# Performance-related aspects in the Big Data Astronomy Era: architects in software optimization

Daniele Tavagnacco - INAF-Observatory of Trieste
on behalf of EUCLID SDC-IT

# Design and Optimization



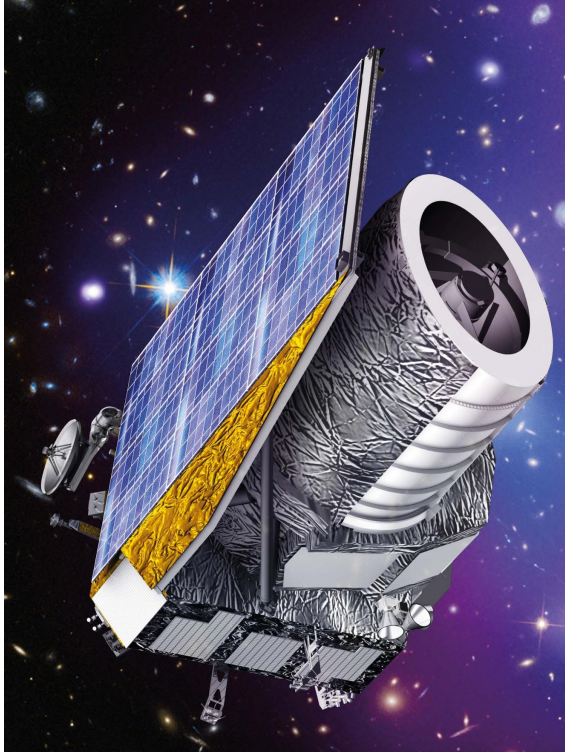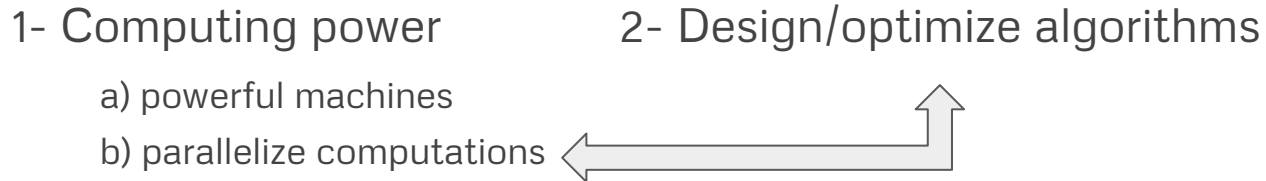image credits: web

# EUCLID mission


image credits: ESA

- ESA medium class space mission

- Universe expansion, dark energy, dark matter, gravity

- launch 2022

- 15,000 deg$^2$ survey, 6 years

- 2 instruments: VISible imager, Near Infrared SPectrograph

- ~10$^9$ observed sources, ~10$^6$ sources with spectrum

- >15 PB data

- lookback ~10 billion years (z~2)

# *Design and Big Data*

EUCLID : **15 PB** of data to be ~~processed~~ **reduced**

technical only      technical + human: understandable

*Big Data: datasets difficult to process in acceptable time frame or cost range*

1- Computing power      2- Design/optimize algorithms

    a) powerful machines
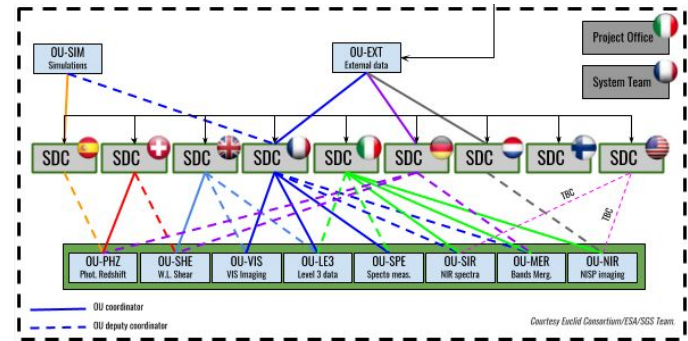
    b) parallelize computations

