

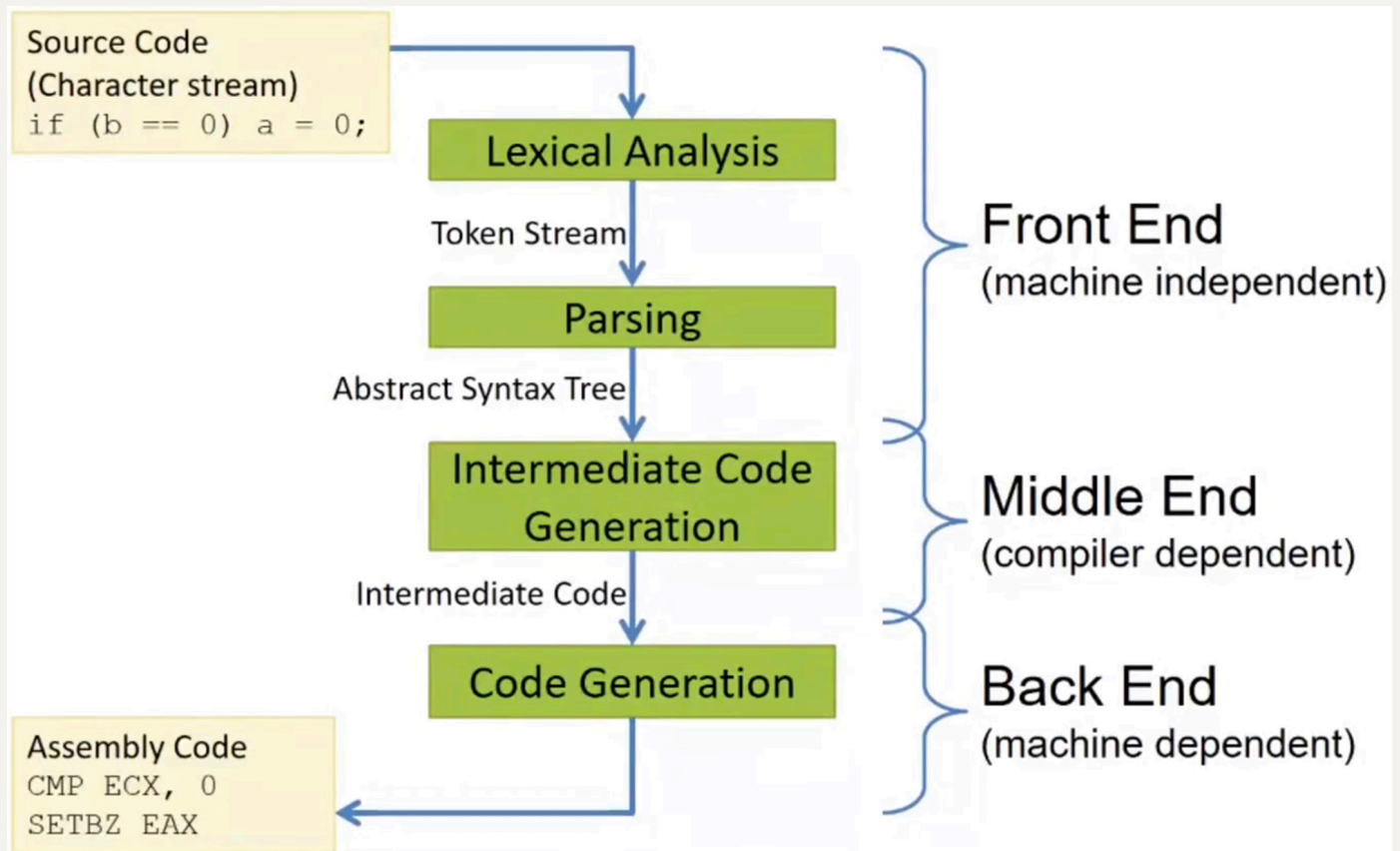
Compiler Design – Lecture note week 2

- Author: Ruben Schenk
- Date: 11.10.2021
- Contact: ruben.schenk@inf.ethz.ch

3. X86 LITE

3.1 Simplified Compiler Structure

A simplified compiler structure looks as follows:



3.2 X86 vs. X86Lite

x86 assembly is very complicated:

