





Memory hierarchies Cache-based dual-core architectures									
*SSE2	Intel Xeon5100	AMD Opteron							
Peak Performance	12.0 GFlop/s	5.2 GFlop/s							

*SSE2		Xeon5100	Opteron	Itanium2
Peak Per Core frequ	formance Jency	12.0 GFlop/s 3.0 GHz	5.2 GFlop/s 2.6 GHz	6.4 GFlop/s 1.6 GHz
#Register	rs	8 / 16 [*]	16 / 32 [*]	128
	Size	32 kB	64 kB	16 KB
L1	BW	96 GB/s	41.6 GB/s	51.2 GB/s
	Latency	2 cycles	3 cycles	1 cycle
	Size	4 MB (2 cores)	1 MB	256 KB
L2	BW	96 GB/s	41.6 GB/s	51.2 GB/s
	Latency	7 cycles	~13 cycles	5-6 cycles
	Size			6 / 12 MB
L3	BW			51.2 GB/s
	Latency			12-13 cycles
Maman	BW	10.6 GB/s	6.4 GB/s	8.5 GB/s
wemory	Latency	~200 ns	< 100 ns	~200 ns
5.2007 (*	13)	hpc@rrze.uni-erla	ngen.de	ParCFD 2007 Short Course

Memory hierarchies: Cache structure



Intel

- Cache line data is always consecutive
 - Cache use is optimal for contiguous access (stride 1)
 - Non-consecutive access reduces performance (worst case: \geq cache line size)
 - Ensure spatial locality by blocking or optimizing data layout
- Caches (~MB) must be mapped to memory locations (~GB)
 - Cache multi-associativity enhances utilization



Memory hierarchies Cache structure



- If one item is needed, the cache line it belongs to is fetched (miss)
- Cache line fetch/load has large latency penalty
- "neighboring" items can then be used from cache н.



Memory hierarchies Cache mapping



Cache mapping: н.

(16)

- Pairing of memory locations with cache line
- e.g. mapping 1 GB of main memory to 512 KB of cache

Directly mapped cache: н.

- Every memory location can be mapped to exactly one cache location
- If cache size=n, i-th memory location can be stored at cache location mod(i,n)
- Memory access with stride=cache size will not allow caching of more than one line of data, i.e. effective cache size is one line!
- No penalty for stride-one access

HPC

Short Course











Lattice Boltzmann method: Basic implementation strategy



F(0:18,x,y,z,0:1)

- Use "full matrix" representation:
 - $F(\alpha, x, y, z, 0) = f_{\alpha}(x, y, z, t) \& F(\alpha, x, y, z, 1) = f_{\alpha}(x, y, z, t+1)$ (t odd; $\vec{x} = (x, y, z)$)
 - For complex geometries, *vector-like representations* might be interesting as they can block out the solid parts
 - However, the connectivity of the cells has to be stored additionally (therefore, you need at least 1/3 of solid to save memory)

Basic Optimizations

- Analyze relaxation step: Many operations can be eliminated (common sub-expressions, zero-velocity components):
 # of floating point operations depends on compiler & optimization level!
- Combine Collide & Stream step in a single loop to minimize data transfer!

hpc@rrze.uni-erlangen.de



Lattice Boltzmann method Performance (1)



 Standard performance measure for LBM: Million fluid cell updates per second: MFLUP/s (or Million Lattice cell updates per second: MLUP/s)

Performance estimation

- Assumption:
 200 Floating point operations per fluid cell update
- If data transfer is infinite fast (all data fits into cache):
 - Max. MFLUP/s = PeakPerformance / (200 Flop/cell update)
 - Peak Performance = 4,000 MFLOP/s \rightarrow Max. 20 MFLUP/s
 - In reality complex kernels do not achieve peak performance and transfer speed is finite even for caches → ~10-15 MFLUP/s

Lattice Boltzmann method Performance (2)



