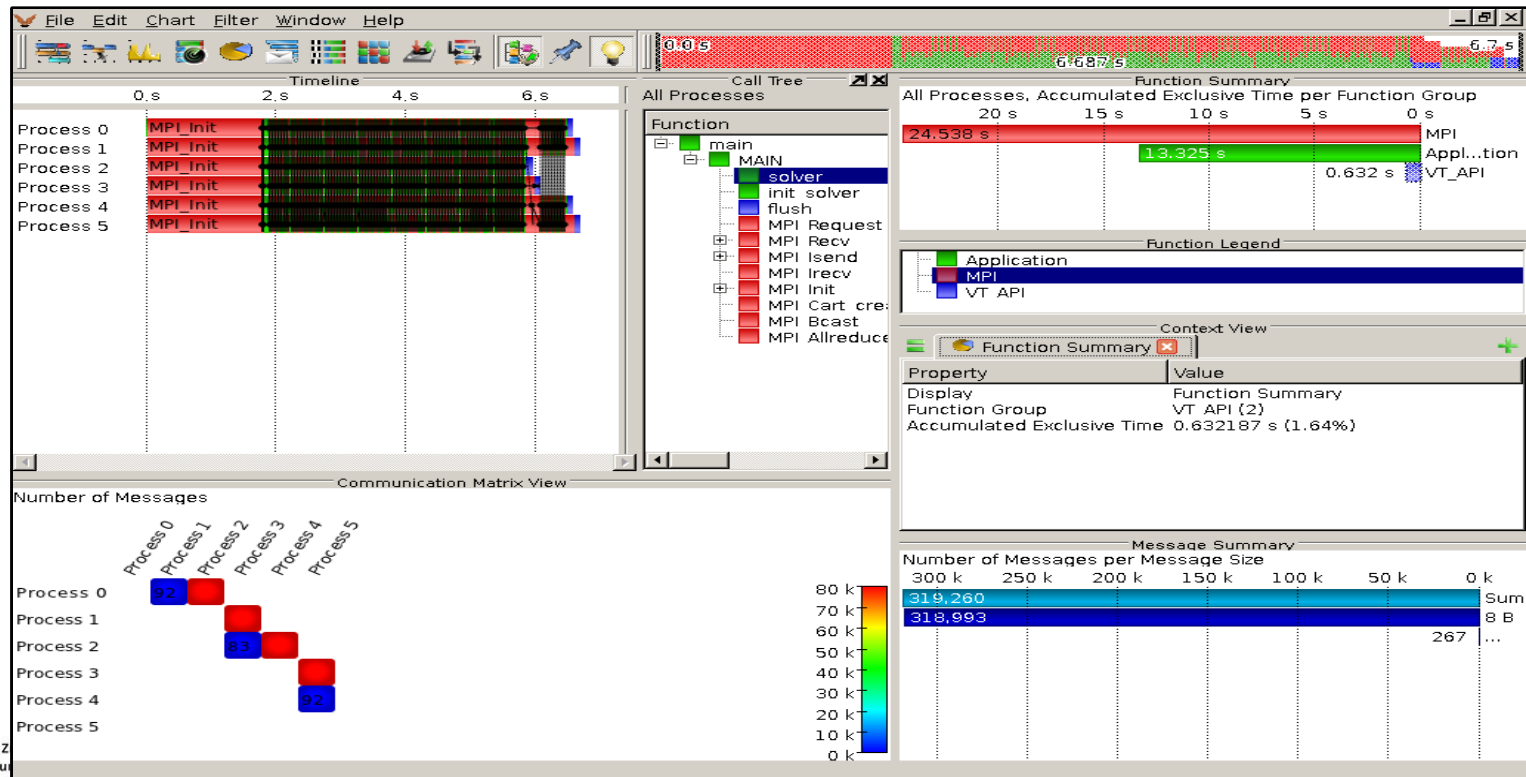
# VAMPIR (http://www.vampir.eu)

```
n1:/home/piccinali/trunk/laplace/src/vampir $    which mpif90-vt
/softs/openmpi-1.4.3/bin/mpif90-vt
```

```
n1:/home/piccinali/trunk/laplace/src $    make clean ; make FCOMPILER=mpif90-vt FFLAGS="-O3 -w" PROGRAM=exe.vampir
