

Intermediate Representation I: Control Structures

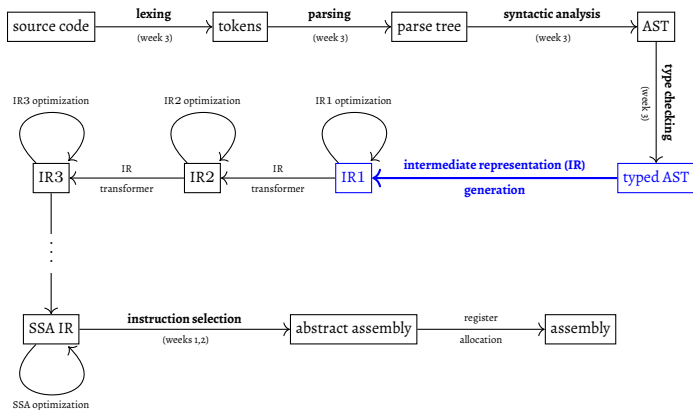
CSE 302 – Compilers – Week 4

Kaustuv Chaudhuri

Inria & École polytechnique

2019-09-30

Where in the compiler are we?



Summary of our Compiler

Things you have seen

- Define the grammar of the language
- Parse source code into an abstract syntax tree (AST)
- Convert an AST for **sequential code** (BX0) directly into machine instructions (instruction selection)

Summary of our Compiler

Things you have seen

- Define the grammar of the language
- Parse source code into an abstract syntax tree (AST)
- Convert an AST for **sequential code** (BX0) directly into machine instructions (instruction selection) by **spilling** all temporaries onto the stack

Summary of our Compiler

Things you have seen

- Define the grammar of the language
- Parse source code into an abstract syntax tree (AST)
- Convert an AST for **sequential code** (BX0) directly into machine instructions (instruction selection) by **spilling** all temporaries onto the stack

What more is necessary to compile BX1?

- Conditionals and loops
- Booleans and boolean expressions
- Comparisons

Summary of our Compiler

Things you have seen

- Define the grammar of the language
- Parse source code into an abstract syntax tree (AST)
- Convert an AST for **sequential code** (BX0) directly into machine instructions (instruction selection) by **spilling** all temporaries onto the stack

What more is necessary to compile BX1?

- Conditionals and loops
- Booleans and boolean expressions
- Comparisons

Common feature: **jumps**

Today's Agenda

- 1 Jumps and comparisons in AMD64
- 2 Register Transfer Language (RTL)
- 3 Compiling control structures to RTL
- 4 Code-generation for RTL

Jumps and Comparisons in AMD64

Labels and Instruction Pointer

- **Pointers:** (virtual) memory addresses as values
Can be dereferenced to read/write memory
- **Label:** addresses of an instruction in compiled code
- **Instruction pointer** (`%rip`): contains the address of the *next* instruction, so setting `%rip` changes the code path
- **Jump:** specialized instructions to set `%rip`.
 - `jmp`: absolute jump
 - `je`, `jne`, `jle`, etc: conditional jumps

Example of Jumps and Conditionals

```
var x = 0, y = 1 : int64;  
var z : int64;  
  
while (x < 200) {  
    print x;  
    z = x;  
    x = y;  
    y = y + z;  
}
```

```
main:  
    pushq %rbp  
    movq %rsp, %rbp  
    subq $24, %rsp  
    movq $0, -8(%rbp)           # x = 0  
    movq $1, -16(%rbp)         # y = 1  
    .L0:                         # while begin  
        cmpq $199, -8(%rbp)    # 199 >= x ?  
        jge .L1  
        movq -8(%rbp), %rdi  
        callq bxl_print_int  
        movq -8(%rbp), %r11  
        movq %r11, -24(%rbp)    # z = x  
        movq -16(%rbp), %r11  
        movq %r11, -8(%rbp)    # x = y  
        movq -24(%rbp), %r11  
        addq %r11, -16(%rbp)    # y = y + z  
        jmp .L0  
    .L1:                         # while end  
        movq %rbp, %rsp  
        popq %rbp  
        xorq %rax, %rax  
        retq
```