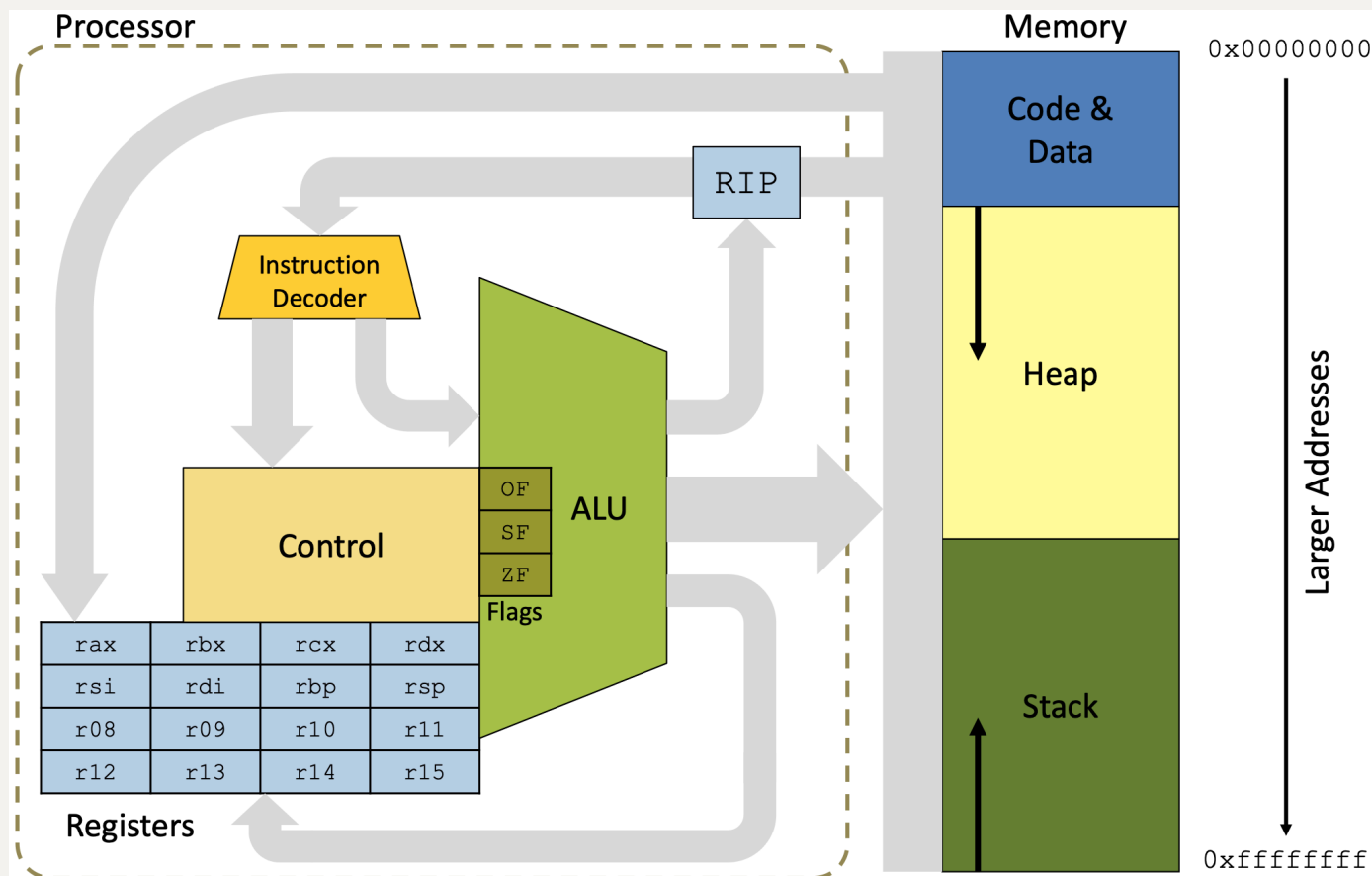
- 8-, 16, 32-, and 64-bit values + floating point, etc.
- Intel 64 and IA 32 have a *huge* number of functions
- For machine code, the instruction range is in size from 1 to 17 bytes

x86Lite assembly is a very simple subset of X86

- Only 64-bit signed integers (no floating point, no 16-bit, no etc.)
- Only about 20 instructions
- Sufficient as a target language for general-purpose computing

3.3 X86 Schematic

The X86 schematic looks as follows:



3.3.1 Registers

There are three special **registers**:

