

# Compilers and computer architecture: The RISC-V architecture

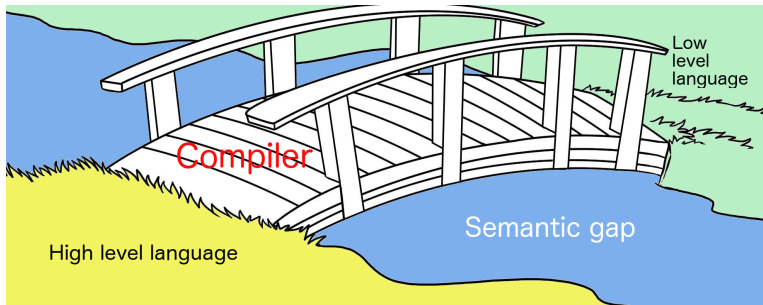
Martin Berger <sup>1</sup>

November 2019

---

<sup>1</sup>Email: [M.F.Berger@sussex.ac.uk](mailto:M.F.Berger@sussex.ac.uk), Office hours: Wed 12-13 in  
Chi-2R312

# Recall the function of compilers



# Introduction

In previous lectures, we focussed on generating code for simple architectures like the stack machine, or accumulator machines.

Now we want to do something more interesting, generating code for a real CPU.

We focus on the RISC-V (family of) processor(s).

# CISC vs RISC

Processors can roughly be classified as

- ▶ **CISC** (complex instruction set computer)
- ▶ **RISC** (reduced instruction set computer)

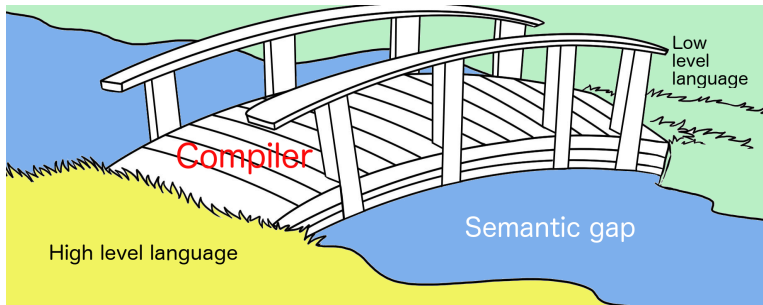
What is the instruction set?

# Instruction set architecture

In a CPU we distinguish between

- ▶ **Instruction set architecture**, that is externally visible aspects like the supported data types (e.g. 32 bit Ints, 80 bit floats etc), instructions, number and kinds of registers, addressing modes, memory architecture, interrupt etc. **This is what the programmer uses to get the CPU to do things.**
- ▶ **Microarchitecture**, which how the instruction set is implemented. The microarchitecture is not visible to the programmer, and CPUs with different microarchitectures can share a common instruction set. For example Intel and AMD support very similar instruction sets (x86 derived) but have very different microarchitectures. **The programmer can ignore this, this is for the hardware people.**

# Recall:



# Instruction set architecture

There is a **semantic gap** between high-level programming languages and machine languages: the former have powerful features (e.g. method invocation) that translate to a large number of simple machine instructions.

In the past it was thought that making machine commands more powerful would close or narrow the semantic gap, making compiled code faster. Examples of powerful machine commands include directly enabling constructs such as procedure calls, or complicated array access in single instructions.

CPUs with such instruction sets are called **CISC** (complex instruction set computer). Complex because the instructions do complicated things and are complex to implement in hardware.

# Instruction set architecture

In some cases CISC architectures led to faster compiled code. But in the 1970s researchers began to study instruction set architecture and compiled code carefully and noticed two things.

- ▶ Compilers rarely make use of the complex instructions provided by CISC machines.
- ▶ Complex operations tended to be slower than a sequence of simpler operations doing the same thing.
- ▶ Often real-world programs spend most of their time executing simple operations.
- ▶ Implementing complex operations leads to complicated CPU architecture that slow down the execution of simple instructions. Worse: simple operations are slower even when you don't use the complex operations.



# RISC

These empirical insights lead to a radical rethink of instruction set architecture.

Optimise as much as possible the few instructions that are used the most in practise. Try and make them exceedingly fast.