# SDC-IT level3 activities
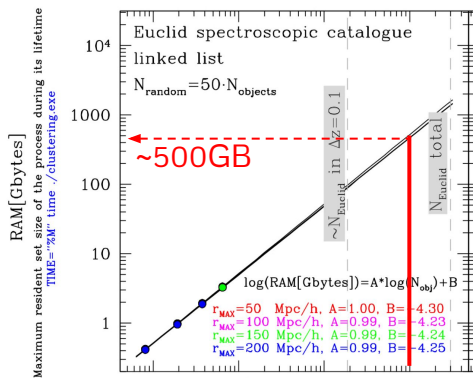
**Euclid computing infrastructure:**

- 3 levels: level1 (collect), level2 (prepare), level3 (science)
- distributed infrastructure (--> rules)
  - common environment for sw
  - minimize effort in production and testing
    (common development tools, test tools, …)
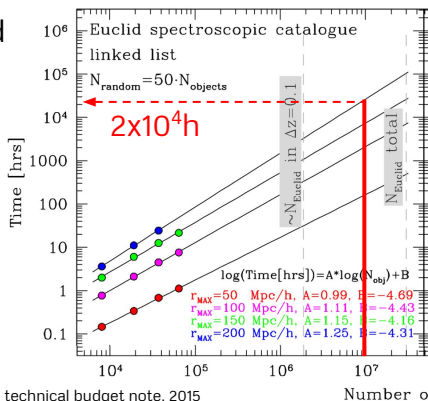
**SDC-IT supervise level3 Galaxy Clustering software:**

- integration in EC framework
  (C++, Python, 3rd party sw like *swarp*, *sextractor*, *h5py*…)
- software porting in **C++** or Python
- support for refactoring and optimization
- deployment in CI environment

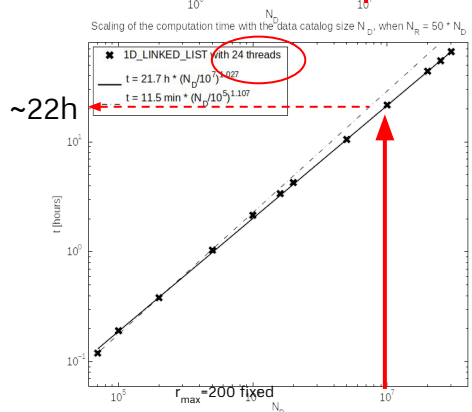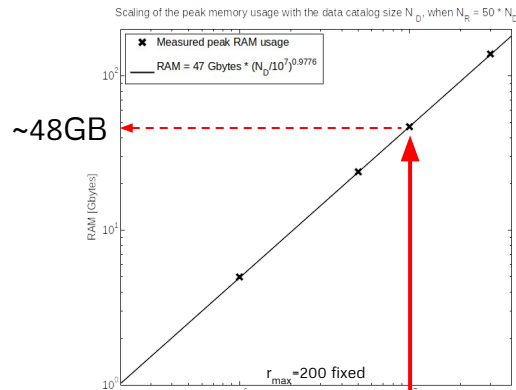# *Two Point Correlation Function GC example*



**Before...**
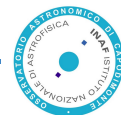on Euclid estimated
size of ~$10^7$ objects

**...after**
measured on
simulated
catalogs

~500GB

$2 \times 10^4$h

~48GB

~22h

# *Software design*

*Activity performed before writing any line of code*

Aimed at reducing:

- rigidity    - any change affects many parts of the system
- fragility    - change breaks unexpected parts of the system
- immobility  - code hard to reuse because it cannot be disentangled

Based on:

- **S**ingle Responsibility    - each entity has only one responsibility
- **O**pen/Close        - entities open for extension, closed for modifications
- **L**iskov's Substitution    - Open/Close applied to behaviour
- **I**nterface Segregation  - avoid general purpose interfaces
- **D**ependency Inversion  - decoupling high-level /low-level modules with interfaces

# *Scientific software: how good is design?*

**Software is a collection modules that:**

- operate in harmony
- have simple APIs
- hide complexity internally

**Requirements change during lifetime:**

- extend functionalities
- maintain reliability when extending
- reuse parts of the code

**The quantity of data to be reduced is increasing:**

- code scalability
- how many data are "big data"

