EECS 583 – Class 12 Superblock Scheduling Intro. to Loop Scheduling

University of Michigan

October 19, 2011

Announcements & Reading Material

- Reminder: HW 2 Speculative LICM
 - » Daya has office hours Thurs and Fri 3-5pm
 - » Deadline extended to Monday 11:59pm
- Class project proposals
 - » Next week: Daya and I will meet with each group for 15 mins to discuss informal project proposals
 - Next wk: No class on Wednes, but we will have class Monday
 - » Signup sheet on my door Thurs (tomorrow)
- Today's class
 - "Sentinel Scheduling for VLIW and Superscalar Processors",
 S. Mahlke et al., ASPLOS-5, Oct. 1992, pp.238-247.
- Next class
 - » "Iterative Modulo Scheduling: An Algorithm for Software Pipelining Loops", B. Rau, MICRO-27, 1994, pp. 63-74.

Homework Problem –	Op	erat	ion Sc	hedul	ing
Machina: 2 jaqua 1 mamory port 1 ALU		Op 1	priorit	у	
Machine. 2 issue, 1 memory port, 1 ALO		1	0 7		
Memory port = 2 cycles, pipelined		3	4		
ALU = 1 cycle		4	5		
•		5	2		
		6	3		
0.1 (1m) (2m) 0.0		7 0	3		
		0 9	2 3		
		10	1		
2,3 (3) (4m) 2,2		RU m	an	Sch	nedule
1 1 2	time		MEM	Time	Dlaced
$25 \left(\frac{1}{5} \right)^{3,4} \left(\frac{1}{5} \right)^{3,4}$		ALU			
3,5 (5) (6) (7) 4,4	0		X	0	2
	1		Х	1	1
1 5.5 (8) (9m) 0.4	2		Х	2	4
$\langle \downarrow_1 \rangle$	3	Χ	Х	3	3, 9
	4	Χ		4	6
6,6 (10)	5	Χ		5	7
	6	Χ		6	5
Calculate height-based priorities	7	Χ		7	8
 Schedule using <u>Operation</u> scheduler 	8	Х		8	10

From Last Time: Generalize Beyond a Basic Block

- Superblock
 - » Single entry
 - » Multiple exits (side exits)
 - » No side entries
- Schedule just like a BB
 - » Priority calculations needs change
 - » Dealing with control deps



Lstart in a Superblock

Not a single Lstart any more

- » 1 per exit branch (Lstart is a vector!)
- » Exit branches have probabilities

Lstart0

Estart

op

1

2

3

4

5

6

Lstart1



Operation Priority in a Superblock

Priority – Dependence height and speculative yield

- » Height from op to exit * probability of exit
- » Sum up across all exits in the superblock

Priority(op) = **SUM**(Probi * (MAX_Lstart – Lstarti(op) + 1))

valid late times for op



Dependences in a Superblock



* Data dependences shown, all are reg flow except $1 \rightarrow 6$ is reg anti

* Dependences define precedence ordering of operations to ensure correct execution semantics

* What about control dependences?

* Control dependences define precedence of ops with respect to branches

Conservative Approach to Control Dependences



* Make branches barriers, nothing moves above or below branches

* Schedule each BB in SB separately

* Sequential schedules

* Whole purpose of a superblock is lost

Upward Code Motion Across Branches

- Restriction 1a (register op)
 - » The destination of op is not in liveout(br)
 - » Wrongly kill a live value
- Restriction 1b (memory op)
 - » Op does not modify the memory
 - Actually live memory is what matters, but that is often too hard to determine
- Restriction 2
 - » Op must not cause an exception that may terminate the program execution when br is taken
 - Op is executed more often than it is supposed to (<u>speculated</u>)
 - » Page fault or cache miss are ok
- Insert control dep when either restriction is violated



Downward Code Motion Across Branches

- Restriction 1 (liveness)
 - » If no compensation code
 - Same restriction as before, destination of op is not liveout
 - » Else, no restrictions
 - Duplicate operation along both directions of branch if destination is liveout
- Restriction 2 (speculation)
 - » Not applicable, downward motion is not speculation
- Again, insert control dep when the restrictions are violated
- Part of the philosphy of superblocks is no compensation code inseration hence R1 is enforced!



Add Control Dependences to a Superblock



Class Problem



Draw the dependence graph

Relaxing Code Motion Restrictions

- Upward code motion is generally more effective
 - Speculate that an op is useful (just like an out-of-order processor with branch pred)
 - » Start ops early, hide latency, overlap execution, more parallelism
- Removing restriction 1
 - » For register ops use register renaming
 - » Could rename memory too, but generally not worth it
- Removing restriction 2
 - » Need hardware support (aka <u>speculation models</u>)
 - Some ops don't cause exceptions
 - Ignore exceptions
 - Delay exceptions



R1: y is not in liveout(1)R2: op 2 will never cause an exception when op1 is taken

Restricted Speculation Model

- Most processors have 2 classes of opcodes
 - » Potentially exception causing
 - load, store, integer divide, floating-point
 - » Never excepting
 - Integer add, multiply, etc.
 - Overflow is detected, but does not terminate program execution
- Restricted model
 - » R2 only applies to potentially exception causing operations
 - Can freely speculate all never exception ops (still limited by R1 however)



We assumed restricted speculation when this graph was drawn.

