

Transcript

Introduction to Operating Systems VO

Version v1.0

Daniel Kocher

Daniel.Kocher@stud.sbg.ac.at

Salzburg, January 20, 2015

Myth Busting

Myth: I prefer programming language X because it makes it more convinient to implement my application.

Answer: Wrong!

The fundamental problem of programming is to establish functional correctness and adequate performance. Different languages provide different tools of automating the process of establishing correctness and performance. The language should be chosen based on that insight.



Figure 1: Process of establishing correctness.

Ultimate goal:

A Compiler checking correctness in such a sense that no exception is thrown while executing in any case. But this is infeasible (mathematical proof possible).

Hardware Exception: Interrupts, a mechanism to stop memory accesses in hardware.

Myth: I like garbage-collected languages like Java because they free me from memory management.

Answer: Wrong!

A garbage collector (GC) provides safe deallocation of unneeded memory but the programmer still needs to say what is unneeded, otherwise the system will run out of memory (memory leak).

General runtime complexity of a garbage collector: the size of the heap.



Figure 2: Unreachable, reachable and needed set of a program.

Multicore

Amdahl's Law: P represents program parallelism on N cores

- $S(N) = N$ if $P = 100\%$. This is ideal multicore scalability.
- In general: $S(N) = \frac{1}{(1-P) \cdot \frac{P}{N}}$

Sequential vs. Parallelized Code

The bottleneck of parallelized is the memory bus. It is limiting execution even if a problem can be perfectly parallelized without any side effects. Thus a parallelization factor of 100% is infeasible. Any shared resource at any level of an architecture creates a limitation.



Figure 3: Amdahl's law plot.

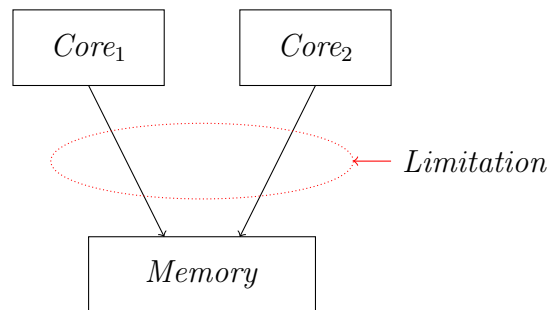


Figure 4: Bottleneck of parallelized code.

Architecture

Von Neumann Architecture

Introduced in 1945. This is a stored program computer: data = program. The Fetch-Decode-Execute cycle (see Figure 7) modifies the state of the machine.

