# "OpenMP Does Not Scale"

**Ruud van der Pas**

**Distinguished Engineer**

*Architecture and Performance*

*SPARC Microelectronics, Oracle*

*Santa Clara, CA, USA*

ORACLE

# Agenda

- **The Myth**
- **Deep Trouble**
- **Get Real**
- **The Wrapping**

# The Myth

http://www.reference.com/search?q=myth  Google

Getting Started  Latest Headlines  Developer Guide

myth - Information from Refe...

# Myth

Wikiped...

Myth

**"A myth, in popular use, is something that is widely believed but false."** ......

Mythology, mythography, or folkloristics. In these academic fields, a myth (*mythos*) is a sacred story concerning the origins of the world or how the world and the creatures in it came to have their present form. The active beings in myths are generally gods and heroes. Myths often are said to take place before recorded history begins. In saying that a myth is a sacred narrative, what is meant is that a myth is believed to be true by people who attach religious or spiritual significance to it. Use of the term by scholars does not imply that the narrative is either true or false. See also legend and tale.

A myth, in popular use, is something that is widely believed but false. This usage, which is often pejorative, arose from labeling the religious stories and beliefs of other cultures as being incorrect, but it has spread to cover non-religious beliefs as well. Because of this usage, many people take offense when the religious narratives they believe to be true are called myths (see Religion and mythology for more information). This usage is frequently confused with fiction, legend, fairy tale, folklore, fable, and urba... distinct meaning in academia.

Phoenix Myth

Myth Nightclub

Myth

P Indicates premium content, which is available only to subscribers.

*(source: www.reference.co...*

# "OpenMP Does Not Scale"

**A Common Myth**

**A Programming Model Can Not "Not Scale"**

**What Can Not Scale:**

**The Implementation**
**The System Versus The Resource Requirements**

**Or ..... You**

# *Hmmm .... What Does That Really Mean ?*

# Some Questions I Could Ask

*"Do you mean you wrote a parallel program, using OpenMP and it doesn't perform?"*

*"I see. Did you make sure the program was fairly well optimized in sequential mode?"*

ORACLE®

# Some Questions I Could Ask

*"Oh. You didn't. By the way, why do you expect the program to scale?"*

*"Oh. You just think it should and you used all the cores. Have you estimated the speed up using Amdahl's Law?"*

*"No, this law is not a new EU financial bail out plan. It is something else."*

# Some Questions I Could Ask

*"I understand. You can't know everything. Have you at least used a tool to identify the most time consuming parts in your program?"*

*"Oh. You didn't. You just parallelized all loops in the program. Did you try to avoid parallelizing innermost loops in a loop nest?"*

*"Oh. You didn't. Did you minimize the number of parallel regions then?"*

*"Oh. You didn't. It just worked fine the way it was.*

ORACLE®

# More Questions I Could Ask

*"Did you at least use the nowait clause to minimize the use of barriers?"*

*"Oh. You've never heard of a barrier. Might be worth to read up on."*

*"Do all threads roughly perform the same amount of work?"*

*"You don't know, but think it is okay. I hope you're right."*

ORACLE®

# I Don't Give Up That Easily

*"Did you make optimal use of private data, or did you share most of it?"*

*"Oh. You didn't. Sharing is just easier. I see.*

# I Don't Give Up That Easily

*"You seem to be using a cc-NUMA system. Did you take that into account?"*

*"You've never heard of that either. How unfortunate. Could there perhaps be any false sharing affecting performance?"*

*"Oh. Never heard of that either. May come handy to learn a little more about both."*

ORACLE®

# The Grass Is Always Greener …

*"So, what did you do next to address the performance ?"*

*"Switched to MPI. I see. Does that perform any better then?"*

*"Oh. You don't know. You're still debugging the code."*

# Going Into Pedantic Mode

*"While you're waiting for your MPI debug run to finish (are you sure it doesn't hang by the way ?), please allow me to talk a little more about OpenMP and Performance."*

# Deep Trouble

# OpenMP And Performance/1

- The transparency and ease of use of OpenMP are a mixed blessing

  → Makes things pretty easy

  → May mask performance bottlenecks

- In the ideal world, an OpenMP application "just performs well"

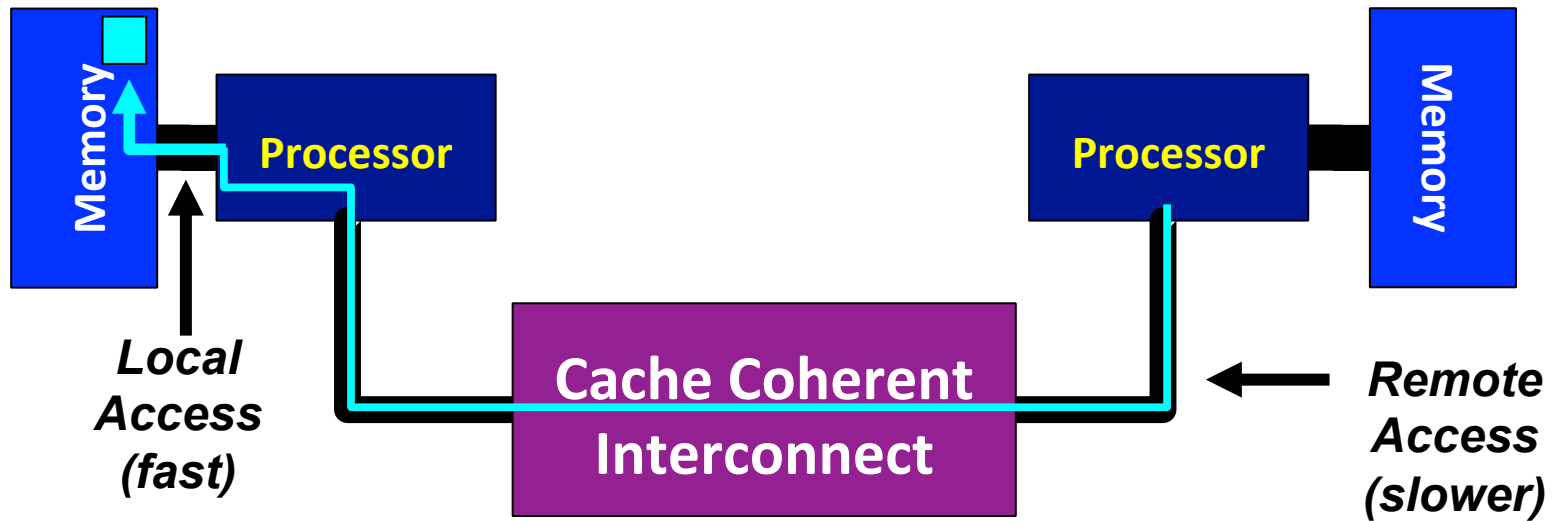- Unfortunately, this is not always the case

# OpenMP And Performance/2

- Two of the more obscure things that can negatively impact performance are cc-NUMA effects and False Sharing

