For final slides and example code see: https://tiny.cc/NLPE-SC19



Node-Level Performance Engineering

Georg Hager, Jan Eitzinger, Gerhard Wellein Erlangen Regional Computing Center (RRZE) and Department of Computer Science

University of Erlangen-Nuremberg

SC19 full-day tutorial November 17, 2019 Denver, CO

slide updated



Agenda



Preliminaries		08:30
Introduction to m	ulticore architecture	
Threads, cores, S	SIMD, caches, chips, sockets, ccNUMA	
Multicore tools (p	oart I)	10:00
Microbenchmarki	ing for architectural exploration	10:30
Streaming bench	marks	
Hardware bottlen	iecks	
Node-level perfor	mance modeling (part I)	
The Roofline Mod	del	12:00
Lunch break		
Multicore tools (p	oart II)	13:30
Node-level perfor	mance modeling (part II)	
Case studies: Jac	cobi solver, sparse MVM, tall & skinny MM	15:00
Optimal resource	utilization	15:30
SIMD parallelism		
ccNUMA		
OpenMP synchro	onization and multicores	17:00
(c) RRZE 2019	Node-Level Performance Engineering	



Prelude: Scalability 4 the win!





(c) RRZE 2019

Node-Level Performance Engineering

Scalability Myth: Code scalability is the key issue







- 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



Introduction: Modern node architecture

A glance at basic core features: pipelining, superscalarity, SMT, SIMD Caches and data transfers through the memory hierarchy Accelerators Bottlenecks & hardware-software interaction

Multi-core today: Intel Xeon Skylake SP (2017)

- LL5JJ
- Xeon "Skylake SP" (Platinum/Gold/Silver/Bronze): Up to 28 cores running at 2+ GHz (+ "Turbo Mode": 3.8+ GHz)
 - Mesh interconnect
 - Reincarnated as "Cascade Lake" in 2018
- Simultaneous Multithreading
 - \rightarrow reports as 56-way chip
- 8 Billion Transistors / 14 nm
- Die size: ~500 mm²





2015: Broadwell architecture

- Ring instead of mesh interconnect
- Cluster on Die (analogous to SNC)
- Up to 24 cores

Optional: "Sub-

General-purpose cache based microprocessor core





Stored-program computer

- 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



LL5J

Idea:

- Split complex instruction into several simple / fast steps (stages)
- Each step takes the same time, e.g., one 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 sufficient # of independent instructions required
 - Requires complex instruction scheduling by compiler/hardware
 - software-pipelining / out-of-order execution
- Pipelining is widely used in modern computer architectures





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

Superscalar Processors – Instruction Level Parallelism



- Multiple units enable use of Instruction Level Parallelism (ILP): Instruction stream is "parallelized" on the fly
- Instructions from different loop iterations retired at the same time



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



SMT principle (2-way example):





- x86 SIMD instruction sets:
 - SSE: register width = 128 Bit \rightarrow 2 double precision floating point operands
 - AVX: register width = 256 Bit \rightarrow 4 double precision floating point operands
 - AXV512: you get it.
- Adding two registers holding double precision floating point operands



(c) RRZE 2019

_ _ _ _ .



T1 T2 T1 T2 T1 T2 P P P P P L1D L1D L1D L1D L2 L2 L2 L2 L2 L2 Memory I	T1 T2 T1 T2 T1 T P P P L1D L1D L1D L2 L2 L2		Maximur P _{core} = / Super-	n floatin = n_{sup}^{FP}	g poin per	t (FP) p n _{FMA} / SIMD	n _{SIM}	ance:	
Typical representatives	n ^{FP} [inst./cy]	n _{FMA}	n _{SIMD} [ops/inst.]	y iac	C	ode	f [Gcy/s]	P _{core} [GF/s]	
Nehalem	2	1	2	Q1/2009	X	5570	2.93	11.7	
Westmere	2	1	2	Q1/2010	X	5650	2.66	10.6	
Sandy Bridge	2	1	4	Q1/2012	E5-	2680	2.7	21.6	
Ivy Bridge	2	1	4	Q3/2013	E5-2	660 v2	2.2	17.6	
Haswell	2	2	4	Q3/2014	E5-2	695 v3	2.3	36.8	
Broadwell	2	2	4	Q1/2016	E5-2	699 v4	2.2	35.2	
Skylake	2	2	8	Q3/2017	Golo	6148	2.4	76.8	
 AMD Zen	2	2	2	Q1/2017	Epy	c 7451	2.3	18.4	_
IBM POWER8	2	2	2	Q2/2014	S8	22LC	2.93	23.4 1	8



s = 0.0do i = 1,N s = s + a(i)

enddo

...In single precision on an AVXcapable core (ADD latency = 3 cy)

How fast can this loop possibly run with data in the L1 cache?

- Loop-carried dependency on summation variable
- Execution stalls at every ADD until previous ADD is complete
- \rightarrow No pipelining?
- \rightarrow No SIMD?



