Compiler Backend Martin Kellogg

PA3 deadline is today (AoE)
 How is it going?

- PA3 deadline is **today** (AoE)
 - How is it going?
- I will hold an extra office hour today 11:30-12:30 for those who would like to see a PA3 test case
 - You may also be able to catch me either between 2 and 2:30 in my office or at the CS seminar this afternoon, but no promises

- PA3 deadline is **today** (AoE)
 - How is it going?
- I will hold an extra office hour today 11:30-12:30 for those who would like to see a PA3 test case
 - You may also be able to catch me either between 2 and 2:30 in my office or at the CS seminar this afternoon, but no promises
- PA4 is still Coming Soon[™] (I'm actually trying to get this autograder right the first time...)

- PA3 deadline is **today** (AoE)
 - How is it going?
- I will hold an extra office hour today 11:30-12:30 for those who would like to see a PA3 test case
 - You may also be able to catch me either between 2 and 2:30 in my office or at the CS seminar this afternoon, but no promises
- PA4 is still Coming Soon[™] (I'm actually trying to get this autograder right the first time...)
- Note that PA4c1's specification is TAC -> TAC
 - that is, the input is also a .cl-tac file
 - PA4c1 is due April 28, and is mostly optional

Agenda

- Interprocedural optimizations (and analysis)
- Compiler backend overview
 - Instruction selection
 - Instruction scheduling
 - Register allocation basics
- Next time:
 - Deep dive into register allocation

Agenda

- Interprocedural optimizations (and analysis)
- Compiler backend overview
 - Instruction selection
 - Instruction scheduling
 - Register allocation basics
- Next time:
 - Deep dive into register allocation

• Dividing the program up into procedures give one big benefit: separate compilation

- Dividing the program up into procedures give one big benefit: separate compilation
 - we can also optimize each procedure independently using global analyses like those we've discussed today

- Dividing the program up into procedures give one big benefit: separate compilation
 - we can also optimize each procedure independently using global analyses like those we've discussed today
- However, procedure calls also introduce significant overhead
 pre-call/post-return bookkeeping, prologue/epilogue, jump

- Dividing the program up into procedures give one big benefit: separate compilation
 - we can also optimize each procedure independently using global analyses like those we've discussed today
- However, procedure calls also introduce significant overhead
 pre-call/post-return bookkeeping, prologue/epilogue, jump
- Calls are also hard to reason about in global optimizations
 - compiler doesn't know what will happen inside the call

- Dividing the program up into procedures give one big benefit: separate compilation
 - we can also optimize each procedure independently using global analyses like those we've discussed today
- However, procedure calls also introduce significant overhead
 pre-call/post-return bookkeeping, prologue/epilogue, jump
- Calls are also hard to reason about in global optimizations
 - compiler doesn't know what will happen inside the call
- These downsides of procedure calls motivate *interprocedural* (or "*whole-program*") optimizations that span procedure boundaries

- Dividing the program up into procedures give one big benefit:
 separate compilation
 - we can also optimize e global analyses like the
- However, procedure calls
 pre-call/post-return bd

Today we will take a brief look at two interprocedural optimizations:

- inlining
- tail call optimization
- Calls are also hard to reason about in global optimizations
 - compiler doesn't know what will happen inside the call
- These downsides of procedure calls motivate *interprocedural* (or "*whole-program*") optimizations that span procedure boundaries

• Consider the following procedure:

int add(int x, int y):
 return x + y

• Consider the following procedure:

```
int add(int x, int y):
 return x + y
```

• Imagine generating code for a call to this procedure:

• Consider the following procedure:

```
int add(int x, int y):
 return x + y
```

Imagine generating code for a call to this procedure:
 the actual procedure body is only one instruction

• Consider the following procedure:

```
int add(int x, int y):
 return x + y
```

- Imagine generating code for a call to this procedure:
 - the actual procedure body is only one instruction
 - the prologue and epilogue dominate, we may have to spill registers at call sites, etc.

• Consider the following procedure:

```
int add(int x, int y):
 return x + y
```

- Imagine generating code for a call to this procedure:
 - the actual procedure body is only one instruction
 - the prologue and epilogue dominate, we may have to spill registers at call sites, etc.
- Key idea of *inlining*: for such a procedure call, replace the call with the procedure's body

• Inlining is useful when:

- Inlining is useful when:
 - the body of the procedure to be inlined is much shorter than the prologue/epilogue

- Inlining is useful when:
 - the body of the procedure to be inlined is much shorter than the prologue/epilogue
 - inlining enables **specialization** (e.g., arguments are constants)

- Inlining is useful when:
 - the body of the procedure to be inlined is much shorter than the prologue/epilogue
 - inlining enables **specialization** (e.g., arguments are constants)
 - inlining enables other optimizations (e.g., part of the body is dead at this particular call site)

