Code Generation

Martin Kellogg

Course Announcements

- PA3c2 (TAC) is due later this week ("before spring break")
 - you should already have started, or you are behind
 - if there is demand from the class, I will consider a short extension on this assignment to e.g., Monday
- I have become aware of a **bug in the reference compiler**'s x86-64 module; a fix is forthcoming. For now, don't trust it.
- Don't forget there is a midterm in this class the week after spring break!

• In this lecture (and, unfortunately, the one on the Monday after spring break...) I am going to show you how to do code generation for a stack machine ("Cool ASM"/"Cool bytecode")

- In this lecture (and, unfortunately, the one on the Monday after spring break...) I am going to show you how to do code generation for a stack machine ("Cool ASM"/"Cool bytecode")
 - Including all of the complexity of handling objects, etc.

- In this lecture (and, unfortunately, the one on the Monday after spring break...) I am going to show you how to do code generation for a stack machine ("Cool ASM"/"Cool bytecode")
 - Including all of the complexity of handling objects, etc.
- A stack machine is simpler than the x86-64 we saw on Monday
 This is on purpose (where's the fun if I make it too easy...)

- In this lecture (and, unfortunately, the one on the Monday after spring break...) I am going to show you how to do code generation for a stack machine ("Cool ASM"/"Cool bytecode")
 - Including all of the complexity of handling objects, etc.
- A stack machine is simpler than the x86-64 we saw on Monday
 This is on purpose (where's the fun if I make it too easy...)
- Hard PA3 task for you: merge...

- In this lecture (and, unfortunately, the one on the Monday after spring break...) I am going to show you how to do code generation for a stack machine ("Cool ASM"/"Cool bytecode")
 - Including all of the complexity of handling objects, etc.
- A stack machine is simpler than the x86-64 we saw on Monday
 This is on purpose (where's the fun if I make it too easy...)
- Hard PA3 task for you: merge...
 - ...what you know about x86-64 (i.e., last lecture)

- In this lecture (and, unfortunately, the one on the Monday after spring break...) I am going to show you how to do code generation for a stack machine ("Cool ASM"/"Cool bytecode")
 - Including all of the complexity of handling objects, etc.
- A stack machine is simpler than the x86-64 we saw on Monday
 This is on purpose (where's the fun if I make it too easy...)
- Hard PA3 task for you: merge...
 - ...what you know about x86-64 (i.e., last lecture)
 - ...with what you know about code generation for a simpler target (this lecture)

- Stack machine basics
 - accumulator, stack pointer

- Stack machine basics
 - accumulator, stack pointer
- Stack discipline, calling convention for our stack machine
 - with a bit of optimization thrown in to give you a taste of the idea

- Stack machine basics
 - accumulator, stack pointer
- Stack discipline, calling convention for our stack machine
 - with a bit of optimization thrown in to give you a taste of the idea
- Object layout
 - \circ $\;$ and its interactions with subtyping $\;$

- Stack machine basics
 - accumulator, stack pointer
- Stack discipline, calling convention for our stack machine
 - with a bit of optimization thrown in to give you a taste of the idea
- Object layout
 - and its interactions with subtyping
- Dispatch tables/vtables
 - this gets us to dynamic dispatch

- Stack machine basics
 - accumulator, stack pointer
- Stack discipline, calling convention for our stack machine
 - with a bit of optimization thrown in to give you a taste of the idea
- Object layout
 - \circ and its interactions with φ
- Dispatch tables/vtables
 - \circ this gets us to dynamic di

Q: What's on the midterm?

- Stack machine basics
 - accumulator, stack pointer
- Stack discipline, calling convention for our stack machine
 - with a bit of optimization thrown in to give you a taste of the idea
- Object layout
 - \circ and its interactions with φ
- Dispatch tables/vtables
 - \circ this gets us to dynamic di

Q: What's on the midterm? A: **Only what we get through today**. I probably won't ask much about it.

- A simple evaluation model
 - No variables or registers

- A simple evaluation model
 - No variables or registers
- A stack of values for intermediate results
 - Imagine a calculator that uses "reverse Polish" notation

- A simple evaluation model
 - No variables or registers
- A stack of values for intermediate results
 - Imagine a calculator that uses "reverse Polish" notation
- Think of it this way:
 - I am going to explain how you'd build "Cool javac"; your compiler is supposed to be "Cool g++".

- A simple evaluation model
 - No variables or registers
- A stack of values for intermediate results
 - Imagine a calculator that uses "reverse Polish" notation
- Think of it this way:
 - I am going to explain how you'd build "Cool javac"; your compiler is supposed to be "Cool g++".
 - "Cool javac" : Cool -> stack machine bytecode
 - A "JVM" is also necessary to interpret the bytecode!

- A simple evaluation model
 - No variables or registers
- A stack of values for intermediate results
 - Imagine a calculator that uses "reverse Polish" notation
- Think of it this way:
 - I am going to explain how you'd build "Cool javac"; your compiler is supposed to be "Cool g++".
 - "Cool javac" : Cool -> stack machine bytecode
 - A "JVM" is also necessary to interpret the bytecode!
 - "Cool g++": Cool -> x86-64

• Consider two instructions:

- Consider two instructions:
 - o push i
 - place the integer i on top of the stack

- Consider two instructions:
 - o push i
 - place the integer i on top of the stack
 - o add
 - pop two elements, add them and put the result back on the stack

- Consider two instructions:
 - o push i
 - place the integer i on top of the stack
 - o add
 - pop two elements, add them and put the result back on the stack
- A program to compute 7 + 5:
 - push 7
 - push 5
 - add

- Consider two instructions:
 - o push i
 - place the integer i on top of the stack
 - o add
 - pop two elements, add them and put the result back on the stack
- A program to compute 7 + 5: push 7 push 5 add

Quick poll: how much does this have in common with the 280 project? (determines how fast I go)

• Each instruction:

- Each instruction:
 - Takes its operands from the top of the stack

- Each instruction:
 - Takes its operands from the top of the stack
 - Removes those operands from the stack