Plain scalar code, no SIMD

do $i = 1, N$ s = s + a(i) enddo
LOAD r1.0 \leftarrow 0
i 🗲 1
loop:
LOAD r2.0 \leftarrow a(i)







 \rightarrow 1/24 of ADD peak

Applicable peak for the sum reduction (II)



s3





SIMD-vectorization (8-way MVE) x pipelining (3-way MVE)

```
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
```

```
do i = 1,N,24
s10=s10+a(i+0); s20=s20+a(i+8); s30=s30+a(i+16)
s11=s11+a(i+1); s21=s21+a(i+9); s31=s31+a(i+17)
s12=s12+a(i+2); s22=s22+a(i+10); s32=s32+a(i+18)
s13=s13+a(i+3); s23=s23+a(i+11); s33=s33+a(i+19)
s14=s14+a(i+4); s24=s24+a(i+12); s34=s34+a(i+20)
s15=s15+a(i+5); s25=s25+a(i+13); s35=s35+a(i+21)
s16=s16+a(i+6); s26=s26+a(i+14); s36=s36+a(i+22)
s17=s17+a(i+7); s27=s27+a(i+15); s37=s37+a(i+23)
```

```
enddo
```

```
s = s + s10 + s11 + ... + s37
```

	s10	s20	s30
loop: LOAD $[r4.0,,r4.7] \leftarrow [a(i),,a(i+7)] \# SIMD LO$	AD s11	s21	s31
LOAD $[r5.0,, r5.7] \leftarrow [a(i+8),, a(i+15)] \# SIMD$	s12	s22	s32
LOAD $[10.0,, 10.7] \subset [a(1+10),, a(1+25)] # SIMD$	s13	s23	s33
ADD r1 \leftarrow r1 + r4 # SIMD ADD ADD r2 \leftarrow r2 + r5 # SIMD ADD	OOV s14	s24	s34
ADD r3 \leftarrow r3 + r6 # SIMD ADD	↑ s15	s25	s35
i+=24 →? loop	s16	s26	s36
result C r1.0+r1.1++r3.6+r3.7	s17	s27	s37



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(:)



Putting the cores & caches together

AMD Epyc 7451 24-Core Processor («Naples»)





24 cores per socket

- 4 chips w/ 6 cores each ("Zeppelin" die)
 - 3 cores share 8MB L3 ("Core Complex", "CCX")
- DDR4-2666 memory interface with 2 channels per chip
 - MemBW per node:
 - 16 ch x 8 byte x 2.666 GHz = 341 GB/s
- Two-way SMT
 - Two 256-bit (actually 4 128-bit) SIMD FP units
 AVX2, 8 flops/cycle
- 32 KiB L1 data cache per core
- 512 KiB L2 cache per core
- 2 x 8 MiB L3 cache per chip
 - 64 MiB L3 cache per socket
- ccNUMA memory architecture
- Infinity fabric between CCX's and between chips



Interlude: A glance at current accelerator technology

Nvidia "Volta" V100 vs. Intel Skylake "Platinum"

Nvidia V100 "Volta" specs

Architecture

- 21.1 B Transistors
- ~ 1.4 GHz clock speed
- ~ 80 "SM" units -
 - 64 SP "cores" each (FMA)
 - 32 DP "cores" each (FMA)
 - 8 "Tensor Cores" each
 - 2:1 SP:DP performance
- ~7 TFlop/s DP peak
- 6 MiB L2 Cache
- 4096-bit HBM2
- MemBW ~ 900 GB/s (theoretical)
- MemBW ~ 830 GB/s (measured)



© Nvidia



Node-Level Performance Engineering

Trading single thread performance for parallelism: GPGPUs vs. CPUs



GPU vs. Cf light spe	⊃U ed estimate	Control	ALU ALU		
(per dev	ice)	Cache			
/lemBW	~ 8-10x	DRAM			DRAM
Peak	~ 5-10x		CPU	GPU	

	2x Intel Xeon Platinum 8160	NVidia Tesla V100 "Volta"
Cores@Clock	2 x 24 @ ≥2.1 GHz	80 SMs @ ~1.4 GHz
SP Performance/core	≥ 134 GFlop/s	~179 GFlop/s
Threads@STREAM	~20	> 20000
SP peak	≥ 6 TFlop/s	~14 TFlop/s
Stream BW (meas.)	2 x 105 GB/s	830 GB/s
Transistors / TDP	~2x8 Billion / 2x150 W	21 Billion/250 W



Node topology and programming models

Parallelism in a modern compute node



 Parallel and shared resources (potential bottlenecks!) within a sharedmemory 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
- Memory bus per socket 7
- Intersocket link 8
- PCIe bus(es) 9
- Other I/O resources 10

How does your application react to all of those details?

Parallel programming models: *Pure MPI*





Parallel programming models:

Pure threading on the node







- Node-level hardware parallelism takes many forms
 - Sockets/devices CPU: 1-8, GPGPU: 1-6
 - Cores moderate (CPU: 4-64) to massive (GPGPU: 10's-100'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
 - Slowly improving, though (OpenMP 4.0, MPI 3.0)
- Apart from overheads, performance features are largely independent of the programming model



Multicore Performance and Tools

Part 1: Topology and affinity



- 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, Intel Amplifier,...
- Performance Profilers
 Intel VtuneTM, likwid-perfctr, PAPI based tools, HPC Toolkit, Linux perf, ...
- Microbenchmarking STREAM, likwid-bench, Imbench



