

Instruction Scheduling

Note by Baris Aktemur:

Our slides are adapted from Cooper and Torczon's slides that they prepared for COMP 412 at Rice.

Copyright 2010, Keith D. Cooper & Linda Torczon, all rights reserved.

Students enrolled in Comp 412 at Rice University have explicit permission to make copies of these materials for their personal use.

Faculty from other educational institutions may use these materials for nonprofit educational purposes, provided this copyright notice is preserved.

What Makes Code Run Fast?

- Many operations have non-zero latencies
- Modern machines can issue several operations per cycle
- Execution time is *order-dependent* (and has been since the 60's)

<u>Operation</u>	<u>Cycles</u>
load	3
store	3
loadI	1
add	1
mult	2
fadd	1
fmult	2
shift	1
branch	0 to 8

- Loads & stores may or may not block on issue
 - > Non-blocking \Rightarrow fill those issue slots
- Branch costs vary with path taken
- Branches typically have delay slots
 - > Fill slots with unrelated operations
 - > Percolates branch upward
- Scheduler should hide the latencies

Assumed latencies for example on next slide.

Example

$$a \leftarrow a * 2 * b * c * d$$

Simple schedule

Start	Operations
1	loadAI rarp,@a ⇒ r1
4	add r1,r1 ⇒ r1
5	loadAI rarp,@b ⇒ r2
8	mult r1,r2 ⇒ r1
10	loadAI rarp,@c ⇒ r2
13	mult r1,r2 ⇒ r1
15	loadAI rarp,@d ⇒ r2
18	mult r1,r2 ⇒ r1
20	storeAI r1 ⇒ rarp,@a

Schedule loads early

Start	Operations
1	loadAI rarp,@a ⇒ r1
2	loadAI rarp,@b ⇒ r2
3	loadAI rarp,@c ⇒ r3
4	add r1,r1 ⇒ r1
5	mult r1,r2 ⇒ r1
6	loadAI rarp,@d ⇒ r2
7	mult r1,r3 ⇒ r1
9	mult r1,r2 ⇒ r1
11	storeAI r1 ⇒ rarp,@a

Reordering operations for speed is called *instruction scheduling*

ALU Characteristics



This data is surprisingly hard to measure accurately

- Value-dependent behavior
- Context-dependent behavior
- Compiler behavior
 - Have seen gcc underallocate & inflate operation costs with memory references (spills)
 - Have seen commercial compiler generate 3 extra ops per divide raising effective cost by 3
- Difficult to reconcile measured reality with the data in the Manuals (e.g. integer divide on Nehalem)

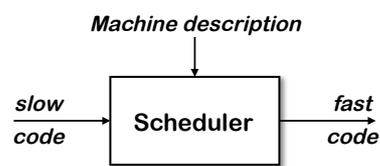
Intel E5530 operation latencies	
Instruction	Cost
64 bit integer subtract	1
64 bit integer multiply	3
64 bit integer divide	41
Double precision add	3
Double precision subtract	3
Double precision multiply	5
Double precision divide	22
Single precision add	3
Single precision subtract	3
Single precision multiply	4
Single precision divide	14

Instruction Scheduling (Engineer's View)

The Problem

Given a code fragment for some target machine and the latencies for each individual operation, reorder the operations to minimize execution time

The Concept



The Task

- Produce correct code
- Minimize wasted cycles
- Avoid spilling registers
- Operate efficiently

Comp 412, Fall 2010

4

Instruction Scheduling (The Abstract View)

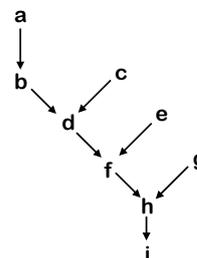
To capture properties of the code, build a precedence graph G

- Nodes $n \in G$ are operations with $type(n)$ and $delay(n)$
- An edge $e = (n_1, n_2) \in G$ if & only if n_2 uses the result of n_1

```

a: loadAI  rarp,@a => r1
b: add     r1,r1 => r1
c: loadAI  rarp,@b => r2
d: mult    r1,r2 => r1
e: loadAI  rarp,@c => r3
f: mult    r1,r2 => r1
g: loadAI  rarp,@d => r2
h: mult    r1,r2 => r1
i: storeAI r1 => rarp,@a
    
```

The Code



The Precedence Graph

Comp 412, Fall 2010

5

Instruction Scheduling (Definitions)

A *correct schedule* S maps each $n \in N$ into a non-negative integer representing its cycle number, and

1. $S(n) \geq 0$, for all $n \in N$, obviously
2. If $(n_1, n_2) \in E$, $S(n_1) + \text{delay}(n_1) \leq S(n_2)$
3. For each type t , there are no more operations of type t in any cycle than the target machine can issue

The *length* of a schedule S , denoted $L(S)$, is

$$L(S) = \max_{n \in N} (S(n) + \text{delay}(n))$$

The goal is to find the shortest possible correct schedule.

