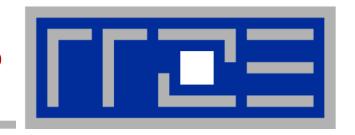
For final slides and example code see:

https://goo.gl/DrGlq5



# Node-Level Performance Engineering



Erlangen Regional Computing Center (RRZE) and Department of Computer Science University of Erlangen-Nuremberg

SC16 full-day tutorial November 13, 2016 Salt Lake City, UT, USA



qrme.com





slide updated from USB version

# **Agenda**



|       | - | Preliminaries   | 08:30 |
|-------|---|---|-------|
| GW    | - | Introduction to multicore architecture                    |       |
|       |   | Threads, cores, SIMD, caches, chips, sockets, ccNUMA      |       |
|       | • | Multicore tools   | 10:00 |
| Н5    | • | Microbenchmarking for architectural exploration           | 10:30 |
|       |   | <ul> <li>Streaming benchmarks</li> </ul>                  |       |
|       |   | <ul><li>Hardware bottlenecks</li></ul>                    |       |
|       | • | Node-level performance modeling (part I)                  |       |
|       |   | The Roofline Model and dense MVM                          | 12:00 |
| MD GH | • | Lunch break   |       |
|       | • | Node-level performance modeling (part II)                 | 13:30 |
|       |   | <ul><li>Case studies: Sparse MVM, Jacobi solver</li></ul> |       |
|       | • | Optimal resource utilization                              |       |
|       |   | <ul> <li>SIMD parallelism</li> </ul>                      | 15:00 |
|       |   | <ul><li>ccNUMA</li></ul>                                  | 15:30 |
|       |   | <ul><li>OpenMP synchronization and multicores</li></ul>   |       |
|       | _ | Pattern-driven performance engineering                    | 17:00 |

#### **A conversation**



From a student seminar on "Efficient programming of modern multi- and manycore processors"

**Student**: I have implemented this algorithm on the GPGPU, and it

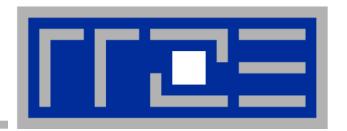
solves a system with 26546 unknowns in 0.12 seconds,

so it is really fast.

**Me**: What makes you think that 0.12 seconds is fast?

**Student**: It is fast because my baseline C++ code on the CPU is about

20 times slower.



# Prelude: Scalability 4 the win!

### Scalability Myth: Code scalability is the key issue



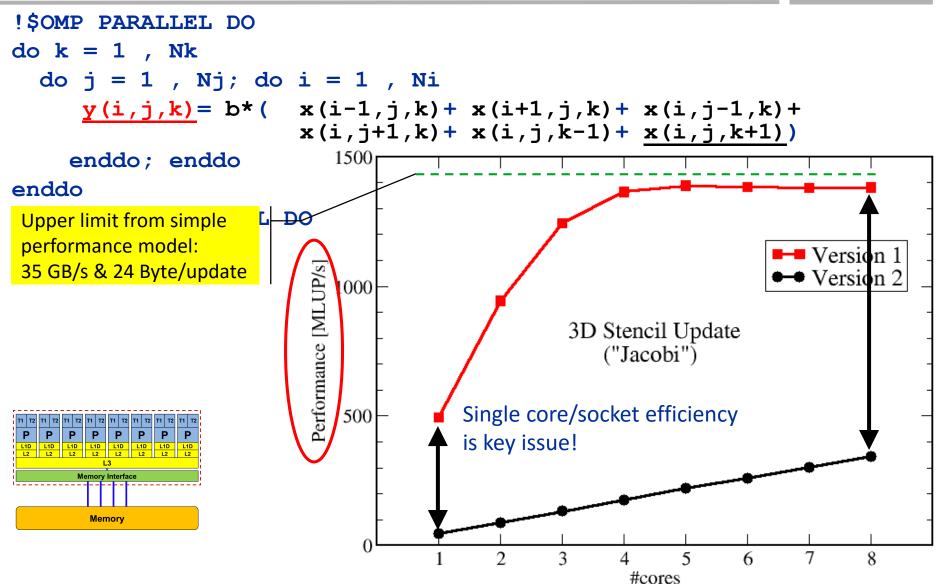
```
!$OMP PARALLEL DO
do k = 1 , Nk
  do j = 1 , Nj; do i = 1 , Ni
      y(i,j,k) = b*(x(i-1,j,k) + x(i+1,j,k) + x(i,j-1,k) +
                        x(i,j+1,k) + x(i,j,k-1) + x(i,j,k+1)
     enddo; enddo
enddo
                                  3D Stencil Update
!$OMP END PARALLEL DO
                                      ("Jacobi")
Changing only the compile
options makes this code
                                   ■ Version 1
                           Speed-Up
                                     Version 2
scalable on an 8-core chip
                                                          Prepared for
                                                          the highly
                                                          parallel era!
                                                           -03 - xAVX
      Memory
```

6

#cores

## Scalability Myth: Code scalability is the key issue





## Questions to ask in high performance computing

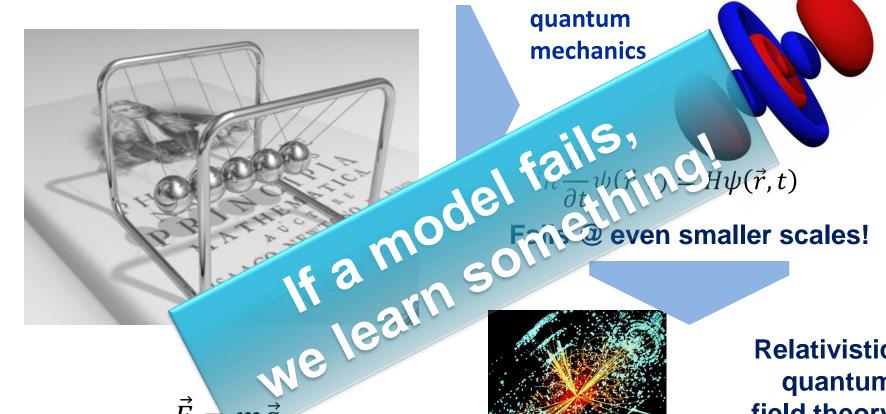


- Do I understand the performance behavior of my code?
  - Does the performance match a model I have made?
- What is the optimal performance for my code on a given machine?
  - High Performance Computing == Computing at the bottleneck
- Can I change my code so that the "optimal performance" gets higher?
  - Circumventing/ameliorating the impact of the bottleneck
- My model does not work what's wrong?
  - This is the good case, because you learn something
  - Performance monitoring / microbenchmarking may help clear up the situation

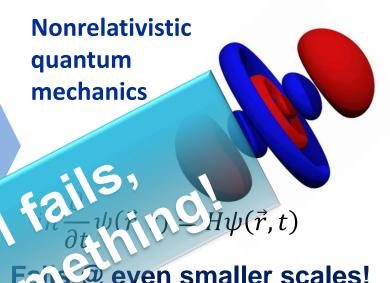
# **How model-building works: Physics**



#### **Newtonian mechanics**



Fails @ small scales!



Relativistic quantum field theory

 $U(1)_{V} \otimes SU(2)_{L} \otimes SU(3)_{C}$ 



# Introduction: Modern node architecture

A glance at basic core features:

pipelining, superscalarity, SMT

Caches and data transfers through the memory hierarchy

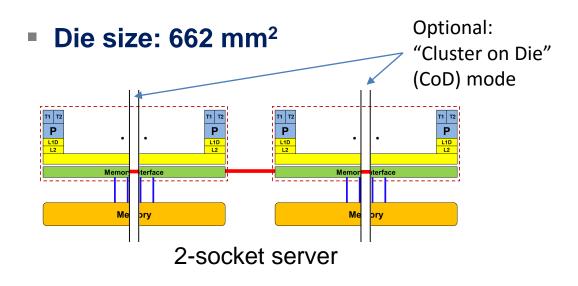
**Accelerators** 

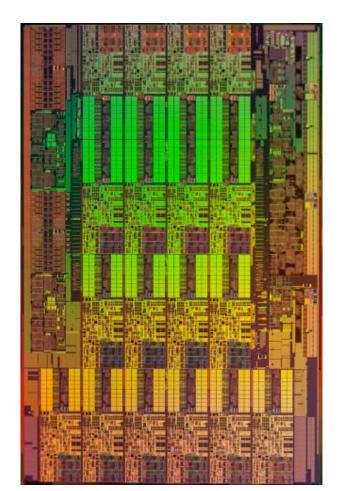
**Bottlenecks & hardware-software interaction** 

## Multi-core today: Intel Xeon 2600v3 (2014)



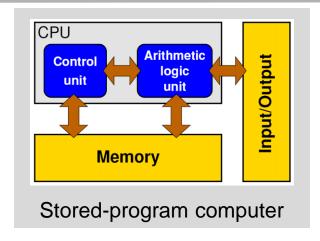
- Xeon E5-2600v3 "Haswell EP":
   Up to 18 cores running at 2+ GHz (+ "Turbo Mode": 3.5+ GHz)
- Simultaneous Multithreading→ reports as 36-way chip
- 5.7 Billion Transistors / 22 nm



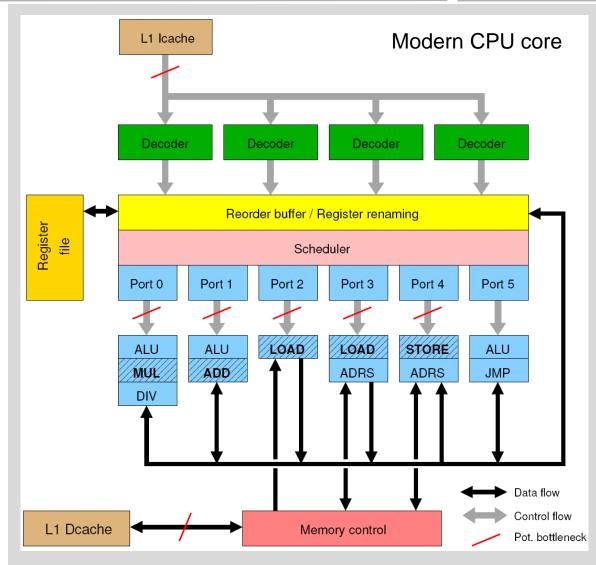


### General-purpose cache based microprocessor core





- Implements "Stored Program Computer" concept (Turing 1936)
- Similar designs on all modern systems
- (Still) multiple potential bottlenecks
- The clock cycle is the "heartbeat" of the core



# Pipelining of arithmetic/functional units



#### Idea:

- Split complex instruction into several simple / fast steps (stages)
- Each step takes the same amount of time, e.g. a single cycle
- Execute different steps on different instructions at the same time (in parallel)

### Allows for shorter cycle times (simpler logic circuits), e.g.:

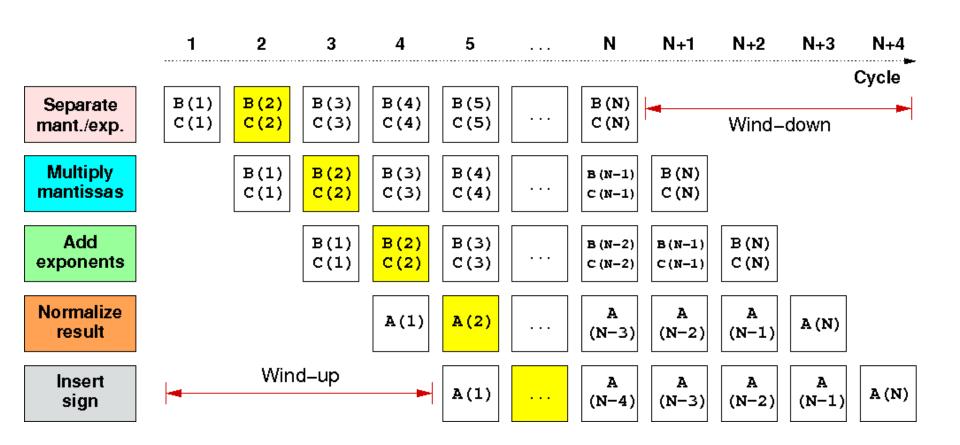
- floating point multiplication takes 5 cycles, but
- processor can work on 5 different multiplications simultaneously
- one result at each cycle after the pipeline is full

#### Drawback:

- Pipeline must be filled startup times (#Instructions >> pipeline steps)
- Efficient use of pipelines requires large number of independent instructions → instruction level parallelism
- Requires complex instruction scheduling by compiler/hardware softwarepipelining / out-of-order
- Pipelining is widely used in modern computer architectures

## 5-stage Multiplication-Pipeline: A(i)=B(i)\*C(i); i=1,...,N





First result is available after 5 cycles (=latency of pipeline)! Wind-up/-down phases: Empty pipeline stages

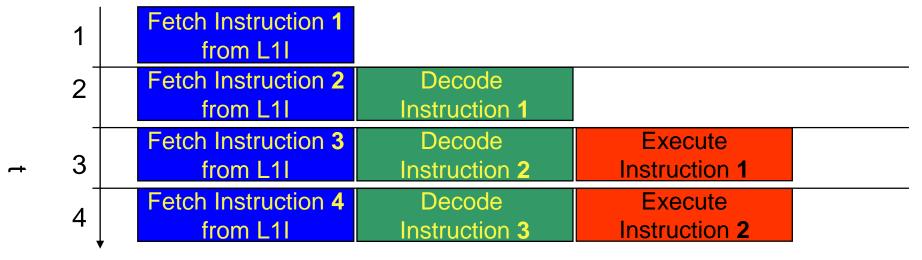
## **Pipelining: The Instruction pipeline**



 Besides arithmetic & functional unit, instruction execution itself is pipelined also, e.g.: one instruction performs at least 3 steps:



Hardware Pipelining on processor (all units can run concurrently):

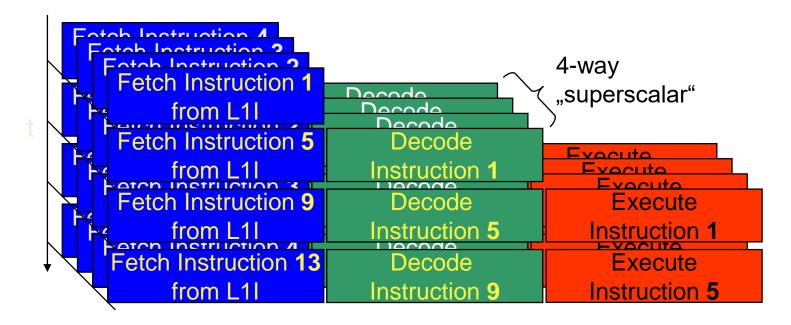


- Branches can stall this pipeline! (Speculative Execution, Predication)
- Each unit is pipelined itself (e.g., Execute = Multiply Pipeline)

#### **Superscalar Processors – Instruction Level Parallelism**



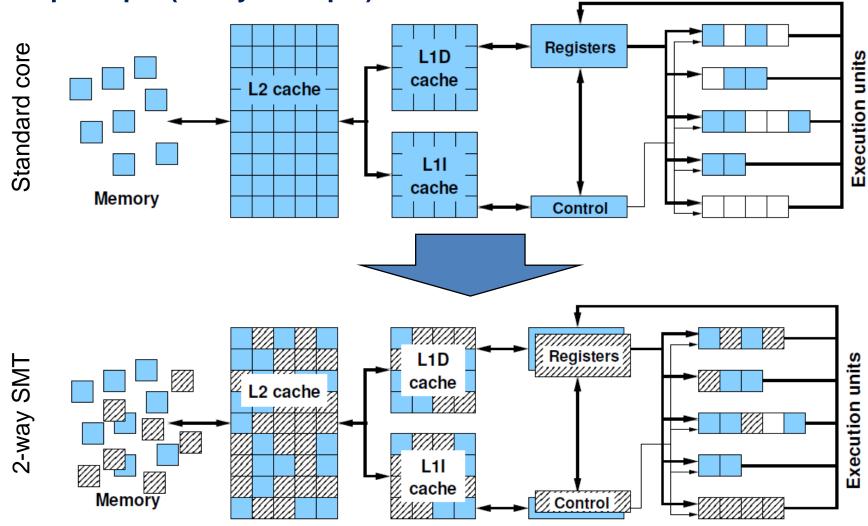
Multiple units enable use of Instrucion Level Parallelism (ILP): Instruction stream is "parallelized" on the fly



- Issuing m concurrent instructions per cycle: m-way superscalar
- Modern processors are 3- to 6-way superscalar & can perform 2 or 4 floating point operations per cycles

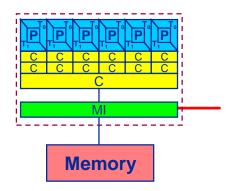






# **SMT** impact

- SMT adds another layer of topology (inside the physical core)
- Caveat: SMT threads share all caches!
- Possible benefit: Better pipeline throughput
  - Filling otherwise unused pipelines
  - Filling pipeline bubbles with other thread's executing instructions:



#### Thread 0:

Dependency → pipeline stalls until previous MULT is over

Thread 1:

Unrelated work in other thread can fill the pipeline bubbles

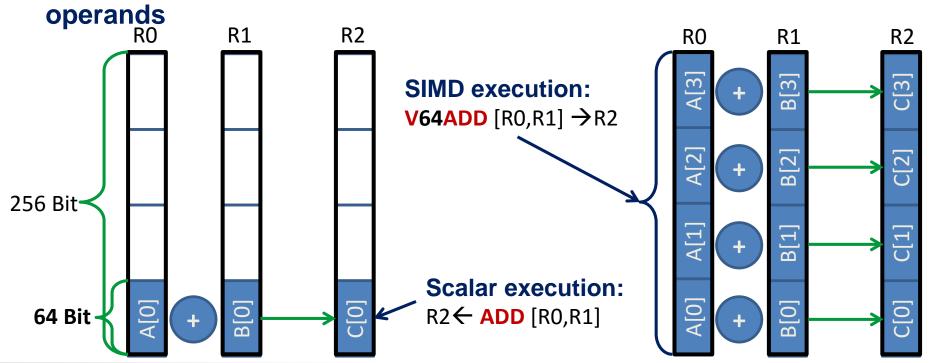
Beware: Executing it all in a single thread (if possible) may achieve the same goal without SMT:

#### Core details: SIMD processing



- Single Instruction Multiple Data (SIMD) operations allow the concurrent execution of the same operation on "wide" registers
- x86 SIMD instruction sets:
  - SSE: register width = 128 Bit → 2 double precision floating point operands
  - AVX: register width = 256 Bit → 4 double precision floating point operands

Adding two registers holding double precision floating point



#### SIMD processing – Basics



Steps (done by the compiler) for "SIMD processing"

```
for(int i=0; i<n;i++)
C[i]=A[i]+B[i];
```

"Loop unrolling"

```
for(int i=0; i<n;i+=4) {
        C[i] =A[i] +B[i];
        C[i+1]=A[i+1]+B[i+1];
        C[i+2]=A[i+2]+B[i+2];
        C[i+3]=A[i+3]+B[i+3];}
//remainder loop handling</pre>
```

```
Load 256 Bits starting from address of A[i] to register R0

Add the corresponding 64 Bit entries in R0 and R1 and store the 4 results to R2

Store R2 (256 Bit) to address starting at C[i]
```

```
LABEL1:

VLOAD R0 ← A[i]

VLOAD R1 ← B[i]

V64ADD[R0,R1] → R2

VSTORE R2 → C[i]

i←i+4

i<(n-4)? JMP LABEL1

//remainder loop handling
```

#### SIMD processing – Basics



No SIMD vectorization for loops with data dependencies:

```
for(int i=0; i<n;i++)
A[i]=A[i-1]*s;
```

"Pointer aliasing" may prevent SIMDfication

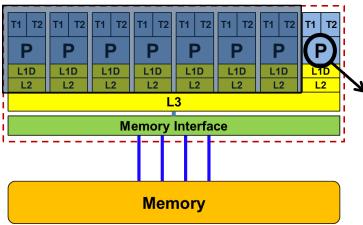
```
void f(double *A, double *B, double *C, int n) {
    for(int i=0; i<n; ++i)
        C[i] = A[i] + B[i];
}</pre>
```

- C/C++ allows that  $A \rightarrow \&C[-1]$  and  $B \rightarrow \&C[-2]$  $\rightarrow C[i] = C[i-1] + C[i-2]$ : dependency  $\rightarrow$  No SIMD
- If "pointer aliasing" does not happen, tell it to the compiler:

```
-fno-alias (Intel), -Msafeptr (PGI), -fargument-noalias (gcc)
-restrict keyword (C only!):
void f(double restrict *A, double restrict *B, double restrict *C, int n) {...}
```

#### There is no single driving force for single core performance!





Maximum floating point (FP) performance:

$$P_{core} = n_{super}^{FP} \cdot n_{FMA} \cdot n_{SIMD} \cdot f$$

| Super-    |  |  |  |  |
|-----------|--|--|--|--|
| scalarity |  |  |  |  |

FMA factor SIMD factor

Clock Speed

| Typical representatives | $n_{super}^{FP}$ inst./cy | $n_{FMA}$ | $n_{SIMD}$ ops/inst. |         | Code       | f [GHz] | P <sub>core</sub><br>[GF/s] |
|-------------------------|---------------------------|-----------|----------------------|---------|------------|---------|-----------------------------|
| Nehalem                 | 2                         | 1         | 2                    | Q1/2009 | X5570      | 2.93    | 11.7                        |
| Westmere                | 2                         | 1         | 2                    | Q1/2010 | X5650      | 2.66    | 10.6                        |
| Sandy Bridge            | 2                         | 1         | 4                    | Q1/2012 | E5-2680    | 2.7     | 21.6                        |
| Ivy Bridge              | 2                         | 1         | 4                    | Q3/2013 | E5-2660 v2 | 2.2     | 17.6                        |
| Haswell                 | 2                         | 2         | 4                    | Q3/2014 | E5-2695 v3 | 2.3     | 36.8                        |
| Broadwell               | 2                         | 2         | 4                    | Q1/2016 | E5-2699 v4 | 2.2     | 35.2                        |
| IBM POWER8              | 2                         | 2         | 2                    | Q2/2014 | S822LC     | 2.93    | 23.4                        |

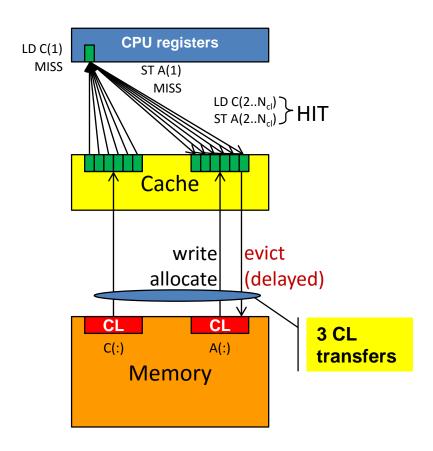
#### Registers and caches: Data transfers in a memory hierarchy



#### How does data travel from memory to the CPU and back?

- Remember: Caches are organized in cache lines (e.g., 64 bytes)
- Only complete cache lines are transferred between memory hierarchy levels (except registers)
- MISS: Load or store instruction does not find the data in a cache level
   → CL transfer required

Example: Array copy A(:)=C(:)

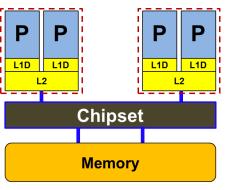


#### Commodity cluster nodes: From UMA to ccNUMA

Basic architecture of commodity compute cluster nodes



### Yesterday (2006): UMA

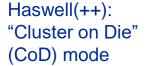


Uniform Memory Architecture (UMA)

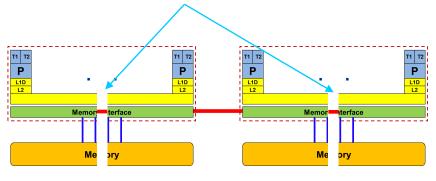
Flat memory; symmetric MPs

But: system "anisotropy"

2-socket server



#### Today: ccNUMA



2-socket server

Cache-coherent Non-Uniform Memory Architecture (ccNUMA)

Two or more **NUMA** domains per node

ccNUMA provides scalable bandwidth but: Where does my data finally end up?



# Interlude: A glance at current accelerator technology

NVidia "Pascal" GP100

VS.

