# Scalable OpenMP Programming

*Dieter an Mey*

*Center for Computing and Communication*
*RWTH Aachen University, Germany*

www.rz.rwth-aachen.de

anmey@rz.rwth-aachen.de

# Overview

- **Why OpenMP**
- **Short OpenMP Introduction**
- **OpenMP on NUMA Machines**
- **OpenMP on Clusters**
- **Conclusion**

**Scalable OpenMP Programming – D.an Mey**

**RWTH**
Center for
Computing and Communication

# Overview

- **Why OpenMP**
- **Short OpenMP Introduction**
- **OpenMP on NUMA Machines**
- **OpenMP on Clusters**
- **Conclusion**

**Scalable OpenMP Programming – D.an Mey**

**RWTH**
Center for
Computing and Communication

# RWTH Aachen University: Key Figures  WT 06/07
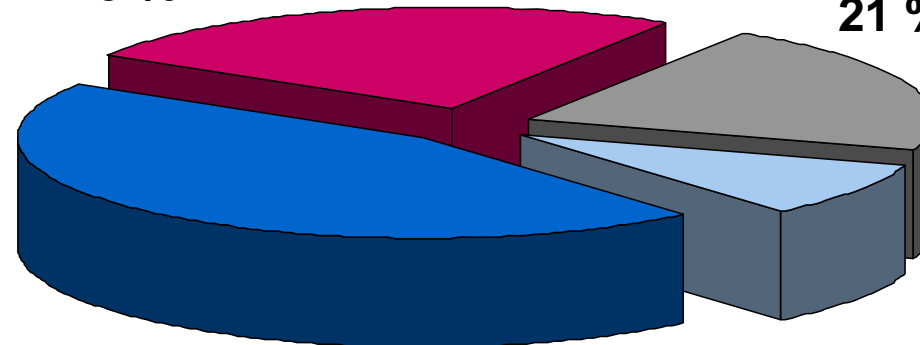
**30180 students**

**428 professorships**

**3600 academic staff**

**Natural Sciences 25 %**

**Humanities, Social Sciences and Economics 21 %**

**260 institutes**

**Engineering 45 %**

**Medicine 9 %**

**Scalable OpenMP Programming – D.an Mey**

RWTH Center for Computing and Communication
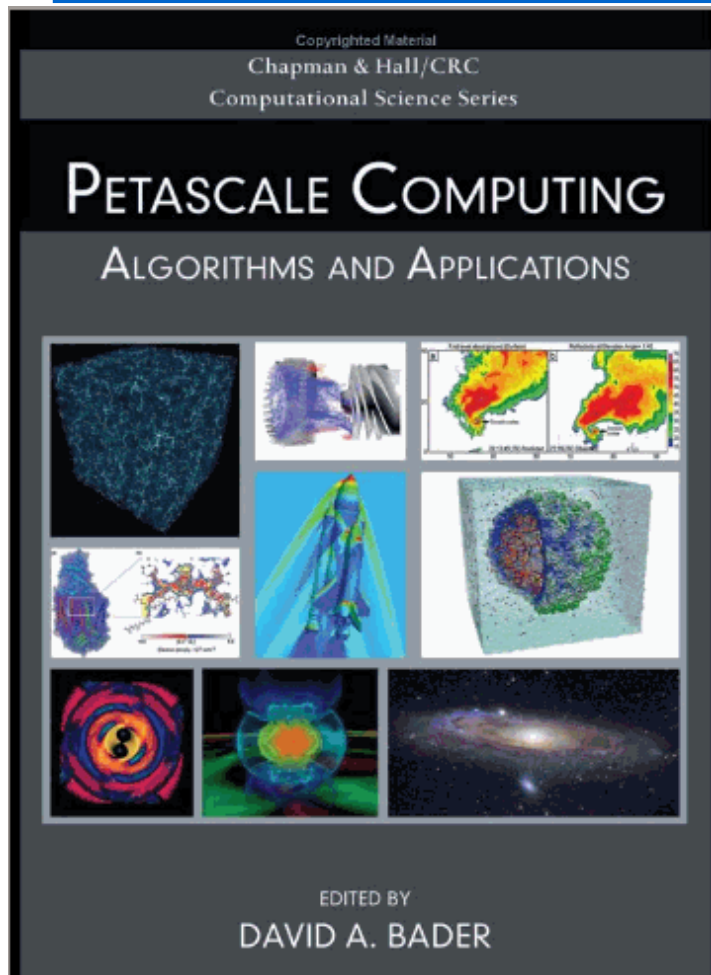
# Why OpenMP?

- Large codes mainly in C++ and Fortran and some C

- Software lifetime measured in decades

- MPI is there to stay on clusters

  - Cannot always be applied easily – if at all

  - Scalability may be limited due to underlying problem (geometry etc.)

  - "MPI only" may not be appropriate for "many cores"
    => MPI + OpenMP (hybrid)

- OpenMP is the alternative and the supplement to MPI

- Scalability of OpenMP limited by current machinery

- So far scalability explored on

  - Sun Fire E25K (144 cores UltraSPARC IV)

  - Sun UltraSPARC T2 (64 threads in one "Niagara 2" chip)

  - Intel Cluster OpenMP

  - ScaleMP "Virtual SMP"