- ***Neither of these are restricted to OpenMP***

  → They come with shared memory programming on modern cache based systems

  → But they might show up because you used OpenMP

  → In any case they are important enough to cover here

24

ORACLE®

# *Considerations for cc-NUMA*
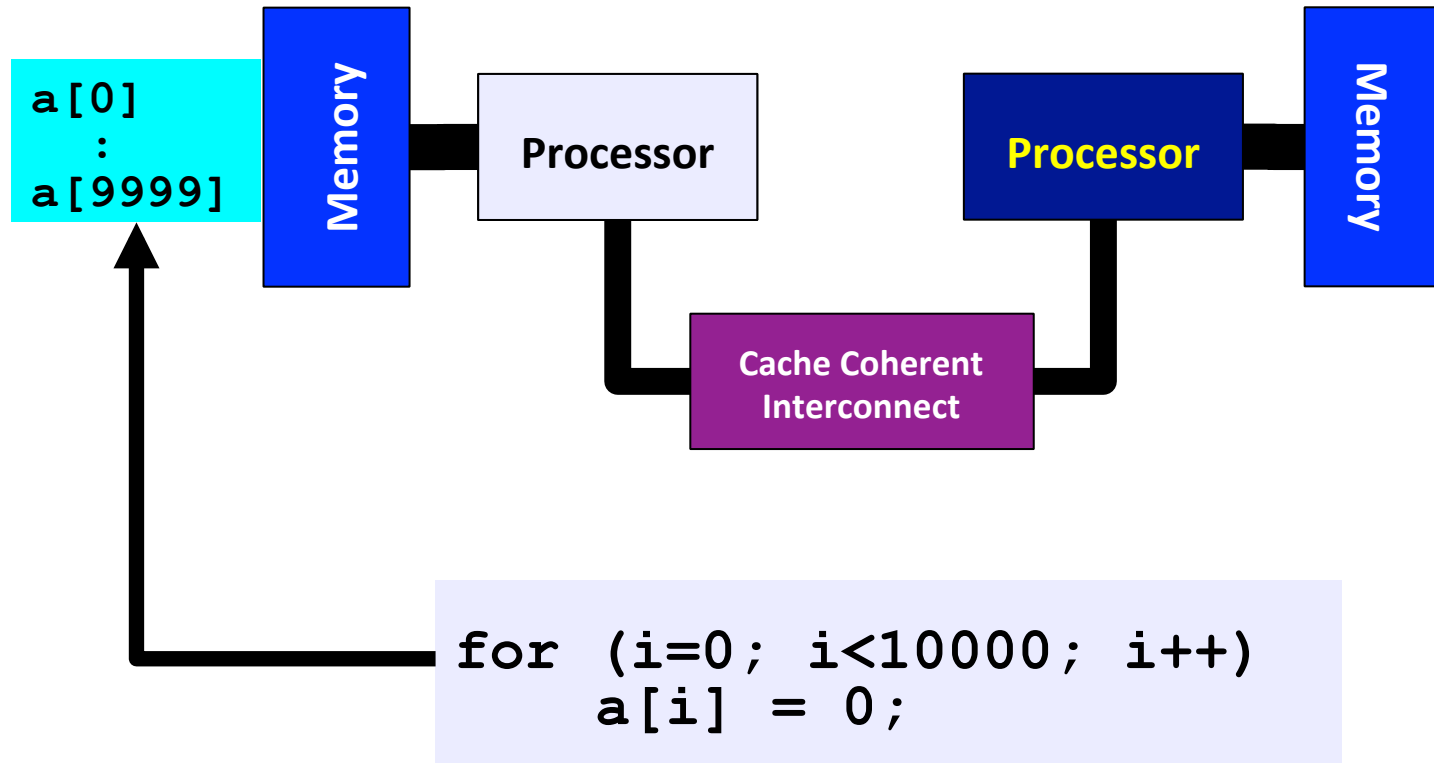
# A Generic cc-NUMA Architecture



**Main Issue: How To Distribute The Data ?**
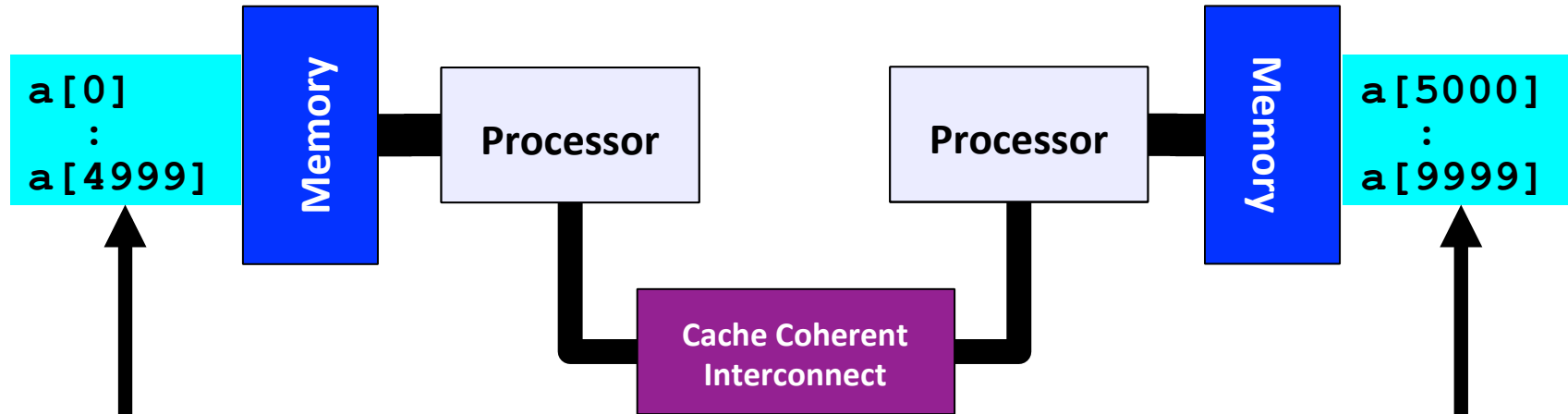
# About Data Distribution

- Important aspect on cc-NUMA systems

    → If not optimal, longer memory access times and hotspots

- OpenMP 4.0 does provide support for cc-NUMA

    → Placement under control of the Operating System (OS)

    → User control through OMP_PLACES

- Windows, Linux and Solaris all use the "First Touch" placement policy by default

    → May be possible to override default (check the docs)

ORACLE®

# Example First Touch Placement/1



```
for (i=0; i<10000; i++)
      a[i] = 0;
```

**First Touch**
*All array elements are in the memory of the processor executing this thread*

# Example First Touch Placement/2

a[0]
:
a[4999]

**Memory**

**Processor**

**Processor**

**Memory**

a[5000]
:
a[9999]

**Cache Coherent Interconnect**

`#pragma omp parallel for num_threads(2)`

```
for (i=0; i<10000; i++)
    a[i] = 0;
```

