(intel)

# Intel® Xeon 5500 Platforms, Integrated Memory Controllers and NUMA

**David Levinthal, Julia Fedorova, Dmitry Ryabtsev**
**SSG/DPD/PAT**

---

## Agenda

**NUMA and Enabling: Overview**

**Topology Overview**

**BIOS Options**

**OS dependent NUMA concerns**

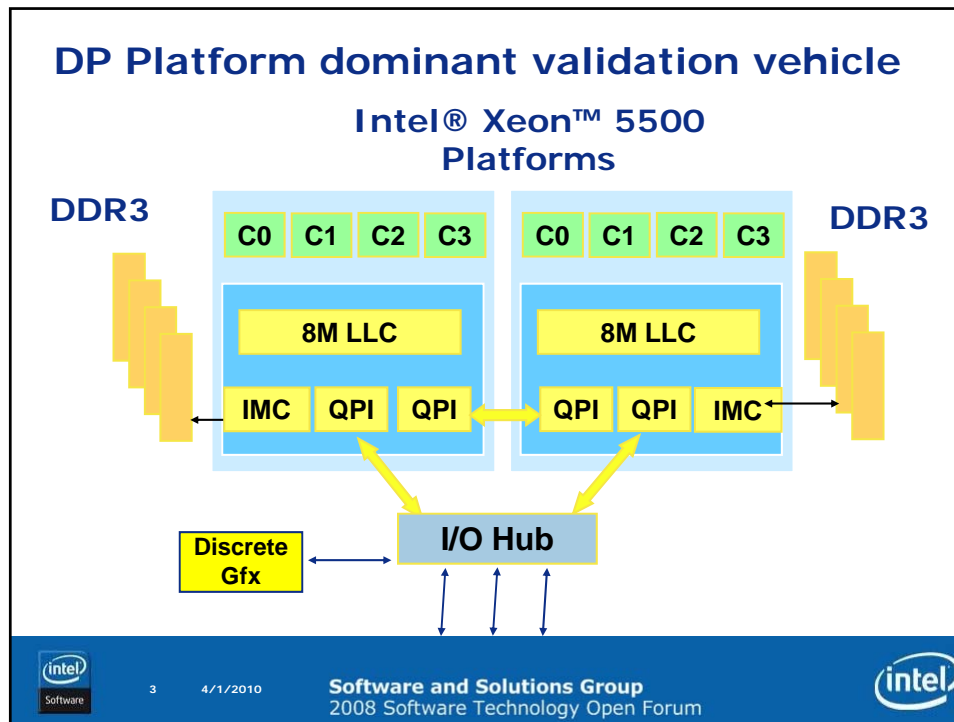**Identifying memory locality (and lack thereof) on Intel® Xeon 5500 processors**

**Summary**

(intel)
Software

2    4/1/2010

**Software and Solutions Group**
2008 Software Technology Open Forum

(intel)

# DP Platform dominant validation vehicle

## Intel® Xeon™ 5500 Platforms

**DDR3**

| C0 | C1 | C2 | C3 |

| C0 | C1 | C2 | C3 |

**DDR3**

**8M LLC**

**8M LLC**

| IMC | QPI | QPI |

| QPI | QPI | IMC |

**Discrete Gfx**

**I/O Hub**

3    4/1/2010    **Software and Solutions Group**
2008 Software Technology Open Forum

(intel)

---

# NUMA, Quickpath and
# Intel® Xeon™ 5500 Platforms

Quickpath Interfaces greatly increase memory bandwidth of our platforms

Integrated memory controllers on each socket access dimms

• Quickpath interconnctions provide cache coherency

• Bandwidth improves by ~4X

Bandwidth improvement comes at a price

• Non uniform memory access

• Latency to dimms on remote sockets is ~2X larger

**Pealing away the Bandwidth layer reveals the NUMA Latency layer**

4    4/1/2010    **Software and Solutions Group**
2008 Software Technology Open Forum

(intel)