**RWTH**
**C**
**C**
**C**
Center for
Computing and Communication

# Statistics from the "First Petascale Book"

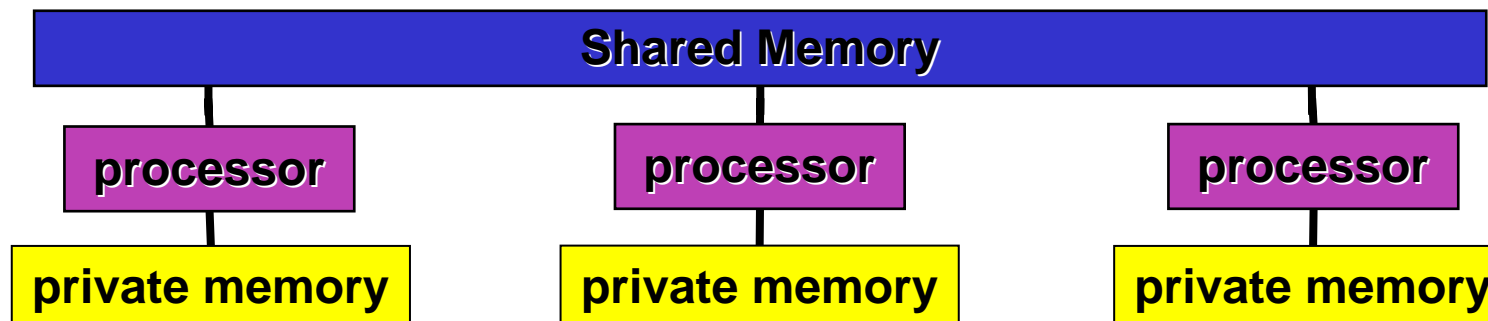| Keyword | Hits | Remarks |
|---|---|---|
| MPI | 612 | since 1994 |
| OpenMP | 150 | since 1997 with some 28 hits in our own chapter about OpenMP |
| threads | 109 | frequently in the context of OpenMP, 57 in our chapter about OpenMP |
| C++ | 87 | since 1983 |
| Fortran | 69 | since 1957 |
| Chapel | 49 | with some 22 hits in Zima's chapter about Chapel |
| UPC | 30 | since 2001 |
| Co-array Fortran | 27 | since 1998 |
| hybrid MPI/OpenMP | ~26 | hard to count |
| C | ~20 | hard to count |
| HPF | 11 | since 1993 |
| X10 | 9 | |
| Fortress | 6 | |
| Java | 5 | since 1995 |
| Titanium | 3 | |
| posix threads | 2 | 1995, Linux since 2003 |

*Petascale Computing: Algorithms and Applications (Chapman & Hall/Crc Comp. Sci. Ser.)*
*edited by David A. Bader , 2007, 528 pages, 24 contributions, 90 contributors*

**Scalable OpenMP Programming – D.an Mey**

RWTH
Center for
Computing and Communication

# Overview

- **Why OpenMP**
- **Short OpenMP Introduction**
- **OpenMP on NUMA Machines**
- **OpenMP on Clusters**
- **Conclusion**

**RWTH** CCC
Center for
Computing and Communication
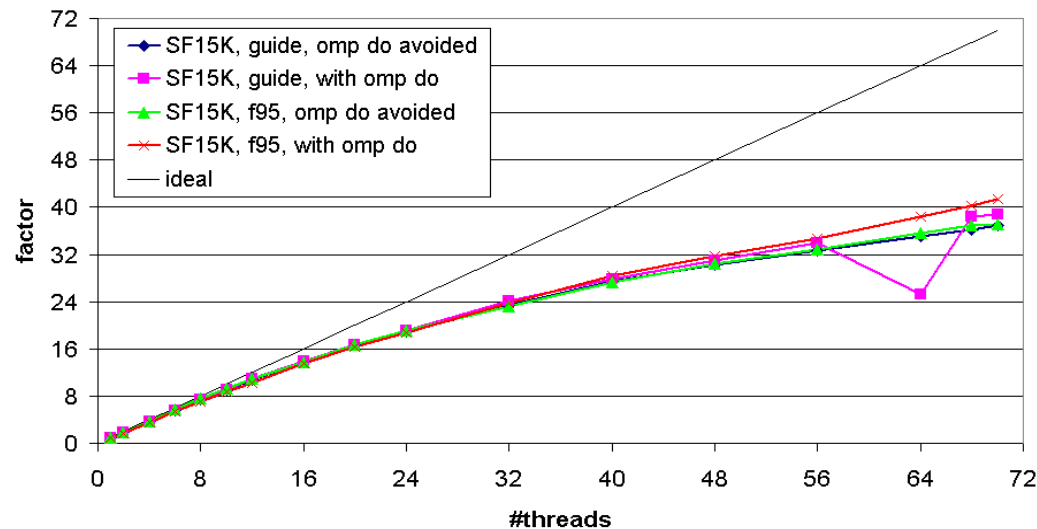
# Memory Model of OpenMP

- OpenMP: Shared-Memory model
  - All threads share a common address space (shared memory)
  - Threads can have private data (explicit user control)

- Relaxed memory consistency
  - Temporary View ("*Caching*"):
    Memory consistency is guaranteed only after synchronization points, namely implicit and explicit `flush`es

    - Each OpenMP `barrier` includes a `flush`
    - Exit from worksharing constructs include barriers by default
    - Entry to and exit from `critical` regions include a `flush`
    - Entry to and exit from lock routines (OpenMP API) include a `flush`

| Shared Memory | | |
|---|---|---|
| processor | processor | processor |
| private memory | private memory | private memory |

RWTH

Center for
Computing and Communication

# Heat Flow Simulation with FEM - ThermoFlow60

*Thomas Haarmann, Wolfgang Koschel,  Jet Propulsion Laboratory, RWTH Aachen University*

- simulation of the heat flow in a rocket combustion chamber

- Finite Element Method

- OpenMP Parallelelization
  - 30000 lines of Fortran
  - 200 OpenMP directives, 69 parallel loops,
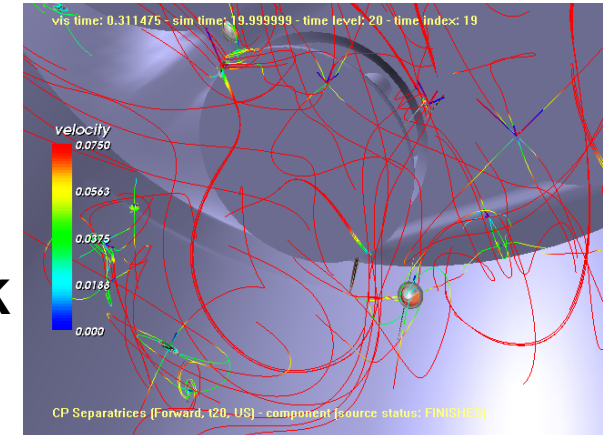  - 1 main parallel region, "*orphaning*"



**speed-up**

- SF15K, guide, omp do avoided
- SF15K, guide, with omp do
- SF15K, f95, omp do avoided
- SF15K, f95, with omp do
- ideal

factor / #threads

**Speedup: ~40 with 68 threads**

**Scalable OpenMP Programming – D.an Mey**

RWTH CCC
Center for Computing and Communication

# Nested OpenMP for Critical Point Computation

*Samuel Sarholz, Andreas Gerndt, Computing and Communication Center, RWTH Aachen University*