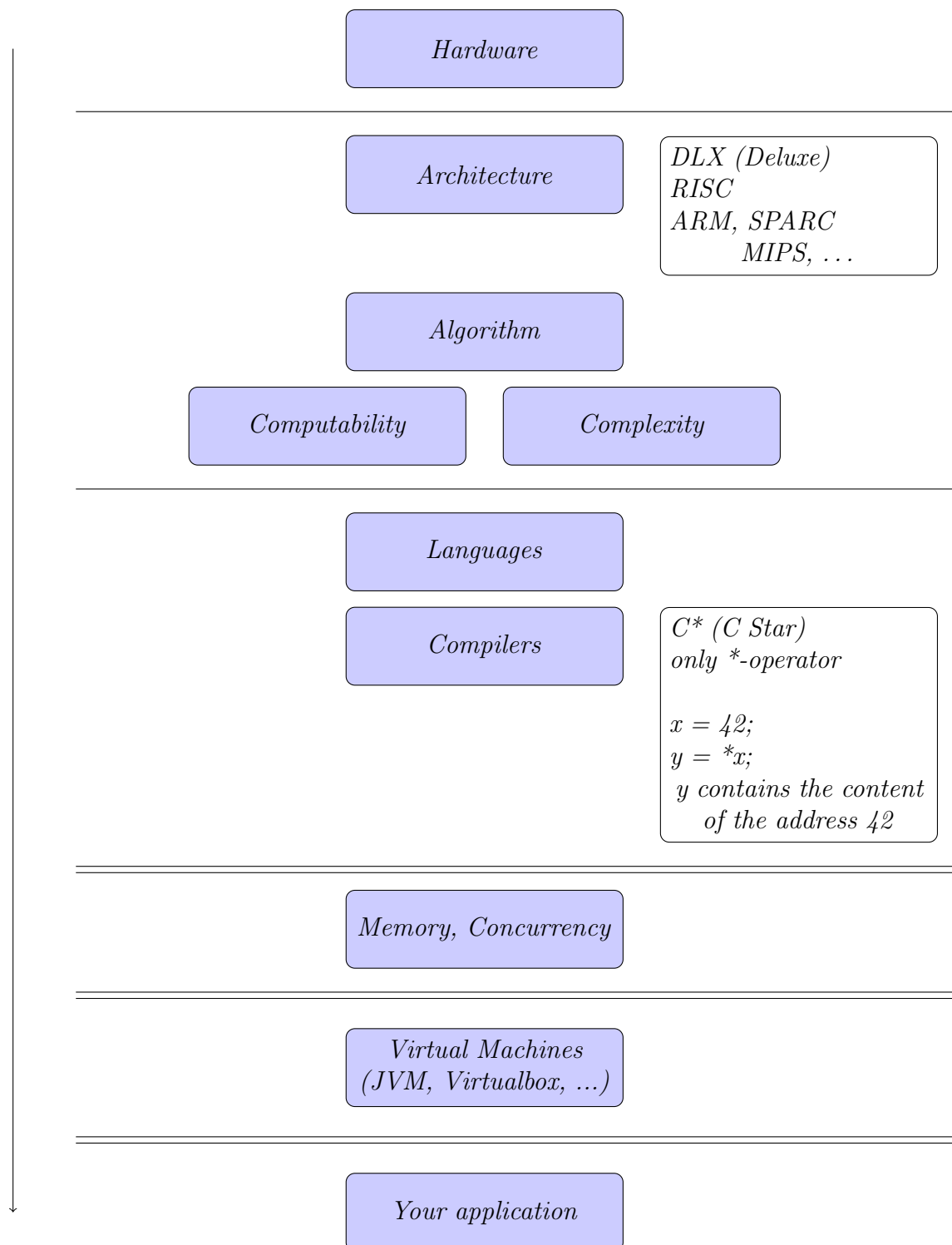


Figure 5: Architecture hierarchy (top-down.

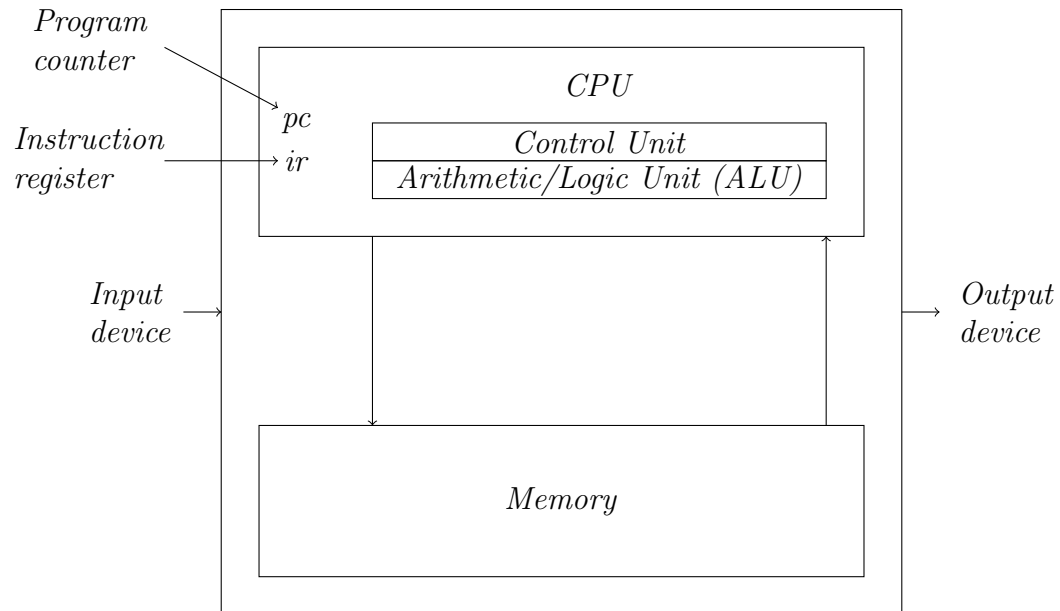


Figure 6: Von Neumann Architecture.

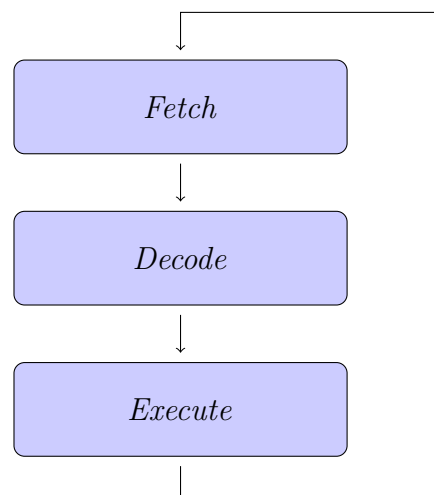


Figure 7: Fetch-Decode-Execute cycle.

DLX Machine

- Control unit:
Instruction register (ir) and program counter (pc).
- Arithmetic unit:
32x 32-bit registers; `reg[0]`, `reg[1]`, ..., `reg[31]`. `reg[0]` always contains the value 0 and `reg[31]` is the link register. Both are reserved by convention and must not be used for any other purpose.
- Memory:
n 32-bit words, byte-addressed (see Figure 8), word-aligned;
`mem[0]`, ..., `mem[n - 1]`.

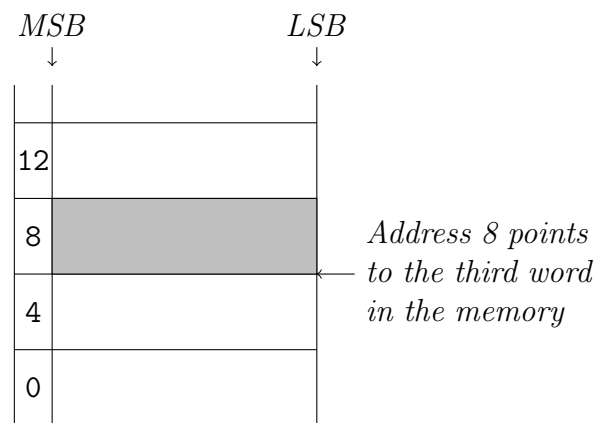


Figure 8: Visualization of a byte-addressed memory of 32-bit words.