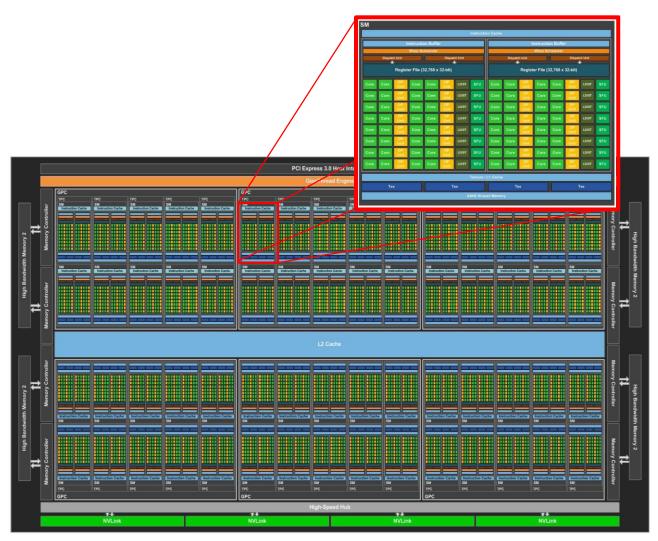
Intel Xeon Phi "Knights Landing"

# NVidia Pascal GP100 block diagram



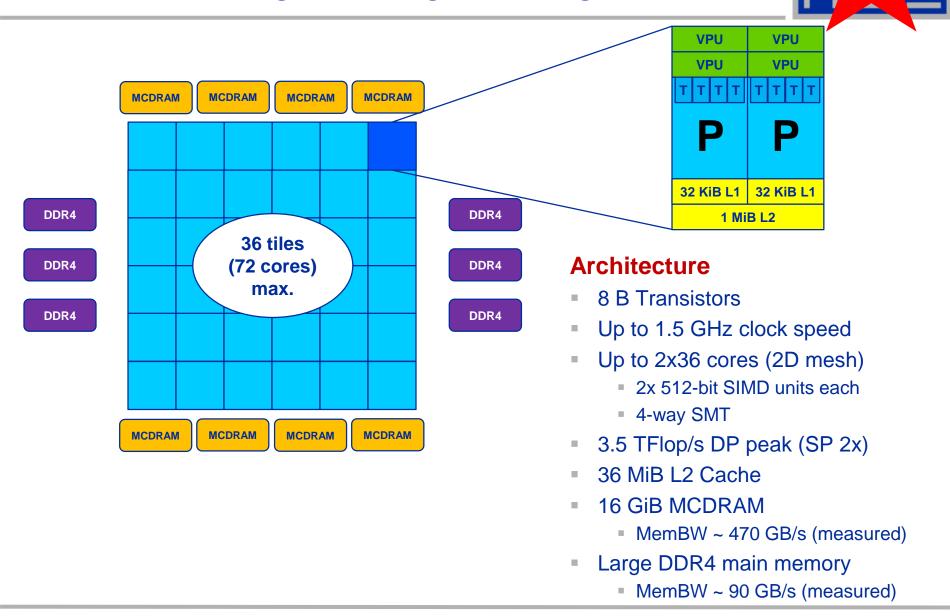
#### **Architecture**

- 15.3 B Transistors
- ~ 1.4 GHz clock speed
- Up to 60 "SM" units
  - 64 (SP) "cores" each
- 5.7 TFlop/s DP peak
- 4 MB L2 Cache
- 4096-bit HBM2
- MemBW ~ 732 GB/s (theoretical)
- MemBW ~ 510 GB/s (measured)
- 2:1 SP:DP performance



© NVIDIA Corp.

# Intel Xeon Phi "Knights Landing" block diagram



updated

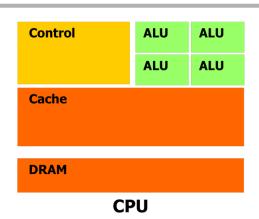
# Trading single thread performance for parallelism:

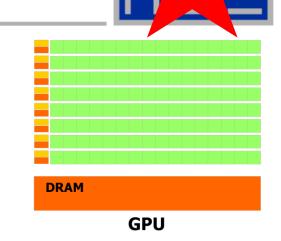
GPGPUs vs. CPUs

GPU vs. CPU light speed estimate (per device)

MemBW  $\sim 5-10x$ 

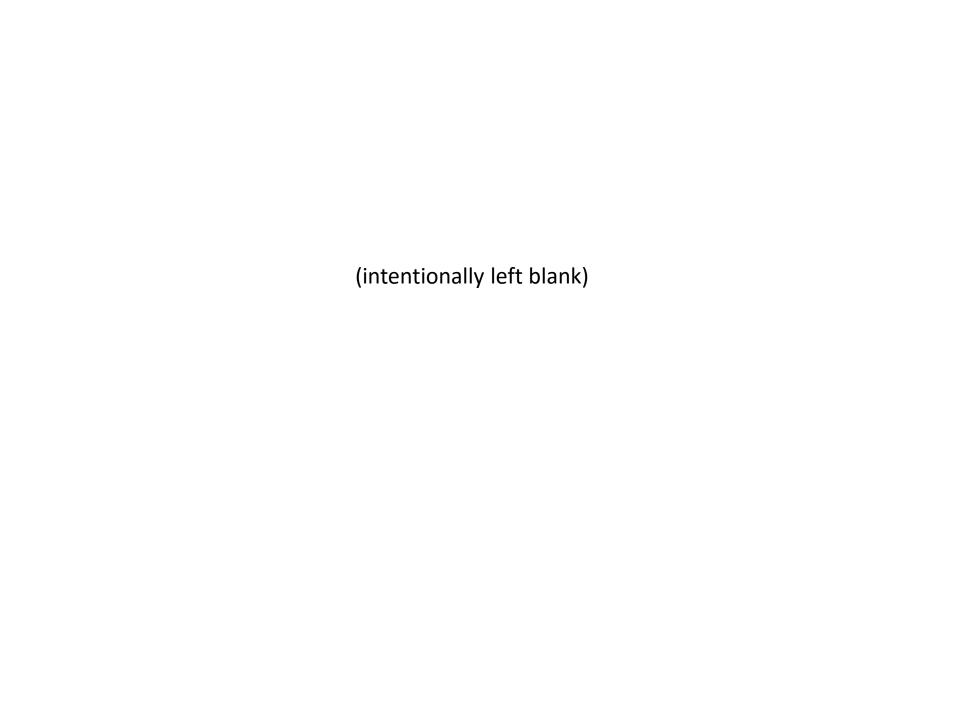
Peak ~ 6-15x

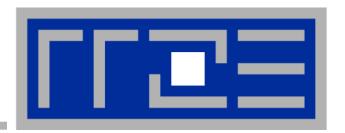




updated

|                     | 2x Intel Xeon E5-<br>2697v4 "Broadwell" | Intel Xeon Phi 7250<br>"Knights Landing" | NVidia Tesla P100<br>"Pascal" |
|---------------------|---|--|-------------------------------|
| Cores@Clock         | 2 x 18 @ ≥2.3 GHz                       | 68 @ 1.4 GHz                             | 56 SMs @ ~1.3 GHz             |
| SP Performance/core | ≥73.6 GFlop/s                           | 89.6 GFlop/s                             | ~166 GFlop/s                  |
| Threads@STREAM      | ~8                                      | ~40                                      | >8000?                        |
| SP peak             | ≥2.6 TFlop/s                            | 6.1 TFlop/s                              | ~9.3 TFlop/s                  |
| Stream BW (meas.)   | 2 x 62.5 GB/s                           | 450 GB/s (HBM)                           | 510 GB/s                      |
| Transistors / TDP   | ~2x7 Billion / 2x145 W                  | 8 Billion / 215W                         | 14 Billion/300W               |



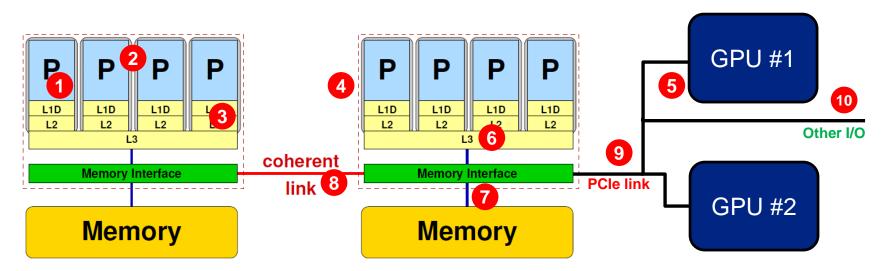


# Node topology and programming models

# Parallelism in a modern compute node



Parallel and shared resources within a shared-memory node



#### Parallel resources:

- Execution/SIMD units 1
- Cores 2
- Inner cache levels 3
- Sockets / ccNUMA domains 4
- Multiple accelerators 5

#### **Shared resources:**

- Outer cache level per socket 6
- Memory bus per socket 7
- Intersocket link 8
- PCle bus(es)
- Other I/O resources 10

How does your application react to all of those details?

#### **Parallel programming models**

#### on modern compute nodes



#### Shared-memory (intra-node)

- Good old MPI
- OpenMP
- POSIX threads
- Intel Threading Building Blocks (TBB)
- Cilk+, OpenCL, StarSs,... you name it

#### "Accelerated"

- OpenMP 4.0+
- CUDA
- OpenCL
- OpenACC

#### Distributed-memory (inter-node)

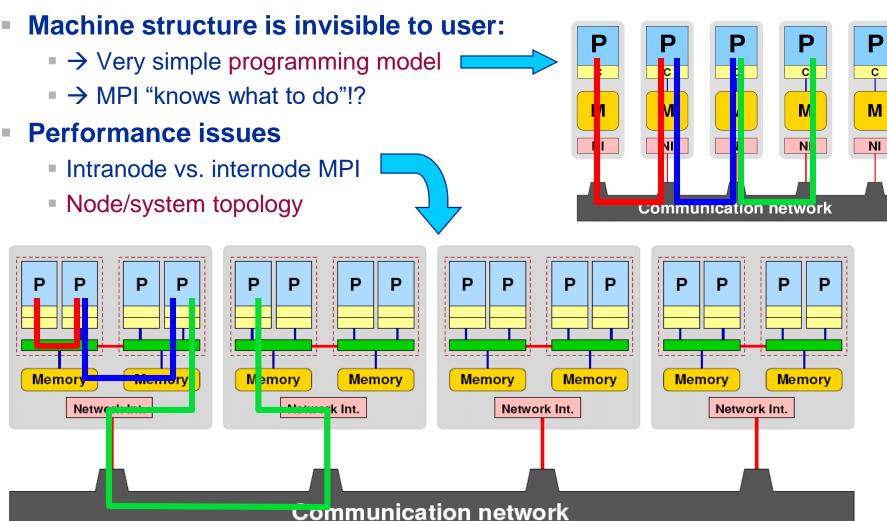
- MPI
- PGAS (CAF, UPC, ...)
- Hybrid
  - Pure MPI + X, X == <you name it>

All models require awareness of topology and affinity issues for getting best performance out of the machine!

#### **Parallel programming models:**

#### Pure MPI





### **Parallel programming models:**

Pure threading on the node

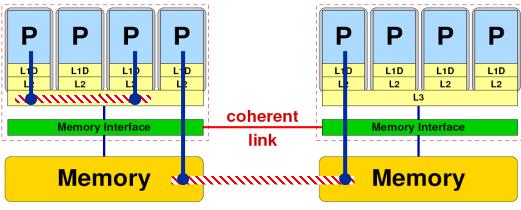


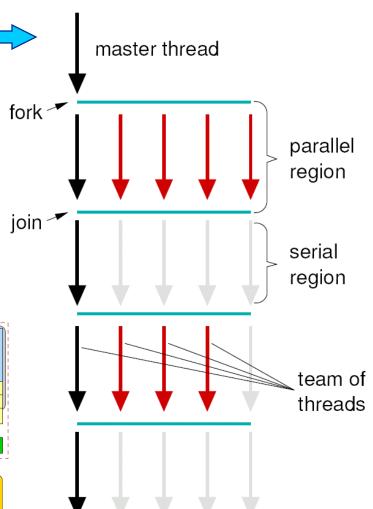
#### Machine structure is invisible to user

- → Very simple programming model
- Threading SW (OpenMP, pthreads, TBB,...) should know about the details

#### Performance issues

- Synchronization overhead
- Memory access
- Node topology





### Parallel programming models: Lots of choices

Hybrid MPI+OpenMP on a multicore multisocket cluster

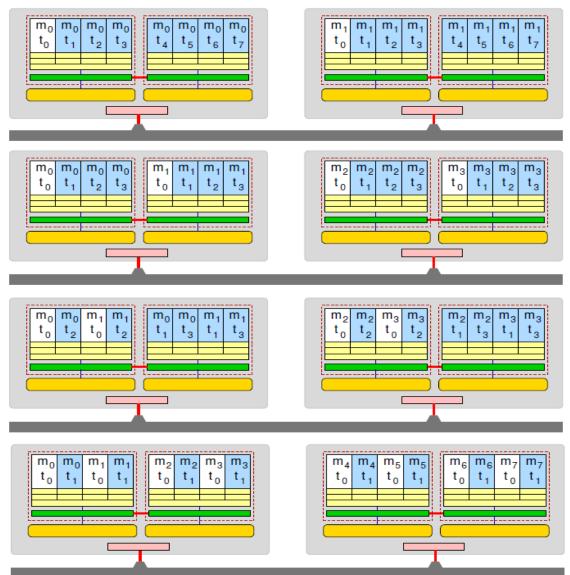


# One MPI process / node

One MPI process / socket: OpenMP threads on same socket: "blockwise"

> OpenMP threads pinned "round robin" across cores in node

Two MPI processes / socket
OpenMP threads
on same socket



#### Conclusions about architecture



- Modern computer architecture has a rich "topology"
- Node-level hardware parallelism takes many forms
  - Sockets/devices CPU: 1-8, GPGPU: 1-6
  - Cores moderate (CPU: 4-16) to massive (GPGPU: 1000's)
  - SIMD moderate (CPU: 2-8) to massive (GPGPU: 10's-100's)
  - Superscalarity (CPU: 2-6)
- Exploiting performance: parallelism + bottleneck awareness
  - "High Performance Computing" == computing at a bottleneck
- Performance of programming models is sensitive to architecture
  - Topology/affinity influences overheads
  - Standards do not contain (many) topology-aware features
  - Apart from overheads, performance features are largely independent of the programming model



# **Multicore Performance and Tools**

# **Tools for Node-level Performance Engineering**



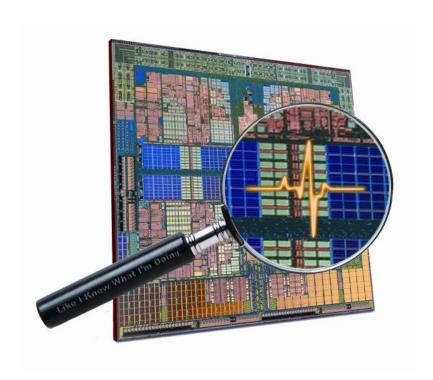
- Gather Node Information
   hwloc, likwid-topology, likwid-powermeter
- Affinity control and data placement
   OpenMP and MPI runtime environments, hwloc, numactl, likwid-pin
- Runtime Profiling
   Compilers, gprof, HPC Toolkit, ...
- Performance Profilers
   Intel Vtune<sup>TM</sup>, likwid-perfctr, PAPI based tools, Linux perf, ...
- Microbenchmarking
   STREAM, likwid-bench, Imbench

### **LIKWID** performance tools



#### **LIKWID** tool suite:

Like
I
Knew
What
I'm
Doing



Open source tool collection (developed at RRZE): https://github.com/RRZE-HPC/likwid

J. Treibig, G. Hager, G. Wellein: *LIKWID: A*lightweight performance-oriented tool suite for x86

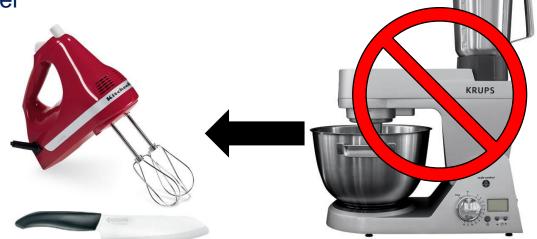
multicore environments. PSTI2010, Sep 13-16, 2010,
San Diego, CA http://arxiv.org/abs/1004.4431

#### **Likwid Tool Suite**



#### Command line tools for Linux:

- easy to install
- works with standard linux kernel
- simple and clear to use
- supports Intel and AMD CPUs



#### Current tools:

- likwid-topology: Print thread and cache topology
- likwid-powermeter: Measure energy consumption
- likwid-pin: Pin threaded application without touching code
- likwid-perfctr: Measure performance counters
- likwid-bench: Microbenchmarking tool and environment
- ... some more

# Output of likwid-topology -g

#### on one node of Intel Haswell-EP



```
CPU name:
              Intel(R) Xeon(R) CPU E5-2695 v3 @ 2.30GHz
              Intel Xeon Haswell EN/EP/EX processor
CPU type:
CPU stepping: 2
Hardware Thread Topology
Sockets:
Cores per socket:
Threads per core:
                                                         Available
HWThread
              Thread
                            Core
                                           Socket
43
44
              ( 0 28 1 29 2 30 3 31 4 32 5 33 6 34 7 35 8 36 9 37 10 38 11 39 12 40 13 41 )
Socket 0:
Socket 1:
              ( 14 42 15 43 16 44 17 45 18 46 19 47 20 48 21 49 22 50 23 51 24 52 25 53 26 54 27 55 )
                                                                                                               All physical
                                                                                                             processor IDs
Cache Topology
Size:
                            32 kB
Cache groups: (0 28) (1 29) (2 30) (3 31) (4 32) (5 33) (6 34) (7 35) (8 36) (9 37) (10 38) (11 39) (12 40) (13 41
) ( 14 42 ) ( 15 43 ) ( 16 44 ) ( 17 45 ) ( 18 46 ) ( 19 47 ) ( 20 48 ) ( 21 49 ) ( 22 50 ) ( 23 51 ) ( 24 52 ) ( 25 53 ) ( 26 54 ) ( 27 55 )
Level:
Size:
                            256 kB
Cache groups: (0 28) (1 29) (2 30) (3 31) (4 32) (5 33) (6 34) (7 35) (8 36) (9 37) (10 38) (11 39) (12 40) (13 41
) ( 14 42 ) ( 15 43 ) ( 16 44 ) ( 17 45 ) ( 18 46 ) ( 19 47 ) ( 20 48 ) ( 21 49 ) ( 22 50 ) ( 23 51 ) ( 24 52 ) ( 25 53 ) ( 26 54 ) ( 27 55 )
Level:
Size:
Cache groups: ( 0 28 1 29 2 30 3 31 4 32 5 33 6 34 ) ( 7 35 8 36 9 37 10 38 11 39 12 40 13 41 ) ( 14 42 15 43 16 44 17 45 18 46 19 47 20 48 )
( 21 49 22 50 23 51 24 52 25 53 26 54 27 55 )
```

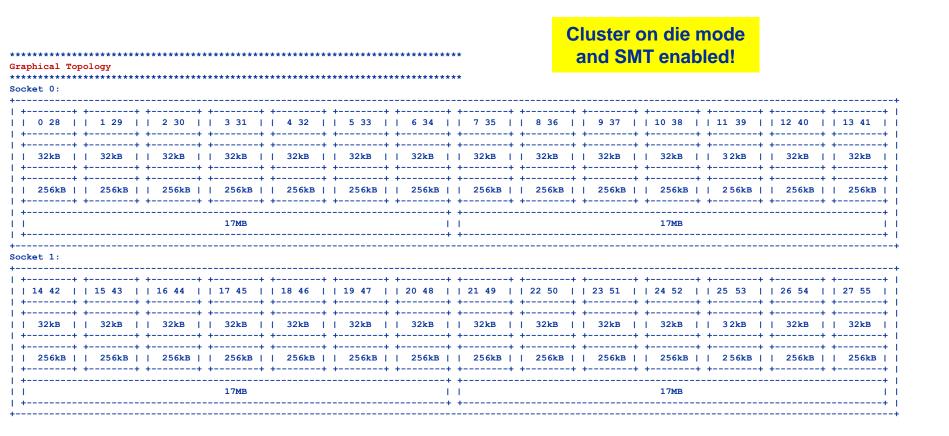
# **Output of likwid-topology continued**



```
NUMA Topology
NUMA domains:
Domain:
Processors:
                ( 0 28 1 29 2 30 3 31 4 32 5 33 6 34 )
Distances:
                                   10 21 31 31
                                  13292.9 MB
Free memory:
                                  15941.7 MB
Total memory:
Domain:
Processors:
                ( 7 35 8 36 9 37 10 38 11 39 12 40 13 41 )
Distances:
                                   21 10 31 31
                                  13514 MB
Free memory:
                                  16126.4 MB
Total memory:
Domain:
Processors:
                 ( 14 42 15 43 16 44 17 45 18 46 19 47 20 48 )
Distances:
                                   31 31 10 21
                                   15025.6 MB
Free memory:
Total memory:
                                   16126.4 MB
Domain:
                ( 21 49 22 50 23 51 24 52 25 53 26 54 27 55 )
Processors:
                                   31 31 21 10
Distances:
Free memory:
                                  15488.9 MB
                                  16126 MB
Total memory:
```

# **Output of likwid-topology continued**







# **Enforcing thread/process-core affinity under the Linux OS**

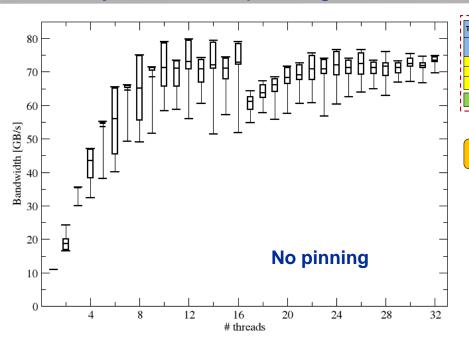
Standard tools and OS affinity facilities under program control

likwid-pin

# **Example: STREAM benchmark on 16-core Sandy Bridge:**

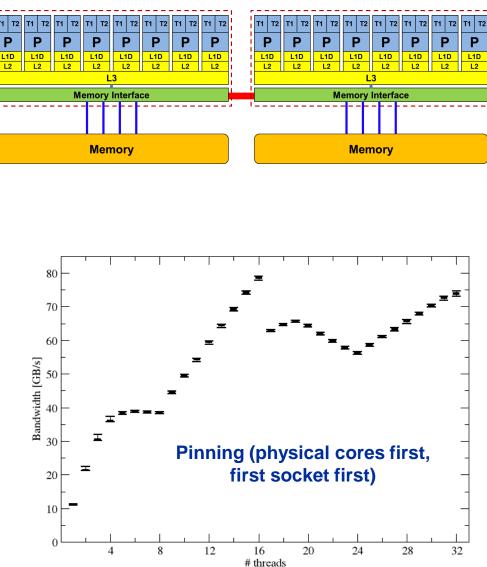
Anarchy vs. thread pinning







- Eliminating performance variation
- Making use of architectural features
- Avoiding resource contention



### More thread/Process-core affinity ("pinning") options



- Highly OS-dependent system calls
  - But available on all systems

Linux: sched\_setaffinity()
Windows: SetThreadAffinityMask()

- Hwloc project (http://www.open-mpi.de/projects/hwloc/)
- Support for "semi-automatic" pinning in some compilers/environments
  - All modern compilers with OpenMP support
  - Generic Linux: taskset, numactl, likwid-pin (see below)
  - OpenMP 4.0 (see OpenMP tutorial)
- Affinity awareness in MPI libraries
  - SGI MPT
  - OpenMPI
  - Intel MPI
  - · ...