- Each instruction:
 - Takes its operands from the top of the stack
 - Removes those operands from the stack
 - Computes the required operation on them

- Each instruction:
 - Takes its operands from the top of the stack
 - Removes those operands from the stack
 - Computes the required operation on them
 - Pushes the result on the stack

stack

• Each instruction:

- Takes its operands from the top of the stack
- Removes those operands from the stack
- Computes the required operation on them
- Pushes the result on the stack



push 7

- Each instruction:
 - Takes its operands from the top of the stack
 - Removes those operands from the stack
 - Computes the required operation on them
 - Pushes the result on the stack



push 7 push 5

- Each instruction:
 - Takes its operands from the top of the stack
 - Removes those operands from the stack
 - Computes the required operation on them
 - Pushes the result on the stack



- Each instruction:
 - Takes its operands from the top of the stack
 - Removes those operands from the stack
 - Computes the required operation on them
 - Pushes the result on the stack

Stack Machines: Good For Compiler Writers

• Stack machines have a very nice property from our perspective as compiler writers:

Stack Machines: Good For Compiler Writers

- Stack machines have a very nice property from our perspective as compiler writers:
 - Each operation takes operands from the same place and puts results in the same place
- Stack machines have a very nice property from our perspective as compiler writers:
 - Each operation takes operands from the same place and puts results in the same place
- This means a **uniform** compilation scheme

- Stack machines have a very nice property from our perspective as compiler writers:
 - Each operation takes operands from the same place and puts results in the same place
- This means a **uniform** compilation scheme
 - We don't have to worry about where operands are ("they're on the stack")

- Stack machines have a very nice property from our perspective as compiler writers:
 - Each operation takes operands from the same place and puts results in the same place
- This means a **uniform** compilation scheme
 - We don't have to worry about where operands are ("they're on the stack")
- And therefore means we can write a very simple compiler

- Stack machines have a very nice property from our perspective as compiler writers:
 - Each operation takes operands from the same place and puts results in the same place
- This means a **uniform** compilation scheme
 - We don't have to worry about where operands are ("they're on the stack")
- And therefore means we can write a very simple compiler
 - Many/most compilers start by targeting a stack machine

- Stack machines have a very nice property from our perspective as compiler writers:
 - Each operation takes operands from the same place and puts results in the same place
- This means a **uniform** compilation scheme
 - We don't have to worry about where operands are ("they're on the stack")
- And therefore means we can write a very simple compiler
 - Many/most compilers start by targeting a stack machine
 - A good idea for you, too! (in PA3)

- Location of the operands is implicit
 - Always on the top of the stack

- Location of the operands is implicit
 - \circ Always on the top of the stack
- No need to specify operands explicitly

- Location of the operands is implicit
 - \circ Always on the top of the stack
- No need to specify operands explicitly
- No need to specify the location of the result

- Location of the operands is implicit
 - Always on the top of the stack
- No need to specify operands explicitly
- No need to specify the location of the result
- Instruction "add" as opposed to "add r1, r2"
 - Smaller encoding of instructions
 - More compact programs (great for network!)

- Location of the operands is implicit
 - \circ Always on the top of the stack
- No need to specify operands explicitly
- No need to specify the location of the result
- Instruction "add" as opposed to "add r1, r2"
 - Smaller encoding of instructions
 - More compact programs (great for network!)
- This is one reason why Java bytecode uses a stack evaluation model

- The add instruction does 3 memory operations
 - Two reads and one write to the stack

- The add instruction does 3 memory operations
 - Two reads and one write to the stack
 - This is very slow!



- The add instruction does 3 memory operations
 - Two reads and one write to the stack
 - This is very slow!
 - The top of the stack is "frequently" accessed (so much)



- The add instruction does 3 memory operations
 - Two reads and one write to the stack
 - This is very slow!
 - The top of the stack is "frequently" accessed (so much)
- We can do one simple optimization to mitigate this disadvantage



- The add instruction does 3 memory operations
 - Two reads and one write to the stack
 - This is very slow!
 - The top of the stack is "frequently" accessed (so much)
- We can do one simple optimization to mitigate this disadvantage
 - Bonus pedagogical benefit: I can mention again that you shouldn't prematurely optimize on PA3



- The add instruction does 3 memory operations
 - Two reads and one write to the stack
 - This is very slow!
 - The top of the stack is "frequently" accessed (so much)
- We can do one simple optimization to mitigate this disadvantage
 - Bonus pedagogical benefit: I can mention again that you shouldn't prematurely optimize on PA3
- Can anyone guess what it is?



• Key insight: keep the top of the stack in a register (called the *accumulator*)

- Key insight: keep the top of the stack in a register (called the *accumulator*)
 - This should remind you of how fold works
 - It's almost like functional programming is relevant to the rest of the class...

- Key insight: keep the top of the stack in a register (called the *accumulator*)
 - This should remind you of how fold works
 - It's almost like functional programming is relevant to the rest of the class...
 - Register accesses are much faster

- Key insight: keep the top of the stack in a register (called the *accumulator*)
 - This should remind you of how fold works
 - It's almost like functional programming is relevant to the rest of the class...
 - Register accesses are much faster
- The add instruction is now acc <- acc + top_of_stack
 - Only one memory operation, much faster!

• The result of computing an expression is **always** in the accumulator

- The result of computing an expression is always in the accumulator
- For an operation op (e₁, ..., e_n), push the accumulator on the stack after computing each of e₁, ..., e_{n-1}

- The result of computing an expression is always in the accumulator
- For an operation op (e₁, ..., e_n), push the accumulator on the stack after computing each of e₁, ..., e_{n-1}
 - \circ e_n's result is in the accumulator before op

- The result of computing an expression is **always** in the accumulator
- For an operation op (e₁, ..., e_n), push the accumulator on the stack after computing each of e₁, ..., e_{n-1}
 - e_n's result is in the accumulator before op
 - After the operation **pop** n-1 values

- The result of computing an expression is **always** in the accumulator
- For an operation op (e₁, ..., e_n), push the accumulator on the stack after computing each of e₁, ..., e_{n-1}
 - e_n's result is in the accumulator before op
 - After the operation **pop** n-1 values
- After computing an expression the stack is as before

- The result of computing an expression is **always** in the accumulator
- For an operation op (e₁, ..., e_n), push the accumulator on the stack after computing each of e₁, ..., e_{n-1}
 - \circ e_n's result is in the accumulator before op
 - After the operation **pop** n-1 values
- After computing an expression the stack is as before

Let's look at an example...