Syntax Formats

General syntax of an instruction: `op a, b, c`.

F1

The length of `a` and `b` allow to address all 32 registers. The two's complement is used here because the implementation of arithmetics is easier (in contrast to the one's complement).

<i>op</i>	$0 \leq a \leq 31$	$0 \leq b \leq 31$	$-2^{15} \leq c \leq 2^{15} - 1$
<i>6 bits</i> 2^6 ops	<i>5 bits</i> $2^5 - 1$	<i>5 bits</i> $2^5 - 1$	<i>16 bits sign-extended to 32 bits</i>

F2

E.g. `R1 = R2 + R3`

<i>op</i>	$0 \leq a \leq 31$	$0 \leq b \leq 31$	<i>unused</i>	$0 \leq c \leq 31$
<i>6 bits</i> 2^6 ops	<i>5 bits</i> $2^5 - 1$	<i>5 bits</i> $2^5 - 1$	<i>11 bits</i>	<i>5 bits</i> $2^5 - 1$

F3

Absolute addressing in memory.

<i>op</i>	$0 \leq c \leq 2^{26} - 1$
<i>6 bits</i> 2^6 ops	<i>26 bits</i> $2^{26} - 1$

Von Neumann Bottleneck: Memory read/write operations limit the performance.

Register Instructions

Arithmetic Instructions

F1

Instruction	Semantics	Additional information
ADDI a, b, c	reg[a]:=reg[b]+c; pc:=pc+4;	Add immediate c is data (a constant)
SUBI a, b, c	reg[a]:=reg[b]-c; pc:=pc+4;	Subtract immediate c is data (a constant)
MULI a, b, c	reg[a]:=reg[b]*c; pc:=pc+4;	Multiply immediate c is data (a constant)
DIVI a, b, c	reg[a]:=reg[b]/c; pc:=pc+4;	Divide immediate c is data (a constant)
MODI a, b, c	reg[a]:=reg[b]%c; pc:=pc+4;	Modulo immediate c is data (a constant)
CMPI a, b, c	reg[a]:=reg[b]-c; pc:=pc+4;	Compare immediate c is data (a constant) reg[a]==0 if reg[b]==c reg[a]>=0 if reg[b]>=c reg[a]>0 if reg[b]>c reg[a]=<0 if reg[b]=<c reg[a]<0 if reg[b]<c reg[a]!=0 if reg[b]!=c

F2

Instruction	Semantics	Additional information
ADD a, b, c	reg[a] := reg[b] + reg[c]; pc := pc + 4;	Add Register addressing reg[c]
SUB a, b, c	reg[a] := reg[b] - reg[c]; pc := pc + 4;	Subtract Register addressing reg[c]
MUL a, b, c	reg[a] := reg[b] * reg[c]; pc := pc + 4;	Multiply Register addressing reg[c]
DIV a, b, c	reg[a] := reg[b] / reg[c]; pc := pc + 4;	Divide Register addressing reg[c]
MOD a, b, c	reg[a] := reg[b] % reg[c]; pc := pc + 4;	Modulo Register addressing reg[c]
CMP a, b, c	reg[a] := reg[b] - reg[c]; pc := pc + 4;	Compare Register addressing reg[c] reg[a] == 0 if reg[b] == reg[c] reg[a] >= 0 if reg[b] >= reg[c] reg[a] > 0 if reg[b] > reg[c] reg[a] <= 0 if reg[b] <= reg[c] reg[a] < 0 if reg[b] < reg[c] reg[a] != 0 if reg[b] != reg[c]

Register Allocation Problem

Registers have to be used in order to execute an instruction. There are 29 registers (in theory) to use. In practice, at least in this course, less registers can be used because some registers are reserved for special purposes (stack pointer, heap pointer, globals pointer, frame pointer, ...). It must be guaranteed that these will not be used and furthermore already allocated registers must not be used for an instruction.

Examples:

C Code	Instructions	Additional information
1 + 2;	ADDI 1, 0, 1 ADDI 2, 0, 2 ADD 1, 1, 2	First, naive solution
1 + 2;	ADDI 1, 0, 1 ADDI 1, 1, 2	Second, better solution
1 + 2;	ADDI 1, 0, 3	Third, best solution Constant folding
if(1 < 2)	ADDI 1, 0, 1 ADDI 2, 0, 2 CMP 1, 0, 2 BGE 1, 0, <loc>	<loc> represents the line of code in the assembly code the CPU continues with

Memory Instructions**F1**

Instruction	Semantics	Additional information
LDW a, b, c	reg[a] := mem[(reg[b]+c)/4]; pc := pc+4;	Load word (from memory) Register-relative addressing
STW a, b, c	mem[(reg[b]+c)/4] := reg[a]; pc := pc+4;	Store word (into memory) Register-relative addressing
POP a, b, c	reg[a] := mem[reg[b]/4]; reg[b] := reg[b]+c; pc := pc+4;	Pop (from stack) c: size of the popped chunk reg[b]: contains stack pointer
PSH a, b, c	reg[b] := reg[b]-c; mem[reg[b]/4] := reg[a]; pc := pc+4;	Push (onto stack) c: size of the popped chunk reg[b]: contains stack pointer

The stack grows from high to low addresses (top-down). In the POP and PSH instructions, $\text{reg}[b] := \text{reg}[b] + c$ and $\text{reg}[b] := \text{reg}[b] - c$ represent the actual removal of the element from the stack, respectively. In case of the C* language c will always be 4 because only a single type (`int`) exists in this particular language. In general c is the amount of data you want to remove from the stack.