# *Optimization within Euclid GC*

**Scientific software:**

- has special life cycle

- mainly developed by scientists

- No a priori requirements

**Refactoring the code:**

- more understandable

- cleaner and tidier

- removing redundacies and unused code

- generalize to allow reuse

- change internal structure

  (smooth flow, avoid nested conditions)

- improve performance

```cpp
 3
 4   void loopVectorA(vector<int> & vec) {
 5
 6       int c = 0;
 7
 8       // simple loop
 9       for (size_t ivec=0; ivec<vec.size(); ivec++) {
10           c += vec[ivec];
11       }
12
13   }
14
15   void loopVectorB(vector<int>  vec) {
16
17       int c = 0;
18
19       // optimized loop
20       for (vector<int>::iterator it = vec.begin(); it != vec.
21
22           c += *it;
23       }
24
25   }
```

Optimized code

```
1  loopVectorA(std::vector<int, std::allocator<int> >&):
2    rep ret
3  loopVectorB(std::vector<int, std::allocator<int> >):
4    rep ret
```

Compiler result

# *Optimization within Euclid GC*

**Scientific software:**

- has peculiar life cycle
- mainly developed by scientists
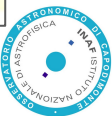- No a priori requirements

**Refactoring the code:**

- more understandable
- cleaner and tidier
- removing redundacies and unused code
- generalize to allow reuse
- change internal structure
  (smooth flow, avoid nested conditions)
- improve performance (**know the tool**)

```
 3
 4   void loopVectorA(vector<int> & vec) {
 5
 6       int c = 0;
 7
 8       // simple loop
 9       for (size_t ivec=0; ivec<vec.size(); ivec++) {
10           c += vec[ivec];
11       }
12
13   }
14
15   void loopVectorB(vector<int>  vec) {
16
17       int c = 0;
18
19       // optimized loop
20       for (vector<int>::iterator it = vec.begin(); it != vec.
21
22           c += *it;
23       }
24
25   }
```

Optimized code

```
1   loopVectorA(std::vector<int, std::allocator<int> >&):
2     rep ret
3   loopVectorB(std::vector<int, std::allocator<int> >):
4     rep ret
```

Compiler result

# *Optimization within Euclid GC*

**Scientific software:**

- has special life cycle
- mainly developed by scientists
- No a priori requirements

**Refactoring the code:**

- more understandable
- cleaner and tidier
- removing redundacies and unused code
- generalize to allow reuse
- change internal structure
  (smooth flow, avoid nested conditions)
- improve performance (**know the tool**)

```python
>>> if (os.path.exists("filename")) :
...     os.remove("filename")
... f = open("filename",'w')

# -------

>>> if not (os.path.exists("dirname")) :
...     os.makedirs("dirname")
```

Python!

```python
f = open("filename","w+")

# -------

os.makedirs("dirname")
```