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

Output of likwid-topology -g

on one node of Intel Haswell-EP

Г	=
•	

														-																
CPU name: CPU type: CPU stepping:	Intel (R) Intel X 2) Xeon(R eon Hasw) CPU E ell EN/	5-269 'EP/EX	5 v3 @ proce	2.3 ssoi	30GHz	****	***	****	****	****	****	*																
Hardware Threa	d Topolog	JY *******	******	****	*****	****	****	****	***	****	****	****	****	*																
Sockets:			2																											
Cores per sock	et:		14																											
Threads per co	re:		2																											
HWThread	Thread		Core		So	cket			Ava	ilab	le			-																
0	0		0		0				*																					
1	0		1		0				*																					
•••																														
43	1		1		1				*																					
44	1		2		1				*																					
Socket 0:	(0 28 3	1 29 2 3	0 3 31	4 32	5 33 6	5 34	7 35	8 3	59	37 :	10 3	8 11	39	12	40 1	13 4	1)	4 0	7 66											
		15 45 1		45 1				40 2.						- 52	2 25	55	20 2)4 Z	/ 55	, ,				A	l pl	าyร	sica	al		
*************	*******	******	******	****	*****	****	****	****	***	****	****	****	****	*										pro	ce	SSC	or I	Ds		
**************************************	*******	******	******	*****	*****	****	****	****	***	****	****	****	****	*																
Level:			1																											
Size:			32 kB																											
Cache groups:	(028)) (129) (2	30)	(3 31	.)	(43	2)	(5	33) (6 34)	(7	35)	(8 36	5)	(9	37)	(1	.0 38	3)	(11	39) (12 4	0)	(134	L
) (14 42) (15 43)	(1644)) (17	45)	(184	6)	(19	47)	(20 4	18)	(2	1 49))	(22	2 50) (23	51) (24 5	2)	(2	5 53) (26	54)	(2'	755)	
Level:			2																											
Size:			256 kB																											
Cache groups:	(028) (129) (2	30)	(3 31	.)	(43	2)	(5	33) (6 34)	(7	35)	(8 36	5)	(9	37)	(1	.0 31	3)	(11	39) (12 4	0)	(134	1
) (14 42) (15 43)	(1644)) (17	45)	(184	6)	(19	47)	(20 4	18)	(2	1 49)	(22	2 50) (23	51) (24 5	2)	(2	5 53) (26	54)	(2'	755)	
Level:			3																											
Size:			17 MB																											
Cache groups: (21 49 22 50 2	(0 28 3 23 51 24	1 29 2 3 52 25 5	0 3 31 3 26 54	4 32 27 5	5336 5)	5 34) (7 35	8 3	369	37	10 3	8 11	39	9 12	40	13 4	11)	(1	.4 42	2 15	43 :	L6 4	4 17	45	18 4	6 19	47	20 48)
														-																

Output of likwid-topology continued



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



Cluster on die mode

***	***************************************													and SMT enabled!																
Gra	aphical Topology ************************************																													
Soc	ket 0):																												
+																														+
	0 2	28	1	L 29	+ +- 	2 30	+ -	++ 3 31	· +	4 32	+ +	5 33	+ - 	634	+ +	7 35	1	8 36	+	+	9 37	+	10 38	+ +	11 39	+ +	12 40	1	13 41	+
		+	+		+ +-		+ -	++	+		+ +		+ -	+	+ +		+	+	+	+	+	+		+ +		+ +	+	+		+ 1
	32k	:B	1 3	32kB		32kB	I I	32kB	i	32kB	i i	32kB		32kB	i	32kB	I	32kB		Ť	32kB	ī	32kB		32kB		32kB	ī	32kB	T I
14		·+	+		+ +- + +-		+ -+ -+ -+ -+ -+ -+ -+ -+ -+ -+ -+ -+ -+	++	• +		+ + + +		+ ·	+	⊦ + ⊦ +		+ -	+	+	+	+	+++++++++++++++++++++++++++++++++++++++		+ + + +		+ + + +	+	· +		·+ -+
i	256	kB	1 2	256kB	i i	256kB	i I	256kB	I	256kB	i i	256kB		256kB	I	256kB	i.	256k	в	Ì	256kB	I	256kB		256kB		256kB	Ì	256kB	
11	++ ++ ++ ++ ++ ++ ++ ++ ++ ++ + ++ ++ ++ ++ ++ +												-+ ++ ++ ++ ++ ++ ++ ++ +														·+ ·+			
ii	17MB												17MB															i i		
1 4	++ +-																											•+		
Soc	ket 1	.:																												
1 4		+	+		+ +-		+ -	++	+		+ +	·	+ -	+	+ +		+	+	+	+	+	+		+ +		+ +	+	+		-+
	14 4	2	15	5 43	I I	16 44		17 45	1	18 46		19 47		20 48	I	21 49	L	22 50	I.	T	23 51	I	24 52		25 53		26 54	I	27 55	11
14		·+	+		+ +- + +-		+ -+ -+ -+ -+ -+ -+ -+ -+ -+ -+ -+ -+ -+	++	• +		+ + + +		+ ·	+	⊦ + ⊦ +		+ -	+	+	+	+	+++++++++++++++++++++++++++++++++++++++		+ + + +		+ + + +	+	· +		·+ -+
j.	32k	в	1 3	32kB	L L	32kB		32kB	I	32kB	I I	32kB		32kB	I	32kB	I.	32kB	1	I	32kB	I	32kB	I I	32kB	I I	32kB	I	32kB	тi
14		·+	+		+ +-		+ -	++	• +		+ + +	+	+ ·	+	⊦ + ⊦ +		+ -	+	+	+-	+	+		+ +	·	+ +	+	· +		·+ -+
	256	ikB	1 2	256kB	i i	256kB	i I	256kB	j	256kB	i i	256kB	i I	256kB	i	256kB	i.	256k	ві	i	256kB	j	256kB	i i	256kB	ii	256kB	i	256kB	ij
	++ ++ ++ ++ ++ +										+	⊦ + ⊦ +		+	+	+	+	+	+		+ +		+ +	+	+		·+ -+			
i i								17MB							i								17MB							i i
· · · · · · · · · · · · · · · · · · ·																										+				