This makes the task of the compiler (much) harder, but the compiler has to compile a program only once, whereas the CPU would have to support complex instructions all the time.

# RISC

RISC processors like RISC-V are the outcome.

One key feature of RISC is that external memory was only accessible by a load or store instruction. All other instructions were limited to internal registers. Hence RISC is also called **load/store architecture**.

This drastically simplifies processor design: allowing instructions to be fixed-length, simplifying pipelines, and isolating the logic for dealing with the delay in completing a memory access (cache miss, etc.) to only two instructions.

## RISC vs CISC today

Despite RISC being technically better, still the most popular desktop/server family of chips (Intel x86) is not RISC. (Phones are dominated by RISCish architectures (ARM).) Reasons:

- ▶ Large amount of legacy x86 code (e.g. Microsoft products), locking PC users into x86 CISC. x86 has been backwards compatible back to the 8080 processor (1974).
- ▶ Intel earns much more money than producers of RISC chips so can spend much more on research, design and manufacturing, keeping CISC chips competitive with RISC.
- ▶ 'Under the hood' modern Intel processors are also RISC: the complicated x86 machine instructions are translated at run-time into much simpler underlying RISC microcode (which is not user-visible).
- ▶ Both ARM and x86 provide dedicated instructions for cryptography (AES).

## Aside: extreme RISC

Sometimes a workload is so extremely skewed towards certain commands that it makes sense to abandon traditional CPU architecture (including RISC) and use special purpose processors.

- ▶ GPUs (graphics processing unit) which explore the unusually high degree of SIMD parallelism of most graphics and image processing algorithms.
- ▶ Bitcoin lead to dedicated chips that compute nothing but SHA256.
- ▶ Google's TPU (Tensor Processing Units) for speeding up computation in neural nets (typically 15x - 30x vs GPU), and more energy efficient (70x vs GPU, 200x vs CPU). Basically matrix multiplication and activation function engine. Currently a **lot** of work in this space.

# The RISC-V processor

The original RISC processor was MIPS, John Hennessy in Stanford. For various reasons it has been replaced by **RISC-V**!

RISC-V is **open source** and has an extremely clean and simple design. For those reasons it has emerged as a serious competitor to ARM.

Several industrial strength compilers to RISC-V exist, including LLVM and GCC.

# The RISC-V processor

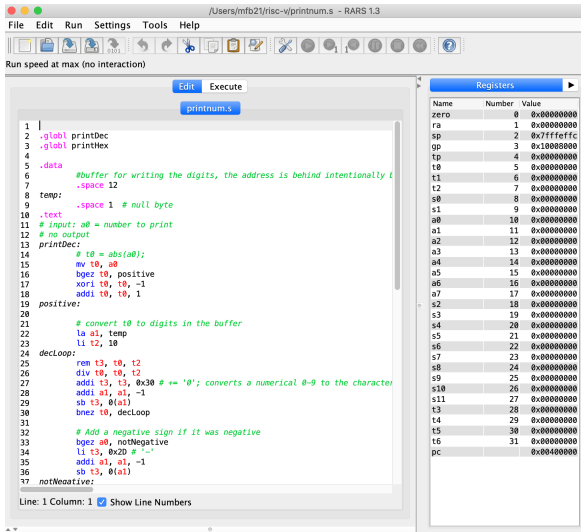
In its basic form, RISC-V has 32 integer registers.

RISC-V is a **load-and-store** architecture, meaning:

Except for memory access instructions, instructions address only registers.

# RARS

I recommend using a RISC-V simulator like RARS for learning RISC-V.



The screenshot shows the RARS simulator window. The title bar indicates the path `/Users/mfb21/risc-v/printnum.s - RARS 1.3`. The menu bar includes File, Edit, Run, Settings, Tools, and Help. The toolbar contains icons for file operations and execution. Below the toolbar, it says "Run speed at max (no interaction)".