- **Analysis of complex and accurate fluid dynamics simulations**

- **Extraction of Critical Points for VR ( Location with velocity = 0 )**

- **25-100% efficiency with 128 threads on Sun Fire E25K (72 UltraSPARC IV dual core) depending on data set**



```
// Loop over time levels
#pragma omp parallel for num_threads(nTimeThreads) schedule(dynamic,1)
for (curT=1; curT<=maxT; ++curT) {
// Loop over Blocks
#pragma omp parallel for num_threads(nBlockThreads) schedule(dynamic,1)
for (curB=1; curB<=maxB; ++curB) {
// Loop over Cells
#pragma omp parallel for num_threads(nCellThreads) schedule(guided)
for (curC=1; curC<=maxC; ++curC) {
FindCriticalPoints (curT, curB, curC); // highly adaptive algorithm (bisectioning)
} } }                                   // huge load imbalances
```

**14**

**RWTH**
Center for
Computing and Communication

# Overview

- **Why OpenMP**
- **Short OpenMP Introduction**
- **OpenMP on NUMA Machines**
- **OpenMP on Clusters**
- **Conclusion**

**RWTH**
Center for
Computing and Communication

# The Earth is Flat

*OpenMP is Hardware agnostic*
*It has no notion of data locality*
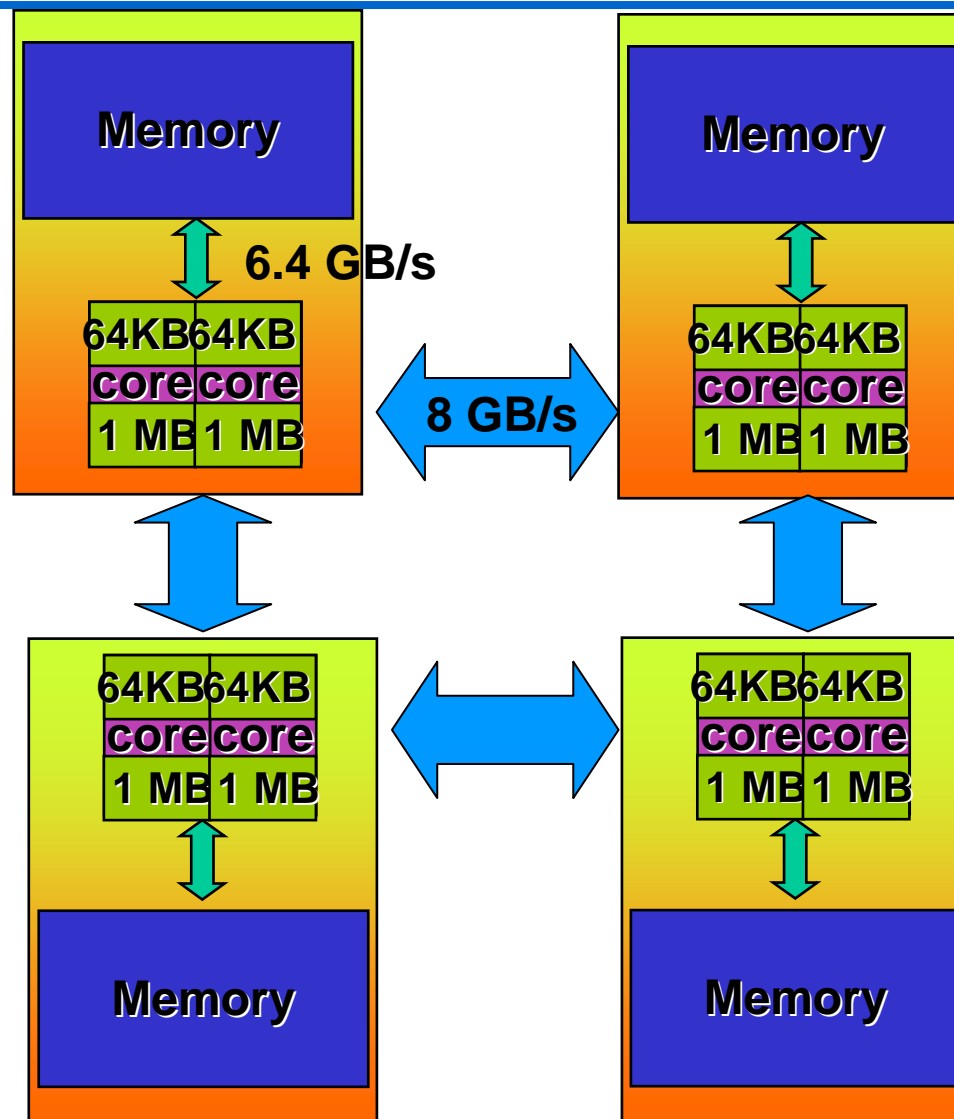
=>

The Affinity Problem**:**

**How to maintaining or improve
the nearness of threads and
their most frequently used data**

Or**:**

**Where to run threads?
Where to place data?**

RWTH
Center for
Computing and Communication

# Sun Fire V40z (w/ dualcore AMD Opteron Chip)

**Memory**

**6.4 GB/s**

| 64KB | 64KB |
|------|------|
| core | core |
| 1 MB | 1 MB |

**Memory**

| 64KB | 64KB |
|------|------|
| core | core |
| 1 MB | 1 MB |

**8 GB/s**

| 64KB | 64KB |
|------|------|
| core | core |
| 1 MB | 1 MB |

| 64KB | 64KB |
|------|------|
| core | core |
| 1 MB | 1 MB |

**Memory**

**Memory**

**4 AMD Opteron 875 dual core processors 2.2 GHz**

*Cache-coherent HyperTransport Connections*

**Scalable OpenMP Programming – D.an Mey**

RWTH
**Center for Computing and Communication**

# Sparse-Matrix-Vector-Multiplication as part of the Navier Stokes Solver DROPS ( C++ )

*C. Terboven (RZ,RWTH), A. Spiegel (RZ,RWTH), D. an Mey (RZ,RWTH),*
*S. Gross(IGPM,RWTH), and V. Reichelt (IGPM,RWTH).*

**4xdualcore Opteron 2.2 GHz, ccNUMA**

**12xdualcore UltraSPARC 1.2 GHz, flat memory**

- SF 2900 (first touch)
- SF V40z (first touch)
- SF 2900 (ignore locality)
- (ignore locality)

19,6 Mio nonzeros
233,334 matrix dimension
225 MB memory footprint

MFLOPS

1400.00
1200.00
1000.00
800.00
600.00
400.00
200.00
0.00

0    5    10    15    20    25

# threads

*IWOMP 2005*

**Performance of a cc-Numa system is very sensitive to data placement.**

18

*Scalable OpenMP Programming – D.an Mey*

RWTH
CCC
Center for
Computing and Communication

# Thread-Data-Affinity (1 of 2)

- In an ideal world the operating system together with the OpenMP runtime system would handle affinity automatically.

