



Massively Parallel Computer Architecture High-Performance Microprocessor Architecture and Programming

Hiroshi Nakashima
(ACCMS, Kyoto University)



Contents

- **Three Major Weapons of High-Performance Microprocessors**
 - **Instruction Pipeline, ILP, and Out-of-Order Execution**
 - basic mechanism of instruction execution
 - hazards caused by dependency and means to cope with them
 - **Branch Prediction and Speculation**
 - predicting branch address and direction
 - means to cope with misprediction
 - **Cache and TLB**
 - basics of configuration
 - cache/TLB-awareness and improvement

2



Overview (1/4)

■ Factor to Determine Performance

- First-level Approximation = #basic-ops (≈machine insts) × 1 / frequency

```
for (x=0; x<nx; x++)
  un[y][x]=u[y][x]+
    (dt/(h*h)) * (u[y][x+1]+u[y][x-1]+
      u[y+1][x]+u[y-1][x]-
      4*u[y][x]);
```

FP register

GP (int) register

- $f=u[\dots] \dots 5$, $un[\dots]=f \dots 1$
- $f=f+f \dots 5$, $f=f*f \dots 3$, $f=f/f \dots 1$
- $g++ \dots 1$, $g<g \dots 1$, $if()goto \dots 1$

→ 18 / frequency?
→ supercomputer in KU (18core Xeon)
≈ 7 / frequency (why so small?)

3



Overview (2/4)

■ Why "#ops / freq" Overestimates = Parallel Execution of Multiple Insts

- in KU's supercomputer
 - $f=mem \times 2$
 - $mem=f \times 1$
 - $f=f+f \times 1$
 - $\{f=f*f, f=f/f\} \times 2$ (except for $f=f/f \times 2$)
 - integer insts $\ni \{g++, g<g, if()goto\} \times 4$

can be executed in parallel

in reality, we have some other restrictions

→ $\max(\{f=mem\} \times 5/2, \{mem=f\} \times 1/1,$
 $\{f=f+f\} \times 5/1, \{f=f*f\} \times 3+ \{f=f/f\} \times 1/2,$
 $\{g++\} \times 1+ \{g<g\} \times 1+ \{if()\} \times 1/4)$
 = 5? (< 7)

4



Overview (3/4)

■ Why max(...) Underestimates = Latency of Depending Instruction

- in KU's supercomputer (unit=cycle)
 - $f=mem \dots 5$ (+7+23+60?) ← 1st / 2nd / 3rd-level cache miss
 - $f=f+f \dots 3$, $f=f*f \dots 5$, $f=f/f \dots 14$ (!!)
 - $\{g++, g<g\} \dots 1$
 - $if()goto \dots 1$ (+15?+7+...)

and we have to wait 5-cycle for next divide

branch misprediction / 1st / ... cache miss

→ avoid division as much as possible!!

- c.f. other representative ops
 - $g=mem \dots 4$ (+...)
 - $g=g+g \dots 1$, $g=g*g \dots 3$, $g=g/g \dots 20\sim 26$

5



Overview (4/4)

Loop-invariant code motion is a fundamental optimization but not always (able to be) done by compilers.

■ Avoiding Division

```
dthh=dt/(h*h);
for (t=0; t<tmax; t++){
  for (y=0; y<ny; y++) for (x=0; x<nx; x++){
    un[y][x]=u[y][x]+
      dthh*(u[y][x+1]+u[y][x-1]+
        u[y+1][x]+u[y-1][x]-
        4*u[y][x]);
    tmp=un; un=u; u=tmp;
  }
}
```

When un, dt, h are global, a compiler might unsure that dt and h are invariant due to update of un[[]].

lesson-1: Do a simple optimization by yourself.
lesson-2: Don't use global variables in a loop if possible. (e.g., move them into local variables)

■ Performance of Modified Code in KU's Supercomputer ≈ 5~6 / frequency

- Reasons why >5 are shown later.
 - Many pitfalls to make >> 6
- Key issues are to avoid them

6

Massively Parallel Comp. Arch. (2) © 2010-2017 H. Nakashima

Basics of Instruction Pipeline (1/6)

- Instruction Pipeline**
 - Split instruction execution flow into multiple stages.
 - All stages work in parallel.
- In Textbooks ...**

→ (up to 5 (= #stages) instructions are executed in parallel.