# *Need to become a code expert?*

Source code

```
2
3   void Test1(int a, int b, int c, int d) {
4
5       int res = a;
6       res += b;
7       res += c;
8       res += d;
9
10  }
11
12  void Test2(int a, int b, int c, int d) {
13
14      int res = a + b + c + d;
15
16  }
17
```

```
1   Test1(int, int, int, int):
2       push rbp
3       mov rbp, rsp
4       mov DWORD PTR [rbp-20], edi
5       mov DWORD PTR [rbp-24], esi
6       mov DWORD PTR [rbp-28], edx
7       mov DWORD PTR [rbp-32], ecx
8       mov eax, DWORD PTR [rbp-20]
9       mov DWORD PTR [rbp-4], eax
10      mov eax, DWORD PTR [rbp-24]
11      add DWORD PTR [rbp-4], eax
12      mov eax, DWORD PTR [rbp-28]
13      add DWORD PTR [rbp-4], eax
14      mov eax, DWORD PTR [rbp-32]
15      add DWORD PTR [rbp-4], eax
16      pop rbp
17      ret
18  Test2(int, int, int, int):
19      push rbp
20      mov rbp, rsp
21      mov DWORD PTR [rbp-20], edi
22      mov DWORD PTR [rbp-24], esi
23      mov DWORD PTR [rbp-28], edx
24      mov DWORD PTR [rbp-32], ecx
25      mov eax, DWORD PTR [rbp-24]
26      mov edx, DWORD PTR [rbp-20]
27      add edx, eax
28      mov eax, DWORD PTR [rbp-28]
29      add edx, eax
30      mov eax, DWORD PTR [rbp-32]
31      add eax, edx
32      mov DWORD PTR [rbp-4], eax
33      pop rbp
34      ret
```
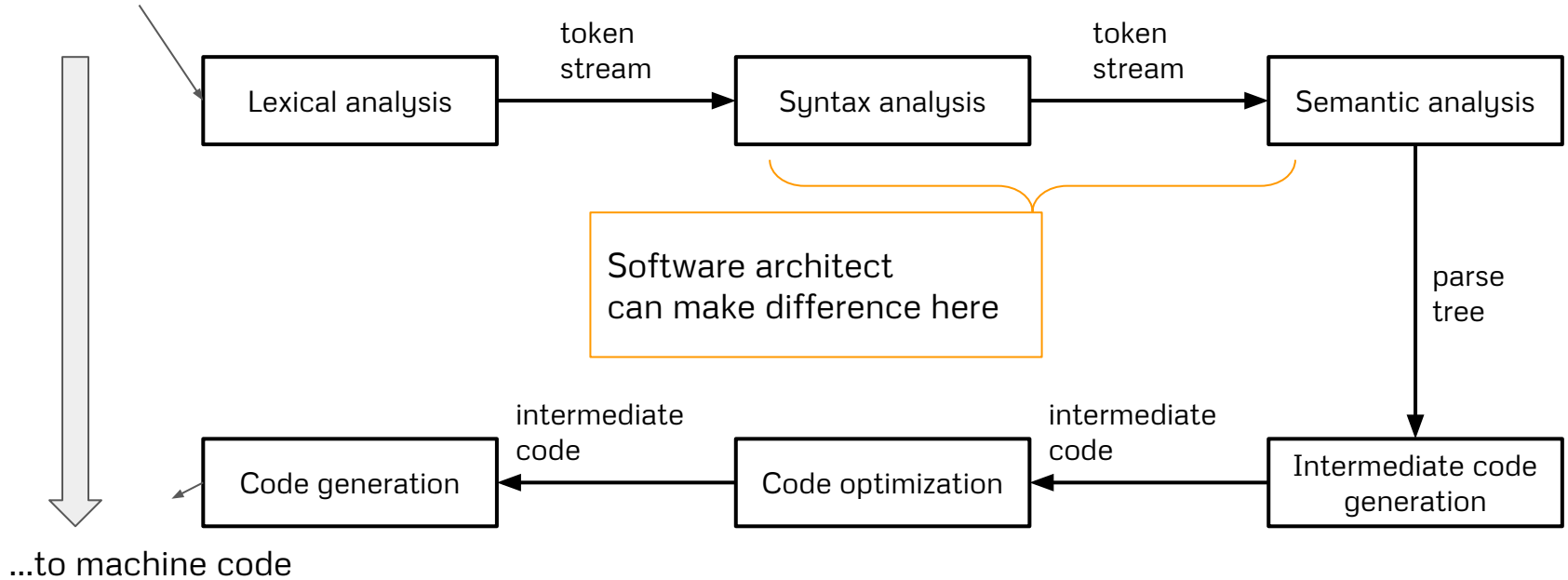Compiled code

1reg

2reg

Optimized code by compiler

```
1   Test1(int, int, int, int):
2       rep ret
3   Test2(int, int, int, int):
4       rep ret
```

# *What a compiler does?*



From source...