- Inlining is useful when:
 - the body of the procedure to be inlined is much shorter than the prologue/epilogue
 - inlining enables **specialization** (e.g., arguments are constants)
 - inlining enables other optimizations (e.g., part of the body is dead at this particular call site)
- Inlining also has risks:

- Inlining is useful when:
 - the body of the procedure to be inlined is much shorter than the prologue/epilogue
 - inlining enables **specialization** (e.g., arguments are constants)
 - inlining enables other optimizations (e.g., part of the body is dead at this particular call site)
- Inlining also has risks:
 - increases code size (may overflow instruction cache)

- Inlining is useful when:
 - the body of the procedure to be inlined is much shorter than the prologue/epilogue
 - inlining enables **specialization** (e.g., arguments are constants)
 - inlining enables other optimizations (e.g., part of the body is dead at this particular call site)
- Inlining also has risks:
 - increases code size (may overflow instruction cache)
 - increases register pressure

- Inlining is useful when:
 - the body of the procedu the prologue/epilogue

Note similar benefits and risks to **loop unrolling**. Deciding whether to inline is similarly complicated!

- inlining enables **specialization** (e.g., arguments are constants)
- inlining enables other optimizations (e.g., part of the body is dead at this particular call site)
- Inlining also has risks:
 - increases code size (may overflow instruction cache)
 - increases register pressure

• In practice, production compilers will use heuristics to decide when/if to inline, such as:

- In practice, production compilers will use heuristics to decide when/if to inline, such as:
 - Is this procedure a **leaf** in the call graph?
 - That is, does it not call any other procedures itself?

- In practice, production compilers will use heuristics to decide when/if to inline, such as:
 - Is this procedure a **leaf** in the call graph?
 - That is, does it not call any other procedures itself?
 - Is the callee procedure significantly smaller than the calling procedure?

- In practice, production compilers will use heuristics to decide when/if to inline, such as:
 - Is this procedure a **leaf** in the call graph?
 - That is, does it not call any other procedures itself?
 - Is the callee procedure significantly smaller than the calling procedure?
 - Static *call count*: the number of distinct sites that call the procedure.
 - Any procedure called just once is a good inlining candidate.

- In practice, production compilers will use heuristics to decide when/if to inline, such as:
 - Is this procedure a **leaf** in the call graph?
 - That is, does it not call any other procedures itself?
 - Is the callee procedure significantly smaller than the calling procedure?
 - Static *call count*: the number of distinct sites that call the procedure.
 - Any procedure called just once is a good inlining candidate.
 - **Profile data**, such as fraction of execution time (if available)

• Consider a procedure that calls another procedure and then immediately returns, like this example:

```
int foo(...):
 ...
 return bar(...)
```

• What will happen as **bar** returns?

• Consider a procedure that calls another procedure and then immediately returns, like this example:

```
int foo(...):
 ...
 return bar(...)
```

- What will happen as **bar** returns?
 - We will execute the epilogues of foo and bar in sequence, with no intervening instructions

• Consider a procedure that calls another procedure and then immediately returns, like this example:

```
int foo(...):
 ...
 return bar(...)
```

- What will happen as **bar** returns?
 - We will execute the epilogues of foo and bar in sequence, with no intervening instructions
 - Including many redundant operations (e.g., resetting %rsp)
• *Tail-call elimination* is an interprocedural optimization that allows a function called as the last instruction in a procedure to return to the calling procedure's caller directly

- *Tail-call elimination* is an interprocedural optimization that allows a function called as the last instruction in a procedure to return to the calling procedure's caller directly
 - This eliminates redundant operations in the epilogue

- *Tail-call elimination* is an interprocedural optimization that allows a function called as the last instruction in a procedure to return to the calling procedure's caller directly
 - This eliminates redundant operations in the epilogue
- This optimization is most important for *tail-recursive* procedures that call themselves as the last operation in their body
 e.g., imagine a naive Fibonacci implementation

- *Tail-call elimination* is an interprocedural optimization that allows a function called as the last instruction in a procedure to return to the calling procedure's caller directly
 - This eliminates redundant operations in the epilogue
- This optimization is most important for *tail-recursive* procedures that call themselves as the last operation in their body
 - e.g., imagine a naive Fibonacci implementation
 - Tail-call elimination often reduces asymptotic stack space requirements from linear to constant for tail-recursive calls

- *Tail-call elimination* is an interprocedural optimization that allows a function called as the last instruction in a procedure to return to the calling procedure's caller directly
 - This eliminates redundant operations in the epilogue
- This optimization is most important for *tail-recursive* procedures that call themselves as the last operation in their body
 - e.g., imagine a naive Fibonacci implementation
 - Tail-call elimination often reduces asymptotic stack space requirements from linear to constant for tail-recursive calls
- Functional languages practically require tail-call elimination

- The biggest difficulty in interprocedural optimization is maintaining support for **separate compilation**
 - Traditional "*compilation unit*" is a procedure or file

- The biggest difficulty in interprocedural optimization is maintaining support for separate compilation
 Traditional "compilation unit" is a procedure or file
 - Traditional "*compilation unit*" is a procedure or file
- If all of the code is definitely available, it suffices to **track dependencies** between procedures from an optimization perspective, and then re-optimize whenever a dependent procedure changes

- The biggest difficulty in interprocedural optimization is maintaining support for **separate compilation**
 - Traditional "*compilation unit*" is a procedure or file
- If all of the code is definitely available, it suffices to **track dependencies** between procedures from an optimization perspective, and then re-optimize whenever a dependent procedure changes
- Alternatively, we can defer interprocedural optimization until *link time*, when a *linker* combines the object files from each compilation unit into a single executable. We'll talk more about this next week.