*First Touch*
*Both memories now have "their own half" of the array*

ORACLE®
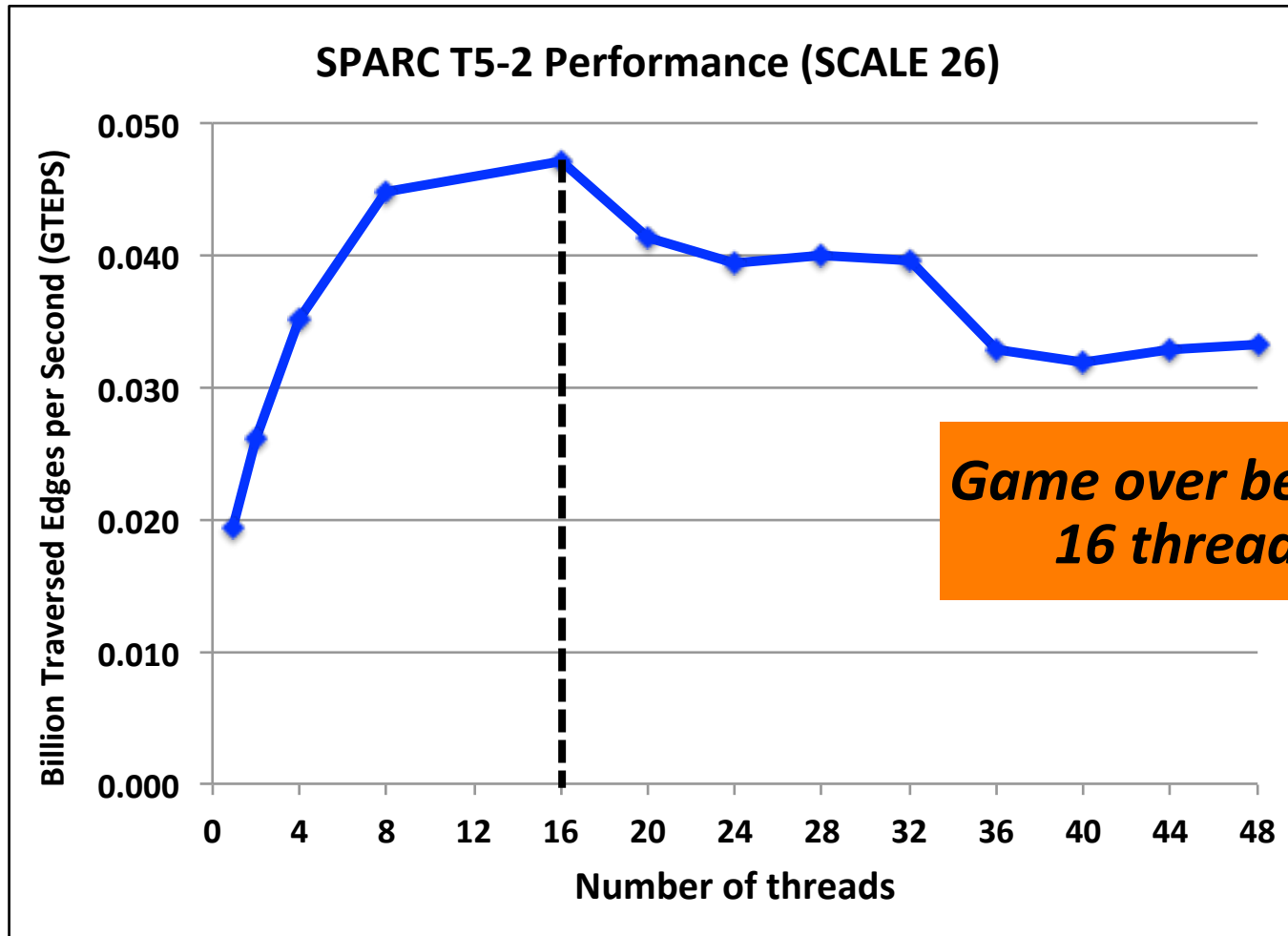
# Get Real

"Don't Try This At Home (yet)"
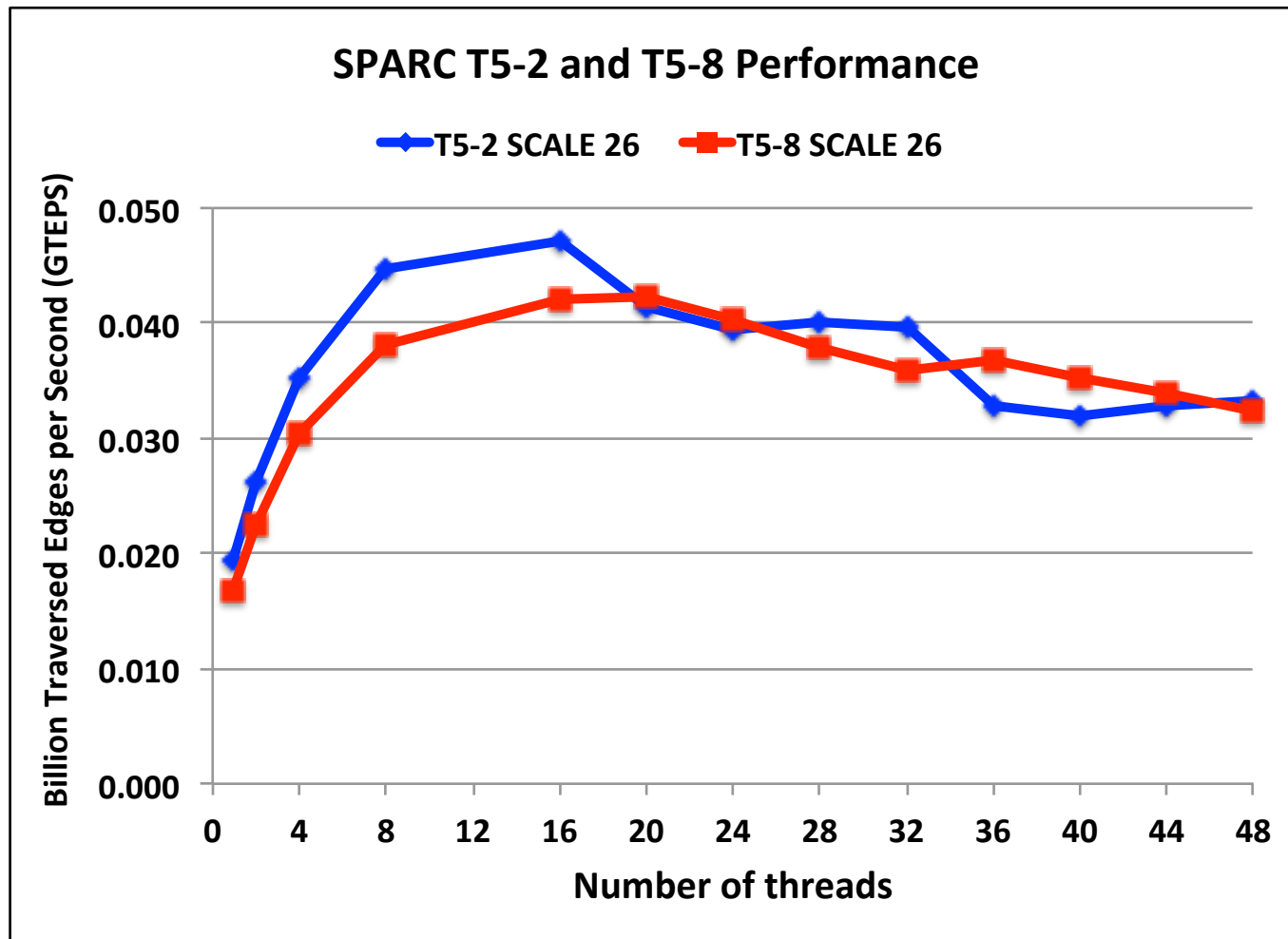
# The Initial Performance (35 GB)



SPARC T5-2 Performance (SCALE 26)

*Game over beyond 16 threads*

**That doesn't scale very well**

**Let's use a bigger machine !**

# Initial Performance (35 GB)



SPARC T5-2 and T5-8 Performance

# Oops! That can't be true

# Let's run a larger graph !

ORACLE®

3

# Initial Performance (280 GB)



SPARC T5-2 and T5-8 Performance

# *Let's Get Technical*

# Total CPU Time Distribution



Total CPU Time Percentage Distribution (Baseline, SCALE 26)

Legend:
- Atomic operations
- OMP-atomic_wait
- OMP-critical_section_wait
- OMP-implicit_barrier
- Other
- Function 1
- Function 2

X-axis: Number of threads

# Bandwidth Of The Original Code



SPARC T5-2 Measured Bandwidth (BASE, SCALE 28, 16 threads)

Less than half of the memory bandwidth is used

# Summary Original Version

- Communication costs are too high
    - Increases as threads are added
    - This seriously limits the number of threads used
    - This is turn affects memory access on larger graphs
- The bandwidth is not balanced
- Fixes:
    - Find and fix many OpenMP inefficiencies
    - Use some efficient atomic functions

ORACLE®

# Methodology

If The Code Does Not Scale Well

*Use A Profiling Tool*

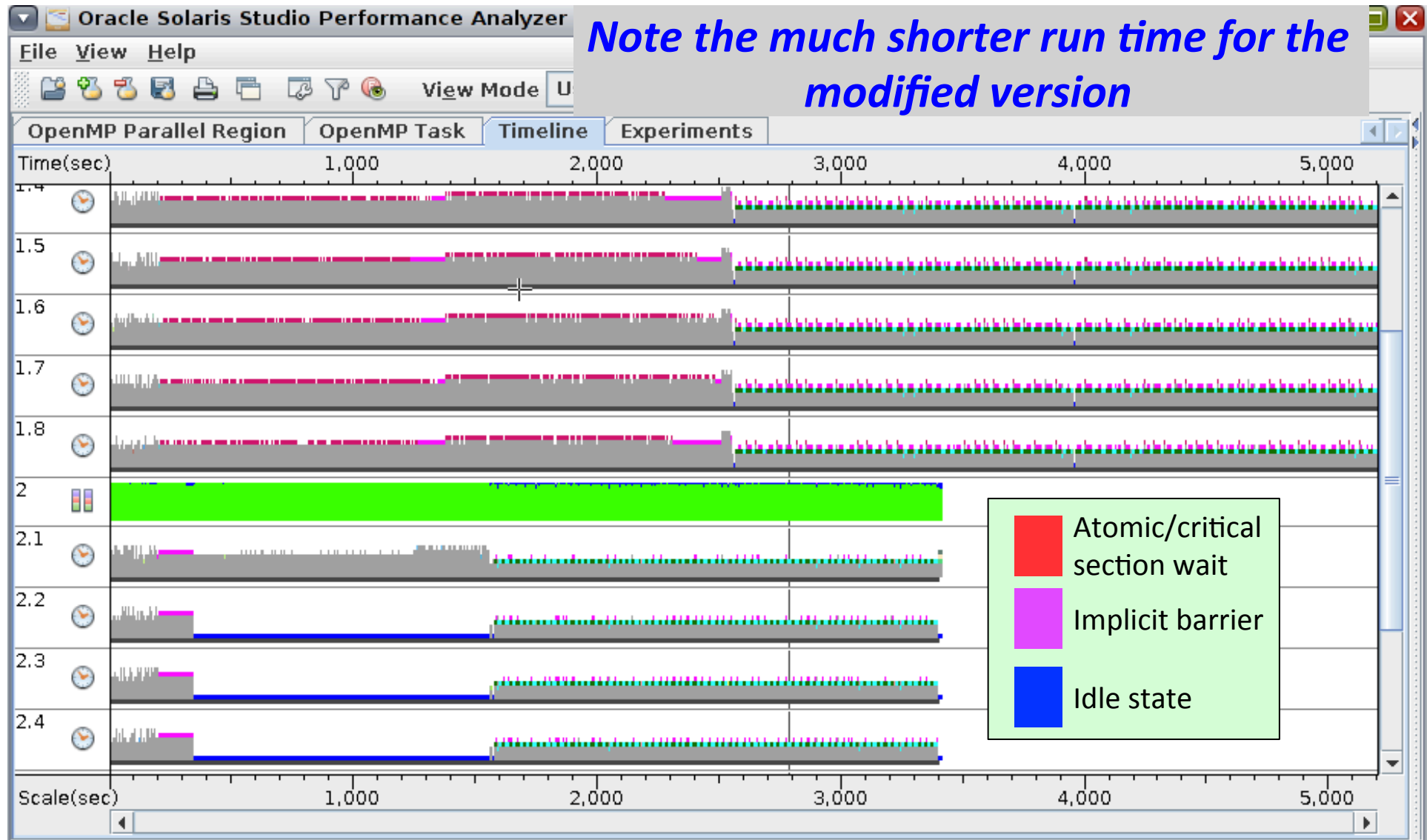Use The Checklist To Identify Bottlenecks

