

Python vs Big Data

"Python?? Why Python?"

<https://www.youtube.com/watch?v=VrutixOEtOM>

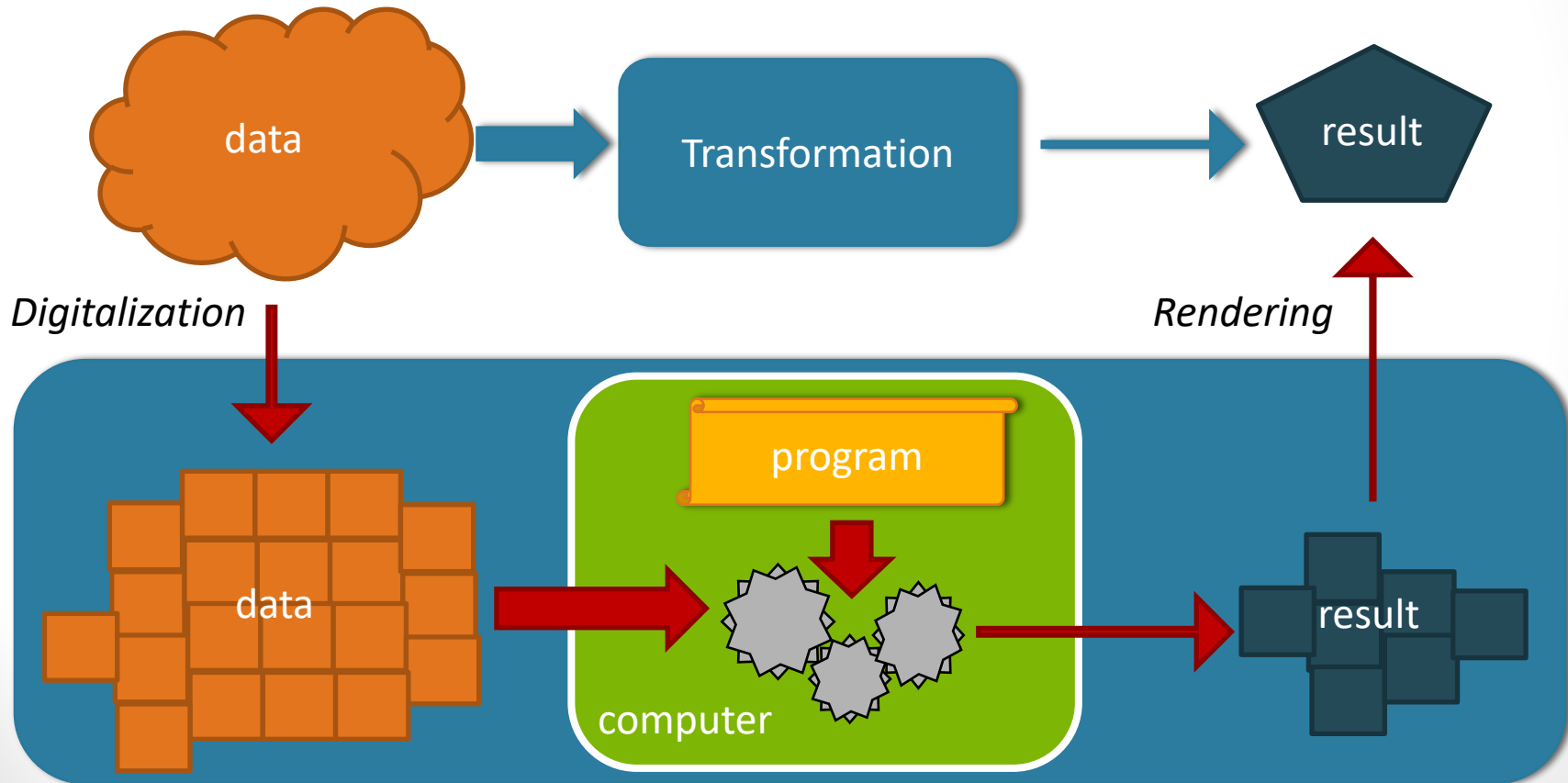
What Is Python

- Python is a **high-level, interpreted and general-purpose dynamic** programming language that focuses on **code readability**.
- The Python is widely used and have a large and active programmer **community**.
- It has a **comprehensive and large standard library** that has automatic memory management and dynamic features.