- In Reality ...**
 - ≈20 stages (or more) and more complicated (as discussed later)

7

Massively Parallel Comp. Arch. (2) © 2010-2017 H. Nakashima

Basics of Instruction Pipeline (2/6)

- Part of DE-solver loop (assume latency=1)**

	#u[y][x+1]	#u[y][x-1]	#u[y][x+1]+ #u[y][x-1]	#u[y+1][x]	#u[y+1][x]
IF	I=PC++;	I=PC++;	I=PC++;	I=PC++;	I=PC++;
ID	A=guy;	A=guy;	A=fuzp;	A=guy;	A=fuzp;
EX	D=A+B+C;	D=A+B-C;	D=A+B;	D=A+B;	D=A+B;
MA	M=D;	M=D;	M=D;	M=D;	M=D;
WB	fuzp=M;	fuzm=M;	fuzp=M;	fuzp=M;	fuzp=M;

8

Massively Parallel Comp. Arch. (2) © 2010-2017 H. Nakashima

Basics of Instruction Pipeline (3/6)

- Hazard due to Instruction Dependence and Remedy for it**
 - true/flow dependence = output of predecessor is used by successor
 - RAW (read-after-write) hazard = generate/use-timing is reversed.

→bypassing / forwarding

avoid hazard by bypassing data through stages

9

Massively Parallel Comp. Arch. (2) © 2010-2017 H. Nakashima

Basics of Instruction Pipeline (4/6)

- Latency ≥ 2 → not always avoidable → Stall**
 - Even when latency of f=mem, f=f+f, f=f*f is 2 ...

cannot be bypassed

detect RAW hazard and stall

10

Massively Parallel Comp. Arch. (2) © 2010-2017 H. Nakashima

Basics of Instruction Pipeline (5/6)

- Remedy=Static Scheduling (1/2)**

anti-dependent: use → define dependency

16+8=24/iteration

11

Massively Parallel Comp. Arch. (2) © 2010-2017 H. Nakashima

Basics of Instruction Pipeline (6/6)

- Remedy=Static Scheduling (2/2)**

removing anti-dependence by register renaming

17+1=18/iteration

12

Massively Parallel Comp. Arch. (2) © 2010-2017 H. Nakashima

ILP (1/13)

- ILP: Instruction Level Parallelism**
 - Execution order of mutually independent instructions can be modified arbitrarily.
 - They can be executed **in parallel**
 - Execute multiple instructions **in parallel** with multiple function units (e.g. ALUs)
 - Variants**
 - super scalar**: execute ordinary instructions
 - static: in-order execution of compiled code (1990's)
 - dynamic: change order at run time (**out-of-order**, now)
 - VLIW**: one instruction performs multiple operations (e.g. Itanium)
 - SIMD**: one instruction operates on multiple data (SSEx: streaming SIMD extension, AVX)

Annotations: "scalar means 'not vector (processor)'" (pointing to super scalar); "very long inst. word" (pointing to super scalar); "single inst. multiple data" (pointing to SIMD).

13

Massively Parallel Comp. Arch. (2) © 2010-2017 H. Nakashima

ILP (2/13)

- KU's Supercomputer**
 - Issue Rate = 7** (#operations which can start at every cycle)
 - #Function Units**
 - {g, f} = mem × 2, mem = {g, f} × 1
 - {f = f * f [+ f], f = f + f} × 1, {f = f * f [+ f], f = f / f} × 1,
 - {g ++, g < g} × 2, {g ++, g < g, if() goto} × 2

Annotation: "Instructions per Cycle" (pointing to Issue Rate = 7)

→ **maximum execution throughput (IPC) = 7**
 → (16inst/7) / iteration = 2.3 cycle / iteration ?