# NUMA Modes on DP Systems Controlled in BIOS
## Non Numa

- **Even/Odd lines assigned to sockets 0/1**
  - Line interleaving

## NUMA mode

- **First Half of memory space on socket 0**
- **Second half on socket 1**
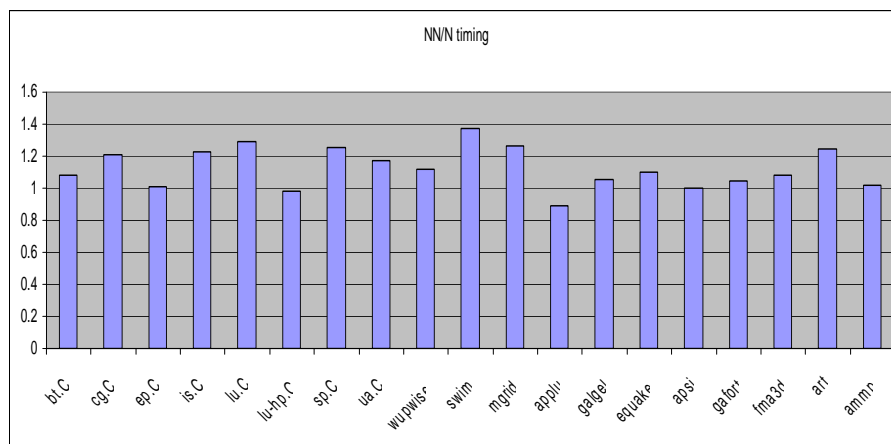- **Default on Intel® Xeon™ 5500 Processors**

(intel) Software    5    4/1/2010    **Software and Solutions Group**
2008 Software Technology Open Forum    (intel)

# NON-NUMA/NUMA Timings for Specomp* and NAS* Parallel Benchmarks

NN/N timing

bt.C, cg.C, ep.C, is.C, lu.C, lu-hp.C, sp.C, ua.C, wupwise, swim, mgrid, applu, galgel, equake, apsi, gafort, fma3d, art, ammp

* Other names and brands may be claimed as the property of others.

(intel) Software    6    4/1/2010    **Software and Solutions Group**
2008 Software Technology Open Forum    (intel)

**3**

# Non Uniform Memory Access and Parallel Execution

**Process parallel is intrinsically NUMA friendly**

- **Affinity pinning maximizes local memory access**
- **MPI**
- **Parallel submission to batch queues**
- **Standard for HPC**

**Shared memory threading is more problematic**

- **Explicit threading, TBB, openMP***
- **NUMA friendly data decomposition (page based) has not been required**
- **OS scheduled thread migration can aggravate situation**

\* Other names and brands may be claimed as the property of others.

7  4/1/2010  **Software and Solutions Group**
2008 Software Technology Open Forum

(intel) Software  (intel)

---

# HPC Applications will see Large Performance Gains due to Bandwidth Improvements

**A remaining performance bottleneck may be due to non uniform memory access latency**

**Intel® PTU data access profiling feature was designed to address NUMA**

- **Intel® Xeon™ 5500 processors events were designed to provide the required data**

8  4/1/2010  **Software and Solutions Group**
2008 Software Technology Open Forum

(intel) Software  (intel)

**Data Access Events on Intel® Xeon™ 5500 processors Reveal NUMA Access Pattern**

**"miss" events are inclusive**
– **Sum over all data sources and their individual latencies**

**Intel® Xeon™ 5500 processor Precise events are exclusive**

**Per data source**

9    4/1/2010    **Software and Solutions Group**
2008 Software Technology Open Forum

---

**Data Access Events Reveal NUMA Access Pattern**



10    4/1/2010    **Software and Solutions Group**
2008 Software Technology Open Forum

## Controlling NUMA Data Locality on Linux* and Windows*

### Linux* assigns physical pages on "first touch"
- ie buffer initialization not malloc
- If each thread initializes its data, things are good
- Can also use numactl or numalib

### Windows assigns physical pages with "allocation"
- VirtualAlloc works like malloc on Linux*
  - Physical pages assigned at first use
- malloc & VirtualAllocExNuma allocation must be parallelized
  - Buffers are no longer contiguous linear address ranges
  - Much MUCH harder

* Other names and brands may be claimed as the property of others.