Agenda

- Interprocedural optimizations (and analysis)
- Compiler backend overview
 - Instruction selection
 - Instruction scheduling
 - Register allocation basics
- Next time:
 - Deep dive into register allocation











• The *backend* of an optimizing compiler converts from optimized IR into optimized assembly for the target machine

- The *backend* of an optimizing compiler converts from optimized IR into optimized assembly for the target machine
- Three major subproblems:

- The *backend* of an optimizing compiler converts from optimized IR into optimized assembly for the target machine
- Three major subproblems:
 - Instruction selection: map IR into assembly code, combine low-level IR operations into machine instructions (take advantage of addressing modes, etc.)

- The *backend* of an optimizing compiler converts from optimized IR into optimized assembly for the target machine
- Three major subproblems:
 - Instruction selection: map IR into assembly code, combine low-level IR operations into machine instructions (take advantage of addressing modes, etc.)
 - Instruction scheduling: Reorder instructions to minimize execution time, hide latencies from processor function units, memory/cache stalls

- The *backend* of an optimizing compiler converts from optimized IR into optimized assembly for the target machine
- Three major subproblems:
 - Instruction selection: map IR into assembly code, combine low-level IR operations into machine instructions (take advantage of addressing modes, etc.)
 - Instruction scheduling: Reorder instructions to minimize execution time, hide latencies from processor function units, memory/cache stalls
 - Register allocation: Map abstract registers to actual registers, add code to spill values to memory and reload as needed, etc

- The *backend* of an optimizir IR into optimized assembly
- Three major subproblems:

 Instruction selection: ma low-level IR operations advantage of addressing modes, etc.)

 and then start on register

 allocation. Wednesday, we'll cover
 register allocation in more detail.
 - Instruction scheduling: Reorder instructions to minimize execution time, hide latencies from processor function units, memory/cache stalls

Today we will briefly cover

instruction selection + scheduling,

 Register allocation: Map abstract registers to actual registers, add code to spill values to memory and reload as needed, etc

- Goal: map IR to assembly code
 - assuming known storage layout and code shape

- Goal: map IR to assembly code
 - assuming known storage layout and code shape
- Problem to solve: given the low-level IR, there are many possible code sequences that implement it correctly

- Goal: map IR to assembly code
 - assuming known storage layout and code shape
- Problem to solve: given the low-level IR, there are many possible code sequences that implement it correctly
 - e.g., how many ways can we "set %rax to zero" in x86-64?

- Goal: map IR to assembly code
 - assuming known storage layout and code shape
- Problem to solve: given the low-level IR, there are many possible code sequences that implement it correctly
 - e.g., how many ways can we "set %rax to zero" in x86-64?

movq	\$0, %rax	salq	64, %rax
subq	%rax, %rax	shrq	64, %rax
xorq	%rax, %rax	imulq	\$0, %rax

- Goal: map IR to assembly code
 - assuming known storage layout and code shape
- Problem to solve: given the low-level IR, there are many possible code sequences that implement it correctly
 - e.g., how many ways can we "set %rax to zero" in x86-64?

movq	\$0, %rax	salq	64, %rax
subq	%rax, %rax	shrq	64, %rax
xorq	%rax, %rax	imulq	\$0, %rax

• Many machine instructions also do **several things at once** (e.g., register arithmetic and effective address calculation in a movq)

• Several possibilities: **fastest**, smallest, minimize power consumption (e.g., don't use a functional unit if leaving it powered-down is a win), reduce memory traffic, etc.

- Several possibilities: fastest, smallest, minimize power consumption (e.g., don't use a functional unit if leaving it powered-down is a win), reduce memory traffic, etc.
 - Typically we want "fastest", but depends on opt. target

- Several possibilities: **fastest**, smallest, minimize power consumption (e.g., don't use a functional unit if leaving it powered-down is a win), reduce memory traffic, etc.
 - Typically we want "fastest", but depends on opt. target
- Sometimes not obvious
 - e.g., if one of the function units in the processor is idle and we can select an instruction that uses that unit, it effectively executes for free, even if that instruction wouldn't be chosen normally
 - (Some interaction with scheduling here...)
 - (and it might consume extra power, so bad if that matters)