14

Massively Parallel Comp. Arch. (2) © 2010-2017 H. Nakashima

ILP (3/13)

- Static Scheduling**
 - Latency of f = mem, f = f {+, *} f = 2

9+6=15/iteration → IPC=17/15=1.13

15

Massively Parallel Comp. Arch. (2) © 2010-2017 H. Nakashima

ILP (3/13)

- loop unrolling**
 - Expand k iterations to have one iteration
 - Improve parallelism in loop body (& reduce loop-control overhead)

```
for (x=0; x<nx; x++) un[y][x]=u[y][x]+...;
```

↓ 4-way unrolling

```
for (x=0; x<(nx/4)*4; x+=4) {
  un[y][x+0]=u[y][x+0]+...;
  un[y][x+1]=u[y][x+1]+...;
  un[y][x+2]=u[y][x+2]+...;
  un[y][x+3]=u[y][x+3]+...;
}
```

```
for (; x<nx; x++) un[y][x]=u[y][x]+...;
```

16

Massively Parallel Comp. Arch. (2) © 2010-2017 H. Nakashima

ILP (5/13)

- Dependence Graph of 4-way Unrolled Loop**

if #registers is sufficient

24/(4*iteration) → IPC=56/24=2.33 (68/24=2.83)

17

Massively Parallel Comp. Arch. (2) © 2010-2017 H. Nakashima

ILP (6/13)

- Limitation of Static Scheduling (incl. Loop Unrolling) = Insufficient Parallelism inside one or k Iterations**

declining parallelism

→ **Scheduling across Iterations = Software Pipelining**

18

Massively Parallel Comp. Arch. (2) © 2010-2017 H. Nakashima

ILP (7/13)

Software Pipelining

combine with next iteration with shift

2way-unroll and perform scheduling a little bit loosely

10(2*iteration) → IPC=30/10=3.0 (34/10=3.4)

Massively Parallel Comp. Arch. (2) © 2010-2017 H. Nakashima

ILP (8/13)

Dynamic Scheduling

= dynamic execution ordering

automatically done by microprocessor

- Out-of-order issue of executable instructions with ready-to-use operands irrespective of their program order
- Dynamic register renaming to remove unnecessary restriction of execution order
- Speculative execution of instructions given by branch predictor
- Reordering for in-order completion to cope with branch misprediction, exception, interrupt, etc.

20

Massively Parallel Comp. Arch. (2) © 2010-2017 H. Nakashima

ILP (8/13)

multiple issue of (e.g.) 4 instructions

out-of-order issue add: wait for loaded data load: not wait for add

register renaming: dynamic mapping from logical registers (f=16/g=16) to physical registers of larger amount

speculation: execute predicted insts before branch target is determined

reordering: complete when all predecessors complete (e.g. up to 4 insts)

```

fuzp = (guy+gx+8)
fuzm = (guy+gx-8)
fuzp = fuzp + fuzm
fuzz = (guy+gx)
fuzp = fuzp - fuzz
fuzm = (guy+gx)
fuzz = fuzz * ffour
fuzp = fuzp - fuzz
fuzm = fuzp * fdtth
fuzp = fuzp + ftmp
*(guy+gx) = fuzp
gx = gx + 8
gx < gx + 8
if (<) goto
fuzp = (guy+gx+8)
fuzm = (guy+gx-8)
fuzp = fuzp + fuzm
fuzz = (guy+gx)
fuzp = fuzp - fuzz
fuzm = (guy+gx)
fuzz = fuzz * ffour
fuzp = fuzp + ftmp
fuzm = (guy+gx)

```

21

Massively Parallel Comp. Arch. (2) © 2010-2017 H. Nakashima

ILP (9/13)

How Out-of-Order Execution Progresses

5/iteration > 4/iteration

→ #instructions to start/complete increases monotonically

→ reservation station (60@KU) for inst to start and/or reorder buffer (192@KU) for inst to complete will be exhausted

→ effective fetch/issue rate becomes 3.2/cycle eventually

actual latency (e.g. in KU) of load/FP-ops > 2

→ buffers exhaust frequently

→ hard to find executable insts

→ likely to be > 5/iteration

4/iteration

5/iteration → IPC=116/5=3.2

22

Massively Parallel Comp. Arch. (2) © 2010-2017 H. Nakashima

ILP (10/13)

Exploiting SIMD (AVX) Mechanism

- ymm: 256-bit media register
 - char X 32, short X 16, int X 8, long long X 4
 - float X 8, double X 4

→ perform a particular operation on multiple data

- 4-way unroll and perform 4 ops by one inst

```

for (x=0; x<nx; x+=4) {
  un[y][x] = u[y][x];
  dtth*(u[y][x+1]+u[y][x-1]+u[y+1][x]+u[y-1][x]-4*u[y][x]);
  un[y][x+1] = u[y][x+1];
  dtth*(u[y][x+2]+u[y][x]+u[y+1][x+1]+u[y-1][x+1]-4*u[y][x+1]);
  un[y][x+2] = u[y][x+2];
  dtth*(u[y][x+3]+u[y][x+1]+u[y+1][x+2]+u[y-1][x+2]-4*u[y][x+2]);
  un[y][x+3] = u[y][x+3];
  dtth*(u[y][x+4]+u[y][x+2]+u[y+1][x+3]+u[y-1][x+3]-4*u[y][x+3]);
}