• Computing 7 + 5 with an accumulator:







•••

• Computing 7 + 5 with an accumulator:





acc <- 7 push acc

• Computing 7 + 5 with an accumulator:



• Computing 7 + 5 with an accumulator:



Stack Machines: Bigger Example: 3 + (7 + 5)

<u>Code</u> <u>Acc</u> <u>Stack</u>

Stack Machines: Bigger Example: 3 + (7 + 5)

<u>Code</u>	<u>Acc</u>	<u>Stack</u>
acc <- 3	3	<init></init>

Stack Machines: Bigger Example: 3 + (7 + 5)

<u>Code</u>	<u>Acc</u>	<u>Stack</u>
acc <- 3	3	<init></init>
push acc	3	3, <init></init>
<u>Code</u>		
-------------	--	--
acc <- 3		
push acc		
acc <- 7		

<u>Stack</u> <init> 3, <init> 3, <init>

<u>Code</u>
acc <- 3
push acc
acc <- 7
push acc

<u>Stack</u> <init> 3, <init> 3, <init> 7, 3, <init>

<u>Code</u>	<u>Acc</u>	<u>Stack</u>
acc <- 3	3	<init></init>
push acc	3	3, <init></init>
acc <- 7	7	3, <init></init>
push acc	7	7, 3, <init></init>
acc <- 5	5	7, 3, <init></init>

<u>Code</u>	<u>Acc</u>	<u>Stack</u>
acc <- 3	3	<init></init>
push acc	3	3, <init></init>
acc <- 7	7	3, <init></init>
push acc	7	7, 3, <init></init>
acc <- 5	5	7, 3, <init></init>
acc <- acc + top_of_stack	12	7, 3, <init></init>

<u>Code</u>	<u>Acc</u>	<u>Stack</u>
acc <- 3	3	<init></init>
push acc	3	3, <init></init>
acc <- 7	7	3, <init></init>
push acc	7	7, 3, <init></init>
acc <- 5	5	7, 3, <init></init>
acc <- acc + top_of_stack	12	7, 3, <init></init>
рор	12	3, <init></init>

<u>Code</u>	<u>Acc</u>	<u>Stack</u>
acc <- 3	3	<init></init>
push acc	3	3, <init></init>
acc <- 7	7	3, <init></init>
push acc	7	7, 3, <init></init>
acc <- 5	5	7, 3, <init></init>
acc <- acc + top_of_stack	12	7, 3, <init></init>
рор	12	3, <init></init>
acc <- acc + top of stack	15	3, <init></init>

<u>Code</u>	<u>Acc</u>	<u>Stack</u>
acc <- 3	3	<init></init>
push acc	3	3, <init></init>
acc <- 7	7	3, <init></init>
push acc	7	7, 3, <init></init>
acc <- 5	5	7, 3, <init></init>
acc <- acc + top_of_stack	12	7, 3, <init></init>
рор	12	3, <init></init>
acc <- acc + top_of_stack	15	3, <init></init>
рор	15	<init></init>

Stack Machines: Accumulator Invariants

- The result of computing an expression is **always** in the accumulator
- For an operation op (e₁, ..., e_n), push the accumulator on the stack after computing each of e₁, ..., e_{n-1}
 - \circ e_n's result is in the accumulator before op
 - After the operation **pop** n-1 values
- After computing an expression the stack is as before
 - It is **CRITICAL** that the stack is preserved across the evaluation of a subexpression

Stack Machines: Accumulator Invariants

- The result of computing an expression is **always** in the accumulator
- For an operation op (e₁, ..., e_n), push the accumulator on the stack after computing each of e₁, ..., e_{n-1}
 - \circ e's result is in the accumulator before op
 - After the operation **pop** n-1 values
- After computing an expression the stack is as before
 - It is **CRITICAL** that the stack is preserved across the evaluation of a subexpression
 - Stack before evaluating 7 + 5 is 3, <init>

Stack Machines: Accumulator Invariants

- The result of computing an expression is **always** in the accumulator
- For an operation op (e₁, ..., e_n), push the accumulator on the stack after computing each of e₁, ..., e_{n-1}
 - \circ e's result is in the accumulator before op
 - After the operation **pop** n-1 values
- After computing an expression the stack is as before
 - It is **CRITICAL** that the stack is preserved across the evaluation of a subexpression
 - Stack before evaluating 7 + 5 is 3, <init>
 - Stack after evaluating 7 + 5 is also 3, <init>

• The "compiler" we will outline in this two-lecture series will generate code for a stack machine with an accumulator

- The "compiler" we will outline in this two-lecture series will generate code for a stack machine with an accumulator
 - You might want to run the resulting code on a processor

- The "compiler" we will outline in this two-lecture series will generate code for a stack machine with an accumulator
 - You might want to run the resulting code on a processor
- To do so, we'll implement stack machine instructions using **Cool-ASM** instructions and registers

- The "compiler" we will outline in this two-lecture series will generate code for a stack machine with an accumulator
 - You might want to run the resulting code on a processor
- To do so, we'll implement stack machine instructions using Cool-ASM instructions and registers
 - Cool-ASM is a bytecode format described in the CRM
 - I will assume access to a Cool-ASM interpreter

- The "compiler" we will outline in this two-lecture series will generate code for a stack machine with an accumulator
 - You might want to run the resulting code on a processor
- To do so, we'll implement stack machine instructions using Cool-ASM instructions and registers
 - Cool-ASM is a bytecode format described in the CRM
 - I will assume access to a Cool-ASM interpreter
 - Fun PA3 activity for you: map this Cool-ASM to x86-64
 - effectively, this would be using Cool-ASM as an IR!