Without register-relative addressing, which is a concrete form of indirect addressing, you could only talk about as many memory cells as your program uses. The only way to get the same expressivity as with register-relative addressing is to write self-modifying code.

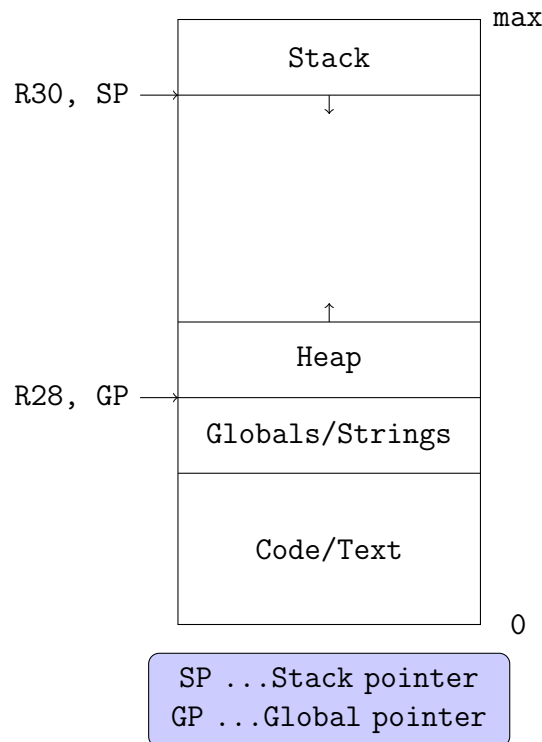


Figure 9: Memory layout when executing a program

Examples:

C Code	Instructions	Additional information
int x;	LDW 1, 28, -4	R1 := x
x = x + 1;	ADDI 2, 0, 1	R2 := 1
	ADD 1, 1, 2	R1 := R1 + R2
	STW 1, 28, -4	mem[x] = R1
*x + 1;	LDW 1, 28, -4	R1 := x
	LDW 1, 1, 0	R1 := mem[R1 + 0]
	ADDI 2, 0, 1	R2 := 1
	ADD 1, 1, 2	R1 := R1 + R2
*(x + 1);	LDW 1, 28, -4	R1 := x
	ADDI 2, 0, 4	R2 := 0 + 4, (4 because it's an address)
	ADD 1, 1, 2	R1 := R1 + R2
	LDW 1, 1, 0	R1 := mem[R1 + 0]
		*(x + i) is equivalent to x[i]

When dealing with $*(x + 1)$, the type of x has to be known as well as its size in the memory. The size of x in the memory determines the offset which is added to the address of x . In case of the C* language this value is always 4 because there is only one type and no support for composed types (**struct**). Nevertheless, the size of the type needs to be added in general.

Formula to compute $*(x + i) = x[i]: x + \text{sizeof}(x) * i$.

Control Instructions

F1

Instruction	Semantics	Additional information
BEQ a, c	if(reg[a]==0) pc:=pc+c*4; else pc:=pc+4;	Branch on equal to zero Conditional branch, pc-relative
BGE a, c	if(reg[a]>=0) pc:=pc+c*4; else pc:=pc+4;	Branch on greater than or equal to zero Conditional branch, pc-relative
BGT a, c	if(reg[a]>0) pc:=pc+c*4; else pc:=pc+4;	Branch on greater than to zero Conditional branch, pc-relative
BLE a, c	if(reg[a]<=0) pc:=pc+c*4; else pc:=pc+4;	Branch on less than or equal to zero Conditional branch, pc-relative
BLT a, c	if(reg[a]<0) pc:=pc+c*4; else pc:=pc+4;	Branch on less than to zero Conditional branch, pc-relative
BNE a, c	if(reg[a]!=0) pc:=pc+c*4; else pc:=pc+4;	Branch on not equal to zero Conditional branch, pc-relative
BR c	pc:=pc+c*4;	Branch (unconditional)
BSR c	reg[31]:=pc+4; pc:=pc+c*4;	Branch to subroutine (unconditional) The link register (R31) is saved to be able to return to the correct instruction

F2

Instruction	Semantics	Additional information
RET c	pc:=reg[c];	c = R31

F3

Instruction	Semantics	Additional information
JSR c	reg[31]:=pc+4; pc:=c;	c = R31 absolute addressing

All branches (conditional as well as unconditional) are pc-relative. The multiplication of *c* by 4 is because the words in the memory are byte-addressed. Branches are useful (especially because they operate pc-relative) because you generate so-called *relocatable code*. *Relocatable code* can be moved in memory anywhere and will not change its behavior because every address is computed relative to the program counter. The drawback of branches is the restriction of the addressable space: there may be more memory than you can use with branches. The solution is absolute addressing used by the F3 format. The F3 format allows you to address a range of 2^{26} .