• One algorithm for instruction selection is tree pattern matching

- One algorithm for instruction selection is tree pattern matching
- Goal: find a sequence of machine instructions that perform the computation described by the program's IR code

- One algorithm for instruction selection is tree pattern matching
- Goal: find a sequence of machine instructions that perform the computation described by the program's IR code
- Algorithm:

- One algorithm for instruction selection is tree pattern matching
- Goal: find a sequence of machine instructions that perform the computation described by the program's IR code
- Algorithm:
 - Describe each machine instruction we want to consider using same low-level IR used for program

- One algorithm for instruction selection is tree pattern matching
- Goal: find a sequence of machine instructions that perform the computation described by the program's IR code
- Algorithm:
 - Describe each machine instruction we want to consider using same low-level IR used for program
 - *Tile* the low-level IR tree with operation (instruction) trees
Instruction Selection: Tree Pattern Matching

- One algorithm for instruction selection is tree pattern matching
- Goal: find a sequence of machine instructions that perform the computation described by the program's IR code
- Algorithm:
 - Describe each machine instruction we want to consider using same low-level IR used for program
 - *Tile* the low-level IR tree with operation (instruction) trees
 - A tiling "*implements*" a tree if it covers every node in the tree and the overlap between any two tiles (trees) is limited to a single compatible node

• Two common algorithms to generate tilings:

- Two common algorithms to generate tilings:
 - Maximal munch:

- Two common algorithms to generate tilings:
 - **Maximal munch**:
 - Top-down tree walk

- Two common algorithms to generate tilings:
 - Maximal munch:
 - Top-down tree walk
 - Find largest tile that fits each node

- Two common algorithms to generate tilings:
 - Maximal munch:
 - Top-down tree walk
 - Find largest tile that fits each node
 - Why largest? Heuristic: One instruction that "does more" is likely cheaper than several that do less

- Two common algorithms to generate tilings:
 - Maximal munch:
 - Top-down tree walk
 - Find largest tile that fits each node
 - Why largest? Heuristic: One instruction that "does more" is likely cheaper than several that do less
 - **Dynamic programming:**
 - Assign costs to each node in the tree using a cost model

- Two common algorithms to generate tilings:
 - Maximal munch:
 - Top-down tree walk
 - Find largest tile that fits each node
 - Why largest? Heuristic: One instruction that "does more" is likely cheaper than several that do less
 - **Dynamic programming:**
 - Assign costs to each node in the tree using a cost model
 - cost = cost of individual node + subtree costs

- Two common algorithms to generate tilings:
 - Maximal munch:
 - Top-down tree walk
 - Find largest tile that fits each node
 - Why largest? Heuristic: One instruction that "does more" is likely cheaper than several that do less
 - **Dynamic programming:**
 - Assign costs to each node in the tree using a cost model
 - cost = cost of individual node + subtree costs
 - Try all possible combinations bottom-up, pick cheapest

- Two common algorithms to generate tilings:
 - Maximal munch:
 - Top-down tree walk
 - Find largest tile that fits each node
 - Why largest? Heuristic: One instruction that "does more" is likely cheaper than several that do less
 - **Dynamic programming:**
 - Assign costs to each node in the tree using a cost model
 - cost = cost of individual node + subtree costs
 - Try all possible combinations bottom-up, pick cheapest
 - Slower, but optimal for a given cost model

Instruction Selection: T.P.M. Example



Tree for a[i] := x

Instruction Selection: T.P.M. Example



Tree for a[i] := x

Given a tiled tree, to generate code:

Given a tiled tree, to generate code:

• Do a postorder tree walk with node-dependant order for children

Given a tiled tree, to generate code:

- Do a postorder tree walk with node-dependant order for children
- Each tile corresponds to a code sequence; emit code sequences in order

Given a tiled tree, to generate code:

- Do a postorder tree walk with node-dependant order for children
- Each tile corresponds to a code sequence; emit code sequences in order
- Connect tiles by using same register (or temporary) name to tie boundaries together

Instruction Selection: T.P.M. Example



LOAD r1 <- M[fp+a_{off}]
 ADDI r2 <- 4 + r0
 MUL r2 <- r2 * r_i
 ADD r1 <- r1 + r2
 LOAD r2 <- M[fp+x_{off}]
 STORE M[r1+0] <- r2

Tree for a[i] := x

Trivia Break: Computer Science

This American engineer, inventor and science administrator joined MIT in 1919. One of his major scientific works was a differential analyzer, a mechanical analog computer with some digital components that could solve differential equations with as many as 18 independent variables. However, he is best-known for his work in scientific administration: he was vice-president and dean of MIT's School of Engineering, and then became the head of the U.S. Office of Scientific Research and Development during the second world war. He was also instrumental in the founding of the National Science Foundation, and he founded the company that eventually became Raytheon while he was at MIT.