- The "compiler" we will outline in this two-lecture series will generate code for a stack machine with an accumulator
 - You might want to run the resulting code on a processor
- To do so, we'll implement stack machine instructions using Cool-ASM instructions and registers
 - Cool-ASM is a bytecode format described in the CRM
 - I will assume access to a Cool-ASM interpreter
 - Fun PA3 activity for you: map this Cool-ASM to x86-64
 - effectively, this would be using Cool-ASM as an IR!
 - you are welcome to do so in your compiler if you want
 - would it come before or after three-address code?

• Cool-ASM is a RISC-style assembly language

- Cool-ASM is a RISC-style assembly language
 - An *assembly language* is an untyped, unsafe, low-level, fast programming language with few-to-no primitives.

- Cool-ASM is a RISC-style assembly language
 - An *assembly language* is an untyped, unsafe, low-level, fast programming language with few-to-no primitives.
- A *register* is a fast-access untyped global variable shared by the entire assembly program.

- Cool-ASM is a RISC-style assembly language
 - An *assembly language* is an untyped, unsafe, low-level, fast programming language with few-to-no primitives.
- A *register* is a fast-access untyped global variable shared by the entire assembly program.
 - Cool-ASM: 8 general registers and 3 special ones (stack pointer, frame pointer, return address)

- Cool-ASM is a RISC-style assembly language
 - An *assembly language* is an untyped, unsafe, low-level, fast programming language with few-to-no primitives.
- A *register* is a fast-access untyped global variable shared by the entire assembly program.
 - Cool-ASM: 8 general registers and 3 special ones (stack pointer, frame pointer, return address)
- An *instruction* is a primitive statement in assembly language that operates on registers
 - Cool-ASM: add, jmp, ld, push, ...

- Cool-ASM is a RISC-style assembly language
 - An *assembly language* is an untyped, unsafe, low-level, fast programming language with few-to-no primitives.
- A *register* is a fast-access untyped global variable shared by the entire assembly program.
 - Cool-ASM: 8 general registers and 3 special ones (stack pointer, frame pointer, return address)
- An *instruction* is a primitive statement in assembly language that operates on registers
 - Cool-ASM: add, jmp, ld, push, ...
- A *load-store architecture*: bring values in to registers from memory to operate on them.

Cool-ASM in Brief: Sample Instructions

add r2 <- r5 r2 lir5 <- 183 dr2 < r1[5]st r1[6] <- r7 my_label: pushr1 subr1 < -r11bnz r1 my label

```
; r2 = r5 + r2
:r5 = 183
: r2 = *(r1+5)
;*(r1+6) = r7
-- dashdash also a comment
;*sp = r1; sp --;
:r1--:
; if (r1 != 0) goto my_label
```

Cool-ASM in Brief: Sample Instructions

add $r_{2} < -r_{5} r_{2}$ lir5 <- 183 dr2 < r1[5]st r1[6] <- r7 my_label: pushr1 subr1 < r11bnz r1 my label

Details of Cool-ASM don't matter much; you're not required to implement this. See the CRM if you want full details.

-- dashdash also a comment
: *sp = r1: sp --:

; if (r1 != 0) goto my_label

• The accumulator is kept in register r1

- The accumulator is kept in register r1
 - This is just a **convention**. You could pick **r2**.

- The accumulator is kept in register r1
 - This is just a convention. You could pick r2.
- The stack is kept in memory

- The accumulator is kept in register r1

 This is just a convention. You could pick r2.
- The stack is kept in memory
- The stack grows towards lower addresses
 - Standard architecture convention (MIPS)

- The accumulator is kept in register r1
 - This is just a convention. You could pick r2.
- The stack is kept in memory
- The stack grows towards lower addresses
 - Standard architecture convention (MIPS)
- The address of the next unused location on the stack is kept in register sp

- The accumulator is kept in register r1
 - This is just a convention. You could pick r2.
- The stack is kept in memory
- The stack grows towards lower addresses
 - Standard architecture convention (MIPS)
- The address of the next unused location on the stack is kept in register sp
 - The top of the stack is always at address sp + 1

- The accumulator is kept in register r1
 - This is just a convention. You could pick r2.
- The stack is kept in memory
- The stack grows towards lower addresses
 - Standard architecture convention (MIPS)
- The address of the next unused location on the stack is kept in register sp
 - The top of the stack is always at address sp + 1
- Cool-ASM "Word Size" = 1 = number of memory cells taken up by one integer/pointer/string

- The accumulator is kept in register r1
 - This is just a convention. You could pick r2.
- The stack is kept in memory
- The stack grows towards lower addresses
 - Standard architecture convention (MIPS)
- The address of the next unused location on the stack is kept in register sp
 - The top of the stack is always at address sp + 1
- Cool-ASM "Word Size" = 1 = number of memory cells taken up by one integer/pointer/string
 - Note intentional simplification vs x86! This is what makes a good IR

Trivia Break: Computer Science

This US Navy Rear Admiral is credited with inventing one of the first compilers for a computer programming language, as well as popularizing machine-independent languages (particularly COBOL). The term "debugging" was also popularized when a logbook that is

sometimes erroneously credited to this person (who was one of the associated project's leaders) physically recorded a moth that had gotten into the relays:



Cool Assembly Example

Stack-machine code for 7 + 5:

Equivalent Cool-ASM:

Cool Assembly Example

Stack-machine code for 7 + 5:

acc <- 7

Equivalent Cool-ASM: li r17

Cool Assembly Example

Stack-machine code for 7 + 5: acc <- 7 push acc

Equivalent Cool-ASM: li r1 7 sw sp[0] <- r1 sub sp <- sp 1
Stack-machine code for 7 + 5: acc <- 7 push acc

acc <- 5

Equivalent Cool-ASM: li r1 7 sw sp[0] <- r1 sub sp <- sp 1 li r1 5

Stack-machine code for 7 + 5: acc <- 7 push acc acc <- 5

acc <- acc + top_of_stack</pre>

Equivalent Cool-ASM: li r1 7 sw sp[0] <- r1 sub sp <- sp 1 li r1 5 lw r2 <- sp[1] add r1 <- r1 r2

Stack-machine code for 7 + 5: acc <- 7 push acc acc <- 5 acc <- acc + top_of_stack

pop

Equivalent Cool-ASM: li r 1 7 sw sp[0] <- r1 sub sp <- sp 1</pre> li r1 5 lw r2 <- sp[1] add r1 < -r1 r2add sp <- sp 1

Stack-machine code for 7 + 5: acc <- 7 push acc acc <- 5 acc <- acc + top_of_stack Equivalent Cool-ASM: li r 1 7 sw sp[0] <- r1 sub sp <- sp 1 li r1 5 lw r2 <- sp[1] add r1 < -r1 r2add sp <- sp 1

рор

We now generalize this to a simple language...