Example:

C Code	Instructions	Additional information
<code>while(x<1) {</code>	<code>LDW 1, 28, -4</code>	<code>R1 := x</code>
<code><body></code>	<code>ADDI 2, 0, 1</code>	<code>R2 := 1</code>
<code>}</code>	<code>CMP 1, 1, 2</code>	<code>R1 := R1 - R2</code>
	<code>BGE 1, 0, 0</code>	c is set to zero, because at this point of time we don't know where to jump to (in general). A FixUp will later update this address.
	<code><body></code>	The body of the function
	<code>BR 0, 0, <top></code>	<code><top></code> is the address offset to branch BEFORE the first generated instruction (<code>LDW 1, 28, -4</code>).

C* is a Turing-complete language:

- arithmetics or integer
- dereferencing operator
- assignment
- while loops

Functions are not needed to satisfy Turing-completeness but reduces the size of code by preventing code duplication.

In general, every if-else-construct can be replaced by while loops, but not the other way round. However, a while loop can be replaced by a combination of functions and if-else-constructs for Turing-completeness.

Functions

Why do we use a stack and no other data structure? Because a stack guarantees proper nesting of functions, constant time access and has no spatial drawback. These facts are a result of the LIFO (Last In First Out) property of a stack.

Example:

C Code	Instructions	Additional information
<code>f(x);</code>	<code>LDW 1, 28, -4</code>	<code>R1 := x</code>
	<code>PSH 1, 30, 4</code>	Push argument on stack
	<code><body of f></code>	
	<code>BSR 0, 0, <faddr></code>	Branch to subroutine at address offset <code><faddr></code>
	<code>ADD 1, 0, 27</code>	R27 is the return register (by convention)

`<faddr>` is often resolved by a part of the compiler called *Linker*. There may be two kinds of references: symbolic references (e.g. `f`, the name of the function) and direct references (e.g. `addr(f)`, the absolute address of the function in memory).

Examples:

C Code	Instructions	Additional information
<code>int f(int x) {</code>	<code>PSH 31, 30, 4</code>	Save link register (push R31 onto stack)
<code><body of f></code>	<code>PSH 29, 30, 4</code>	Save frame pointer (push R29 onto stack)
<code>}</code>	<code>ADD 29, 0, 30</code>	Set new frame pointer to current stack pointer See Figure 10
<code>return x;</code>	<code>LDW 1, 29, 12</code>	Load argument <code>x</code> from memory. The offset is 12 here because <code>x</code> is pushed onto the stack before link and frame register.
	<code>ADD 27, 0, 1</code>	Save <code>x</code> into return register
	<code>POP 29, 30, 4</code>	Restore caller's frame (register)
	<code>POP 31, 30, 8</code>	Pop link register from stack as well as <code>x</code> . This could also be done by two <code>POP</code> instructions each popping 4 bytes. In general the size to pop is $4 + \#args * sizeof(arg_i)$.
	<code>RET 0, 0, 31</code>	Return to where function was invoked

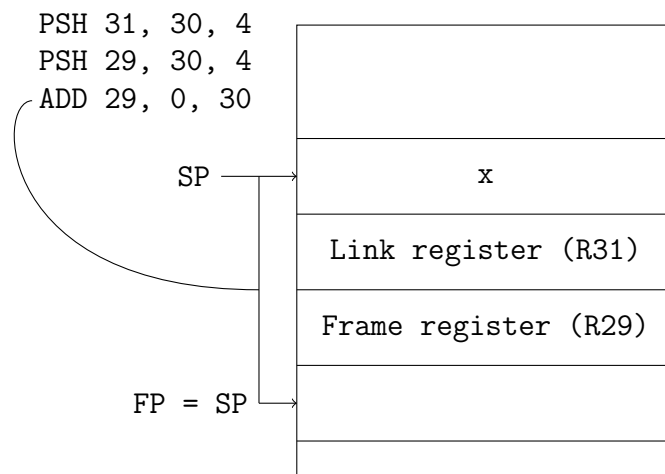


Figure 10: Stack snapshot of a function with one argument `x`.

Caller vs. Callee:

The caller of a function `f` invokes the function, e.g. `f(x);`.

The callee is the function itself, e.g. `int f(int x) { ... }`.

The callee can again be a caller if it invokes another function in its body, e.g. `int f(int x) { ... g(x); ... }`.

Declaration vs. Definition:

A declaration is inserted into the symbol table of a compiler, e.g. `int x;`. The declaration tells the compiler that a variable or function exists.

A definition is the actual value assignment, e.g. `x = 2;`. A definition can be done multiple times in a source code.