```

→ but compilers often fail to SIMD-vectorize

23

Massively Parallel Comp. Arch. (2) © 2010-2017 H. Nakashima

ILP (11/13)

Why Compilers Fail to SIMD-vectorize

```

for (x=0; x<nx; x+=4) {
  un[y][x] = u[y][x] + dtth*(...);
  ...
  un[y][x+3] = u[y][x+3] + dtth*(...);
}

```

→ mutually dependent

actually they are independent definitely but compilers are often ignorant of it

this transformation does not preserve semantics if dependent

```

for (x=0; x<nx; x+=4) {
  double t0 = u[y][x] + dtth*(...);
  ...
  double t3 = u[y][x+3] + dtth*(...);
  un[y][x] = t0;
  ...
  un[y][x+3] = t3;
}

```

this code means they are independent

SIMD-vectorizable if this transformation is correct

24

Massively Parallel Comp. Arch. (2) © 2010-2017 H. Nakashima

ILP (12/13)

- How to Make Loop SIMD-Vectorizable

```

for (t=0; t<tmax; t++){
  double (*restrict uu)[NX]=u,
         (*restrict uun)[NX]=un;
  for (y=0; y<ny; y++) for (x=0; x<nx; x++){
    uun[y][x]=uu[y][x]+
    dthh*(uu[y][x+1]+uu[y][x-1]+
    uu[y+1][x]+uu[y-1][x]-4*uu[y][x]);
    tmp=un; un=u; u=tmp;
  }
}

```

may run incorrectly if assurance is incorrect

- *restrict p
- =*p or *(p+exp) is not be updated by a pointer other than p
- programmer assures &uun[y][x] != &uun[y'][x']
- SIMD-vectorizable → 1.8 cycle/iteration (> lower bound=1.25)

25

Massively Parallel Comp. Arch. (2) © 2010-2017 H. Nakashima

ILP (13/13)

- Fine Tuning and Caution

```

uun[y][x]=uu[y][x]+
dthh*(uu[y][x+1]+uu[y][x-1]+
uu[y+1][x]+uu[y-1][x]-4*uu[y][x]);

```

This transformation does not preserve program semantics
 ← FP operations are not associative

```

uun[y][x]=uu[y][x]+
dthh*((uu[y][x+1]+uu[y][x-1])+
(uu[y+1][x]+uu[y-1][x])-4*uu[y][x]);

```

reduce total latency of RHS
 → reduce #-of in-flight instructions and possibility of buffer exhaustion
 → increase chance to find executable instructions

operation result of tuned code may be slightly different from original one

26

Massively Parallel Comp. Arch. (2) © 2010-2017 H. Nakashima

Branch Prediction (1/8)

- Control Dependency
- Target of branch instruction (& its successors) depends on the branch.
- Control hazard makes pipeline stalled during fetch and decode of branch target which takes place after calculation of branch address and direction.

```

gx=gx+8
gx<gx+8
if(<)goto
fuzp=* (guy+gx+8)

```

branch address and direction are determined

fetch and decode of branch target ≈ 15 cycle

27

Massively Parallel Comp. Arch. (2) © 2010-2017 H. Nakashima

Branch Prediction (2/8)

- Predicting Branch Address
- Correspondence of (address of) branch instruction and target address is usually invariant (if taken).
 0x124C:if(<)goto xloop(=PC-94=0x11E0)
- Branch Target Buffer (BTB)
 - A kind of cache to keep (inexact) correspondence of branch instruction address and target address.
 - On the fetch of (branch) instruction, simultaneously lookup BTB to get "maybe target address".
 - Branch latency = 1 if success.
 - BTB@KU: 4K(??) entry (+ for indirect branches)
- Return Address Stack (RAS)
 - Stack to predict return address: push by call / pop by return
 - RAS@KU: 16(??) entry

28

Massively Parallel Comp. Arch. (2) © 2010-2017 H. Nakashima

Branch Prediction (3/8)

