CS4212: Compiler Design

Week 2: SIMPLE interpreter and x86Lite

Ilya Sergey

ilya@nus.edu.sg

ilyasergey.net/CS4212/

Plan for today

- An interpreter for simple imperative language
- A translator to OCaml
- Basics of x86 assembly

• How to represent programs as data structures.

• How to write programs that process programs.



Everyone's Favorite Function

language that we'll call "SIMPLE"

X = 6;ANS = 1;whileNZ (X) { ANS = ANS * X;X = X + -1;}

- We need to describe the constructs of this SIMPLE

 - Semantics: what is the meaning (behavior) of a legal "program"?

• Consider this implementation of factorial in a hypothetical programming

(Simple IMperative Programming LanguagE):

Syntax: which sequences of characters count as a legal "program"?

"Object" vs. "Meta" language

Object language:

the language (syntax / semantics) being described or manipulated

Today's example: SIMPLE

Course project: $OAT \Rightarrow LLVM \Rightarrow x86asm$

Clang compiler: $C/C++ \Rightarrow LLVM \Rightarrow x86asm$

Metacircular interpreter: lisp Metalanguage:

the language (syntax / semantics) used to *describe* some object language

interpreter written in OCaml

compiler written in OCaml

compiler written in C++

interpreter written in lisp

Grammar for a Simple Language

BNF grammars are themselves domain-specific metalanguages for describing the syntax of other languages... oncrete syntax (grammar) for the Simple language: Vritten in "Backus-Naur form"

exp> and <cmd> are nonterminals

:=', '|', and <...> symbols are part of the metalanguage eywords, like 'skip' and 'ifNZ' and symbols, like '{' and '+' re part of the *object language*

eed to represent the *abstract syntax* e. hide the irrelevant of the concrete syntax)

Implement the *operational semantics* (i.e. define the behavior, or meaning, of the program)

Demo: Interpreters in OCaml

- Interpreting expressions

• <u>https://github.com/cs4212/week-01-simple-2024</u>

• Translating Simple programs to OCaml programs

Week 2: x86Lite

(the target architecture for this module)

Intel's X86 Architecture

- 1978: Intel introduces 8086 \bullet
- 1982: 80186, 80286 \bullet
- 1985: 80386
- 1989: 80486 (100MHz, 1μm)
- 1993: Pentium
- 1995: Pentium Pro
- 1997: Pentium II/III \bullet
- 2000: Pentium 4
- 2003: Pentium M, Intel Core \bullet
- 2006: Intel Core 2 \bullet
- 2008: Intel Core i3/i5/i7 \bullet
- 2011: SandyBridge / IvyBridge \bullet
- 2013: Haswell \bullet
- 2014: Broadwell \bullet
- 2015: Skylake (core i3/i5/i7/i9) (2.4GHz, 14nm) \bullet
- 2016: Xeon Phi





Intel Processor Transistor Count



Transistor Count (log scale)

X86 vs. X86lite

- X86 assembly is *very* complicated:
 - 8-, 16-, 32-, 64-bit values + floating points, etc.
 - Intel 64 and IA 32 architectures have a *huge* number of functions
 - "CISC" complex instructions
 - Machine code: instructions range in size from 1 byte to 17 bytes
 - Lots of hold-over design decisions for backwards compatibility
 - instruction-selection level
- X86lite is a *very* simple subset of X86:
 - Only 64 bit signed integers (no floating point, no 16bit, no ...)
 - Only about 20 instructions
 - Sufficient as a target language for general-purpose computing

– Hard to understand, there is a large book about optimizations at just the

X86 Schematic



Larger Addresses

X86lite Machine State: Registers

- Register File: 16 64-bit registers
 - general purpose accumulator rax —
 - rbx base register, pointer to data
 - counter register for strings & loops rcx
 - rdx data register for I/O
 - pointer register, string source register - rsi
 - pointer register, string destination register - rdi
 - base pointer, points to the stack frame rbp ____
 - stack pointer, points to the top of the stack rsp
 - r08-r15 general purpose registers
- a "virtual" register, points to the current instruction rip rip is manipulated only indirectly via jumps and return. ____