Enforcing thread/process-core affinity

likwid-pin OpenMP affinity mechanisms
Example: OpenMP STREAM benchmark on 16-core system:

FFBE

Anarchy vs. thread pinning



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 also be used as a superior replacement for taskset
- Supports logical core numbering within a node

- Simple usage with physical (kernel) core IDs:
 - likwid-pin -c 0-3,4,6 ./myApp parameters
 - OMP_NUM_THREADS=4 likwid-pin -c 0-9 ./myApp parameters
- Simple usage with logical core IDs ("thread groups"):
 - likwid-pin -c S0:0-7 ./myApp parameters
 - likwid-pin -c C1:0-2 ./myApp parameters

Likwid Currently available thread domains







Running the STREAM benchmark with likwid-pin:

\$ likwid-pin -c S0:0-3 ./stream





Processor: smallest entity able to run a thread or task (hardware thread) Place: one or more processors \rightarrow thread pinning is done place by place Free migration of the threads on a place between the processors of that place.

OMP_PLACES	Place ==
threads	Hardware thread (hyper-thread)
cores	All HW threads of a single core
sockets	All HW threads of a socket
abstract_name(num_places)	Restrict # of places available

abstract name

Or use explicit numbering, e.g. 8 places, each consisting of 4 processors:

- OMP_PLACES="{0,1,2,3}, {4,5,6,7}, {8,9,10,11}, ... {28,29,30,31}"
- OMP_PLACES="{0:4}, {4:4}, {8:4}, ... {28:4}"
- OMP_PLACES="{0:4}:8:4"

<lower-bound>:<number of entries>[:<stride>]

Caveat: Actual behavior is implementation defined!



Determines how places are used for pinning:

OMP_PROC_BIND	Meaning
FALSE	Affinity disabled
TRUE	Affinity enabled, implementation defined strategy
CLOSE	Threads bind to consecutive places
SPREAD	Threads are evenly scattered among places
MASTER	Threads bind to the same place as the master thread that was running before the parallel region was entered

 If there are more threads than places, consecutive threads are put into individual places ("balanced")

Some simple OMP_PLACES examples



- Intel Xeon w/ SMT, 2x10 cores, 1 thread per physical core, fill 1 socket OMP_NUM_THREADS=10 OMP_PLACES=cores OMP_PROC_BIND=close

 Always prefer abstract places instead of HW thread IDs!
 - Intel Xeon Phi with 72 cores, 32 cores to be used, 2 threads per physical core OMP_NUM_THREADS=64 OMP_PLACES=cores(32) OMP_PROC_BIND=close # spread will also do
- Intel Xeon, 2 sockets, 4 threads per socket (no binding within socket!) OMP_NUM_THREADS=8 OMP_PLACES=sockets OMP_PROC_BIND=close # spread will also do
- Intel Xeon, 2 sockets, 4 threads per socket, binding to cores OMP_NUM_THREADS=8 OMP_PLACES=cores OMP_PROC_BIND=spread

likwid-mpirun MPI startup and Hybrid pinning



- How do you manage affinity with MPI or hybrid MPI/threading?
- In the long run a unified standard is needed
- Till then, likwid-mpirun provides a portable/flexible solution
- The examples here are for Intel MPI/OpenMP programs, but are also applicable to other threading models

```
Pure MPI:
likwid-mpirun -np 16 -nperdomain S:2 ./a.out
Hybrid:
likwid-mpirun -np 16 -pin S0:0,1 S1:0,1 ./a.out
```

likwid-mpirun 1 MPI process per node



likwid-mpirun -np 2 -pin N:0-11 ./a.out



Intel MPI+compiler:

OMP_NUM_THREADS=12 mpirun -ppn 1 -np 2 -env KMP_AFFINITY scatter ./a.out

(c) RRZE 2019

likwid-mpirun 1 MPI process per socket



likwid-mpirun -np 4 -pin S0:0-5_S1:0-5 ./a.out



Intel MPI+compiler:

OMP_NUM_THREADS=6 mpirun -ppn 2 -np 4 \ -env I_MPI_PIN_DOMAIN socket -env KMP_AFFINITY scatter ./a.out



Microbenchmarking for architectural exploration (and more)

Probing of the memory hierarchy Saturation effects in cache and memory

Latency and bandwidth in modern computer environments





(c) RRZE 2019



- 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 architectures, ever!

