# Mestrado em Engenharia Informática

2024/25

*João Luís Sobral*

# *Course foundations*

- **Lecturing/teaching team**
  - *Lectures: João Sobral*
  - *Lab classes: João Sobral, Rui Silva, João Barbosa, Miguel Braga*

- **Documentation & oral lectures/classes**
  - *Documentation (slides, lab guides, books, reports…): English, <u>on BB</u>*
  - *Oral presentations and work defence: optional, but usually Portuguese*

- **Expected background**
  - *Imperative Programming and Computer Architecture*
  - *Lab skills in Unix, C programming and debugging*

- **Assessment**
  - *1x written test/exam (test 20th Dec, exam 24th Jan, weight: 30%)*
  - *2x intermediate works w/ reports (submission 21th Oct & 19th Nov, weight: 2x15%)*
  - *1x final work w/ report & oral defence (end of 2024 & 6-10 Jan, weight: 40%)*
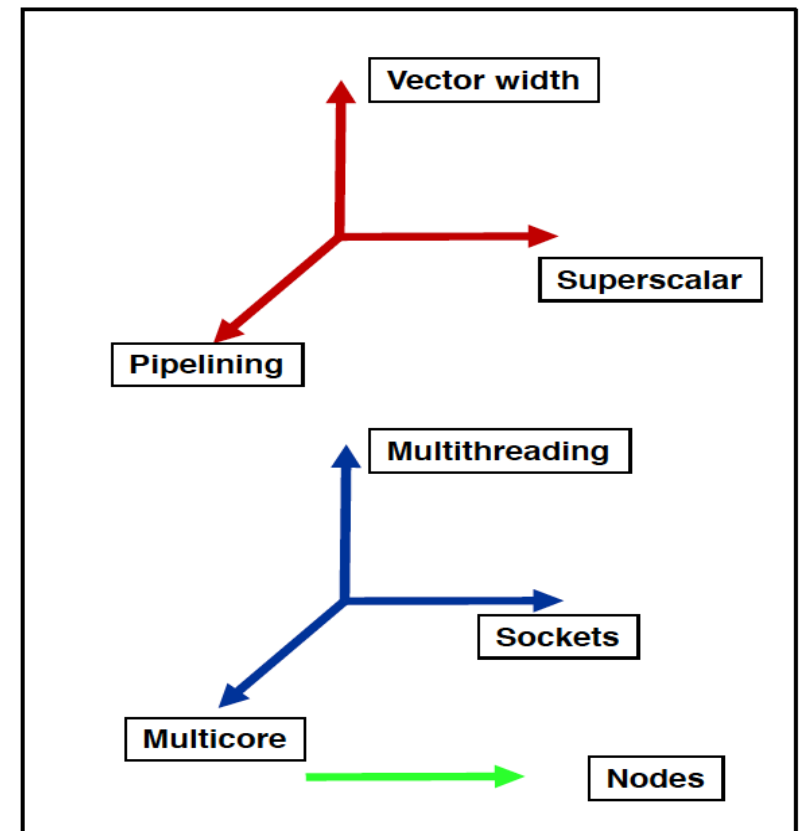
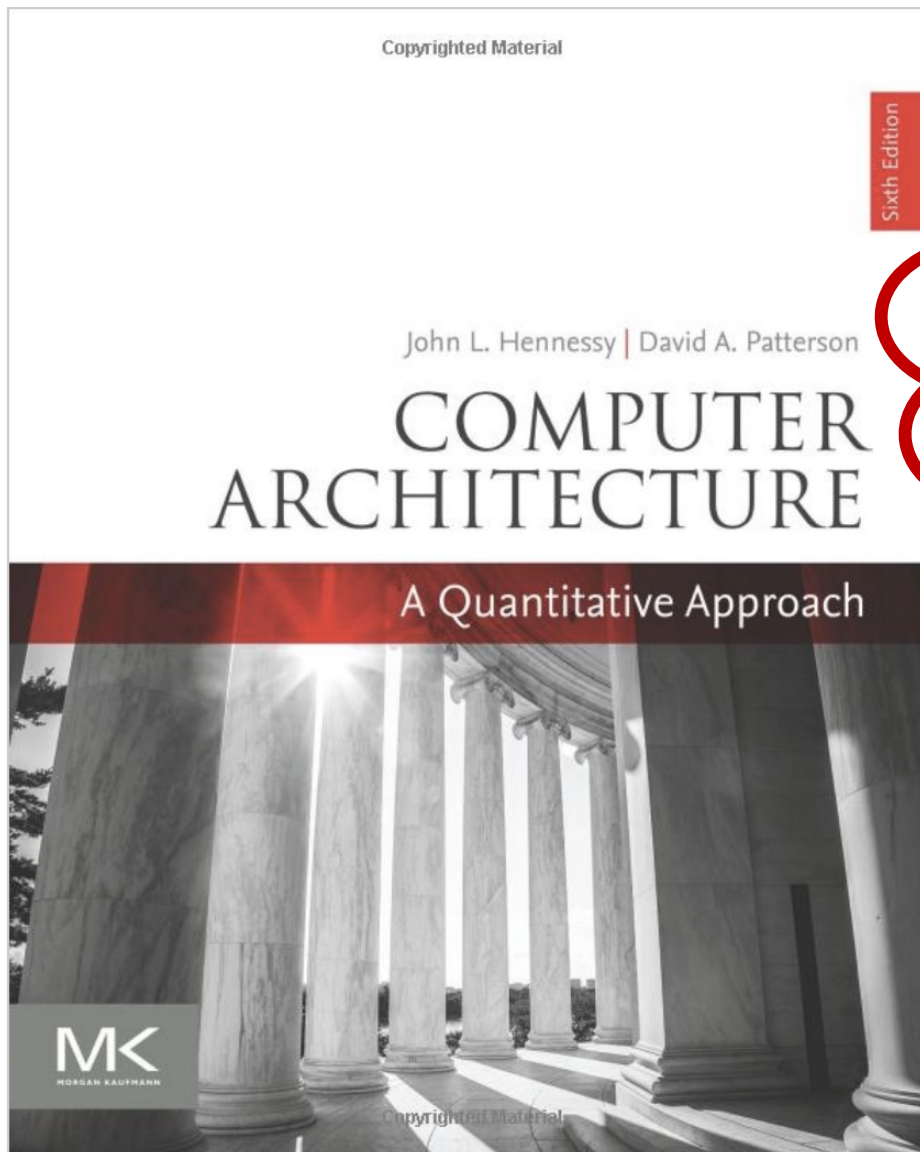# *Focus on this course: <u>performance engineering</u>*