The main window is divided into two panes. The left pane, titled "printnum.s", contains assembly code for a program that prints the absolute value of a number in decimal. The right pane, titled "Registers", shows a table of 32 registers.

```
1 |
2 .globl printDec
3 .globl printHex
4
5 .data
6     #buffer for writing the digits, the address is behind intentionally l
7     .space 12
8 temp:
9     .space 1 # null byte
10
11 .text
12 # input: a0 = number to print
13 # no output
14 printDec:
15     # t0 = abs(a0);
16     mv t0, a0
17     bgez t0, positive
18     xori t0, t0, -1
19     addi t0, t0, 1
20 positive:
21     # convert t0 to digits in the buffer
22     la a1, temp
23     li t2, 10
24 decLoop:
25     rem t3, t0, t2
26     div t0, t0, t2
27     addi t3, t3, 0x30 # += '0'; converts a numerical 0-9 to the character
28     addi a1, a1, -1
29     sb t3, 0(a1)
30     bnez t0, decLoop
31
32     # Add a negative sign if it was negative
33     bgez a0, notNegative
34     li t3, 0x2D # '-'
35     addi a1, a1, -1
36     sb t3, 0(a1)
37 notNegative:
```

The "Registers" window displays the following table:

Name	Number	Value
zero	0	0x00000000
ra	1	0x00000000
sp	2	0x7ffffeffc
gp	3	0x10000000
tp	4	0x00000000
t0	5	0x00000000
t1	6	0x00000000
t2	7	0x00000000
s0	8	0x00000000
s1	9	0x00000000
a0	10	0x00000000
a1	11	0x00000000
a2	12	0x00000000
a3	13	0x00000000
a4	14	0x00000000
a5	15	0x00000000
a6	16	0x00000000
a7	17	0x00000000
s2	18	0x00000000
s3	19	0x00000000
s4	20	0x00000000
s5	21	0x00000000
s6	22	0x00000000
s7	23	0x00000000
s8	24	0x00000000
s9	25	0x00000000
s10	26	0x00000000
s11	27	0x00000000
t3	28	0x00000000
t4	29	0x00000000
t5	30	0x00000000
t6	31	0x00000000
pc		0x00400000

At the bottom of the left pane, it says "Line: 1 Column: 1" and "Show Line Numbers" is checked.

# Getting RARS

Homepage: `https://github.com/TheThirdOne/rars`

Download \*.jar at `https://github.com/TheThirdOne/rars/releases/download/v1.3.1/rars1\_3\_1.jar`

Easy to launch: `java -jar rars1_3_1.jar`

Has useful online help.



# RARS

You need to learn RARS by yourself and in tutorials.

# RISC-V

Here is a basic overview of RISC-V. We'll only cover issues that are relevant for code generation. You are expected to familiarise yourself with RISC-V programming on your own.

This should not be difficult with RARS, as RISC-V is an exceptionally clean architecture.

# RISC-V registers

x0 / zero
x1
x2
x3
x4
x5
x6
x7
x8
x9
x10
x11
x12
x13
x14
x15
x16
x17
x18
x19
x20
x21
x22
x23
x24
x25
x26
x27
x28
x29
x30
x31
XLEN
XLEN-1 0
pc

RISC-V has the following registers (all are 32 bits).

- ▶ 32 general purpose registers
- ▶ A program counter (PC)

# RISC-V registers

x0 / zero
x1
x2
x3
x4
x5
x6
x7
x8
x9
x10
x11
x12
x13
x14
x15
x16
x17
x18
x19
x20
x21
x22
x23
x24
x25
x26
x27
x28
x29
x30
x31
XLEN
XLEN-1
0
pc

CPU general-purpose registers have assigned functions:

- ▶ x0 is hard-wired to 0, and can be used as target register for any instruction whose result is to be discarded. x0 can also be used as a source of 0 if needed.
- ▶ x1-x31 are general purpose registers. The 32 bit integers they hold are interpreted, depending on the instruction that access the registers. (Examples: Boolean values, two's complement signed binary integers or unsigned binary integers, stack pointer or return address).

# RISC-V registers: PC

x0 / zero
x1
x2
x3
x4
x5
x6
x7
x8
x9
x10
x11
x12
x13
x14
x15
x16
x17
x18
x19
x20
x21
x22
x23
x24
x25
x26
x27
x28
x29
x30
x31
XLEN
XLEN-1 0
pc

The program counter (PC) register, points to the instruction to be executed next. The PC cannot directly be written or read using load/store instructions. It can only be influenced by executing instructions which change the PC as a side-effect.