Tackle Them One By One

This Is An Incremental Approach

But Very Rewarding
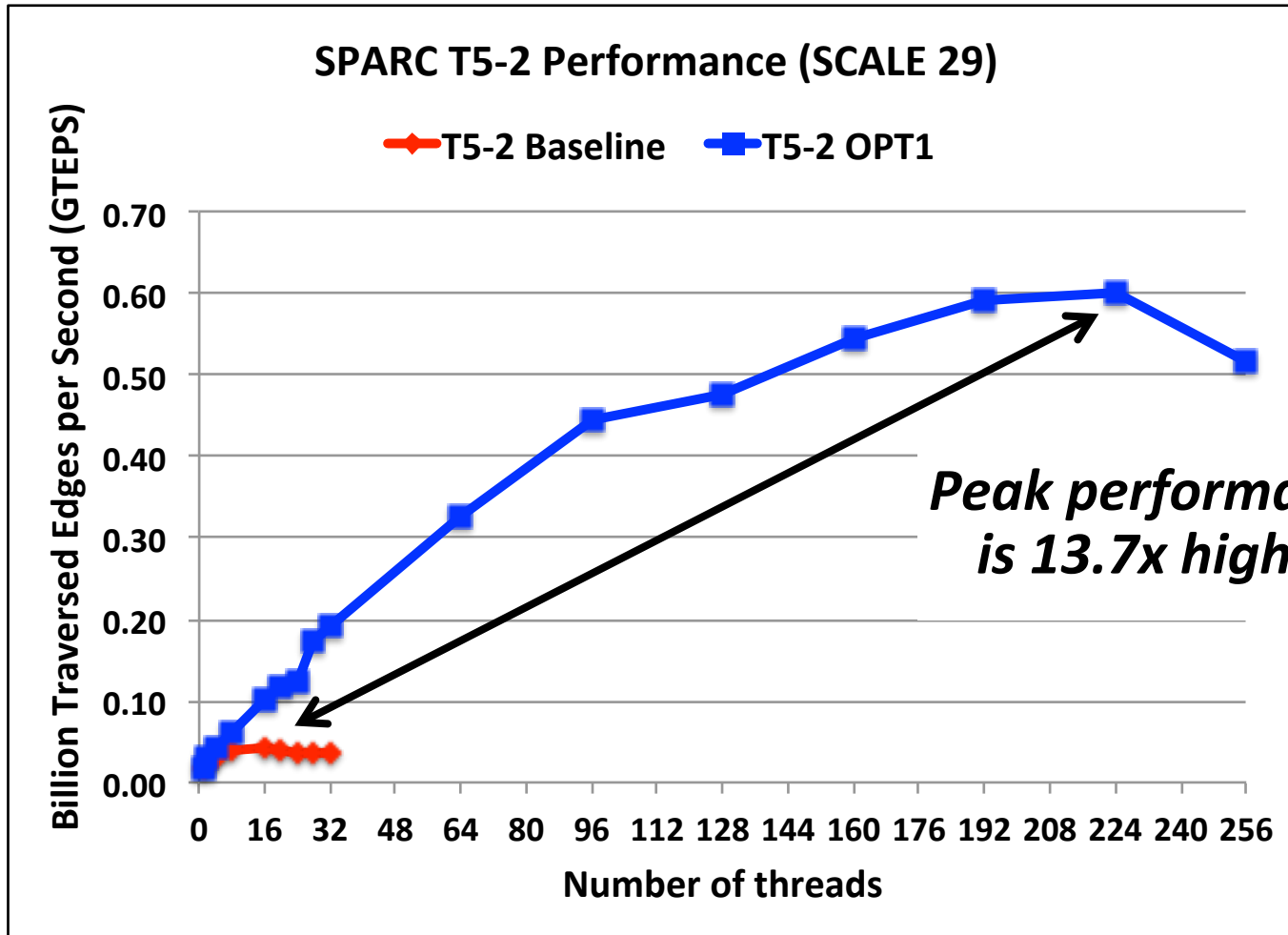
ORACLE®

**Secret Sauce**

BO → BO

4

# Comparison Of The Two Versions



*Note the much shorter run time for the modified version*

Legend:
- **Atomic/critical section wait** (red)
- **Implicit barrier** (magenta)
- **Idle state** (blue)

# Performance Comparison



**SPARC T5-2 Performance (SCALE 29)**
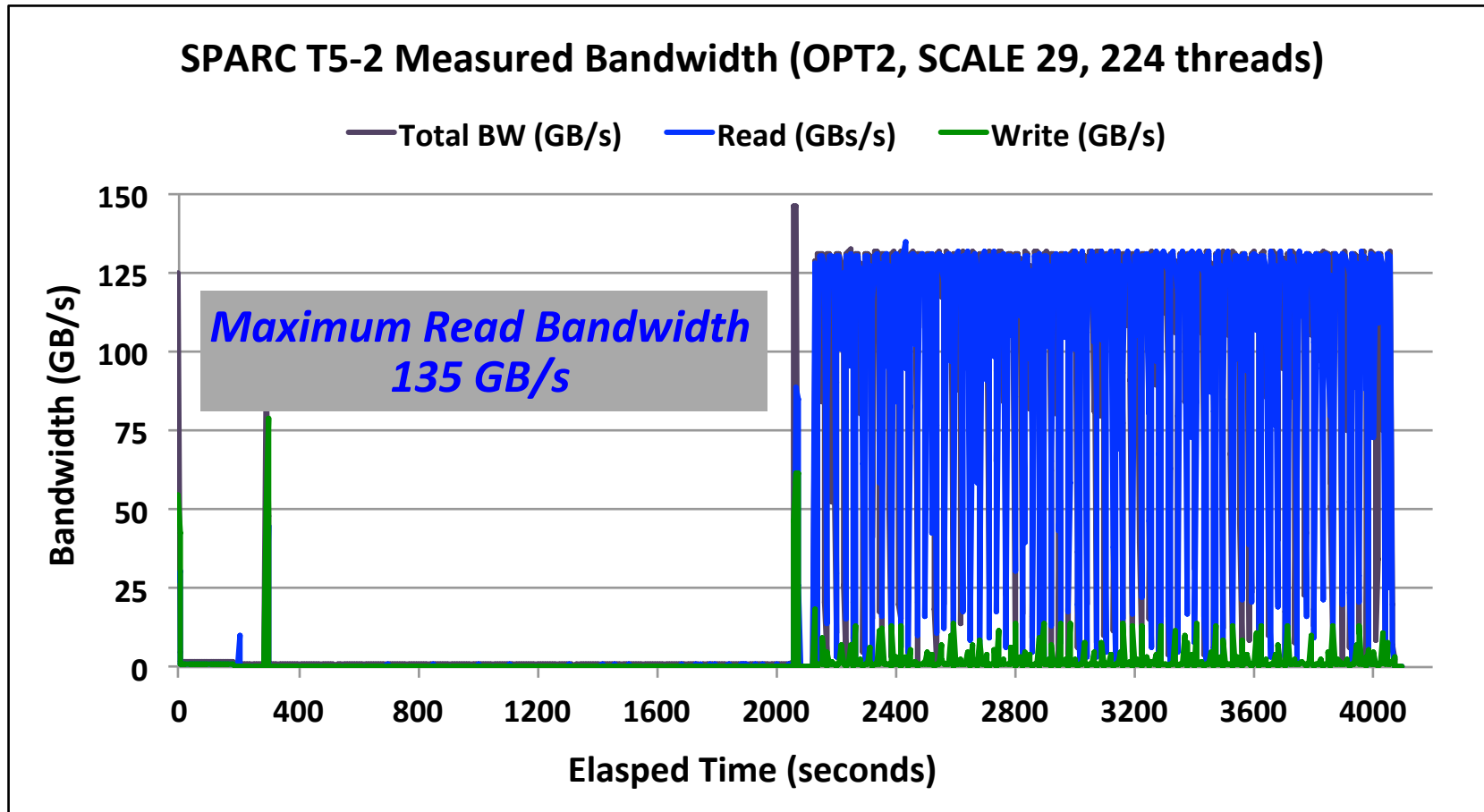
*Peak performance is 13.7x higher*

# Observations

First Touch Placement Is Not Used

The Code Does Not Exploit Large Pages

But Needs It ….

