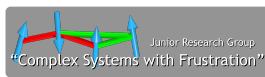


Simulating Spin Models on GPU

Martin Weigel

Institut für Physik, Johannes-Gutenberg-Universität Mainz, Germany

Monte Carlo Algorithms in Statistical Physics
Melbourne, July 26–28, 2010



GPU computing



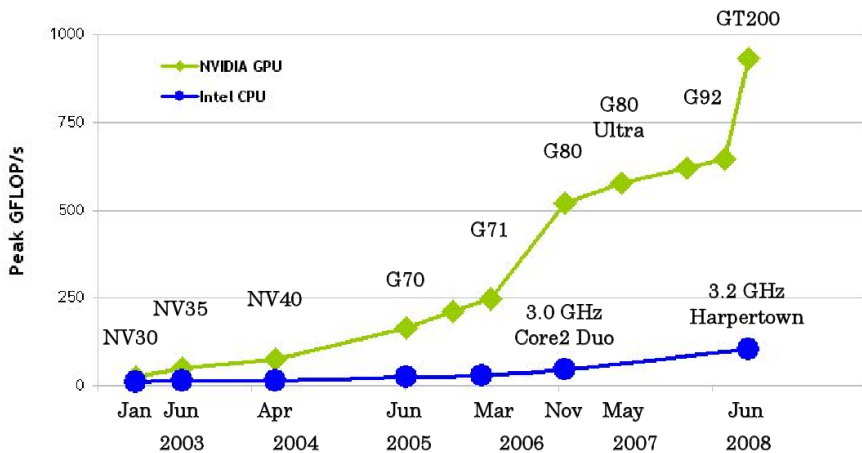
traditional interpretation of GPU computing

GPU computing

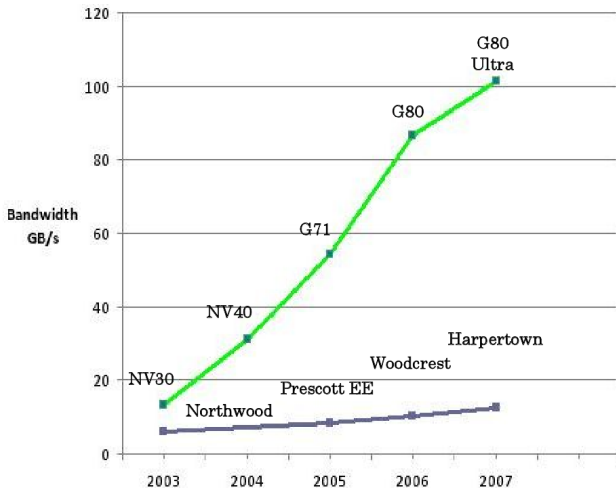


traditional interpretation of GPU computing

GPU computing



GPU computing



GPU computation frameworks

GPGPU = General Purpose Computation on Graphics Processing Unit

“Old” times: use original graphics primitives

- OpenGL
- DirectX

Vendor specific APIs for GPGPU:

- NVIDIA CUDA: library of functions performing computations on GPU (C, C++, Fortran), additional preprocessor with language extensions
- ATI/AMD Stream: similar functionality for ATI GPUs

Device independent schemes:

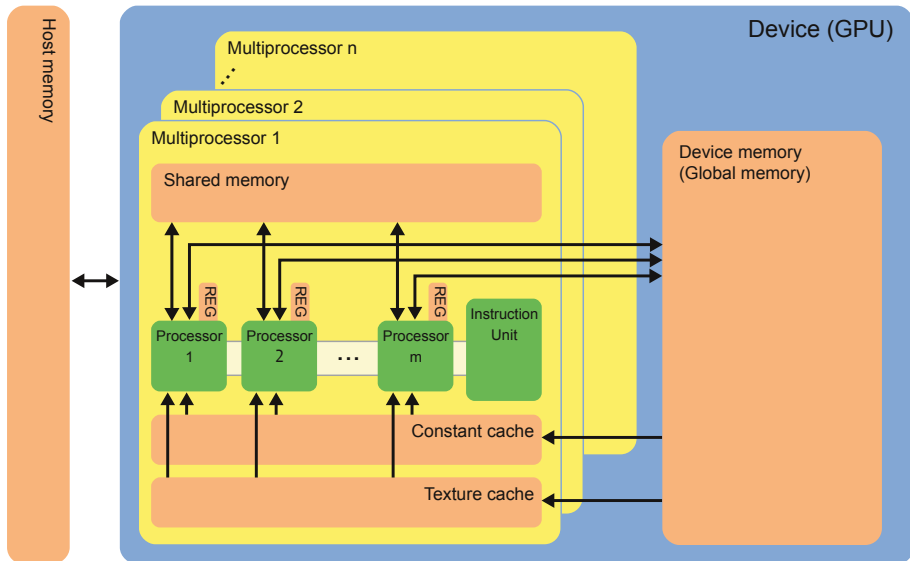
- BrookGPU (Stanford University): compiler for the “Brook stream program language” with backends for different hardware; now merged with AMD Stream
- Sh (University of Waterloo): metaprogramming language for programmable GPUs
- OpenCL (Open Computing Language): open framework for parallel programming across a wide range of devices, ranging from CPUs, Cell processors and GPUs to handheld devices

NVIDIA Tesla C1060

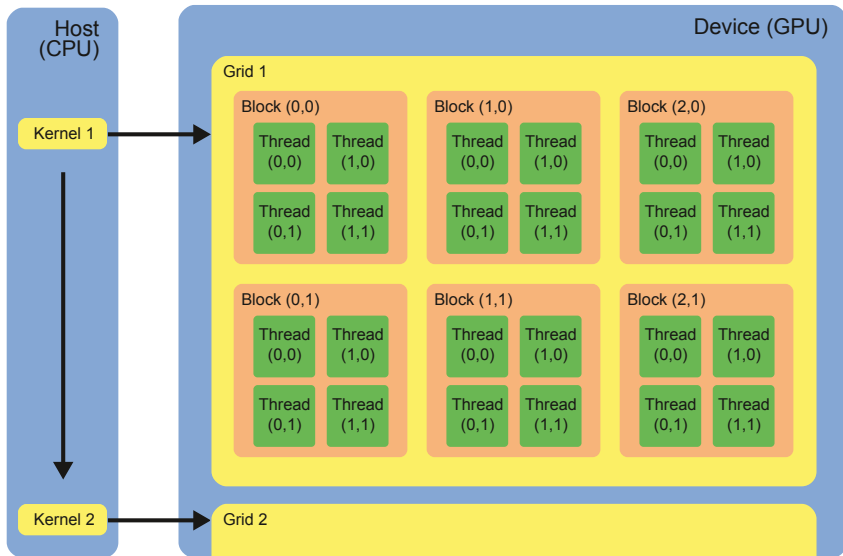
- 240 streaming processor cores
- 1.3 GHz clock frequency
- single precision peak performance 933 GFLOP/s
- double precision peak performance 78 GFLOP/s
- 4 GB GDDR3 RAM
- memory bandwidth 102 GB/s



NVIDIA architecture



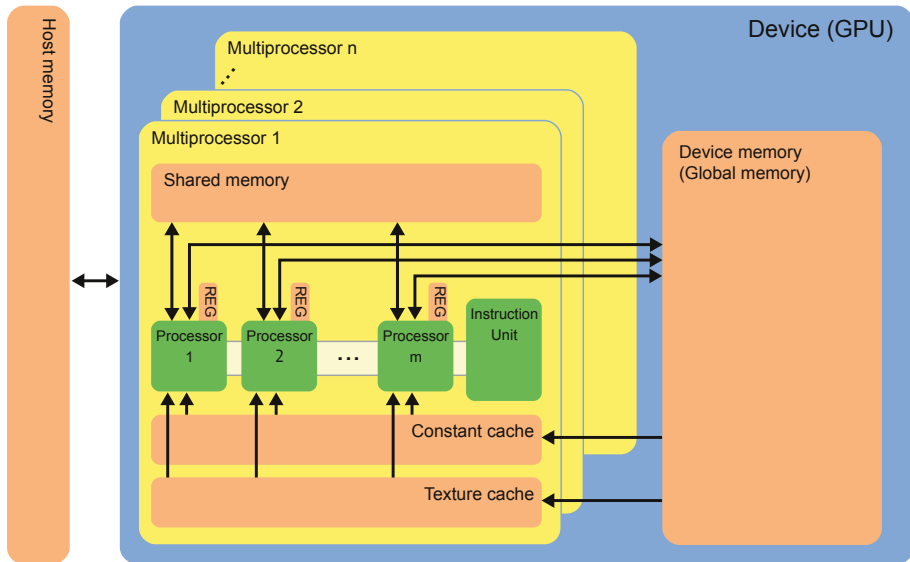
NVIDIA architecture



Compute model:

- all GPU calculations are encapsulated in dedicated functions (“kernels”)
- two-level hierarchy of a “grid” of thread “blocks”
- mixture of vector and parallel computer:
 - different threads execute the same code on different data (branching possible)
 - different blocks run independently
- threads on the same multiprocessor communicate via shared memory, different blocks are not meant to communicate
- atomic operations available
- coalescence of memory accesses

NVIDIA architecture



Memory layout:

- *Registers*: each multiprocessor is equipped with several thousand registers with local, zero-latency access
- *Shared memory*: processors of a multiprocessor have access a small amount (16 KB for Tesla, 48 KB for Fermi) of on chip, very small latency shared memory
- *Global memory*: large amount (currently up to 4 GB) of memory on separate DRAM chips with access from every thread on each multiprocessor with a latency of several hundred clock cycles
- *Constant and texture memory*: read-only memories of the same speed as global memory, but cached
- *Host memory*: cannot be accessed from inside GPU functions, relatively slow transfers

Design goals:

- a large degree of locality of the calculations, reducing the need for communication between threads
- a large coherence of calculations with a minimum occurrence of divergence of the execution paths of different threads
- a total number of threads significantly exceeding the number of available processing units
- a large overhead of arithmetic operations and shared memory accesses over global memory accesses

Spin models

Consider classical spin models with nn interactions, in particular

Ising model

$$\mathcal{H} = -J \sum_{\langle ij \rangle} s_i s_j + H \sum_i s_i, \quad s_i = \pm 1$$

Heisenberg model

$$\mathcal{H} = -J \sum_{\langle ij \rangle} \vec{S}_i \cdot \vec{S}_j + \vec{H} \cdot \sum_i \vec{S}_i, \quad |\vec{S}_i| = 1$$

Edwards-Anderson spin glass

$$\mathcal{H} = - \sum_{\langle ij \rangle} J_{ij} s_i s_j, \quad s_i = \pm 1$$

Random number generation

Monte Carlo simulations need (many) random numbers:

- want to update the spins on a lattice with many threads in parallel \Rightarrow each thread need its own instance of RNG
- stream of random numbers for each thread must be uncorrelated from all the others, i.e.
 - all threads work on the same global sequence *or*
 - threads produce non-overlapping sub-sequences
- should store everything in shared memory \Rightarrow need RNGs with small state

Linear congruential generators

Simplest choice satisfying these requirements is linear congruential generator (LCG):

$$x_{n+1} = (ax_n + b) \bmod m.$$

- for $m = 2^{32}$ or $2^{32} - 1$, the maximal period is of the order $p \approx m \approx 10^9$, much too short for large-scale simulations
- one should actually use at most \sqrt{p} numbers of the sequence
- for $m = 2^{32}$, modulo can be implemented as overflow, but then period of lower rank bits is only 2^k
- has poor statistical properties, e.g., k -tuples of (normalized) numbers lie on hyper-planes
- state is just 4 bytes per thread
- can easily skip ahead via $x_{n+t} = (a_t x_n + b_t) \bmod m$, where $a_t = a^t \bmod m$, $b_t = \sum_{i=1}^t a^i c \bmod m$.
- could be improved by choosing $m = 2^{64}$ and truncation to 32 most significant bits, period $p = m \approx 10^{18}$ and 8 bytes per thread

Lagged Fibonacci generators

Better statistical properties are achieved with lagged Fibonacci generators, for instance, of the three-term form

$$x_n = \alpha x_{n-r} \oplus \beta x_{n-s} \bmod m.$$

- good results with $\oplus = +$ and sufficiently large lags
- long period $p = 2^{31}(2^r - 1)$ (e.g., $p \approx 10^{400}$ for $r = 1279$)
- can be implemented exactly in floating-point arithmetic

$$u_n = \alpha u_{n-r} + \beta u_{n-s} \bmod 1.$$

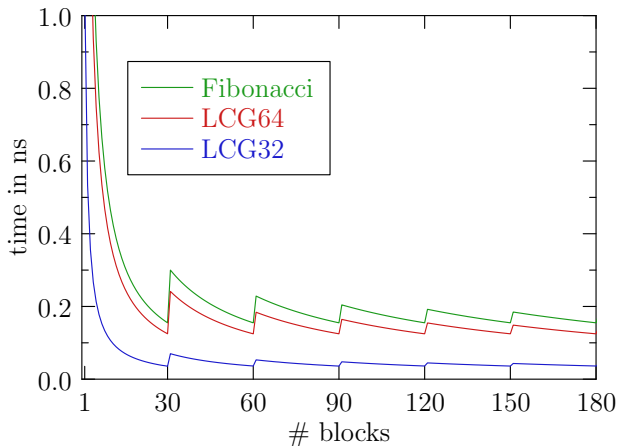
- if number of threads T is less than smaller lag, s numbers out of the same sequence can be generated in one vectorized call
- requires s/T words of local storage (in single precision)
- seeding for independent sub-sequences is possible

Other generators

There is, of course, a wealth of further generators which might be suitable

- shift-register or Tausworthe generators \Rightarrow special case of generalized lagged Fibonacci generators
- RANLUX generator \Rightarrow expensive, but good
- Mersenne twister has good properties, but requires large state per thread (624 words)
- ...

Performance



maximum speed-up between 70 and 100

Metropolis simulations

Computations need to be organized to suit the GPU layout for maximum performance:

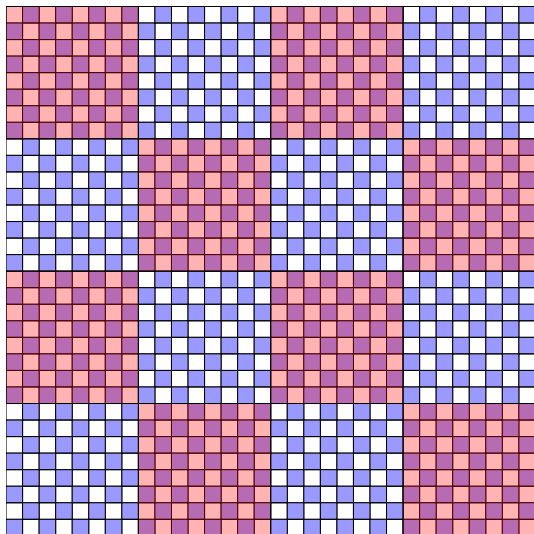
- many threads/blocks required to alleviate memory latencies
- per block shared memory much faster than global memory \Rightarrow define sub-problems that fit in to the (currently 16KB) shared memory
- threads within a block can synchronize and reside on the same processor, blocks must be independent (arbitrary execution order)

Consequences for (Metropolis) simulations:

- best to use an independent RNG per thread \Rightarrow need to make sure that sequences are uncorrelated
- divide system into tiles that can be worked on independently \Rightarrow level-1 checkerboard
- each tile should fit into shared memory
- divide tile (again) in checkerboard fashion for parallel update with different threads \Rightarrow level-2 checkerboard