- It easily extensible by other programming language
 - <https://jakevdp.github.io/WhirlwindTourOfPython/>
 - <https://www.python.org/>
 - [https://en.wikipedia.org/wiki/Python_\(programming_language\)](https://en.wikipedia.org/wiki/Python_(programming_language))

Why Python... some step back...

It's a dirty job, but someone have to do it

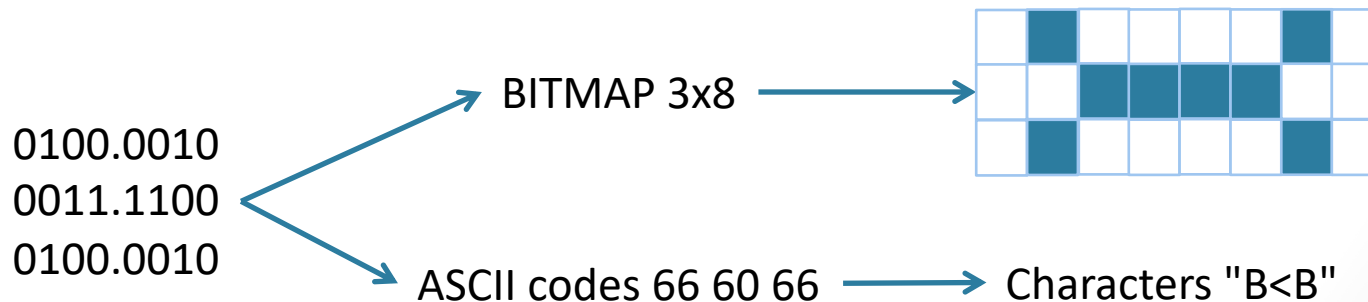
- People needs to elaborate **data** in order to extract **results**