- If data must be transferred from and to main memory in each time step:
 - Assumption: full use of each cache line loaded
 - Data to be transferred for a single fluid cell update: (2+1)*19*8 Byte → 456 Bytes/(fluid cell update)
 - Max. MFLUP/s= MemoryBandwidth / (456 Bytes/(fluid cell update))
 - MemoryBandwidth = 6,000 MByte/s \rightarrow Max. 13 MFLUP/s
 - In reality even simple kernels only achieve 50%-60% of theoretical main memory bandwidth → ~ 6-8 MFLUP/s
- Crossover between cache and memory bound computations: 2 *19 * $xMax^3$ * 8Byte ~ L2/L3 cache size \rightarrow xMax ~ 14-15 (for 1 MB cache)











Literature



Performance of LBM:

T. Pohl, N. Thürey, F. Deserno, U. Rüde, P. Lammers, G. Wellein, and T. Zeiser, In: IEEE/ACM: Proceedings of the IEEE/ACM SC2004 Conference (2004), pp. 1-13. Performance Evaluation of Parallel Large-Scale Lattice Boltzmann Applications on Three Supercomputing Architectures

T. Pohl, M. Kowarschik, J. Wilke, K. Iglberger, and U. Rüde, Parallel Processing Letters. Vol. 13, No. 4 (2003), pp. 549-560.

Optimization and Profiling of the Cache Performance of Parallel Lattice Boltzmann Codes

G. Wellein, T. Zeiser, G. Hager, and S. Donath, Computers & Fluids, Vol. 35 (2006), pp. 910-919.

On the Single Processor Performance of Simple Lattice Boltzmann Kernels

M. Schulz, M. Krafczyk, J. Tölke, E. Rank, in: M. Breuer, F. Durst, C. Zenger (Eds.), High Performance Scientific and Engineering Computing, Springer, Berlin (2001), pp. 115–122.

Parallelization strategies and efficiency of CFD computations in complex geometries using lattice Boltzmann methods on high performance computers

23.05.2007 (45) hpc@rrze.uni-erlangen.de







Literature Performance of LBM: J. Habich, Bachelor Thesis, Erlangen 2006. Improving computational efficiency of Lattice Boltzmann methods on complex geometries S. Donath, Bachelor Thesis, Erlangen 2004. On Optimized Implementations of the Lattice Boltzmann Method on Contemporary Architectures K. Iglberger, Bachelor Thesis, Erlangen 2003. Performance Analysis and Optimization of the Lattice Boltzmann Method in 3D Find these and more interesting Bachelor and Master Theses at: http://www10.informatik.uni-erlangen.de/en/Publications/Theses/ ParCFD 2007 23.05.2007 (46) hpc@rrze.uni-erlangen.de Short Course **Efficient implementation of** lattice Boltzmann kernels – Part II

T. Zeiser, G. Wellein, G. Hager, S. Donath HPC Services Regionales Rechenzentrum Erlangen Friedrich-Alexander-University Erlangen-Nuremberg Germany



International Conference on Parallel Computational Fluid Dynamics 2007 – Short Course





Principle idea of "patch" approach



- Divide the domain into small patches
 - Use full arrays within the each patch
 - Use ghost cells for exchange between patches
 - only allocate patches which contain fluid areas
 - a patch should be much smaller than the total domain to block out most of the solid
 - but large enough to allow efficient processing
- Very much like parallelization with domain decomposition (but many patches go to the same processor)







References: IWTM/FHG, LSS/Uni-Erlangen, ...

23.05.2007 (53)

hpc@rrze.uni-erlangen.de



Optimization of data locality



 Processing of 1-D list no longer depends on lexicographic ordering owing to (necessary) use of adjacency list

→ added flexibility, e.g.

- "implicit" blocking (without loop overhead)
- use of space filling curves

· · · · ·

- \rightarrow No need for the solver to know about the ordering!
- → Automatic handling of bounce back by adjacency list!