## How:

– **plan and code efficient <u>algorithms</u> & <u>data structures</u>**

– **<u>profiling & evaluating</u>** (through measurements) **the execution efficiency**

– **understanding the organization of a computer system (its <u>architecture</u>)**

## Where in the hardware:

– **in sequential code with ILP**
  - **pipelining**
  - **superscalar w/ out-of-order exec**
  - **vector processing**

– **in a memory hierarchy**
  - **multi-level caches (+ distributed memory)**

– **w/ support for thread parallelism**
  - **multithreading in-core (SMT)**
  - **multithreading in multicore**

– **w/ computing accelerators** (GPU, …)

– **w/ support for process-based parallelism**
  - **multiprocessing in a computer cluster**

**Sixth Edition**

Copyrighted Material

John L. Hennessy | David A. Patterson

# COMPUTER ARCHITECTURE

## A Quantitative Approach

Copyrighted Material

## Table of Contents

# Key textbook (2)

THE OPENMP COMMON CORE

Making OpenMP Simple Again

Timothy G. Mattson, Yun (Helen) He, and Alice E. Koniges

Taylor & Francis Group
an **informa** business

**PARALLEL COMPUTING ARCHITECTURES AND APIs**

IoT Big Data Stream Processing

**Vivek Kale**

CRC Press
Taylor & Francis Group

A CHAPMAN & HALL BOOK

Structured Parallel Programming

*Patterns for Efficient Computation*

Michael McCool, Arch D. Robison, James Reinders

MK
MORGAN KAUFMANN

# Recommended textbook (1)

David B. Kirk
Wen-mei W. Hwu

**THIRD EDITION**

Programming Massively
Parallel Processors

A Hands-on Approach

**Contents**

**GPU & CUDA**

1. Introduction
2. Data parallel computing
3. Scalable parallel execution
4. Memory and data locality
5. Performance considerations
6. Numerical considerations
7. Parallel patterns: Convolution
8. Parallel patterns: Prefix Sum
9. Parallel patterns : Parallel Histogram Computation
10. Parallel patterns: Sparse Matrix Computation
11. Parallel patterns: Merge Sort
12. Parallel patterns: Graph Searches
13. CUDA dynamic parallelism
14. Application case study—non-Cartesian magnetic …
15. Application case study—molecular visualization …
16. Application case study—machine learning
17. Parallel programming and computational thinking
18. Programming a heterogeneous computing cluster
19. Parallel programming with OpenACC
20. More on CUDA and graphics processing computing
21. Conclusion and outlook

Appendix A. An introduction to OpenCL
Appendix B. THRUST: a productivity-oriented library for
                    CUDA

# *Background concepts from a basic Computer Systems course*

## Some notes/comments for this course:

– *some slides are borrowed from*

*and some from*

# A New Golden Age for Computer Architecture

BY JOHN L. HENNESSY AND DAVID A. PATTERSON

# Growth of computer performance
*(single-core)*

Figure 6. Growth of computer performance using integer programs (SPECintCPU).

# Current Trends in Architecture

- Cannot continue to leverage Instruction-Level Parallelism (ILP)
    - Single processor performance improvement ended in 2003

- New models for performance:
    - Data-level **parallelism** (DLP)
    - Thread-level **parallelism** (TLP)
    - Process-level **parallelism**
    - **Parallelism** with accelerators

- These require explicit restructuring of the application

# Performance optimization

- **Typical optimisation cycle**

# Performance evaluation

- **Can be performed at different levels**

**High Level Language**

*compiler*

**Assembly Language**

*assembler*

**Machine Language**

*Machine interpretation*

Control signals

```
temp = v[k];
v[k] = v[k+1];
v[k+1] = temp;
```

```
lw $15, 0($2)
lw $16, 4($2)
sw $16, 0($2)
sw $15, 4($2)
```

```
0000 1001 1100 0110
1010 1111 0101 1000
1100 0110 1010 1111
0101 1000 0000 1001
```

```
ALUOP[0:3] <=
InstReg[9:11] & MASK
```

**Applications**

**Programming languages**

**Compiler**

**ISA**

**Datapath Control**

**Functional units**

**Transistors    connections**

**Execution time (or Solutions per hour)**

**millions instructions per second: MIPS millions of FP operations per second: MFLOP/s**

**Megabytes per second**

**Cycles per second (*clock rate*)**

# Bandwidth and Latency

- Bandwidth or throughput
    - Total work done in a given time
    - 32,000-40,000x improvement for processors
    - 300-120x improvement for memory and disks

- Latency or response time
    - Time between start and completion of an event
    - 50-90x improvement for processors
    - 6-8x improvement for memory and disks

# Measuring Performance

- Typical performance metrics:
  - Response time
  - Throughput

- Speedup of X relative to Y
  - Execution time$_Y$ / Execution time$_X$

- Execution time
  - Wall clock time:  includes all system overheads
  - CPU time:  only computation time

- Benchmarks
  - Kernels (e.g. matrix multiply)
  - Toy programs (e.g. sorting)
  - Synthetic benchmarks (e.g. Dhrystone)
  - Benchmark suites (e.g. SPEC06fp, TPC-C)

# Principles of Computer Design

- The Processor Performance Equation

$$CPU\ time = CPU\ clock\ cycles\ for\ a\ program \times Clock\ cycle\ time$$

$$CPU\ time = \frac{CPU\ clock\ cycles\ for\ a\ program}{Clock\ rate}$$

$$CPI = \frac{CPU\ clock\ cycles\ for\ a\ program}{Instruction\ count}$$

$$CPU\ time = Instruction\ count \times Cycles\ per\ instruction \times Clock\ cycle\ time$$

$$\frac{Instructions}{Program} \times \frac{Clock\ cycles}{Instruction} \times \frac{Seconds}{Clock\ cycle} = \frac{Seconds}{Program} = CPU\ time$$

# Measuring the CPU Time (single-core): an example

Consider the following piece of executable x86 code:

```
    movl 10, %eax
    movl 0, %ecx
loop:
    addl %eax, %ecx
    decl %eax
    jnz loop
```

Estimate its execution time in a computer system with a
2 GHz clock frequency and a CPI=1.5.