- Predicting Branch Direction
- Directions of a conditional branch (are expected to) have temporal locality
 e.g. loop terminator: #taken=n-1 / #not-taken=1
- Simple prediction mechanism = 1-bit predictor
 - 1-bit array indexed by a segment of branch inst address
 - Record 0 (not-taken) / 1 (taken) for each branch inst
 - predict (not) taken if it was (not) taken last time
 - sensitive to noise (e.g. miss twice for loop terminator)
- 2-bit saturate-counter predictor
 - 2-bit array indexed by a segment of branch inst address
 - not-taken → -1 / taken → +1 (min=0 / max=3)
 - 0 or 1 → predict as not-taken / 2 or 3 → predict as taken
 - one miss does not change its mind (e.g. #miss of loop terminator = 1)

29

Massively Parallel Comp. Arch. (2) © 2010-2017 H. Nakashima

Branch Prediction (4/8)

- 1bit vs 2bit

```

for (y=0; y<4; y++) for (x=0; x<4; j++)
  un[y][x]=u[y][x]+dthh*(...);

```

x	actual	prediction	
		1bit	2bit
1	T	N	c=2 → T
2	T	T	c=3 → T
3	T	T	c=3 → T
4	N	T	c=3 → T

T : taken
 N : not-taken

accuracy = 50% accuracy = 75%

30

Massively Parallel Comp. Arch. (2) © 2010-2017 H. Nakashima

Branch Prediction (5/8)

- Very Dangerous Code

```

for(i=0; i<N; i++) {
  if(i&1) {for odd i}
  else {for even i}
}

```

accuracy for TNTNTN... = 1/2 at most, 0 at worst

→ for(i=0; i<N; i+=2) {for even i}
 for(i=1; i<N; i+=2) {for odd i}

or

```

for(i=0; i<N; i+=2) {
  {for even i}
  {for odd i+1}
}

```

recode as them if possible

31

Massively Parallel Comp. Arch. (2) © 2010-2017 H. Nakashima

Branch Prediction (6/8)

- Using Global history (GH)

- 2-bit prediction is too local
 - spatial: only aware of a particular branch
 - temporal: remember almost nothing of actions except for last four
- GH: remember all directions of last n branches (e.g. $n = 8$)
- Mix GH into index of 2-bit predictor
 - gselect (64K)
`pred_table[gh][branch_inst_addr&0xFF]`
 - gshare (64K)
`pred_table[gh^(branch_inst_addr&0xFFFF)]`

in KU, gshare or its variant is used ($n=32?$)

32

Massively Parallel Comp. Arch. (2) © 2010-2017 H. Nakashima

Branch Prediction (7/8)

- Effect of GH (1/2)

```

for(t=0; t<tmax; t++) {
  for(y=0; y<4; y++) for(x=0; x<4; j++)
    un[y][x]=u[y][x]+dthh*(...);
  tmp=un; un=u; u=tmp;
}

```

for(t): T T T T T T T T T T

for(y): T T T T N T T T T N

for(x): TTTN TTTN TTTN TTTN TTTN TTTN TTTN TTTN

y=0				y=1					
x	GH	for(x)	pred	actual	x	GH	for(x)	pred	actual
1	TTTNT	TTT...	T	T	1	TTTTN	TTT...	T	T
2	TTTNT	TTT...	T	T	2	TTTNT	TTT...	T	T
3	TNTTT	TTT...	T	T	3	TNTTT	TTT...	T	T
4	NNTTT	NNN...	N	N	4	TNTTT	NNN...	N	N

33

Massively Parallel Comp. Arch. (2) © 2010-2017 H. Nakashima

Branch Prediction (8/8)

- Effect of GH (2/2)

```

for(p=h; p!=NULL; p=p->next){
  if(p->data==x) break;
}
if(p==NULL) ...

```

- Truth value of $p==NULL$
 - last branch is for $p->data==x$ (true) → false
 - last branch is for $p!=NULL$ (false) → true
 - predictable as $GH=?????T \rightarrow false / ????N \rightarrow true$
- $if(i&1)\{...\}else\{...\}$ is also predictable but ...