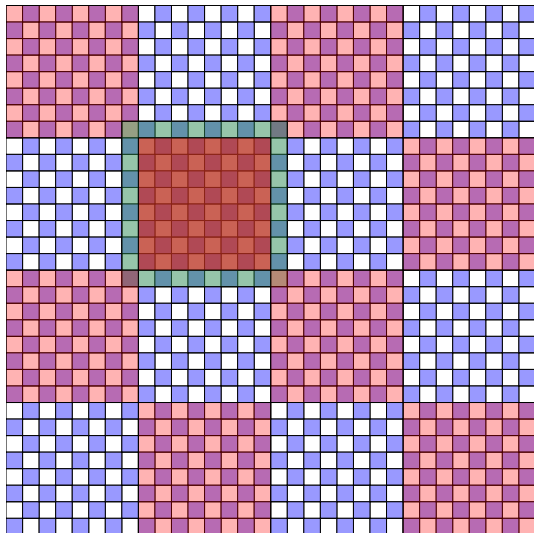
Checkerboard decomposition

- (red) large tiles: thread blocks
- (red) small tiles: individual threads
- load one large tile (plus boundary) into shared memory
- perform several spin updates per tile



Checkerboard decomposition

- (red) large tiles: thread blocks
- (red) small tiles: individual threads
- load one large tile (plus boundary) into shared memory
- perform several spin updates per tile



Implementation

```
__global__ void metro_sublattice_shared(spin_t *s, int *ranvec, int offset)
{
    int n = threadIdx.y*(BLOCKL/2)+threadIdx.x;

    int xoffset = (2*blockIdx.x+(blockIdx.y+offset)%2)*BLOCKL;
    int yoffset = blockIdx.y*BLOCKL;

    __shared__ spin_t sS[(BLOCKL+2)*(BLOCKL+2)];

    sS[(threadIdx.y+1)*(BLOCKL+2)+2*threadIdx.x+1] = s[(yoffset+threadIdx.y)*L+xoffset+2*threadIdx.x];
    sS[(threadIdx.y+1)*(BLOCKL+2)+2*threadIdx.x+2] = s[(yoffset+threadIdx.y)*L+xoffset+2*threadIdx.x+1];
    if(threadIdx.y == 0) {
        if(blockIdx.y == 0) {
            sS[2*threadIdx.x+1] = s[(L-1)*L+xoffset+2*threadIdx.x];
            sS[2*threadIdx.x+2] = s[(L-1)*L+xoffset+2*threadIdx.x+1];
        } else {
            sS[2*threadIdx.x+1] = s[(yoffset-1)*L+xoffset+2*threadIdx.x];
            sS[2*threadIdx.x+2] = s[(yoffset-1)*L+xoffset+2*threadIdx.x+1];
        }
    }
    if(threadIdx.y == BLOCKL-1) {
        if(blockIdx.y == GRIDL-1) {
            sS[(BLOCKL+1)*(BLOCKL+2)+2*threadIdx.x+1] = s[xoffset+2*threadIdx.x];
            sS[(BLOCKL+1)*(BLOCKL+2)+2*threadIdx.x+2] = s[xoffset+2*threadIdx.x+1];
        } else {
            sS[(BLOCKL+1)*(BLOCKL+2)+2*threadIdx.x+1] = s[(yoffset+BLOCKL)*L+xoffset+2*threadIdx.x];
            sS[(BLOCKL+1)*(BLOCKL+2)+2*threadIdx.x+2] = s[(yoffset+BLOCKL)*L+xoffset+2*threadIdx.x+1];
        }
    }
    if(threadIdx.x == 0) {
        if(xoffset == 0) sS[(threadIdx.y+1)*(BLOCKL+2)] = s[(yoffset+threadIdx.y)*L+(L-1)];
        else sS[(threadIdx.y+1)*(BLOCKL+2)] = s[(yoffset+threadIdx.y)*L+xoffset-1];
    }
    if(threadIdx.x == BLOCKL/2-1) {
        if(xoffset == L-BLOCKL) sS[(threadIdx.y+1)*(BLOCKL+2)+BLOCKL+1] = s[(yoffset+threadIdx.y)*L];
        else sS[(threadIdx.y+1)*(BLOCKL+2)+BLOCKL+1] = s[(yoffset+threadIdx.y)*L+xoffset+BLOCKL];
    }
}
```

Implementation

```
__shared__ int ranvecS[THREADS];
ranvecS[n] = ranvec[(blockIdx.y*(GRIDL/2)+blockIdx.x)*THREADS+n];

__syncthreads();

int x1 = (threadIdx.y%2)+2*threadIdx.x;
int x2 = ((threadIdx.y+1)%2)+2*threadIdx.x;
int y = threadIdx.y;

for(int i = 0; i < SWEEPS_LOCAL; ++i) {
    int ide = sS(x1,y)*(sS(x1-1,y)+sS(x1+1,y)+sS(x1,y-1)+sS(x1,y+1));
    if(ide <= 0 || RAN(ranvecS[n]) < boltzD[ide]) sS(x1,y) = -sS(x1,y);

    __syncthreads();

    ide = sS(x2,y)*(sS(x2-1,y)+sS(x2+1,y)+sS(x2,y-1)+sS(x2,y+1));
    if(ide <= 0 || RAN(ranvecS[n]) < boltzD[ide]) sS(x2,y) = -sS(x2,y);

    __syncthreads();
}

s[(yoffset+threadIdx.y)*L+xoffset+2*threadIdx.x] = sS[(threadIdx.y+1)*(BLOCKL+2)+2*threadIdx.x+1];
s[(yoffset+threadIdx.y)*L+xoffset+2*threadIdx.x+1] = sS[(threadIdx.y+1)*(BLOCKL+2)+2*threadIdx.x+2];
ranvec[(blockIdx.y*(GRIDL/2)+blockIdx.x)*THREADS+n] = ranvecS[n];
}
```


How to assess performance?

- what to compare to (one CPU core, whole CPU, SMP system, ...) here: Tesla C1060 vs. Intel QuadCore (Yorkfield) @ 3.0 GHz/6 MB
- for really fair comparison: optimize CPU code for cache alignment, use SSE instructions etc.
- ignore measurements, since spin flips per μs , (ns, ps) is well-established unit for spin systems

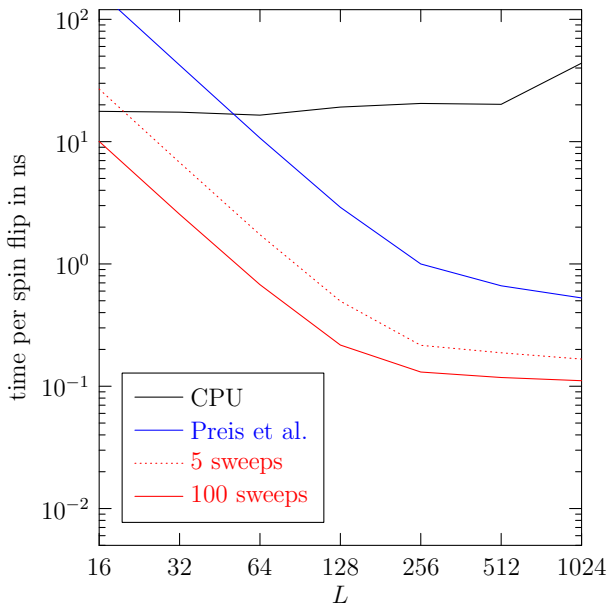
How to assess performance?