Trivia Break: Computer Science

Vannevar Bush wrote a 1945 article about this hypothetical electromechanical device for interacting with microform documents. Bush envisioned that individuals would compress and store all of their books, records, and communications in this device, "mechanized so that it may be consulted with exceeding speed and flexibility". The individual was supposed to use it as an automatic personal filing system, making the it "an enlarged intimate supplement to his memory." The concept influenced the development of early hypertext systems, eventually leading to the creation of the World Wide Web, and personal knowledge base software.

- Goal: reorder instructions to minimize execution time given instruction and operand latencies
 - Assume fixed program code at this point

- Goal: reorder instructions to minimize execution time given instruction and operand latencies
 - \circ $% \ensuremath{\mathsf{Assume}}$ Assume fixed program code at this point
- Why?

- Goal: reorder instructions to minimize execution time given instruction and operand latencies
 - Assume fixed program code at this point
- Why?
 - Many operations have non-zero latencies
 - Modern machines can issue several operations per cycle
 - Want to take advantage of multiple function units on chip
 - Loads and stores may or may not block
 - Branch costs vary
 - Want to help out the branch predictor, if we can

Instruction Scheduling: Example

Instruction Scheduling: Example: Latencies

Operation	Cycles
LOAD	3
STORE	3
ADD	1
MULT	2
SHIFT	1
BRANCH	0 TO 8

(Note that these are just simplified examples; a real machine's behavior will be much more complex!)

Instruction Scheduling: Example

Consider two schedules for the expression w*2*x*y*z: Operation Cycles LOAD 3 3 STORE ADD 1 2 MULT SHIFT 1 BRANCH 0 TO 8

Instruction Scheduling: Example

Consider two schedules for the expression w*2*x*y*z:

Operation	Cycles
LOAD	3
STORE	3
ADD	1
MULT	2
SHIFT	1
BRANCH	0 TO 8

Simple schedule			
1	LOAD	r1 <- w	
4	ADD	r1 <- r1,r1	
5	LOAD	r2 <- x	
8	MULT	r1 <- r1,r2	
9	LOAD	r2 <- y	
12	2 MULT	r1 <- r1,r2	
13	B LOAD	r2 <- z	
16	5 MULT	r1 <- r1,r2	
18	STORE	w <- r1	
21	r1 free		
2 registers, 20 cycles			

Instruction

Quick in-class exercise: turn to someone near you and use your big human brain to come up with a **better schedule**.

Consider two schedules for the expression w*2*x*y*z:

Operation	Cycles
LOAD	3
STORE	3
ADD	1
MULT	2
SHIFT	1
BRANCH	0 TO 8

1 LOAD	r1 <- w	
4 ADD	r1 <- r1,r1	
5 LOAD	r2 <- x	
8 MULT	r1 <- r1,r2	
9 LOAD	r2 <- y	
12 MULT	r1 <- r1,r2	
13 LOAD	r2 <- z	
16 MULT	r1 <- r1,r2	
18 STORE	w <- r1	
21 r1 free		
2 registers, 20 cycles		

Instruction Scheduling: Example

Consider two schedules for the expression w*2*x*y*z:

Operation	Cycles
LOAD	3
STORE	3
ADD	1
MULT	2
SHIFT	1
BRANCH	0 TO 8

Simple schedule		
1	LOAD	r1 <- w
4	ADD	r1 <- r1,r1
5	LOAD	r2 <- x
8	MULT	r1 <- r1,r2
9	LOAD	r2 <- y
12	MULT	r1 <- r1,r2
13	LOAD	r2 <- z
16	MULT	r1 <- r1,r2
18	STORE	w <- r1
21	r1 free	
2 registers, 20 cycles		

Loads early

	1	LOAD	r1 <- w
	2	LOAD	r2 <- x
	3	LOAD	r3 <- y
	4	ADD	r1 <- r1,r1
	5	MULT	r1 <- r1,r2
	6	LOAD	r2 <- z
	7	MULT	r1 <- r1,r3
	9	MULT	r1 <- r1,r2
	11	STORE	w <- r1
14 r1 is free			
3 registers, 13 cycles			

List Scheduling

• List scheduling constructs a schedule, one cycle at a time, by doing a (mildly modified) topological sort of a precedence graph

- List scheduling constructs a schedule, one cycle at a time, by doing a (mildly modified) topological sort of a precedence graph
 - in the precedence graph:
 - nodes are operations
 - edges are dependencies: if the operation at node n2 uses the result of the operation at n1, then (n1, n2) is an edge

- List scheduling constructs a schedule, one cycle at a time, by doing a (mildly modified) topological sort of a precedence graph
 - in the precedence graph:
 - nodes are operations
 - edges are dependencies: if the operation at node n2 uses the result of the operation at n1, then (n1, n2) is an edge
 - use *priority* to choose among ready (=in-degree 0) operations
 - priority = number of cycles on critical path to the end (usually)