Used A Smarter Memory Allocator

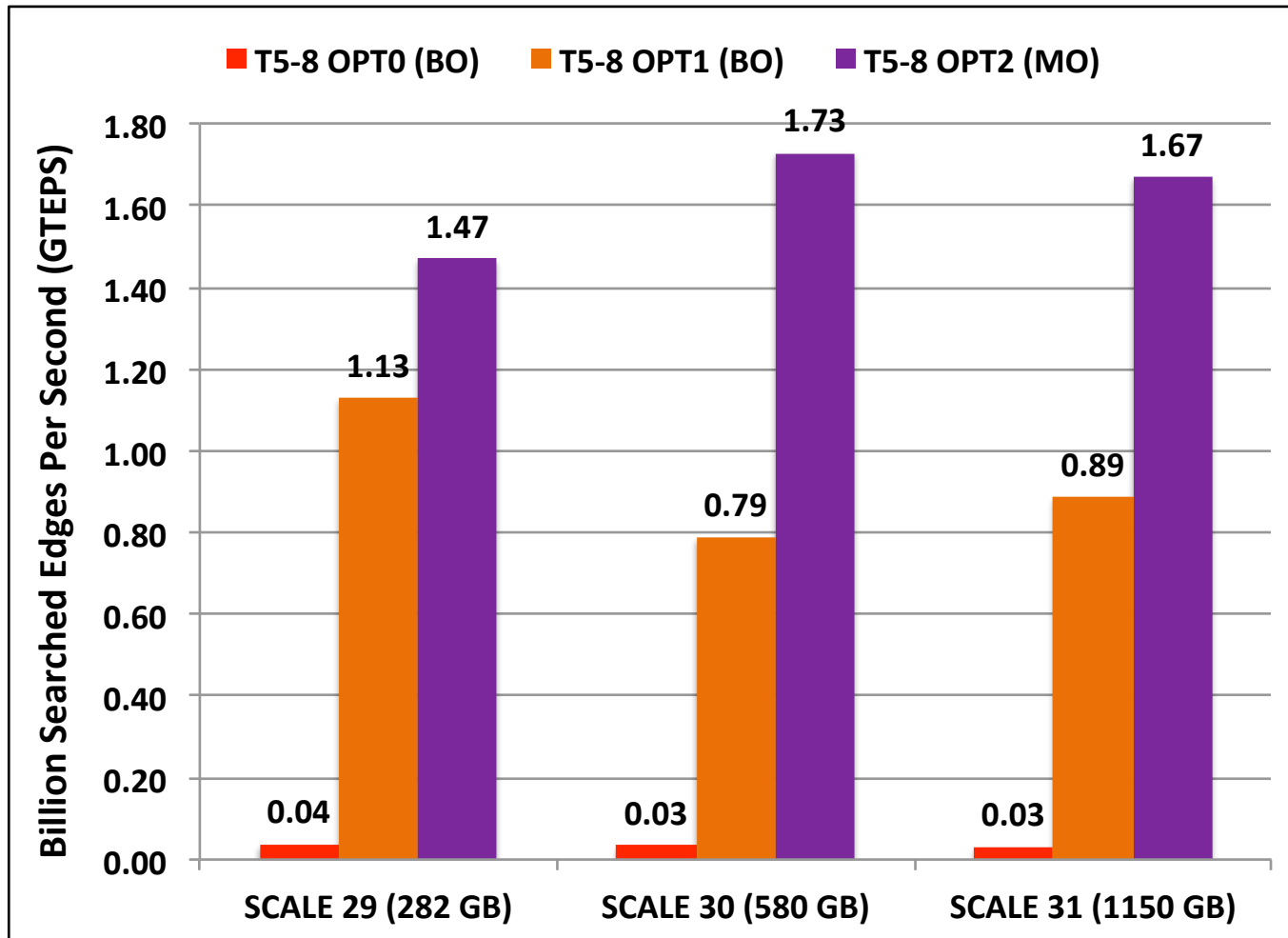# Bandwidth Of The New Code



SPARC T5-2 Measured Bandwidth (OPT2, SCALE 29, 224 threads)

Maximum Read Bandwidth 135 GB/s

# The Result

# Bigger Is Definitely Better!

**72x speed up !**

SPARC T5-8 Graph Search Performance
(SCALE 30 - Memory Requirements: 580 GB)

■ Search Time (minutes)　—■— Speed Up

*Search time reduced from 12 hours to 10 minutes*

**Elapsed Time (minutes)**

800
700
600
500
400
300
200
100
0

80
70
60
50
40
30
20
10
0

1　2　4　8　16　32　64　128　256　384　512

**Number of threads**

ORACLE®

# Tuning Benefit Breakdown



SPARC T5-8 Speed Up Over OPT0

OPT1    OPT2

Bigger is better

Somewhat diminishing return

# Different
# Secret Sauce

**MO** → **MOBO**

# A Simple OpenMP Change

**57-75x improvement**



Chart: Billion Searched Edges Per Second (GTEPS)

Legend: T5-8 OPT0 (BO), T5-8 OPT1 (BO), T5-8 OPT2 (MO), T5-8 OPT2 (MOBO)

- SCALE 29 (282 GB): 0.04, 1.13, 1.47, 2.15
- SCALE 30 (580 GB): 0.03, 0.79, 1.73, 2.43
- SCALE 31 (1150 GB): 0.03, 0.89, 1.67, 2.42

ORACLE

53

# *"I Value My Personal Space"*

# My Favorite Simple Algorithm

```
void mxv(int m,int n,double *a,double *b[],
         double *c)
{
   for (int i=0; i<m; i++)           parallel loop
   {
      double sum = 0.0;
      for (int j=0; j<n; j++)
          sum += b[i][j]*c[j];
      a[i] = sum;
   }
}
```
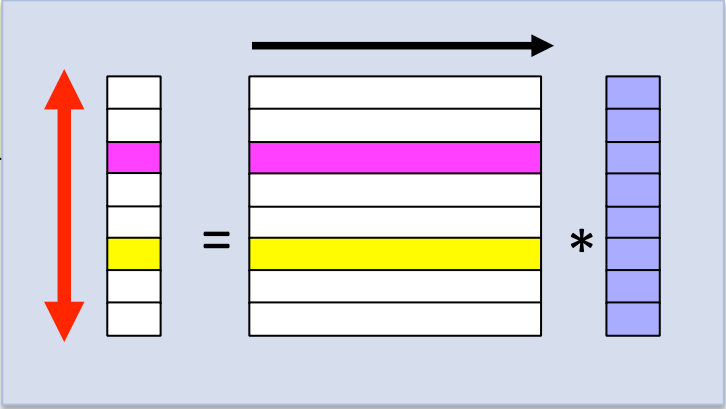
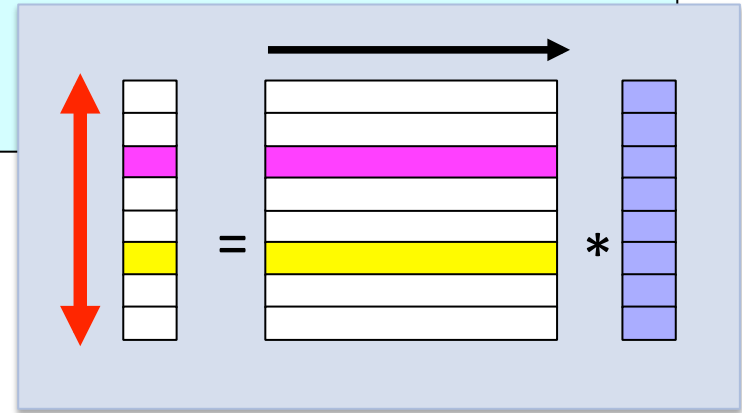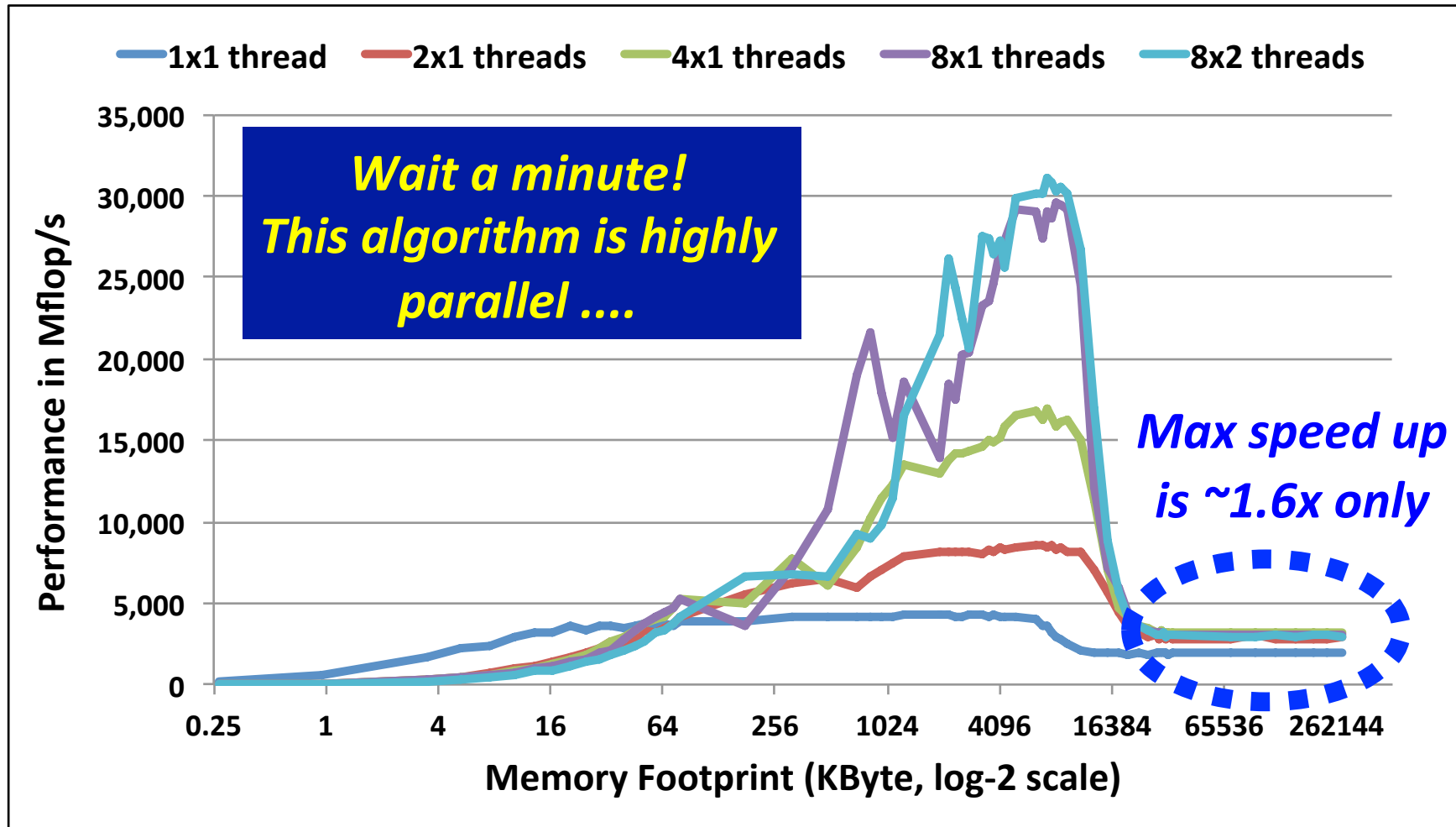# The OpenMP Source

```
#pragma omp parallel for default(none) \
        shared(m,n,a,b,c)
for (int i=0; i<m; i++)
{
  double sum = 0.0;
  for (int j=0; j<n; j++)
     sum += b[i][j]*c[j];
  a[i] = sum;
}
```