Easy eproducible on different set of data

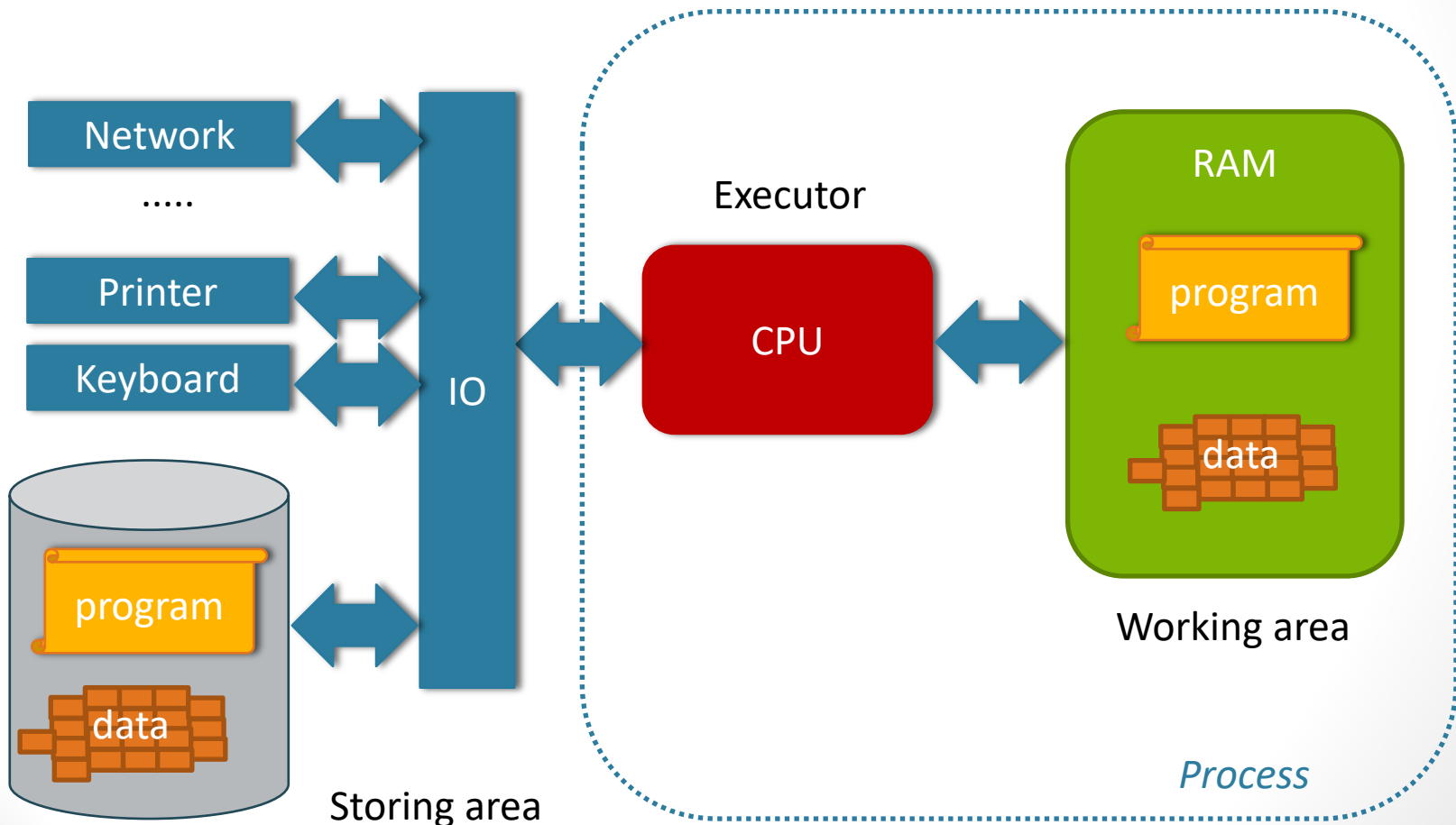
Data coding

- Digital computers can handle only binary signals: sequences of 0 and 1 (**bit** = binary digit)
- In order to transform data by digital computers, it needs to **digitalize** data, i.e. transform real samples (images, sound, etc.) into sequences of bits, packed for technological and historical reasons into group of 8 bit, called **bytes**.
- The meaning of a sequence is given by the **format** used to code and interpreter the sequence, eg. ASCII, bitmap, mp3.



Computers at hardware level

A very schematic and simplified draft of a digital computer



Coding transformations

- A classical digital computer transforms digital data by following a **program**, i.e. a **sequence of commands** that describes the transformations to be applied to data in order to reach results.
- A program can be written using various **Hi-Level** programming languages, i.e. language for humans, eg. *ADA, C, C++, Perl, Python, Java, Pascal, Basic*.
- Computers, at hardware level, understand only a very trivial set of commands, the *Assembly*, a **Low-Level** programming language, a language for CPUs.

Hi-Level languages

BASIC:

```
10 INPUT "Your name?: ", NAME$
20 PRINT "Hello "; NAME$
```

C:

```
#include <stdio.h>
char * name[100];
int main() {
    printf("Your name?: ");
    scanf("%s",name);
    printf("Hello %s\n", name);
    return 0;
}
```

Python:

```
name=input("Your name?: ")
print("Hello",name)
```

Java:

```
package stringvariables ;
import java.util.Scanner;
public class StringVariables {
    Scanner user_input = new Scanner( System.in );
    String name;
    System.out.print("Your name?");
    name = user_input.next( );
    System.out.print("Hello "+name);
}
```