- In simple situations things might work well:

  - Exclusive access to the compute node

  - Single level of parallelism

  - Static program behaviour concerning thread and data usage

  - Initialization of data by the same thread which later uses the data („first touch policy")

**RWTH** Center for Computing and Communication
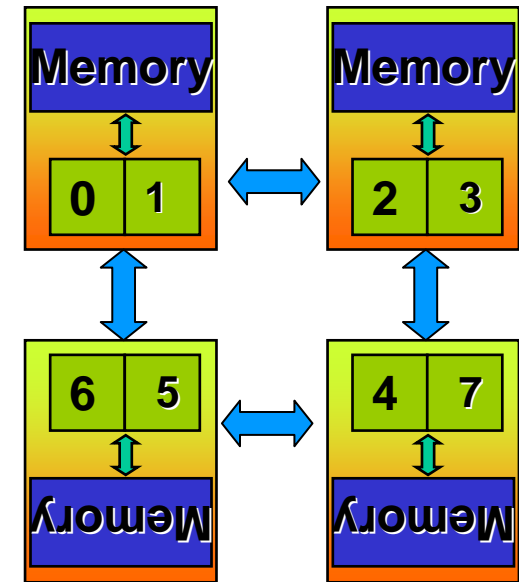
# Thread-Data-Affinity (2 of 2)

- In more complicated situations, you may want to

  - Bind threads explicitly
    (How about multi user mode? Hybrid parallelization?)

  - Carefully initialize data

  - If necessarry and possible, migrate data (or threads)

    - Solaris MPO madvise() implements „next touch strategy"

    - Linux 2.6.18: move_pages()
      can be used to implement „next touch strategy"
      ( RWTH: prototype by RZ, better solution by LfBS )

    - Windows: Migration is not yet supported)

- Nested OpenMP is implemented with thread pool

  - Inner teams' threads loose affinity to their data

  - Sun Studio on Solaris: SUNW_MP_THR_AFFINITY=TRUE

**RWTH**
CCC
**Center for
Computing and Communication**

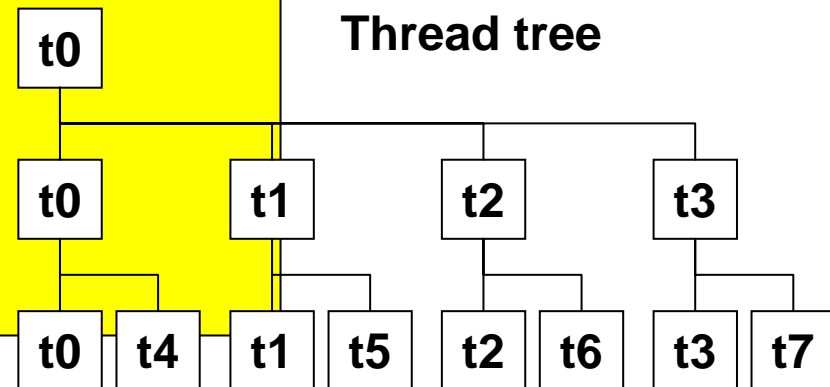# OpenMP nested, here: 4x2 threads

```fortran
!$omp  parallel private(me) num_threads(4)
      me = omp_get_thread_num()
      CALL stream(a(1,me),b(1,me),c(1,me))
!$omp end parallel
      ...
      subroutine stream (a,b,c)
      double precision a(*),b(*),c(*)
      ...
!$omp parallel do num_threads(2)
        do 50 j = 1,n
            c(j) = a(j)+x*b(j)
  50      continue
!$omp end parallel do
      ...
```
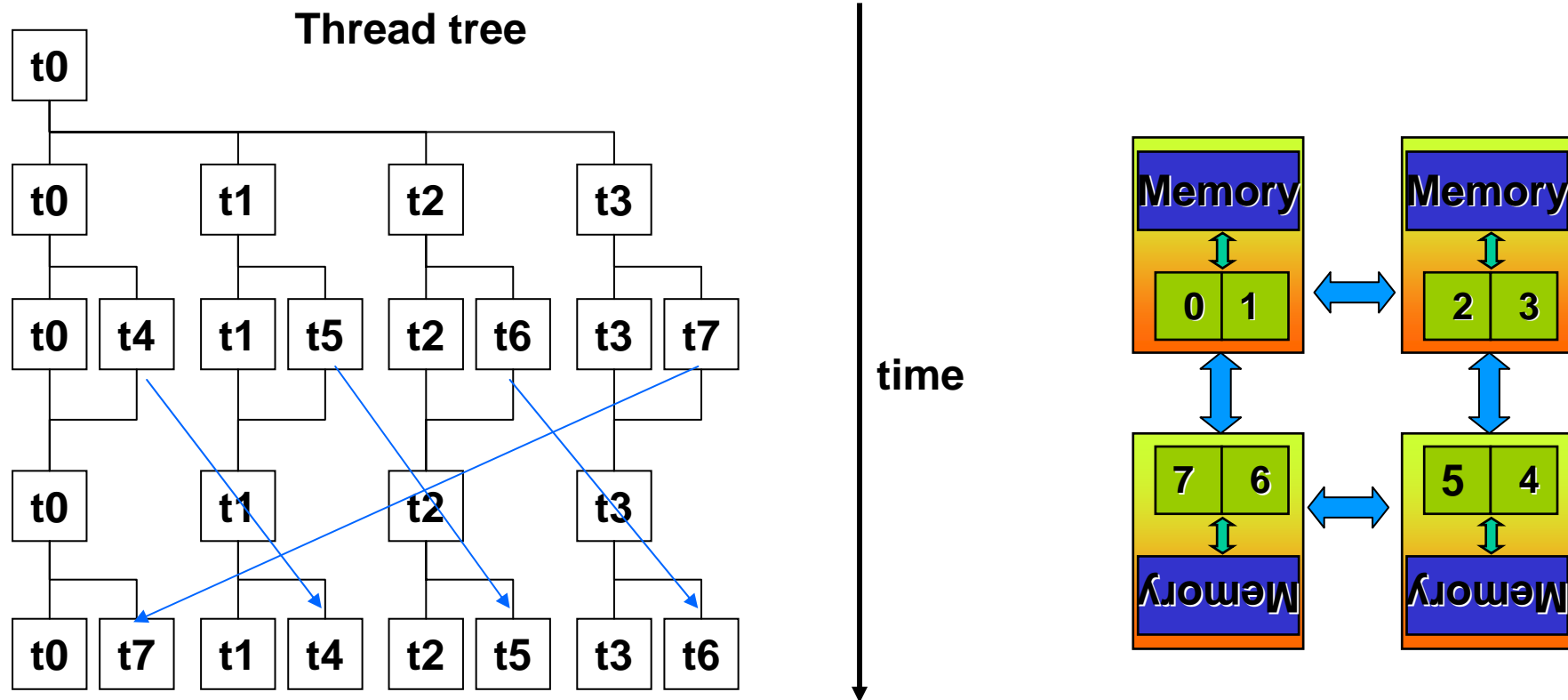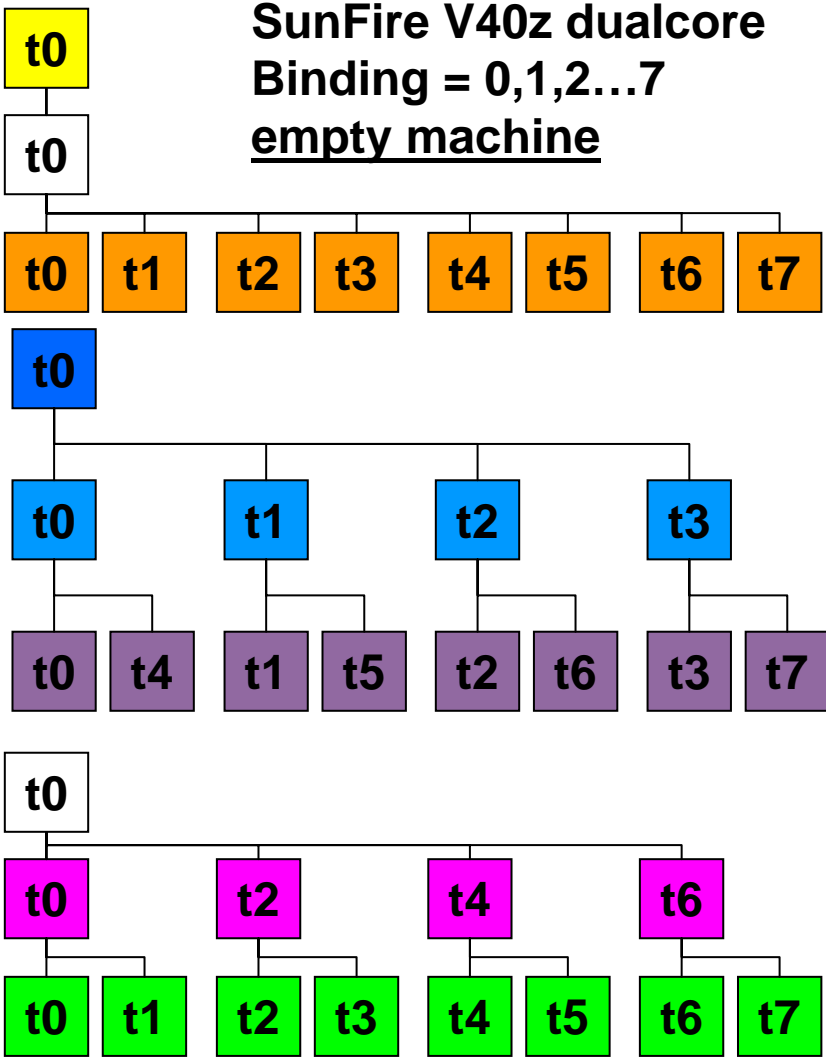
**Memory**   **Memory**