(intel) Software

11     4/1/2010     **Software and Solutions Group**
2008 Software Technology Open Forum

(intel)

---

## Data Locality, Threaded Applications and Bandwidth

**Consider a threaded triad**

```
int triad(int len, double *a, double *b,
                              double *c, double *x);
     int i,bytes = 24;
     #pragma omp parallel
     {
     #pragma omp for private (i)
     #pragma vector nontemporal
     for(i=0;i<len;i++)a[i]=b[i]+x*c[i];
     }
     return bytes
```

**Parallelizes the work**
> function called 1000 times, len=8192000
> ~ 1B cachelines written NT, 2B read

(intel) Software

12     4/1/2010     **Software and Solutions Group**
2008 Software Technology Open Forum

(intel)

**Data Locality, Threaded Applications and Bandwidth**

**Run an OpenMP\* triad under my usual mini_app driver, the resulting BW is only**

> **~ 5bytes/cycle for 8 threads**

**Running in Non Numa Mode results in**

> **~8.5 Bytes/cycle**

**Why?**

**Default Version Allocates Buffers on Thread 0**
**Using only one Memory Controller**

\* Other names and brands may be claimed as the property of others.

(intel)
Software
13    4/1/2010    **Software and Solutions Group**
2008 Software Technology Open Forum
(intel)

---

**Performance Events and NUMA Sources**

- **Offcore_Response_0**
  **8 flavors of Request Type X 8 flavors of $line Source**
  - **+ all combinations..**
    **(~65K possible programmings)**

- **One "gotcha"...**
  **NT stores to local Dram**
  **appear to go to another core's cache**
  **(data source = 2 instead of 0x40)**

(intel)
Software
14    4/1/2010    **Software and Solutions Group**
2008 Software Technology Open Forum
(intel)

## PTU Display Shows Local and Remote Access for OpenMP Triad

**Software and Solutions Group**
2008 Software Technology Open Forum

---

## Need to Distribute "Allocation" "Allocate" on First Touch

### Original allocation

```
buf1 = (char *) malloc(DIM*(sizeof (double))+1024);
buf2 = (char *) malloc(DIM*(sizeof (double))+1024);
buf3 = (char *) malloc(DIM*(sizeof (double))+1024);
a = (double *) buf1;
b = (double *) buf2;
c = (double *) buf3;
for(num=0;num<len;num++)
{
        a[num]=10.;
        b[num]=10.;
        c[num]=10.;
}
```

## Initialization must also be done in Parallel

* Other names and brands may be claimed as the property of others.

**Software and Solutions Group**
2008 Software Technology Open Forum

## Parallel "Allocation" for Linux*
## Requires Parallel Initialization

Parallel allocation

```
        buf1 = (char *) malloc(DIM*(sizeof (double))+1024);
        buf2 = (char *) malloc(DIM*(sizeof (double))+1024);
        buf3 = (char *) malloc(DIM*(sizeof (double))+1024);
        a = (double *) buf1;
        b = (double *) buf2;
        c = (double *) buf3;
#pragma omp parallel
{
#pragma omp for private(num)
        for(num=0;num<len;num++)
        {
                a[num]=10.;
                b[num]=10.;
                c[num]=10.;
        }

}
```

* Other names and brands may be claimed as the property of others.

17   4/1/2010   **Software and Solutions Group**
2008 Software Technology Open Forum

| Event | Triad_omp | |
|---|---|---|
| CPU_CLK_UNHALTED.THREAD | 2.23E+11 | |
| CPU_CLK_UNHALTED.THREAD;Socket 0 | 7.51E+10 | |
| CPU_CLK_UNHALTED.THREAD;Socket 1 | 1.48E+11 | |
| OFFCORE_RESPONSE_0.ANY_REQUEST.ANY_LOCATION | 3.13E+09 | |
| OFFCORE_RESPONSE_0.ANY_REQUEST.ANY_LOCATION;Socket 0 | 1.56E+09 | |
| OFFCORE_RESPONSE_0.ANY_REQUEST.ANY_LOCATION;Socket 1 | 1.56E+09 | |
| OFFCORE_RESPONSE_0.ANY_REQUEST.LOCAL_CACHE_DRAM | 1.56E+09 | |
| OFFCORE_RESPONSE_0.ANY_REQUEST.LOCAL_CACHE_DRAM; Socket 0 | 1.55E+09 | |
| OFFCORE_RESPONSE_0.ANY_REQUEST.LOCAL_CACHE_DRAM; Socket 1 | 8000000 | |
| OFFCORE_RESPONSE_0.ANY_REQUEST.REMOTE_DRAM | 1.55E+09 | |
| OFFCORE_RESPONSE_0.ANY_REQUEST.REMOTE_DRAM;Socket 0 | 1.55E+09 | |
| OFFCORE_RESPONSE_0.ANY_REQUEST.REMOTE_DRAM;Socket 1 | 100000 | |