#### Likwid-pin

#### Overview



- Pins processes and threads to specific cores without touching code
- Directly supports pthreads, gcc OpenMP, Intel OpenMP
- Based on combination of wrapper tool together with overloaded pthread library → binary must be dynamically linked!
- Can be used as a superior replacement for taskset
- Supports logical core numbering within a node

#### Usage examples:

```
likwid-pin -c 0-3,4,6 ./myApp parameters
```

- likwid-pin -c S0:0-7 ./myApp parameters
- likwid-pin -c N:0-15 ./myApp parameters
- OMP\_NUM\_THREADS is set by the tool if not set explicitly

### **LIKWID** terminology

#### Thread group syntax

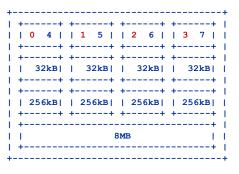


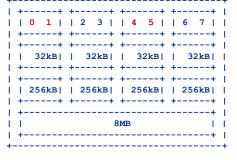
- The OS numbers all processors (hardware threads) on a node
- The numbering is enforced at boot time by the BIOS and may have nothing to do with topological entities
- LIKWID concept: thread group consisting of HW threads sharing a topological entity (e.g., socket, shared cache,...)
- A thread group is defined by a single character + index
- Example: likwid-pin -c S1:0-3,6,7 ./a.out
- Group expression chaining with @: likwid-pin -c S0:0-3@S1:0-3 ./a.out
- Alternative expression based syntax: likwid-pin -c E:S0:4:2:4 ./a.out

E:<thread domain>:<num threads>:<chunk size>:<stride>

Expression syntax is convenient for Xeon Phi: likwid-pin -c E:N:120:2:4 ./a.out numbering across physical cores first within the group

compact numbering within the group

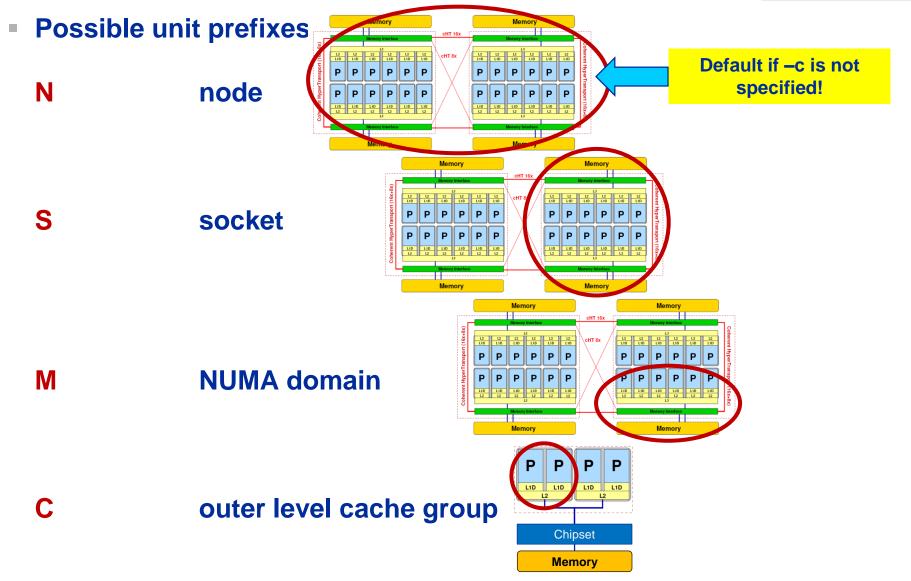




#### Likwid

#### Currently available thread domains





#### Likwid-pin

Example: Intel OpenMP



#### Running the STREAM benchmark with likwid-pin:

```
$ likwid-pin -c S0:0-3 ./stream
Double precision appears to have 16 digits of accuracy
Assuming 8 bytes per DOUBLE PRECISION word
Array size = 20000000
Offset
                      32
                                                            Main PID always
 The total memory requirement is 457 MB
                                                                 pinned
 You are running each test 10 times
 The *best* time for each test is used
                                                             Skip shepherd
 *EXCLUDING* the first and last iterations
                                                            thread if necessary
[pthread wrapper] MAIN -> 0
[pthread wrapper] PIN MASK: 0->1 1->2 2->3
[pthread wrapper] SKIP MASK: 0x1
        threadid 140668624234240 -> SKIP
        threadid 140668598843264 -> core 1 - OK
        threadid 140668594644992 -> core 2 - OK
        threadid 140668590446720 -> core 3 - OK
                                                             Pin all spawned
                                                              threads in turn
  [... rest of STREAM output omitted ...]
```

# Intel KMP\_AFFINITY environment variable



• KMP\_AFFINITY=[<modifier>,...]<type>[,<permute>][,<offset>]

- modifier
  - granularity=<specifier> takes the following specifiers: fine, thread, and core
  - norespect
  - noverbose
  - proclist={proc-list>}
  - respect
  - verbose

- type (required)
  - compact
  - disabled
  - explicit (GOMP\_CPU\_AFFINITY)
  - none
  - scatter

**OS processor IDs** 

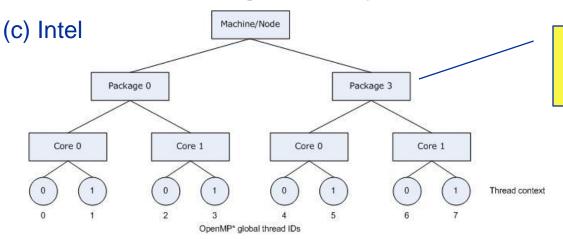
Respect an OS affinity mask in place

- Default:
  - noverbose,respect,granularity=core
- KMP\_AFFINITY=verbose, none to list machine topology map

### Intel KMP AFFINITY examples

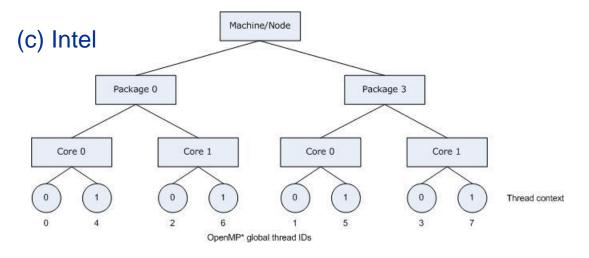


# KMP\_AFFINITY=granularity=fine,compact



Package means chip/socket

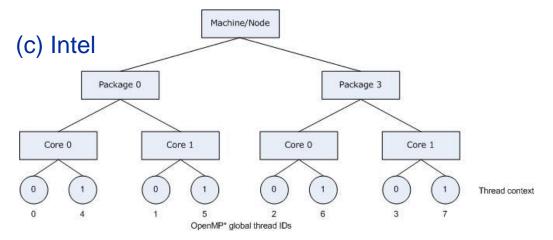
# KMP\_AFFINITY=granularity=fine,scatter



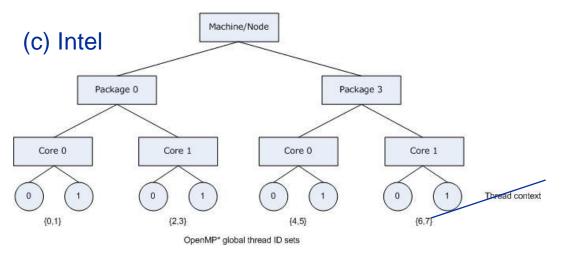
### Intel KMP AFFINITY permute example



KMP\_AFFINITY=granularity=fine,compact,1,0



KMP\_AFFINITY=granularity=core,compact

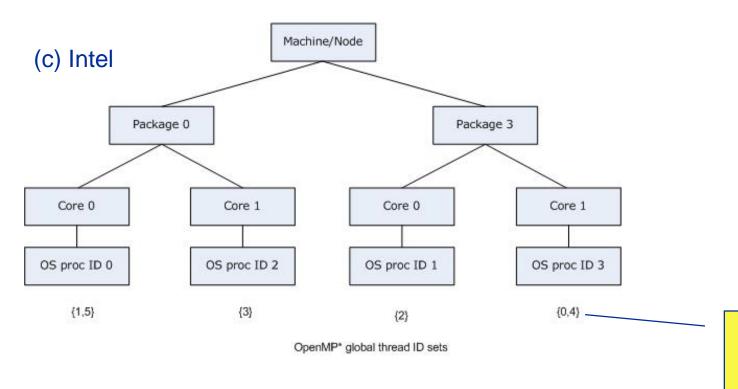


Threads may float within core

### **GNU GOMP AFFINITY**



■ GOMP\_AFFINITY=3,0-2 used with 6 threads



Round robin oversubscription

Always operates with OS processor IDs



# Multicore performance tools: Probing performance behavior

likwid-perfctr

#### Basic approach to performance analysis



- Runtime profile / Call graph (gprof): Where are the hot spots?
- 2. Instrument hot spots (prepare for detailed measurement)
- 3. Find performance signatures

#### **Possible signatures:**

- Bandwidth saturation
- Instruction throughput limitation (real or language-induced)
- Latency impact (irregular data access, high branch ratio)
- Load imbalance
- ccNUMA issues (data access across ccNUMA domains)
- Pathologic cases (false cacheline sharing, expensive operations)

likwid-perfctr can help here

**Goal**: Come up with educated guess about a performance-limiting motif (Performance Pattern)

# **Probing performance behavior**



- How do we find out about the performance properties and requirements of a parallel code?
  - Profiling via advanced tools is often overkill
- A coarse overview is often sufficient
  - likwid-perfctr (similar to "perfex" on IRIX, "hpmcount" on AIX, "lipfpm" on Linux/Altix)
  - Simple end-to-end measurement of hardware performance metrics
  - "Marker" API for starting/stopping counters
  - Multiple measurement region support
  - Preconfigured and extensible metric groups, list with

```
likwid-perfctr -a
```

BRANCH: Branch prediction miss rate/ratio

CACHE: Data cache miss rate/ratio

CLOCK: Clock of cores

DATA: Load to store ratio

FLOPS\_DP: Double Precision MFlops/s FLOPS\_SP: Single Precision MFlops/s

FLOPS\_X87: X87 MFlops/s

L2: L2 cache bandwidth in MBytes/s
L2CACHE: L2 cache miss rate/ratio
L3: L3 cache bandwidth in MBytes/s
L3CACHE: L3 cache miss rate/ratio

MEM: Main memory bandwidth in MBytes/s

TLB: TLB miss rate/ratio

#### Example usage with preconfigured metric group (shortened)



```
$ likwid-perfctr -C N:0-3 -g FLOPS DP ./stream.exe
                Intel(R) Xeon(R) CPU E5-2660 v2 @ 2.20GHz
CPU name:
CPU type:
                Intel Xeon IvyBridge EN/EP/EX processor
CPU clock:
                2.20 GHz
                                                                         Configured metrics
[... YOUR PROGRAM OUTPUT ...]
                                                  Always
                                                                            (this group)
                                                 measured
Group 1: FLOPS DP
                                         Counter
                                                       Core 0
                                                                     Core 1
                                                                                  Core 2
                                                                                                Core
                  Event
                                            EXXC0
            INSTR RETIRED ANY
                                                      521332883
                                                                    523904122
                                                                                 519696583
                                                                                               519193
          CPU CLK UNHALTED CORE
                                                                                              1376447
                                            FIXC1
                                                     1379625927
                                                                   1381900036 I
                                                                                1378355460
          CPU CLK UNHALTED REF
                                            FIXC2
                                                     1389460886
                                                                   1393031508 I
                                                                                1387504228
                                                                                              1385276
  FP COMP OPS EXE SSE FP PACKED DOUBLE
                                             PMC0
                                                      176216849
                                                                    176176025
                                                                                 177432054
                                                                                               176367
  FP COMP OPS EXE SSE FP SCALAR DOUBLE
                                                                       599
                                                        1554
                                                                                     72
                                                                                                  27
                                             PMC1
        SIMD FP 256 PACKED DOUBLE
                                             PMC2
                                                          0
                                                                        0
                                                                                     O
                            Core 0
         Metric
                                        Core 1
                                                     Core 2
                                                                   Core 3
                                                     0.6856
                                                                   0.6856
                            0.6856
                                        0.6856
   Runtime (RDTSC) [s] |
  Runtime unhalted [s] |
                            0.6270
                                        0.6281
                                                     0.6265
                                                                  0.6256
                                                                                            Derived
                                                                 2186.2243
       Clock [MHz]
                          2184.6742 | 2182.6664 |
                                                   2185.7404
                                                                                            metrics
                                                                   2.6511
           CPI
                            2.6463
                                        2.6377
                                                     2.6522
         MFLOP/s
                           514.0890 I
                                       513.9685 I
                                                    517.6320 I
                                                                  514.5273
       AVX MFLOP/s
                                           0
                                                       0
                                                                      0
     Packed MUOPS/s
                                                    258.8160
                           257.0434
                                       256.9838 I
                                                                  257,2636
     Scalar MUOPS/s
                            0.0023
                                        0.0009
                                                     0.0001
                                                                3.938426e-05
```

#### Marker API (C/C++ and Fortran)



- A marker API is available to restrict measurements to code regions
- The API only turns counters on/off. The configuration of the counters is still done by likwid-perfctr
- Multiple named region support, accumulation over multiple calls
- Inclusive and overlapping regions allowed

```
#include <likwid.h>
                                    // must be called from serial region
LIKWID MARKER INIT;
#pragma omp parallel
                                    // only reqd. if measuring multiple threads
  LIKWID MARKER THREADINIT;
LIKWID MARKER START ("Compute");
                                             Activate macros with -DLIKWID PERFMON
                                             Run likwid-perfctr with -m option to
LIKWID MARKER STOP("Compute");
                                             activate markers
LIKWID MARKER START ("Postprocess");
LIKWID MARKER STOP("Postprocess");
                                    // must be called from serial region
LIKWID MARKER CLOSE;
```

#### Best practices for runtime counter analysis



#### Things to look at (in roughly this order)

- Excess work
- Load balance (flops, instructions, BW)
- In-socket memory BW saturation
- Flop/s, loads and stores per flop metrics
- SIMD vectorization
- CPI metric
- # of instructions,
   branches, mispredicted branches

#### **Caveats**

- Load imbalance may not show in CPI or # of instructions
  - Spin loops in OpenMP barriers/MPI blocking calls
  - Looking at "top" or the Windows Task Manager does not tell you anything useful
- In-socket performance saturation may have various reasons
- Cache miss metrics are sometimes misleading



# **Measuring energy consumption** with LIKWID

# Measuring energy consumption

likwid-powermeter and likwid-perfctr-g ENERGY



- Implements Intel RAPL interface (Sandy Bridge)
- RAPL = "Running average power limit"

CPU name: Intel Core SandyBridge processor

CPU clock: 3.49 GHz

\_\_\_\_\_

Base clock: 3500.00 MHz Minimal clock: 1600.00 MHz

Turbo Boost Steps:

C1 3900.00 MHz

C2 3800.00 MHz

C3 3700.00 MHz

C4 3600.00 MHz

\_\_\_\_\_

Thermal Spec Power: 95 Watts

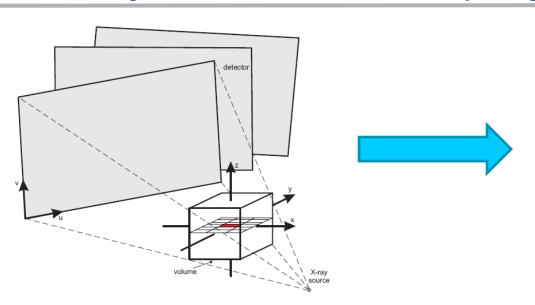
Minimum Power: 20 Watts
Maximum Power: 95 Watts

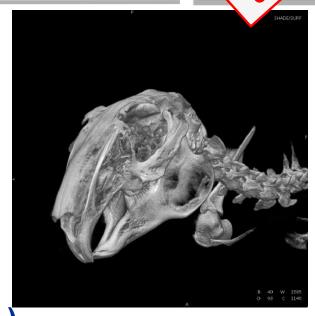
Maximum Time Window: 0.15625 micro sec

\_\_\_\_\_

# **Example:**

A medical image reconstruction code on Sandy Bridge





Sandy Bridge EP (8 cores, 2.7 GHz base freq.)

| Test case          | Runtime [s] | Power [W] |                    | Energy [J] |
|--------------------|-------------|-----------|--------------------|------------|
| 8 cores, plain C   | 90.43       | 90        | Fas<br><b>→</b> le | 8110       |
| 8 cores, SSE       | 29.63       | 93        | ter o              | 2750       |
| 8 cores (SMT), SSE | 22.61       | 102       | code               | 2300       |
| 8 cores (SMT), AVX | 18.42       | 111       |                    | 2040       |

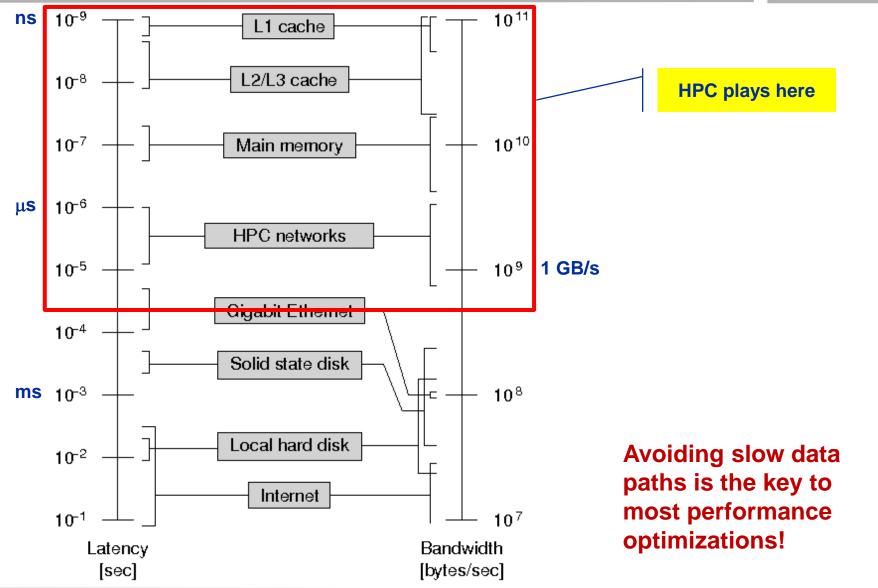


# Microbenchmarking for architectural exploration (and more)

Probing of the memory hierarchy
Saturation effects in cache and memory

# Latency and bandwidth in modern computer environments

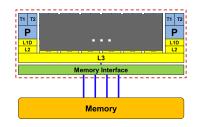




# **Intel Xeon E5 multicore processors**



| Microarchitecture               | SandyBridge-EP          | IvyBridge-EP            | Haswell-EP              |  |  |  |
|---------------------------------|-------------------------|-------------------------|-------------------------|--|--|--|
| Shorthand                       | SNB                     | IVB                     | HSW                     |  |  |  |
| Xeon Model                      | E5-2680                 | E5-2690 v2              | E5-2695 v3              |  |  |  |
| Year                            | 03/2012                 | 09/2013                 | 09/2014                 |  |  |  |
| Clock speed (fixed)             | 2.7 GHz                 | 2.2 GHz                 | 2.3 GHz                 |  |  |  |
| Cores/Threads                   | 8/16                    | 10/20                   | 14/28                   |  |  |  |
| Load/Store throughput per cycle |                         |                         |                         |  |  |  |
| AVX(2)                          | 1 LD & 1/2 ST           | 1 LD & 1/2 ST           | 2 LD & 1 ST             |  |  |  |
| SSE/scalar                      | 2 LD    1 LD & 1 ST     | 2 LD    1 LD & 1 ST     | 2 LD & 1 ST             |  |  |  |
| L1 port width                   | $2\times16+1\times16$ B | $2\times16+1\times16$ B | $2\times32+1\times32$ B |  |  |  |
| ADD throughput                  | 1 / cy                  | 1 / cy                  | 1 / cy                  |  |  |  |
| MUL throughput                  | 1 / cy                  | 1 / cy                  | 2 / cy                  |  |  |  |
| FMA throughput                  | n/a                     | n/a                     | 2 / cy                  |  |  |  |
| L2-L1 data bus                  | 32 B                    | 32 B                    | 64 B                    |  |  |  |
| L3-L2 data bus                  | 32 B                    | 32 B                    | 32 B                    |  |  |  |
| LLC size                        | $20\mathrm{MiB}$        | 25 MiB                  | 35 MiB                  |  |  |  |
| Main memory                     | 4×DDR3-1600             | 4×DDR3-1866             | 4×DDR4-2133             |  |  |  |
| Peak memory BW                  | 51.2 GB/s               | 51.2 GB/s               | 68.3 GB/s               |  |  |  |
| Load-only BW                    | 43.6 GB/s (85%)         | 46.1 GB/s (90%)         | 60.6 GB/s (89%)         |  |  |  |
| $T_{\rm L3Mem}$ per CL          | 3.96 cy                 | 3.05 cy                 | 2.43 cy                 |  |  |  |



FP instructions throughput per core

Max. data transfer per cycle between caches

Peak main memory bandwidth

#### The parallel vector triad benchmark

A "swiss army knife" for microbenchmarking



#### Simple streaming benchmark:

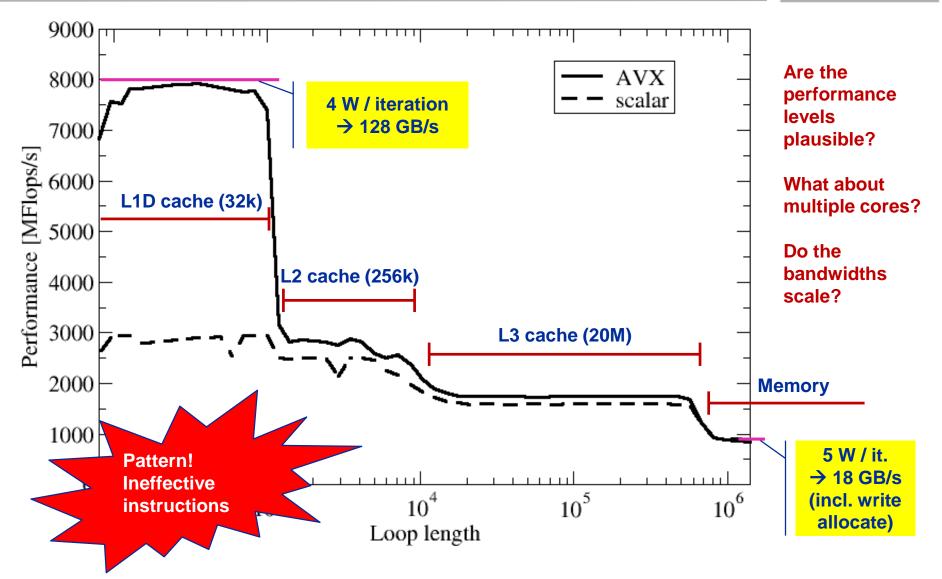
```
double precision, dimension(N) :: A,B,C,D
A=1.d0; B=A; C=A; D=A

do j=1,NITER
    do i=1,N
        A(i) = B(i) + C(i) * D(i)
    enddo
    if(.something.that.is.never.true.) then
        call dummy(A,B,C,D)
    endif
enddo
```

- Report performance for different N
- Choose NITER so that accurate time measurement is possible
- This kernel is limited by data transfer performance for all memory levels on all current architectures!

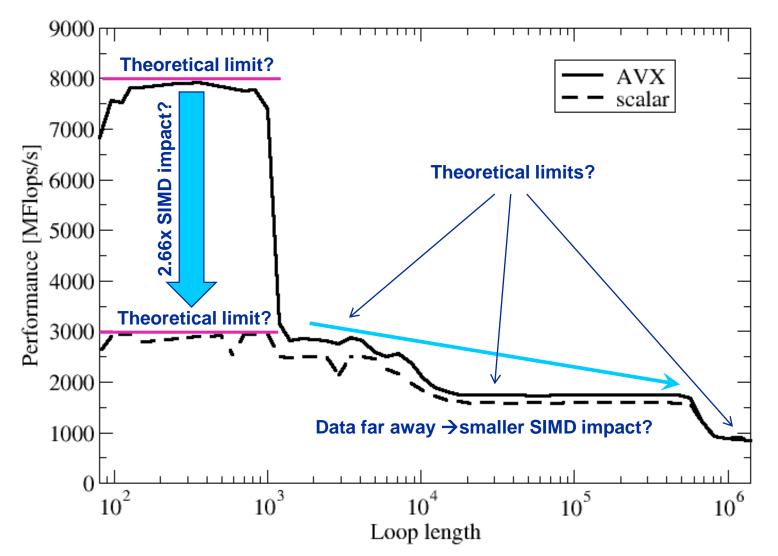
# A(:)=B(:)+C(:)\*D(:) on one Sandy Bridge core (3 GHz)





# A(:)=B(:)+C(:)\*D(:) on one Sandy Bridge core (3 GHz): Observations and further questions





See later for answers!