Example:

C Code	Instructions	Additional information
<code>y = x + f(x);</code>	<code>LDW 1, 28, -4</code>	Load <code>x</code>
	<code>PSH 1, 30, 4</code>	Save context of <code>f</code>
	<code>LDW 1, 28, -4</code>	Load <code>x</code>
	<code>PSH 1, 30, 4</code>	Push actual argument(s) for <code>f</code> onto stack
	<code>BSR 0, 0, <faddr></code>	Branch to <code>f</code>
	<code>...</code>	Execute function <code>f</code>
	<code>POP 1, 30, 4</code>	Restore context of <code>f</code>
	<code>ADD 2, 0, 27</code>	Store return value into <code>R2</code>
	<code>ADD 1, 1, 2</code>	Actual addition of <code>x</code> and <code>f(x)</code>
	<code>STW 1, 28, -8</code>	Store result.
		-8 because it's the second global variable

Heap

Dynamically allocated memory is located on the heap. In C there is the `malloc/free` combination to accomplish this, in Java there is the `new` keyword to allocate memory dynamically but no explicit mechanism to deallocate it. Java has a garbage collector which safely deallocates unneeded memory. However, Java provides implicit deallocation mechanisms. The first is to set a reference explicitly to null:

```
x = new A(); // dynamic allocation
...
...
x = null;
// if no other reference to x exists here
// x will be deallocated by the garbage collector
```

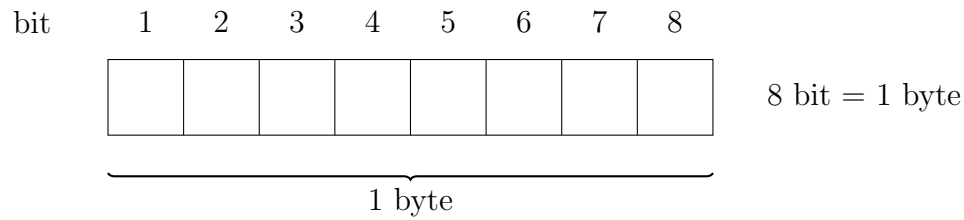
The second approach uses namespaces. When jumping out from a namespace, every namespace-local variables are unneeded if and only if there are no references from outside. If there are no such references the garbage collector will deallocate these variables:

```
{
    x = new A();
    ...
    ...
}
```

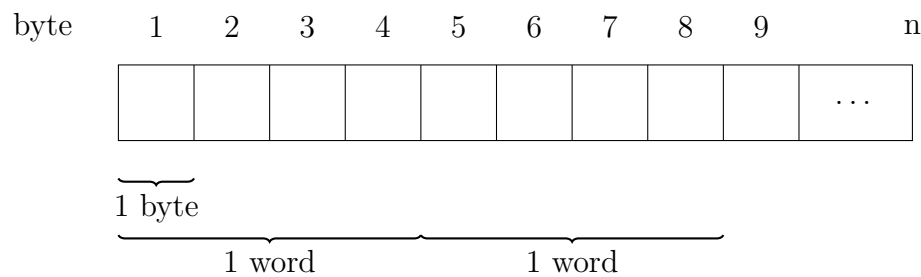
The simplest (but most unintelligent) implementation of a dynamic heap allocator is a so-called *bump pointer allocator*. A bump pointer is a pointer which only counts upwards, never downwards. In this case, there is no `free` to deallocate because it only grows upwards. However, the implementation of `malloc` is pretty simple and straigh-forward:

```
void* malloc(int s) {
    LDW 1, 29, 12
    ADD 27, 0, 26
    ADD 26, 26, 1 // R26 = Heap pointer
}
```

So every time `malloc` is invoked, these three DLX instructions are executed.



Byte-addressing and word-alignment:



Wie funktioniert der Zugriff auf den Stack?
 Wieso wird R1, R28, -4 gemacht? **Latency**.

Implementation von void* malloc(int s):

```
LDW 1, 29, 12
ADD 27, 0, 26
ADD 26, 26, 1
```

s: R1 = sum[R29 + 12]

Da man in C* keinen Zugriff auf Register hat, muss obiges Verhalten per Assembler implementiert werden.

malloc ist in der (g)libc enthalten und ist heutzutage de facto ein Betriebssystem-feature.

Implementation von int getchar():

```
RDC 0, 0, 27
```

Implementation von void putchar(int c):

```
LDW 1, 29, 12
WRC 0, 0, 1
```

Liest den nächsten Character aus einem Stream.

```
R27 = getchar();
```

Streams: `stdin`, `stdout`, `stderr`

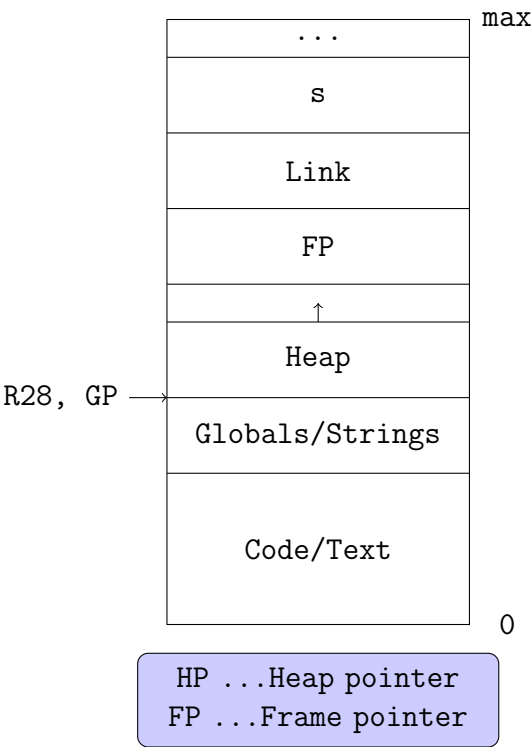


Figure 11: Memory layout revisited

F2

Instruction	Semantics	Additional information
RDC c	reg[c] = getchar(); pc:=pc+4;	
WRC c	putchar(reg[c]); pc:=pc+4;	

Why is it so hard to implement RDC and WRC on application level? Because operating systems stuff is self-referential!

