

Lecture Notes on Calling Conventions

15-411: Compiler Design
Frank Pfenning*

Lecture 11
October 6, 2015

1 Introduction

In Lab 3 you will be adding functions to the arithmetic language with loops and conditionals. Compiling functions creates some new issues in the front end and the back end of the compiler. In the front end, we need to make sure functions are called with the right number of arguments, and arguments of the right type. In the back end, we need to create assembly code that respects the *calling conventions* of the machine architecture. Strict adherence to the calling conventions is crucial so that your code can interoperate with library routines, and the environment can call functions that you define.

Calling conventions are rather machine-specific and often quite arcane. You must carefully read the Section 3.2 of the AMD64 ABI [MHJM09]¹. The 15-213 course's lecture slides Machine-Level Programming (Procedures), as of Fall 2015, cover x86-64 exclusively. You can find more detail in Chapter 3 of Bryant and O'Hallaron's textbook [BO16].

2 Intermediate representations

We have already seen our discussion of intermediate representations that function calls should take pure arguments in order to easily guarantee the left-to-right evaluation order prescribed by our language semantics. Moreover, they should be lifted to the level of commands rather than remain embedded inside expressions because functions may have side-effects.

*with contributions by Rob Simmons

¹Available at <http://refspecs.linuxfoundation.org/elf/x86-64-abi-0.99.pdf>

We will discuss two lower-level forms of intermediate representation for function calls. The first, higher-level, form is one we will treat as an addition to language of quads (three-address code), an instruction

$$d \leftarrow f(s_1, \dots, s_n)$$

where each s_i is a source operand and d is a destination operand. In the language of quads, we can extend our generic $\text{def}(l, x)$, $\text{use}(l, x)$ and $\text{succ}(l, l')$ to handle this new form of instruction.

$$\frac{l : d \leftarrow f(s_1, \dots, s_n)}{\begin{array}{l} \text{def}(l, d) \\ \text{use}(l, s_i) \quad (1 \leq i \leq n) \\ \text{succ}(l, l + 1) \end{array}} J_8$$

The second lower-level intermediate representation, which will do register allocation on, will also be treated a three-address code language, but it may fit better with two-address code in your compiler. In this lower-level intermediate language, we will simply have calls of the form

call f

All arguments and return values will be designated by special registers.

3 x86-64 Calling Conventions

In x86-64, the first six arguments are passed in registers, the remaining arguments are passed on the stack. The result is returned in a specific return register `%rax`. These conventions do not count floating point arguments and results, which are passed in the dedicated floating point registers `%xmm0` to `%xmm7` and on the stack only if there are more than eight floating point parameters. Fortunately, our language has only integers at the moment, so you do not have to worry about the conventions for floating point numbers.

In the IA32 instruction set that preceded x86-64, stack frames were required to have a *frame* or *base* pointer `%ebp` which had to be saved and restored in a very specific fashion with each function call. It provided a reliable pointer to the beginning of a stack frame for easy calculation of frame offsets to handle references to arguments and local variables. It also made it easier for tools such as `gdb` to print backtraces of the stack. On the x86-64 this information is maintained elsewhere and a frame pointer is no longer required.

The general organization of stack frames at the time a procedure is called, will be as follows.

| Address | Contents | Frame |
|----------|----------------|--------|
| ... | ... | Caller |
| 16(%rsp) | argument 8 | |
| 8(%rsp) | argument 7 | |
| (%rsp) | return address | |

Note that all arguments take 8 bytes of space on the stack, even if the type of argument would indicate that only 4 bytes need to be passed.