Unstructured approach: LB-DC implementation



Parallelization by domain decomposition



- Cartesian cutting is no longer appropriate
- "unstructured grid" → use graph pratitioner (e.g. metis)
- or just cut the 1-D list into equal chunks
- → avoiding load imbalance is more important than minimizing communication
- → Trying to overlap communication and computation is generally not worth the effort as most MPI implementation do not support it out of the box (i.e. data transfer only takes place while within an MPI call).













Intel Trace Collector/Analyzer (I)



C

Summary



Intel Trace Collector/Analyzer (III)







Intel Trace Collector/Analyzer (IV)



Eile S	tyle <u>W</u> indows	Help F1	*				X
View Cr	ians <u>N</u> avigate	Advanced Lay	Tota	l Tima lel (Proc	acc by Collective	Oneration)	
ñ.	UDI Denet	NBI Badwaa		Maan	RMDan	speranony	
80	WPI_Bcast	11 199er 3	11. 402e-3	E 701er3	5 499o. 2		
P0 D4	10 6040-3	11.1090-5	10.6430-3	5.701e-3	5.28250-3		19.8+-3
Po	9,9020-3	550.6	9.9570-3	4.97850-3	4.92350-3		
De	6.2754.3	560-6	6.331e-3	3.1655a-3	3.10954-3		
PA	5.5860.3	510-6	5.6370-3	2.8185013	2.76750-3		17.6e-3
P5	4.873e-3	580-6	4.931e-3	2.4655e-3	2.4075e-3		
P6	1.395e-3	53e-6	1.4480-3	724e-6	671e-6		
P7	605e-6	55e-6	660e-6	330e-6	275e-6		15.4e-3
PB	6354-6	2136-6	8486-6	4240-6	211e-6		
P9	525e+6	81e-6	606e-6	303e-6	222e-6		
P10	671e-6	281e-6	952e-6	476e-6	195e-6		13.28=3
P11	505e-6	163e-6	668e-6	334e-6	171e-6		
P12	1.725e-3	539e-6	2.264e-3	1.132e-3	593e-6		11e-3
P13	20.946e-3	4.384e-3	25.33e-3	12.665e-3	8.281e-3		
P14	20.42e+3	4.339e-3	24.759e-3	12.3795e-3	8.0405e-3		
P15	20.887e-3	4.44e-3	25.327e-3	12.6635e-3	8.2235e-3		8.8e-3
Sum	105.767e-3	25.996e-3	131.763e-3				
Mean	6.61044e-3	1.62475e-3		4.11759e-3			
StdDev	7.5181e-3	2.97052e-3			6.235964-3		6.68-3
							4.40-3
							2.24-3
							0
13.210 8	932, 13.269 732	: 0.058 799	680.	All_Proce	66665	Application expanded in {	.) Tag Filter
0.07	(00)						ParCFD 2007
JU7	(66)		hpc@)rrze.uni-	erlangen.d	9	Short Course
							Short Course

Intel VTune (II) General Events EXT SNOOP Ratios Event groups: All Events Event Code: 0x7 Available events Mask: Combined mask from tables below. Category: Multi-Core Events:External Bus Events L2_ST.SELF.MESI Definition: External snoops. L2_ST.SELF.M_STATE Description: This event counts the snoop resp tion: This event counts the shoop responses to bus transactions. Responses can ted separately by type and by bus agent. With the 'THIS_AGENT' mask the event L2_ST.SELF.S_STATE Event description Explain... L2 store requests Click the "Explain... " button for more details on the -1selected event. Selected events Counts bus transactions initiated by any agent on the bus, in systems where each processor i attached to a different bus, each core counts Event Name Sample After CPU CLK UNHALTED.CORE 2992000 INST RETIRED.ANY 2992000 SIMD_INST_RETIRED.ANY 1000000 RUS TRANS MEM ALL AGENTS 100000 Counts snoop CLEAN responses. They occur when the responder's cache does not have the snooped address. Run Information Calibration runs: 1 Counts snoop HITM responses. They occur when the snooped address at the responder's cache is in a Modified state. Sampling runs: 1 OK Cancel ANY 0x28 Apply 1.0 ParCFD 2007 23.05.2007 (68) hpc@rrze.uni-erlangen.de Short Course C