Assembly

instruction in memory used by CPU

08048918

08048919

0804891b

0804891e

08048925

08048929

0804892b

0804892d

0804892e

0804892f

08048930

08048933

08048934

08048939

0804893e

08048941

08048944

08048946

08048947

08048948

instruction transliterated for humans

pushl %ebp

movl %esp, %ebp

subl \$0x4, %esp

movl \$0x0, 0xffffffffc(%ebp)

cmpl \$0x63, 0xffffffffc(%ebp)

jle 08048930

jmp 08048948

nop

nop

nop

movl 0xffffffffc(%ebp), %eax

pushl %eax

pushl \$0x8049418

call 080487c0 <printf>

addl \$0x8, %esp

incl 0xffffffffc(%ebp)

jmp 08048925

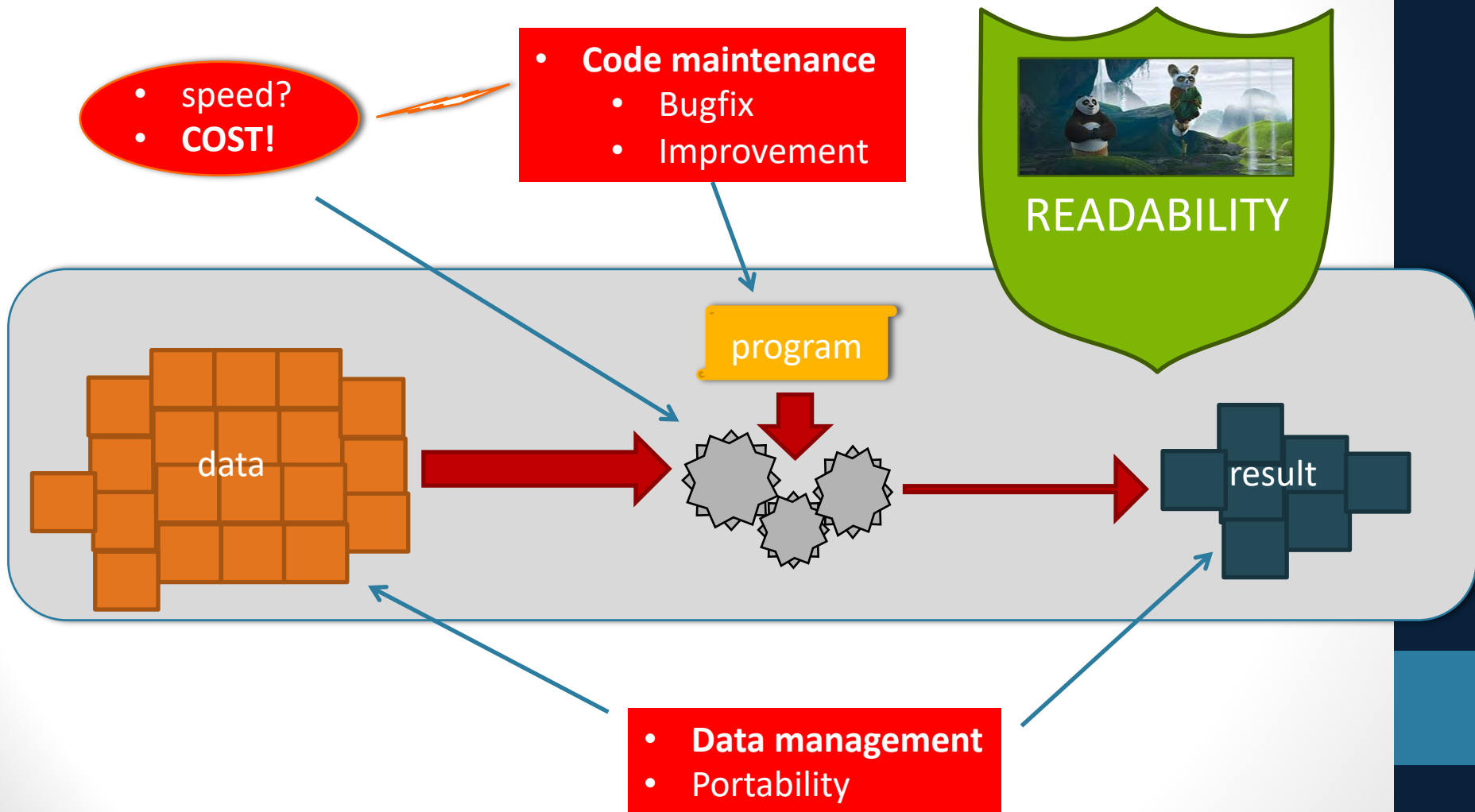
nop

nop

xorl %eax, %eax

Use computers? Start problems!