## RISC-V registers: do you notice something?

x0 / zero
x1
x2
x3
x4
x5
x6
x7
x8
x9
x10
x11
x12
x13
x14
x15
x16
x17
x18
x19
x20
x21
x22
x23
x24
x25
x26
x27
x28
x29
x30
x31
XLEN
XLEN-1
pc

No explicit SP (and no `push`, no `pop` commands)



Can be simulated with general purpose register and normal commands.

# RISC-V registers: usage conventions

Register name	Symbolic name	Description
32 integer registers		
x0	Zero	Always zero
x1	ra	Return address
x2	sp	Stack pointer
x3	gp	Global pointer
x4	tp	Thread pointer
x5	t0	Temporary / alternate return address
x6–7	t1–2	Temporary
x8	s0/fp	Saved register / frame pointer
x9	s1	Saved register
x10–11	a0–1	Function argument / return value
x12–17	a2–7	Function argument
x18–27	s2–11	Saved register
x28–31	t3–6	Temporary

Interesting for us mostly x1 (ra), x2 (sp), x8 (s0)

Note that those are **usage conventions**. I recommend adhering to them if you want to interface with other RISC-V software (e.g. assembler).

# RISC-V Datatypes

RISC-V has address space of  $2^{32}$  bytes for all memory accesses.

Address space is **circular**, so that the byte at address  $2^{32} - 1$  is adjacent to the byte at address zero.

Memory is byte-addressable.

A word of memory is defined as 32 bits (4 bytes).

A halfword is 16 bits (2 bytes).

A byte 8 bits.



# RISC-V memory alignment

RISC-V has fixed-length 32-bit instructions that **must** be aligned on 32-bit boundaries (i.e. at memory locations divisible by 4).

**But ...** Accessed memory addresses **need not** be aligned, **but accesses to aligned addresses may be faster**; for example, simple CPUs may implement unaligned accesses with slow software emulation driven from an alignment failure interrupt.

The assembler will help you with alignment. We come back to this.

# RISC-V instructions

CPU instructions are organized into the following functional groups:

- ▶ Load and store (memory access)
- ▶ Immediates (handling of constants)
- ▶ Computational (e.g. integer arithmetic and boolean logic)
- ▶ Jump and branch (conditional and unconditional)
- ▶ Many others (e.g. SIMD, vectoring)

Each instruction is 32 bits long in memory.

Important: RISC-V processors use a simple load/store architecture; all operations (e.g. addition, comparison) are performed on operands held in processor registers.

Main memory is accessed **only** through load and store instructions.

# RISC-V instructions

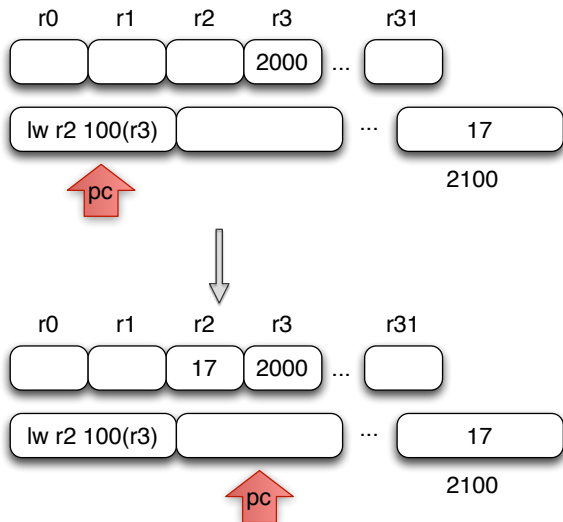
The command

```
lw reg1 offset(reg2)
```

(where `offset` is a 16-bit integer) adds the content `reg2` and the 16 bit value `offset`, obtaining a new number  $n$ , and then looks up the 32 bit value stored in memory at  $n$ . That value is then loaded into register `reg1` as a signed integer.

The sum of `reg2` and `offset` must be word aligned (i.e. the two least significant bits must be 0), otherwise an error will occur.