S is *time-optimal* if $L(S) \leq L(S_i)$, for all other schedules S_i

A schedule might also be optimal in terms of registers, power, or space....

Instruction Scheduling (What's so difficult?)

Critical Points

- All operands must be available
- Multiple operations can be *ready*
- Moving operations can lengthen register lifetimes
- Placing uses near definitions can shorten register lifetimes
- Operands can have multiple predecessors

Together, these issues make scheduling *hard* (NP-Complete)

Local scheduling is the simple case

- Restricted to straight-line code
- Consistent and predictable latencies

Instruction Scheduling: The Big Picture

1. Build a precedence graph, P
2. Compute a priority function over the nodes in P
3. Use list scheduling to construct a schedule, 1 cycle at a time
 - a. Use a queue of operations that are ready
 - b. At each cycle
 - I. Choose the highest priority ready operation & schedule it
 - II. Update the ready queue

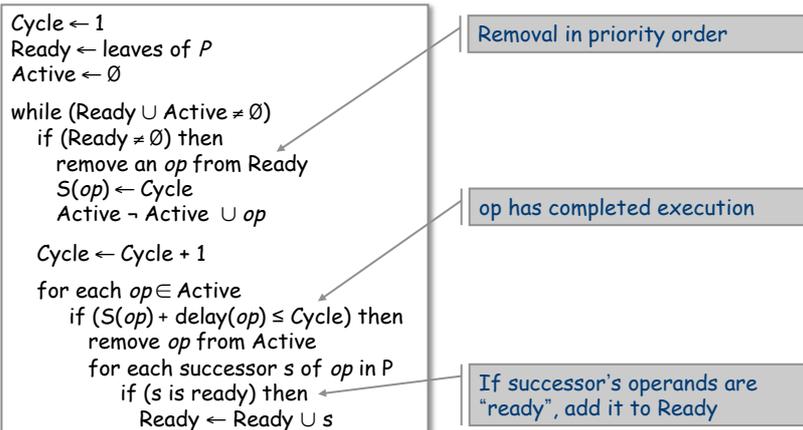
Local list scheduling

- The dominant algorithm for thirty years
- A greedy, heuristic, local technique

Comp 412, Fall 2010

*8

Local List Scheduling



Comp 412, Fall 2010

9

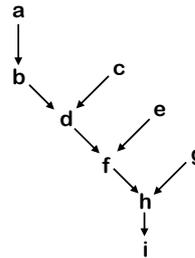
Scheduling Example

1. Build the precedence graph

```

a: loadAI  rarp,@a => r1
b: add     r1,r1 => r1
c: loadAI  rarp,@b => r2
d: mult    r1,r2 => r1
e: loadAI  rarp,@c => r3
f: mult    r1,r2 => r1
g: loadAI  rarp,@d => r2
h: mult    r1,r2 => r1
i: storeAI r1      => rarp,@a
    
```

The Code



The Precedence Graph

Scheduling Example

1. Build the precedence graph 2. Determine priority

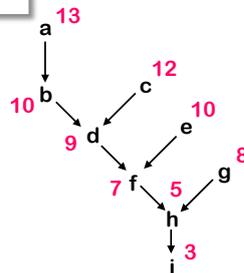
Operation	Cycles
load	3
store	3
loadI	1
add	1
mult	2
fadd	1
fmult	2
shift	1
branch	0 to 8

weighted path

```

a: loadAI  rarp,@a => r1
b: add     r1,r1 => r1
c: loadAI  rarp,@b => r2
d: mult    r1,r2 => r1
e: loadAI  rarp,@c => r3
f: mult    r1,r2 => r1
g: loadAI  rarp,@d => r2
h: mult    r1,r2 => r1
i: storeAI r1      => rarp,@a
    
```

The Code



The Precedence Graph

Scheduling Example

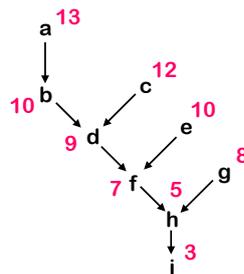
1. Build the precedence graph
2. Determine priorities: longest latency-weighted path
3. Perform list scheduling

```

1) a: loadAl    r0,@w  => r1
2) c: loadAl    r0,@x  => r2
3) e: loadAl    r0,@y  => r3
4) b: add       r1,r1   => r1
5) d: mult     r1,r2   => r1
6) g: loadAl    r0,@z  => r2
7) f: mult     r1,r3   => r1
9) h: mult     r1,r2   => r1
11) i: storeAl  r1      => r0,@w
    
```

The Code

Used a new register name



The Precedence Graph

More List Scheduling

List scheduling breaks down into two distinct classes

Forward list scheduling

- Start with available operations
- Work forward in time
- Ready \Rightarrow all operands available