- List scheduling constructs a schedule, one cycle at a time, by doing a (mildly modified) topological sort of a precedence graph
 - in the precedence graph:
 - nodes are operations
 - edges are dependencies: if the operation at node n2 uses the result of the operation at n1, then (n1, n2) is an edge
 - use *priority* to choose among ready (=in-degree 0) operations
 - priority = number of cycles on critical path to the end (usually)
- Note: may need to rename registers to avoid false dependencies and conflicts

- List scheduling constructs Full algorithm if you want to implement it yourself: (mildly modified) topolog P = precedence graph:
 - in the precedence gra
 - nodes are operati
 - edges are depend the result of the c
 - use *priority* to choose
 - priority = number
 (usually)
- Note: may need to renam and conflicts

P = precedence graph; Cycle = 1; Ready = leaves of P; Active = empty; while (Ready and/or Active are not empty) if (Ready is not empty) remove an op from Ready; S(op) = Cycle;Active = Active È op; Cycle++: for each op in Active if $(S(op) + delay(op) \le Cycle)$ remove op from Active; for each successor s of op in P if (s is ready - i.e., all operands available) add s to Ready

ga

List Scheduling: Example

Code

a LOAD r1 <- w r1 <- r1,r1 b ADD c LOAD r2 <- x d MULT r1 <- r1,r2 e LOAD r2 <- y f MULT r1 <- r1,r2 g LOAD r2 <- z r1 <- r1,r2 h MULT i STORE w <- r1


g ⁸



Code 13 a a LOAD r1 <- w 10_h c 12 r1 <- r1,r1 b ADD 9 e 10 c LOAD r2 <- x d MULT r1 <- r1,r2 g ⁸ e LOAD r2 <- y 5 h f MULT r1 <- r1,r2 g LOAD r2 <- z 3 h MULT r1 <- r1,r2 cycle: 1 i STORE w <- r1 ready: a c e g active: -

instr done1 a LOAD 4

Code		¹³ a
a LOAD	r1 <- w	10 12
b ADD	r1 <- r1,r1	
c LOAD	r2 <- x	⁹ d e ¹⁰
d MULT	r1 <- r1,r2	7.5 - 8
e LOAD	r2 <- y	r g c
f MULT	r1 <- r1,r2	⁵ h
g LOAD	r2 <- z	3
h MULT	r1 <- r1,r2	
i STORE	w <- r1	cycle: 1 2 ready: a c e g active: a

- # instr done
- 1 a LOAD 4
- 2 cLOAD 5

Code		¹³ a
a LOAD	r1 <- w	10 12
b ADD	r1 <- r1,r1	
c LOAD	r2 <- x	⁹ d e ¹⁰
d MULT	r1 <- r1,r2	7. 8
e LOAD	r2 <- y	f g
f MULT	r1 <- r1,r2	⁵ h
g LOAD	r2 <- z	3
h MULT	r1 <- r1,r2	
i STORE	w <- r1	cycle: 1 2 3
		ready: a c e g
		active: a c

- # instr done
- 1 a LOAD 4
- 2 cLOAD 5
- 3 e LOAD 6

	¹³ a
r1 <- w	10 12
r1 <- r1,r1	
r2 <- x	⁹ d e ¹⁰
r1 <- r1,r2	7.5 - 8
r2 <- y	r gc
r1 <- r1,r2	⁵ h
r2 <- z	3
r1 <- r1,r2	
w <- r1	cycle: 1 2 3 4
	ready: a c e g b
	active: a c e
	r1 <- w r1 <- r1,r1 r2 <- x r1 <- r1,r2 r2 <- y r1 <- r1,r2 r2 <- z r1 <- r1,r2 w <- r1

- # instr done
- 1 a LOAD 4
- 2 c LOAD 5 3 e LOAD 6
- 4 b ADD 5

Code		¹³ a
a LOAD	r1 <- w	10 12
b ADD	r1 <- r1,r1	
c LOAD	r2 <- x	⁹ d e ¹⁰
d MULT	r1 <- r1,r2	7
e LOAD	r2 <- y	r g c
f MULT	r1 <- r1,r2	⁵ h
g LOAD	r2 <- z	3
h MULT	r1 <- r1,r2	
i STORE	w <- r1	cycle: 1 2 3 4 5
		ready: acegbd
		active: a c e b

- # instr done
- 1 a LOAD 4
- 2 c LOAD 5 3 e LOAD 6
- 4 b ADD 5
- 5 d MULT 7

Code		¹³ a
a LOAD	r1 <- w	10 12
b ADD	r1 <- r1,r1	
c LOAD	r2 <- x	⁹ d e ¹⁰
d MULT	r1 <- r1,r2	7 8
e LOAD	r2 <- y	f g
f MULT	r1 <- r1,r2	⁵ h
g LOAD	r2 <- z	3
h MULT	r1 <- r1,r2	
i STORE	w <- r1	cycle: 123456
		ready: acegbd
		active: a c e b d

- # instr done
- 1 a LOAD 4
- 2 c LOAD 5 3 e LOAD 6
- 4 b ADD 5
- 5 d MULT 7
 - 6 gLOAD 9