- what to compare to (one CPU core, whole CPU, SMP system, ...) here: Tesla C1060 vs. Intel QuadCore (Yorkfield) @ 3.0 GHz/6 MB
- for really fair comparison: optimize CPU code for cache alignment, use SSE instructions etc.
- ignore measurements, since spin flips per μs , (ns, ps) is well-established unit for spin systems

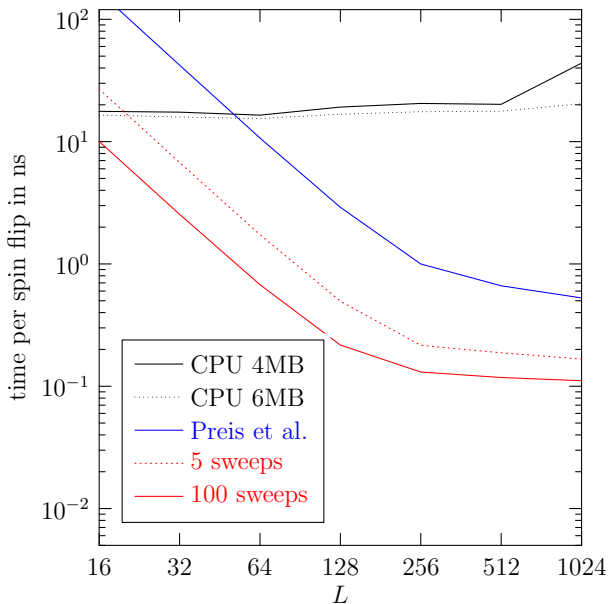
Example: Metropolis simulation of 2D Ising system

- use 32-bit linear congruential generator
- no neighbor table since integer multiplies and adds are very cheap (4 instructions per clock cycle and processor)
- need to play with tile sizes to achieve best throughput

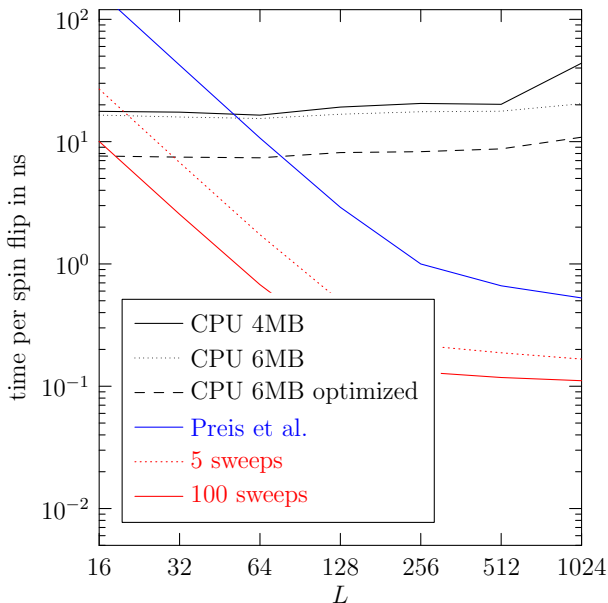
2D Ising ferromagnet



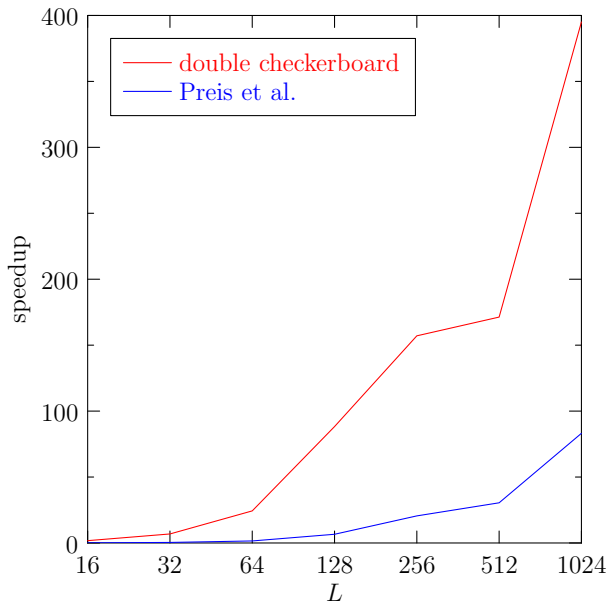
2D Ising ferromagnet



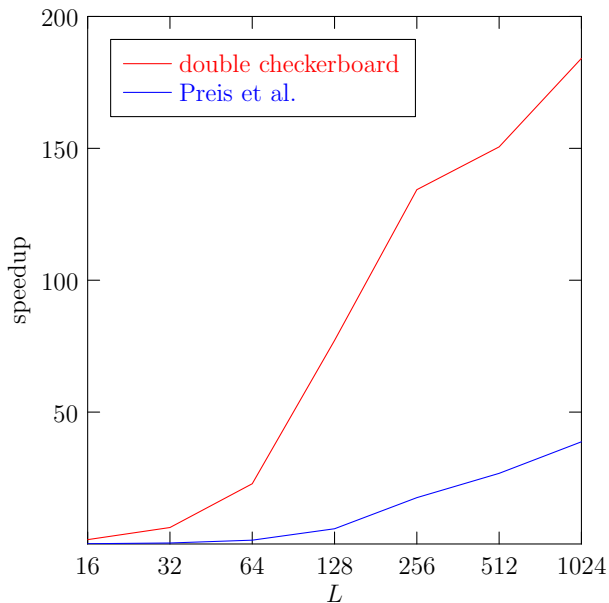
2D Ising ferromagnet



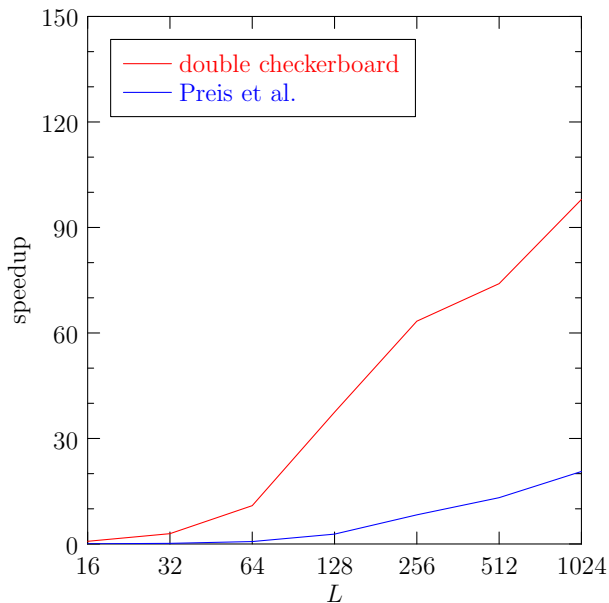
2D Ising ferromagnet



2D Ising ferromagnet



2D Ising ferromagnet

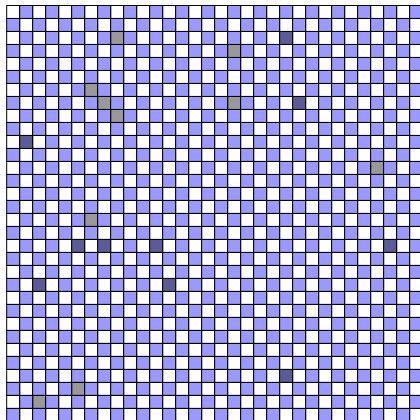


Is it correct?

Detailed balance,

$$T(\{s'_i\} \rightarrow \{s_i\})p_\beta(\{s'_i\}) = T(\{s_i\} \rightarrow \{s'_i\})p_\beta(\{s_i\}),$$

only holds for *random* updates.



Is it correct?

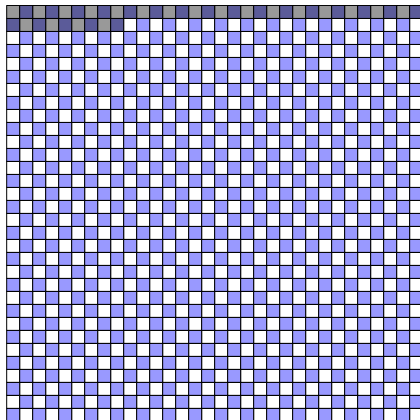
Detailed balance,

$$T(\{s'_i\} \rightarrow \{s_i\})p_\beta(\{s'_i\}) = T(\{s_i\} \rightarrow \{s'_i\})p_\beta(\{s_i\}),$$

only holds for *random* updates.