Comparison Instruction

```
cmpq    $42, %rax  
      src1    src2  
    (reg/imm) (reg/mem)
```

- Sets the status flags register (%**eflags**) based on `src1 - src2`
- Note: does not affect any of the other registers
- Overlays:
 - `cmpl $42, %eax`
 - `cmpw $42, %ax`
 - `cmpb $42, %al`

Jump Instructions

```
jmp    .L42
jcnd  .L42
      label
```

- Jumps unconditionally (**jmp**) or when a particular condition flag is set/unset based on an earlier **cmpq**
- Some conditional jumps: (note synonyms)

Instruction	Condition	wrt: [cmpq src1, src2]
je, jz	ZF=1	src1 - src2 == 0
jne, jnz	ZF=0	src1 - src2 != 0
jle, jng	ZF=1 or SF≠OF	src1 - src2 <= 0
jl, jnge	ZF=0 and SF≠OF	src1 - src2 < 0
jge, jnl	ZF=1 or SF=OF	src1 - src2 >= 0
jg, jnle	ZF=0 and SF=OF	src1 - src2 > 0

Register Transfer Language

Intermediate Languages

- The AST retains much of the structure of the original program
 - Control structures
 - Variables
 - Aggregate types, abstract types, polymorphism, etc. (to come)
- An **intermediate language** eliminates most of these structures and abstractions
 - Only unstructured control (jumps, implicit or explicit)
 - **Pseudo-registers** of only primitive types + pointers
 - Simple family of instructions, close to assembly
 - Structure access is explicit (to come)
- **Control Flow Graph** (CFG): the purpose of all intermediate languages is to yield a graph representation of the flow of control through a program, which can be used for analysis and optimization in later phases

Register Transfer Language (RTL)

Our first intermediate language

- Register Transfer Language
 - Can use an infinite number of **pseudo-registers** (aka “pseudos”)
 - Will be written #1, #2, etc.
 - Overlay suffixes: q (64-bit), d (32-bit), w (16-bit), b (8-bit)
 - Special pseudo ## for “don’t care”
 - Line comments will now start with ;
 - Pervasive use of **labels**, written L0, L1, etc. (globally unique)
 - Each RTL **instruction** has:
 - A unique **source** label (aka “in label”)
 - A list of **successor** labels (aka “out label”)
 - No global relative ordering of instructions
 - For procedures: (to come)
 - Two unique labels: *enter*, *exit*
 - Certain pseudo-registers marked as *input* / *output*
 - Instructions know which procedure they are for
- Assembly:
 - Can only use a finite number of **hardware registers**
 - Has a fixed ordering of instructions

RTL Instructions by Example

(Inspired by, but not the same as, AMD64 assembly)

L42:	binop	add,	#7q,	#10q	→	L57
in label	instruction	opcode	arg1	arg2		out label

- All arguments must be registers (except for **move**)
- Some instructions have an additional **opcode**, which unites many instructions into the same category
- Must mention both in label and out label(s)
- Multiple out labels separated by ,

RTL Instructions: Summary

RTL instruction	Description
Li: move n, r1 \rightarrow Lo	move imm., r1 = n
Li: copy r1, r2 \rightarrow Lo	copy regs, r2 = r1
Li: unop op, r1 \rightarrow Lo	r1 = op r1
Li: binop op, r1, r2 \rightarrow Lo	r2 = r2 op r1
Li: ubbranch op, r1 \rightarrow Lo1, Lo2	unary branch: op r1
Li: bbranch , r1, r2 \rightarrow Lo1, Lo2	binary branch: r1 op r2
Li: goto \rightarrow Lo	unconditional jump
Li: call f(r1, ..., rn), r \rightarrow Lo	function call, result in r
Li: return r	return (note: no out label)