rm -f exe mem_ntk.o Laplace_mpi.o
mpicc -O3 -w -c mem_ntk.c
mpif90-vt -O3 -w -c Laplace_mpi.F90
mpif90-vt -O3 -w -o exe.vampir mem_ntk.o Laplace_mpi.o
```
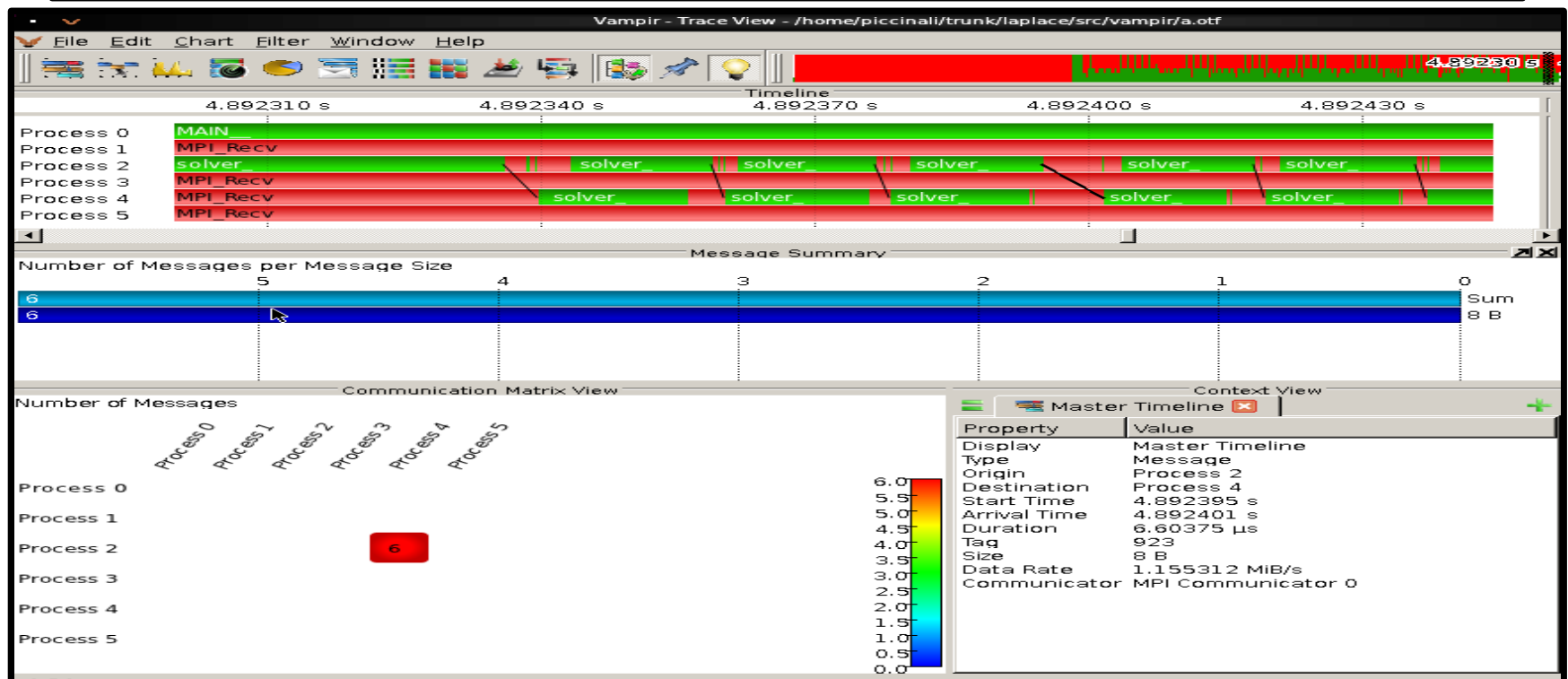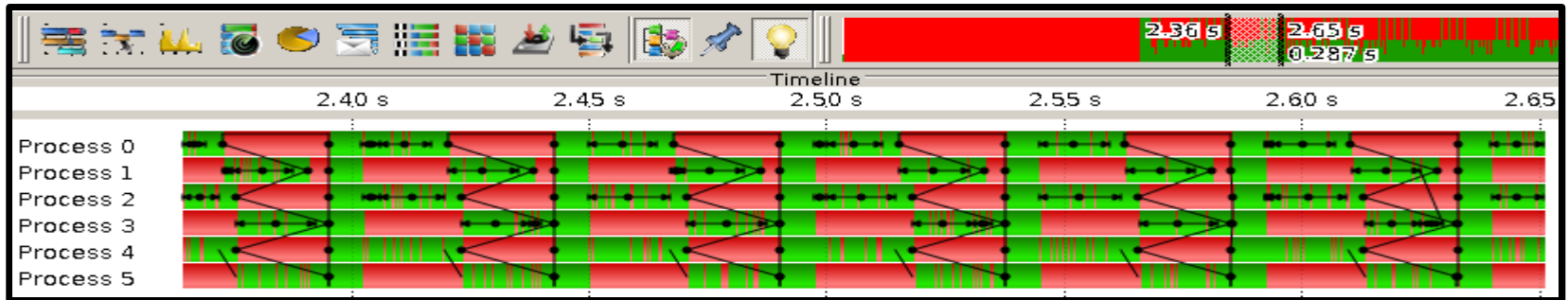
```
sbatch.sh ../exe.vampir 6 6 1 "1920 1920 200 1.0d-5" "" "-bind-to-core"
```