Со	de		¹³ a
а	LOAD	r1 <- w	10 12
b	ADD	r1 <- r1,r1	
С	LOAD	r2 <- x	⁹ d e ¹⁰
d	MULT	r1 <- r1,r2	7.6 - 8
е	LOAD	r2 <- y	r g c
f	MULT	r1 <- r1,r2	⁵ h
g	LOAD	r2 <- z	3
h	MULT	r1 <- r1,r2	
i	STORE	w <- r1	cycle: $\frac{123456}{7}$ ready: $\frac{123456}{7}$ active: $\frac{123456}{7}$

- # instr done1 a LOAD 42 c LOAD 5
- 3 e LOAD 6 4 b ADD 5
- 5 d MULT 7
- 6 gLOAD 9 7 fMULT 9

Code		¹³ a
a LOAD	r1 <- w	10 12
b ADD	r1 <- r1,r1	
c LOAD	r2 <- x	⁹ d e ¹⁰
d MULT	r1 <- r1,r2	7.6 8
e LOAD	r2 <- y	r g c
f MULT	r1 <- r1,r2	⁵ h
g LOAD	r2 <- z	3
h MULT	r1 <- r1,r2	
i STORE	w <- r1	cycle: 1234567 8 ready: acegbdf active: acebd gf

1	a LOAD	4
2	c LOAD	5
3	e LOAD	6
4	b ADD	5
5	d MULT	7
6	g LOAD	9
7	f MULT	9

7 f MULT 8 – (stall)

instr

#

done

Code		¹³ a
a LOAD	r1 <- w	10, 12
b ADD	r1 <- r1,r1	
c LOAD	r2 <- x	⁹ d e ¹⁰
d MULT	r1 <- r1,r2	7.6
e LOAD	r2 <- y	r g
f MULT	r1 <- r1,r2	⁵ h
g LOAD	r2 <- z	3
h MULT	r1 <- r1,r2	
i STORE	w <- r1	cycle: 12345678 9 ready: acegbdf h active: acebdgf

1	a LOAD
2	c LOAD
3	e LOAD
4	b ADD
5	d MULT
6	g LOAD
7	f MULT
8	– (stall)
9	h MULT

#

instr

done

Code		¹³ a
a LOAD	r1 <- w	10, 12
b ADD	r1 <- r1,r1	
c LOAD	r2 <- x	⁹ d e ¹⁰
d MULT	r1 <- r1,r2	7.6 - 8
e LOAD	r2 <- y	, r g
f MULT	r1 <- r1,r2	⁵ h
g LOAD	r2 <- z	3
h MULT	r1 <- r1,r2	
i STORE	w <- r1	cycle: 123456789 10 ready: acegbdfh active: acebdgf h

#	instr	done
1	a LOAD	4
2	c LOAD	5
3	e LOAD	6
4	b ADD	5
5	d MULT	7
6	g LOAD	9
7	f MULT	9
8	– (stall)	
9	h MUIT	11

10 - (stall)

Code		¹³ a	2	c LOAD
a LOAD	r1 <- w	10 12	3	e LOAD
b ADD	r1 <- r1,r1		4	b ADD
c LOAD	r2 <- x	⁹ d e ¹⁰	5	d MULT
d MULT	r1 <- r1,r2	7.6	6	g LOAD
e LOAD	r2 <- y	r g	7	f MULT
f MULT	r1 <- r1,r2	⁵ h	8	– (stall)
g LOAD	r2 <- z	3	9	h MULT
h MULT	r1 <- r1,r2		10	– (stall)
i STORE	w <- r1	cycle: 12345678910 11 ready: acegbdfh i active: acebdgfh	11	i STORE

#

instr

a LOAD

done

• Alternative: **backwards** list scheduling

- Alternative: **backwards** list scheduling
 - Work from the root to the leaves

- Alternative: **backwards** list scheduling
 - Work from the root to the leaves
 - Schedules instructions from end to beginning of the block

- Alternative: **backwards** list scheduling
 - Work from the root to the leaves
 - Schedules instructions from end to beginning of the block
- In practice, production compilers typically **try both** and pick the result that minimizes costs

- Alternative: **backwards** list scheduling
 - Work from the root to the leaves
 - Schedules instructions from end to beginning of the block
- In practice, production compilers typically **try both** and pick the result that minimizes costs
 - Little extra expense since the precedence graph and other information can be reused

- Alternative: **backwards** list scheduling
 - Work from the root to the leaves
 - Schedules instructions from end to beginning of the block
- In practice, production compilers typically **try both** and pick the result that minimizes costs
 - Little extra expense since the precedence graph and other information can be reused
 - Different directions win in different cases

• It's possible to do regional or even global list scheduling

- It's possible to do regional or even global list scheduling
 - However, it's much more complicated and gains are often small