n	immediate (no \$)	r1, r2, ...	pseudo-registers
op	opcode	Li, Lo, ...	labels

RTL: Example

```
var x = 0, y = 1 : int64;  
var z : int64;  
  
while (x < 200) {  
    print x;  
    z = x;  
    x = y;  
    y = y + z;  
}
```

```
enter: L0  
exit: L9  
----  
L0: move 0, #0q      → L1  
L1: move 1, #1q      → L2  
L2: move 200, #2q    → L3  
L3: bbranch jl, #0q, #2q → L4, L8  
L4: call bxl_print(#0q), ## → L5  
L5: copy #0q, #3q    → L6  
L6: copy #1q, #0q    → L7  
L7: binop add, #3q, #1q → L3  
L8: move 0, #0q      → L9  
L9: return #0q
```

RTL Generation

Generating RTL

(From the typed AST – for now)

- Basic idea: adapt **maximal munch** to deal with labels, control structures, new expression forms, etc.
 - Next few slides present the *top down* variant
- The **BX1** AST is still simple enough that we can directly generate RTL
- Later in the course, the AST will first need to be *simplified* or **elaborated** before the RTL generation

RTL Generation: Integer Expressions

Simple cases

$$L_i = \text{RTL}_i(e, r_d, L_o)$$

- Given: e (an `int64` expression), r_d (a destination pseudo-register), L_o (an out-label)
- Returns: L_i (an in-label)
- Side-effect: emits RTL instructions

e	$L_i = \text{RTL}_i(e, r_d, L_o)$	fresh
42	L_i : <code>move 42, rd \rightarrow Lo</code>	L_i
<code>ri</code>	L_i : <code>copy ri, rd \rightarrow Lo</code>	L_i
$e_1 + e_2$	$L_i = \text{RTL}_i(e_1, r_t, L1)$ $L1 = \text{RTL}_i(e_2, r_d, L2)$ $L2$: <code>binop addq, rt, rd \rightarrow Lo</code>	$L2, r_t$

Read bottom to top

RTL Generation: Boolean Expressions

$$L_i = \text{RTL}_b(e, L_t, L_f)$$

- Given: e (a **bool** expression), L_t (an out-label for the **true** case), L_f (out-label for the **false** case)
- Returns: L_i (an in-label)
- Side-effect: emits RTL instructions

e	$L_i = \text{RTL}_b(e, L_t, L_f)$	fresh
true	L_i is just L_t	—
false	L_i is just L_f	—
$! e_1$	$L_i = \text{RTL}_b(e_1, L_f, L_t)$	—
$e_1 \ \&\& \ e_2$	$L_i = \text{RTL}_b(e_1, L1, L_f)$ $L1 = \text{RTL}_b(e_2, L_t, L_f)$	—
$e_1 \ \ e_2$	$L_i = \text{RTL}_b(e_1, L_t, L1)$ $L1 = \text{RTL}_b(e_2, L_t, L_f)$	—
$e_1 < e_2$	$L_i = \text{RTL}_i(e_1, r1, L1)$ $L1 = \text{RTL}_i(e_2, r2, L2)$ $L2: \text{bbranch } j_l, r1, r2 \rightarrow L_t, L_f$	$L2, r1, r2$

RTL Generation: Boolean Comparison

When $e_1, e_2 : \text{bool}$, you can use the following equivalences:

$$(e_1 == e_2) \equiv (e_1 \&\& e_2) \mid \mid ! (e_1 \mid \mid e_2)$$

$$(e_1 != e_2) \equiv ! (e_1 \&\& e_2) \&\& (e_1 \mid \mid e_2)$$

Note: each subexpression e_1 and e_2 must only be computed **once**!