The function that is called, the *callee*, should set up its stack frame, reserving space for local variables, spilled temps that could not be assigned to registers, and arguments passed to functions it calls in turn. We recommend calculating the total space needed and then decrementing the stack pointer `%rsp` by the appropriate amount only one time within the function. By changing the stack pointer only once, at the beginning, references to parameters and local variables remain constant throughout the function's execution. The stack then looks as follows, where the size of the callee's stack frame is n .

| Position | Contents | Frame |
|-----------------|---|----------|
| ... | ... | Caller |
| $n + 16(\%rsp)$ | argument 8 | |
| $n + 8(\%rsp)$ | argument 7 | |
| $n + 0(\%rsp)$ | return address | |
| | local variables | Callee |
| | ... | |
| | argument build area for function calls | |
| | ... | |
| (%rsp) | end of frame | |
| | 128 bytes | red zone |

Note that `%rsp` should be aligned $0 \bmod 16$ before another function is called, and may be assumed to be aligned $8 \bmod 16$ on function entry. This happens because the `call` instruction saves the 64-bit return address on the stack.

The area below the stack pointer is called the *red zone* and may be used by the callee as temporary storage for data that is not needed across function calls or even to build arguments to be used before a function call. The ABI states that the red zone "shall not be modified by signal or interrupt handlers." This can be tricky, however, because, for example, Linux kernel code may not respect the red zone and overwrite this area. We therefore suggest not using the red zone.

4 Register Convention

We extract from [MHJM09] the relevant information on register usage. In the first column is a suggested numbering for the purpose of register allocation.

| Abstract form | x86-64 Register | Usage | Preserved across function calls |
|--------------------------|-----------------|---------------|---------------------------------|
| <i>res</i> ₀ | %rax | return value* | No |
| <i>arg</i> ₁ | %rdi | argument 1 | No |
| <i>arg</i> ₂ | %rsi | argument 2 | No |
| <i>arg</i> ₃ | %rdx | argument 3 | No |
| <i>arg</i> ₄ | %rcx | argument 4 | No |
| <i>arg</i> ₅ | %r8 | argument 5 | No |
| <i>arg</i> ₆ | %r9 | argument 6 | No |
| <i>ler</i> ₇ | %r10 | caller-saved | No |
| <i>ler</i> ₈ | %r11 | caller-saved | No |
| <i>lee</i> ₉ | %rbx | callee-saved | Yes |
| <i>lee</i> ₁₀ | %rbp | callee-saved* | Yes |
| <i>lee</i> ₁₁ | %r12 | callee-saved | Yes |
| <i>lee</i> ₁₂ | %r13 | callee-saved | Yes |
| <i>lee</i> ₁₃ | %r14 | callee-saved | Yes |
| <i>lee</i> ₁₄ | %r15 | callee-saved | Yes |
| | %rsp | stack pointer | Yes |

The starred registers have a potentially relevant alternative use. %a1 (the lower 8 bits of %rax) contains the number of floating point arguments on the stack in a call to varargs functions. %rbp is the frame pointer for the stack frame, in an x86-like calling convention (which is optional for the x86-64).

5 Typical Calling Sequence

If we have 6 or fewer arguments, a typical calling sequence for 32-bit arguments with an instruction

$$d \leftarrow f(s_1, s_2, s_3)$$

will have the following form:

```

arg3 ← s3
arg2 ← s2
arg1 ← s1
call f
d ← res0

```

First we move the temps into the appropriate argument registers, then we call the function f (represented by a symbolic label), and then we move the result register into the desired destination.

This organization, perhaps just before register allocation, has the advantage that the live ranges of fixed registers (called *precolored nodes* in register allocation) is minimized. This is important to avoid potential conflict. We have already applied a similar technique in the implementation of `div` and `mod` operations, which expect their arguments in fixed registers.

Let us state this as a fundamental principle of code generation that you should strive to adhere to:

The live range of precolored registers should be as short as possible!

We can now see a problem with our previous calculation of `def` and `use` information: the above sequence to actually implement the function call will overwrite the argument registers `%edx`, `%esi`, `%edi` as well as the result register `%eax` (the lower 32bits of the return register `%rax`)! In fact, any of the argument registers, the result register, as well as `%r10` and `%r11`, may not be preserved across function calls and therefore have to be considered to be *defined* by the call. If we represent this in the low-level intermediate language, we would *add* to the rule R_8 the following rule R'_8 :

$$\frac{l : \text{call } f \quad \text{caller-save}(r)}{\text{def}(l, r)} J'_8$$

where `caller-save(r)` is true of register r among `%rax`, `%rdi`, `%rsi`, `%rdx`, `%rcx`, `%r8`, `%r9`, `%r10`, and `%r11`.