# Performance On Intel Nehalem



System: Intel X5570 with 2 sockets,
8 cores, 16 threads at 2.93 GHz

Notation: Number of cores x
number of threads within core

57

**?**

# Let's Get Technical

# A Two Socket Nehalem System



socket 0

| | | | hw thread 0 |
| --- | --- | --- | --- |
| | caches | core 0 | hw thread 1 |
| | caches | core 1 | hw thread 0 |
| shared cache | | | hw thread 1 |
| | caches | core 2 | hw thread 0 |
| | | | hw thread 1 |
| | caches | core 3 | hw thread 0 |
| | | | hw thread 1 |

socket 1

| | | | hw thread 0 |
| --- | --- | --- | --- |
| | caches | core 0 | hw thread 1 |
| | caches | core 1 | hw thread 0 |
| shared cache | | | hw thread 1 |
| | caches | core 2 | hw thread 0 |
| | | | hw thread 1 |
| | caches | core 3 | hw thread 0 |
| | | | hw thread 1 |

QPI Interconnect

memory

memory

60

ORACLE

# Data Initialization Revisited

```
#pragma omp parallel default(none) \
        shared(m,n,a,b,c) private(i,j)
{
#pragma omp for
    for (j=0; j<n; j++)
        c[j] = 1.0;

#pragma omp for
    for (i=0; i<m; i++)
    {
        a[i] = -1957.0;
        for (j=0; j<n; j++)
            b[i][j] = i;
    } /*-- End of omp for --*/

} /*-- End of parallel region --*/
```
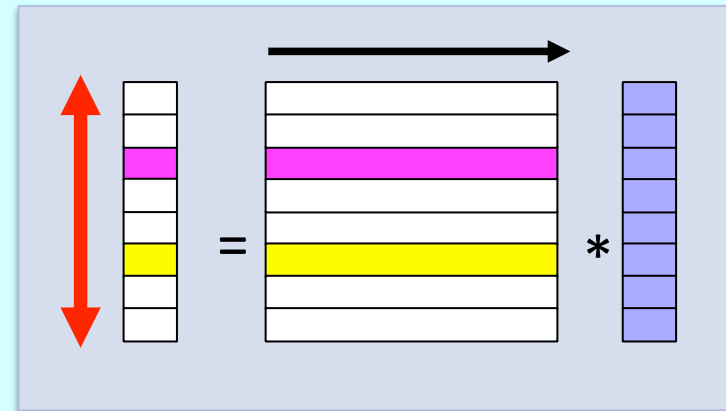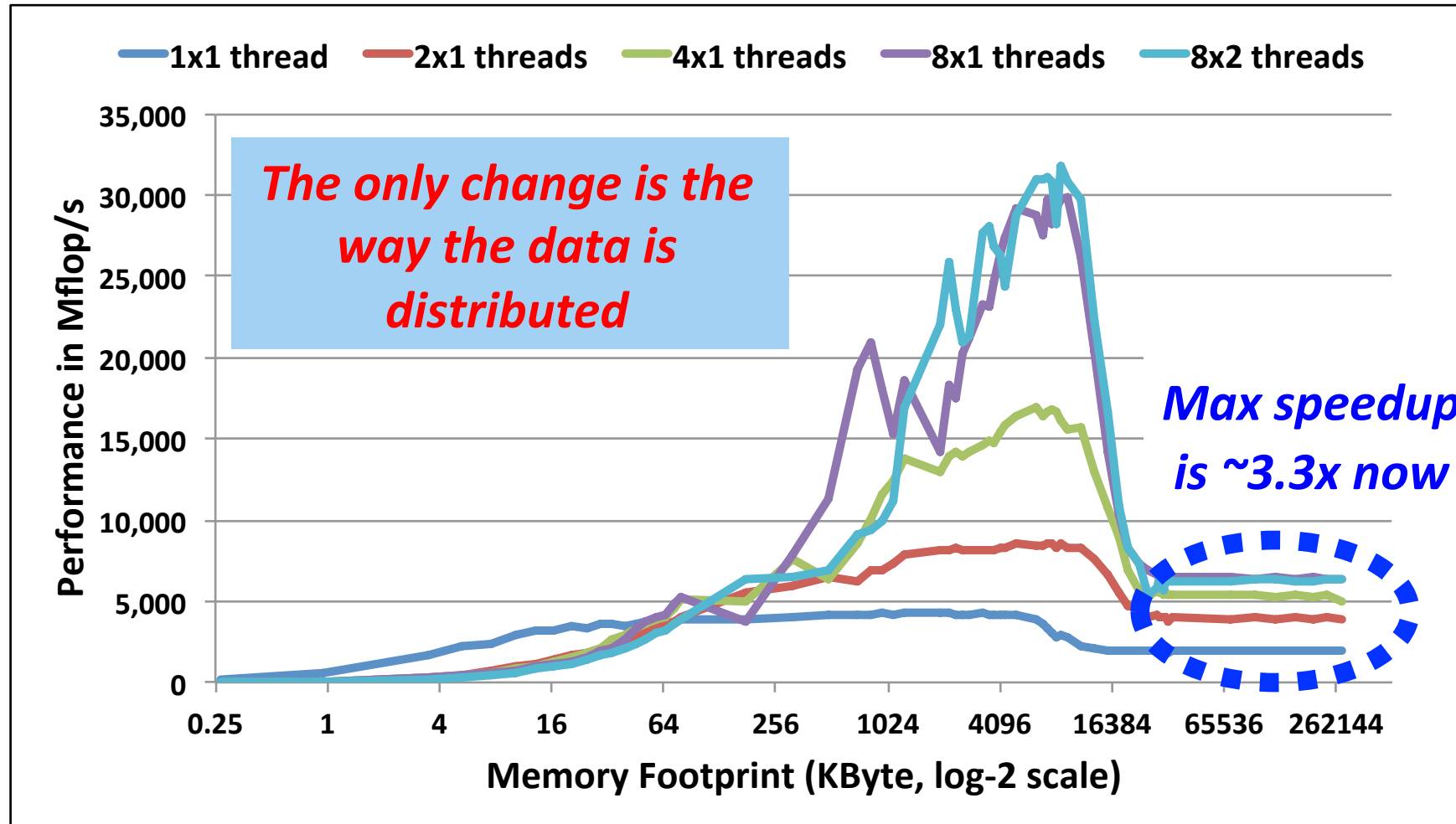
ORACLE