| 0 | 1 |   | 2 | 3 |

| 6 | 5 |   | 4 | 7 |

**Memory**   **Memory**

**Thread tree**

t0

t0   t1   t2   t3

t0 t4 t1 t5 t2 t6 t3 t7

**Scalable OpenMP Programming – D.an Mey**

RWTH
Center for
Computing and Communication

# OpenMP nested

**Thread tree**



**time**

Typically OS threads are organized in a pool
and may be allocated variably,
thus loosing data affinity !

**Scalable OpenMP Programming – D.an Mey**

RWTH
Center for
Computing and Communication

# OpenMP nested

**SunFire V40z dualcore**
**Binding = 0,1,2…7**
**empty machine**



| #thrds | First touch | Affinity | min MB/s | max MB/s |
|--------|-------------|----------|----------|----------|
| 1 x 8 | Initial thread | n.a. | 2525 | 2534 |
| 1 x 8 | All inner threads | n.a. | 11786 | 11869 |
| 4 x 2 | Initial thread | no | 4x629 | 4x631 |
| 4 x 2 | Initial thread | yes | 4x628 | 4x631 |
| 4 x 2 | Inner master | no | 4x1312 | 4x1332 |
| 4 x 2 | Inner master | yes | 4x1329 | 4x1334 |
| 4 x 2 | Inner master | Yes+sort | 4x2881 | 4x2943 |
| 4 x 2 | All inner threads | no | 4x2640 | 4x2948 |
| 4 x 2 | All inner threads | yes | 4x2893 | 4x2950 |
| 4 x 2 | All inner threads | Yes + sort | 4x2922 | 4x2934 |