- `rip`: The *instruction pointer*, holds the address of the next instruction
- `rbp`: The *base pointer*, used for call-stack manipulation
- `rsp`: The *stack pointer*, used for call-stack manipulation

3.3.2 Memory

The memory consists of three parts:

- **Code & Data**: Holds the actual program instructions as well as program constants and globals
- **Stack**: Used for function calls and local variables
- **Heap**: Dynamically allocated memory, e.g. via calls to `malloc()`

3.4 Instructions

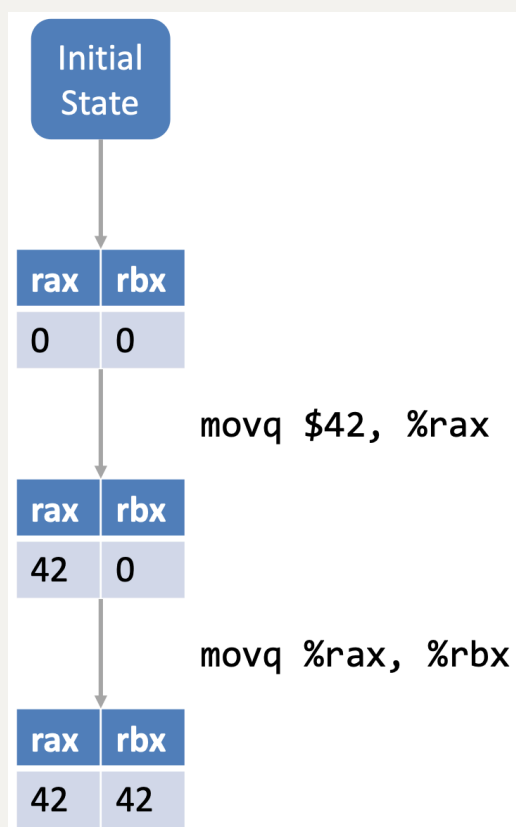
3.4.1 `mov`

The `mov` instructions is of the following form:

```
1  movq SRC, DEST
```

Here, `SRC` and `DEST` are *operands*. `DEST` is treated as a location, either a register or a memory address. `SRC` is treated as a value and is the content of either a register or a memory address or an immediate constant or a label.

Example of a `mov` instruction:



A Note About Instruction Syntax

The most important note is that we have the source *before* the destination. Furthermore:

- Immediate values are prefixed with `$`
- Registers are prefixed with `%`
- Mnemonic suffixed (`movq` vs ``mov`):

- `q` -> quadword (4 words)
- `l` -> long (2 words)
- `w` -> word (16 bits)
- `b` -> byte (8 bits)

3.4.2 X86 Operands

| TYPE | DESCRIPTION | EXAMPLE |
|------|---|--|
| Imm | 64-bit literal signed integer ("immediate") | <code>move \$4, %rax</code> |
| Lbl | a "label" representing a machine address | <code>call FOO</code> |
| Reg | one of the 16 registers | <code>move %rbx, %rax</code> |
| Ind | machine address: [base:Reg][index:reg] disp:int32 | <code>move 12(%rax, %rcx), %rbx</code> |

3.4.3 Arithmetic Instructions

| INSTRUCTION | DESCRIPTION | EXAMPLE | NOTES |
|-----------------------------|--|------------------------------|--|
| <code>negs DEST</code> | 2's complement negation | <code>negs %rax</code> | |
| <code>add SRC, DEST</code> | <code>DEST <- DEST + SRC</code> | <code>add %rbx, %rax</code> | |
| <code>Subq SRC, DEST</code> | <code>DEST <- DEST - SRC</code> | <code>subq \$4, %rsp</code> | |
| <code>Imulq SRC, Reg</code> | <code>Reg <- Reg * Src</code> (truncated 128-bit mult.) | <code>imulq \$2, %rax</code> | Reg must be a register, not a memory address |