#### The throughput-parallel vector triad benchmark



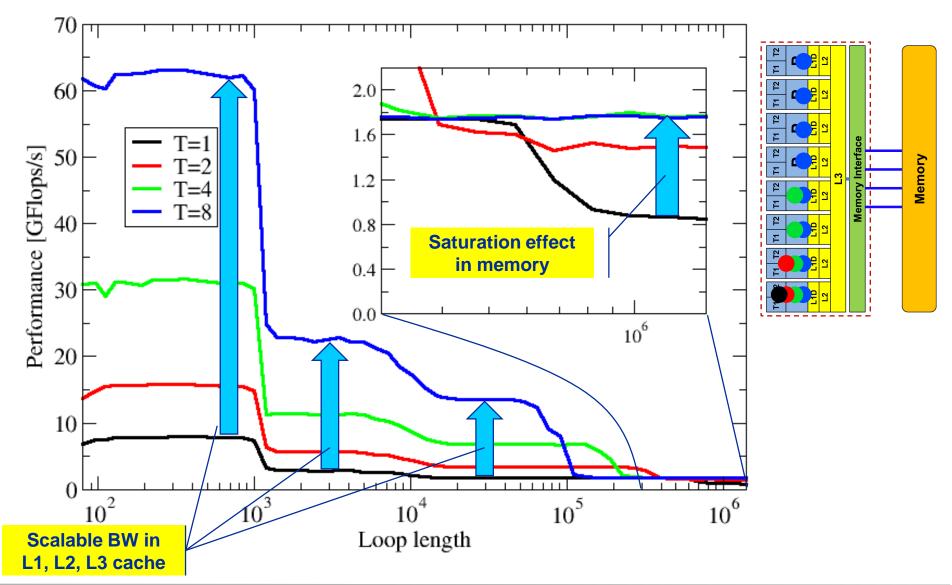
#### **Every core runs its own, independent triad benchmark**

```
double precision, dimension(:), allocatable :: A,B,C,D
!$OMP PARALLEL private(i,j,A,B,C,D)
allocate (A(1:N), B(1:N), C(1:N), D(1:N))
A=1.d0; B=A; C=A; D=A
do j=1,NITER
  do i=1,N
    A(i) = B(i) + C(i) * D(i)
  enddo
  if(.something.that.is.never.true.) then
    call dummy (A,B,C,D)
  endif
enddo
!$OMP END PARALLEL
```

### > pure hardware probing, no impact from OpenMP overhead

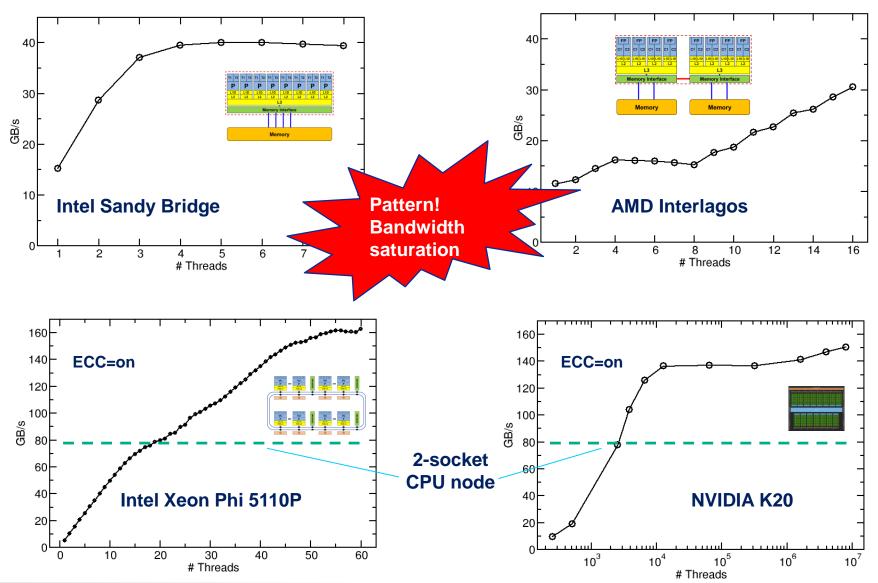
# Throughput vector triad on Sandy Bridge socket (3 GHz)





# Attainable memory bandwidth: Comparing architectures





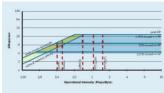
#### **Conclusions from the microbenchmarks**



- Affinity matters!
  - Almost all performance properties depend on the position of
    - Data
    - Threads/processes
  - Consequences
    - Know where your threads are running
    - Know where your data is

Bandwidth bottlenecks are ubiquitous







## "Simple" performance modeling: The Roofline Model

Loop-based performance modeling: Execution vs. data transfer

**Example: array summation** 

**Example: dense & sparse matrix-vector multiplication** 

Example: a 3D Jacobi solver Model-guided optimization

R.W. Hockney and I.J. Curington:  $f_{1/2}$ : A parameter to characterize memory and communication bottlenecks. Parallel Computing 10, 277-286 (1989). DOI: 10.1016/0167-8191(89)90100-2

W. Schönauer: <u>Scientific Supercomputing: Architecture and Use of Shared and Distributed Memory Parallel Computers</u>. Self-edition (2000)

S. Williams: <u>Auto-tuning Performance on Multicore Computers</u>. UCB Technical Report No. UCB/EECS-2008-164. PhD thesis (2008)

#### Prelude: Modeling customer dispatch in a bank



Revolving door throughput: b<sub>S</sub> [customers/sec]















Processing capability:

Pmax [tasks/sec]



Intensity:

/ [tasks/customer]





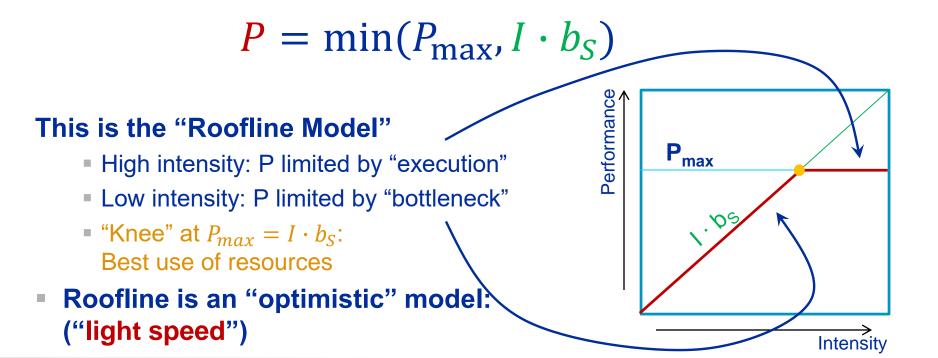
#### Prelude: Modeling customer dispatch in a bank



#### How fast can tasks be processed? P [tasks/sec]

#### The bottleneck is either

- The service desks (max. tasks/sec):  $P_{\text{max}}$
- The revolving door (max. customers/sec):  $I \cdot b_S$



#### The Roofline Model



[Byte/s]

- 1. P<sub>max</sub> = Applicable peak performance of a loop, assuming that data comes from the level 1 cache (this is not necessarily P<sub>peak</sub>)
   → e.g., P<sub>max</sub> = 176 GFlop/s
- 2. I = Computational intensity ("work" per byte transferred) over the slowest data path utilized (code balance  $B_C = I^{-1}$ )  $\rightarrow$  e.g., I = 0.167 Flop/Byte  $\rightarrow B_C = 6$  Byte/Flop
- 3.  $b_S$  = Applicable peak bandwidth of the slowest data path utilized  $\rightarrow$  e.g.,  $b_S$  = 56 GByte/s

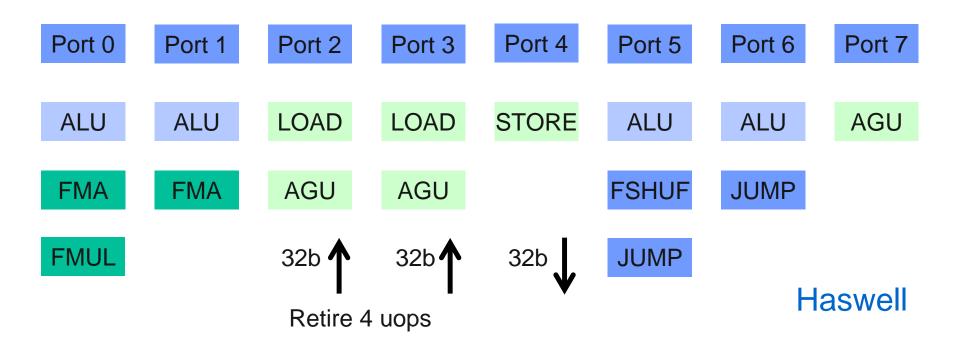
### Expected performance:

$$P = \min(P_{\text{max}}, I \cdot b_S) = \min\left(P_{\text{max}}, \frac{b_S^{\prime}}{B_C}\right)$$
[Byte/Flop]

## Preliminary: Estimating $P_{\text{max}}$



Every new CPU generation provides incremental improvements.



#### Example: Estimate $P_{\text{max}}$ of vector triad on Haswell



```
double *A, *B, *C, *D;
for (int i=0; i<N; i++) {
    A[i] = B[i] + C[i] * D[i];
}</pre>
```

Minimum number of cycles to process one AVX-vectorized iteration (one core)?

→ Equivalent to 4 scalar iterations

```
Cycle 1: LOAD + LOAD + STORE
```

Cycle 2: LOAD + LOAD + FMA + FMA

Cycle 3: LOAD + LOAD + STORE

**Answer: 1.5 cycles** 

## Example: Estimate $P_{\text{max}}$ of vector triad on Haswell (2.3 GHz)



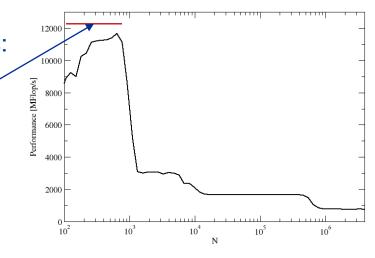
```
double *A, *B, *C, *D;
for (int i=0; i<N; i++) {
    A[i] = B[i] + C[i] * D[i];
}</pre>
```

## What is the performance in GFlops/s per core and the bandwidth in GBytes/s?

One AVX iteration (1.5 cycles) does  $4 \times 2 = 8$  flops:

$$\frac{2.3 \cdot 10^9 \text{ cy/s}}{1.5 \text{ cy}} \cdot 4 \text{ updates} \cdot \frac{2 \text{ flops}}{\text{update}} = 12.27 \frac{\text{Gflops}}{\text{s}}$$

$$6.13 \cdot 10^9 \frac{\text{updates}}{\text{s}} \cdot 32 \frac{\text{bytes}}{\text{update}} = 196 \frac{\text{Gbyte}}{\text{s}}$$



## $P_{\text{max}}$ + bandwidth limitations: The vector triad



Vector triad A(:)=B(:)+C(:)\*D(:) on a 2.3 GHz 14-core Haswell chip

Consider full chip (14 cores):

Memory bandwidth:  $b_S = 50$  GB/s

Code balance (incl. write allocate):

 $B_c = (4+1) \text{ Words } / 2 \text{ Flops} = 20 \text{ B/F} \rightarrow / = 0.05 \text{ F/B}$ 

 $\rightarrow$   $l \cdot b_s = 2.5$  GF/s (0.5% of peak performance)

 $P_{\text{peak}}$  / core = 36.8 Gflop/s ((8+8) Flops/cy x 2.3 GHz)

 $P_{\text{max}}$  / core = 12.27 Gflop/s (see prev. slide)

 $\rightarrow$   $P_{\text{max}}$  = 14 \* 12.27 Gflop/s =172 Gflop/s (33% peak)

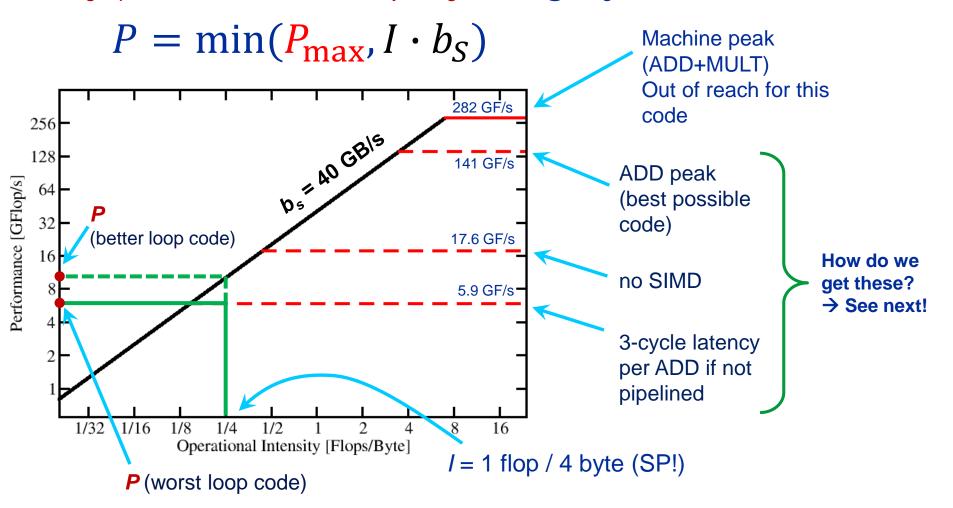
 $P = \min(P_{\text{max}}, I \cdot b_S) = \min(172,2.5) \text{ GFlop/s}$ 

## A not so simple Roofline example



Example: do i=1,N; s=s+a(i); enddo

in single precision on a 2.2 GHz Sandy Bridge socket @ "large" N



### Applicable peak for the summation loop

#### Plain scalar code, no SIMD

```
LOAD r1.0 ← 0

i ← 1

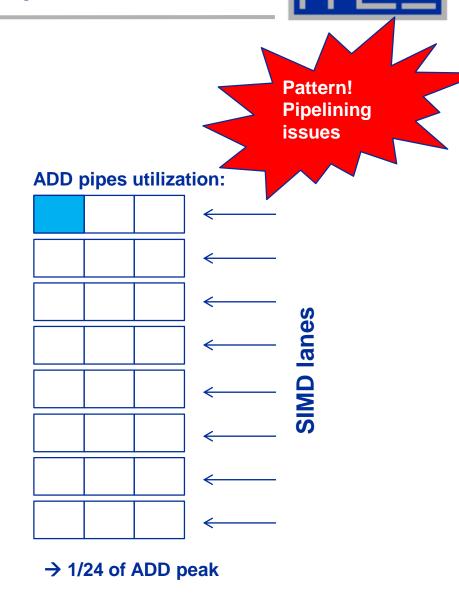
loop:

LOAD r2.0 ← a(i)

ADD r1.0 ← r1.0+r2.0

++i →? loop

result ← r1.0
```



#### Applicable peak for the summation loop



#### Scalar code, 3-way unrolling

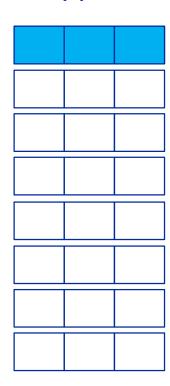
LOAD r1.0 
$$\leftarrow$$
 0  
LOAD r2.0  $\leftarrow$  0  
LOAD r3.0  $\leftarrow$  0  
i  $\leftarrow$  1

#### loop:

LOAD r4.0 
$$\leftarrow$$
 a(i)  
LOAD r5.0  $\leftarrow$  a(i+1)  
LOAD r6.0  $\leftarrow$  a(i+2)

ADD 
$$r1.0 \leftarrow r1.0 + r4.0$$
  
ADD  $r2.0 \leftarrow r2.0 + r5.0$   
ADD  $r3.0 \leftarrow r3.0 + r6.0$ 

#### **ADD pipes utilization:**



→ 1/8 of ADD peak

### Applicable peak for the summation loop

#### SIMD-vectorized, 3-way unrolled

LOAD 
$$[r1.0,...,r1.7] \leftarrow [0,...,0]$$
  
LOAD  $[r2.0,...,r2.7] \leftarrow [0,...,0]$   
LOAD  $[r3.0,...,r3.7] \leftarrow [0,...,0]$   
 $i \leftarrow 1$ 

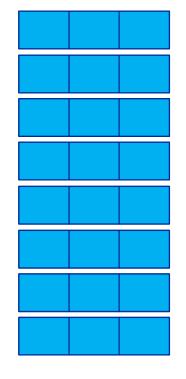
Pattern! ALU saturation

#### ADD pipes utilization:

#### loop:

LOAD 
$$[r4.0,...,r4.7] \leftarrow [a(i),...,a(i+7)]$$
  
LOAD  $[r5.0,...,r5.7] \leftarrow [a(i+8),...,a(i+15)]$   
LOAD  $[r6.0,...,r6.7] \leftarrow [a(i+16),...,a(i+23)]$ 

ADD r1 
$$\leftarrow$$
 r1 + r4  
ADD r2  $\leftarrow$  r2 + r5  
ADD r3  $\leftarrow$  r3 + r6

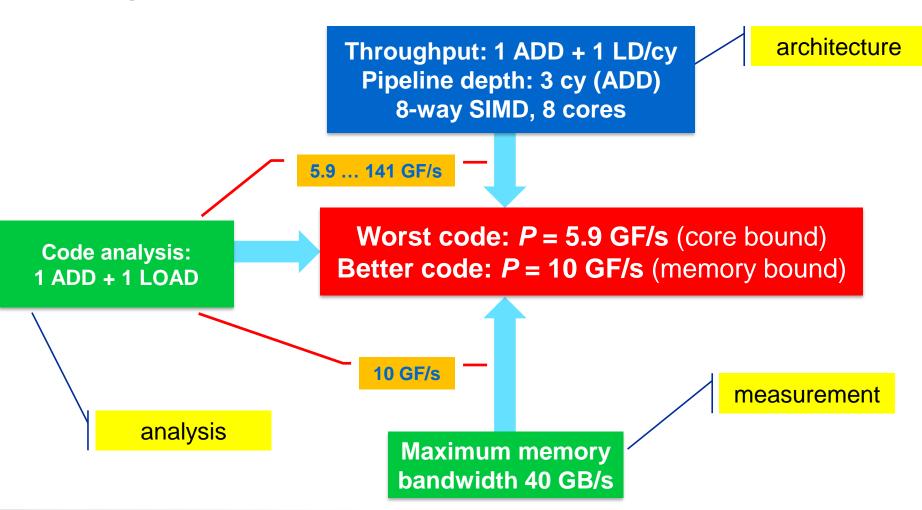


→ ADD peak

#### Input to the roofline model



... on the example of in single precision



#### **Prerequisites for the Roofline Model**



- The roofline formalism is based on some (crucial) assumptions:
  - There is a clear concept of "work" vs. "traffic"
    - "work" = flops, updates, iterations...
    - "traffic" = required data to do "work"
  - Attainable bandwidth of code = input parameter! Determine effective bandwidth via simple streaming benchmarks to model more complex kernels and applications
  - Data transfer and core execution overlap perfectly!
    - Either the limit is core execution or it is data transfer.
  - Slowest limiting factor "wins"; all others are assumed to have no impact
  - Latency effects are ignored: perfect data streaming, "steady-state" execution, no start-up effects





# **Case study: Dense Matrix Vector Multiplication**

Where does the data come from?

### **Example: Dense matrix-vector multiplication in DP (AVX)**



do 
$$c = 1$$
 ,  $C$   
do  $r = 1$  ,  $R$   

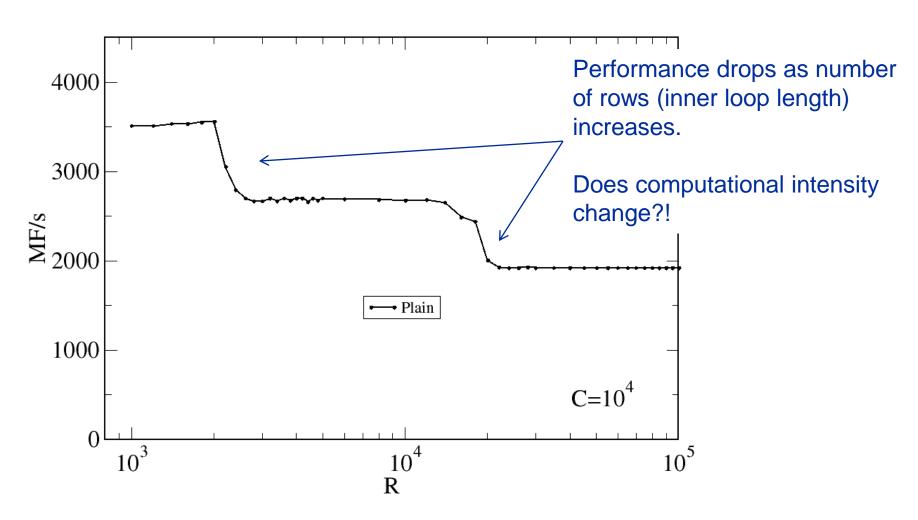
$$y(r) = y(r) + A(r,c) * x(c)$$
enddo
enddo
$$do c = 1$$

- Assume C = R ≈ 10,000
- Applicable peak performance?
- Relevant data path?
- Computational Intensity?

```
do c = 1 , C
  tmp=x(c)
  do r = 1 , R
     y(r)=y(r) + A(r,c)* tmp
  enddo
enddo
```

#### DMVM (DP) - Single core performance vs. column height





Intel Xeon E5 2695 v3 (Haswell-EP), 2.3 GHz, CoD mode, Core  $P_{peak}$ =18.4 GF/s, Caches: 32 KB / 256 KB / 35 MB, PageSize: 2 MB; ifort V15.0.1.133;  $b_S$  = 32 Gbyte/s

#### **DMVM** data traffic analysis

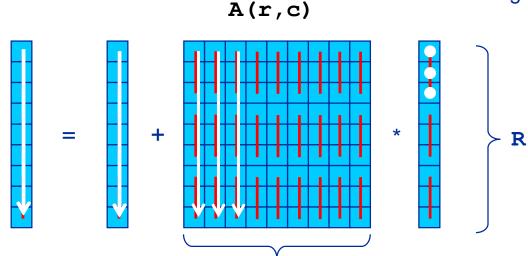


tmp stays in a register during inner loop

A(:,:) is loaded from memory – no data reuse

y(:) is loaded and stored in each outer iteration  $\rightarrow$  for c>1 update y(:) in cache

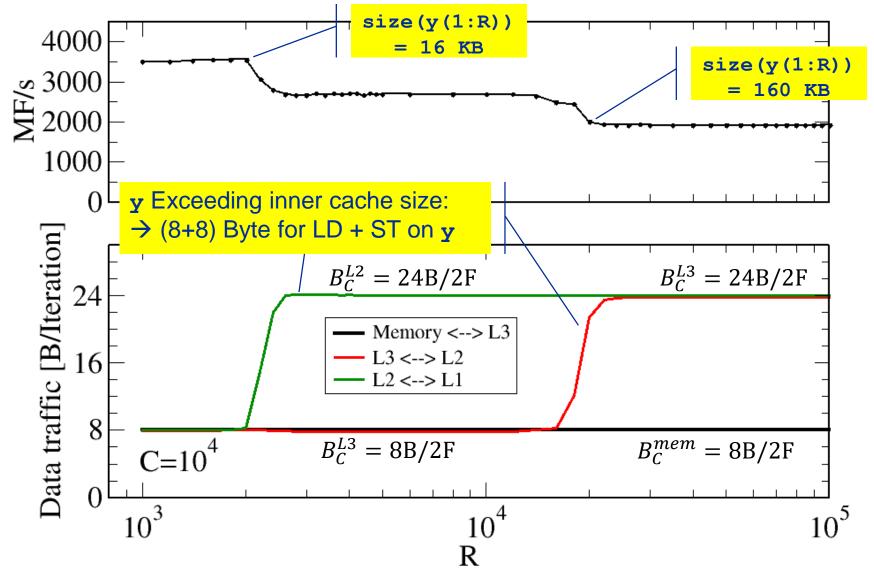
y(:) may not fit in innermost cache → more traffic from lower level caches for larger R



Roofline analysis: Distinguish code balance in memory  $(B_C^{mem})$  from code balance in relevant cache level(s)  $(B_C^{L3}, B_C^{L2},...)$ !