# RISC-V instructions



# RISC-V instructions

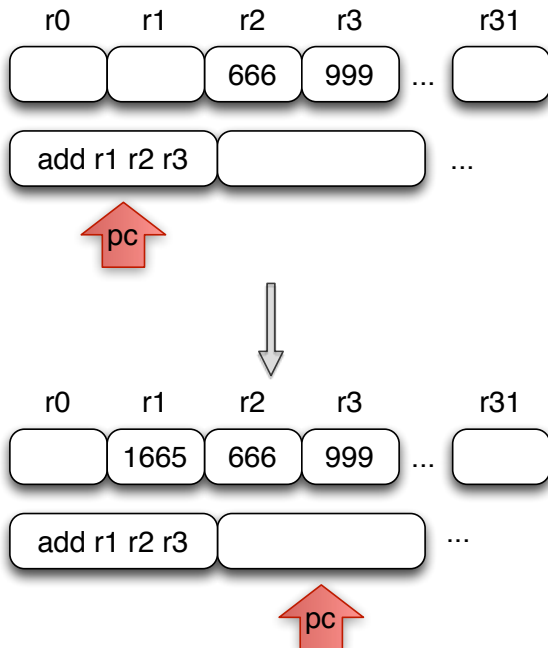
The command

```
add reg1 reg2 reg3
```

Adds the contents of registers `reg2` and `reg3`, and stores the result in `reg1`.

Note that the `reg1`, `reg2` and `reg3` don't have to be distinct.

## RISC-V instructions



# RISC-V instructions

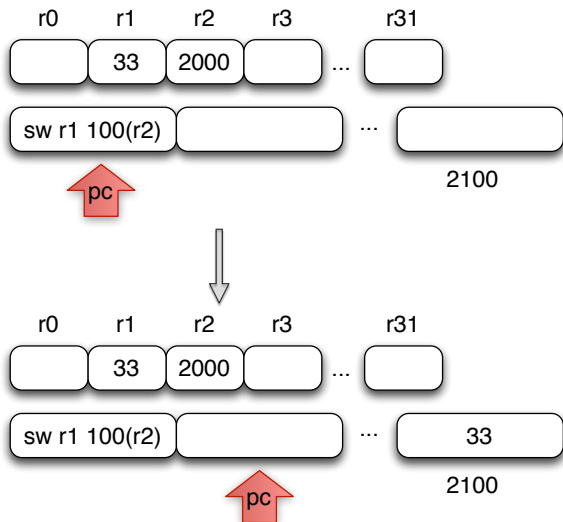
The command

```
sw reg1 offset(reg2)
```

(where `offset` is an integer) stores the 32 bit word currently in `reg1` at the address obtained by adding the 16 bit value `offset` to the content of register `reg2`.

The sum of `reg2` and `offset` must be word aligned (i.e. the two least significant bits must be 0), otherwise an error will occur.

# RISC-V instructions





# RISC-V instructions

The command

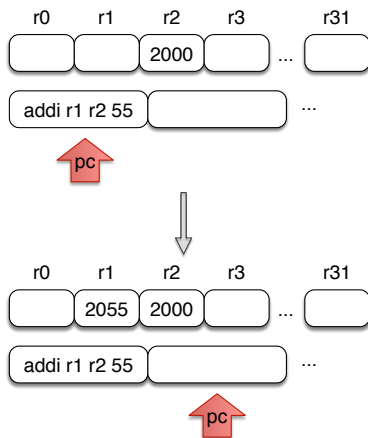
```
addi reg1 reg2 imm
```

Adds the 16 bit **signed** integer `imm` to the word currently in `reg2`, storing the result in register `reg1`. Here the 'u' in `addiu` means **unsigned**. In first approximation that means overflow is not checked when adding (no error is caused by overflowing).

Not checking overflow is useful e.g. when you want 'wrap around' a sum at 0 or  $2^{32} - 1$ . You want this e.g. when doing cryptography. In addition we consider e.g. the SP an unsigned integer.

But `imm` is signed, so we can increment and decrement e.g. the SP.