• We have these Cool-ASM instructions:

• We have these Cool-ASM instructions:

push rX st sp[0] <- rX</th> sub sp <- sp 1</td>

• We have these Cool-ASM instructions:

push rX	st sp[0] <- rX sub sp <- sp 1
pop rX	ld rX <- sp[1] add sp <- sp 1

• We have these Cool-ASM instructions:

push rX	st sp[0] <- rX sub sp <- sp 1
pop rX	ld rX <- sp[1] add sp <- sp 1
rX <- top	Id rX <- sp[1]

A Small Language

• We will consider a source language with integers and integer operations (only, for now)

A Small Language

- We will consider a source language with integers and integer operations (only, for now)
- Full grammar:

 $\begin{array}{l} \mathsf{P} -> \mathsf{D} \ ; \mathsf{P} \ | \ \mathsf{D} \\ \mathsf{D} -> \ def \ id(\mathsf{ARGS}) = \mathsf{E} \\ \mathsf{ARGS} -> \ id, \ \mathsf{ARGS} \ | \ id \\ \mathsf{E} -> \ int \ | \ id \ | \ if \ \mathsf{E}_1 = \mathsf{E}_2 \ then \ \mathsf{E}_3 \ else \ \mathsf{E}_4 \\ & \quad | \ \mathsf{E}_1 + \mathsf{E}_2 \ | \ \mathsf{E}_1 - \mathsf{E}_2 \ | \ id(\mathsf{E}_1, ..., \mathsf{E}_n) \end{array}$

Reminder of how to read this:

- capital letters are non-terminals
- italics = lexemes
- "|" separates options

A Small Language (continued)

• The first function definition **f** is the "main" routine

A Small Language (continued)

- The first function definition **f** is the "main" routine
- Running the program on input i means computing f(i)

A Small Language (continued)

- The first function definition **f** is the "main" routine
- Running the program on input i means computing f(i)
- Program for computing the Fibonacci numbers:

```
def fib(x) =
    if x = 1 then
    else
        if x = 2 then
         else
             fib(x - 1) + fib(x - 2)
```

• For each expression **e** we generate Cool- ASM code that:

For each expression e we generate Cool- ASM code that:
 Computes the value of e in r1 (our accumulator)

- For each expression e we generate Cool- ASM code that:
 - Computes the value of e in r1 (our accumulator)
 - Preserves sp and the contents of the stack

- For each expression e we generate Cool- ASM code that:
 - Computes the value of e in r1 (our accumulator)
 - Preserves sp and the contents of the stack
- We define a *code generation function* cgen(e) whose result is the code generated for e

- For each expression **e** we generate Cool- ASM code that:
 - Computes the value of e in r1 (our accumulator)
 - Preserves sp and the contents of the stack
- We define a *code generation function* cgen(e) whose result is the code generated for e
- Like type rules, our code generation function is *syntax-directed*

- For each expression **e** we generate Cool- ASM code that:
 - Computes the value of e in r1 (our accumulator)
 - Preserves sp and the contents of the stack
- We define a *code generation function* cgen(e) whose result is the code generated for e
- Like type rules, our code generation function is *syntax-directed*
 - in other words, it is a big case statement, based on the structure of the input

- For each expression **e** we generate Cool- ASM code that:
 - Computes the value of e in r1 (our accumulator)
 - Preserves sp and the contents of the stack
- We define a *code generation function* cgen(e) whose result is the code generated for e
- Like type rules, our code generation function is *syntax-directed*
 - in other words, it is a big case statement, based on the structure of the input
 - we will have one code generation rule for constants, one for addition, one for subtraction, etc.
 - one for each kind of expression!

Code Generation: Constants

Code Generation: Constants

• The code to evaluate a constant simply copies it into the accumulator:

cgen(123) = li r1 123

Code Generation: Constants

- The code to evaluate a constant simply copies it into the accumulator:
 cgen(123) = li r1 123
- Note that this also preserves the stack, as required

 $cgen(e_1 + e_2) =$

$$cgen(e_1 + e_2) = cgen(e_1)$$

$$cgen(e_{1} + e_{2}) = cgen(e_{1})$$

$$push r1$$

$$cgen(e_{2}) ;; e_{2} now in r1$$

```
cgen(e<sub>1</sub> + e<sub>2</sub>) =
    cgen(e<sub>1</sub>)
    push r1
    cgen(e<sub>2</sub>) ;; e<sub>2</sub> now in r1
    pop t1
```

```
cgen(e<sub>1</sub> + e<sub>2</sub>) =
    cgen(e<sub>1</sub>)
    push r1
    cgen(e<sub>2</sub>) ;; e<sub>2</sub> now in r1
    pop t1
    add r1 <- t1 r1</pre>
```

```
cgen(e<sub>1</sub> + e<sub>2</sub>) =
    cgen(e<sub>1</sub>)
    push r1
    cgen(e<sub>2</sub>) ;; e<sub>2</sub> now in r1
    pop t1
    add r1 <- t1 r1</pre>
```

• Possible optimization: put the result of e₁ directly in register t1?

Code Generation Mistake: Addition Alternative

• Unsafe optimization: put the result of e₁ directly in register t1?