<https://www.youtube.com/watch?v=tiq6v39YliQ>



Develop Code: a job for teams

- Code should (must?) be:
 - **readable**: projects pass thorough many hands and may live, from change to change, for many years
 - **easy to develop**:
 - **easy syntax** → fast learning
 - **not error-prone**: syntax should aid to avoid errors i.e. a good programming style
 - **with a lot of *already made wheels***: a wide library collection of good functions aid to build up good code rapidly (*dont reinvent the wheel*)
 - **Cool**: a large connected community of geeks that codes with your programming language probably have already solved all of your possible problems.

A troubled kid: Speed

- Speed generally conflicts with code maintenance.

Fast codes require the full control the flow of the instructions in order to obtain hi speed of execution...

but (usually):

- is coded using a "raw" programming language (eg. C,C++) thus it result often unreadable.
- it don't use "abstractions" for implementing algorithm and managing data thus it became easy to make mistakes and bugs
- libraries are implemented from scratch in order to optimize code or remove unused part of code, thus "new code, new bugs"



"Don't run if there is not needs"



Interpreter vs Compiler

- The process of translate from HI to Low Level can be made in two way: translate the program with a **compiler** o execute the program with an **interpreter**
- Compilers:
 - take a lot of time for compile phase but the result, *the executable*, run fast on CPU.
 - Any new release of the code have to be compiled again
 - there no easy ways to run the code step by step for test (you have to use a *debugger*)
- Interpreters:
 - designed for interactive mode: easy to debug code
 - code is executed by an agent, not directly by CPU
 - easy to *port* to new kind of computer
 - Not so fast: each line have to be translated anytime is executed

Speed

C

```
char* aword=malloc(sizeof(char)*10);
scanf("%s",aword);
for (i=0;strlen(aword);i++){
    printf("%c\n",aword[i]);
}
free(aword);
```

- + fast: compiled for the running CPU
- + small binary
- + no need other piece of code
- unreadable
- memory mgmt is our duty
- easy to make mistakes on syntax