Usually applied sequential update
merely satisfies (global) *balance*,

$$\sum_{\{s'_i\}} T(\{s'_i\} \rightarrow \{s_i\})p_\beta(\{s'_i\}) =$$
$$\sum_{\{s'_i\}} T(\{s_i\} \rightarrow \{s'_i\})p_\beta(\{s_i\}) = p_\beta(\{s_i\})$$



Is it correct?

Detailed balance,

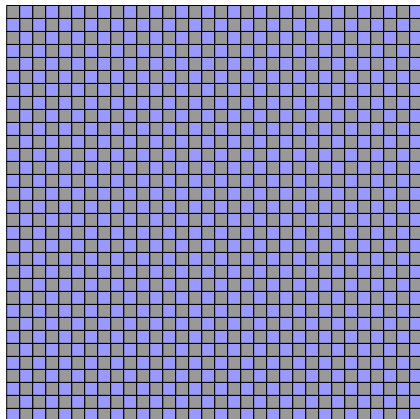
$$T(\{s'_i\} \rightarrow \{s_i\})p_\beta(\{s'_i\}) = T(\{s_i\} \rightarrow \{s'_i\})p_\beta(\{s_i\}),$$

only holds for *random* updates.

Usually applied sequential update
merely satisfies (global) *balance*,

$$\sum_{\{s'_i\}} T(\{s'_i\} \rightarrow \{s_i\})p_\beta(\{s'_i\}) = \sum_{\{s'_i\}} T(\{s_i\} \rightarrow \{s'_i\})p_\beta(\{s_i\}) = p_\beta(\{s_i\})$$

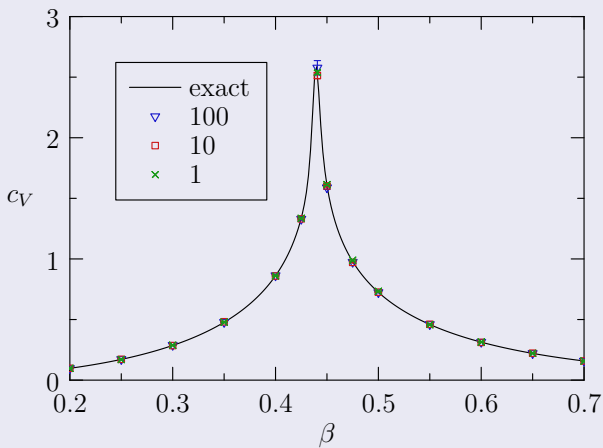
Similarly for checkboard update. Could
restore detailed balance on the level of
several sweeps, though:



$$AAAA(M)AAAABBBB(M)BBBBAAAA(M)AAAABBBB(M)BBBB \dots$$

A closer look

Comparison to exact results:



A closer look

Random number generators: significant deviations from exact result for test case of 1024×1024 system at $\beta = 0.4$, 10^7 sweeps

- checkerboard update uses random numbers in different way than sequential update
- linear congruential generators can skip ahead: “right” way uses non-overlapping sub-sequences
- “wrong” way uses sequences from random initial seeds, many of which must overlap

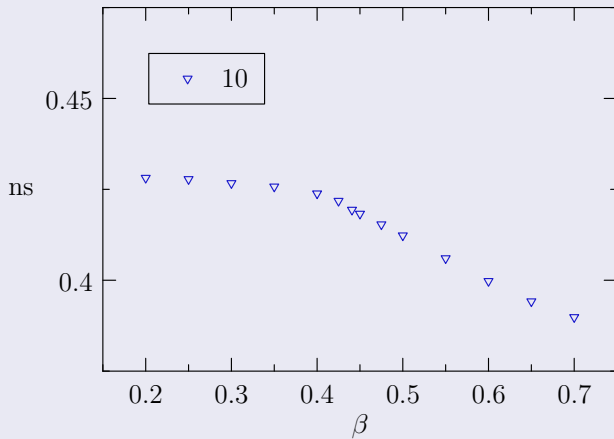
A closer look

Random number generators: significant deviations from exact result for test case of 1024×1024 system at $\beta = 0.4$, 10^7 sweeps

method	e	$\frac{E-E_{\text{exact}}}{\sigma(E)}$	C_V	$\frac{C_V-C_{V,\text{exact}}}{\sigma(C_V)}$
exact	1.106079207	0	0.8616983594	0
sequential update				
LCG32	1.1060788(15)	-0.26	0.83286(45)	-63.4
LCG64	1.1060800(17)	0.49	0.86101(59)	-1.14
Fibonacci	1.1060788(17)	-0.18	0.86131(59)	-0.64
checkerboard update				
LCG32, sequential	1.0910932(13)	-11451	0.73668(43)	-288
LCG32, random	1.1060776(17)	0.92	0.86162(62)	-0.12
LCG64, sequential	1.106066(16)	-0.81	5.321(14)	310
LCG64, random	1.1060776(17)	0.92	0.86162(62)	-0.12
Fibonacci	1.1060799(19)	0.39	0.86083(52)	-1.64

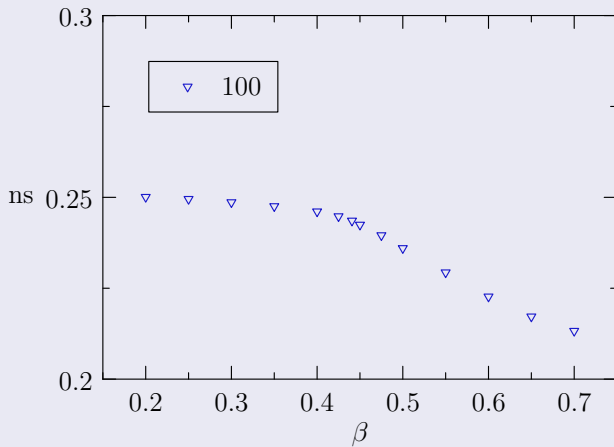
A closer look

Temperature dependent spin flip times:



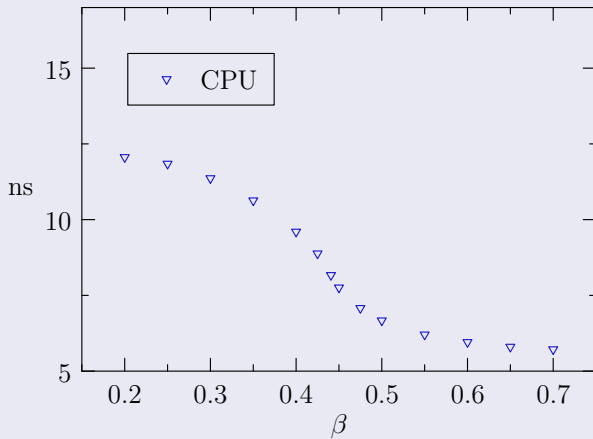
A closer look

Temperature dependent spin flip times:



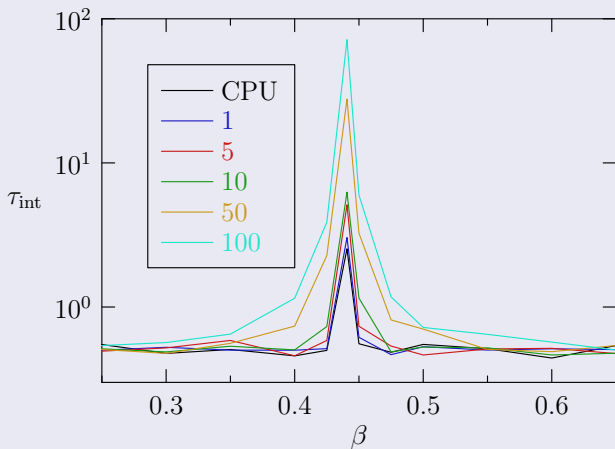
A closer look

Temperature dependent spin flip times:



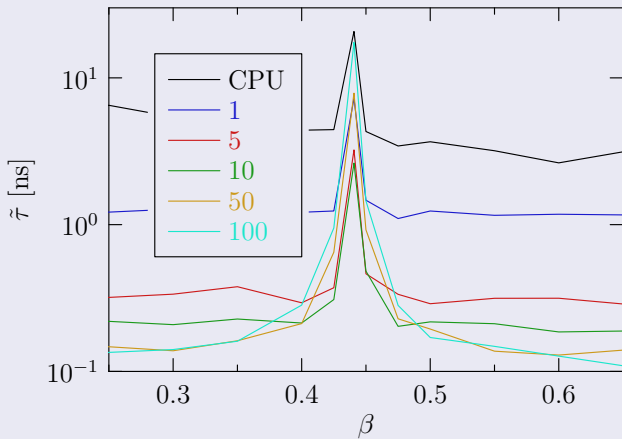
A closer look

Autocorrelation times:



A closer look

Real time to create independent spin configuration:



Heisenberg model

Maximum performance around 100 ps per spin flip for Ising model (vs. around 10 ns on CPU). What about continuous spins, i.e., float instead of int variables?