RTL Generation: Statements

$$L_i = \text{RTL}_s(s, L_o)$$

- Given: s a statement, L_o an out-label
- Returns: L_i , an in-label
- Side-effect: emits RTL instructions

s	$L_i = \text{RTL}_s(s, L_o)$	proviso
$x = e; \quad (\text{int64})$	$L_i = \text{RTL}_i(e, r, L_o)$	$r = \text{lookup}(x)$
$x = e; \quad (\text{bool})$	$L_i = \text{RTL}_b(e, Lt, Lf)$ $Lt: \text{move } 1, r \rightarrow L_o$ $Lf: \text{move } 0, r \rightarrow L_o$	$r = \text{lookup}(x)$ Lt, Lf fresh

Handle print similarly

RTL Generation: Blocks, Conditionals, Loops

s	$Li = RTL_s(s, Lo)$	fresh
$\{ \}$	Li is just Lo	—
$\{s_1 \ s_2 \ \cdots \ s_n\}$	$Li = RTL_s(s_1, L2)$ $L2 = RTL_s(s_2, L3)$ \vdots $Ln = RTL_s(s_n, Lo)$	—
<u>if</u> (e_c) s_t <u>else</u> s_f	$Li = RTL_b(e_c, Lt, Lf)$ $Lt = RTL_s(s_t, Lo)$ $Lf = RTL_s(s_f, Lo)$	—
<u>while</u> (e_c) s_b	$Li = RTL_b(e_c, Lt, Lo)$ $Lt = RTL_s(s_b, Lend)$ $Lend: \text{goto } Li$	Lend

RTL to Assembly

Stupid Algorithm (That Works)

- For all non-jumping instructions, jump to out-labels.



Stupid Algorithm (That Works)

Step 2/3

- For jumping instructions, jump to the false out-label

L10: **ubbranch** **jz**, #3 → L20, L30



```
.L10:
    ### TODO: load #3 into %rax
    cmpq $0, %rax
    je .L20           # true case
    jmp .L30          # fall-through false case
```

Stupid Algorithm (That Works)

Step 3/3

- Line up all the generated instructions in some order (say top-to-bottom)
- Remove redundant `jumps`.

```
### ...  
addq %esi, %eax  
movq %eax, -8(%rbp)  
jmp .L10  
.L10:  
movq -16(%rbp), %eax  
cmpq $200, %eax  
### ...
```

⇒

```
### ...  
addq %esi, %eax  
movq %eax, -8(%rbp)  
# --- jmp removed ---  
.L10:  
movq -16(%rbp), %eax  
cmpq $200, %eax  
### ...
```

Stupid Algorithm (That Works)

Step 3/3

- Line up all the generated instructions in some order (say top-to-bottom)
- Remove redundant `jmps`.

```
### ...  
addq %esi, %eax  
movq %eax, -8(%rbp)  
jmp .L10  
.L10:  
movq -16(%rbp), %eax  
cmpq $200, %eax  
### ...
```

⇒

```
### ...  
addq %esi, %eax  
movq %eax, -8(%rbp)  
# --- jmp removed ---  
.L10:  
movq -16(%rbp), %eax  
cmpq $200, %eax  
### ...
```

- Warning! Don't remove labels (unless you can prove that nothing else jumps to that label)

BX and RTL: The Road Ahead

- Pointless to build a more complex algorithm *at this point* for linearizing RTL to AMD64
- Next step: RTL with **explicit**s (ERTL):
 - Instructions with both pseudo-registers and machine registers
 - Break down complex instructions (e.g., `idivq`) into simpler instruction sequences
 - Coalesce linear sequences into **basic blocks**
- BX2 will add *functions* and *procedures*, so we will add to ERTL:
 - Input and output pseudo-registers
 - Calling conventions and stack frames
- We will then move on to Static Single Assignment (SSA) form

BX1 and Lab 3/Week 2

- At the end of Week 1 (this Thursday) you should:
 - Parse the **BX1** syntax
 - Have written and submitted 5 example **BX1** programs
These programs will be given to everyone as challenges
 - Test your programs with the provided **BX1** interpreter
- Week 2 tasks:
 - Implement the RTL instruction format
 - Transform **BX1** AST to RTL
 - Transform RTL to **AMD64**
 - We will give you an RTL interpreter, which you can use to test your **BX1** to RTL phase