Simplest instruction: mov

- movq SRC, DEST
- Here, DEST and SRC are operands
- DEST is treated as a *location*
 - A location can be a register or a memory address
- SRC is treated as a value
 - A value is the *contents* of a register or memory address
 - A value can also be an *immediate* (constant) or a label
- movg %rbx, %rax // move the contents of rbx into rax

copy SRC into DEST

• movg \$4, %rax // move the 64-bit immediate value 4 into rax

A Note About Instruction Syntax

- X86 presented in *two* common syntax formats
- AT&T notation: source before lacksquare
 - Prevalent in the Unix/Mac ecosys
 - Immediate values prefixed with '
 - Registers prefixed with '%'
 - Mnemonic suffixes: movq vs. mov
 - q = quadword (4 words)
 - l = long (2 words)
 - w = word
 - **b** = byte
- Intel notation: destination befo
 - Used in the Intel specification /
 - Prevalent in the Windows ecosys
 - Instruction variant determined by

destination			
stems	movq	\$5 ,	%rax
\$'	movl	\$5,	%eax
7		src	dest

Note: X86lite uses the AT&T notation and the 64-bit only version of the instructions and registers.

ore source	mov	rax,	5
manuals	mov	eax,	5
stem		1	
y register name		dest	Src

X86lite Arithmetic instructions

- two's complement negation negq DEST
- addg SRC, DEST \leftarrow DEST \leftarrow DEST + SRC
- subg SRC, DEST \leftarrow DEST \leftarrow DEST SRC
- imulq SRC, Reg Reg ← Reg * SRC (truncated 128-bit mult.)

Examples as written in:

addq	%rbx,	%rax	// rax	+
subq	\$4, r	sp	// rsp	←

ullet

- rax + rbx
- -rsp-4
- Note: Reg (in imulq) must be a register, not a memory address

X86lite Logic/Bit manipulation Operations

- notq DEST logical negation
- andq SRC, DEST D
- orq SRC, DEST \leftarrow DEST \leftarrow
- xorq SRC, DEST DEST \leftarrow
- sarq Amt, DEST → DEST → DEST >> amt (arithmetic shift right)
- shlq Amt, DEST DEST DEST << amt (arithmetic shift left)
- shrq Amt, DEST → DEST → DEST >>> amt (bitwise shift right)

The difference between arithm preserves the sign.

- DEST ← DEST && SRC
- DEST ← DEST || SRC
- DEST ← DEST xor SRC

The difference between arithmetic and bitwise shift is that the former

X86 Operands

•	Imm	64-bit literal signed ir
		42, 0x3de7
•	Lbl	a "label" representing
		the assembler/linker/
		_factorial, .L2
•	Reg	One of the 16 register
		%rax, %r04
•	Ind	[base:Reg][index:Reg
		"Indirect" machine ad
		(%rax), -8(%rbp)

Operands are the values operated on by the assembly instructions

nteger "immediate"

g a machine address, to be resolved by 'loader

rs, the value of a register is of its contents

g,scale:int32][disp] dress (see next slide)

X86 Addressing

- In general, there are three components of an indirect address
 - a machine address stored in a register – Base:
 - Index * scale: a variable offset from the base
 - a constant offset (displacement) from the base – Disp:
- addr(ind) = Base + [Index * scale] + Disp
 - When used as a *location*, ind denotes the address addr(ind)
 - When used as a *value*, ind denotes Mem[addr(ind)], the contents of the memory address
- Example: -4(%rsp) denotes address: rsp 4
- Example: (%rax, %rcx, 4) denotes address: rax + 4*rcx
- Example: 12(%rax, %rcx, 4) denotes address: rax + 4*rcx +12
- Note: Index cannot be rsp

Note: X86lite does not need this full generality. It does not use index * scale.