**Note socket 0/1 switch between PTU runs**

18   4/1/2010   **Software and Solutions Group**
2008 Software Technology Open Forum

| Event | Triad_omp | Triad_NUMA |
|---|---|---|
| CPU_CLK_UNHALTED.THREAD | 2.23E+11 | 1.17E+11 |
| CPU_CLK_UNHALTED.THREAD;Socket 0 | 7.51E+10 | 5.84E+10 |
| CPU_CLK_UNHALTED.THREAD;Socket 1 | 1.48E+11 | 5.83E+10 |
| OFFCORE_RESPONSE_0.ANY_REQUEST.ANY_LOCATION | 3.13E+09 | 3.11E+09 |
| OFFCORE_RESPONSE_0.ANY_REQUEST.ANY_LOCATION;Socket 0 | 1.56E+09 | 1.56E+09 |
| OFFCORE_RESPONSE_0.ANY_REQUEST.ANY_LOCATION;Socket 1 | 1.56E+09 | 1.55E+09 |
| OFFCORE_RESPONSE_0.ANY_REQUEST.LOCAL_CACHE_DRAM | 1.56E+09 | 3.11E+09 |
| OFFCORE_RESPONSE_0.ANY_REQUEST.LOCAL_CACHE_DRAM; Socket 0 | 1.55E+09 | 1.55E+09 |
| OFFCORE_RESPONSE_0.ANY_REQUEST.LOCAL_CACHE_DRAM; Socket 1 | 8000000 | 1.55E+09 |
| OFFCORE_RESPONSE_0.ANY_REQUEST.REMOTE_DRAM | 1.55E+09 | 400000 |
| OFFCORE_RESPONSE_0.ANY_REQUEST.REMOTE_DRAM;Socket 0 | 1.55E+09 | 300000 |
| OFFCORE_RESPONSE_0.ANY_REQUEST.REMOTE_DRAM;Socket 1 | 100000 | 100000 |

(intel) Software

19

**5.1 B/cyc vs 8.5 B/cyc vs 12.5 B/cyc
on a poorly tuned machine**
2008 Software Technology Open Forum

(intel)

---

**OpenMP and Core Affinity Pinning**

**Export KMP_AFFINITY=compact,0,verbose
will pin affinity of threads**

**Just not reproducibly (per socket) on Red Hat
5.1 from run to run**

**Causing problems in multi run PTU collections**

**Problem is that an app does not use OMP
runtime libs to pin affinity until there is a
#pragma parallel {}**

**You must add this around first instruction to pin
affinity of Main thread**

(intel) Software

20      4/1/2010      **Software and Solutions Group**
2008 Software Technology Open Forum

(intel)

**Multi-thread Scaling and NUMA**

**When measuring scaling between 4 and 8 threads (assuming no SMT) the affinity of the 4 threads matters**

**4 threads all on one socket has the same LLC cache size/core as 8 threads**

**BUT**

**2 threads/socket has closer to the same memory BW as the 8 thread run**

**Thus 4->8 scaling will always have a non scaling contribution due to one of these 2 effects**
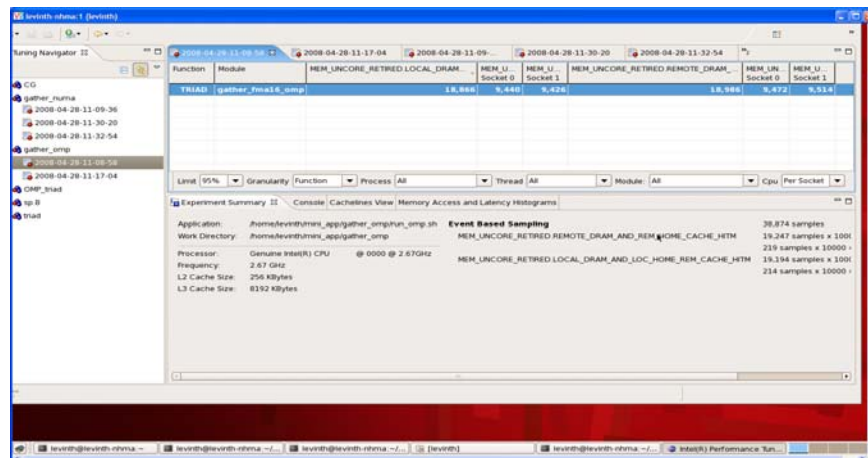
21    4/1/2010    **Software and Solutions Group**
2008 Software Technology Open Forum