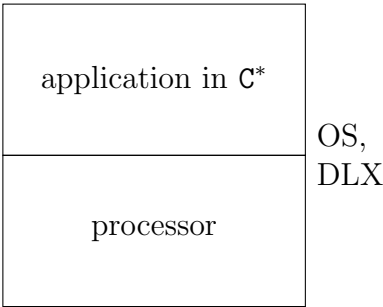


Figure 12: Application and processor

How does I/O work on a processor?

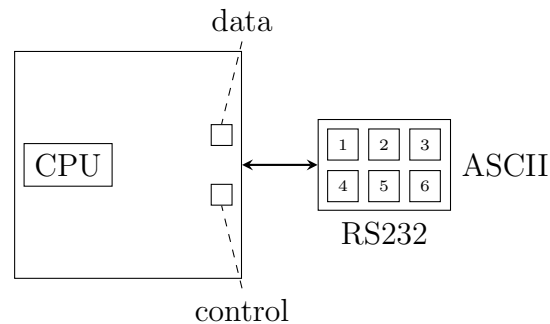


Figure 13: I/O on a processor

Computability & Complexity of algorithms

What is the minimal machine that is still universal?

- RISC: Reduced Instruction Set Computer
Has separate instructions to load/store & compute.
E.g. ARM, MIPS, SPARC, ...
- CISC: Complex Instruction Set Computer
More complex instructions which load, compute & store in a single instruction.
E.g. x86

RISC vs. CISC \leftrightarrow Compiler vs. processor.

OISC/URISC: One/Ultimate Reduced Instruction Set Computer

SUBLEQ a, b, c: Subtraction less or equal

```
mem[b] := mem[b] - mem[a];
if mem[b] ≤ 0 goto c;
else pc := pc + 4;
```

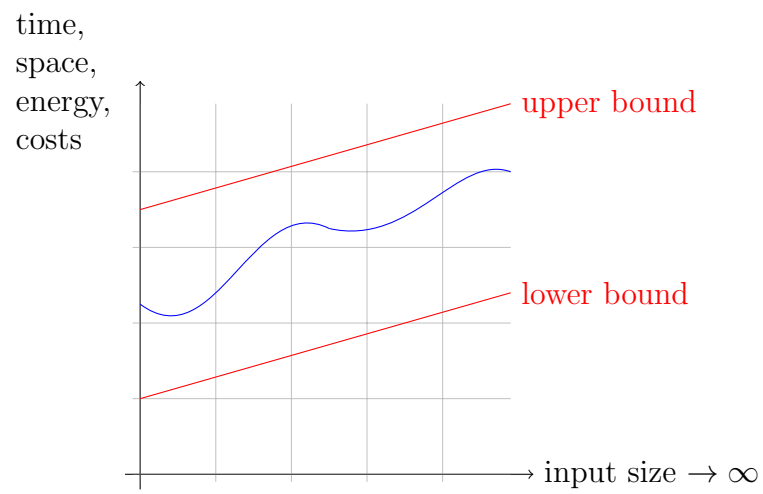


Figure 14: Asymptotic computational complexity [1]

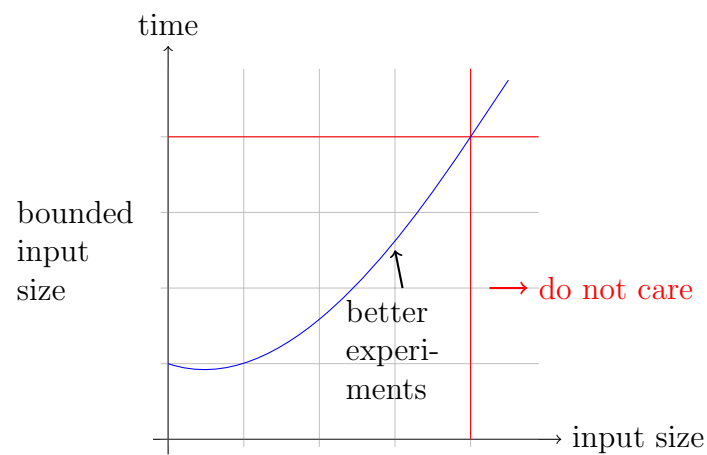


Figure 15: Asymptotic computational complexity [2] (for operating systems)