X86lite Memory Model

- The X86lite memory consists of 2⁶⁴ bytes.
- X86lite treats the memory as consisting of 64-bit (8-byte) quadwords.
- Therefore: legal X86lite memory addresses consist of 64-bit, quadword-aligned pointers.
 - All memory addresses are evenly divisible by 8
- leag Ind, DEST DEST ← addr(Ind) loads a pointer into DEST
- By convention, there is a stack that grows from high addresses to low addresses
- The register rsp points to the top of the stack
 - pushq SRC $rsp \leftarrow rsp 8; Mem[rsp] \leftarrow SRC$
 - DEST \leftarrow Mem[rsp]; rsp \leftarrow rsp + 8 – popq DEST



Comparisons and Conditioning

X86lite State: Flags & Condition Codes

- X86 instructions set flags as a side effect
- X86lite has only 3 flags: •
 - **OF**: "overflow" set when the result is too big/small to fit in 64-bit reg. —
 - **SF**: "sign" set to the sign or the result (0=positive, 1 = negative)
 - **ZF**: "zero" set when the result is 0
- From these flags, we can define *Condition Codes* \bullet
 - You can think of Cond. Codes as of additional registers, whose value changes depending on the current flags
- E.g., cmpq SRC1, SRC2 computes SRC1 SRC2 to set the flags
- Now we can check conditional codes:
 - eq equality holds when **ZF** is set
 - holds when (not **ZF**) – neq inequality
 - lt less than holds when SF <> OF
 - Equivalently: ((SF && not OF) || (not SF && OF))
 - ge greater or equal holds when (not It) holds, i.e. (SF = OF)
 - le than or equal holds when SF <> OF or ZF
 - gt greater than holds when (not le) holds,
 - i.e. (SF = OF) && not(ZF)



Conditional Instructions

- cmpq SRC1, SRC2 Corr
- setb CC DEST DES
- jCC SRC rip
- Example:
 - cmpq %rcx, %rax Com
 jeq __truelbl If ra

Compute SRC2 – SRC1, set condition flags

DEST's lower byte ← if CC then 1 else 0

rip ← if CC then SRC else do nothing

Compare rax to ecx If rax = rcx then jump to truelbl

Code Blocks & Labels

X86 assembly code is organized into *labeled blocks*: \bullet



Labels indicate code locations that can be jump targets (either through conditional branch instructions or function calls).

- Labels are translated away by the linker and loader instructions live in the heap in the "code segment" \bullet
- An X86 program begins executing at a designated code label (usually "main"). lacksquare

```
<instruction>
<instruction>
<instruction>
<instruction>
<instruction>
<instruction>
```

Basic Control Flow

Jumps, Calls, and Return

- jmp SRC **rip** ← SRC Jump to location in SRC
- callq SRC Push rip; rip ← SRC
 - Call a procedure: Push the program counter to the stack (decrementing rsp) and then jump to the machine instruction at the address given by SRC.
- Pop into rip retq
 - Return from a procedure: Pop the current top of the stack into rip ____ (incrementing rsp).
 - This instruction effectively jumps to the address at the top of the stack



Loop-based Factorial in Assembly

.globl _program		
_program:		
movq	\$1, %rax	
movq	\$6, %rdi	
loop:		
cmpq	\$0, %rdi	
je exi	.t	
imulq	%rdi, %rax	
decq	%rdi	
jmp loc	p	
exit:		
retq		

```
int program() {
    int acc = 1;
    int n = 6;
    while (0 < n) {
        acc = acc * n;
        n = n - 1;
    }
    return acc;
}</pre>
```

Demo: Hand-Coded x86Lite

- https://github.com/cs4212/week-02-x86lite
- Basic definitions: x86.ml lacksquare
- Linking with assembly: test.c
- Example program, simple output, factorial

Compiling, Linking, Running

- To use hand-coded X86:
 - by running make
 - 2. Run it, redirecting the output to some .s file, e.g.: ./main1.native >> prog.s
 - 3. clang -o test test.c prog.s

One M1/M2 mac, use the following flags:

4. ./test

Compile OCaml program main1.ml to the executable

```
Use clang (or gcc) to compile & link with test.c:
```

```
clang -arch x86 64 -o test prog.s test.c
```

You should be able to run the resulting executable:

Next lecture

- Compiling simple programs to X86lite
- Intermediate Representations