python

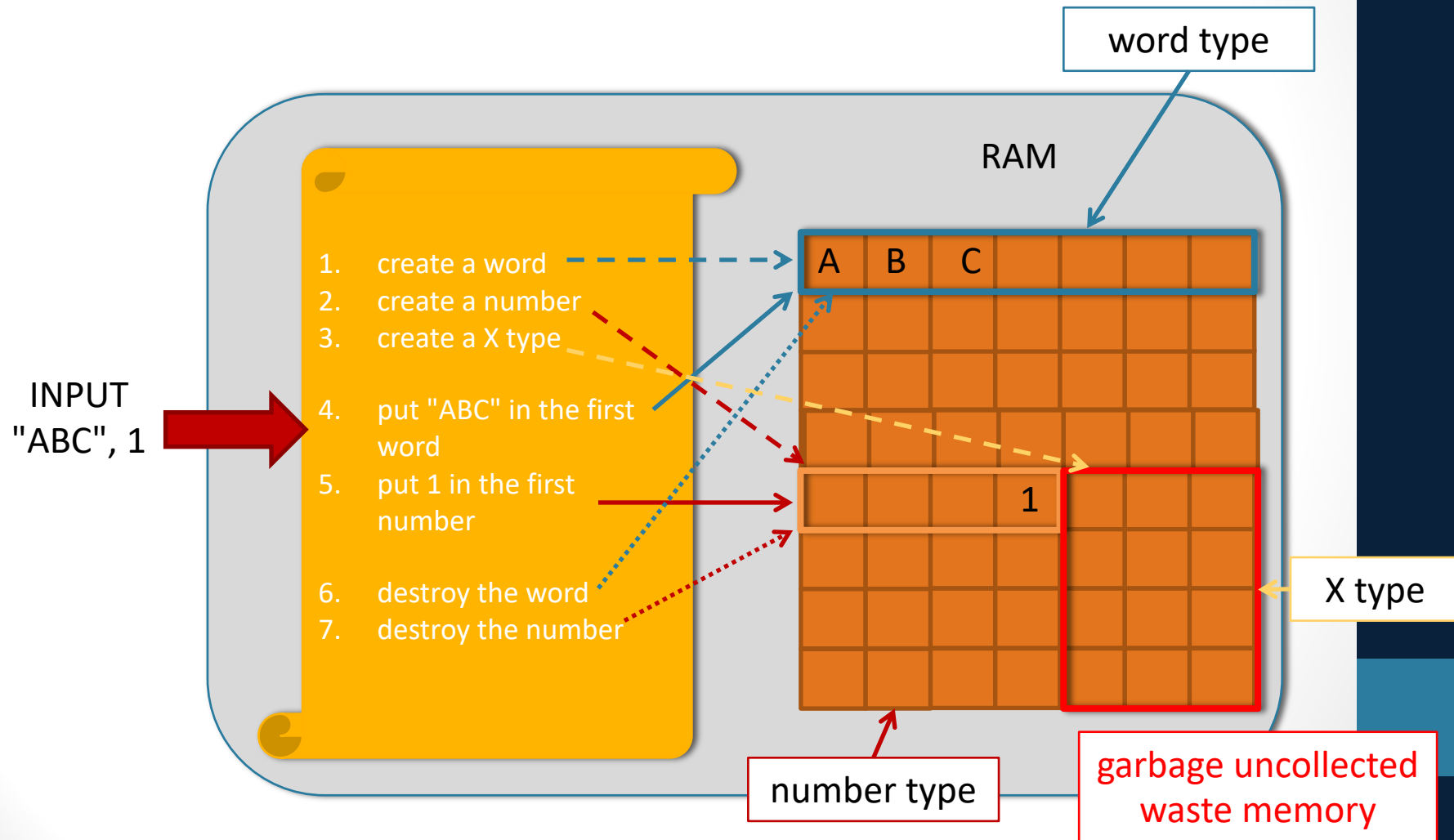
```
aword=input()
for c in aword:
    print(c)
```

- +easy to undestand
- +easy to find errors
- +memory mgmt is delagated to system
- +portable: code is not executed directly by CPU but its translated by the interpreter
- not so fast: it is interpreted
- sometime its fuzzy: managing object requires a background process that sink some cpu time in a unpredictable time
- require a software to be executed: python interpreter