**Per Socket Display + Data Source events Show NUMA /Cross Socket Traffic**



22    4/1/2010    **Software and Solutions Group**
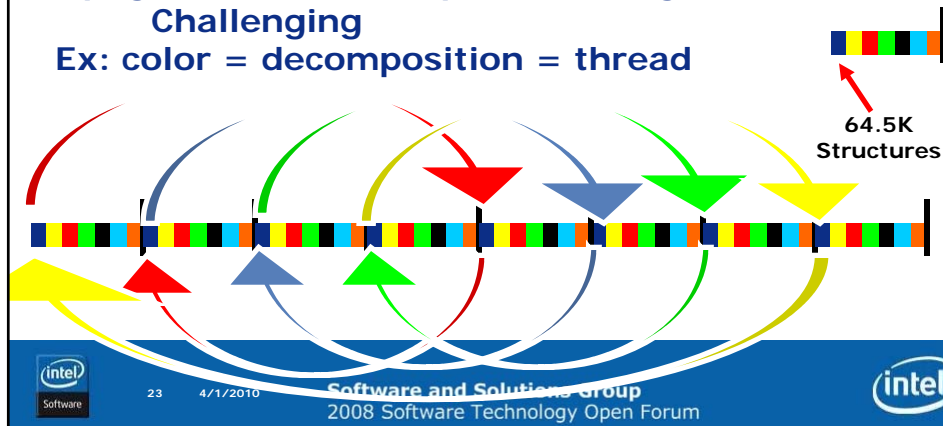2008 Software Technology Open Forum

# Indirect Addressing, Locality and Latency (Diff Eq on Non Uniform Grid, Oil Res)

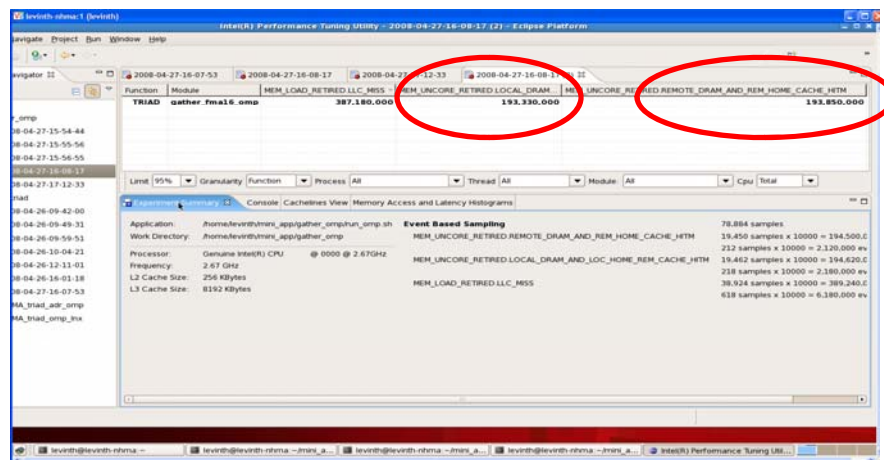**Multi-dimensional array access can cause large address gaps in data decomposition.**

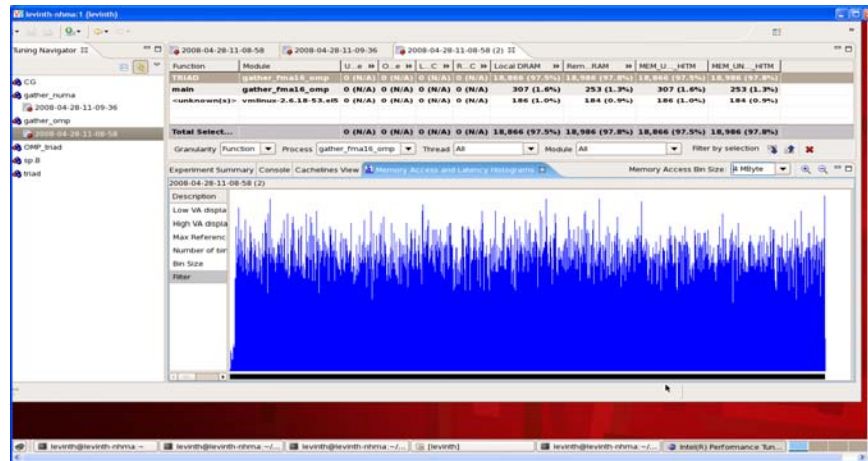**This can make mapping NUMA home node->pages->data decomposition ranges Challenging**
**Ex: color = decomposition = thread**