3.4.4 Logical/Bit Manipulation Instructions

| Instruction | Explanation | Example | Notes |
|-----------------------------|---|------------------------------|------------------------|
| <code>notq DEST</code> | logical negation | <code>notq %rax</code> | bitwise not |
| <code>andq SRC, DEST</code> | <code>DEST <- DEST & SRC</code> | <code>andq %rbx, %rax</code> | bitwise and |
| <code>orq SRC, DEST</code> | <code>DEST <- DEST SRC</code> | <code>orq \$4, %rsp</code> | bitwise or |
| <code>xorq SRC, DEST</code> | <code>DEST <- DEST xor SRC</code> | <code>xorq \$2, %rax</code> | bitwise xor |
| <code>sarq Amt, DEST</code> | <code>DEST <- DEST >> Amt</code> | <code>sarq \$4, %rax</code> | arithmetic shift right |
| <code>shlq Amt, DEST</code> | <code>DEST <- DEST <<< Amt</code> | <code>shlq %rbx, %rax</code> | logical shift left |
| <code>shrq Amt, DEST</code> | <code>DEST <- DEST >>> Amt</code> | <code>shrq \$1, %rsp</code> | logical shift right |

3.4.5 Condition Flags & Codes

Some X86 instructions set flags as side effects:

- **OF**: *overflow* is set when the result is too big/small to fit in a 64-bit register
- **SF**: *sign* is set to the sign of the result (0 means positive, 1 means negative)
- **ZF**: *zero* is set when the result is 0

From these three flags, we can define **condition codes**. If we want to compare `SRC1` to `SRC2`, we compute `SRC1 - SRC2`. We can then define the following condition codes based on the resulting condition flags:

| CODE | CONDITION |
|-------------------------|------------------------|
| e (equality) | ZF is set |
| ne (inequality) | (not ZF) |
| g (greater than) | (not ZF) and (SF = OF) |
| l (less than) | SF <> OF |
| ge (greater or equal) | (SF = OF) |
| le (less than or equal) | SF <> OF or ZF |

3.4.6 Conditional Instructions

We might write conditional instructions in the following way in X86:

```
1 // Conditional instruction in C
2 if (a == b) {
3     // something
4 } else {
5     // somethingElse
6 }
7 // commonCode
```

```
1 ; Conditional instruction in x86
2 cmpq %rax, %rbx
3 je
4
5 somethingElse:
6     <instruction>
7     ;...
8     jmp commonCode
9
```

```

10 something:
11     <instruction>
12     ;...
13
14 commonCode:
15     <instruction>
16     ;...

```

We support the following three **conditional instructions**:

| INSTRUCTION | DESCRIPTION |
|------------------------------|--|
| <code>cmpq SRC2, SRC1</code> | Compute <code>SRC1 - SRC2</code> , set condition flags |
| <code>setbCC DEST</code> | <code>DEST</code> 's lower byte <code><- if CC then 1 else 0</code> |
| <code>jCC SRC</code> | <code>rip <- if CC then SRC else fallthrough</code> |

3.4.7 Code Blocks and Labels

x86 assembly code is organized into **labeled blocks**. Labels indicate code locations than can be jump targets. Labels are translated away by the linker and loader -- instructions live in the *code segment*.

An x86 program begins executing at a designated code label (usually `main`).

3.4.8 Jumps, Call and Return

We might code function calls in the following way in x86:

```

1 void bar() {
2     // ...
3 }
4
5 void foo() {
6     // ...
7     bar();
8 }

```

```

1  bar:
2      <instruction>
3      ;...
4      <instruction>
5      ret
6
7  foo:
8      <instruction>
9      ; ...
10     <instruction>
11     call bar
12     ;...

```

The different instructions one might use are given by the following table:

| INSTRUCTION | DESCRIPTION | NOTES |
|-----------------------|---|---|
| <code>jmp SRC</code> | <code>rip <- SRC</code> | Jump to location in <code>SRC</code> |
| <code>call SRC</code> | Push <code>rip</code> , <code>rip <- SRC</code> (call a procedure) | Push the program counter to the stack (decrementing <code>rsp</code>), and then jump to the machine instruction at the address given by <code>SRC</code> |
| <code>ret</code> | Pop into <code>rip</code> (return from procedure) | Pop the current top of the stack into <code>rip</code> (incrementing <code>rsp</code>). This instruction effectively jumps to the address at the top of the stack. |

3.5 x86Lite Addresses

3.5.1 x86Lite Addressing

We show how **addressing** in x86Lite works with the following simple example:

```

1  long a[0, 42, 2020;]
2
3  long b = (long)a;      // b = address(a)
4  long b = *a;           // b = a[0] = 0
5  long b = *(a+2);       // b = a[2] = 2020
6
7  long c = 1;
8  long b = a[c];         // b = 42
9  long b = a[c+1];       // b = 2020