Code Generation Mistake: Addition Alternative

• Unsafe optimization: put the result of e₁ directly in register t1?

```
cgen(e<sub>1</sub> + e<sub>2</sub>) =

cgen(e<sub>1</sub>)

mov t1 <- r1

cgen(e<sub>2</sub>) ;; e<sub>2</sub> now in r1

add r1 <- t1 r1
```

Code Generation Mistake: Addition Alternative

• Unsafe optimization: put the result of e₁ directly in register t1?

```
cgen(e<sub>1</sub> + e<sub>2</sub>) =

cgen(e<sub>1</sub>)

mov t1 <- r1

cgen(e<sub>2</sub>) ;; e<sub>2</sub> now in r1

add r1 <- t1 r1
```

• To see why this is incorrect, try to generate code for : 3 + (7 + 5)

Code Generation Notes
The code that we generated for addition is a template with "holes" for code for evaluating e₁ and e₂

- The code that we generated for addition is a template with "holes" for code for evaluating e₁ and e₂
- Stack-machine code generation is recursive

- The code that we generated for addition is a template with "holes" for code for evaluating e₁ and e₂
- Stack-machine code generation is recursive
 - Code for $e_1 + e_2$ consists of code for e_1 and e_2 glued together
 - with a bit of code to actually do the addition

- The code that we generated for addition is a template with "holes" for code for evaluating e₁ and e₂
- Stack-machine code generation is recursive
 - Code for e₁ + e₂ consists of code for e₁ and e₂ glued together
 with a bit of code to actually do the addition
- Implication: code generation can be written as a recursive-descent tree walk of the AST
 - \circ At least for expressions

Code Generation: Subtraction

 $cgen(e_1 - e_2) =$

Code Generation: Subtraction

```
cgen(e<sub>1</sub> - e<sub>2</sub>) =
    cgen(e<sub>1</sub>)
    push r1
    cgen(e<sub>2</sub>) ;; e<sub>2</sub> now in r1
    pop t1
    sub r1 <- t1 r1</pre>
```

Code Generation: Subtraction

- Almost identical to addition!
 - only difference is the sub instruction, which does what you'd expect (subtraction)

• To generate code for if, we need control flow instructions

- To generate code for if, we need control flow instructions
- Cool-ASM has two that will be useful:

- To generate code for if, we need control flow instructions
- Cool-ASM has two that will be useful:

```
beq r1 r2 label
conditional branch to label if r1 = r2
```

- To generate code for if, we need control flow instructions
- Cool-ASM has two that will be useful:

beq r1 r2 label conditional branch to label if r1 = r2

jmp label unconditional jump to label

- To generate code for if, we need control flow instructions
- Cool-ASM has two that will be useful:

beq r1 r2 label conditional branch to label if r1 = r2

jmp label unconditional jump to label

Note the similarity to what x86 provides! How easy would it be to convert these Cool-ASM instructions into x86?

cgen(if e1 = e2 then e3 else e4) =

```
cgen(if e1 = e2 then e3 else e4) =
cgen(e1)
```

```
cgen(if e1 = e2 then e3 else e4) =
cgen(e1)
push r1
```

cgen(if e1 = e2 then e3 else e4) =
 cgen(e1)
 push r1
 cgen(e2)

cgen(if e1 = e2 then e3 else e4) =
 cgen(e1)
 push r1
 cgen(e2)
 pop t1

cgen(if e1 = e2 then e3 else e4) =
 cgen(e1)
 push r1
 cgen(e2)
 pop t1
 beq r1 t1 true_branch ;; else fall through

```
cgen(if e1 = e2 then e3 else e4) =
   cgen(e1)
   push r1
   cgen(e2)
   pop t1
   beq r1 t1 true_branch ;; else fall through
   cgen(e4)
```

```
cgen(if e1 = e2 then e3 else e4) =
   cgen(e1)
   push r1
   cgen(e2)
   pop t1
   beq r1 t1 true_branch ;; else fall through
   cgen(e4)
   jmp end_if
```

```
cgen(if e1 = e2 then e3 else e4) =
    cgen(e1)
    push r1
    cgen(e2)
    popt1
    beq r1 t1 true_branch ;; else fall through
    cgen(e4)
    jmp end_if
 true_branch:
```

```
cgen(if e1 = e2 then e3 else e4) =
    cgen(e1)
    push r1
    cgen(e2)
    popt1
    beq r1 t1 true_branch ;; else fall through
    cgen(e4)
    jmp end_if
 true_branch:
    cgen(e3)
```

```
cgen(if e1 = e2 then e3 else e4) =
    cgen(e1)
    pushr1
    cgen(e2)
    popt1
    beq r1 t1 true_branch ;; else fall through
    cgen(e4)
    jmp end_if
 true_branch:
    cgen(e3)
 end if:
```

Trivia Break: Art History

This style of visual arts, architecture, and product design first appeared in Paris in the early 1910s and flourished in the US and Europe during the 1920s. It is characterized by rare and expensive materials, such as ebony and ivory, exquisite craftsmanship, and the inclusion of "modern" materials like chrome plating, stainless steel, and plastic. In New York City, the Empire State Building, Chrysler Building, and other buildings from the 1920s and 1930s are monuments to the style.



Trivia Break: Geography

This lake in Russia is the world's seventh-largest by surface area (it is just larger than Belgium). However, unlike most other freshwater lakes in the world, it is a *rift lake*: a geologically active rift in the Earth's crust at the bottom of the lake is being pulled apart at a rate of about 4mm per year. This has two important implications for this lake: it is both the oldest (~25-30 million years old) and deepest (avg. 744m, deepest 1,642m) lake in the world. Because of its depth, it contains about 22% of the world's freshwater.

• Recall that an *activation record* (or *stack frame*) stores calling context information on the stack during a function call.