This is why there is no cdep between $4 \rightarrow 6$ and $4 \rightarrow 8$

General Speculation Model

- 2 types of exceptions
 - » Program terminating (traps)
 - Div by 0, illegal address
 - Fixable (normal and handled at run time)
 - Page fault, TLB miss
- General speculation
 - Processor provides nontrapping versions of all operations (div, load, etc)
 - Return some bogus value (0)
 when error occurs
 - » R2 is completely ignored, only R1 limits speculation
 - » Speculative ops converted into non-trapping version
 - Fixable exceptions handled as usual for non-trapping ops



Programming Implications of General Spec

- Correct program
 - » No problem at all
 - Exceptions will only result when branch is taken
 - Results of excepting speculative operation(s) will not be used for anything useful (R1 guarantees this!)
- Program debugging
 - Non-trapping ops make this almost impossible
 - Disable general speculation during program debug phase



Class Problem



 Starting with the graph assuming restricted speculation, what edges can be removed if general speculation support is provided?
 With more renaming, what dependences could be removed?

Sentinel Speculation Model

- Ignoring all speculative exceptions is painful
 - » Debugging issue (is a program ever fully correct?)
- Also, handling of all fixable exceptions for speculative ops can be slow
 - » Extra page faults
- Sentinel speculation
 - » Mark speculative ops (opcode bit)
 - Exceptions for speculative ops are noted, but not handed immediately (return garbage value)
 - Check for exception conditions in the "home block" of speculative potentially excepting ops



Delaying Speculative Exceptions

- ✤ 3 things needed
 - » Record exceptions
 - » Check for exceptions
 - » Regenerate exception
 - Re-execute ops including dependent ops
 - Terminate execution or process exception
- Recording them
 - Extend every register with an extra bit
 - Exception tag (or NAT bit)
 - Reg data is garbage when set
 - Bit is set when either
 - Speculative op causes exception
 - Speculative op has a NAT'd source operand (exception propagation)



Delaying Speculative Exceptions (2)

- Check for exceptions
 - Test NAT bit of appropriate register (last register in dependence chain) in home block
 - » Explicit checks
 - Insert new operation to check NAT
 - » Implicit checks
 - Non-speculative use of register automatically serves as NAT check
- Regenerate exception
 - » Figure out the exact cause
 - » Handle if possible
 - » Check with NAT condition branches to "recovery code"
 - Compiler generates the recovery code specific to each check



Delaying Speculative Exceptions (3)

In recovery code, the exception condition Recovery code consists of chain will be regenerated as the excepting op of operations starting with a is re-executed with the same inputs potentially excepting speculative op up to its corresponding check If the exception can be handled, it is, all dependent ops are re-executed, and execution 2': y = *x3': z = y + 41: branch x == 0is returned to point after the check If the exception is a program error, execution is terminated in the recovery code branch NAT(z) fixup done: 4: *w = z Recovery code fixup: 2": y = *x3": z = y + 4jump done

Implicit vs Explicit Checks

- Explicit
 - » Essentially just a conditional branch
 - » Nothing special needs to be added to the processor
 - » Problems
 - Code size
 - Checks take valuable resources
- Implicit
 - » Use existing instructions as checks
 - » Removes problems of explicit checks
 - » However, how do you specify the address of the recovery block?, how is control transferred there?
 - » Hardware table
 - Indexed by PC
 - Indicates where to go when NAT is set
- IA-64 uses explicit checks

Class Problem



- 1. Move ops 5, 6, 8 as far up in the SB as possible assuming sentinel speculation support
- 2. Insert the necessary checks and recovery code (assume ld, st, and div can cause exceptions)

Change Focus to Scheduling Loops



Basic Approach – List Schedule the Loop Body



Total time = 6 * n

Unroll Then Schedule Larger Body



Total time = 7 * n/2

Problems With Unrolling

- Code bloat
 - » Typical unroll is 4-16x
 - » Use profile statistics to only unroll "important" loops
 - » But still, code grows fast
- Barrier after across unrolled bodies
 - » I.e., for unroll 2, can only overlap iterations 1 and 2, 3 and 4, ...
- Does this mean unrolling is bad?
 - » No, in some settings its very useful
 - Low trip count
 - Lots of branches in the loop body
 - » But, in other settings, there is room for improvement

Overlap Iterations Using Pipelining



A Software Pipeline