# RISC-V instructions



# RISC-V instructions

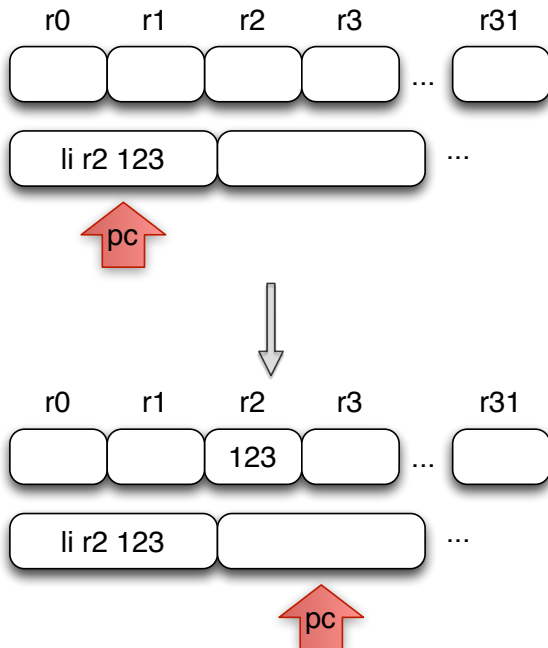
## The **pseudo instruction**

```
li reg imm
```

Stores the **32 bit** integer `imm` in register `reg`.

It is a pseudo instruction in that there is no RISC-V assembly command that directly implements this (MPIS cannot load 32 bit words directly), instead the RISC-V assembler will automatically expand `li reg imm` into a sequence of real assembler commands. When compiling you can easily treat pseudo instructions as real instructions.

## RISC-V instructions



# Our first RISC-V program

Let's write the program  $7+5$ , we want the result in register `r5`.

```
li r6 7
```

```
li r5 5
```

```
add r5 r5 r6
```

## Our second RISC-V program

Let's write  $7+5$ , in accumulator machine form.

- ▶ One argument is in the accumulator.
- ▶ Remaining arguments on the stack.
- ▶ Result should be in accumulator.

Recall RISC-V doesn't have an explicit SP.

Also: no explicit accumulator!

We must simulate both, and that is easy: each register can be used as SP or as accumulator.

## Our second RISC-V program

Let's write  $7+5$ , in accumulator machine form.

Convention: we use register `x2` as stack pointer, and register `x10` as accumulator.

It is customary in RISC-V assembly to write `sp` for the stack pointer (`r29`) and `a0` for register `r4`.

## Our second RISC-V program

Recall that in the accumulator machine model, memory operations work only via the accumulator. With this in mind, here is the program  $7+5$  we are seeking to translate to RISC-V in pseudo-code.

```
acc ← 7
push acc
acc ← 5
acc ← acc + top of stack
pop
```



## Our second RISC-V program

To translate

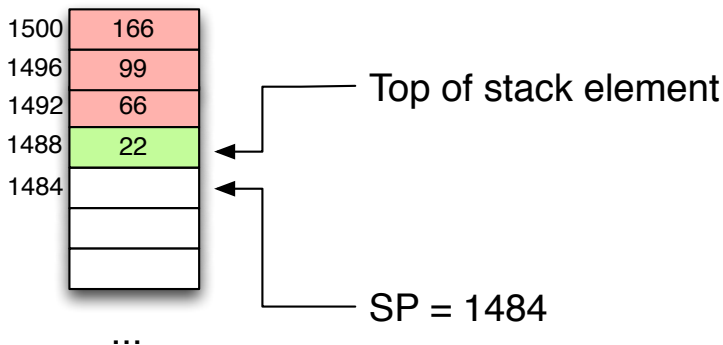
```
acc <- 7
push acc
acc <- 5
acc <- acc + top of stack
pop
```

into RISC-V we adhere to the conventions that

- ▶ The stack grows downwards (i.e. from high to low addresses).
- ▶ The stack pointer `sp` points to the first free memory cell below (in terms of addresses) the top of the stack.

## Our second RISC-V program

- ▶ The stack grows downwards (i.e. from high to low addresses).
- ▶ The stack pointer  $sp$  points to the first free memory cell below (in terms of addresses) the top of the stack.



## Our second RISC-V program

```
acc <- 7
push acc
acc <- 5
acc <- acc+topOfStack
pop
```

```
li a0 7
sw a0 0(sp)
addi sp sp -4
li a0 5
lw t1 4(sp)
add a0 a0 t1
addi sp sp 4
```

Note that the program on the right is really doing almost exactly what we did a few weeks ago when we looked at the accumulator machine, except that

- ▶ we use a temporary `t1`
- ▶ we use RISC-V assembly
- ▶ we have to adjust the stack 'by hand', rather than using built-in `push` and `pop`

# RISC-V

We will soon write a compiler that compiles a simple language with procedures to RISC-V code.