34

Massively Parallel Comp. Arch. (2) © 2010-2017 H. Nakashima

Predicated Instruction (1/2)

- Conditional Branch & SIMD (1/2)

```

for(i=0; i<n; i++){
  if(a[i]<0) a[i]=-a[i]; /* a[i]=abs(a[i]);
}

```

definitely correct transformation

```

for(i=0; i<n; i+=4){
  if(a[i+0]<0) a[i+0]=-a[i+0];
  if(a[i+1]<0) a[i+1]=-a[i+1];
  if(a[i+2]<0) a[i+2]=-a[i+2];
  if(a[i+3]<0) a[i+3]=-a[i+3];
}

```

How SIMD-vectorize these branches??

SIMD-vectorizable of course

35

Massively Parallel Comp. Arch. (2) © 2010-2017 H. Nakashima

Predicated Instruction (1/2)

- Conditional Branch & SIMD (1/2)

```

for(i=0; i<n; i++){
  if(a[i]<0) a[i]=-a[i]; /* a[i]=abs(a[i]);
}

```

```

for(i=0; i<n; i+=4){
  int c0=a[i+0]<0; double t0=-a[i+0];
  int c1=a[i+1]<0; double t1=-a[i+1];
  int c2=a[i+2]<0; double t2=-a[i+2];
  int c3=a[i+3]<0; double t3=-a[i+3];

  a[i+0]=c0 ? t0 : a[i+0];
  a[i+1]=c1 ? t1 : a[i+1];
  a[i+2]=c2 ? t2 : a[i+2];
  a[i+3]=c3 ? t3 : a[i+3];
}

```

conditional move

- transform latency from branch to arithmetic
- no misprediction because not predicting → efficient (usually)

36

Massively Parallel Comp. Arch. (2) © 2010-2017 H. Nakashima

Basics of Cache (1/7)

- Overview of **Cache** (≠ Cash)
 - High-speed memory to fill the access speed gap between CPU and main memory
CPU : cache : memory = 1 : 1~20 : 50~100
 - Keep data to be **likely** accessed in near future
 - exploiting/expecting **locality**
 - temporal: repeated access to data accessed recently
 - spatial: access to data proximate to recently accessed data
 - accesses against expectation result in very poor performance
 - e.g. skipping accesses in a huge array

37

Massively Parallel Comp. Arch. (2) © 2010-2017 H. Nakashima

Basics of Cache (2/7)

- Configuration Parameter: $L \times S \times W = \text{capacity (e.g. 32KB)}$
 - line size (≈32~256B)
line: unit of bookkeeping & data transfer

38

Massively Parallel Comp. Arch. (2) © 2010-2017 H. Nakashima

Basics of Cache (3/7)

- $L > 8B (= \text{sizeof}(\text{double}))$
 - Exploiting/expecting spatial locality
 - access to $u[x]$ → accesses to $u[x+1], u[x+2], \dots$
 - Efficient transfer from/to memory by contiguous access
 - latency = 50~100cycle, bandwidth = 5~10B/cycle

```

for(y=0;y<ny;y++) for(x=0;x<nx;x++)
    un[y][x]=u[y][x]*dthh*(...);
for(x=0;x<nx;x++) for(y=0;y<ny;y++)
    un[y][x]=u[y][x]*dthh*(...);
    
```

miss rate = $8/L$

better

worse

```

for(i=0;i<N;i++){ j=f(i);
    obj[j].x+=obj[j].vx; obj[j].y+=obj[j].vy;
}
for(i=0;i<N;i++){ j=f(i);
    obj.x[j]+=obj.vx[j]; obj.y[j]+=obj.vy[j];
}
    
```

better

worse

39

Massively Parallel Comp. Arch. (2) © 2010-2017 H. Nakashima

Basics of Cache (4/7)

- $L = 2^l, S = 2^s \rightarrow L \times S = 2^{l+s}$
 - Data pair **conflict** with each other if their address difference is $2^{l+s}n$.
 - Arrays of power-of-two size are **dangerous**.

```

#define N 1024
double u[N][N], un[N][N];
...
for(x=1;x<N-1;x++){
    un[y][x]=u[y][x]+dthh*(u[y][x+1]+u[y][x-1]+
        u[y+1][x]+u[y-1][x]+4*u[y][x]);
}
struct {double x[N],y[N],z[N],
        vx[N],vy[N],vz[N];} obj;
...
for(i=0;i<N;i++){
    obj.x[i]+=obj.vx[i]; obj.y[i]+=obj.vy[i];
    obj.z[i]+=obj.vz[i];
}
    
```

conflict with each other
 $L \times S < 8192$ & $W < 4$
→ misses always

$L \times S < 8192$ & $W < 6$
→ misses always

40

Massively Parallel Comp. Arch. (2) © 2010-2017 H. Nakashima

Basics of Cache (5/7)