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





The throughput-parallel vector triad benchmark



Every core runs its own, independent triad benchmark

 \rightarrow pure hardware probing, no impact from OpenMP overhead

```
double precision, dimension(:), allocatable :: A,B,C,D
!$OMP PARALLEL private(i,j,A,B,C,D)
allocate(A(1:N), B(1:N), C(1:N), D(1:N))
A=1.d0; B=A; C=A; D=A
!SOMP SINGLE
stime = timestamp()
!$OMP END SINGLE
do j=1,NITER
  do i=1,N
    A(i) = B(i) + C(i) * D(i)
  enddo
  <<obscure dummy call>>
enddo
!$OMP SINGLE
etime = timestamp()
!$OMP END SINGLE
!$OMP END PARALLEL
Mflops = (2.d0*NITER)*N*num threads / (etime-stime) / 1.0e6
```

Throughput vector triad on Sandy Bridge socket (3 GHz)





Attainable memory bandwidth: Comparing architectures





(c) RRZE 2019

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

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)



Simplistic view of the hardware:



Simplistic view of the software:

```
! may be multiple levels
do i = 1,<sufficient>
        <complicated stuff doing
        N flops causing
        V bytes of data transfer>
enddo
```

Computational intensity $I = \frac{N}{V}$ \rightarrow Unit: flop/byte

A simple performance model for loops



How fast can tasks be processed? P [flop/s]

The bottleneck is either

- The execution of work:
- The data path:

 P_{peak} [flop/s] $I \cdot b_S$ [flop/byte x byte/s]





 $P_{peak} \left| \frac{F}{s} \right|$

 $b_S\left[\frac{B}{s}\right]$

Apply the naive Roofline model in practice

- Machine parameter #1:
- Machine parameter #2:
- Code characteristic:

Peak performance:

Memory bandwidth:

Computational Intensity: $I = \frac{F}{R}$



The Roofline Model – Basics



Compare capabilities of different machines



- RLM always provides upper bound but is it realistic?
- If code is not able to reach this limit (e.g. contains add operations only) machine parameter need to redefined (e.g., $P_{peak} \rightarrow P_{peak}/2$)



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





Multiple ceilings may apply

- Different P_{max}
 → different flat ceilings
- Different bandwidths / data paths
 - \rightarrow different inclined ceilings





Haswell port scheduler model:





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

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

```
\rightarrow 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_{max} of vector triad on Haswell (2.3 GHz)

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



*P*_{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: $b_{\rm S} = 50$ GB/s Code balance (incl. write allocate): $B_{\rm c} = (4+1)$ Words / 2 Flops = 20 B/F \rightarrow / = 0.05 F/B

 \rightarrow *I* · *b*_S = 2.5 GF/s (0.5% of peak performance)

 P_{peak} / core = 36.8 Gflop/s ((8+8) Flops/cy x 2.3 GHz) P_{max} / core = 12.27 Gflop/s (see prev. slide)

→ *P*_{max} = 14 * 12.27 Gflop/s =172 Gflop/s (33% peak)

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

A not so simple Roofline example



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

in single precision on an 8-core 2.2 GHz Sandy Bridge socket @ "large" N







The roofline formalism is based on some (crucial) prerequisites:

- 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 saturated bandwidth of the chip via simple streaming benchmarks to model more complex kernels and applications

Assumptions behind the model:

- 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, i.e. perfect streaming mode
- "Steady state" code execution (no wind-up/-down effects)





.

Node-Level Performance Engineering

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)
- 4. Hit the core bottleneck by good serial code
 (e.g., -fno-alias → see later)
- Shift P_{max} by accessing additional hardware features or using a different algorithm/implementation (e.g., scalar → SIMD)







Multicore performance tools: Probing performance behavior

likwid-perfctr

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
```
likwid-perfctr

Example usage with preconfigured metric group (shortened)





(c) RRZE 2019

(c) RRZE 2019

Node-Level Performance Engineering

98

likwid-perfctr *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-marker.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;
```



likwid-perfctr *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

Using Roofline for monitoring "live" jobs on a cluster

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





(c) RRZE 2019



Case study: A Jacobi smoother

The basic performance properties in 2D

Layer conditions Optimization by spatial blocking



- 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: $y(:) \leftarrow x(:)$

Swap y
$$\leftrightarrow$$
 x



enddo

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







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

> (here: Multiply by 4 FLOP/LUP to get FLOP/s rate) Node-Level Performance Engineering





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





(c) RRZE 2019



Case study: A Jacobi smoother

The basics in two dimensions Layer conditions Optimization by spatial blocking





(c) RRZE 2019



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

			miss				
		hit		miss			
			miss				



(c) RRZE 2019

Analyzing the data flow



Reduce inner (j-) loop dimension successively



x(0:jmax1+1,0:kmax+1)

		miss				
	hit		miss			
		miss				

Best case: 3 "layers" of grid fit into the cache!

k



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

(c) RRZE 2019





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)





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

\rightarrow "Spatial Blocking" of j-loop:

enddo

enddo enddo

New layer condition (blocking) 3 * jblock * 8B < CacheSize/2

→ Determine for given CacheSize an appropriate jblock value:

jblock < CacheSize / 48 B

Establish the layer condition by blocking



Split up						
domain into						
subblocks:						
o a block						
size = 5						
0120 - 0						

(c) RRZE 2019

Establish the layer condition by blocking





(c) RRZE 2019

Establish layer condition by spatial blocking





(c) RRZE 2019

Validating the hypotheis: Measure memory code balance





(c) RRZE 2019

OpenMP parallelization of the blocked 2D stencil







 Caveat: LC must be fulfilled per thread → shared cache causes smaller blocks!







OpenMP parallelization and blocking for shared cache





Stencil shapes and layer conditions



- a) Long-range r = 2: 5 layers (2r + 1)
- b) Long-range r = 2 with gaps: 6 layers (2 per populated row)
- c) Asymmetric: 3 layers
- d) 2D box: 3 layers



- 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?

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: <u>arXiv:1410.5010</u>



2

3

4

cores

5

Node-Level Performance Engineering

7

6



Case study: Sparse Matrix Vector Multiplication

Sparse Matrix Vector Multiplication (SpMV)

- **FFEE**
- Key ingredient in some matrix diagonalization algorithms
 - Lanczos, Davidson, Jacobi-Davidson
- Store only N_{nz} nonzero elements of matrix and RHS, LHS vectors with N_r (number of matrix rows) entries
- "Sparse": N_{nz} ~ N_r



SpMVM characteristics



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





- **val[]** stores all the nonzeros (length N_{nz})
- col_idx[] stores the column index
 of each nonzero (length N_{nz})
- **row_ptr[]** stores the starting index of each new row in **val[]** (length: N,)





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