```



```

1 ; Array [0, 42, 2020]
2 ; Array address 0xBEEF
3
4 movq %0xBEEF, %rax
5
6 movq %rax, %rbx ; rbx = 0xBEEF
7 movq (%rax), %rbx ; rbx = 0
8 movq 16(%rax), %rbx ; rbx = 2020
9
10 movq $1, %rcx
11 movq (%rax, %rcx), %rbx ; rbx = 42
12 movq 8(%rax, %rcx), %rbx ; rbx = 2020

```

In general, there are three components to an **indirect address**:

- *Base*: a machine address stored in a register
- *Index*: a variable offset from the base
- *Disp*: a constant offset (displacement) from the base

We therefore have: $\text{addr}(\text{ind}) = \text{Base} + [\text{Index} * 8] + \text{Disp}$. When used as a location, `ind` denotes the address $\text{addr}(\text{ind})$. When used as a value, `ind` denotes $\text{Mem}[\text{addr}(\text{ind})]$, the contents of the memory address.

Examples:

| EXPRESSION | ADDRESS |
|----------------------------|--------------------------------|
| <code>-8(%rsp)</code> | <code>rsp - 8</code> |
| <code>(%rax, %rcx)</code> | <code>rax + 8 * rcx</code> |
| <code>8(%rax, %rcx)</code> | <code>rax + 8 * rcx + 8</code> |

3.5.2 x86Lite Memory Model

The x86Lite memory consists of 2^{64} bytes numbered `0x00000000` through `0xffffffff`. The memory is treated as consisting of 64-bit (8 byte) words. Therefore: *legal x86Lite memory addresses consists of 64-bit, quadword-aligned pointers*. This means, that all memory addresses are evenly divisible by 8.

To load a pointer into `DEST`, we use `leaq Ind, DEST` (`DEST <- addr(Ind)`).

By convention, the stack grows from high addresses to low addresses.

The register `rsp` points to the top of the stack:

- `pushq SRC: rsp <- rsp - 8; Mem[rsp] <- SRC`
- `popq DEST: DEST <- Mem[rsp]; rsp <- rsp + 8`

Here is a nice website to explore assembly code given some code snippet in another language:
[Compiler Explorer](#)

3.6 Example: Handcoding x86Lite

Let's look at how we would implement the **factorial** function in **x86Lite**:

```

1  ; long factorial(long i) {
2  ;   if (i > 11) {
3  ;     return i * factorial(i-11);
4  ;   }
5  ;   return 11;
6  ;}
7
8  .text
9  .global factorial
10
11 factorial:
12   ; i is in %rdi
13
14   ; boilerplate
15   pushq   %rbp
16   movq    %rsp, %rbp
17
18   ; if (i > 11)
19   cmpq    $1, %rdi      ; computes %rdi - 1
20   jle     .BASECASE     ; if (i <= 1)
21
22   ; (i > 11) holds at this point
23   pushq   %rdi          ; stores the current value of i on top of the
calls tack
24
25   subq    $1, %rdi
26   callq   factorial
27
28   ; %rax holds factorial(i - 1)
29   popq    %rdi          ; %rdi holds again the original value of i
30   imulq   %rdi, %rax     ; i * factorial(i - 1)
31
32   jmp     .EXIT
33

```

```

34 .BASECASE:
35     movq    $1, %rax
36
37 .EXIT:
38     ; rest of boilerplate
39     movq    %rbp, %rsp
40     popq    %rbp
41     ret     ; return the value in %rax
42
43 .data

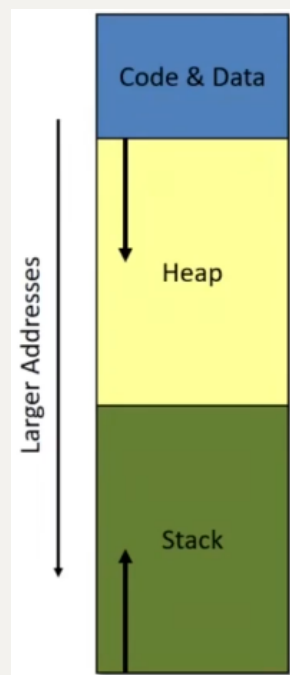
```

Remark: By convention, compilers often use a `.` in front of a label that is internal, i.e. not a global label (compare `factorial` to `.EXIT` in the code above).