$CPU_{time} = \text{\#Instr} * CPI * Clk_{cycle}$

$PU_{time} = 32 * 1.5 * 1/(2*10^{-9})$ sec $= 24* 10^{-9}$ s $= $ **24 ns**

# *Lei de Amdahl*

**O ganho no desempenho – *speedup* –**
obtido com a melhoria do tempo de execução de uma parte do sistema,
está limitado pela fração de tempo dessa parte.

$$Speedup_{overall} = \frac{Tempo\_exec_{antigo}}{Tempo\_exec_{novo}} = \frac{1}{f / s}$$

**f** - fração com melhoria **s**
**s** - *speedup* da fração

**Ex.1:** Se **10%** de um programa executa **90x** mais rápido

*Overall speedup* = 1.11

**Ex.2:** Se **90%** de um prog executa **90x** mais rápido

*Overall speedup* = 9.09

**Paralelismo:**
se $N_{proc} \equiv speedup$, trocar **s** por $N_{proc}$



Amdahl's Law

Parallel Portion
— 50%
— 75%
— 90%
— 95%

Speedup

Number of Processors

$$\text{Overall Speedup} = \frac{\text{Old execution time}}{\text{New execution time}}$$

$$= \frac{1}{\left( ( 1 - \text{Fraction}_{enhanced} ) + \frac{\text{Fraction}_{enhanced}}{\text{Speedup}_{enhanced}} \right)}$$

$$Speedup = \frac{1}{(1 - p) + p/N}$$

$$\text{Speedup}(N) = \frac{1}{(1-P)+\frac{P}{N}}$$

Serial part of job =
1 (100%) - Parallel part

Parallel part is divided
up by N workers

# CPU profiling using sampling

- Record stacks at a timed interval: simple and effective
  - Pros: Low (deterministic) overhead
  - Cons: Coarse accuracy, but usually sufficient
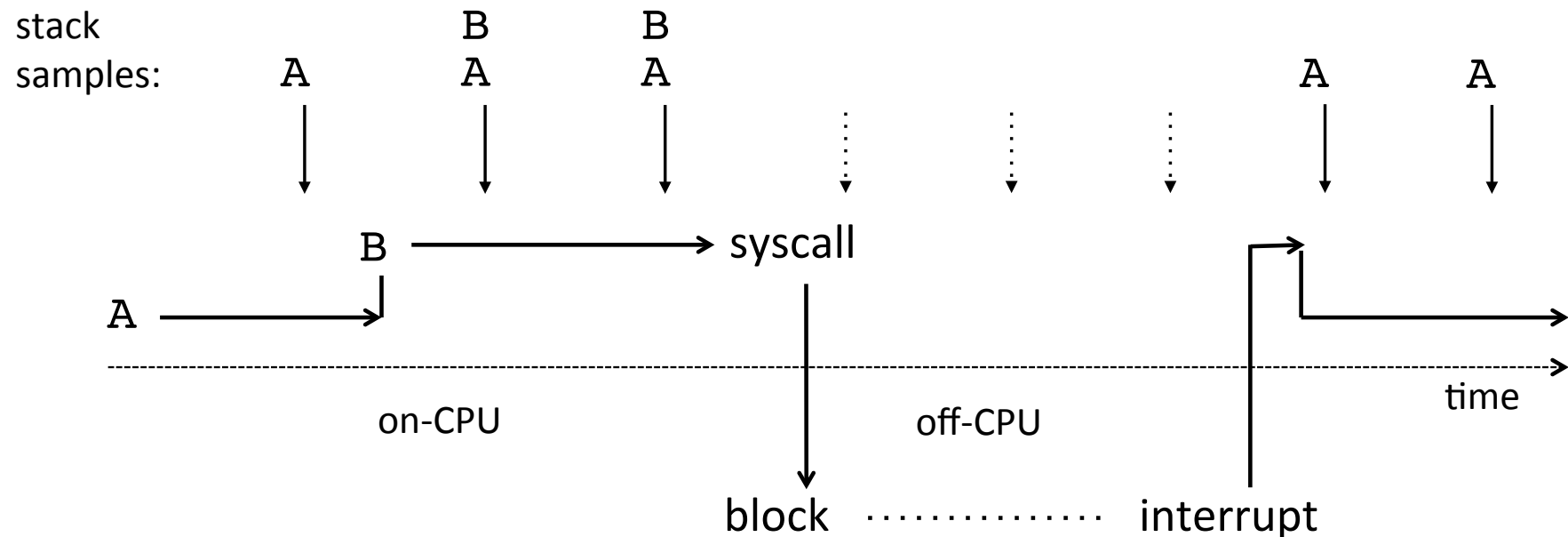
stack
samples:

```
                 B           B                                       A        A
    A            A           A
    |            |           |           ┊          ┊          ┊     |        |
    ↓            ↓           ↓           ↓          ↓          ↓     ↓        ↓

             B ────────────────→ syscall                    ┌──────→
    A ───────┘                      |                        │
                                    |                        │
  ──────────────────────────────────────────────────────────┴──────────────→
  ┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄┄→
           on-CPU                      ↓        off-CPU                    time

                                     block ┄┄┄┄┄┄┄┄┄┄ interrupt
```
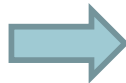
# Profiling code execution