```
!$OMP parallel do schedule(???)
do i = 1,Nr
  do j = row_ptr(i), row_ptr(i+1) - 1
    c(i) = c(i) + val(j) * b(col_idx(j))
  enddo
enddo
enddo
!$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
- Can we explain this?
- Is there a "light speed" for SpMV?
- Optimization?





Example: SpMV 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

$$B_{c}^{DP,CRS} = \frac{8 + 4 + 8\alpha + 20/N_{nzr}}{2} \frac{B}{F} = \left(6 + 4\alpha + \frac{10}{N_{nzr}}\right) \frac{B}{F}$$

Absolute minimum code balance:
$$B_{\min} = 6 \frac{B}{F}$$

 $\Rightarrow I_{\max} = \frac{1}{6} \frac{F}{B}$

Hard upper limit for in-memory performance: b_S/B_{\min}

The " α effect"



 $B_c^{DP,CRS}(\alpha) = \frac{8+4+8\alpha+20/N_{nzr}}{2}\frac{B}{F}$

 $=\left(6+4\alpha+\frac{10}{N_{max}}\right)\frac{B}{F}$

DP CRS code balance

- α quantifies the traffic for loading the RHS
 - $\alpha = 0 \rightarrow \text{RHS}$ is in cache
 - $\alpha = 1/N_{nzr} \rightarrow RHS$ loaded once
 - $\alpha = 1 \rightarrow$ no cache
 - $\alpha > 1 \rightarrow$ Houston, we have a problem!
- "Target" performance = b_S/B_c
- Caveat: Maximum memory BW may not be achieved with spMVM (see later)

Can we predict α ?

- Not in general
- Simple cases (banded, block-structured): Similar to layer condition analysis

\rightarrow Determine α by measuring the actual memory traffic

Determine α (RHS traffic quantification)



$$B_{c}^{DP,CRS} = \left(6 + 4\alpha + \frac{10}{N_{nzr}}\right)\frac{B}{F} = \frac{V_{meas}}{N_{nz} \cdot 2 F}$$

 V_{meas} is the measured overall memory data traffic (using, e.g., likwidperfctr)

Solve for
$$\alpha$$
:

$$\alpha = \frac{1}{4} \left(\frac{V_{meas}}{N_{nz} \cdot 2 \text{ bytes}} - 6 - \frac{10}{N_{nzr}} \right)$$

Example: kkt_power matrix from the UoF collection on one Intel SNB socket

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

- $V_{meas} \approx 258 \text{ MB}$
- $\rightarrow \alpha = 0.36, \alpha N_{nzr} = 2.5$
- → RHS is loaded 2.5 times from memory



$$\frac{B_c^{DP,CRS}(\alpha)}{B_c^{DP,CRS}(1/N_{nzr})} = 1.11$$

11% extra traffic → optimization potential!



(c) RRZE 2019



Benchmark system: Intel Xeon Ivy Bridge E5-2660v2, 2.2 GHz, $b_S = 46.6 \text{ GB/s}$

Matrix	Ν	N _{nzr}	B_c^{opt} [B/F]	P_{opt} [GF/s]
DLR1	278,502	143	6.1	7.64
scai1	3,405,035	7.0	8.0	5.83
kkt_power	2,063,494	7.08	8.0	5.83











kkt_power
Now back to the start...





Node-Level Performance Engineering

Investigating the load imbalance with kkt_power







- Conclusion from the 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
 - Do not forget about load balancing!
- Consequence: Modeling is not always 100% predictive. It's all about learning more about performance properties!