- Recall that an *activation record* (or *stack frame*) stores calling context information on the stack during a function call.
 - Code for function calls/definitions depends on the layout of the activation record

- Recall that an *activation record* (or *stack frame*) stores calling context information on the stack during a function call.
 - Code for function calls/definitions depends on the layout of the activation record
- A very simple AR suffices for this language:

- Recall that an *activation record* (or *stack frame*) stores calling context information on the stack during a function call.
 - Code for function calls/definitions depends on the layout of the activation record
- A very simple AR suffices for this language:
 - The result is always in the accumulator
 - Thus, no need to store the result in the AR

- Recall that an *activation record* (or *stack frame*) stores calling context information on the stack during a function call.
 - Code for function calls/definitions depends on the layout of the activation record
- A very simple AR suffices for this language:
 - The result is always in the accumulator
 - Thus, no need to store the result in the AR
 - The activation record holds actual parameters
 - For $f(x_1,...,x_n)$, push $x_1,...,x_n$ on the stack
 - These are the only variables in this language

• The *calling convention* (or *stack discipline*) guarantees that on function exit sp is the same as it was on entry

- The *calling convention* (or *stack discipline*) guarantees that on function exit sp is the same as it was on entry
 - So, no need to save sp explicitly

- The *calling convention* (or *stack discipline*) guarantees that on function exit sp is the same as it was on entry
 - So, no need to save sp explicitly
 - Our simple AR maintains stack discipline (why?)

- The *calling convention* (or *stack discipline*) guarantees that on function exit sp is the same as it was on entry
 - So, no need to save sp explicitly
 - Our simple AR maintains stack discipline (why?)
- We do need the return address

- The *calling convention* (or *stack discipline*) guarantees that on function exit sp is the same as it was on entry
 - So, no need to save sp explicitly
 - Our simple AR maintains stack discipline (why?)
- We do need the return address
- Also, it's handy to have a pointer to the start of the current activation

- The *calling convention* (or *stack discipline*) guarantees that on function exit sp is the same as it was on entry
 - So, no need to save sp explicitly
 - Our simple AR maintains stack discipline (why?)
- We do need the return address
- Also, it's handy to have a pointer to the start of the current activation
 - This pointer lives in register fp ("frame pointer")
Calling Convention

- The *calling convention* (or *stack discipline*) guarantees that on function exit sp is the same as it was on entry
 - So, no need to save sp explicitly
 - Our simple AR maintains stack discipline (why?)
- We do need the return address
- Also, it's handy to have a pointer to the start of the current activation
 - This pointer lives in register **fp** ("frame pointer")
 - Reason for frame pointer will be clear shortly

• Summary: for this language, an AR with

- Summary: for this language, an AR with
 - 1. the caller's frame pointer

- Summary: for this language, an AR with
 - 1. the caller's frame pointer
 - 2. the actual parameters, and

- Summary: for this language, an AR with
 - 1. the caller's frame pointer
 - 2. the actual parameters, and
 - 3. the return address

- Summary: for this language, an AR with
 - 1. the caller's frame pointer
 - 2. the actual parameters, and
 - 3. the return address

suffices.

• Visual learners: consider a call to f(x,y). The AR will be:

- Summary: for this language, an AR with
 - 1. the caller's frame pointer
 - 2. the actual parameters, and
 - 3. the return address

suffices.

• Visual learners: consider a call to f(x,y). The AR will be:



• The *calling sequence* is the instructions (of both caller and callee) to set up a function invocation

- The *calling sequence* is the instructions (of both caller and callee) to set up a function invocation
- Requires one new instruction: call label

- The *calling sequence* is the instructions (of both caller and callee) to set up a function invocation
- Requires one new instruction: call label
 - Jump to label, save address of next instruction in the special-purpose register ra
 - On other architectures (like x86!) the return address is stored on the stack by the "call" instruction
 - (This is also called "*branch and link*".)

cgen(f(e₁,...,e_n)) =

 $cgen(f(e_1,...,e_n)) =$

push fp

• The caller saves its value of the frame pointer

cgen(f(e₁,...,e_n)) = push fp cgen(e₁) push r1

- The caller saves its value of the frame pointer
- Then it saves the actual arguments in order

cgen(f(e₁,...,e_n)) = push fp cgen(e₁) push r1

...

- The caller saves its value of the frame pointer
- Then it saves the actual arguments in order

cgen(f(e₁,...,e_n)) = push fp cgen(e₁) push r1

```
...
cgen(e<sub>n</sub>)
push r1
```

- The caller saves its value of the frame pointer
- Then it saves the actual arguments in order

cgen(f(e₁,...,e_n)) = push fp cgen(e₁) push r1

•••

cgen(e_n) push r1 call f_entry

- The caller saves its value of the frame pointer
- Then it saves the actual arguments in order
- The caller saves the return address in register ra (via the call instruction)
- The AR so far is n+1 bytes long

cgen(f(e₁,...,e_n)) = push fp cgen(e₁) push r1

•••

cgen(e_n) push r1 call f_entry pop fp

- The caller saves its value of the frame pointer
- Then it saves the actual arguments in order
- The caller saves the return address in register ra (via the call instruction)
- The AR so far is n+1 bytes long
- Caller restores fp

 $cgen(def f(x_1,...,x_n) = e) =$

cgen(def f(x₁,...,x_n) = e) = f_entry:

```
cgen(def f(x<sub>1</sub>,...,x<sub>n</sub>) = e) =
f_entry:
mov fp<-sp
```

cgen(def f(x₁,...,x_n) = e) = f_entry: mov fp<-sp push ra

```
cgen(def f(x<sub>1</sub>,...,x<sub>n</sub>) = e) =
f_entry:
mov fp<-sp
push ra
cgen(e)
```

```
cgen(def f(x<sub>1</sub>,...,x<sub>n</sub>) = e) =

f_entry:

mov fp<-sp

push ra

cgen(e)

ra <- top
```

```
cgen(def f(x<sub>1</sub>,...,x<sub>n</sub>) = e) =
f_entry:
mov fp<-sp
push ra
cgen(e)
ra <- top
add sp <- sp z
```

 $\operatorname{cgen}(\operatorname{def} f(x_1,...,x_n) = e) =$ f_entry: mov fp<-sp push ra cgen(e) ra <- top add sp <- sp z return

New instruction: return

 Jump to address in register ra

 $cgen(def f(x_1,...,x_n) = e) =$ f entry: mov fp<-sp push ra cgen(e) ra <- top add sp <- sp z return

- New instruction: return

 Jump to address in register ra
- Note: the frame pointer points to the top, not bottom of the frame