# Data Placement Matters!



**Performance in Mflop/s** vs **Memory Footprint (KByte, log-2 scale)**

Legend: 1x1 thread, 2x1 threads, 4x1 threads, 8x1 threads, 8x2 threads

*The only change is the way the data is distributed*
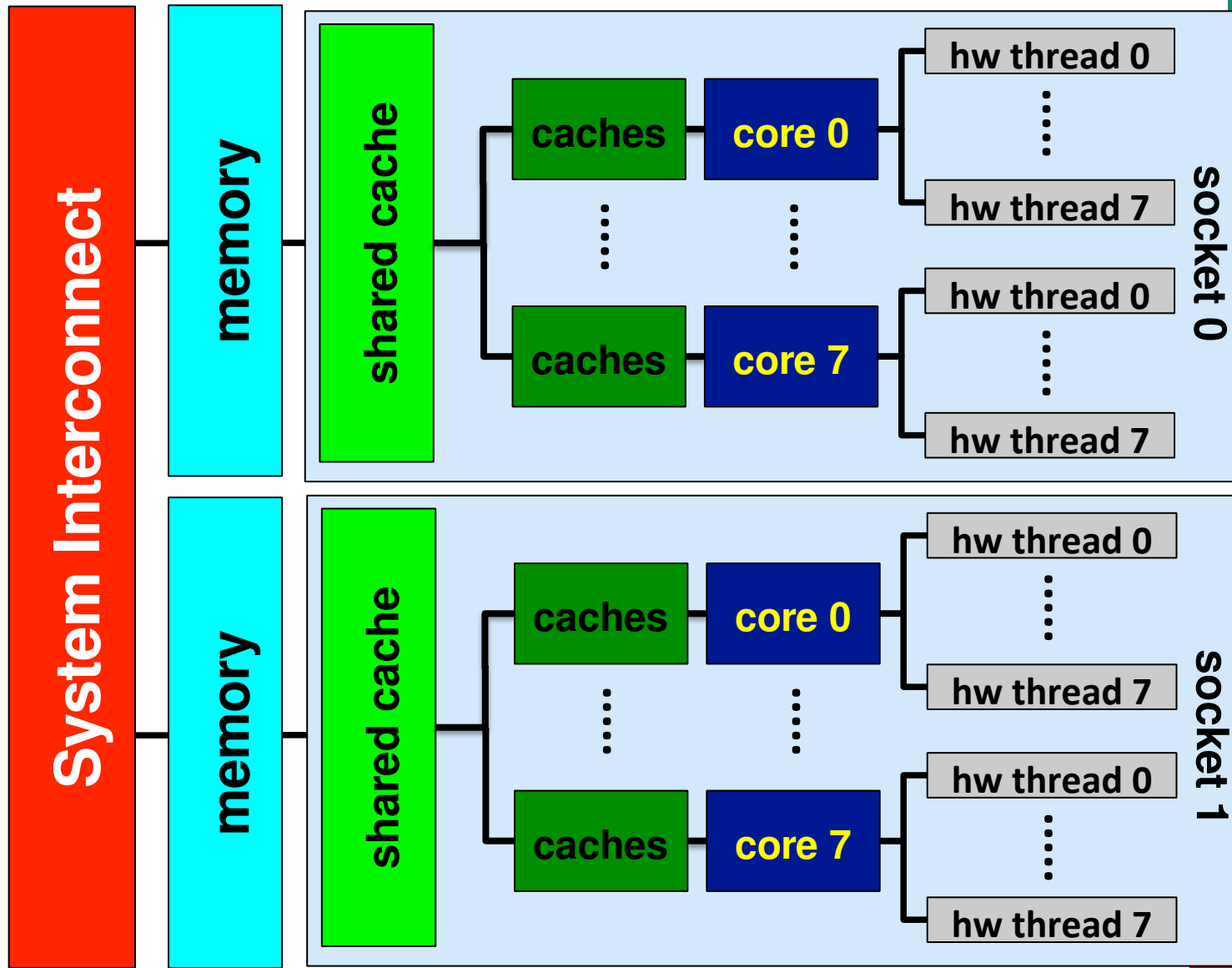
*Max speedup is ~3.3x now*

*System: Intel X5570 with 2 sockets, 8 cores, 16 threads at 2.93 GHz*

*Notation: Number of cores x number of threads within core*

62

# A Two Socket SPARC T4-2 System

ORACLE®

# Performance On SPARC T4-2



**Scaling on larger matrices is affected by cc-NUMA effects (similar as on Nehalem)**

*Note that there are no idle cycles to fill here*

*Max speed up is ~5.8x*

Legend: 1x1 Thread, 8x1 Threads, 16x1 Threads, 16x2 Threads

Y-axis: Performance in Mflop/s — 0, 5,000, 10,000, 15,000, 20,000, 25,000, 30,000, 35,000

X-axis: Memory Footprint (KByte, log-2 scale) — 0.25, 1, 4, 16, 64, 256, 1024, 4096, 16384, 65536, 262144