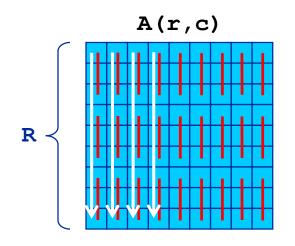
### DMVM (DP) – Single core data traffic analysis





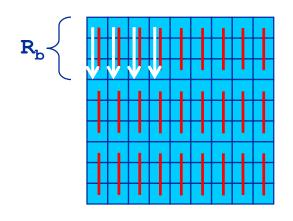
#### Reducing traffic by blocking





```
do c = 1 , C
  tmp=x(c)
  do r = 1 , R
     y(r)=y(r) + A(r,c)* tmp
  enddo
enddo
```

y(:) may not fit into some cache → more traffic for lower level

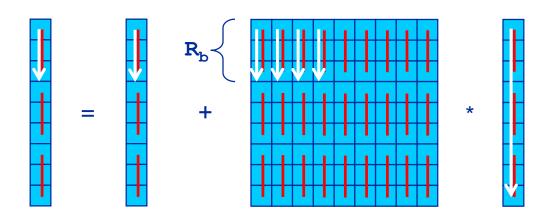


```
do rb = 1 , R , R<sub>b</sub>
  rbS = rb
  rbE = min((rb+R<sub>b</sub>-1), R)
  do c = 1 , C
    do r = rbS , rbE
      y(r)=y(r) + A(r,c)*x(c)
  enddo
enddo
enddo
```

y (rbs:rbE)
may fit into
some cache if
R<sub>b</sub> is small
enough
→ traffic
reduction

### Reducing traffic by blocking

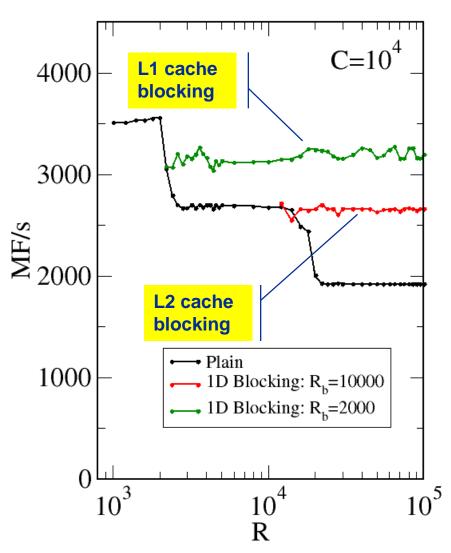




- LHS only updated once in some cache level if blocking is applied
- Price: RHS is loaded multiple times instead of once!
  - How often? → R / R<sub>b</sub> times
- Consequence: For R=C, traffic reduction of LHS & RHS by a factor of R / (1+R/[2R<sub>b</sub>])
  - Still a large reduction if block size can be made larger than about 10

### DMVM (DP) – Reducing traffic by inner loop blocking





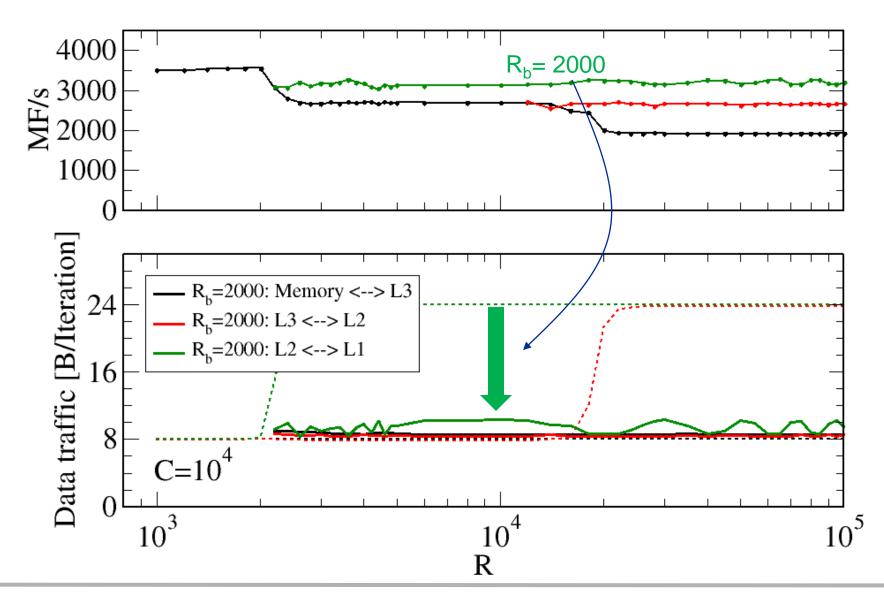
- "1D blocking" for inner loop
- Blocking factor R<sub>b</sub> ←→ cache level

```
do rb = 1 , R , R_b
 rbS = rb
 rbE = min((rb+R_b-1), R)
 do c = 1 , C
   do r = rbS , rbE
      y(r) = y(r) + A(r,c) *x(c)
   enddo
 enddo
enddo
```

10<sup>5</sup> → Fully reuse subset of y (rbS:rbE) from L1/L2 cache

### DMVM (DP) – Validation of blocking optimization





#### **DMVM (DP) – OpenMP parallelization**



```
!$omp parallel do reduction(+:y)
do c = 1 , C
  do r = 1 , R
     y(r) = y(r) + A(r,c) * x(c)
enddo ; enddo
!$omp end parallel do plain code
```

```
!$omp parallel do private(rbS,rbE)
do rb = 1 , R , R<sub>b</sub>
rbS = rb
rbE = min((rb+R<sub>b</sub>-1), R)
do c = 1 , C
   do r = rbS , rbE
      y(r) = y(r) + A(r,c) * x(c)
enddo ; enddo
!$omp end parallel do blocked code
```

#### DMVM (DP) – OpenMP parallelization & saturation "Using more cores can heal bad single-core performance" Roofline limit **Bandwidth** 8000 $B_C = 4$ Byte/Flop $b_{\rm S} = 32 {\rm GB/s}$ memory traffic unchanged 6000 → saturation unchanged! ¥ 4000 saturation influenced by serial performance • Plain So, is blocking 1D blocking: R<sub>b</sub>=2000; par. blocks 2000 useless? → NO (see later) 0 3 5 6 Can we do

blocking good for single thread performance (reduced in-cache traffic)

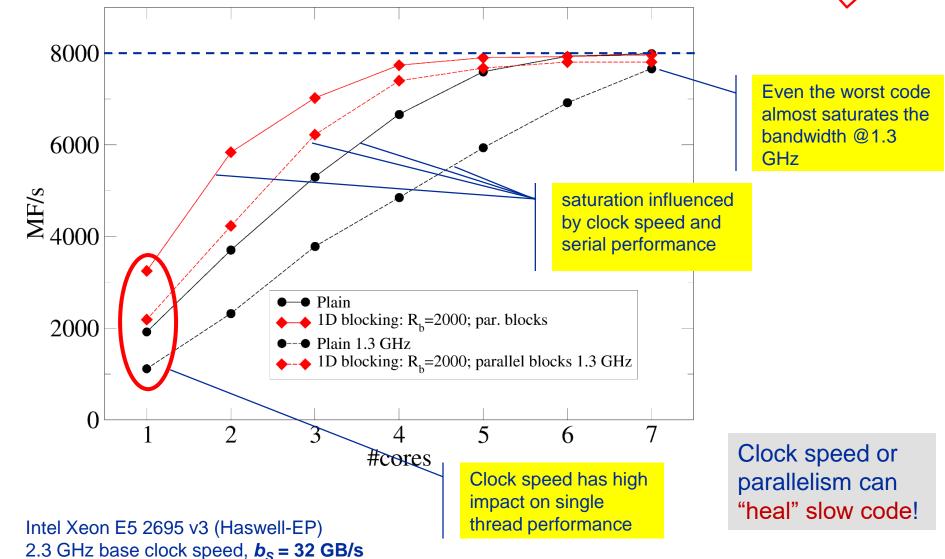
blocking good for single thread performance (reduced in-cache traffic)

Can we do anything to improve  $B_C^{mem}$ ?

→ NO, not here

### DMVM (DP) - saturation & clock speed





## **Conclusions from the dMVM example**

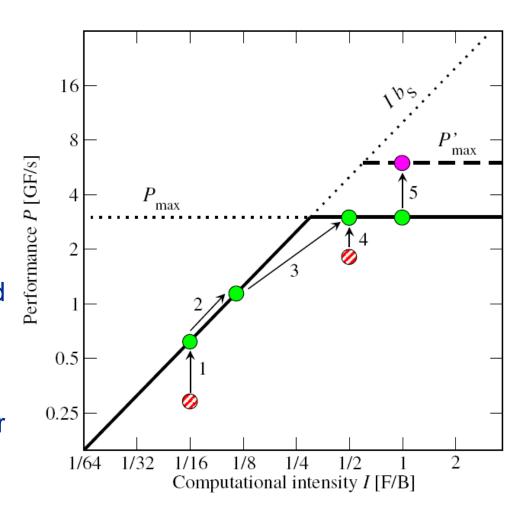
- We have found the reasons for the breakdown of single-core performance with growing number of matrix rows
  - LHS vector fitting in different levels of the cache hierarchy
  - Validated theory by performance counter measurements
- Inner loop blocking was employed to improve code balance in L3 and/or L2
  - Validated by performance counter measurements
- Blocking led to better single-threaded performance

- Saturated performance unchanged (as predicted by Roofline)
  - Because the problem is still small enough to fit the LHS at least into the L3 cache

#### Typical code optimizations in the Roofline Model



- Hit the BW bottleneck by good serial code (e.g., Perl → Fortran)
- Increase intensity to make better use of BW bottleneck (e.g., loop blocking → see later)
- 3. Increase intensity and go from memory-bound to core-bound (e.g., temporal blocking)
- Hit the core bottleneck by good serial code
   (e.g., -fno-alias → see later)
- Shift P<sub>max</sub> by accessing additional hardware features or using a different algorithm/implementation (e.g., scalar → SIMD)



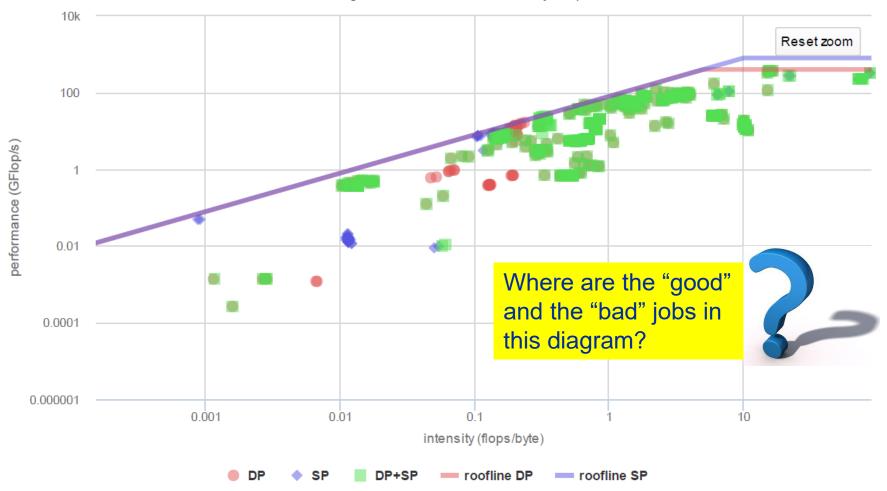
#### Using Roofline for monitoring "live" jobs on a cluster

Based on measured BW and Flop/s data via likwid-perfctr



#### Ganglia Data / Roofline (04. Feb. 2016 - 14:12:24)

Click and drag to zoom in. Hold down shift key to x-pan.





## Case study: A Jacobi smoother

The basic performance properties in 2D

Layer conditions

**Optimization by spatial blocking** 

#### Stencil schemes



- Stencil schemes frequently occur in PDE solvers on regular lattice structures
- Basically it is a sparse matrix vector multiply (spMVM) embedded in an iterative scheme (outer loop)
- but the regular access structure allows for matrix free coding

```
do iter = 1, max_iterations

Perform sweep over regular grid: y(:) 
Swap y 
x
```

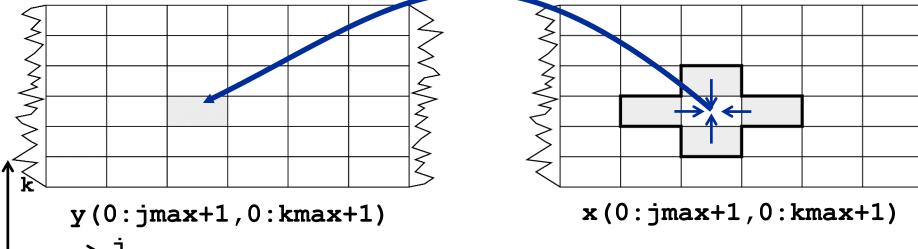
enddo

- Complexity of implementation and performance depends on
  - stencil operator, e.g. Jacobi-type, Gauss-Seidel-type, ...
  - spatial extent, e.g. 7-pt or 25-pt in 3D,...

## Jacobi-type 5-pt stencil in 2D



Lattice Update (LUP)



Appropriate performance metric: "Lattice Updates per second" [LUP/s]

(here: Multiply by 4 FLOP/LUP to get FLOP/s rate)

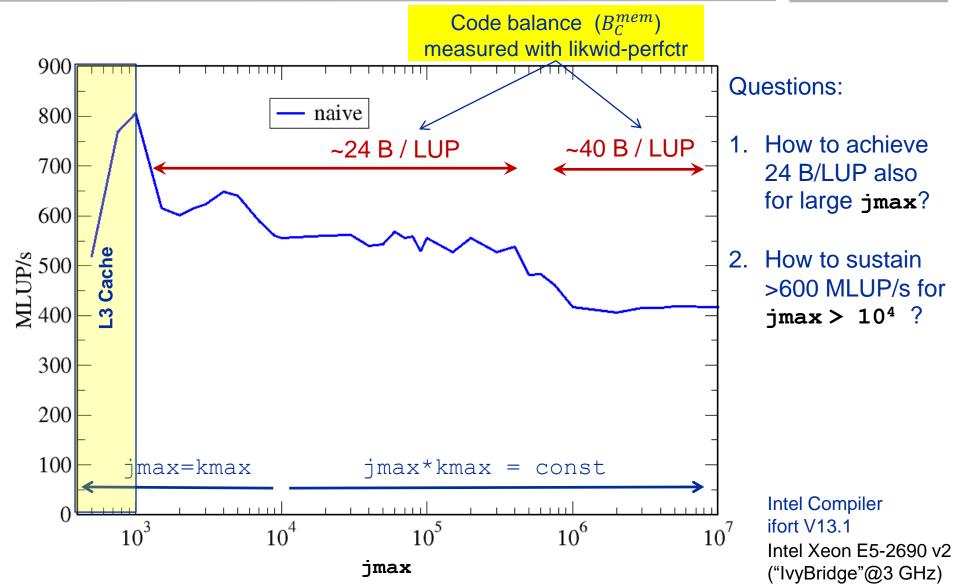
## Jacobi 5-pt stencil in 2D: data transfer analysis

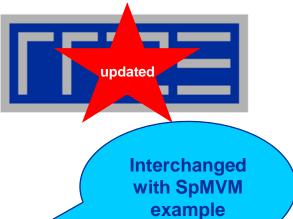


```
Available in cache
    LD+ST y(i,k)
                                       (used 2 updates before) ,
      (incl. write
       allocate)
                                                            LD x(j+1,k)
    do k=1,kmax
      do j=1,jmax
        y(j,k) = const * ( x(j-1,k) + x(j+1,k) &
                                   x(j,k-1) + x(j,k+1))
      enddo
    enddo
                                         LD \times (j,k-1)
                                                            LD x(j,k+1)
Naive balance (incl. write allocate):
x(:,:):3LD+
y(:,:):1 ST+1LD
\rightarrow B<sub>c</sub> = 5 Words / LUP = 40 B / LUP (assuming double precision)
```

### Jacobi 5-pt stencil in 2D: Single core performance







## Case study: A Jacobi smoother

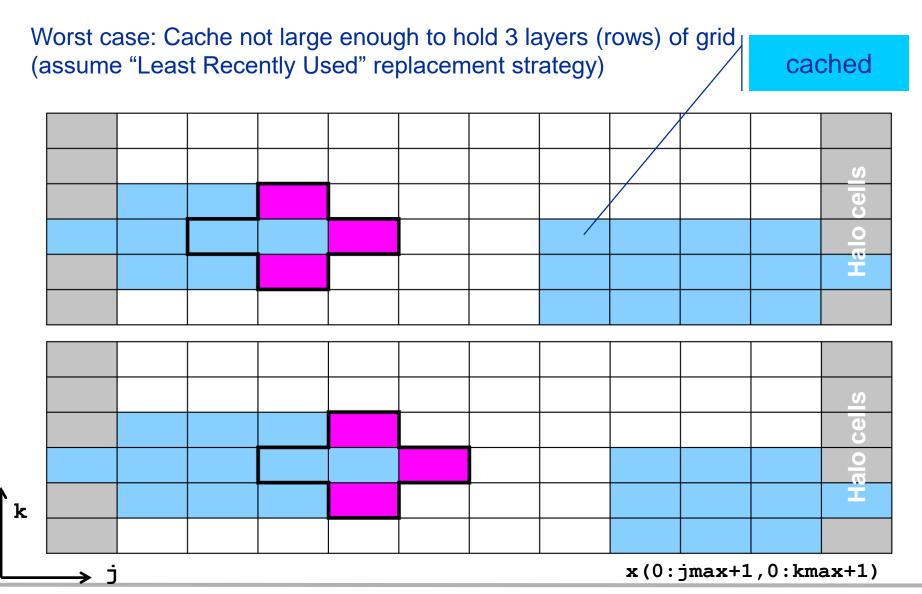
The basics in two dimensions

Layer conditions

Optimization by spatial blocking

#### **Analyzing the data flow**

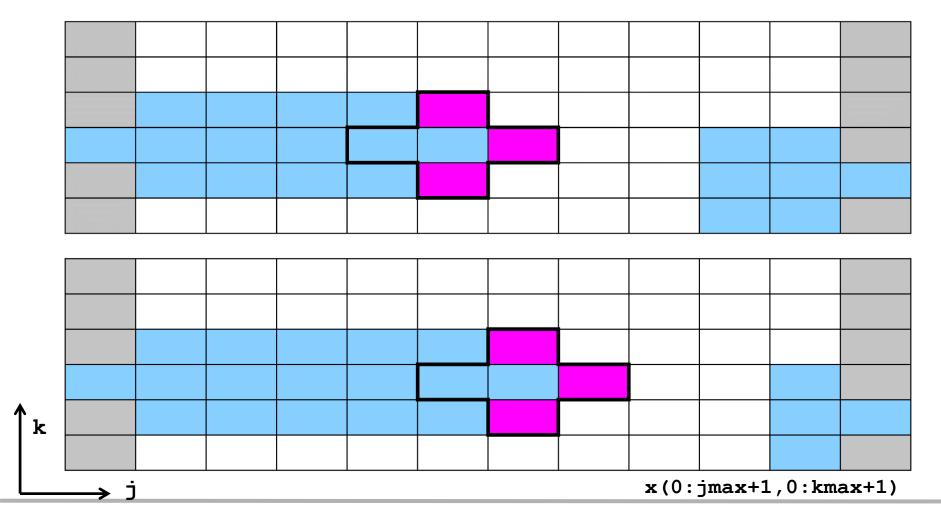




### **Analyzing the data flow**

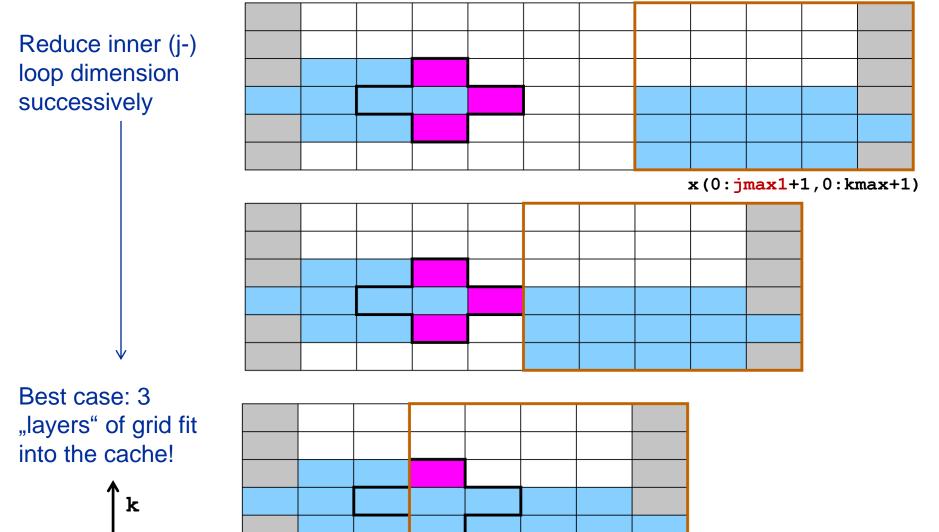


Worst case: Cache not large enough to hold 3 layers (rows) of grid (+assume "Least Recently Used" replacement strategy)



## **Analyzing the data flow**

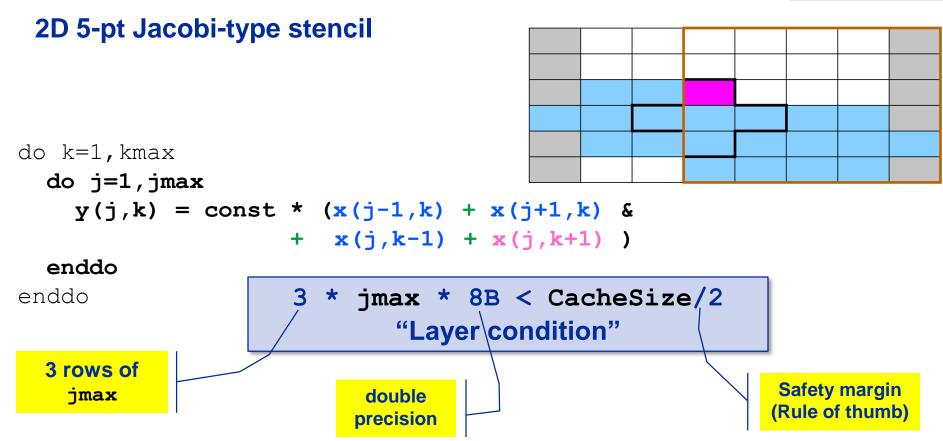




x(0:jmax2+1,0:kmax+1)

## Analyzing the data flow: Layer condition





#### Layer condition:

- Does not depend on outer loop length (kmax)
- No strict guideline (cache associativity data traffic for y not included)
- Needs to be adapted for other stencils (e.g., 3D 7-pt stencil)

## Analyzing the data flow: Layer condition (2D 5-pt Jacobi)



```
y: (1 LD + 1 ST) / LUP
                                                                x: 1 LD / LUP
           do k=1, kmax
             do j=1,jmax
               y(j,k) = const * (x(j-1,k) + x(j+1,k) &
                               + x(j,k-1) + x(j,k+1))
YES
             enddo
                                                              B_C = 24 B / LUP
           enddo
 3 * jmax * 8B < CacheSize/2
      Layer condition fulfilled?
                                      y: (1 LD + 1 ST) / LUP
           do k=1, kmax
             do j=1,jmax
NO
               y(j,k) = const * (x(j-1,k) + x(j+1,k) &
                               + x(j,k-1) + x(j,k+1))
             enddo
                           x: 3 LD / LUP
                                                             B_C = 40 B / LUP
           enddo
```

### Analyzing the data flow: Layer condition (2D 5-pt Jacobi)



- Establish layer condition for all domain sizes
- Idea: Spatial blocking
  - Reuse elements of x () as long as they stay in cache
  - Sweep can be executed in any order, e.g. compute blocks in j-direction

### → "Spatial Blocking" of j-loop:

→ Determine for given CacheSize an appropriate jblock value:

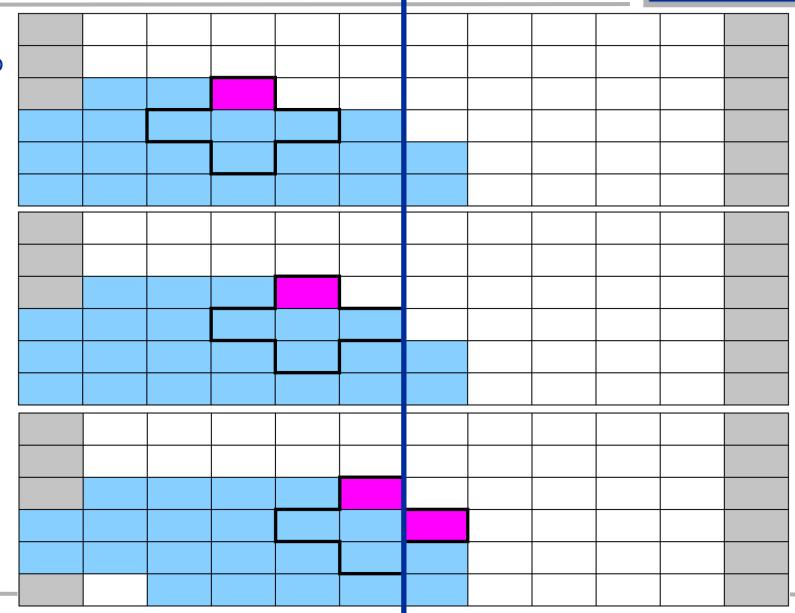
```
jblock < CacheSize / 48 B</pre>
```

## Establish the layer condition by blocking



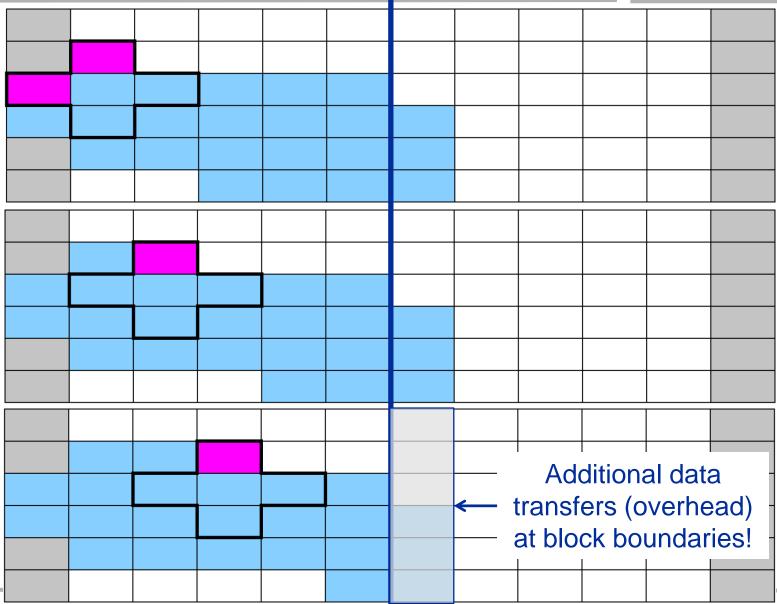
Split up domain into subblocks:

e.g. block size = 5



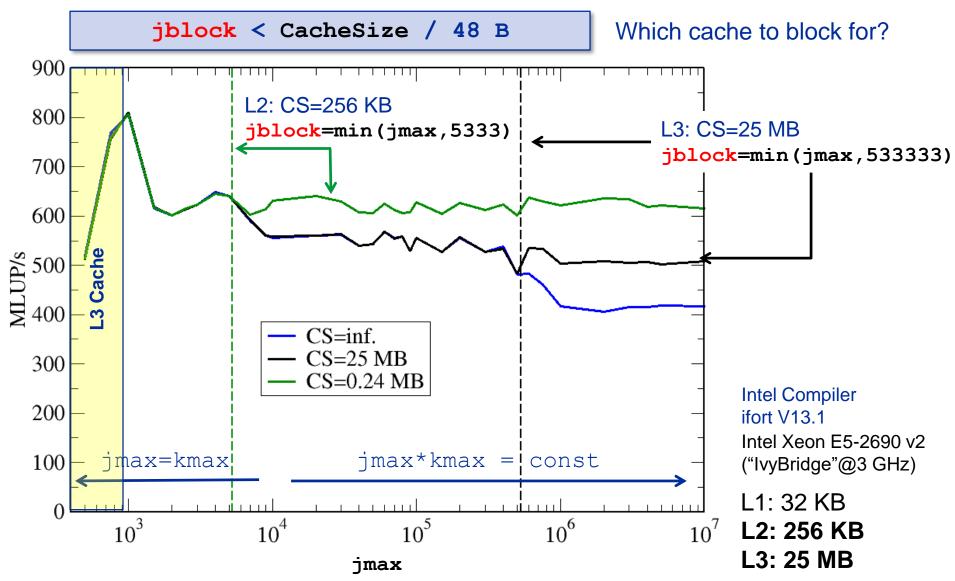
## Establish the layer condition by blocking





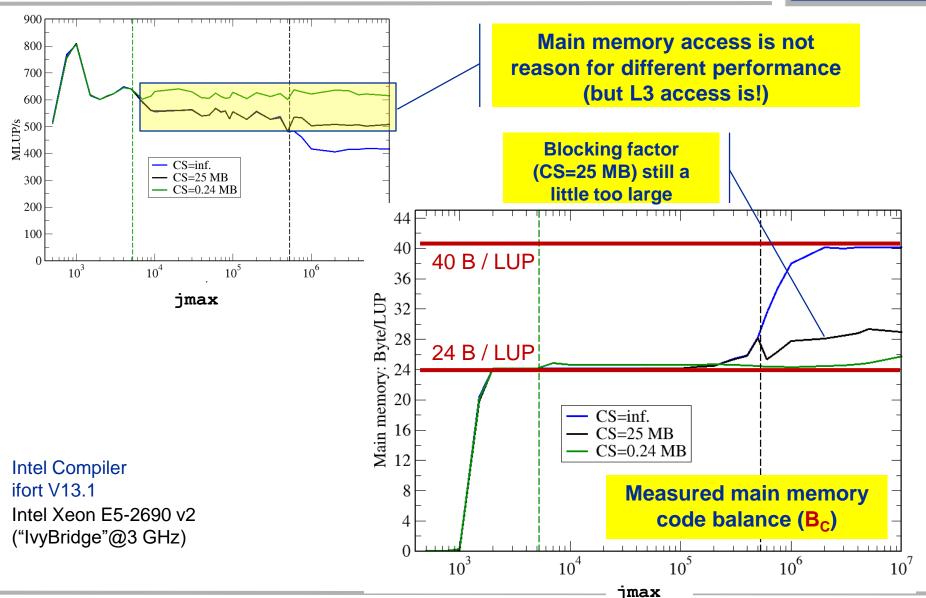
# Establish layer condition by spatial blocking





## Layer condition & spatial blocking: Memory code balance





## Jacobi Stencil - OpenMP parallelization



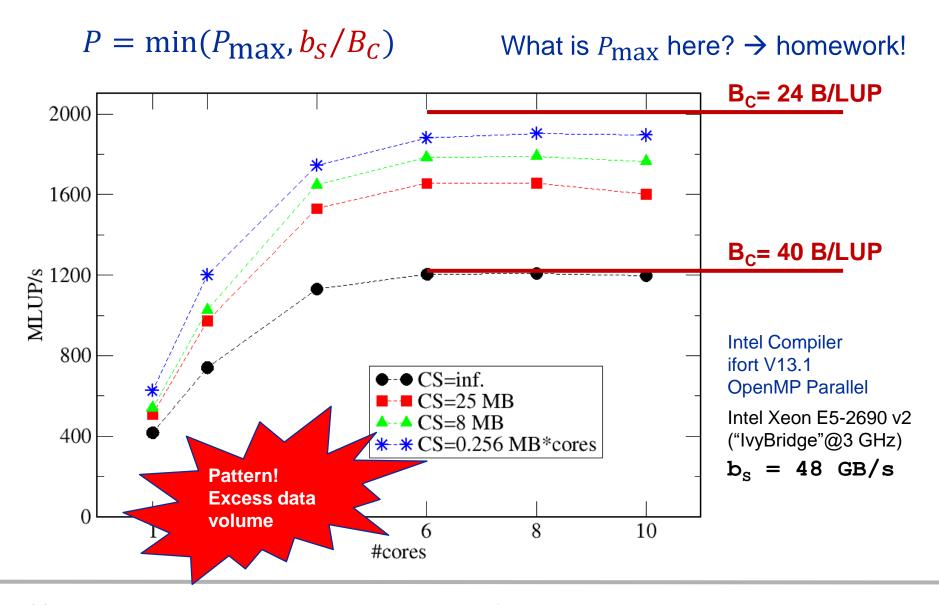
Equally large chunks in k-direction→ "Layer condition" for each thread

```
"Layer condition" for shared cache:
nthreads * 3 *imax * 8B < CS/2
```

Homework: what about if the cache to block for is not shared among the threads?

## Socket scaling – validate Roofline model





#### **Conclusions from the Jacobi example**



- We have made sense of the memory-bound performance vs. problem size
  - "Layer conditions" lead to predictions of code balance
  - "What part of the data comes from where" is a crucial question
  - The model works only if the bandwidth is "saturated"
    - In-cache modeling is more involved
- Avoiding slow data paths == re-establishing the most favorable layer condition
- Improved code showed the speedup predicted by the model
- Optimal blocking factor can be estimated
  - Be guided by the cache size the layer condition
  - No need for exhaustive scan of "optimization space"
- Food for thought
  - Multi-dimensional loop blocking would it make sense?
  - Can we choose a "better" OpenMP loop schedule?
  - What would change if we parallelized inner loops?



# **Case study: Sparse Matrix Vector Multiplication**

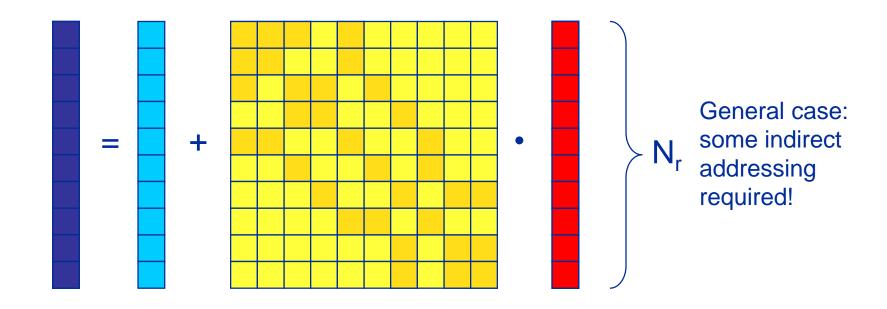
An example for what to do if the Roofline model isn't quite so predictive

# **Sparse Matrix Vector Multiplication**



(spMVM) Key ingredient in some matrix diagonalization algorithms

- Lanczos, Davidson, Jacobi-Davidson
- Store only N<sub>nz</sub> nonzero elements of matrix and RHS, LHS vectors with N<sub>r</sub> (number of matrix rows) entries
- "Sparse": N<sub>nz</sub> ~ N<sub>r</sub>



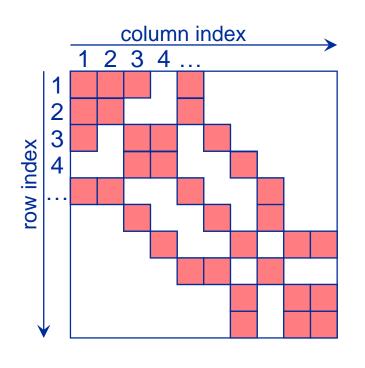
# **SpMVM** characteristics



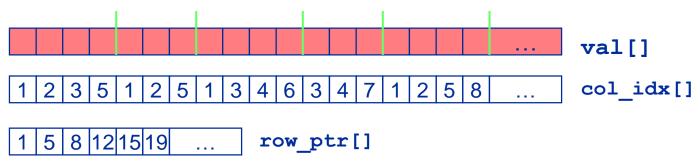
- For large problems, spMVM is inevitably memory-bound
  - Intra-socket saturation effect on modern multicores
- SpMVM is easily parallelizable in shared and distributed memory
- Data storage format is crucial for performance properties
  - Most useful general format on CPUs: Compressed Row Storage (CRS)
  - Depending on compute architecture

# **CRS** matrix storage scheme





- val[] stores all the nonzeros (length N<sub>nz</sub>)
- col\_idx[] stores the column index of each nonzero (length N<sub>nz</sub>)
- row\_ptr[] stores the starting index of each new row in val[] (length: N<sub>r</sub>)



# **Case study: Sparse matrix-vector multiply**



- Strongly memory-bound for large data sets
  - Streaming, with partially indirect access:

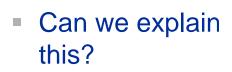
```
!$OMP parallel do
do i = 1,Nr
do j = row_ptr(i), row_ptr(i+1) - 1
   c(i) = c(i) + val(j) * b(col_idx(j))
   enddo
enddo
!$OMP end parallel do
```

- Usually many spMVMs required to solve a problem
- Now let's look at some performance measurements...

## **Performance characteristics**

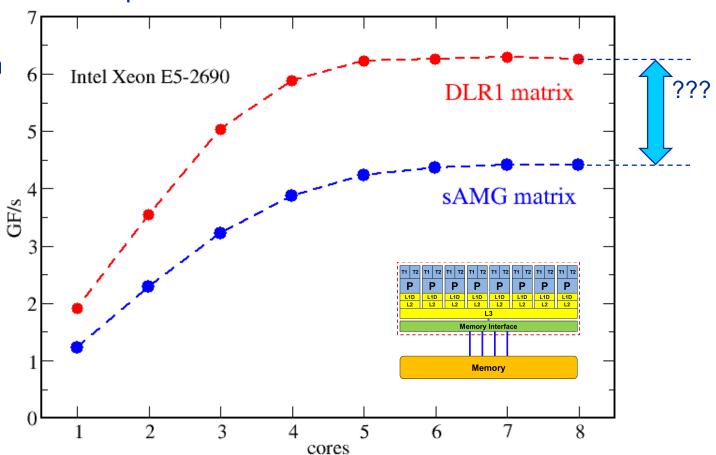


- Strongly memory-bound for large data sets → saturating performance across cores on the chip
- Performance seems to depend on the matrix



Is there a "light speed" for spMVM?

Optimization?



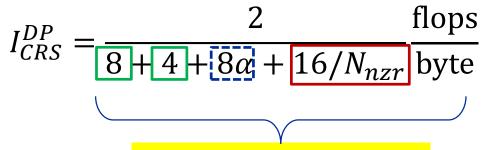
## **Example: SpMVM node performance model**



Sparse MVM in double precision w/ CRS data storage:

do i = 1,
$$N_r$$
  
do j = row\_ptr(i), row\_ptr(i+1) - 1  
 $C(i) = C(i) + val(j) * B(col_idx(j))$   
enddo  
enddo

- DP CRS comp. intensity
  - α quantifies traffic for loading RHS
    - $\alpha = 0 \rightarrow RHS$  is in cache
    - $\alpha = 1/N_{nzr} \rightarrow RHS$  loaded once
    - $\alpha = 1 \rightarrow$  no cache
    - $\alpha > 1 \rightarrow$  Houston, we have a problem!
  - "Expected" performance = b<sub>S</sub> x I<sub>CRS</sub>
  - Determine α by measuring performance and actual memory traffic
    - Maximum memory BW may not be achieved with spMVM



"best" in-memory intensity:  $I = \frac{1}{6}$  F/B, or  $B_C^{mem} = 6$  B/F

Node-Level Performance Engineering

### **Determine RHS traffic**



$$I_{CRS}^{DP} = \frac{2}{8+4+8\alpha+16/N_{nzr}} \frac{\text{flops}}{\text{byte}} = \frac{N_{nz} \cdot 2 \text{ flops}}{V_{meas}}$$

- $V_{meas}$  is the measured overall memory data traffic (using, e.g., likwid-perfctr)
- Solve for  $\alpha$ :

$$\alpha = \frac{1}{4} \left( \frac{V_{meas}}{N_{nz} \cdot 2 \text{ bytes}} - 6 - \frac{8}{N_{nzr}} \right)$$

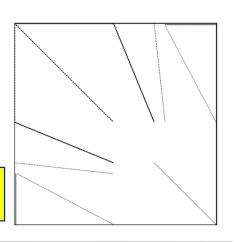
 Example: kkt\_power matrix from the UoF collection on one Intel SNB socket

$$N_{nz} = 14.6 \cdot 10^6$$
,  $N_{nzr} = 7.1$ 

- $V_{meas} \approx 258 \text{ MB}$
- $\rightarrow \alpha = 0.43, \, \alpha N_{nzr} = 3.1$
- → RHS is loaded 3.1 times from memory
- and:

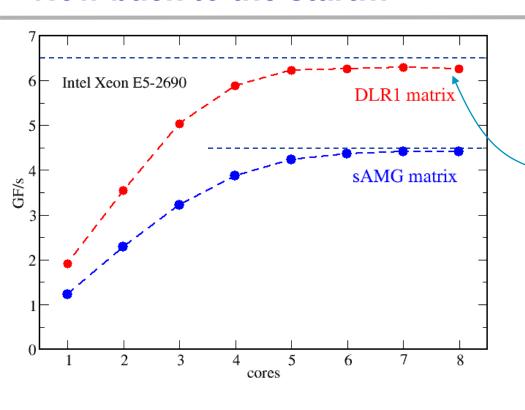
$$\frac{I_{CRS}^{DP}(1/N_{nzr})}{I_{CRS}^{DP}(\alpha)} = 1.15$$

15% extra traffic → optimization potential!



### Now back to the start...





$$b_{\rm S} = 39 \, {\rm GB/s}$$

$$B_c^{min} = 6 \, B/F$$

Maximum spMVM performance:

$$P_{max} = 6.5 \,\mathrm{GF/s}$$

- → DLR1 causes minimum code balance!
- sAMG matrix code balance:

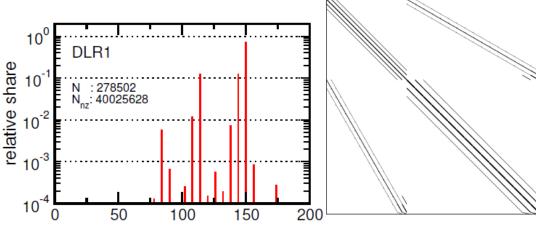
$$B_c \le \frac{b_S}{4.5 \text{ GF/s}} = 8.7 \text{ B/F}$$

- Why is this only an upper limit?
- What is the next step?
- Could we have predicted this qualitative difference?

# **Sparse matrix testcases**

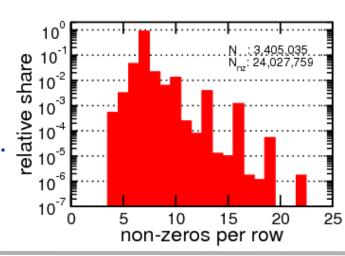
## "DLR1" (A. Basermann, DLR)

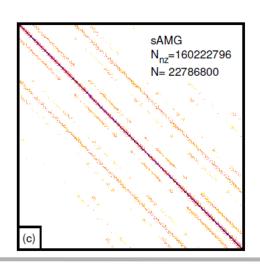
Adjoint problem computation (turbulent transonic flow over a wing) with the TAU CFD system of the German Aerospace Center (DLR) Avg. non-zeros/row ~150



## "sAMG" (K. Stüben, FhG-SCAI)

Matrix from FhG's adaptive multigrid code sAMG for the irregular discretization of a Poisson problem on a car geometry. Avg. non-zeros/row ~ 7





#### Roofline analysis for spMVM



#### Conclusion from Roofline analysis

- The roofline model does not "work" for spMVM due to the RHS traffic uncertainties
- We have "turned the model around" and measured the actual memory traffic to determine the RHS overhead
- Result indicates:
  - 1. how much actual traffic the RHS generates
  - 2. how efficient the RHS access is (compare BW with max. BW)
  - 3. how much optimization potential we have with matrix reordering

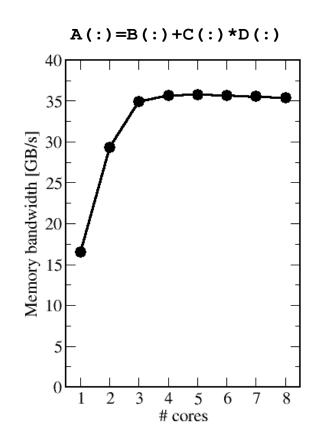
 Consequence: Modeling is not always 100% predictive. It's all about *learning more* about performance properties!

### Shortcomings of the roofline model



- Saturation effects in multicore chips are not explained
  - Reason: "saturation assumption"
  - Cache line transfers and core execution do sometimes not overlap perfectly
  - It is not sufficient to measure single-core STREAM to make it work
  - Only increased "pressure" on the memory interface can saturate the bus
     → need more cores!
- In-cache performance is not correctly predicted
- The ECM performance model gives more insight:

H. Stengel, J. Treibig, G. Hager, and G. Wellein: *Quantifying performance bottlenecks of stencil computations using the Execution-Cache-Memory model.* Proc. <u>ICS15</u>, the 29th International Conference on Supercomputing, June 8-11, 2015, Newport Beach, CA. <u>DOI: 10.1145/2751205.2751240</u>. Preprint: arXiv:1410.5010



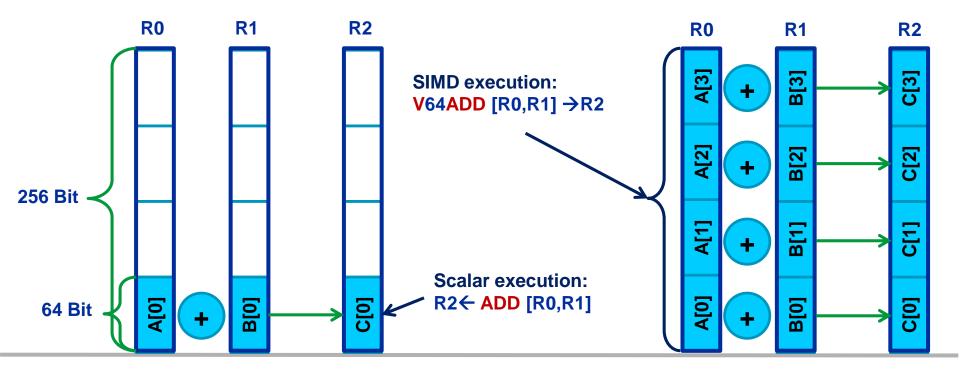


# Coding for SingleInstructionMultipleData processing

# SIMD processing – Basics



- Single Instruction Multiple Data (SIMD) operations allow the concurrent execution of the same operation on "wide" registers.
- x86 SIMD instruction sets:
  - SSE: register width = 128 Bit → 2 double precision floating point operands
  - AVX: register width = 256 Bit → 4 double precision floating point operands
- Adding two registers holding double precision floating point operands



### **Vectorization compiler options (Intel)**



- The compiler will vectorize starting with -02.
- To enable specific SIMD extensions use the –x option:
  - -xSSE2 vectorize for SSE2 capable machines Available SIMD extensions: SSE2, SSE3, SSSE3, SSE4.1, SSE4.2, AVX
  - -xAVX (2) on Sandy Bridge (Haswell) processors

#### Recommended option:

-xHost will optimize for the architecture you compile on

On AMD Opteron: use plain -o3 as the -x options may involve CPU type checks.

#### **Vectorization compiler options**



#### Controlling non-temporal stores (part of the SIMD extensions)

-opt-streaming-stores always|auto|never

**always** use NT stores, assume application is memory

bound (use with caution!)

auto compiler decides when to use NT stores

**never** do not use NT stores unless activated by

source code directive

#### **Vectorization source code directives**



- Fine-grained control of loop vectorization
- Use !DEC\$ (Fortran) or #pragma (C/C++) sentinel to start a compiler directive
- #pragma vector always
  vectorize even if it seems inefficient (hint!)
- #pragma novector do not vectorize even if possible
- #pragma vector nontemporal use NT stores when allowed (i.e. alignment conditions are met)
- #pragma vector aligned specifies that all array accesses are aligned to 16-byte boundaries (DANGEROUS! You must not lie about this!)

#### **User mandated vectorization**



- Since Intel Compiler 12.0 the simd pragma is available
- #pragma simd enforces vectorization where the other pragmas fail
- Prerequesites:
  - Countable loop
  - Innermost loop
  - Must conform to for-loop style of OpenMP worksharing constructs
- There are additional clauses: reduction, vectorlength, private
- Refer to the compiler manual for further details

```
#pragma simd reduction(+:x)
  for (int i=0; i<n; i++) {
      x = x + A[i];
}</pre>
```

 NOTE: Using the #pragma simd the compiler may generate incorrect code if the loop violates the vectorization rules!