**Scalable OpenMP Programming – D.an Mey**

**RWTH**
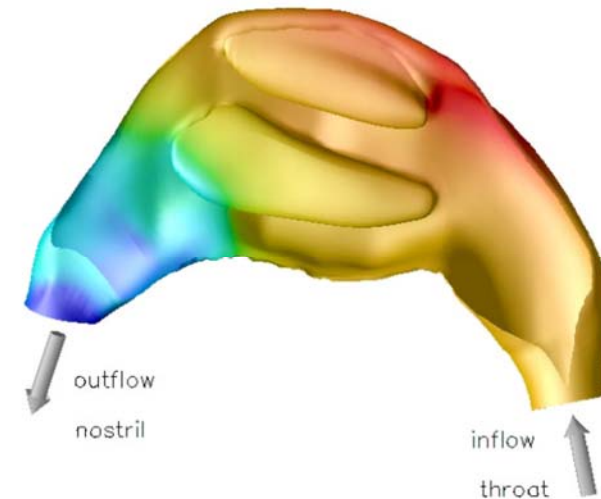Center for
Computing and Communication

# Simulating the Flow through the Human Nose TFS on Solaris

*S. Johnson (PSP), C. Ierotheou (PSP),*
*A. Spiegel (RZ,RWTH), D. an Mey (RZ,RWTH),*
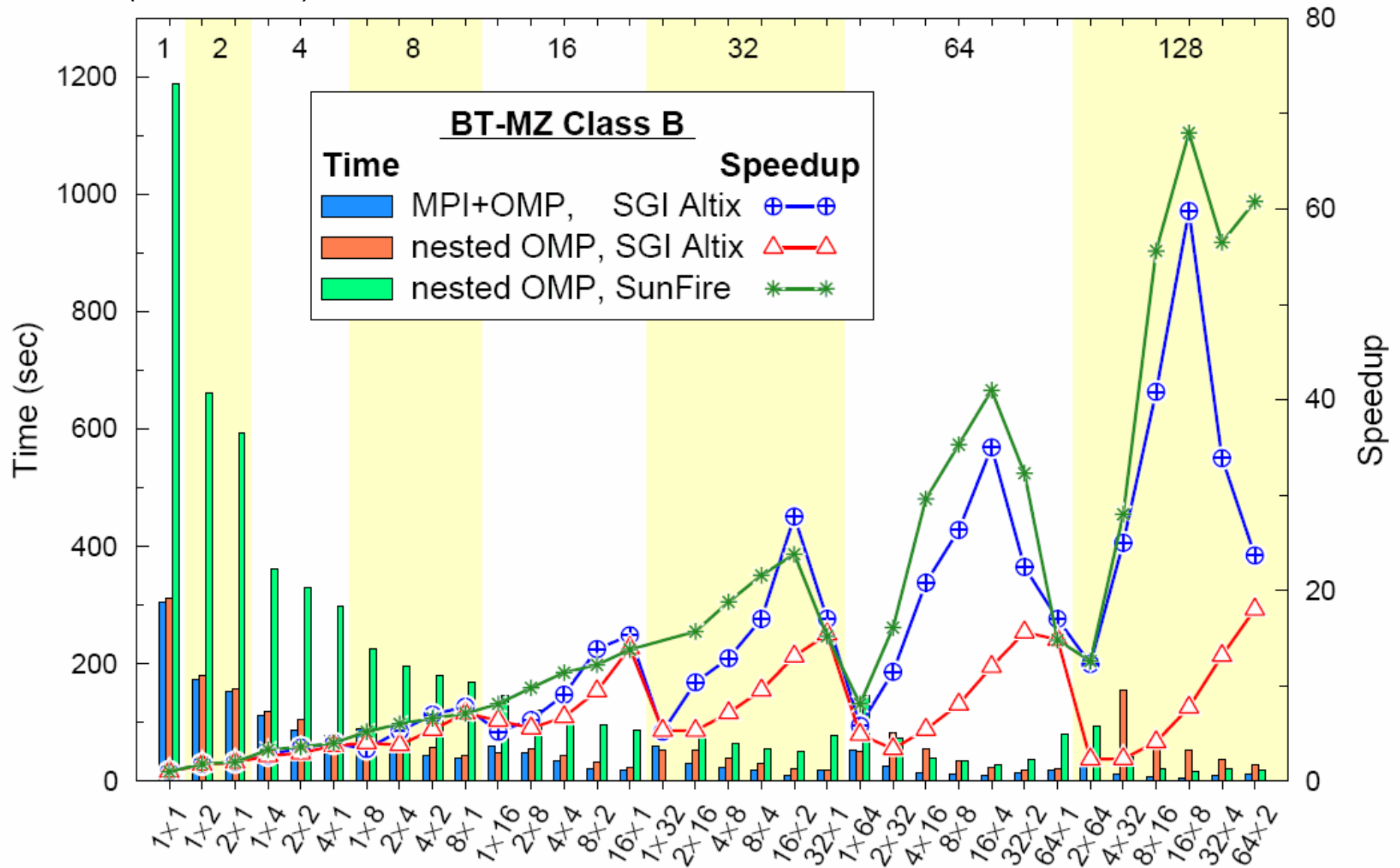*I. Hörschler (AIA, RWTH)*

**SUNW_MP_THR_AFFINITY=TRUE**

**Thread affinity + processor binding + data migration improved the performance by ~25 % on a Sun Fire E 25K**

outflow

nostril

inflow

throat

| Before | | Improved thread affinity | | |
|---|---|---|---|---|
| #threads | Speed-up | #threads | Speed-up | Strategy (best effort) |
| 64 | **20** | 64 | **25** | thread balancing 2-11 threads per team, static schedule, 16 threads in outer team |
| 121 | **20** | 128 | **27** | block grouping, 16 threads in outer team |

RWTH
Center for
Computing and Communication

# NPB Benchmark BT-MZ Class B

*H. Jin (Nasa Ames), et. al.*

# Overview

- **Why OpenMP**
- **Short OpenMP Introduction**
- **OpenMP on NUMA Machines**
- **OpenMP on Clusters**
- **Conclusion**

RWTH
Center for
Computing and Communication

# OpenMP on Clusters

o Multiple Approaches (based on MPI, on DSM …)
     so far not very successful or uncomplete.

o Intel Cluster OpenMP on Commodity Infiniband Cluster

   o Based on TreadMarks (twin pages, sending diffs,…)

   o Integrated in commercial compiler (C++ and F95)

   o Profits from OpenMP's memory model
     (relaxed consistency, temporary view of shared data, consistency
     enforced at well defined synchronization points.)

   o Need to explicitly mark some shared variables (`sharable` directive)

o ScaleMP – Versatile SMP™ Architecture

   o Aggregation of multiple x86 boards into one larger system

   o Cache coherent connection through InfiniBand

   o Modified IB stack and BIOS, caching strategies

   o Single system image, virtual SMP machine

   o Aggregation of all I/O resources to the OS

o Affinity matters!