3.7 Programming in x86Lite

3.7.1 Three parts of the C memory model

We want to quickly revisit the three different parts of the C memory model, shown in the picture below.



- The **code & data** (or `.text`) segment: contains compile code, constant strings, etc.
- The **heap**: stores dynamically allocated objects, is allocated via `malloc` and deallocated via `free`
- The **stack**: stores local variables, the return address of a function and other bookkeeping information

3.7.2 Local vs. Temporal Variable Storage

We somehow need space to store things like global variables, values passed as arguments to procedures, and local variables. The processor provides two options for storing stuff:

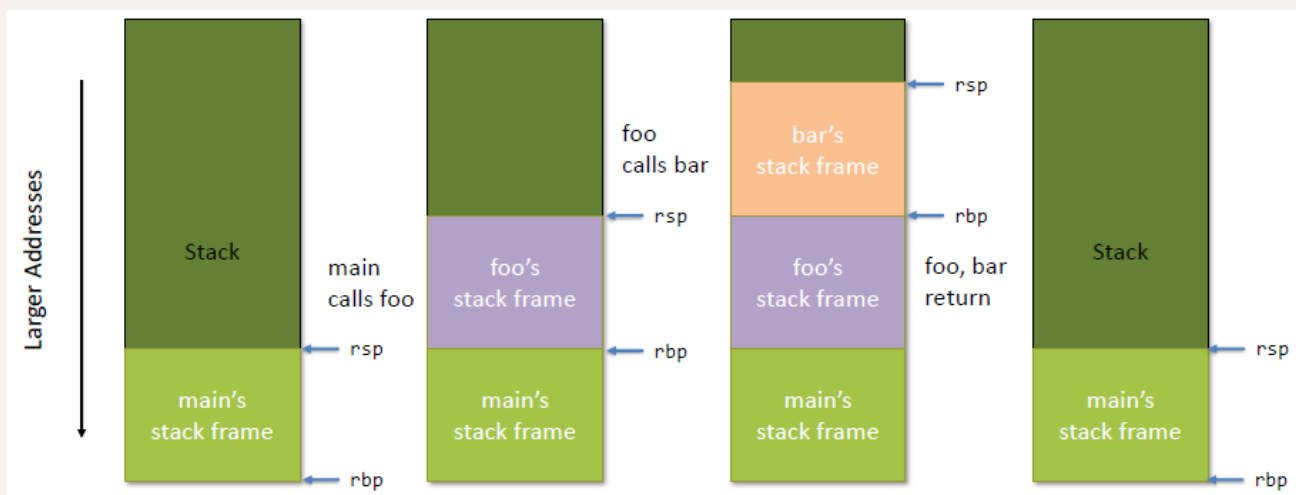
- *Registers*: fast, small size, very limited number
- *Memory (Stack)*: slow, very large amount of space

Example:

```
1  ; int i = 5
2
3  ; Option 1:
4  ; store to a register
5  ; register is "blocked"
6  movq    $5, %rax
7
8  ; Option 2:
9  ; store on the stack
10 subq    $8, %rsp
11 movq    $5, (%rsp)
12 ;...
13 movq    (%rsp), %rax
```

3.7.3 The Stack

The following picture shows how we use the stack in a program with different calls. This corresponds to the "boilerplate" code in the previous example with the `factorial`. We adjust the pointers to the bottom and the top of the stack before and after calling a "function", such that the function has its own **stack frame**.



3.7.4 Calling Conventions

The **calling conventions** cover three main topics:

- Specify the locations of arguments
 - Passed to a function, and
 - Returned by a function
- Designate registers as either
 - Caller Save -- e.g. freely usable by the called code
 - Callee Save -- e.g. must be restored by the called code
- Define the protocol for deallocating stack-allocated arguments, either
 - Caller cleans up
 - Callee cleans up

The widely used calling conventions for x86-64 systems are as follows:

- Callee save registers: `rbp, rbx, r12-r15`
- Caller save: all others
- Call parameters:
 - Parameter 1-6: `rdi, rsi, rdx, rcx, r8, r9`
 - Parameter 7+: on the stack (in right-to-left order), thus, for `n > 6`, the nth argument is at `((n - 7) * 8 + rbp)`
- Return value is in `rax`
- 128-byte "red zone" -- scratchpad for the callee's data