Lexical analysis → token stream → Syntax analysis → token stream → Semantic analysis

Software architect can make difference here

parse tree

Code generation ← intermediate code ← Code optimization ← intermediate code ← Intermediate code generation
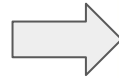
...to machine code

# *Maintaining the code*

*Revisit the code adopting new features provided by language evolution*
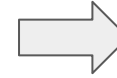
```
 3   // Sum first 1000 numbers
 4   int sum1() {
 5
 6       // set counter to 0
 7       int sum = 0;
 8
 9       // loop 1000 numbers adding to result
10       for (int i=0; i<=1000; ++i) {
11
12           sum +=i;
13       }
14
15       // return number
16       return sum;
17   }
18
19
20   // ------------------------ //
21
22
23   // Sum first 1000 numbers
24   int sum2() {
25
26       // set last number of the sum
27       int last = 1000;
28
29       // use series sum
30       int sum = last*(last + 1)/2;
31
32       // return result
33
34       return sum;
35   }
```

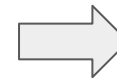code

```
 1   sum1():
 2       push rbp
 3       mov rbp, rsp
 4       mov DWORD PTR [rbp-4], 0
 5       mov DWORD PTR [rbp-8], 0
 6       jmp .L2
 7   .L3:
 8       mov eax, DWORD PTR [rbp-8]
 9       add DWORD PTR [rbp-4], eax
10       add DWORD PTR [rbp-8], 1
11   .L2:
12       cmp DWORD PTR [rbp-8], 1000
13       jle .L3
14       mov eax, DWORD PTR [rbp-4]
15       pop rbp
16       ret
17   sum2():
18       push rbp
19       mov rbp, rsp
20       mov DWORD PTR [rbp-4], 1000
21       mov eax, DWORD PTR [rbp-4]
22       add eax, 1
23       imul eax, DWORD PTR [rbp-4]
24       mov edx, eax
25       shr edx, 31
26       add eax, edx
27       sar eax
28       mov DWORD PTR [rbp-8], eax
29       mov eax, DWORD PTR [rbp-8]
30       pop rbp
31       ret
```

compiled

```
 1   sum1():
 2       xor edx, edx
 3       xor eax, eax
 4   .L3:
 5       add eax, edx
 6       add edx, 1
 7       cmp edx, 1001
 8       jne .L3
 9       rep ret
10   sum2():
11       mov eax, 500500
12       ret
```

optimized C++11

```
 1   sum1():
 2       mov eax, 500500
 3       ret
 4   sum2():
 5       mov eax, 500500
 6       ret
```

optimized C++14

# *The role of human (scientist) architect*

Design code *properly* (scalability, maintenance, extension,…)

Consider performance when *designing code* and *picking algorithms*

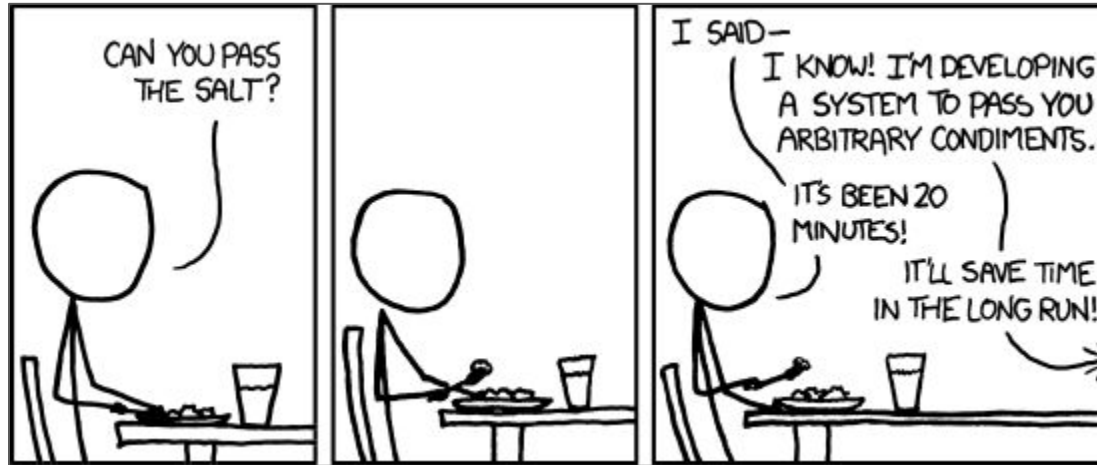Adopt *optimized* features provided by language evolution

Know the **tool**: C++ is not Fortran, Python is not IDL

When optimize don't rely only on "tips&tricks"

If you use C++, trust the compiler, it contains 45 years of improvements…

# Q&A



http://xkcd.com/974