 $cgen(def f(x_1,...,x_n) = e) =$ f entry: mov fp<-sp push ra cgen(e) ra <- top add sp <- sp z return

- New instruction: return

 Jump to address in register ra
- Note: the frame pointer points to the top, not bottom of the frame
- The callee pops the return address, the actual arguments and the saved value of the frame pointer

 $cgen(def f(x_1,...,x_n) = e) =$ f entry: mov fp<-sp push ra cgen(e) ra <- top add sp <- sp z return

- New instruction: return

 Jump to address in register ra
- Note: the frame pointer points to the top, not bottom of the frame
- The callee pops the return address, the actual arguments and the saved value of the frame pointer
- z = n + 2 (so far)



Before call

SP FP







Code Generation: Variables

• Variable references are the last construct

Code Generation: Variables

- Variable references are the last construct
- The "variables" of a function are just its **parameters**
- Variable references are the last construct
- The "variables" of a function are just its parameters
 They are all in the AR

- Variable references are the last construct
- The "variables" of a function are just its **parameters**
 - \circ $\,$ They are all in the AR $\,$
 - Pushed by the caller

- Variable references are the last construct
- The "variables" of a function are just its **parameters**
 - \circ $\,$ They are all in the AR $\,$
 - Pushed by the caller
- Problem: because the stack grows when intermediate results are saved, the variables are **not at a fixed offset** from sp

- Variable references are the last construct
- The "variables" of a function are just its **parameters**
 - \circ $\,$ They are all in the AR $\,$
 - Pushed by the caller
- Problem: because the stack grows when intermediate results are saved, the variables are **not at a fixed offset** from sp
 - Challenge question: what are they at a fixed offset from?

- Variable references are the last construct
- The "variables" of a function are just its **parameters**
 - \circ They are all in the AR
 - Pushed by the caller
- Problem: because the stack grows when intermediate results are saved, the variables are **not at a fixed offset** from sp
 - Challenge question: what are they at a fixed offset from?
 - Answer: the frame pointer

- Variable references are the last construct
- The "variables" of a function are just its **parameters**
 - \circ They are all in the AR
 - Pushed by the caller
- Problem: because the stack grows when intermediate results are saved, the variables are **not at a fixed offset** from sp
 - Challenge question: what are they at a fixed offset from?
 - Answer: the frame pointer
 - Always points to the return address on the stack
 - = the value of sp on function entry

- Variable references are the last construct
- The "variables" of a function are just its **parameters**
 - \circ $\,$ They are all in the AR $\,$
 - Pushed by the caller
- Problem: because the stack grows when intermediate results are saved, the variables are **not at a fixed offset** from sp
 - Challenge question: what are they at a fixed offset from?
 - Answer: the frame pointer
 - Always points to the return address on the stack
 - = the value of sp on function entry
 - It doesn't move => args on the stack are at a fixed offset

 Example: For a function def f(x₁, x₂) = e the activation and frame pointer are set up as follows:

 Example: For a function def f(x₁, x₂) = e the activation and frame pointer are set up as follows:



 Example: For a function def f(x₁, x₂) = e the activation and frame pointer are set up as follows:



• x₁ (first parameter) is at **fp + 2**

 Example: For a function def f(x₁, x₂) = e the activation and frame pointer are set up as follows:



- x₁ (first parameter) is at **fp + 2**
- x_2^{-} (second parameter) is at fp + 1

 Example: For a function def f(x₁, x₂) = e the activation and frame pointer are set up as follows:



- x₁ (first parameter) is at **fp + 2**
- x₂ (second parameter) is at fp + 1
- Thus:

cgen(x_i) = Id r1 <- fp[z]

 Example: For a function def f(x₁, x₂) = e the activation and frame pointer are set up as follows:



• The activation record must be **designed together** with the code generator

- The activation record must be **designed together** with the code generator
- Code generation can be done by **recursive traversal** of the AST

- The activation record must be **designed together** with the code generator
- Code generation can be done by **recursive traversal** of the AST
- As you write your compiler, we recommend starting with a stack machine (simpler!)

- The activation record must be **designed together** with the code generator
- Code generation can be done by **recursive traversal** of the AST
- As you write your compiler, we recommend starting with a stack machine (simpler!)
 - ./cool –asm generates Cool-ASM stack machine code for Cool
 - use this to help you with PA3!

- The activation record must be **designed together** with the code generator
- Code generation can be done by **recursive traversal** of the AST
- As you write your compiler, we recommend starting with a stack machine (simpler!)
 - ./cool -asm generates Cool-ASM stack machine code for Cool
 - use this to help you with PA3!
- Production compilers do different things:

- The activation record must be **designed together** with the code generator
- Code generation can be done by **recursive traversal** of the AST
- As you write your compiler, we recommend starting with a **stack machine** (simpler!)
 - ./cool -asm generates Cool-ASM stack machine code for Cool
 - use this to help you with PA3!
- Production compilers do different things:
 - keep as many values as possible in registers, etc

- The activation record must be **designed together** with the code generator
- Code generation can be done by **recursive traversal** of the AST
- As you write your compiler, we recommend starting with a **stack machine** (simpler!)
 - ./cool -asm generates Cool-ASM stack machine code for Cool
 - use this to help you with PA3!
- Production compilers do different things:
 - keep as many values as possible in registers, etc
 - save this stuff for PA4

After Break

- We will continue our discussion of this stack machine language by discussing the **"hard" features of Cool**
 - object layout for OOP
 - dynamic dispatch

After Break

- We will continue our discussion of this stack machine language by discussing the "hard" features of Cool
 - object layout for OOP
 - dynamic dispatch
- You should start PA3c3 over break
 - I assure you that you will regret it if you do not

After Break

- We will continue our discussion of this stack machine language by discussing the "hard" features of Cool
 - object layout for OOP
 - dynamic dispatch
- You should start PA3c3 over break
 - I assure you that you will regret it if you do not
- Don't forget about the midterm on 3/26
 - I will hold a review session on 3/24 or 3/25 in the evening, keep an eye on Discord for a poll