**64.5K Structures**

23    4/1/2010    Software and Solutions Group
2008 Software Technology Open Forum

# Default Initialization Breaks Array into 8 Contiguous Pieces➔ 50% Non Local Access

24    4/1/2010    Software and Solutions Group
2008 Software Technology Open Forum

# Address Histogram for all Dram Accesses

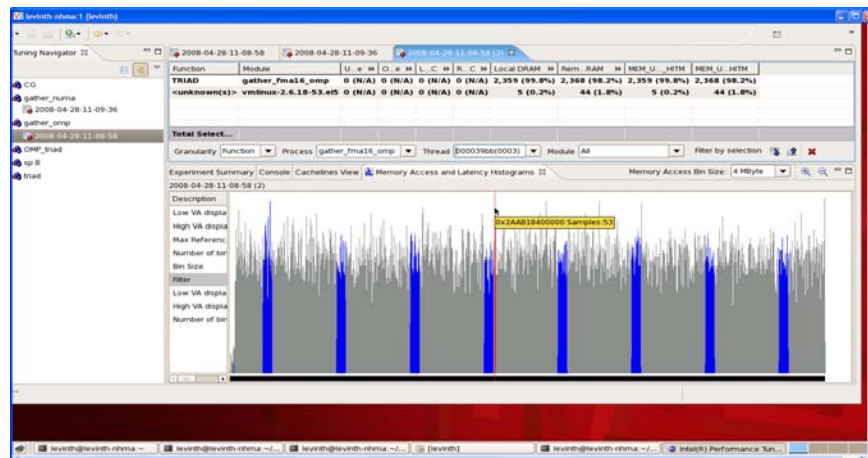**Software and Solutions Group**
2008 Software Technology Open Forum

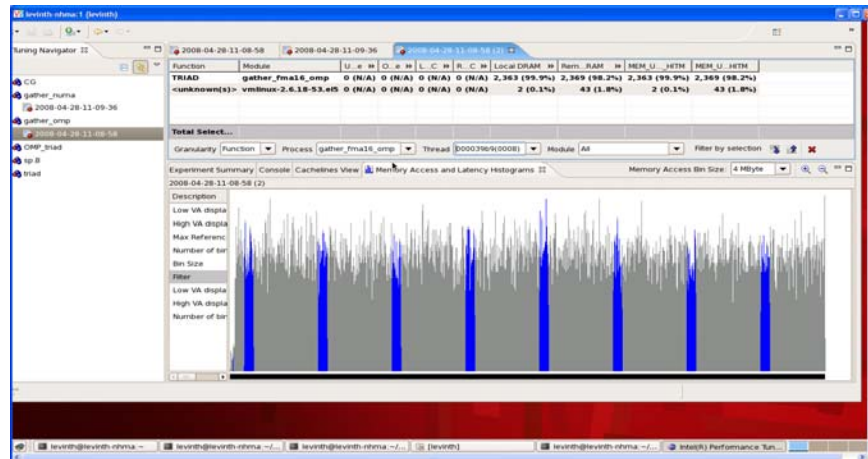# Filtering to a Single Thread Displays the Data Decomposition