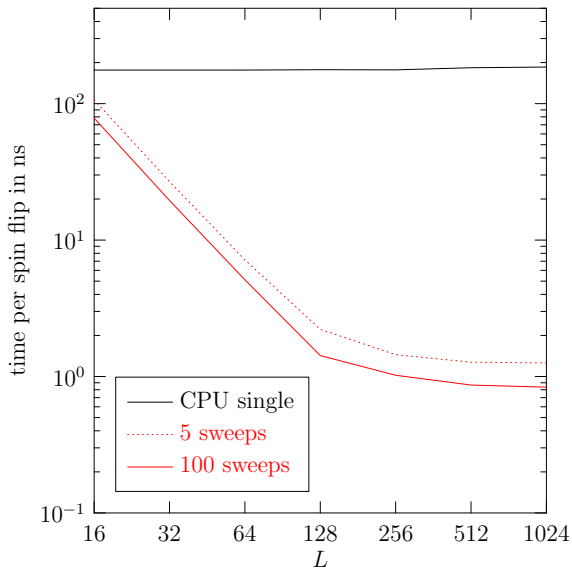
Heisenberg model

Maximum performance around 100 ps per spin flip for Ising model (vs. around 10 ns on CPU). What about continuous spins, i.e., float instead of int variables?

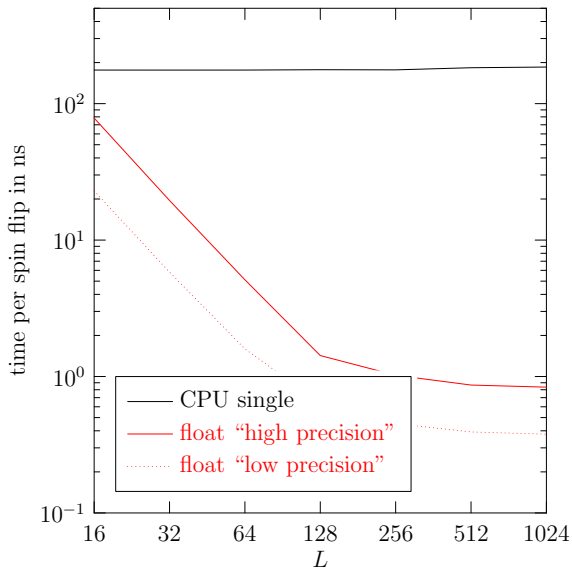
⇒ use same decomposition, but now floating-point computations are dominant:

- CUDA is not 100% IEEE compliant
- single-precision computations are supposed to be fast, double precision (supported since recently) much slower
- for single precision, normal (“high precision”) and extra-fast, device-specific versions of sin, cos, exp etc. are provided

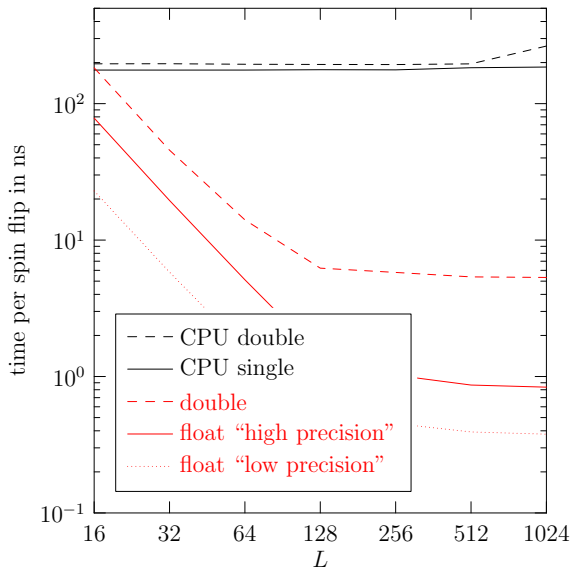
Heisenberg model: performance



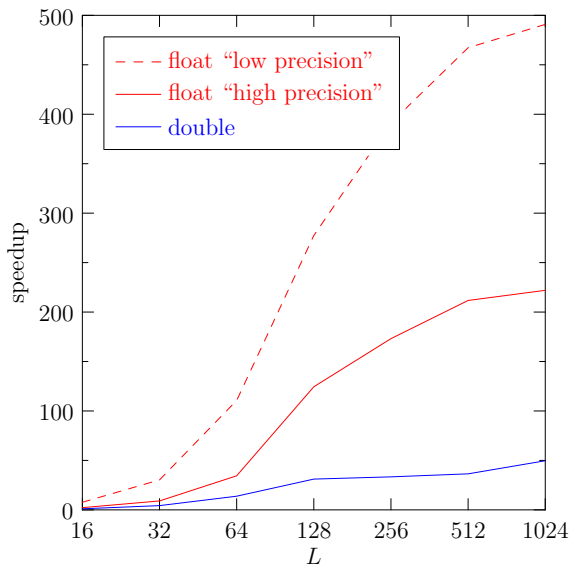
Heisenberg model: performance



Heisenberg model: performance



Heisenberg model: performance



Heisenberg model: stability

Performance results:

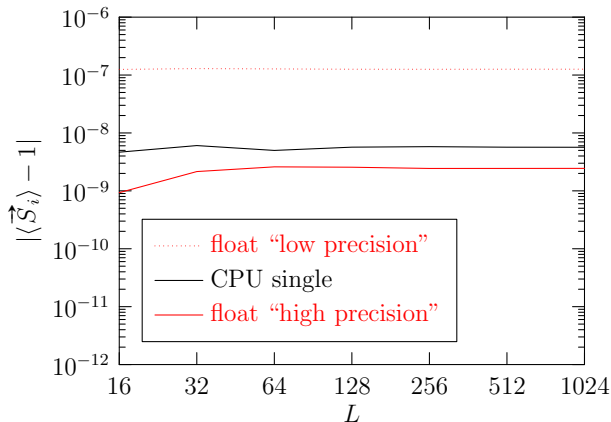
- CPU: 185 ns (single) resp. 264 (double) per spin flip
- GPU: 0.8 ns (single), 0.4 ns (fast single) resp. 5.3 ns (double) per spin flip

Heisenberg model: stability

Performance results:

- CPU: 185 ns (single) resp. 264 (double) per spin flip
- GPU: 0.8 ns (single), 0.4 ns (fast single) resp. 5.3 ns (double) per spin flip

How about stability?

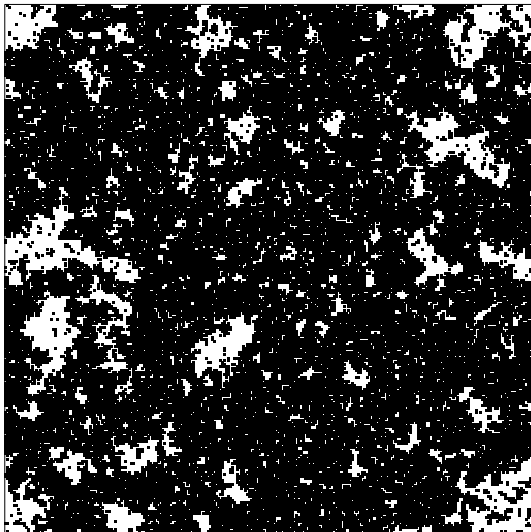


Cluster algorithms

Would need to use cluster algorithms for efficient equilibrium simulation of spin models at criticality:

- 1 Activate bonds between like spins with probability $p = 1 - e^{-2\beta J}$.
- 2 Construct (Swendsen-Wang) spin clusters from domains connected by active bonds.
- 3 Flip independent clusters with probability 1/2.
- 4 Goto 1.