Here we assume that register aliasing is handled correctly, that is, the register allocator understands that, for example, `%eax` constitutes the lower 32 bits of `%rax`.

Note that all *argument registers* and the *result register* are *caller-save*. This is justified by the fact that we often compute a value for the purposes of passing it into a function, but we do not require that value afterwards. Of course, the result register has to be caller-save, since it will be defined by the called function before it returns.

We refer to argument registers more abstractly as $arg_1, arg_2, \dots, arg_6$ and ler_7 and ler_8 for the other two caller-save registers (even if they are not used for passing arguments to a function). We refer to the result register `%rax` as res_0 .

Now if a temp t is live after a function call, we have to add an inference edge connecting t with any of the fixed registers noted above, since the value of those registers are not preserved across a function call.

The other fixed use of argument registers is of course at the beginning of a function. Again, we should be careful to generate code that keeps the live ranges of precolored registers short. We can accomplish this by moving the argument

registers into temps. Under some heuristics in register allocation and coalescing, these moves can sometimes be eliminated. A function $f(x, y, z)$ might then start with

$$\begin{aligned} f : \\ & x \leftarrow arg_1 \\ & y \leftarrow arg_2 \\ & z \leftarrow arg_3 \end{aligned}$$

One more note: if it is possible that the function f is a function accepting a variable number of arguments, some additional considerations apply. For example, the low 8 bits of `%rax`, called `%a1` hold the number of floating point arguments passed to the function. One therefore sometimes sees `xorl %eax, %eax` before a function call to define zero variable arguments.

6 Callee-Save Registers

The typical calling sequence above takes care of treating caller-save registers correctly. But what about callee-save registers, namely `%rbx`, `%rbp`, `%r12`, `%r13`, `%r14` and `%r15`? In compiling a function we are required that the generated code preserves all the callee-save registers. We generically refer to these registers as lee_i where $9 \leq i \leq 14$.

The standard approach is to save those that are needed onto the stack in the function prologue and restore them from the stack in the function epilogue, just before returning. Of course, saving and restoring them all is safe, but may be overkill for small functions that do not require many registers.

Remember that callee-save registers are essentially live throughout the body of a function, since their value at the return instruction matters. This violates our general rule to keep the live ranges of precolored registers short—in fact, they are maximal!

One simple way to deal with this is by listing them last among the registers to be assigned by register allocation. If we need more than the available number of caller-save registers, we assign callee-save registers before we resort to spilling, but make sure the save them at the beginning of a function and restore them at the end. This is generally more efficient than the usual register spilling since such temps still live in a register throughout the function execution. We use this technique in the example in Section 7.

Another solution is to let register allocation together with register coalescing do the job for us. We can move the contents of all the callee-save registers into temps at the beginning of a function and then move them back at the end. If it turns out these temps are spilled, then they will be saved onto the stack. If not, they may be moved from one register to another and then back at the end. However, this only works well with the right heuristics for assigning registers or using

register coalescing.² Register coalescing consults the interference graph to check if we can assign the same register for variable-to-variable moves. Another optimization that can eliminate register-to-register moves is copy propagation, covered in a later lecture. However, copy propagation requires care because it might extend the live range of variables, possibly undoing the care we applied to keep precolored registers contained.

With this technique, the general shape of the code for a function f before register allocation would be

```

f :
    t1 ← lee9
    t2 ← lee10
    ...
    function body
    ...
    lee10 ← t2
    lee9 ← t1
    ret

```

One complication with this approach is that we need to be sure to spill the full 64-bit registers, while registers holding 32-bit integer values might be saved and restored (or directly used as operands) using only 32 bits. Looking ahead, we see that we will need both 32 bit and 64 bit registers and spill slots in the next lab, so we might decide to introduce this complication now. Or we can still treat callee-save registers specially and switch over to a more uniform treatment in the next lab.