- 3C of Cache Misses
 - Compulsory miss
 - data is not in cache on its first access
→ inevitable
 - Capacity miss
 - #-of lines accessed after last access > capacity
← poor temporal (and/or spatial) locality
 - Conflict miss
 - #-of competing lines accessed after last access > #ways
← e.g. discontinuous accesses (esp. with 2^n stride)

41

Massively Parallel Comp. Arch. (2) © 2010-2017 H. Nakashima

Basics of Cache (6/7)

@KU's supercomputer

- (Representative) Memory Hierarchy (≈100)
 - lat = 1~4(+)
 - 64? X 64? X 8? = 32KB (+1.5K? μop)
 - lat = 5~30
 - L2=12(+) L3=34(++)
 - lat = 50~200
 - memory (1GB~) 128GB

Write Through: immediately write to lower level
Write Back: write to lower level at replacement

L2: 64 X 512 X 8 = 256KB/core
L3: 64 X 2048? X 20? X 18 = 45MB/18core

42

Massively Parallel Comp. Arch. (2) © 2010-2017 H. Nakashima

Basics of Cache (7/7)

- Prefetch = Load a line to cache if the line is expected to be accessed in near future**
 - Hardware Prefetch**
 - = prefetch based on detection of particular access pattern e.g. in KU's supercomputer
 - contiguous access to L1D → prefetch next line
 - last k addresses of a load = $b + s \times i$ ($0 \leq i < k$)
 - prefetch line including $b + s \times k$ to L1D
 - a 64B-line in L2 is loaded to L1
 - prefetch next line to L2 if absent
 - watch 32(?) contiguous access streams
 - if L1 require a line in a stream, prefetch next line to L2/L3
 - Software Prefetch**
 - prefetch a line to L1 by a dedicated instruction
 - compiler inserts prefetch instructions estimating access patterns and effectiveness

43

Massively Parallel Comp. Arch. (2) © 2010-2017 H. Nakashima

TLB (1/2)

- Virtual/Physical Address Translation**

44

Massively Parallel Comp. Arch. (2) © 2010-2017 H. Nakashima

TLB (2/2)

- Translation Look-aside Buffer**
 - Cache of PA translated from VA**
 - key = virtual page number (high-order bits of VA)
 - data = physical page number (high-order bits of PA)
 - Relatively small ≈ 10~1000 entries**
 - @KU: L1-ITLB = $32 \times 4\text{way} = 128 (+\alpha)$
 - L1-DTLB = $16 \times 4\text{way} = 64 (+\alpha+\beta)$
 - L2-TLB = $256 \times 6\text{way} = 1536 (+\alpha)$
 - miss → page table lookup → large overhead**
 - many accesses with long stride → miserable

```
for (x=0; x<1000; x++) for (y=0; y<1000; y++)
  un[y][x]=u[y][x]+dthh*(...);
```

2000 pages are accessed in inner loop → 2 misses/iteration

- L2/L3 cache is accessed by PA**
 - contiguous in VA ≠ contiguous in L2/L3

45

Massively Parallel Comp. Arch. (2) © 2010-2017 H. Nakashima

Performance with Cache & TLB (1/7)

- Memory Access Pattern**

```
for (y=0; y<ny; y++) for (x=0; x<nx; x++)
  un[y][x]=u[y][x]+
  dthh*(u[y][x+1]+u[y][x-1]+
  u[y+1][x]+u[y-1][x]-4*u[y][x]);
```

46

Massively Parallel Comp. Arch. (2) © 2010-2017 H. Nakashima

Performance with Cache & TLB (2/7)

47

Massively Parallel Comp. Arch. (2) © 2010-2017 H. Nakashima

Performance with Cache & TLB (3/7)

- Optimization for Large nx : x-dim Tiling**

```
for (xx=0; xx<nx; xx+=TX) {
  int nxcx=(xx+TX<nx)?xx+TX:nx;
  for (y=0; y<ny; y++) for (x=xx; x<nxcx; x++) {
    uu[y][x]=uu[y][x]+
    dthh*(uu[y][x+1]+uu[y][x-1]+
    uu[y+1][x]+uu[y-1][x]-4*uu[y][x]);
  }
}
```
- large nx: #misses / 8-iteration = 4 → 2**