RWTH

Center for
Computing and Communication

# EPCC OpenMP Micro-Benchmarks

J. M. Bull. Measuring Synchronization and Scheduling Overheads in OpenMP. 1999.

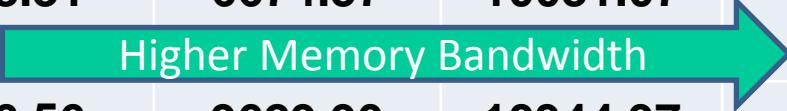| | Tigerton | Opteron | CLOMP | ScaleMP (MEG) |
|---|---|---|---|---|
| **PARALLEL FOR** 2 threads 16 threads | 1.31 5.01 | 1.36 7.17 | 723.77 4342.82 | 264.83 717.77 |
| **BARRIER** 2 threads 16 threads | 0.75 2.55 | 0.58 2.64 | 598.82 4062.67 | 144.45 429.35 |
| **REDUCTION** 2 threads 16 threads | 1.56 5.68 | 2.05 25.77 | 932.18 4686.00 | 298.06 801.91 |

About three orders of magnitude

**Overhead in microseconds [us].**

Binding: 1 Thread/board for CLOMP and ScaleMP(MEG)
8 Threads/board for CLOMP and ScaleMP(MEG)

**Scalable OpenMP Programming – D.an Mey**

# Stream Benchmark

| # threads | Tigerton | Opteron (*) | CLOMP | ScaleMP (RWTH) |
|-----------|----------|-------------|-------|----------------|
| 1 | 2080.78 | 1882.24 | 3321.08 | 2674.13 |
| 2 | 4033.88 | 3665.35 | 6495.34 | 5330.22 |
| 4 | 7008.31 | 6674.57 | 10031.07 | 10439.76 |
| 8 | 7156.56 | 9629.90 | 10344.97 | 17478.77 |
| 16 | 7508.01 | 8787.33 | 10473.24 | 18666.49 |

Higher Memory Bandwidth →

**Bandwidth in MB/s. Scattered Binding.**

(*) We see better performance on our 4-socket Opteron machine running Solaris

**Scalable OpenMP Programming – D.an Mey**
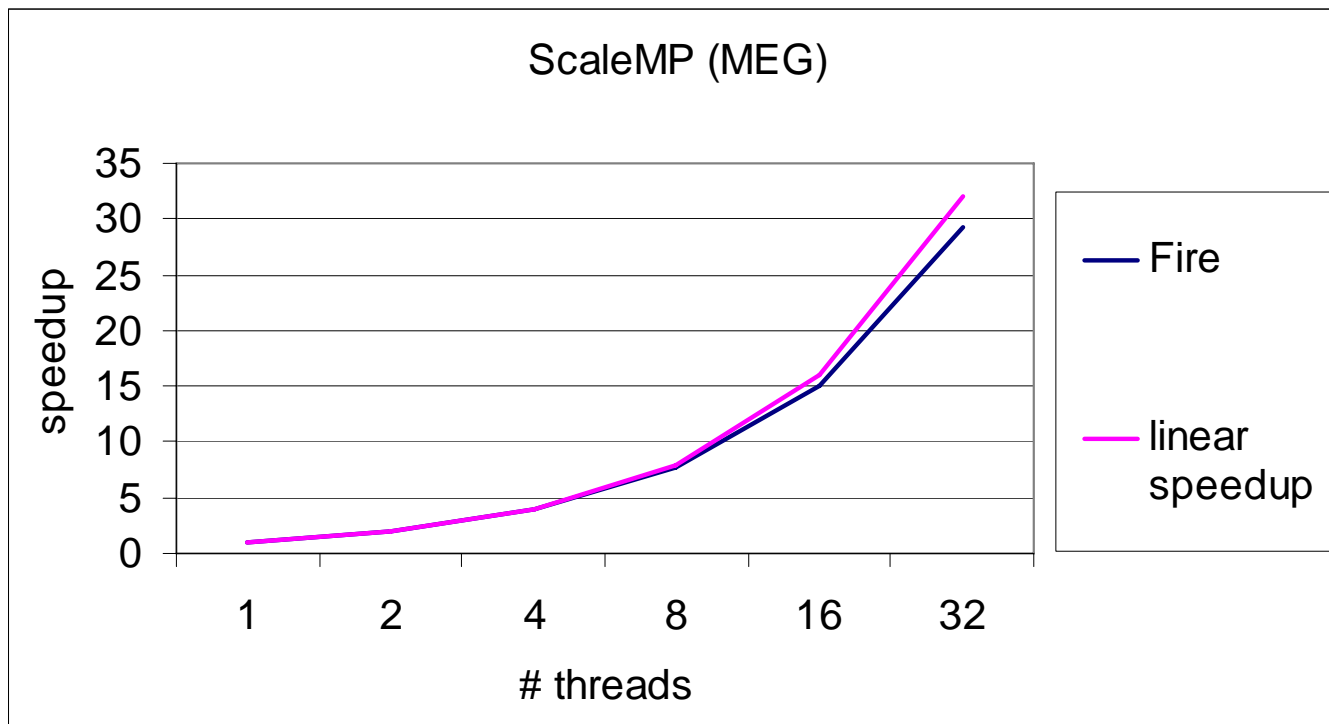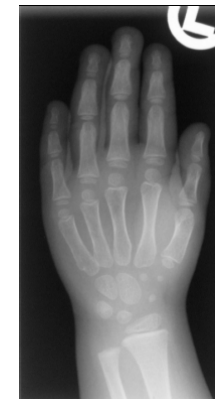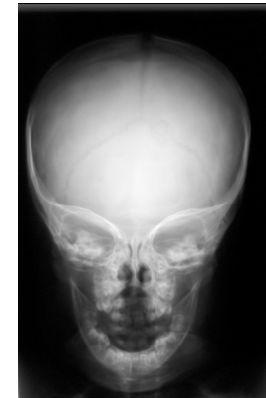
RWTH
Center for
Computing and Communication

FIRE = Flexible Image Retrieval Engine

– Compare the performance of common features on different databases

– Analysis of correlation of different features

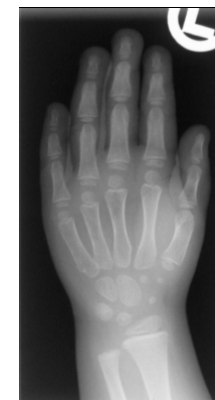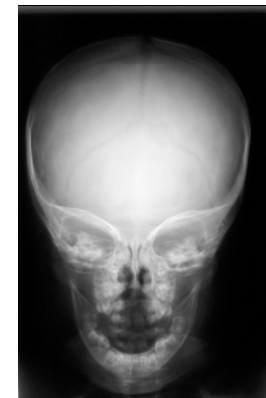*Thomas Deselaers and Daniel Keysers, RWTH I6:*
*Chair for Human Language Technology and Pattern Recognition*