- It's possible to do regional or even global list scheduling
 - However, it's much more complicated and gains are often small
 - e.g., if you schedule an entire EBB, you need a heuristic to estimate the cost of "infinite" paths around loops
 - could e.g., assume all loops execute 10 times

- It's possible to do regional or even global list scheduling
 - However, it's much more complicated and gains are often small
 - e.g., if you schedule an entire EBB, you need a heuristic to estimate the cost of "infinite" paths around loops
 - could e.g., assume all loops execute 10 times
- Some compilers use **profiling information** for scheduling
 - i.e., run the code and see how many cycles different schedules actually take

- It's possible to do regional or even global list scheduling
 - However, it's much more complicated and gains are often small
 - e.g., if you schedule an entire EBB, you need a heuristic to estimate the cost of "infinite" paths around loops
 - could e.g., assume all loops execute 10 times
- Some compilers use **profiling information** for scheduling
 - i.e., run the code and see how many cycles different schedules actually take
 - downside: need to actually run the code
 - e.g., where do you get inputs?

- IR code typically assumes **infinitely-many registers** are available
 - but real machines only have a small number of registers :(

- IR code typically assumes **infinitely-many registers** are available
 - but real machines only have a small number of registers :(
- Task of the *register allocator*: create a mapping from IR's abstract registers to physical registers

- IR code typically assumes **infinitely-many registers** are available
 - but real machines only have a small number of registers :(
- Task of the *register allocator*: create a mapping from IR's abstract registers to physical registers
 - Or, if that's not possible, to memory locations
 - this happens when there aren't enough physical registers

- IR code typically assumes **infinitely-many registers** are available
 - but real machines only have a small number of registers :(
- Task of the *register allocator*: create a mapping from IR's abstract registers to physical registers
 - Or, if that's not possible, to memory locations
 - this happens when there aren't enough physical registers
 - Insert code to move values between registers and memory if needed ("*spill code*")

- IR code typically assumes **infinitely-many registers** are available
 - but real machines only have a small number of registers :(
- Task of the *register allocator*: create a mapping from IR's abstract registers to physical registers
 - Or, if that's not possible, to memory locations
 - this happens when there aren't enough physical registers
 - Insert code to move values between registers and memory if needed ("*spill code*")
 - Typically we will re-run the instruction scheduler if we ever have to spill a register

• Even your PA3 implementation should probably have a simple "register allocator" somewhere

- Even your PA3 implementation should probably have a simple "register allocator" somewhere
 - my suggestion: "everything goes in memory"

- Even your PA3 implementation should probably have a simple "register allocator" somewhere
 - my suggestion: "everything goes in memory"
 - then you can insert load code before any operation that requires its operands in registers without worrying about what's in those registers when you do

- Even your PA3 implementation should probably have a simple "register allocator" somewhere
 - my suggestion: "everything goes in memory"
 - then you can insert load code before any operation that requires its operands in registers without worrying about what's in those registers when you do
- However, one of the first optimizations that I suggest you implement for PA4 is a simple register allocator

- Even your PA3 implementation should probably have a simple "register allocator" somewhere
 - my suggestion: "everything goes in memory"
 - then you can insert load code before any operation that requires its operands in registers without worrying about what's in those registers when you do
- However, one of the first optimizations that I suggest you implement for PA4 is a simple register allocator
 - avoiding memory operations is extremely profitable 44 44 44

- Even your PA3 implementation should probably have a simple "register allocator" somewhere
 - my suggestion: "everything goes in memory"
 - then you can insert load code before any operation that requires its operands in registers without worrying about what's in those registers when you do
- However, one of the first optimizations that I suggest you implement for PA4 is a simple register allocator
 - avoiding memory operations is extremely profitable 44 44 44
 - I recommend waiting until later in PA4 to try more complex schemes
- Register allocation: correctness or optimization
 - depends on memory model

- Register allocation: correctness or optimization
 - depends on memory model
- Local register allocation
 - two approaches: top-down and bottom-up

- Register allocation: correctness or optimization
 - depends on memory model
- Local register allocation
 - two approaches: top-down and bottom-up
- Regional register allocation + complications
 - **liveness analysis** strikes again!

- Register allocation: correctness or optimization
 - depends on memory model
- Local register allocation
 - two approaches: top-down and bottom-up
- Regional register allocation + complications
 - **liveness analysis** strikes again!
- Global register allocation via reduction to graph coloring
 o including a union-find algorithm for fun

Course Announcements

- PA3 deadline is **today** (AoE)
 - How is it going?
- I will hold an extra office hour today 11:30-12:30 for those who would like to see a PA3 test case
 - You may also be able to catch me either between 2 and 2:30 in my office or at the CS seminar this afternoon, but no promises
- PA4 is still Coming Soon[™] (I'm actually trying to get this autograder right the first time...)
- Note that PA4c1's specification is TAC -> TAC
 - that is, the input is also a .cl-tac file
 - PA4c1 is due April 28, and is mostly optional