**System: SPARC T4 with 2 sockets, 16 cores, 128 threads at 2.85 GHz**

**Notation: Number of cores x number of threads within core**

64

ORACLE

# Data Placement Matters!



System: SPARC T4 with 2 sockets, 16 cores, 128 threads at 2.85 GHz
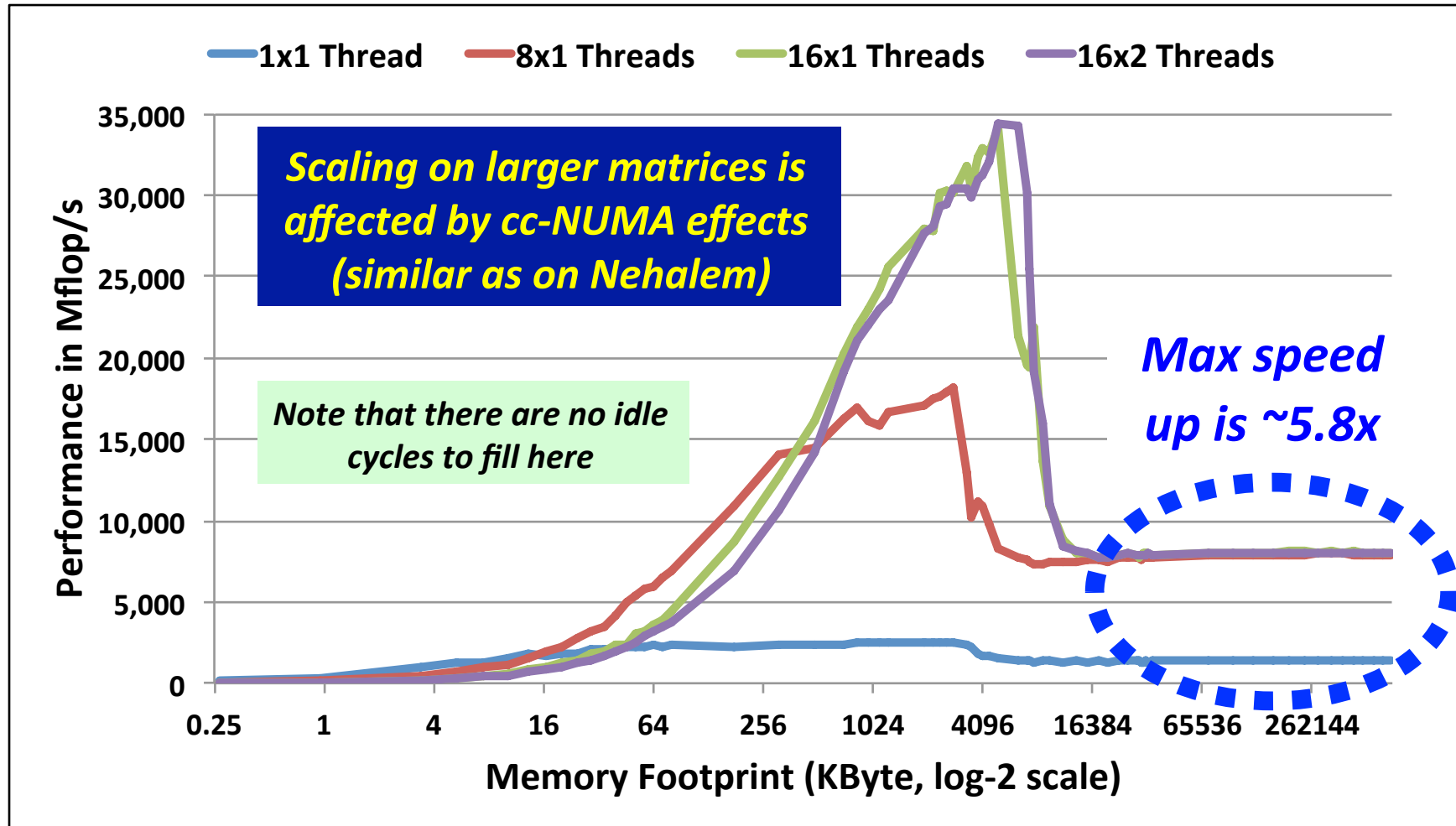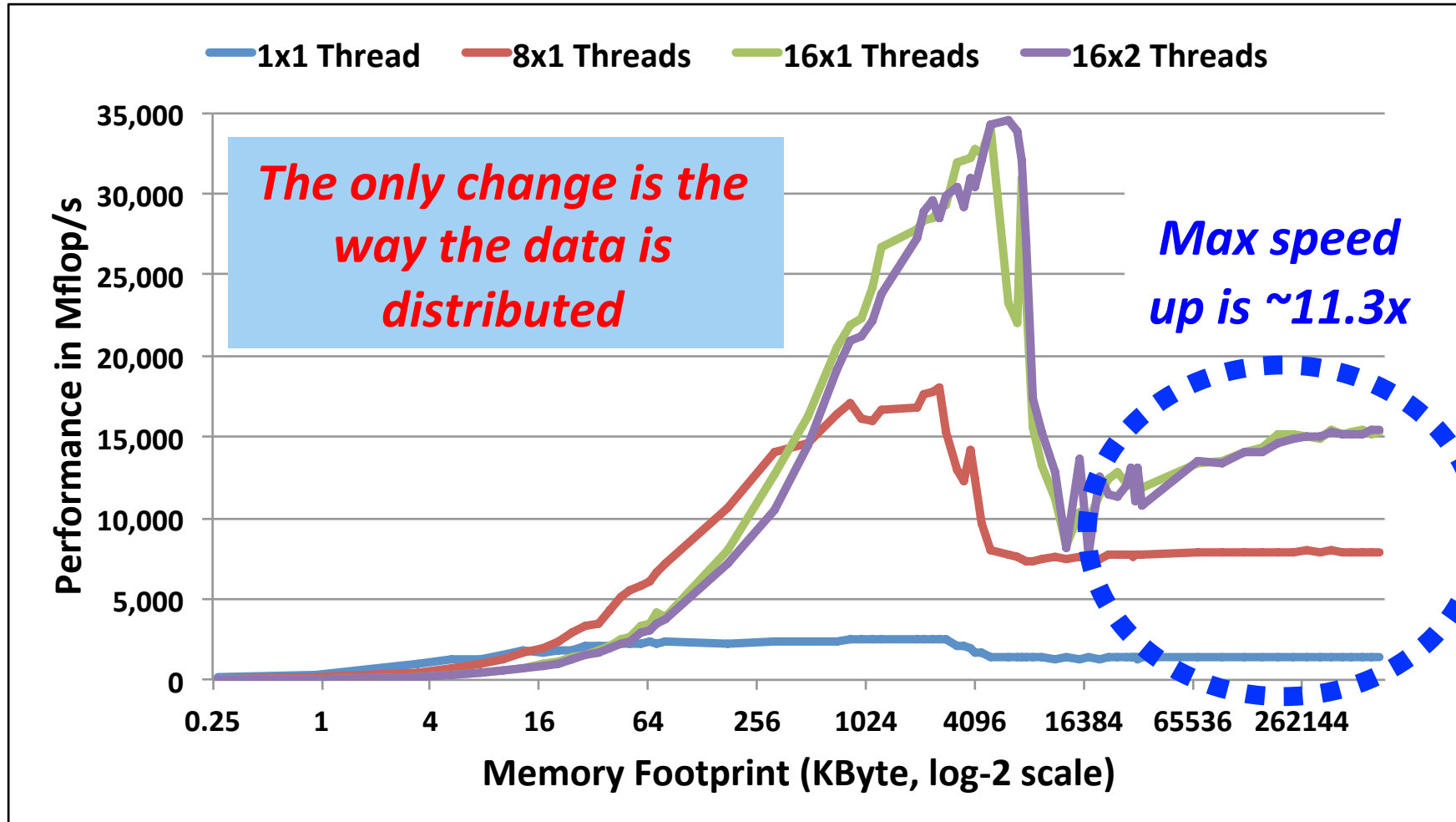
Notation: Number of cores x number of threads within core

# Summary Matrix Times Vector

# The Wrapping

# Wrapping Things Up

*"While we're still waiting for your MPI debug run to finish, I want to ask you whether you found my information useful."*

*"Yes, it is overwhelming. I know."*

*"And OpenMP is somewhat obscure in certain areas. I know that as well."*

ORACLE®

# Wrapping Things Up

*"I understand. You're not a Computer Scientist and just need to get your scientific research done."*

*"I agree this is not a good situation, but it is all about Darwin, you know. I'm sorry, it is a tough world out there."*

ORACLE®

# It Never Ends

*"Oh, your MPI job just finished! Great."*

*"Your program does not write a file called 'core' and it wasn't there when you started the program?"*

*"You wonder where such a file comes from? Let's talk, but I need to get a big and strong coffee first."*

*"WAIT! What did you just say?"*

# It Really Never Ends

*"Somebody told you WHAT ??"*

*"You think GPUs and OpenCL will solve all your problems?"*

*"Let's make that an XL Triple Espresso. I'll buy"*

# Thank You And ..... Stay Tuned !

*ruud.vanderpas@oracle.com*

**Bad OpenMP Does Not Scale**

ORACLE

# *DTrace*
# *Why It Can Be Good For You*

**Ruud van der Pas**

**Distinguished Engineer**

*Architecture and Performance*

*SPARC Microelectronics, Oracle*

*Santa Clara, CA, USA*

ORACLE®

# Motivation

- DTrace is a **Dynamic Tracing Tool**

    → Supported on Solaris, Mac OS X and some Linux flavours

- Monitors the Operating System (OS)

- Through "probes", the user can see what the OS is doing

- Main target: OS related performance issues

- Surprisingly, it can also greatly help to find out what your application is doing though *

*) *A regular profiling tool should be used first*

ORACLE®

# How DTrace Works

- A DTrace probe is written by the user

```
provider:module:function:name
```

- When the probe "fires", the code in the probe is executed

- The probes are based on "providers"

- The providers are pre-defined

  → Example: "sched" provider to get info from the scheduler

  → You can also instrument your own code to have DTrace probes, but there is little need for that

ORACLE®

# Example – Thread Affinity/1

```
sched:::off-cpu
/ pid == $target && self->on_cpu == 1 /        ← "predicate"
{
   self->time_delta = (timestamp - ts_base)/1000;

   @thread_off_cpu [tid-1]     = count();
   @total_thr_state[probename] = count();

   printf("Event %8u %4u %6u %6u    %-16s %8s\n",
          self->time_delta, tid-1, curcpu->cpu_id,
          curcpu->cpu_lgrp, probename, probefunc);

   self->on_cpu     = 0;
   self->time_delta = 0;
}
```

ORACLE®

# Example – Thread Affinity/2

```
pid$target::processor_bind:entry
/ (processorid_t) arg2 >= 0 /
{
   self->time_delta       = (timestamp - ts_base)/1000;
   self->target_processor = (processorid_t) arg2;

   @proc_bind_info[tid-1, curcpu->cpu_id,
                   self->target_processor] = count();

   printf("Event %8u %4u %6u %6u   %9s/%-6d %8s\n",
          self->time_delta, tid-1, curcpu->cpu_id,
          curcpu->cpu_lgrp, "proc_bind", self->target_processor,
          probename);

   self->time_delta       = 0;
   self->target_processor = 0;
}
```

# Example – Example Code

```
#pragma omp parallel
{
   #pragma omp single
   {
     printf("Single region executed by thread %d\n",
            omp_get_thread_num());
   } // End of single region
} // End of parallel region
```

```
export OMP_NUM_THREADS=4
export OMP_PLACES=cores
export OMP_PROC_BIND=close
./affinity.d -c ./omp-par.exe
```

# Example – Result

```
====================================================================
                          Affinity Statistics
====================================================================


Thread      On HW Thread        Lgroup      Created Thread      Count
     0               787              7                   1          1
     0               787              7                   2          1
     0               787              7                   3          1


Thread    Running on HW Thread        Bound to HW Thread
     0                     787                       784
     1                     771                       792
     2                     848                       800
     3                     813                       808
```

# *Thank You And ..... Stay Tuned !*

*ruud.vanderpas@oracle.com*



**Bad OpenMP Does Not Scale**

ORACLE®