Critical configuration



Cluster algorithms

Would need to use cluster algorithms for efficient equilibrium simulation of spin models at criticality:

- 1 Activate bonds between like spins with probability $p = 1 - e^{-2\beta J}$.
- 2 Construct (Swendsen-Wang) spin clusters from domains connected by active bonds.
- 3 Flip independent clusters with probability 1/2.
- 4 Goto 1.

Steps 1 and 3 are local \Rightarrow Can be efficiently ported to GPU.

What about step 2? \Rightarrow Domain decomposition into tiles.

labeling *inside* of domains

- Hoshen-Kopelman
- breadth-first search
- self-labeling
- union-find algorithms

relabeling *across* domains

- self-labeling
- hierarchical approach
- iterative relaxation

BFS or Ants in the Labyrinth

56	57	58	59	60	61	62	63
48	49	50	51	52	53	54	55
40	41	42	43	19	45	46	47
32	33	19	19	19	19	38	39
24	25	26	19	19	29	30	31
16	17	18	19	20	21	22	23
8	9	10	11	12	13	14	15
0	1	2	3	4	5	6	7

BFS or Ants in the Labyrinth

56	57	58	59	60	61	62	63
48	49	50	51	52	53	54	55
40	41	42	43	19	45	46	47
32	33	19	19	19	19	38	39
24	25	26	19	19	29	30	31
16	17	18	19	20	21	22	23
8	9	10	11	12	13	14	15
0	1	2	3	4	5	6	7

only wave-front vectorization would be possible \Rightarrow many idle threads

Self-labeling

56	57	58	59	60	61	62	63
48	49	50	51	52	53	54	55
40	41	42	43	44	45	46	47
32	33	34	35	36	37	38	39
24	25	26	27	28	29	30	31
16	17	18	19	20	21	22	23
8	9	10	11	12	13	14	15
0	1	2	3	4	5	6	7

Diagram illustrating a self-labeling process on an 8x8 grid. The grid contains numbers 0 to 63. Red arrows indicate connections between adjacent cells (horizontally and vertically). A yellow box highlights a 2x2 subgrid centered around the number 19, which is labeled with a red '19' above it. The number 20 is also labeled with a red '19' above it.

Self-labeling

56	57	58	59	60	61	62	63
48	49	50	51	52	53	54	55
40	41	42	43	44	45	46	47
32	33	34	35	36	37	38	39
24	25	26	27	28	29	30	31
16	17	18	19	20	21	22	23
8	9	10	11	12	13	14	15
0	1	2	3	4	5	6	7

effort is $O(L^3)$ at the critical point, but can be vectorized with $O(L^2)$ threads

Union-find

56	57	58	59	60	61	62	63
48	41	41	51	52	53	54	55
40	32	41	41	19	45	46	47
32	32	19	30	30	30	38	39
24	25	26	19	30	30	13	31
16	17	18	19	20	21	13	23
8	9	10	11	12	13	13	15
0	1	2	3	4	5	6	7

Union-find

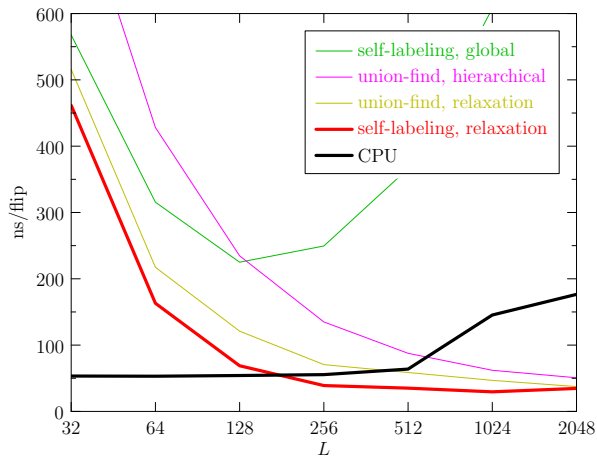
56	57	58	59	60	61	62	63
48	41	41	51	52	53	54	55
40	32	41	41	19	45	46	47
32	32	19	30	30	30	38	39
24	25	26	19	30	30	13	31
16	17	18	19	20	21	13	23
8	9	10	11	12	13	13	15
0	1	2	3	4	5	6	7

tree structure with two optimizations:

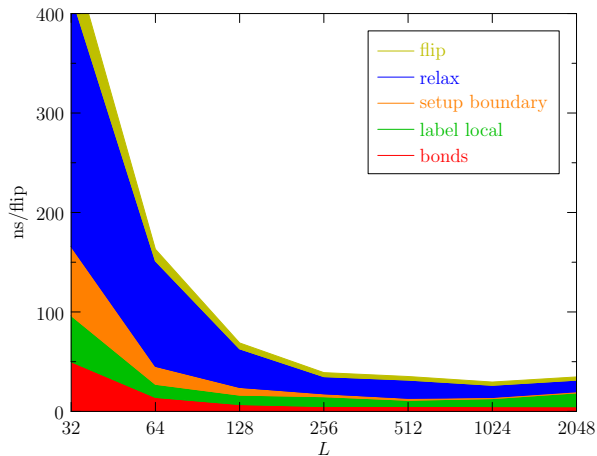
- balanced trees
- path compression

⇒ root finding and cluster union
essentially $O(1)$ operations

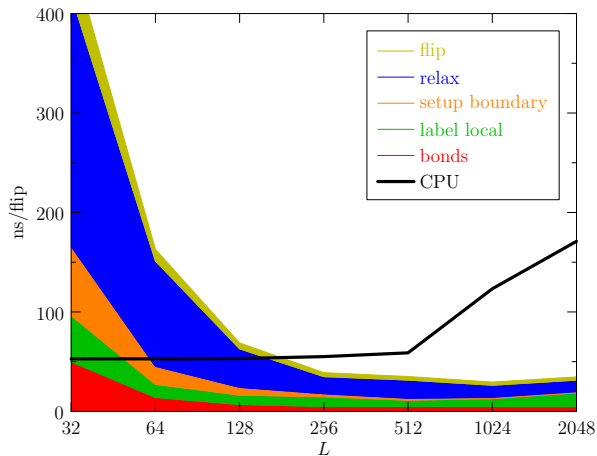
Performance



Performance



Performance



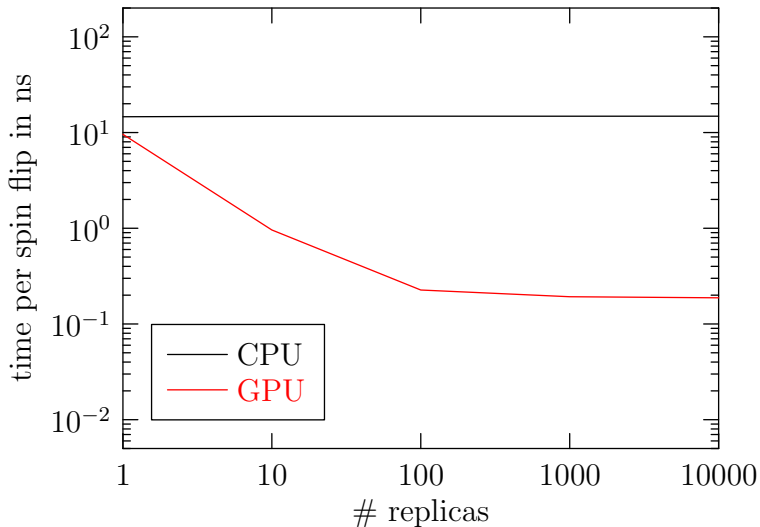
Problems with cluster labeling on GPU:

- overhead from parallelization (relaxation steps)
- lack of thread-level parallelism
- idle threads in hierarchical schemes
- best performance about 29 ns per spin flip, improvements possible
- problems *not* due to type of computations: 2.9 ns per spin flip for SW simulations of several systems in parallel

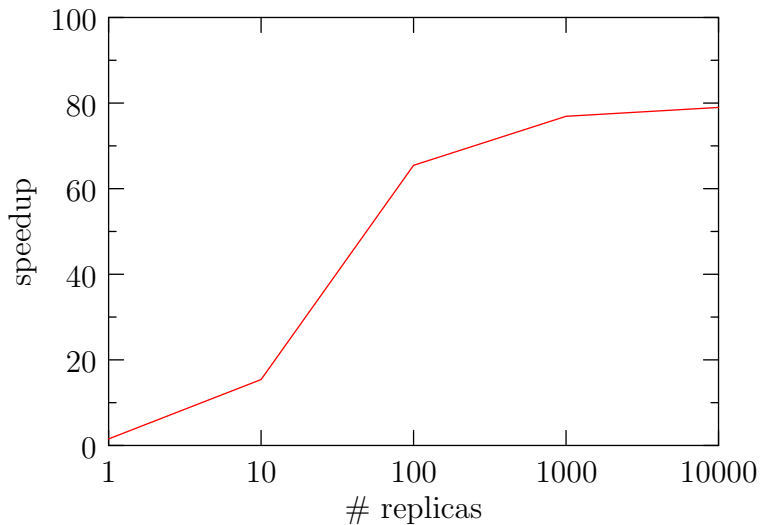
Simulate Edwards-Anderson model on GPU:

- same domain decomposition (checkerboard)
- slightly bigger effort due to non-constant couplings
- higher performance due to larger independence?
- very simple to combine with parallel tempering

Spin glass: performance



Spin glass: performance



Spin glasses: continued

Seems to work well with

- 15 ns per spin flip on CPU
- 180 ps per spin flip on GPU

but not better than ferromagnetic Ising model.

Spin glasses: continued

Seems to work well with

- 15 ns per spin flip on CPU
- 180 ps per spin flip on GPU

but not better than ferromagnetic Ising model.

Further improvement: use multi-spin coding

- Synchronous multi-spin coding: different spins in a single configurations in one word
- Asynchronous multi-spin coding: spins from different realizations in one word

Implementation

```
for(int i = 0; i < SWEEPS_LOCAL; ++i) {
    float r = RAN(ranvecS[n]);
    if(r < boltzD[4]) sS(x1,y) = ~sS(x1,y);
    else {
        p1 = JSx(x1m,y) ^ sS(x1,y) ^ sS(x1m,y); p2 = JSx(x1,y) ^ sS(x1,y) ^ sS(x1p,y);
        p3 = JSy(x1,ym) ^ sS(x1,y) ^ sS(x1,ym); p4 = JSy(x1,y) ^ sS(x1,y) ^ sS(x1,yp);
        if(r < boltzD[2]) {
            ido = p1 | p2 | p3 | p4;
            sS(x1,y) = ido ^ sS(x1,y);
        } else {
            ido1 = p1 & p2; ido2 = p1 ^ p2;
            ido3 = p3 & p4; ido4 = p3 ^ p4;
            ido = ido1 | ido3 | (ido2 & ido4);
            sS(x1,y) = ido ^ sS(x1,y);
        }
    }
}

__syncthreads();

r = RAN(ranvecS[n]);
if(r < boltzD[4]) sS(x2,y) = ~sS(x2,y);
else {
    p1 = JSx(x2m,y) ^ sS(x2,y) ^ sS(x2m,y); p2 = JSx(x2,y) ^ sS(x2,y) ^ sS(x2p,y);
    p3 = JSy(x2,ym) ^ sS(x2,y) ^ sS(x2,ym); p4 = JSy(x2,y) ^ sS(x2,y) ^ sS(x2,yp);
    if(r < boltzD[2]) {
        ido = p1 | p2 | p3 | p4;
        sS(x2,y) = ido ^ sS(x2,y);
    } else {
        ido1 = p1 & p2; ido2 = p1 ^ p2;
        ido3 = p3 & p4; ido4 = p3 ^ p4;
        ido = ido1 | ido3 | (ido2 & ido4);
        sS(x2,y) = ido ^ sS(x2,y);
    }
}

__syncthreads();
}
```

Spin glasses: continued

Seems to work well with

- 15 ns per spin flip on CPU
- 180 ps per spin flip on GPU

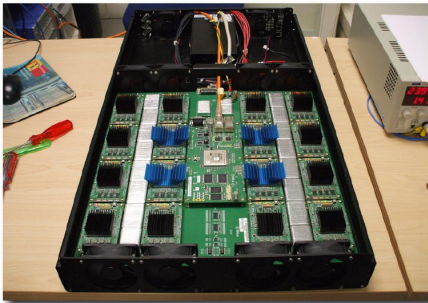
but not better than ferromagnetic Ising model.

Further improvement: use multi-spin coding

- Synchronous multi-spin coding: different spins in a single configurations in one word
- Asynchronous multi-spin coding: spins from different realizations in one word

⇒ brings us down to about 15 ps per spin flip

JANUS, a modular massively parallel and reconfigurable FPGA-based computing system.



JANUS, a modular massively parallel and reconfigurable FPGA-based computing system.

MODEL	Algorithm	JANUS		PC		
		Max size	perfs	AMSC	SMSC	NO MSC
3D Ising EA	Metropolis	96^3	16 ps	45×	190×	
3D Ising EA	Heat Bath	96^3	16 ps	60×		
$Q = 4$ 3D Glassy Potts	Metropolis	16^3	64 ps	1250×	1900×	
$Q = 4$ 3D disordered Potts	Metropolis	88^3	32 ps	125×		1800×
$Q = 4$, $C_m = 4$ random graph	Metropolis	24000	2.5 ns	2.4×		10×

JANUS, a modular massively parallel and reconfigurable FPGA-based computing system.

MODEL	Algorithm	JANUS		PC		
		Max size	perfs	AMSC	SMSC	NO MSC
3D Ising EA	Metropolis	96^3	16 ps	45×	190×	
3D Ising EA	Heat Bath	96^3	16 ps	60×		
$Q = 4$ 3D Glassy Potts	Metropolis	16^3	64 ps	1250×	1900×	
$Q = 4$ 3D disordered Potts	Metropolis	88^3	32 ps	125×		1800×
$Q = 4, C_m = 4$ random graph	Metropolis	24000	2.5 ns	2.4×		10×

Costs:

- Janus: 256 units, total cost about 700,000 Euros
- Same performance with GPU: 64 PCs (2000 Euros) with 2 GTX 295 cards (500 Euros) \Rightarrow 200,000 Euros
- Same performance with CPU only (assuming a speedup of ~ 50): 800 blade servers with two dual Quadcore sub-units (3500 Euros) \Rightarrow 2,800,000 Euros

Conclusions:

- GPGPU promises significant speedups at moderate coding effort
 - Requirements for good performance:
 - large degree of locality \Rightarrow domain decomposition
 - suitability for parallelization (blocks) *and* vectorization (threads)
 - total number of threads much larger than processing units (memory latency)
 - opportunity for using shared memory \Rightarrow performance is memory limited
 - ideally continuous variables
- \Rightarrow maximum speed-up ~ 500
- effort significantly smaller than for special-purpose machines
 - GPGPU might be a fashion, but CPU computing goes the same way

Outlook:

- parallel tempering no problem
- generalized-ensemble techniques might be more of a challenge
- ideal application:
 - spin glasses, disordered systems
 - dynamical MC: nucleation, domain growth