Case study: Tall & Skinny Matrix-Transpose Times Tall & Skinny Matrix (TSMTTSM) Multiplication

TSMTTSM Multiplication



- Block of vectors \rightarrow Tall & Skinny Matrix (e.g. 10⁷ x 10¹ dense matrix)
- Row-major storage format (see SpMVM)
- Block vector subspace orthogonalization procedure requires, e.g. computation of scalar product between vectors of two blocks
- TSMTTSM Mutliplication





General rule for dense matrix-matrix multiply: Use vendor-optimized GEMM, e.g. from Intel MKL¹:



¹Intel Math Kernel Library (MKL) 11.3

TSMTTSM Roofline model

#flops #bytes (slowest data path)

 $C = \alpha \qquad A^T \qquad * B + \beta C$ Optimistic model (minimum data transfer) assuming $M = N \ll K$ and

double precision:

Computational intensity

=

 $I_d \approx \frac{2KMN}{8(KM+KN)}\frac{F}{B} = \frac{M}{8}\frac{F}{B}$

 $= \alpha$

complex double:

$$I_z \approx \frac{8KMN}{16(KM+KN)}\frac{F}{B} = \frac{M}{4}\frac{F}{B}$$



 $\left| * \frac{\beta}{\beta} + \beta \right|$



Now choose
$$M = N = 16 \rightarrow I_d \approx \frac{16}{8} \frac{F}{B}$$
 and $I_z \approx \frac{16}{4} \frac{F}{B}$, i.e. $B_d \approx 0.5 \frac{B}{F}$, $B_z \approx 0.25 \frac{B}{F}$

Intel Xeon E5 2660 v2 ($b_S = 52 \frac{GB}{s}$) $\rightarrow P = 104 \frac{GF}{s}$ (double) Measured (MKL): 16.6 $\frac{GF}{s}$

Intel Xeon E5 2697 v3 ($b_S = 65 \frac{GB}{s}$) $\rightarrow P = 240 \frac{GF}{s}$ (double complex) Measured (MKL): 22.8 $\frac{GF}{s}$

→ Potential speedup: 6–10x vs. MKL

Can we implement a better TSMTTSM kernel than Intel?





Not shown: Inner Loop boundaries (n,m) known at compile time (kernel generation) k assumed to be even



TS: M=N=16; K=10,000,000

System	P _{peak} / b _S	Version	Performance	RLM Efficiency	
Intel Xeon E5 2660 v2 10c@2.2 GHz	176 GF/s	TS OPT	98 GF/s	94 %	
	52 GB/s	TS MKL	16.6 GF/s	16 %	
Intel Xeon E5 2697 v3	582 GF/s 65 GB/s	TS OPT	159 GF/s	66 %	
14c@2.6GHz		TS MKL	22.8 GF/s	9.5 %	



ERLANGEN REGIONAL COMPUTING CENTER



Single Instruction Multiple Data (SIMD)



FRIEDRICH-ALEXANDER UNIVERSITÄT ERLANGEN-NÜRNBERG

SIMD terminology

A word on terminology

- SIMD == "one instruction → several operations"
- "SIMD width" == number of operands that fit into a register
- No statement about parallelism among those operations
- Original vector computers: long registers, pipelined execution, but no parallelism (within the instruction)
 R0
 R1

Today

- x86: most SIMD instructions fully parallel
 - "Short Vector SIMD"
 - Some exceptions on some archs (e.g., vdivpd)
- NEC Tsubasa: 32-way parallelism but SIMD width = 256 (DP)



Example: Data types in 32-byte SIMD registers (AVX[2])

Supported data types depend on actual SIMD instruction set

double		ible dou		double		double double		double doub		double
float	float	float	float	float	float	float float				
int	int	int	int	int	int	int int				





Scalar slot



RIEDRICH-ALEXANDER INIVERSITÄT RLANGEN-NÜRNBERG

SIMD processing – Basics

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



SIMD processing – Basics

No SIMD vectorization for loops with data dependencies:

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

"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" is not used, tell 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) {...}

How to leverage SIMD: your options

Options:

- The compiler does it for you (but: aliasing, alignment, language, abstractions)
- 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:

- xmmintrin.h (SSE)
- pmmintrin.h (SSE2)
- immintrin.h (AVX)
- x86intrin.h (all extensions)

User-mandated vectorization (OpenMP 4)

- Since OpenMP 4.0 SIMD features are a part of the OpenMP standard
- #pragma omp simd enforces vectorization
- Essentially a standardized "go ahead, no dependencies here!"
 - Do not lie to the compiler here!
- Prerequesites:

}

}

- Countable loop
- Innermost loop
- Must conform to for-loop style of OpenMP worksharing constructs
- There are additional clauses:

```
reduction, vectorlength, private, collapse, ...
for (int j=0; j<n; j++) {
    #pragma omp simd reduction(+:b[j:1])
    for (int i=0; i<n; i++) {
        b[j] += a[j][i];</pre>
```

SIMD is an in-core feature!







- 1. Inner loop
- 2. 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)
- 2. 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:

- 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?



How much bandwidth does nonlocal access cost?