#### **x86 Architecture:**

### SIMD and Alignment



#### Alignment issues

- Alignment of arrays with SSE (AVX) should be on 16-byte (32-byte) boundaries to allow packed aligned loads and NT stores (for Intel processors)
  - AMD has a scalar nontemporal store instruction
- Otherwise the compiler will revert to unaligned loads and not use NT stores – even if you say vector nontemporal
- Modern x86 CPUs have less (not zero) impact for misaligned LD/ST, but Xeon Phi (KNC) relies heavily on it!
- How is manual alignment accomplished?
- Dynamic allocation of aligned memory (align = alignment boundary):



# Reading x86 assembly code and exploiting SIMD parallelism

Understanding SIMD execution by inspecting assembly code

SIMD vectorization how-to

Intel compiler options and features for SIMD

### Why and how?



#### Why check the assembly code?

- Sometimes the only way to make sure the compiler "did the right thing"
  - Example: "LOOP WAS VECTORIZED" message is printed, but Loads & Stores may still be scalar!
- Get the assembler code (Intel compiler):

```
icc -S -masm=intel -O3 -xHost triad.c -o a.out
```

Disassemble Executable:

```
objdump -d ./a.out | less
```

#### The x86 ISA is documented in:

Intel Software Development Manual (SDM) 2A and 2B AMD64 Architecture Programmer's Manual Vol. 1-5

#### Basics of the x86-64 ISA



- Instructions have 0 to 4 operands
- Operands can be registers, memory references or immediates
- Opcodes (binary representation of instructions) vary from 1 to 17 bytes
- There are two syntax forms: Intel (left) and AT&T (right)
- Addressing Mode: BASE + INDEX \* SCALE + DISPLACEMENT
- C: A[i] equivalent to \*(A+i) (a pointer has a type: A+i\*8)

```
movaps [rdi + rax*8+48], xmm3
add rax, 8
js 1b
```

```
movaps %xmm4, 48(%rdi,%rax,8)
addq $8, %rax
js ..B1.4
```

```
401b9f: 0f 29 5c c7 30 movaps %xmm3,0x30(%rdi,%rax,8)
401ba4: 48 83 c0 08 add $0x8,%rax
401ba8: 78 a6 js 401b50 <triad_asm+0x4b>
```

#### Basics of the x86-64 ISA



```
16 general Purpose Registers (64bit):
```

```
rax, rbx, rcx, rdx, rsi, rdi, rsp, rbp, r8-r15
alias with eight 32 bit register set:
eax, ebx, ecx, edx, esi, edi, esp, ebp
```

#### Floating Point SIMD Registers:

```
xmm0-xmm15 SSE (128bit) alias with 256-bit registers ymm0-ymm15 AVX (256bit)
```

#### SIMD instructions are distinguished by:

AVX (VEX) prefix: v

Operation: mul, add, mov

Modifier: nontemporal (nt), unaligned (u), aligned (a), high (h)

Width: scalar (s), packed (p)

Data type: single (s), double (d)

### Case Study: Vector Triad (DP) on IvyBridge-EP



```
for (int i = 0; i < length; i++) {
    A[i] = B[i] + D[i] * C[i];
}
To get
objdu
execu-</pre>
```

To get object code use objdump -d on object file or executable or compile with -s

#### Assembly code (-O1):

```
LBB0 3
movsd
       xmm0, [rdx]
mulsd
       xmm0, [rcx]
       xmm0, [rsi]
addsd
       [rax], xmm0
movsd
add
        rsi, 8
add
        rdx, 8
add
        rcx, 8
add
        rax, 8
dec
        edi
        LBB0 3
jne
```

```
..B1.6:

movsd xmm0, [r12+rax*8]

mulsd xmm0, [r13+rax*8]

addsd xmm0, [r14+rax*8]

movsd [r15+rax*8], xmm0

inc rax

cmp rax, rbx

jl ..B1.6
```

7 instructions per loop iteration

```
.L4:
movsd xmm0,[rbx+rax]
mulsd xmm0,[r12+rax]
addsd xmm0,[r13+0+rax]
movsd [rbp+0+rax],xmm0
add rax, 8
cmp rax, r14
jne .L4
```

#### Case Study: Vector Triad (DP) -O3 -xHost

+ #pragma vector aligned

```
updated
```

```
..B1.7:
          ymm0, [r15+rcx*8]
vmovupd
vmovupd
          ymm4, [32+r15+rcx*8]
          ymm7, [64+r15+rcx*8]
vmovupd
          ymm10, [96+r15+rcx*8]
vmovupd
vmulpd
          ymm2, ymm0, [rdx+rcx*8]
vmulpd
          ymm5, ymm4, [32+rdx+rcx*8]
vmulpd
          ymm8, ymm7, [64+rdx+rcx*8]
vmulpd
          ymm11, ymm10, [96+rdx+rcx*8]
          ymm3, ymm2, [r14+rcx*8]
vaddpd
vaddpd
          ymm6, ymm5, [32+r14+rcx*8]
          ymm9, ymm8, [64+r14+rcx*8]
vaddpd
vaddpd
          ymm12, ymm11, [96+r14+rcx*8]
         [r13+rcx*8], ymm3
vmovupd
vmovupd
         [32+r13+rcx*8], ymm6
          [64+r13+rcx*8], ymm9
vmovupd
         [96+r13+rcx*8], ymm12
vmovupd
          rcx, 16
add
cmp
          rcx, rsi
          ..B1.7
ήb
```

# 1.19 instructions per loop iteration

# Case Study: Vector Triad (DP) -O3 -xHost #pragma vector aligned on Haswell-EP



```
..B1.7:
         ymm2, [r15+rcx*8]
vmovupd
         ymm4, [32+r15+rcx*8]
vmovupd
        ymm6, [64+r15+rcx*8]
vmovupd
vmovupd ymm8, [96+r15+rcx*8]
vmovupd ymm0, [rdx+rcx*8]
vmovupd ymm3, [32+rdx+rcx*8]
vmovupd vmm5, [64+rdx+rcx*8]
vmovupd ymm7, [96+rdx+rcx*8]
vfmadd213pd ymm2, ymm0, [r14+rcx*8]
vfmadd213pd ymm4, ymm3, [32+r14+rcx*8]
vfmadd213pd ymm6, ymm5, [64+r14+rcx*8]
vfmadd213pd ymm8, ymm7, [96+r14+rcx*8]
vmovupd [r13+rcx*8], ymm2
vmovupd [32+r13+rcx*8], ymm4
vmovupd [64+r13+rcx*8], ymm6
vmovupd [96+r13+rcx*8], ymm8
add
         rcx, 16
         rcx, rsi
cmp
          ..B1.7
jb
```

On X86 ISA instruction are converted to so-called µops (elementary ops like load, add, mult). For performance the number of µops is important.

23 uops vs. 27 μops (AVX)

# 1.19 instructions per loop iteration

#### **SIMD** processing – The whole picture



| <b>SIMD</b> influences instruction |
|------------------------------------|
| execution in the core - other      |
| runtime contributions stay the     |
| same!                              |

#### Comparing total execution time (cycles):

| same!                              | !                             |        | Execution |           | Memory    |
|------------------------------------|-------------------------------|--------|-----------|-----------|-----------|
|                                    | See P <sub>max</sub> analysis | Scalar | 12        | 15        | <b>21</b> |
| AVX example (Haswell-EP): analysis |                               | SSE    | 6         | <b>15</b> | <b>21</b> |
| Execution                          | (3 cy)                        | AVX    | 3         | 15        | 21        |

Caches

**Memory** 

15 cy

21 cy

Per-cacheline (8 iterations) cycle counts

Total runtime with data loaded from memory:

Scalar 48 cy SSE 42 cy AVX 39 cy

SIMD is only effective if runtime is dominated by instruction execution!

#### How to leverage SIMD: your options



#### **Alternatives:**

- The compiler does it for you (but: aliasing, alignment, language)
- Compiler directives (pragmas)
- Alternative programming models for compute kernels (OpenCL, ispc)
- Intrinsics (restricted to C/C++)
- Implement directly in assembler

#### To use intrinsics the following headers are available:

```
mmintrin.h (SSE)
pmmintrin.h (SSE2)
immintrin.h (AVX)

immintrin.h (AVX)

x86intrin.h (all extensions)

for (int j=0; j<size; j+=16) {
   t0 = _mm_loadu_ps(data+j);
   t1 = _mm_loadu_ps(data+j+4);
   t2 = _mm_loadu_ps(data+j+4);
   t3 = _mm_loadu_ps(data+j+12);
   sum0 = _mm_add_ps(sum0, t0);
   sum1 = _mm_add_ps(sum1, t1);
   sum2 = _mm_add_ps(sum2, t2);
   sum3 = _mm_add_ps(sum3, t3);
}</pre>
```

#### **Rules for vectorizable loops**



- Inner loop
- Countable (loop length can be determined at loop entry)
- 3. Single entry and single exit
- 4. Straight line code (no conditionals)
- 5. No (unresolvable) read-after-write data dependencies
- 6. No function calls (exception intrinsic math functions)

#### **Better performance with:**

- 1. Simple inner loops with unit stride (contiguous data access)
- Minimize indirect addressing
- 3. Align data structures to SIMD width boundary
- 4. In C use the **restrict** keyword and/or **const** qualifiers and/or compiler options to rule out array/pointer aliasing



# Efficient parallel programming on ccNUMA nodes

Performance characteristics of ccNUMA nodes
First touch placement policy

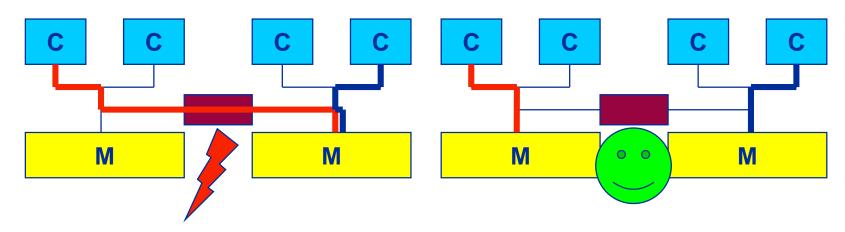
#### ccNUMA performance problems

"The other affinity" to care about



#### ccNUMA:

- Whole memory is transparently accessible by all processors
- but physically distributed
- with varying bandwidth and latency
- and potential contention (shared memory paths)
- How do we make sure that memory access is always as "local" and "distributed" as possible?



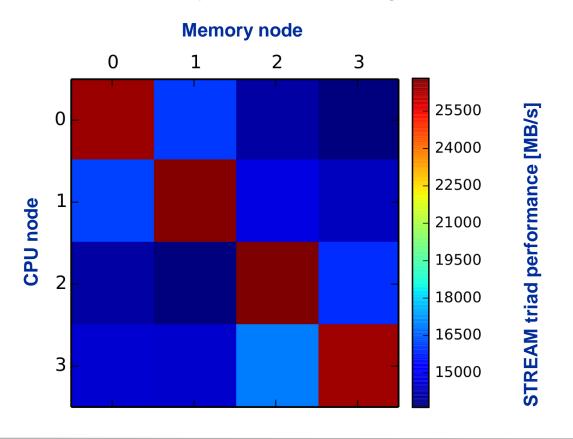
 Page placement is implemented in units of OS pages (often 4kB, possibly more)

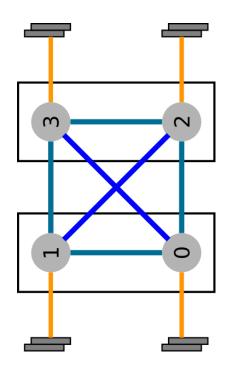
#### Intel Broadwell EP node





- ccNUMA map: Bandwidth penalties for remote access
  - Run 11 threads per ccNUMA domain (half chip)
  - Place memory in different domain → 4x4 combinations
  - STREAM copy benchmark using standard stores





#### numactl as a simple ccNUMA locality tool:

How do we enforce some locality of access?

numactl can influence the way a binary maps its memory pages:

#### Examples:

```
for m in `seq 0 3`; do
    for c in `seq 0 3`; do
    env OMP_NUM_THREADS=8 \
        numactl --membind=$m --cpunodebind=$c ./stream
    enddo
enddo
```

But what is the default without numactl?

### ccNUMA default memory locality



"Golden Rule" of ccNUMA:

# A memory page gets mapped into the local memory of the processor that first touches it!

- Except if there is not enough local memory available
- This might be a problem, see later
- Caveat: "touch" means "write", not "allocate"
- Example:

Memory not mapped here yet

It is sufficient to touch a single item to map the entire page

## Coding for ccNUMA data locality



#### Most simple case: explicit initialization

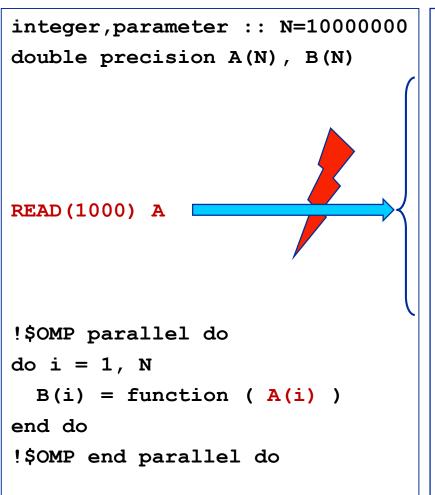
```
integer, parameter :: N=10000000
double precision A(N), B(N)
A=0.d0
!$OMP parallel do
do i = 1, N
 B(i) = function (A(i))
end do
!$OMP end parallel do
```

```
integer,parameter :: N=10000000
double precision A(N),B(N)
!$OMP parallel
!$OMP do schedule(static)
do i = 1, N
 A(i) = 0.d0
end do
!$OMP end do
!$OMP do schedule(static)
do i = 1, N
 B(i) = function (A(i))
end do
!$OMP end do
!$OMP end parallel
```

## Coding for ccNUMA data locality



 Sometimes initialization is not so obvious: I/O cannot be easily parallelized, so "localize" arrays before I/O



```
integer, parameter :: N=10000000
double precision A(N),B(N)
!$OMP parallel
!$OMP do schedule(static)
do i = 1, N
 A(i) = 0.d0
end do
!$OMP end do
!$OMP single
READ (1000) A
!$OMP end single
!$OMP do schedule(static)
do i = 1, N
 B(i) = function (A(i))
end do
!$OMP end do
!$OMP end parallel
```

# **Coding for Data Locality**



- Required condition: OpenMP loop schedule of initialization must be the same as in all computational loops
  - Only choice: static! Specify explicitly on all NUMA-sensitive loops, just to be sure...
  - Imposes some constraints on possible optimizations (e.g. load balancing)
  - If dynamic scheduling/tasking is unavoidable, more advanced methods may be in order
- How about global objects?
  - Better not use them
  - If communication vs. computation is favorable, might consider properly placed copies of global data
- C++: Arrays of objects and std::vector<> are by default initialized sequentially
  - STL allocators provide an elegant solution

#### **Coding for Data Locality:**

Placement of static arrays or arrays of objects



Don't forget that constructors tend to touch the data members of an object. Example:

```
class D {
  double d;
public:
   D(double _d=0.0) throw() : d(_d) {}
  inline D operator+(const D& o) throw() {
    return D(d+o.d);
  }
  inline D operator*(const D& o) throw() {
    return D(d*o.d);
  }
...
};
```

```
→ placement problem with
  D* array = new D[1000000];
```

#### **Coding for Data Locality:**

#### Parallel first touch for arrays of objects



Solution: Provide overloaded D::operator new[]

```
void* D::operator new[](size t n) {
  char *p = new char[n];  // allocate
                                                 parallel first
                                                 touch
  size t i,j;
#pragma omp parallel for private(j) schedule(...)
  for (i=0; i < n; i += size of (D))
    for(j=0; j<sizeof(D); ++j)</pre>
      p[i+j] = 0;
  return p;
void D::operator delete[](void* p) throw() {
  delete [] static cast<char*>p;
```

Placement of objects is then done automatically by the C++ runtime via "placement new"

#### **Coding for Data Locality:**

NUMA allocator for parallel first touch in std::vector<>



```
template <class T> class NUMA Allocator {
public:
  T* allocate(size_type numObjects, const void
               *localityHint=0) {
    size type ofs,len = numObjects * sizeof(T);
    void *m = malloc(len);
    char *p = static cast<char*>(m);
    int i,pages = len >> PAGE BITS;
#pragma omp parallel for schedule(static) private(ofs)
    for(i=0; i<pages; ++i) {</pre>
      ofs = static cast<size t>(i) << PAGE BITS;</pre>
      p[ofs]=0;
    return static cast<pointer>(m);
           Application:
           vector<double, NUMA Allocator<double> > x(10000000)
```

(c) RRZE 2016

## **Diagnosing Bad Locality**

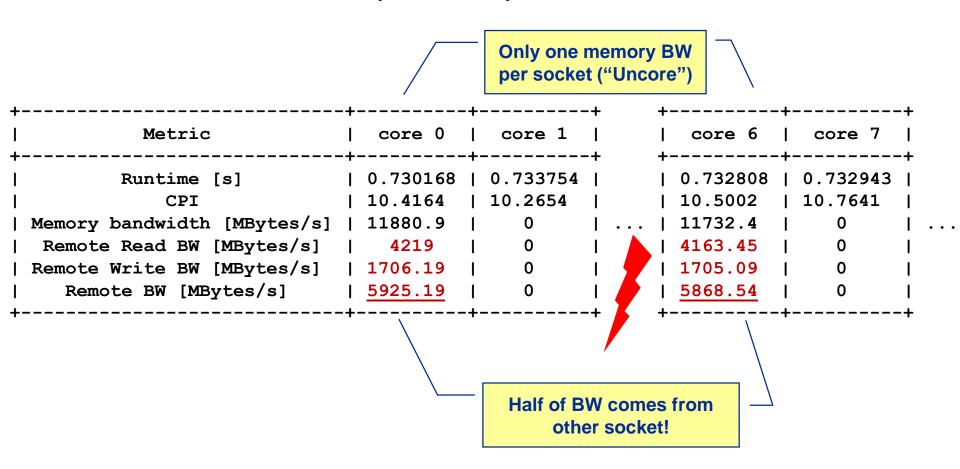
- If your code is cache bound, you might not notice any locality problems
- Otherwise, bad locality limits scalability at very low CPU numbers (whenever a node boundary is crossed)
  - If the code makes good use of the memory interface
- Running with numactl --interleave might give you a hint
  - See later
- Consider using performance counters
  - LIKWID-perfctr can be used to measure nonlocal memory accesses
  - Example for Intel Westmere dual-socket system (Core i7, hex-core):

```
env OMP_NUM_THREADS=12 likwid-perfctr -g MEM -C N:0-11 ./a.out
```

# Using performance counters for diagnosing bad ccNUMA access locality



Intel Westmere EP node (2x6 cores):



#### If all fails...

- Even if all placement rules have been carefully observed, you may still see nonlocal memory traffic. Reasons?
  - Program has erratic access patters → may still achieve some access parallelism (see later)
  - OS has filled memory with buffer cache data:

# numactl --hardware # idle node!

```
available: 2 nodes (0-1)
node 0 size: 2047 MB
node 0 free: 906 MB
node 1 size: 1935 MB
node 1 free: 1798 MB

top - 14:18:25 up 92 days, 6:07, 2 users, load average: 0.00, 0.02, 0.00
Mem: 4065564k total, 1149400k used, 2716164k free, 43388k buffers
Swap: 2104504k total, 2656k used, 2101848k free, 1038412k cached
```

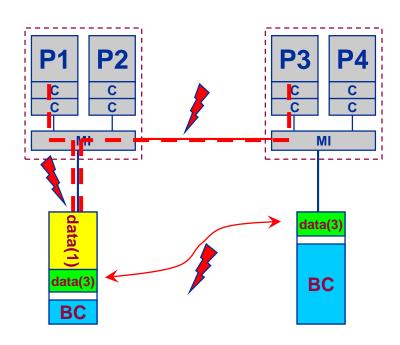
#### ccNUMA problems beyond first touch:

#### Buffer cache

# optional

#### OS uses part of main memory for disk buffer (FS) cache

- If FS cache fills part of memory, apps will probably allocate from foreign domains
- → non-local access!
- "sync" is not sufficient to drop buffer cache blocks



#### Remedies

- Drop FS cache pages after user job has run (admin's job)
  - seems to be automatic after aprun has finished on Crays
- User can run "sweeper" code that allocates and touches all physical memory before starting the real application
- numact1 tool or aprun can force local allocation (where applicable)
- Linux: There is no way to limit the buffer cache size in standard kernels

#### ccNUMA problems beyond first touch:

#### Buffer cache



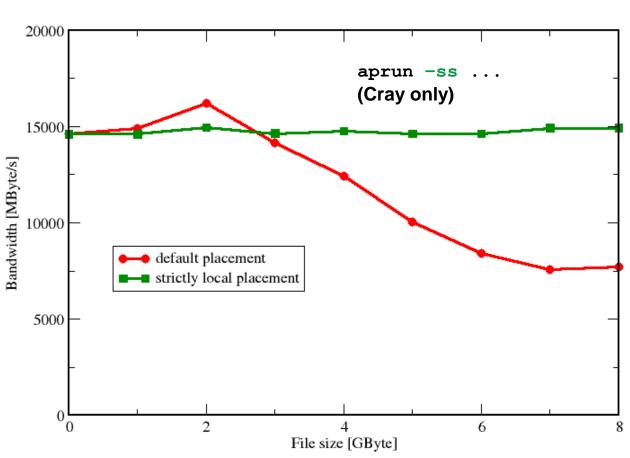
# Real-world example: ccNUMA and the Linux buffer cache

#### **Benchmark:**

- Write a file of some size from LD0 to disk
- Perform bandwidth benchmark using all cores in LD0 and maximum memory installed in LD0

Result: By default, Buffer cache is given priority over local page placement

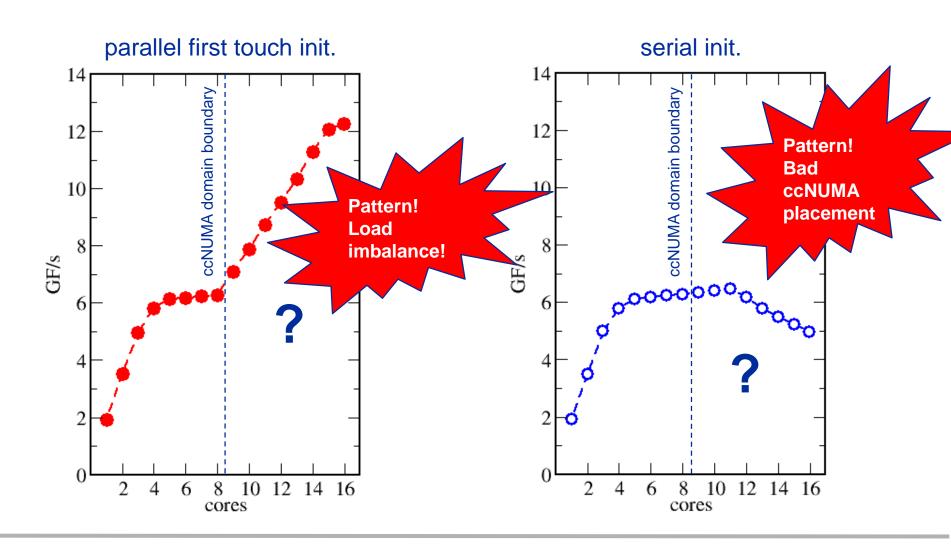
restrict to local domain if possible!



## Back to the spMVM code: a little riddle



#### DLR1 matrix on 2x 8-core Sandy Bridge node

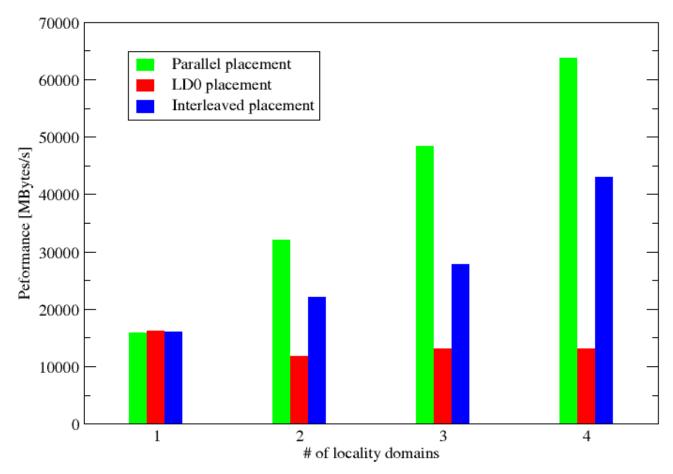


#### The curse and blessing of interleaved placement:

OpenMP STREAM on a Cray XE6 Interlagos node



- Parallel init: Correct parallel initialization
- LD0: Force data into LD0 via numactl -m 0
- Interleaved: numactl --interleave <LD range>

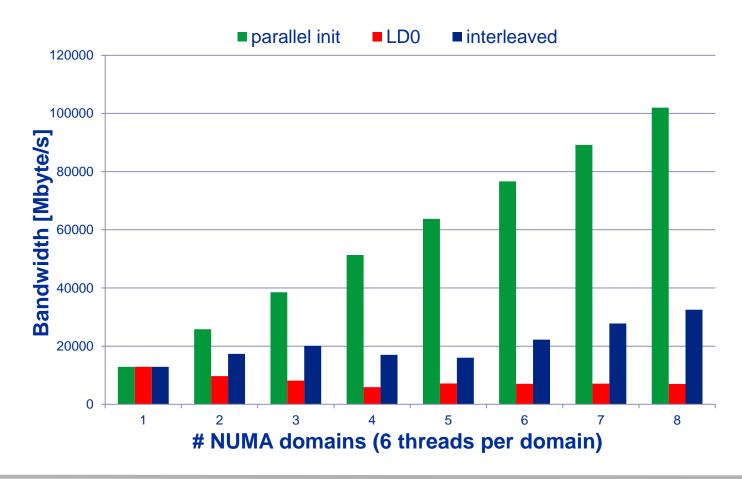


#### The curse and blessing of interleaved placement:

OpenMP STREAM triad on 4-socket (48 core) Magny Cours node



- Parallel init: Correct parallel initialization
- LD0: Force data into LD0 via numactl -m 0
- Interleaved: numactl --interleave <LD range>



#### **Summary on ccNUMA issues**



- Identify the problem
  - Is ccNUMA an issue in your code?
  - Simple test: run with numactl --interleave
- Apply first-touch placement
  - Look at initialization loops
  - Consider loop lengths and static scheduling
  - C++ and global/static objects may require special care
- If dynamic scheduling cannot be avoided
  - Distribute the data anyway, just do not use sequential placement!