With either of the techniques for using callee-save registers, the one additional rule (R'_8) is not enough. We should also note that all *callee-save* registers should be considered *live* at the return instruction.

$$\frac{l : \text{ret } s \quad \text{callee-save}(r)}{\text{use}(l, r)} J'_2$$

We already know, by prior rule, that s itself is live at l . The new rule J'_2 correctly flags all callee-save registers as live throughout the function body, unless they are assigned somewhere. The code pattern above achieves exactly that, cutting their live ranges down to a minimum.

7 An Extended Example

We use the recursive version of the power function as an example to illustrate register allocation in the presence of function calls. The C0 source is on the left; the

²One technique for register coalescing is briefly described in Section 8 of [Lecture 3](#).

abstract assembly on the right. We use a different form of conditional jump in this example, and just proceed to the next instruction if the conditional evaluates to false.

```

int pow(int b, int e)          pow(b,e):
//@requires e >= 0;           if (e == 0) then done else recurse
{                               done:
  if (e == 0)                 ret 1
    return 1;                 recurse:
  else                         t0 <- e - 1
    return b * pow(b, e-1);    t1 <- pow(b, t0)
}                               t2 <- b * t1
                               ret t2

```

First, we convert it to SSA form. Looking at the right, we see it is already in static single assignment form! Looking on the left, we see a purely functional program. Since purely functional programs do not perform assignment, they must already be in SSA form!

Next, we perform liveness analysis. We proceed backward through the program to compute the following information.

| program | def | use | live-in |
|--|----------------|-------------------|-------------------|
| pow(b, e) : | b, e | | |
| if (e == 0) then done else recurse | | e | b, e |
| done : | | | |
| ret 1 | | | |
| recurse : | | | b, e |
| t ₀ ← e - 1 | t ₀ | e | b, e |
| t ₁ ← pow(b, t ₀) | t ₁ | b, t ₀ | b, t ₀ |
| t ₂ ← b * t ₁ | t ₂ | b, t ₁ | b, t ₁ |
| ret t ₂ | | t ₂ | t ₂ |

Next, we move to a lower-level representation, making the precolored registers explicit with the code pattern in Section 5. Keeping in line with our understanding of function calls, every argument and caller-saved register is defined by the call

instruction.

| program | def | use | live-in |
|--|---|----------------|-------------------|
| pow : | arg_1, arg_2 | | |
| $b \leftarrow arg_1$ | b | arg_1 | arg_1, arg_2 |
| $e \leftarrow arg_2$ | e | arg_2 | b, arg_2 |
| if ($e == 0$) then done else recurse | | | b, e |
| done : | | | |
| $res_0 \leftarrow 1$ | res_0 | | |
| ret | | res_0 | res_0 |
| recurse : | | | b, e |
| $t_0 \leftarrow e - 1$ | t_0 | e | b, e |
| $arg_2 \leftarrow t_0$ | arg_2 | t_0 | b, t_0 |
| $arg_1 \leftarrow b$ | arg_1 | b | b, arg_2 |
| call pow | $res_0, arg_1, arg_2,$ $arg_3, arg_4, arg_5,$ arg_6, ler_7, ler_8 | arg_1, arg_2 | b, arg_1, arg_2 |
| $t_1 \leftarrow res_0$ | t_1 | res_0 | b, res_0 |
| $t_2 \leftarrow b * t_1$ | t_2 | b, t_1 | b, t_1 |
| $res_0 \leftarrow t_2$ | res_0 | t_2 | t_2 |
| ret | | res_0 | res_0 |

We have not made any callee-save registers explicit yet, in the hope we will not need them. After all, there are only two variables and three temps in the program, but we have eight caller-save registers.