a.0.def
a.1.events
a.2.events
a.3.events
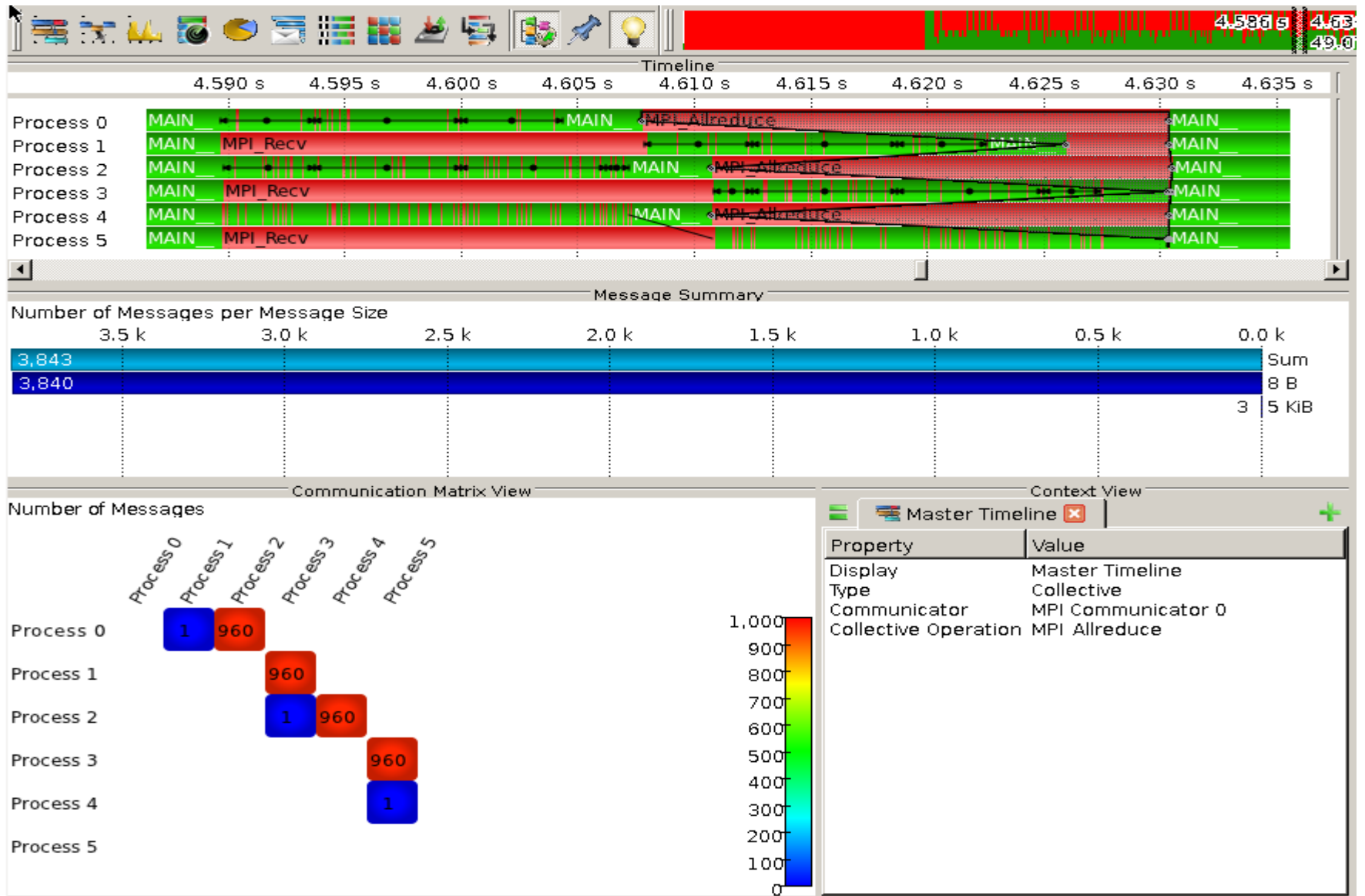a.4.events
a.5.events
a.otf
a.6.events

module load vampir

# VAMPIR

ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

CSCS
Swiss National Supercomputing Centre

# VAMPIR

# Summary



**Performance report**

**Load balance issue ?**

**Computation > 50% runtime ?**

**Communication > 50% runtime ?**

**I/O > 50% runtime ?**

Is computation causing the imbalance ?

Is communication causing the imbalance ?

**Collectives calls ? P2P calls ?**

File sizes ? (stripe) Read/Write sizes ? MPI/IO ?

• Is domain decomposition appropriate ?
• Would RANK_REORDER help ?
• Are MPI_RECV preposted ?
• OPENMP may help to spread the workload

• MPI messages sizes ?
• Are MPI_RECV preposted ?
• OPENMP may help to spread the workload

• Is application vectorized ?
• What is MEMORY cache utilization ? TLB misses ?
• What optimized libs routines are being used ?
• OPENMP may help to spread the workload

ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

ANGD 10/2011

CSCS
Swiss National Supercomputing Centre