**ScaleMP (MEG)**

**Scalable OpenMP Programming – D.an Mey**

**RWTH**
Center for
Computing and Communication

On the new 13 node system:
13 nodes with 2 Harpertown Processors at 2.5 GHz

| | Speed-up | | | |
|---|---|---|---|---|
| #threads | outer level | inner level | nested best effort | #threads on inner x outer level |
| 1 | 1,0 | 1,0 | **1,0** | 1 x 1 |
| 2 | 2,1 | 2,1 | **2,1** | 2 x 1 |
| 4 | 4,0 | 3,9 | **4,1** | 2 x 2 |
| 8 | 7,8 | 7,1 | **8,0** | 2 x 4 |
| 16 | 14,8 | 12,6 | **15,6** | 2 x 8 |
| 32 | 25,5 | | **29,9** | 4 x 8 |
| 64 | 45,4 | | **53,2** | 4 x 16 |
| 104 | | | **67,1** | 4 x 26 |

**Scalable OpenMP Programming – D.an Mey**

RWTH
Center for
Computing and Communication

# FIRE: Image Retrieval System on 144-core SF E25K
## Nested OpenMP improves scalability

| Speedup of FIRE | Sun Fire E25K, 72 dual-core UltraSPARC-IV processors | | |
|---|---|---|---|
| # Threads | Only outer level | Only inner level | Nested OpenMP |
| 4 | --- | 3.8 | --- |
| 8 | --- | 7.6 | --- |
| 16 | 14.8 | 14.1 | 15.4 |
| 32 | 29.6 | 28.9 | 30.6 |
| 72 | 56.5 | --- | 67.6 |
| 144 | --- | --- | 133.3 |

**Scalable OpenMP Programming – D.an Mey**

RWTH
Center for
Computing and Communication

# Sparse Matrix-Vector-Multiplication [Mflop/s]
## Apply Suitable Strategy!



parallel loop over #rows,
**dynamic loop sched**.

#nonzeros
**statically partitioned**

Scalable OpenMP Programming – D.an Mey

- Simulation of Coupled Flow, Heat Transfer and Transport Interaction
- BiCGStab Solver with ILU0 Preconditioner
- Nested Parallelization with OpenMP
- Explicite binding in all inner parallel regions



*SGI F1200 vSMP: 4 nodes with 2x4 cores*
*Dell 140 vSMP: 4 nodes with 2x2 cores*
*Harpertown: 1 nodes with 2x4 cores*

**ScaleMP: ~ 3x on 4 nodes**

Legend:
- SGI F1200 4x2x2 3 GHz
- Dell140 4x2x4 3.16 GHz
- Harpertown 2x4 3 GHz

X-axis: **#threads outer team x #threads inner team**
Y-axis: Seconds

# SHEMAT on ScaleMP (2 of 2)

Cranking up the probem size for the new 13 node system:
13 nodes with 2 Harpertown Processors at 2.5 GHz

| #boards | 1 | 2 | 4 | 8 | 10 |
|---|---|---|---|---|---|
| best effort timing | 8664,8 | 3357,9 | 2264,9 | 1281,8 | 981,6 |
| #threads outer level | 1 | 4 | 8 | 16 | 20 |
| #threads inner level | 4 | 1 | 2 | 2 | 2 |
| #core used per board | 4 | 2 | 4 | 4 | 4 |
| speed-up on board level | **1,0** | **2,6** | **3,8** | **6,8** | **8,8** |
| speed-up versus 1 thread | 1,4 | 3,6 | 5,3 | 9,4 | 12,2 |

Speed-up across the nodes good (threads on outer level don't interact much)
Speed-up withing the nodes bad (limited memory bandwidth)

**Scalable OpenMP Programming – D.an Mey**

RWTH
Center for
Computing and Communication

# Overview

- **Why OpenMP**
- **Short OpenMP Introduction**
- **OpenMP on NUMA Machines**
- **OpenMP on Clusters**
- **Conclusion**

**Scalable OpenMP Programming – D.an Mey**

**RWTH**
Center for
Computing and Communication

# Conclusion

- Scalable applications may need multiple levels of parallelization
- OpenMP suitable for a growing number of cores per node
- Combining MPI and OpenMP is getting more popular
- OpenMP on Clusters an alternative, if MPI is too hard to apply.

- Thread/Data Affinity is essential for OpenMP performance on ccNUMA machines and even more on Clusters
- OpenMP is hardware agnostic
- Needs control of thread and data placement
- Needs data migration, explicite and/or automatic for irregular, adaptive problems

**RWTH**
Center for
Computing and Communication

RZ - Parallel Programming in Computational Engineering and Science (PPCES) March 2009 - HPC Tutorials - - Mozilla Firefox

Datei  Bearbeiten  Ansicht  Chronik  Lesezeichen  Extras  Hilfe

http://www.rz.rwth-aachen.de/go/id/sms/?lang=en

Google

Meistbesuchte Seiten  Erste Schritte  Aktuelle Nachrichten -...  Marder Alarm  http://openmp.org/tw...  http://www.rz.rwth-a...  RZ Aktuelles

**RWTH**
Rechen- und
Kommunikationszentrum

Search   Help   RZ Internal

A-Z   Feedback   RWTH

Students | Faculties and Staff | Projects and Cooperations | MATSE | Press |

News

About Us

Our services

High Performance
Computing

Virtual Reality

Events

Help and Bugs

## Parallel Programming in Computational Engineering and Science (PPCES) March 2009 - HPC Tutorials -

Monday, March 23 - Friday, March 27, 2009

(intel)   **Microsoft**   **Sun** microsystems

Kindly supported by:

**Further information**

Questionnaire:
We hope that you enjoyed
the PPCES. Your feedback is
welcome on the
questionnaire you find here
>>>

| Date | Time | Location |
|---|---|---|
| Monday, March 23<br>Tuesday, March 24<br>Wednesday, March 25<br>Thursday, March 26<br>Friday, March 27 | 14:00 - 17:30 (*)<br>09:00 - 17:30<br>09:00 - 17:30<br>09:00 - 17:30<br>09:00 - 12:30 | Center for Computing and Communication<br>RWTH Aachen University<br>Seffenter Weg 23<br>52074 Aachen |

(*) We like to draw your attention to a presentation by **Horst Simon** (University of California Berkeley) on Monday morning on Future Directions in High Performance Computing 2009-2018 in the SuperC

| | | | | | | |
|---|---|---|---|---|---|---|
| ▪ Introduction | ▪ Related Events | ▪ Sponsors | ▪ Participants | ▪ Flyer & Poster | ▪ Agenda | ▪ Course Material |
| ▪ **Questionnaire!!!** | ▪ Links | ▪ Travel Information | ▪ Contact | | | |

## Introduction