Example: AMD "Epyc" 2-socket system (8 chips, 2 sockets, 48 cores): *STREAM Triad bandwidth measurements* [Gbyte/s]

CPU nod	e 0	1	2	3	4	5	6	7
MEM node 0	32.4	21.4	21.8	21.9	10.6	10.6	10.7	10.8
1	21.5	32.4	21.9	21.9	10.6	10.5	10.7	10.6
2	21.8	21.9	32.4	21.5	10.6	10.6	10.8	10.7
3	21.9	21.9	21.5	32.4	10.6	10.6	10.6	10.7
4	10.6	10.7	10.6	10.6	32.4	21.4	21.9	21.9
5	10.6	10.6	10.6	10.6	21.4	32.4	21.9	21.9
6	10.6	10.7	10.6	10.6	21.9	21.9	32.3	21.4
7	10.7	10.6	10.6	10.6	21.9	21.9	21.4	32.5



Node-Level Performance Engineering

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:

```
numactl --membind=<nodes> a.out  # map pages only on <nodes>
    --preferred=<node> a.out  # map pages on <node>
    # and others if <node> is full
    --interleave=<nodes> a.out  # map pages round robin across
    # all <nodes>
```

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
   done
   done
```

But what is the default without numactl?



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: "to touch" means "to write", not "to allocate"
- Example:

Memory not mapped here yet

double *huge = (double*)malloc(N*sizeof(double));

```
for(i=0; i<N; i++) // or i+=PAGE_SIZE/sizeof(double)
huge[i] = 0.0;
Mapping takes
place here</pre>
```

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

Coding for ccNUMA data locality



Most simple case: explicit initialization



Coding for ccNUMA data locality



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



Node-Level Performance Engineering

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)
 - Presupposes that all worksharing loops with the same loop length have the same thread-chunk mapping
 - If dynamic scheduling/tasking is unavoidable, more advanced methods may be in order
 - OpenMP 5.0 will have rudimentary memory affinity functionality
- 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

Diagnosing bad locality



- If your code is cache bound, you might not notice any locality problems
- Otherwise, bad locality limits scalability (whenever a ccNUMA node boundary is crossed)
 - Just an indication, not a proof yet
- 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 dual-socket system (IvyBridge, 2x10-core):

likwid-perfctr -g NUMA -C M0:0-4@M1:0-4 ./a.out

Using performance counters for diagnosing bad ccNUMA access locality



 Intel Ivy Bridge EP node (running 2x5 threads): measure NUMA traffic per core

likwid-perfctr -g NUMA -C M0:0-4@M1:0-4 ./a.out

Summary output:

L			L	±	┸
Metric	Sum	Min	Max	Avg	T I
Runtime (RDTSC) [s] STAT	4.050483	0.4050483	0.4050483	0.4050483	+
Runtime unhalted [s] STAT	3.03537	0.3026072	0.3043367	0.303537	I
Clock [MHz] STAT	32996.94	3299.692	3299.696	3299.694	I
CPI STAT	40.3212	3.702072	4.244213	4.03212	I
Local DRAM data volume [GByte] STAT	7.752933632	0.735579264	0.823551488	0.7752933632	I
Local DRAM bandwidth [MByte/s] STAT	19140.761	1816.028	2033.218	1914.0761	I
Remote DRAM data volume [GByte] STAT	9.16628352	0.86682464	0.957811776	0.916628352	I
Remote DRAM bandwidth [MByte/s] STAT	22630.098	2140.052	2364.685	2263.0098	I
Memory data volume [GByte] STAT	16.919217152	1.690376128	1.69339104	1.6919217152	I
Memory bandwidth [MByte/s] STAT	41770.861	4173.27	4180.714	4177.0861	I

 Caveat: NUMA metrics vary strongly between CPU models About half of the overall memory traffic is caused by remote domain!

The curse and blessing of interleaved placement: OpenMP STREAM triad on a dual AMD Epyc 7451 (6 cores per LD)



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



Node-Level Performance Engineering

rr@=

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

NUMA balancing is active on many Linux systems today

- Automatic page migration
- Slow process, may take many seconds (configurable)
- Still a good idea to to parallel first touch

If dynamic scheduling cannot be avoided

- Consider round-robin placement as a quick (but non-ideal) fix
- OpenMP 5.0 has some data affinity support



OpenMP performance issues on multicore

Barrier synchronization overhead Topology dependence



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)







!\$OMP PARALLEL ...

\$0MP 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	ntel 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	Overhead grows
Node (20 cores)	4999	7783	with thread count
Node +SMT	5981	9897	•

- Strong dependence on compiler, CPU and system environment!
- OMP_WAIT_POLICY=ACTIVE can make a big difference





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

3.2x cores (20 vs 64) on Phi 4x 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 $1.9 \cdot 12.8 \approx 24x$ more pain \odot .

Tutorial conclusion

- Multicore architecture == multiple complexities
 - Affinity matters \rightarrow pinning/binding is essential
 - Bandwidth bottlenecks \rightarrow inefficiency is often made on the chip level
 - Topology dependence of performance features \rightarrow know your hardware!

Put cores to good use

- Bandwidth bottlenecks \rightarrow surplus cores \rightarrow 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
- I learn more, especially when they do not work!

Simple tools get you 95% of the way

e.g., with the LIKWID tool suite

Node-Level Performance Engineering







Moritz Kreutzer Markus Wittmann Thomas Zeiser Michael Meier Holger Stengel Thomas Gruber Faisal Shahzad Christie L. Alappat Ayesha Afzal Dominik Ernst Julian Hammer Jan Laukemann

THANK YOU.







Bundesministerium für Bildung und Forschung

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 https://blogs.fau.de/hager 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



- SC19 full-day tutorial: Node-Level Performance Engineering
- Presenter(s): Georg Hager, 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.