Next, we build the interference graph. For each line l and each temp t defined at l , we create an edge between t and any variable live in the successor. The only exception is a move $t \leftarrow s$, where we don't create an edge between t and s because they could be consistently be assigned to the same register. We find that only b interferes with other temps and precolored registers.

| temp | interfering with |
|-------|--|
| b | $res_0, arg_1, arg_2, arg_3, arg_4, arg_5, arg_6, ler_7, ler_8, e, t_0, t_1$ |
| e | b |
| t_0 | b |
| t_1 | b |
| t_2 | |

All precolored registers implicitly interfere with each other, so we don't include that in the interference graph. We are left with no caller-saved registers where it is possible to store b . Rather than storing it on the stack, we can store it in one of our callee-saved registers lee_9 , which we will therefore need to push when we enter into the function and pop before returning. We also add an epilouge, `exitpow`, to keep us from having to write the pop operations multiple times in the program.

| program | live-in |
|--|-----------------------|
| pow : | arg_1, arg_2, lee_9 |
| push lee_9 | arg_1, arg_2, lee_9 |
| $b \leftarrow arg_1$ | arg_1, arg_2 |
| $e \leftarrow arg_2$ | b, arg_2 |
| if ($e == 0$) then done else recurse | b, e |
| done : | |
| $res_0 \leftarrow 1$ | |
| goto exitpow | res_0 |
| recurse : | b, e |
| $t_0 \leftarrow e - 1$ | b, e |
| $arg_2 \leftarrow t_0$ | b, t_0 |
| $arg_1 \leftarrow b$ | b, arg_2 |
| call pow | b, arg_1, arg_2 |
| $t_1 \leftarrow res_0$ | b, res_0 |
| $t_2 \leftarrow b * t_1$ | b, t_1 |
| $res_0 \leftarrow t_2$ | t_2 |
| goto exitpow | res_0 |
| exitpow : | res_0 |
| pop lee_9 | res_0 |
| ret | lee_9, res_0 |

While the callee-save $lee_{10}, \dots, lee_{14}$ are still (implicitly) live through this function because they are all needed at the final return, after the rewrite lee_9 no longer is. Therefore, it no longer interferes with any temps.

We can construct a *simplicial elimination ordering*, from the interference graph, such as:

$$b, e, t_0, t_1, t_2$$

We order the colors (machine registers) as

$$res_0, arg_1, \dots, arg_6, ler_7, ler_8, lee_9$$

with the idea that caller-save registers come first (including argument registers which we will likely need anyway), followed by the only callee-save register we are currently permitted to use. If we needed more, we would first have to spill and restore them.

From this we construct the assignment

$$\begin{aligned} b &\mapsto lee_9 \\ e &\mapsto res_0 \\ t_0 &\mapsto res_0 \\ t_1 &\mapsto res_0 \\ t_2 &\mapsto res_0 \end{aligned}$$

Applying the substitutions:

```

pow :
    push lee9
    lee9 ← arg1
    res0 ← arg2
    if (res0 == 0) then done else recurse
done :
    res0 ← 1
    goto exitpow
recurse :
    res0 ← res0 - 1
    arg2 ← res0
    arg1 ← lee9                (redundant)
    call pow
    res0 ← res0                (redundant)
    res0 ← lee9 * res0
    res0 ← res0                (redundant)
    goto exitpow
exitpow :
    pop lee9
    ret

```

There are now some redundant instructions that can be eliminated. The self-moves are obvious, and one line becomes a self-move if we notice that it is moving a value from *lee₉* into *arg₁* that was definitely already there (copy propagation). Using GNU (AT&T) assembler format for x86-64, we end up with:

```

pow:    pushq    %rbx
        movl    %edi, %ebx
        movl    %esi, %eax
        cmpl   $0, %eax
        jne    L1
        movl   $1, %eax
        goto   L2
L1:    subl   $1, %eax
        movl   %eax, %esi
        call  pow
        imull  %ebx, %eax
L2:    popq   %rbx
        ret

```

References

- [BO16] Randal E. Bryant and David R. O'Hallaron. *Computer Systems, A Programmer's Perspective (Third Edition)*. Pearson, 2016.
- [MHJM09] Michael Matz, Jan Hubička, Andreas Jaeger, and Mark Mitchell. System V application binary interface, AMD64 architecture processor supplement. Available at <http://refspecs.linuxfoundation.org/elf/x86-64-abi-0.99.pdf>, May 2009. Draft 0.99.