**Software and Solutions Group**
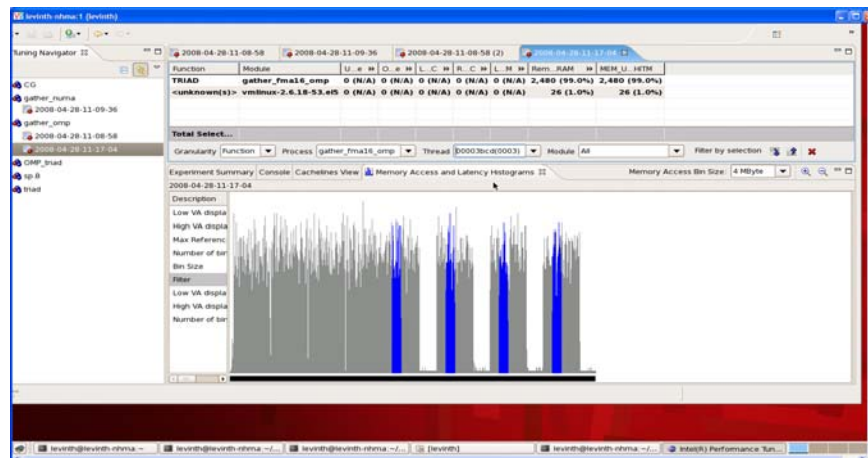2008 Software Technology Open Forum

**A Different Thread**

**Using Only Precise Remote Dram Event**
**Only Half the entries shown**
**Gaps due to lack of events are suppressed**

**Using Only Precise Remote Dram Event**
**Only Half the entries shown**
**Gaps due to lack of events are suppressed**



29    4/1/2010    **Software and Solutions Group**
2008 Software Technology Open Forum

---

# Change Initialization to Follow Work Access Pattern

**Thread initialization with same access sequence as work**

**Expect ~33% improvement**
   – **1/2 of accesses get lower latency by 2**

**Simple OMP ran in 14.3 cycles/cell**

**NUMA initialized version ran in 11.2 cycles/cell**

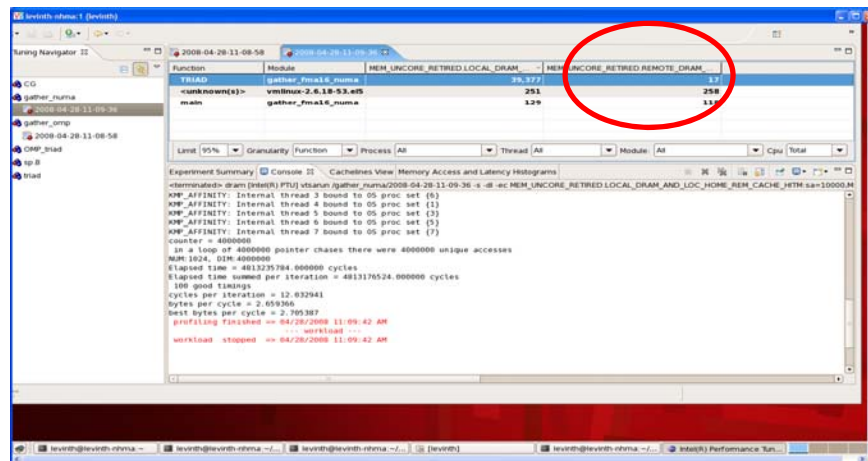**Every access has serious DTLB issues, which don't change with the improved NUMA layout**

30    4/1/2010    **Software and Solutions Group**
2008 Software Technology Open Forum

## Slide 1

# Sampling View for Correctly Initialized Array has no Remote Access



31    4/1/2010    **Software and Solutions Group**
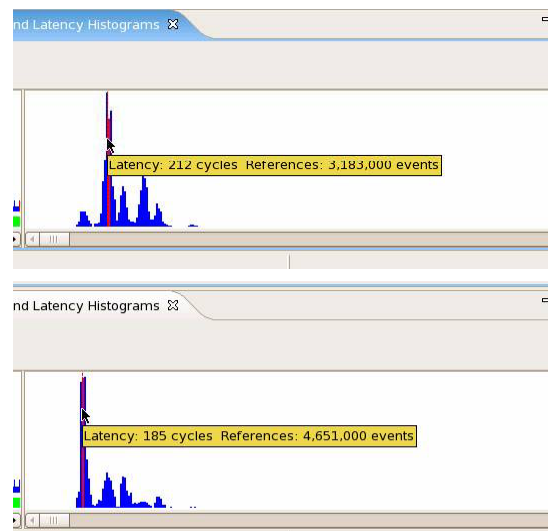2008 Software Technology Open Forum

## Slide 2

# Page Allocation Order Matters

**Serially initialized/allocated**

**Accessed with complex pattern avg Lat =230**

Latency: 212 cycles  References: 3,183,000 events

**Initialized/Allocated and Accessed with complex pattern avg Lat = 209**

Latency: 185 cycles  References: 4,651,000 events

32    4/1/2010    **Software and Solutions Group**
2008 Software Technology Open Forum

**Conclusions**

**NUMA will add complexity to software performance analysis and optimization**

**We have the infrastructure to manage this**

intel)
Software

33   4/1/2010

**Software and Solutions Group**
2008 Software Technology Open Forum

(intel)

---

**Backup**

(intel)
Software

34   4/1/2010

**Software and Solutions Group**
2008 Software Technology Open Forum

(intel)