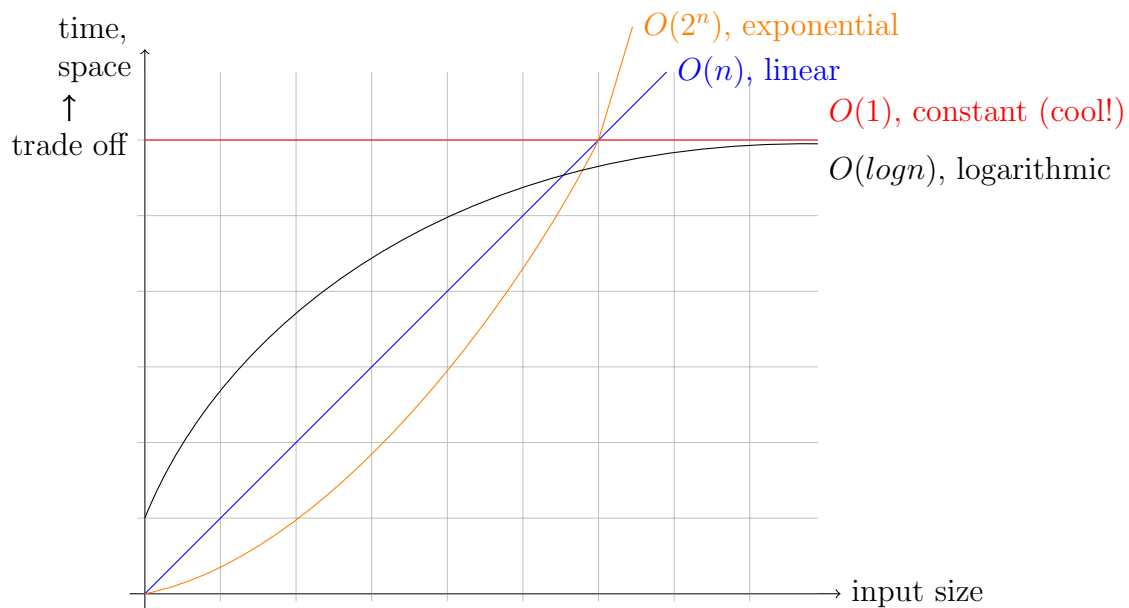


Figure 16: Big O notation (upper bound)

Big O notation (upper bound): $O(\#bits)$

$$O(\#bits) = O(\#bytes) = O(\#char) = O(\#instructions) = O(\#objects)$$

Numeric representation

unary: $\underbrace{|||||}_{16}0$, binary: $\underbrace{10000}_{16}$, octal: $\underbrace{20}_{16}$, decimal: $\underbrace{16}_{16}$, hexadecimal: $\underbrace{10}_{16}$

Exponential reduction from unary to binary representation: $1 + n = \lfloor \log_2 n \rfloor + 1$

The main difference is the denseness of the representations.

```
fib(n) {
  if n ≤ 2
    return n;
  else
    return (fib(n - 1) + fib(n - 2));
}
```

This implementation has a worst-case time complexity of $O(2^n)$. Thus, `fib(7)` will terminate correctly but something like `fib(1024)` will not terminate correctly on most systems.

There exists another implementation of the Fibonacci numbers using dynamic programming. This implementation has a logarithmic worst-case time complexity.

$$\log(n) \cdot \log_d(2) = \log_d(n)$$

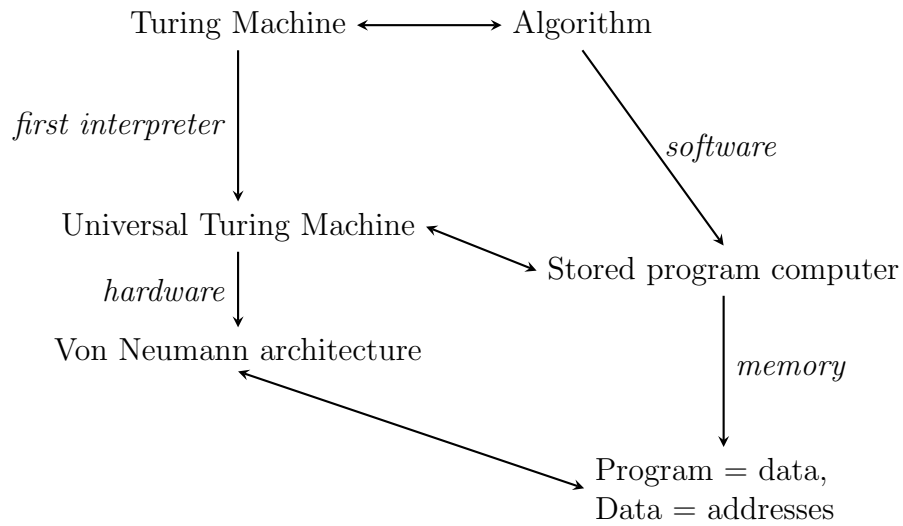
From Turing to digital:

Figure 17: From Turing to digital

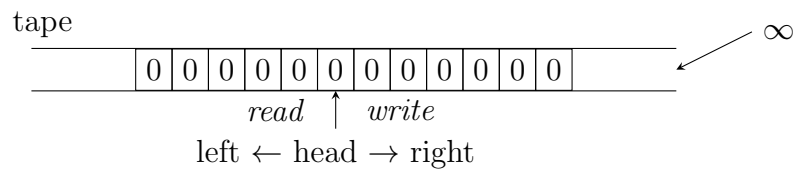
TM:

Figure 18: TM read/write

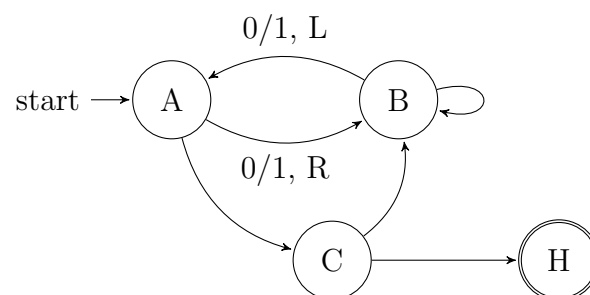
Program:

Figure 19: Busy beaver (3 states)