Speed constrains

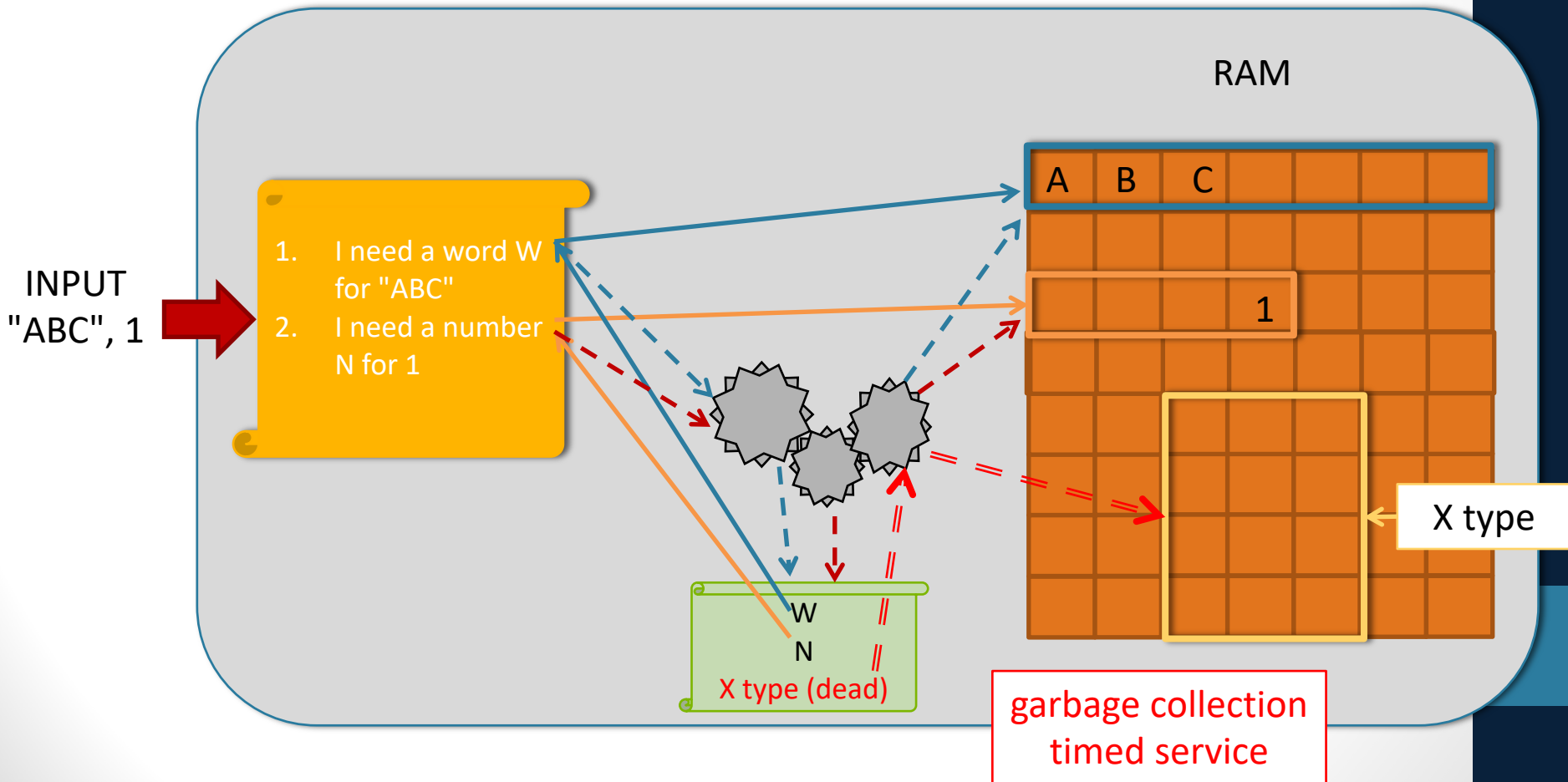
- Speed depends mainly:
 - data management:
 - how objects for data are create and, more important, destroyed.
 - how access to data is made respect to the layered cached memory
 - CPU parallelism:
 - modern CPUs are **superscalar**: can do many steps at the same time, concurrently, if the code permits it.
- On python speed depends mainly:
 - How many you code is “lined”, i.e. command are coded in a single statement that avoid for-cycle
 - How many your code delegates execution to imported library

Data management a *do-it-yourself* view (C style)



Data management

a data-as-service view (Java style)



Python spec

- General purpose language
- Focused on readability
- Interpreted
- Modular
- Dynamic
- Object-oriented
- Portable
- Extensible in C++ & C

Snakify

- Snakify is a platform for e-Learning Python 3
 - Connect to <https://snakify.org/>
 - Sign up using
 - your @unimi.it email as username (dont use your private email, if possible)
 - a password **DIFFERENT** from the one used for email
 - flag the option "**I have a teacher**"
 - put "massimo.marchi@unimi.it" in the field "**Teacher's email**"