Intel VTune (III)



Related Topics		200	1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1	A. @	三 ()					A REAL PROPERTY				
About Sampling Hotspot	Name	CPU_C sample	INST_RE samples	SIMD_IN samples	BUS_TR samples	Clocks per	CPU_CL	INST_R %	SIMD_I	BUS_TR	CPU_CLK_UN = events	-	Events	
The Motron training displays	mod_flow_trt_mp_flow	162 351	7 121,041	213,045	171,536	2.526	94,73%	98.94%	99.42%	94.77%	366,580,350	DUCE 1	TRANS NO	M
function names associated	mod flow mp flow ou	5,701	1 599	841	6,418	17.927	3.33%	0.49%	0.39%	3.55%	12,872,094	DUS_	TRANS_ME	
with selected modules that	mod_flow_mp_flow_ini	2,182	2 370	315	2,194	11.108	1.27%	0.30%	0.15%	1.21%	4,926,663	BUS_	IRANS_ME	M
have symbol information	calc_dens_ux_	1,014	4 .97	87	696	19.690	0.59%	0.08%	0.04%	0.38%	2,289,476	CRU	S per msp.	
available.	MAIN_	64	4 91	0	57	1.325	0.04%	0.07%	0.00%	0.03%	144,503	CPU	CLA UNHA	
You can group the hotspots	read_casefile_arrays_	41	1 126	0	31	0.613	0.02%	0.10%	0,00%	0.02%	92,572	CPU	CLK_UNHA	
Address (RVA). Source Files.	check_	24	4 5	5	30	9.041	0.01%	0.00%	0.00%	0.02%	54,188	LINCT .	OFTIOTO A	
and by Class.	check_cs_flow_	21	1 8	1	26	4.944	0.01%	0.01%	0.00%	0.01%	47,415	I INST	RETIRED.A	
You can view the results in a	out_x_		3 4		5	3,767	0.00%	0.00%	0.00%	0.00%	18,062	INST.	RETIRED.A	
Table or Horizontal Bar Char	mp_lib_mp_mp_wtime_	1	1 0	0	0	0.000	0.00%	0.00%	0.00%	0.00%	2,257	1451	NET OFT	
by right-clicking the view and	for	1	1 0	0	0	0.000	0.00%	0.00%	0.00%	0.00%	2.257	SIMU	INST RET	
or View as Bar Chart from the	for_flush_readahead	1	1 0	0	0	0.000	0.00%	0.00%	0.00%	0.00%	2.257	SIMU	INSI REI	
pop-up menu.	4 H H											1 SIMD	INSI_REI	
From the Hotspot view you	Autority ID		A CARLON AND A							able of Max		cours	D HAN	
Go To:	ACOVITY ID	CONTRACTOR OF	ACTIVITY He	SUIT.	-	10	tai samp	es Dura	tion Ma	cnine war	ne	CPUT	0 10 10	June
All Topics Search	en Sau way eo	12 34 25 1	2007 - 341	family dis-	Inte Late 4, 0	0.11 50.	12006	1000	and here	100	2 T COMPCO	÷		0.00
Bookmarks = Index												2	0	0.00
Tuning Bro II Navigator ** D												-	3	0.00
	1											-		0.00
Sampling Activity	•			1000							-	1		
· Standary concentration	Processes Threads	Modul	les Hot	spots										
v n Run 1	Console 11											14	0.68	• • • •
BINST RETIRED ANY BIND INST RETIRED ANY BSIMD INST RETIRED AN BBUS TRANS MEM ALL	Tuning Console Sun May 20 12: 32:07 200 Sun May 20 12: 32:07 200 Sun May 20 12: 34:23 200	7 127.0.0. 7 127.0.0. 7 127.0.0.	1 (Run 0) C 1 (Run 0) C 1 (Run 0) S	Collection I PU_CLK_U Campling d	or the folio NHALTED. ata was su	core, in core, in	ent(s) is b ST_RETIR ly collecti	eing perf ED.ANY, S ed.	formed: 5IMD_INS	T_RETIRE	D.ANY, BUS_TRA	NS_MEM	1.ALL_AGE	NTS.