48

Massively Parallel Comp. Arch. (2) © 2010-2017 H. Nakashima

Performance with Cache & TLB (4/7)

■ Super Optimization = Temporal Tiling

tile S cache
→ almost hit

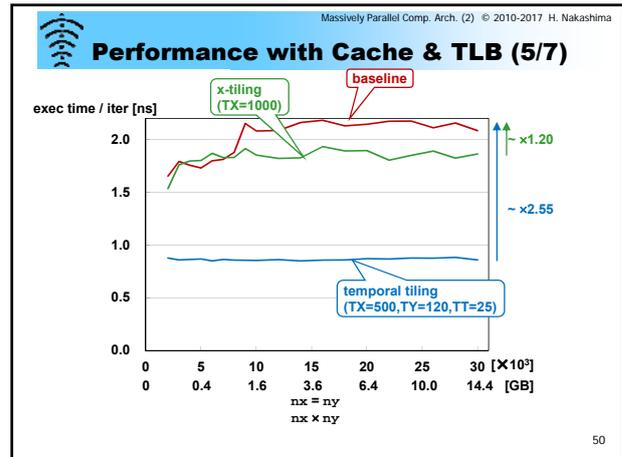
perimeter+0 (or +1)
→ $t \rightarrow t+1$ (&+0)

perimeter+1 (or +2)
→ $t \rightarrow t+2$ (&+1)

perimeter+2 (or +3)
→ $t \rightarrow t+3$ (&+2)

perimeter+3 (or +4)
→ $t \rightarrow t+4$ (&+3)

49



Massively Parallel Comp. Arch. (2) © 2010-2017 H. Nakashima

Performance with Cache & TLB (6/7)

■ Contiguous vs Discontiguous Accesses

```

for (y=0; y<ny; y++) for (x=0; x<nx; x++)
  uun[y][x]=uu[y][x]+
  dthh*(uu[y][x+1]+uu[y][x-1]+
  uu[y+1][x]+uu[y-1][x]-4*uu[y][x])

```

■ #cache-misses = 2~4 / 8-iteration
#TLB-misses = 2~4 / 256K-iteration

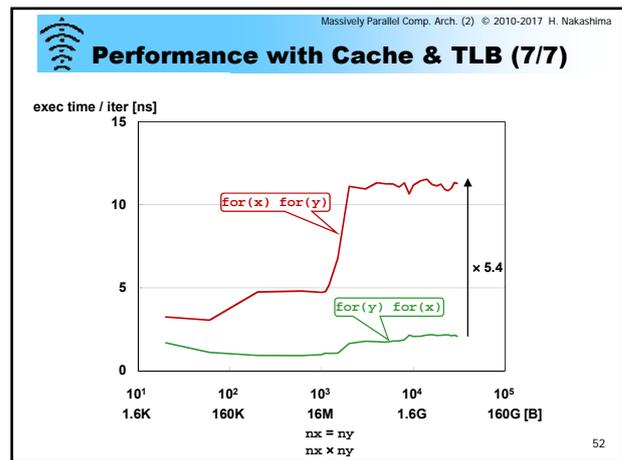
```

for (x=0; x<nx; x++) for (y=0; y<ny; y++)
  uun[y][x]=uu[y][x]+
  dthh*(uu[y][x+1]+uu[y][x-1]+
  uu[y+1][x]+uu[y-1][x]-4*uu[y][x])

```

■ #cache-misses = 2.25 / 1-iteration
#TLB-misses = 2 / 1-iteration

51



Massively Parallel Comp. Arch. (2) © 2010-2017 H. Nakashima

Conclusion: Key Issues of HPC

- Pipeline / ILP / out-of-order / SIMD
 - Relying on compiler/hardware is almost safe.
 - Remember restrict in C programming
- Branch Prediction / Speculation
 - Loops with many iterations are usually OK.
 - Some if-else in loops are really dangerous.
- Cache / TLB
 - Dominant performance factor.
 - Don't let down your stupid hardware.
 - Small care can improve performance significantly.
 - As much contiguous as possible.
 - Avoid access conflict if possible.

53