To understand this, you need to familiarise yourself with RISC-V in the tutorials and in self study.

RISC-V machine code is really straightforward, and not really different from the pseudo machine code we used a few weeks back, except that the assembler syntax is slightly different.

## Interlude on (RISC-V) assembler

Assembler language is a programming language that is close to machine language but not the same.

Why bother with yet another language? Why not program straight in machine language?

# That's why

```
0010011110111101111111111111000001010111110111111000000
000001010010101111101001000000000000100000101011111010
0101000000000010010010101111101000000000000000001100010
1011111010000000000000000011100100011111010111000000000
0001110010001111101110000000000000001100000000001110011
1000000000000110010010010111001000000000000000000010010
100100000001000000000110010110101111101010000000000000
01110000000000000000000001111000000100100000001100001111
110010000010000100010100001000001111111111110111101011
1110111001000000000000110000011110000000100000100000000
000010001111101001010000000000011000000011000001000000
000000111011000010010010000100000001000011000010001111
1011111100000000000010100001001111011110100000000001000
000000001111100000000000000000010000000000000000000001
000000100001
```

## That's why

Here is same code written in assembly language, but no symbolic labels are used as name of registers or memory locations.

```
addi $29, $29, -32      mflo $15
sw $31, 20($29)         addu $25, $24, $15
sw $4, 32($29)          bne $1, $0, -9
sw $5, 36($29)          sw $25, 24($29)
sw $0, 24($29)          lui $4, 4096
sw $0, 28($29)          lw $5, 24($29)
lw $14, 28($29)         jal 1048812
lw $24, 24($29)         addi $4, $4, 1072
multu $14, $14          lw $31, 20($29)
addi $8, $14, 1         addi $29, $29, 32
slti $1, $8, 101        jr $31
sw $8, 28($29)          move $2, $0
```

# That's why

It gets even better with **symbolic** names such as `sp` or `loop`.

```
.text                                addu t0, t6, 1
.align 2                             sw t0, 28(sp)
.globl main                          ble t0, 100, loop
main:                                 la a0, str
    subu sp, sp, 32                  lw a1, 24(sp)
    sw ra, 20(sp)                    jal printf
    sd a0, 32(sp)                    move v0, 0
    sw 0, 24(sp)                     lw ra, 20(sp)
    sw 0, 28(sp)                     addu sp, sp, 32
loop:                                jr ra
    lw t6, 28(sp)                    .data
    mul t7, t6, t6                    .align 0
    lw t8, 24(sp)                    str:
    addu t9, t8, t7                   .asciz "The sum from
    sw t9, 24(sp)                     0 .. 100 is %d"
```



# Assembler vs assembly language

We must carefully distinguish between

- ▶ **Assembly language**, the symbolic representation of a computer's binary machine language.
- ▶ **Assembler**, a program (a mini-compiler) that translates assembly language into real machine code (long sequences of 0s and 1s).

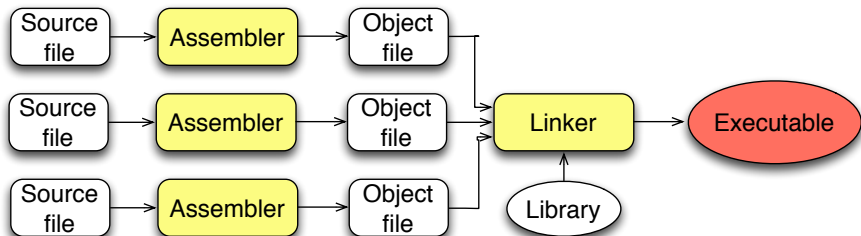
# Assembler, the program

The assembler primarily does two things.

- ▶ Translate commands in assembly language like `addi t3 t6 t8` into machine code.
- ▶ Convert symbolic addresses such as `main` or `loop` into machine addresses such as `100011010011010011010011010101001`. This task is sometimes deferred to the linker.