Elle Edit Navigate Segrch Project Tunin	ig Bun Window Help				
	- 8-8-9-0-0-			-	VTunet
Start 🗅 Help 🗉 🚽 🐨 🗖 🔂 Welce	ome Call Graph Results (127.0.0.1) - Su mod flow trt.F90 Sun May 20 12:1	34:23 2007 - Sampli	- mos	flow trt.	0 - 0
Related Topics	AST CONTACT OF A C				
- About Source/Assembly	Source	CPU_CLK	INST_RE	SIMD_IN	BUS_TRA
Right-click the events for	f local density				
display options. Use the	loc_dens = dd_tmp_0 &	0.26%	0.36%	0.43%	0.26%
toolbar buttons to:	+ dd_tmp_NE + dd_tmp_N + dd_tmp_NW + dd_tmp_W &	0.35%	0.33%	0.29%	0.30%
- Switch	+ dd_tmp_SW + dd_tmp_S + dd_tmp_SE + dd_tmp_E &	0.20%	0.38%	0.48%	0.18%
disassembly views.	+ dd_tmp_T + dd_tmp_TE + dd_tmp_TN + dd_tmp_TW &	0.25%	0.37%	0.87%	0.20%
- Navigate to	+ dd tmp TS + dd tmp B + dd tmp BE + dd tmp BN &	0.29%	0.29%	0.20%	0.20%
the function start or to the	+ dd_tmp_BW + dd_tmp_BS	0.09%	0.28%	0.39%	0.10%
next function. 313					
the code lines that took a long 315	I be a dd trop NE a dd trop SE a dd trop E a dd trop TE a dd trop PE E	0.04%	0.10%	0.05%	0.03%
time to execute. Click a 316	dd two NW, dd two W, dd two SW, dd two TW, dd two PW	0.31%	1.04%	0.62%	0.26%
column to select the event and 317	. og "uuh"uu . og "uuh"u . og "uuh"uu . og "uuh"uu . og "uuh"uu				
click the "largest", "previous", 318	A Tractal to confine the				
- Select code 219	to the presenting of the second	1.40%	1 77%	1 38%	1 76%
Ge Te		<u>.</u>			1 1
All Topics 7 Search Size	Name CPU_C_INST_R_SIMD_I_BUS_TClo	ocks per Instructions Re	stired - CPI		
Bookmarks ill Index	Selected Range 1.44% 2.71% 3.47% 1.30%				
0x10AF	mod flow trt mp flow advrei 94.71% 98.94% 99.42% 94.77%				
Tuning Bro II Navigator					
# Sampling Activity					
♥ Sun May 20 12:34:23 2007					
™ 素Run 1 □Const	ole #			a 1 3	0.0.0
ECPU_CLK_UNHALTED.C/ Tuning I	Console				
DINST_RETIRED.ANY Sun May	y 20 12:32:07 2007 127.0.0.1 (Run 0) Collection for the following event(s) is being performed:	CT AFTIRED ANY DUE	TRANC M	-	
21 SIMD_INST_RETIRED.AN Sun May	y 20 12:34:23 2007 127.0.0.1 (Run 0) Sampling data was successfully collected.	an _ actimic parts BUS	C. Sherry M.	and make pri	-
BUS_TRANS_MEM.ALL					
		DarC		07	1000
2007 (70)	hno@rrzo uni orlongon do	Farur	- J ZU	07	CX