OS file buffer cache may impact proper placement



# **OpenMP performance issues on multicore**

Barrier synchronization overhead Topology dependence

#### The OpenMP-parallel vector triad benchmark

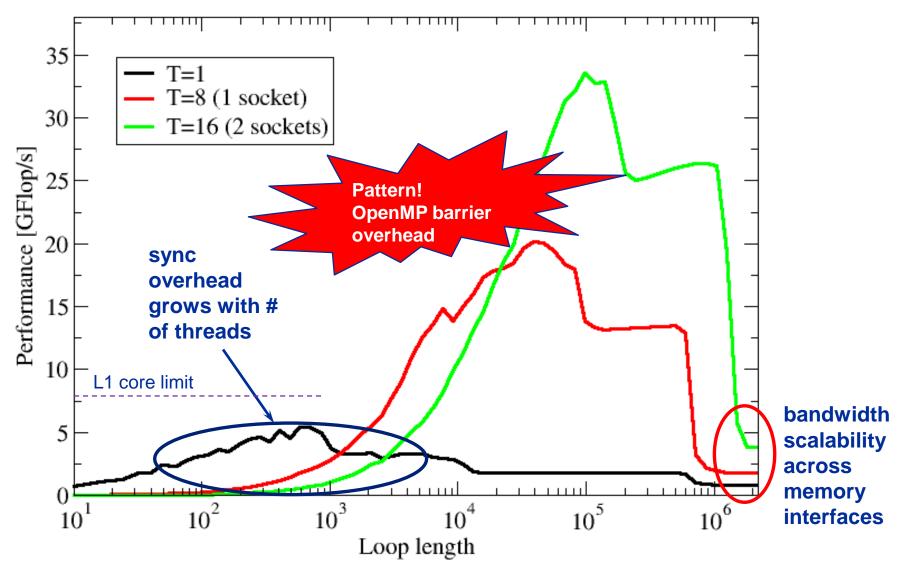


#### OpenMP work sharing in the benchmark loop

```
double precision, dimension(:), allocatable :: A,B,C,D
allocate (A(1:N), B(1:N), C(1:N), D(1:N))
A=1.d0; B=A; C=A; D=A
!$OMP PARALLEL private(i,j)
do j=1,NITER
!$OMP DO
  do i=1,N
    A(i) = B(i) + C(i) * D(i)
  enddo
                          Implicit barrier
!SOMP END DO
  if (.something.that.is.never.true.) then
    call dummy (A,B,C,D)
  endif
enddo
!$OMP END PARALLEL
```

#### **OpenMP vector triad on Sandy Bridge sockets (3 GHz)**





#### Welcome to the multi-/many-core era

Synchronization of threads may be expensive!



!\$OMP PARALLEL ...

!\$OMP BARRIER

!\$OMP DO

!\$OMP ENDDO

!\$OMP END PARALLEL

Threads are synchronized at **explicit** AND **implicit** barriers. These are a main source of overhead in OpenMP progams.

Determine costs via modified OpenMP

Microbenchmarks testcase (epcc)

#### On x86 systems there is no hardware support for synchronization!

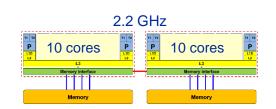
- Next slides: Test OpenMP Barrier performance...
- for different compilers
- and different topologies:
  - shared cache
  - shared socket
  - between sockets
- and different thread counts
  - 2 threads
  - full domain (chip, socket, node)

#### Thread synchronization overhead on IvyBridge-EP

Barrier overhead in CPU cycles



| 2 Threads    | Intel 16.0 | GCC 5.3.0                   |  |
|--------------|------------|-----------------------------|--|
| Shared L3    | 599        | 425                         |  |
| SMT threads  | 612        | 423                         |  |
| Other socket | 1486       | 1067                        |  |
|              |            | Strong topology dependence! |  |



| Full domain       | Intel 16.0 | GCC 5.3.0 |
|-------------------|------------|-----------|
| Socket (10 cores) | 1934       | 1301      |
| Node (20 cores)   | 4999       | 7783      |
| Node +SMT         | 5981       | 9897      |



- Strong dependence on compiler, CPU and system environment!
- OMP\_WAIT\_POLICY=ACTIVE can make a big difference

#### Thread synchronization overhead on Xeon Phi 7210 (64-core)

Barrier overhead in CPU cycles (Intel C compiler 16.03)



2 threads on distinct cores: 730

| 0         | SMT1 | SMT2 | SMT3 | SMT4  |
|-----------|------|------|------|-------|
| One core  | n/a  | 963  | 1580 | 2240  |
| Full chip | 5720 | 8100 | 9900 | 11400 |

Still the pain may be much larger, as more work can be done in one cycle on Phi compared to a full lvy Bridge node

3.2x cores (20 vs 64) on Phi4x more operations per cycle per core on Phi

 $\rightarrow$  4 · 3.2 = 12.8x more work done on Xeon Phi per cycle

1.9x more barrier penalty (cycles) on Phi (11400 vs. 6000)

→ One barrier causes 1.9 · 12.8 ≈ 24x more pain ③.



# Pattern-driven Performance Engineering

Basics of Benchmarking Performance Patterns Signatures

## **Basics of optimization**



- Define relevant test cases
- 2. Establish a sensible performance metric
- Acquire a runtime profile (sequential)
- 4. Identify hot kernels (Hopefully there are any!)
- 5. Carry out optimization process for each kernel



#### **Motivation:**

- Understand observed performance
- Learn about code characteristics and machine capabilities
- Deliberately decide on optimizations

## **Best practices for benchmarking**



#### Preparation

- Reliable timing (minimum time which can be measured?)
- Document code generation (flags, compiler version)
- Get access to an exclusive system
- System state (clock speed, turbo mode, memory, caches)
- Consider to automate runs with a script (shell, python, perl)

#### Doing

- Affinity control
- Check: Is the result reasonable?
- Is result deterministic and reproducible?
- Statistics: Mean, Best ?
- Basic variants: Thread count, affinity, working set size

## Thinking in bottlenecks



- A bottleneck is a performance limiting setting
- Microarchitectures expose numerous bottlenecks

#### **Observation 1:**

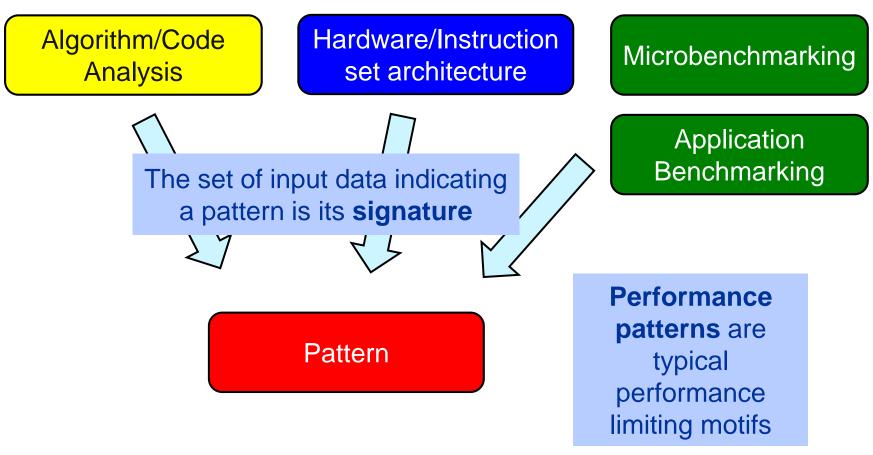
Most applications face a single bottleneck at a time!

#### **Observation 2:**

There is a limited number of relevant bottlenecks!

# **Performance Engineering Process: Analysis**





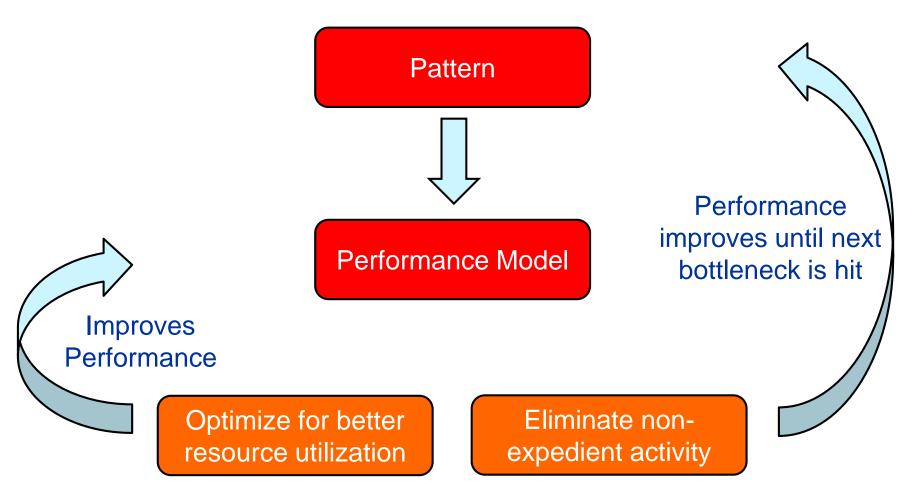
Step 1 Analysis: Understanding observed performance

# **Performance Engineering Process: Modeling** Qualitative view Pattern Wrong pattern Quantitative view **Performance Model Validation** Traces/HW metrics

Step 2 Formulate Model: Validate pattern and get quantitative insight

# **Performance Engineering Process: Optimization**





Step 3 Optimization: Improve utilization of available resources

# **Performance pattern classification**



1. Maximum resource utilization (computing at a bottleneck)

Hazards (something "goes wrong")

3. Work related (too much work or too inefficiently done)

# Patterns (I): Bottlenecks & hazards

| Pattern                               |                      | Performance behavior   | Metric signature, LIKWID performance group(s)  |
|---------------------------------------|----------------------|--|--|
| Bandwidth saturation Jacobi           |                      | Saturating speedup across cores sharing a data path                                      | Bandwidth meets BW of suitable streaming benchmark (MEM, L3)   |
| ALU saturation In-L1 sum optimal code |                      | Throughput at design limit(s)  | Good (low) CPI, integral ratio of cycles to specific instruction count(s) (FLOPS_*, DATA, CPI)                             |
| Inefficient<br>data<br>access         | Excess data volume   | Simple bandwidth performance model much too optimistic                                   | Low BW utilization / Low cache hit ratio, frequent CL evicts or replacements (CACHE, DATA, MEM)                            |
|                                       | Latency-bound access |  |  |
| Micro-architectural anomalies         |                      | Large discrepancy from simple performance model based on LD/ST and arithmetic throughput | Relevant events are very hardware-specific, e.g., memory aliasing stalls, conflict misses, unaligned LD/ST, requeue events |

# Patterns (II): Hazards

| Pattern   | Performance behavior  | Metric signature, LIKWID performance group(s)  |
|---|---|--|
| False sharing of cache lines                          | Large discrepancy from performance model in parallel case, bad scalability                  | Frequent (remote) CL evicts (CACHE)  |
| No parallel initialization  Bad ccNUMA page placement | Bad or no scaling across NUMA domains, performance improves with interleaved page placement | Unbalanced bandwidth on memory interfaces / High remote traffic (MEM)                                  |
| In-L1 sum w/o unrolling  Pipelining issues            | In-core throughput far from design limit, performance insensitive to data set size          | (Large) integral ratio of cycles to specific instruction count(s), bad (high) CPI (FLOPS_*, DATA, CPI) |
| Control flow issues                                   | See above   | High branch rate and branch miss ratio (BRANCH)  |

# Patterns (III): Work-related

| Pattern  |                          | Performance behavior   | Metric signature, LIKWID performance group(s)   |
|--|--------------------------|--|---|
| Load imbalance / serial fraction SpMVM scaling |                          | Saturating/sub-linear speedup  | Different amount of "work" on the cores (FLOPS_*); note that instruction count is not reliable!                     |
| Synchronization overhead                       |                          | Speedup going down as more cores are added / No speedup with small problem sizes / Cores busy but low FP performance | Large non-FP instruction count (growing with number of cores used) / Low CPI (FLOPS_*, CPI)                         |
| Instruction overhead                           |                          | Low application performance, good scaling across cores, performance insensitive to problem size                      | Low CPI near theoretical limit / Large non-FP instruction count (constant vs. number of cores) (FLOPS_*, DATA, CPI) |
| Code<br>composition                            | Expensive instructions   | Similar to instruction overhead  C/C++ aliasing problem  | Many cycles per instruction (CPI) if the problem is large-latency arithmetic  |
|  | Ineffective instructions |  | Scalar instructions dominating in data-parallel loops (FLOPS_*, CPI)  |

## **Patterns conclusion**



- Pattern signature = performance behavior + hardware metrics
- Patterns are applies hotspot (loop) by hotspot
- Patterns map to typical execution bottlenecks
- Patterns are extremely helpful in classifying performance issues
  - The first pattern is always a hypothesis
  - Validation by tanking data (more performance behavior, HW metrics)
  - Refinement or change of pattern
- Performance models are crucial for most patterns
  - Model follows from pattern

## **Tutorial conclusion**



## Multicore architecture == multiple complexities

- Affinity matters → pinning/binding is essential
- Bandwidth bottlenecks → inefficiency is often made on the chip level
- Topology dependence of performance features → know your hardware!

## Put cores to good use

- Bandwidth bottlenecks → surplus cores → functional parallelism!?
- Shared caches → fast communication/synchronization → better implementations/algorithms?

## Simple modeling techniques and patterns help us

- ... understand the limits of our code on the given hardware
- ... identify optimization opportunities
- ... learn more, especially when they do not work!

## Simple tools get you 95% of the way

e.g., with the LIKWID tool suite





Moritz Kreutzer
Markus Wittmann
Thomas Zeiser
Michael Meier
Holger Stengel
Thomas Röhl
Faisal Shahzad
Salah Saleh







# THANK YOU.

## **Presenter Biographies**



**Georg Hager** holds a PhD in computational physics from the University of Greifswald. He is a senior research scientist in the HPC group at Erlangen Regional Computing Center (RRZE). Recent research includes architecture-specific optimization for current microprocessors, performance modeling on processor and system levels, and the efficient use of hybrid parallel systems. His textbook "Introduction to High Performance Computing for Scientists and Engineers" is required or recommended reading in many HPC-related courses around the world. See his blog at <a href="http://blogs.fau.de/hager">http://blogs.fau.de/hager</a> for current activities, publications, and talks.



Jan Eitzinger (formerly Treibig) holds a PhD in Computer Science from the University of Erlangen. He is now a postdoctoral researcher in the HPC Services group at Erlangen Regional Computing Center (RRZE). His current research revolves around architecture-specific and low-level optimization for current processor architectures, performance modeling on processor and system levels, and programming tools. He is the developer of LIKWID, a collection of lightweight performance tools. In his daily work he is involved in all aspects of user support in High Performance Computing: training, code parallelization, profiling and optimization, and the evaluation of novel computer architectures.



**Gerhard Wellein** holds a PhD in solid state physics from the University of Bayreuth and is a professor at the Department for Computer Science at the University of Erlangen. He leads the HPC group at Erlangen Regional Computing Center (RRZE) and has more than ten years of experience in teaching HPC techniques to students and scientists from computational science and engineering programs. His research interests include solving large sparse eigenvalue problems, novel parallelization approaches, performance modeling, and architecture-specific optimization.



#### **Abstract**



- SC16 tutorial: Node-Level Performance Engineering
- Presenter(s): Georg Hager, (Jan Eitzinger), Gerhard Wellein

#### ABSTRACT:

The advent of multi- and manycore chips has led to a further opening of the gap between peak and application performance for many scientific codes. This trend is accelerating as we move from petascale to exascale. Paradoxically, bad node-level performance helps to "efficiently" scale to massive parallelism, but at the price of increased overall time to solution. If the user cares about time to solution on any scale, optimal performance on the node level is often the key factor. We convey the architectural features of current processor chips, multiprocessor nodes, and accelerators, as far as they are relevant for the practitioner. Peculiarities like SIMD vectorization, shared vs. separate caches, bandwidth bottlenecks, and ccNUMA characteristics are introduced, and the influence of system topology and affinity on the performance of typical parallel programming constructs is demonstrated. Performance engineering and performance patterns are suggested as powerful tools that help the user understand the bottlenecks at hand and to assess the impact of possible code optimizations. A cornerstone of these concepts is the roofline model, which is described in detail, including useful case studies, limits of its applicability, and possible refinements.



#### Book:

G. Hager and G. Wellein: <u>Introduction to High Performance Computing for Scientists and Engineers</u>.
 CRC Computational Science Series, 2010. ISBN 978-1439811924
 <a href="http://www.hpc.rrze.uni-erlangen.de/HPC4SE/">http://www.hpc.rrze.uni-erlangen.de/HPC4SE/</a>

#### Papers:

- J. Hofmann, D. Fey, M. Riedmann, J. Eitzinger, G. Hager, and G. Wellein: Performance analysis of the Kahan-enhanced scalar product on current multi- and manycore processors. Accepted for publication in Concurrency & Computation: Practice & Experience. Preprint: <u>arXiv:1604.01890</u>
- M. Röhrig-Zöllner, J. Thies, M. Kreutzer, A. Alvermann, A. Pieper, A. Basermann, G. Hager, G. Wellein, and H. Fehske: Increasing the performance of the Jacobi-Davidson method by blocking. SIAM Journal on Scientific Computing, 37(6), C697–C722 (2015). DOI: 10.1137/140976017, Preprint: <a href="http://elib.dlr.de/89980/">http://elib.dlr.de/89980/</a>
- T. M. Malas, G. Hager, H. Ltaief, H. Stengel, G. Wellein, and D. E. Keyes: Multicore-optimized wavefront diamond blocking for optimizing stencil updates. SIAM Journal on Scientific Computing 37(4), C439-C464 (2015). DOI: 10.1137/140991133, Preprint: arXiv:1410.3060
- J. Hammer, G. Hager, J. Eitzinger, and G. Wellein: Automatic Loop Kernel Analysis and Performance Modeling With Kerncraft. Proc. <u>PMBS15</u>, the 6th International Workshop on Performance Modeling, Benchmarking and Simulation of High Performance Computer Systems, in conjunction with ACM/IEEE Supercomputing 2015 (<u>SC15</u>), November 16, 2015, Austin, TX. <u>DOI: 10.1145/2832087.2832092</u>, Preprint: <u>arXiv:1509.03778</u>



#### Papers continued:

- M. Kreutzer, G. Hager, G. Wellein, A. Pieper, A. Alvermann, and H. Fehske: Performance Engineering of the Kernel Polynomial Method on Large-Scale CPU-GPU Systems. Proc. <u>IPDPS15</u>. DOI: 10.1109/IPDPS.2015.76, Preprint: arXiv:1410.5242
- M. Wittmann, G. Hager, T. Zeiser, J. Treibig, and G. Wellein: Chip-level and multi-node analysis of energy-optimized lattice-Boltzmann CFD simulations. Concurrency and Computation: Practice and Experience (2015). <u>DOI: 10.1002/cpe.3489</u> Preprint: <u>arXiv:1304.7664</u>
- H. Stengel, J. Treibig, G. Hager, and G. Wellein: Quantifying performance bottlenecks of stencil computations using the Execution-Cache-Memory model. Proc. <u>ICS15</u>, <u>DOI: 10.1145/2751205.2751240</u>, Preprint: <u>arXiv:1410.5010</u>
- M. Kreutzer, G. Hager, G. Wellein, H. Fehske, and A. R. Bishop: A unified sparse matrix data format for modern processors with wide SIMD units. SIAM Journal on Scientific Computing 36(5), C401–C423 (2014). DOI: 10.1137/130930352, Preprint: arXiv:1307.6209
- G. Hager, J. Treibig, J. Habich and G. Wellein: Exploring performance and power properties of modern multicore chips via simple machine models. Computation and Concurrency: Practice and Experience (2013). DOI: 10.1002/cpe.3180, Preprint: arXiv:1208.2908
- J. Treibig, G. Hager and G. Wellein: Performance patterns and hardware metrics on modern multicore processors: Best practices for performance engineering. Workshop on Productivity and Performance (PROPER 2012) at Euro-Par 2012, August 28, 2012, Rhodes Island, Greece. <u>DOI: 10.1007/978-3-642-36949-0\_50</u>. Preprint: <u>arXiv:1206.3738</u>



#### Papers continued:

- M. Wittmann, T. Zeiser, G. Hager, and G. Wellein: Comparison of Different Propagation Steps for Lattice Boltzmann Methods. Computers & Mathematics with Applications (Proc. ICMMES 2011). Available online, <u>DOI: 10.1016/j.camwa.2012.05.002</u>. Preprint:<u>arXiv:1111.0922</u>
- M. Kreutzer, G. Hager, G. Wellein, H. Fehske, A. Basermann and A. R. Bishop: Sparse Matrix-vector Multiplication on GPGPU Clusters: A New Storage Format and a Scalable Implementation. Workshop on Large-Scale Parallel Processing 2012 (LSPP12), DOI: 10.1109/IPDPSW.2012.211
- J. Treibig, G. Hager, H. Hofmann, J. Hornegger and G. Wellein: Pushing the limits for medical image reconstruction on recent standard multicore processors. International Journal of High Performance Computing Applications, (published online before print).
   DOI: 10.1177/1094342012442424
- G. Wellein, G. Hager, T. Zeiser, M. Wittmann and H. Fehske: Efficient temporal blocking for stencil computations by multicore-aware wavefront parallelization. Proc. COMPSAC 2009. <u>DOI:</u> 10.1109/COMPSAC.2009.82
- M. Wittmann, G. Hager, J. Treibig and G. Wellein: Leveraging shared caches for parallel temporal blocking of stencil codes on multicore processors and clusters. Parallel Processing Letters 20 (4), 359-376 (2010).
  DOI: 10.1142/S0129626410000296. Preprint: arXiv:1006.3148
- J. Treibig, G. Hager and G. Wellein: LIKWID: A lightweight performance-oriented tool suite for x86 multicore environments. Proc. <u>PSTI2010</u>, the First International Workshop on Parallel Software Tools and Tool Infrastructures, San Diego CA, September 13, 2010. <u>DOI: 10.1109/ICPPW.2010.38</u>. Preprint: <u>arXiv:1004.4431</u>



#### Papers continued:

- G. Schubert, H. Fehske, G. Hager, and G. Wellein: Hybrid-parallel sparse matrix-vector multiplication with explicit communication overlap on current multicore-based systems. Parallel Processing Letters 21(3), 339-358 (2011).
  - DOI: 10.1142/S0129626411000254
- J. Treibig, G. Wellein and G. Hager: Efficient multicore-aware parallelization strategies for iterative stencil computations. Journal of Computational Science 2 (2), 130-137 (2011). <u>DOI</u> 10.1016/j.jocs.2011.01.010
- K. Iglberger, G. Hager, J. Treibig, and U. Rüde: <u>Expression Templates Revisited: A Performance Analysis of Current ET Methodologies</u>. SIAM Journal on Scientific Computing **34**(2), C42-C69 (2012). <u>DOI: 10.1137/110830125</u>, Preprint: <u>arXiv:1104.1729</u>
- K. Iglberger, G. Hager, J. Treibig, and U. Rüde: High Performance Smart Expression Template Math Libraries. 2nd International Workshop on New Algorithms and Programming Models for the Manycore Era (<u>APMM 2012</u>) at <u>HPCS 2012</u>, July 2-6, 2012, Madrid, Spain. <u>DOI: 10.1109/HPCSim.2012.6266939</u>
- J. Habich, T. Zeiser, G. Hager and G. Wellein: Performance analysis and optimization strategies for a D3Q19 Lattice Boltzmann Kernel on nVIDIA GPUs using CUDA. Advances in Engineering Software and Computers & Structures 42 (5), 266–272 (2011). DOI: 10.1016/j.advengsoft.2010.10.007
- J. Treibig, G. Hager and G. Wellein: Multicore architectures: Complexities of performance prediction for Bandwidth-Limited Loop Kernels on Multi-Core Architectures.
   DOI: 10.1007/978-3-642-13872-0 1, Preprint: arXiv:0910.4865.



#### Papers continued:

- G. Hager, G. Jost, and R. Rabenseifner: Communication Characteristics and Hybrid MPI/OpenMP Parallel Programming on Clusters of Multi-core SMP Nodes. In: Proceedings of the Cray Users Group Conference 2009 (CUG 2009), Atlanta, GA, USA, May 4-7, 2009. <u>PDF</u>
- R. Rabenseifner and G. Wellein: Communication and Optimization Aspects of Parallel Programming Models on Hybrid Architectures. International Journal of High Performance Computing Applications 17, 49-62, February 2003.

DOI:10.1177/1094342003017001005