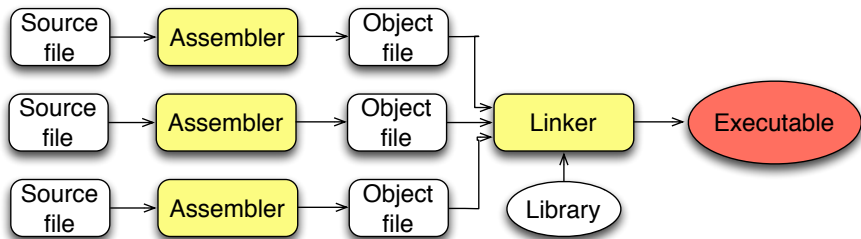
The symbolic addresses in assembly language name commonly occurring bit patterns, such as opcodes and register names, so humans can read and remember them. In addition, assembly language permits programmers to use labels to identify and name particular memory words that hold instructions or data, or that the program can jump to.

# Assembler, the program



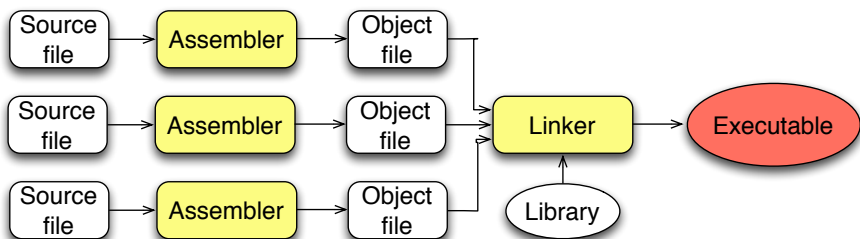
Source files are produced by a compiler. They may contain labels that are not defined in the source file, reference to external code (e.g. `print`).

# Assembler, the program



Assembler translates source files to object files, which are machine code, but contains 'holes' (basically references to external code). Because of holes, object files cannot be executed directly. The holes arise because the assembler translates each file separately.

# Assembler, the program



The linker gets all object files and libraries and puts the right addresses into holes, yielding an executable.

## Assembler, the program

Here is an example of using names: `main` is a **global** name in the sense that other programs can use it. OTOH `loop` is a **local** name: it can only be used (jumped to) inside this program.

```
.text
.align 2
.globl main
main:
    subu sp, sp, 32
    sw ra, 20(sp)
    ...
loop:
    lw t6, 28(sp)
    ...
    ble t0, 100, loop
```

It is the declaration (assembler directive) `.globl main` that makes `main` global.

# Assembler, the program

The assembler processes a source file line by line, translating assembly commands. It keeps track of the size of each command.

```
loop:
    subu sp, sp, 32
    sw ra, 20(sp)
```

When the assembler encounters a line starting with a label, like `loop:` . . . it calculates what address in memory the command just below would be at, and stores the pair of label and address in its symbol table. If it encounters this label later, e.g. `ble t0, 100, loop`, the assembler replaces the label with the address (if local, otherwise the linker does this).

# Helpers

Assembly languages typically offer various features making assembly programming easier. Here are some RISC-V examples.

- ▶ Data layout directives
- ▶ Pseudo instructions
- ▶ Alignment instructions



# Data layout directives

Data layout directives describe data in a more concise and natural manner than its binary representation. Example:

```
.asciz "The sum from 0 .. 100 is %d\n"
```

stores characters from the string in memory. Alternatively we can use the `.byte` directive to obtain the same effect.

```
.byte 84, 104, 101, 32, 115, 117, 109, 32  
.byte 102, 114, 111, 109, 32, 48, 32, 46  
.byte 46, 32, 49, 48, 48, 32, 105, 115  
.byte 32, 37, 100, 10, 0
```

The `.asciz` directive is easier to read for text strings.

# Pseudo instructions

You remember

```
li reg imm?
```

Turns out that `li` is not a RISC-V assembly command.

RISC-V cannot load a 32 bit word in one instruction. Two instructions are needed (one for the lower 16 bits of `imm`, and another one for the upper 16 bits.)

Instead it is a pseudo command: the assembler replaces each occurrence of `li` appropriately.

The RARS simulator shows pseudo instructions and the instructions that the former translate to together.