Backward list scheduling

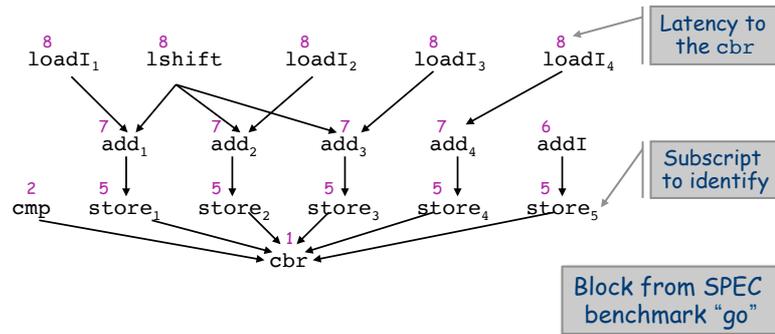
- Start with no successors
- Work backward in time
- Ready \Rightarrow latency covers uses

Variations on list scheduling

- Prioritize critical path(s)
- Schedule last use as soon as possible
- Depth first in precedence graph (minimize registers)
- Breadth first in precedence graph (minimize interlocks)
- Prefer operation with most successors

Local Scheduling

Forward and backward can produce different results



Operation	load	loadI	add	addI	store	cmp
Latency	1	1	2	1	4	1

Comp 412, Fall 2010

14

Local Scheduling

Forward Schedule

	Int	Int	Mem
1	loadI ₁	lshift	
2	loadI ₂	loadI ₃	
3	loadI ₄	add ₁	
4	add ₂	add ₃	
5	add ₄	addI	store ₁
6	cmp		store ₂
7			store ₃
8			store ₄
9			store ₅
10			
11			
12			
13	cbr		

Backward Schedule

	Int	Int	Mem
1	loadI ₄		
2	addI	lshift	
3	add ₄	loadI ₃	
4	add ₃	loadI ₂	store ₅
5	add ₂	loadI ₁	store ₄
6	add ₁		store ₃
7			store ₂
8			store ₁
9			
10			
11	cmp		
12	cbr		

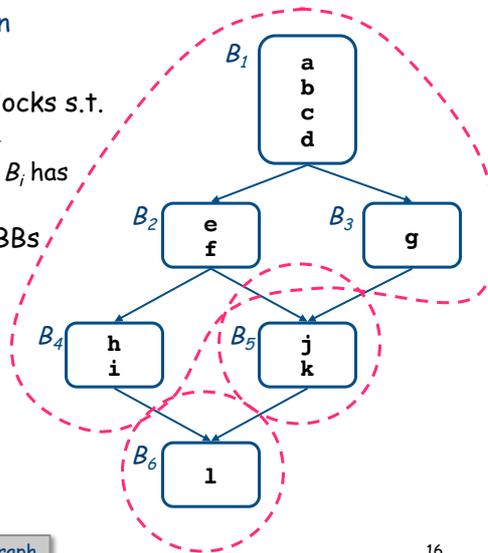
Comp 412, Fall 2010

Using "latency to root" as the priority function 15

Scheduling Larger Regions

One step beyond a block is an Extended Basic Block (EBB)

- EBB is a maximal set of blocks s.t.
 - Set has a single entry, B_i
 - Each block B_j other than B_i has exactly one predecessor
- Example CFG has three EBBs



Comp 412, Fall 2010

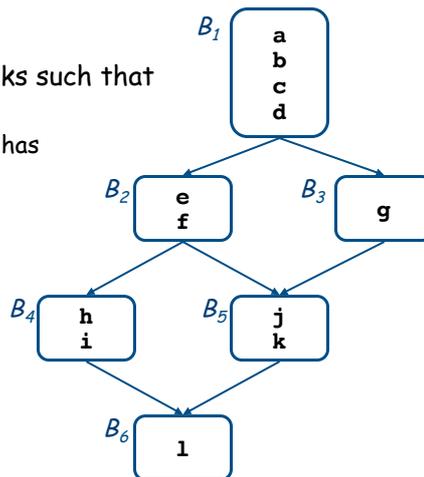
CFG = Control Flow Graph

16

Scheduling Larger Regions

One step beyond a block is an Extended Basic Block (EBB)

- EBB is a maximal set of blocks such that
 - Set has a single entry, B_i
 - Each block B_j other than B_i has exactly one predecessor
- Example has three EBBs
 - Big EBB has two paths
 - $\{B_1, B_2, B_4\}$ & $\{B_1, B_3\}$
- Many optimizations operate on EBBs (including scheduling)



Comp 412, Fall 2010

17

Scheduling Larger Regions

Superlocal Scheduling

- Schedule entire paths through EBBs
- Example has four EBB paths