## Techniques to measure the application time-profile (profiling)

### Polling (sampling)

- the application is periodically interrupted to collect performance data
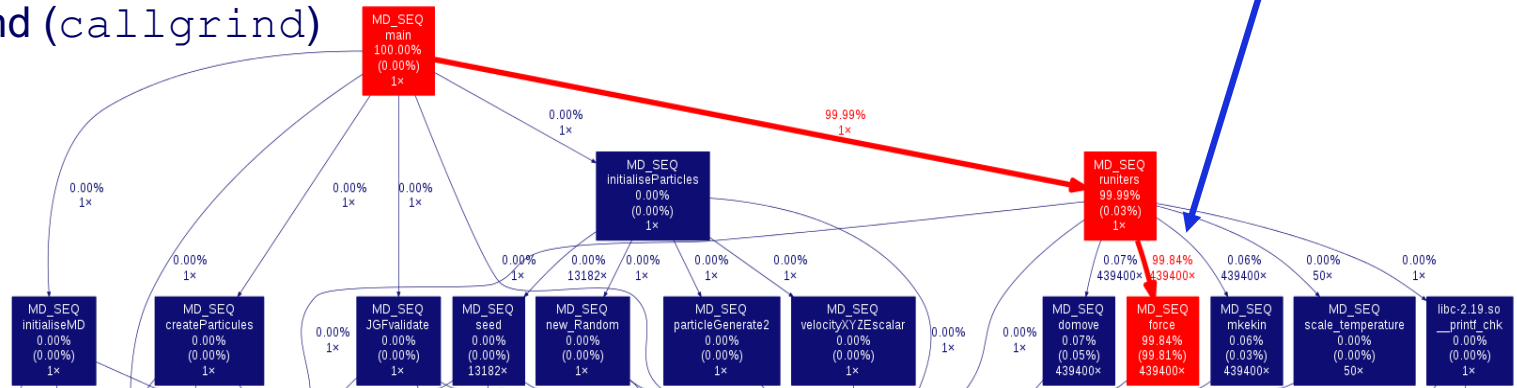
- Example: `gprof`

Flat view ⟹

| % time | cumulative seconds | self seconds | calls | self s/call | total s/call | name |
|--------|--------------------|--------------|-------|-------------|--------------|------|
| 50.00 | 2.47 | 2.47 | 2 | 1.24 | 1.24 | matSqrt |
| 24.70 | 3.69 | 1.22 | 1 | 1.22 | 1.22 | matCube |
| 24.70 | 4.91 | 1.22 | 1 | 1.22 | 1.22 | sysCube |
| 0.61 | 4.94 | 0.03 | 1 | 0.03 | 4.94 | main |
| 0.00 | 4.94 | 0.00 | 2 | 0.00 | 0.00 | vecSqrt |
| 0.00 | 4.94 | 0.00 | 1 | 0.00 | 1.24 | sysSqrt |
| 0.00 | 4.94 | 0.00 | 1 | 0.00 | 0.00 | vecCube |

**Hot-spot**

### Instrumentation

- programmer (or tools) introduce code to collect performance data on useful events
  - may produce better results but also produces more interference (e.g., overhead)
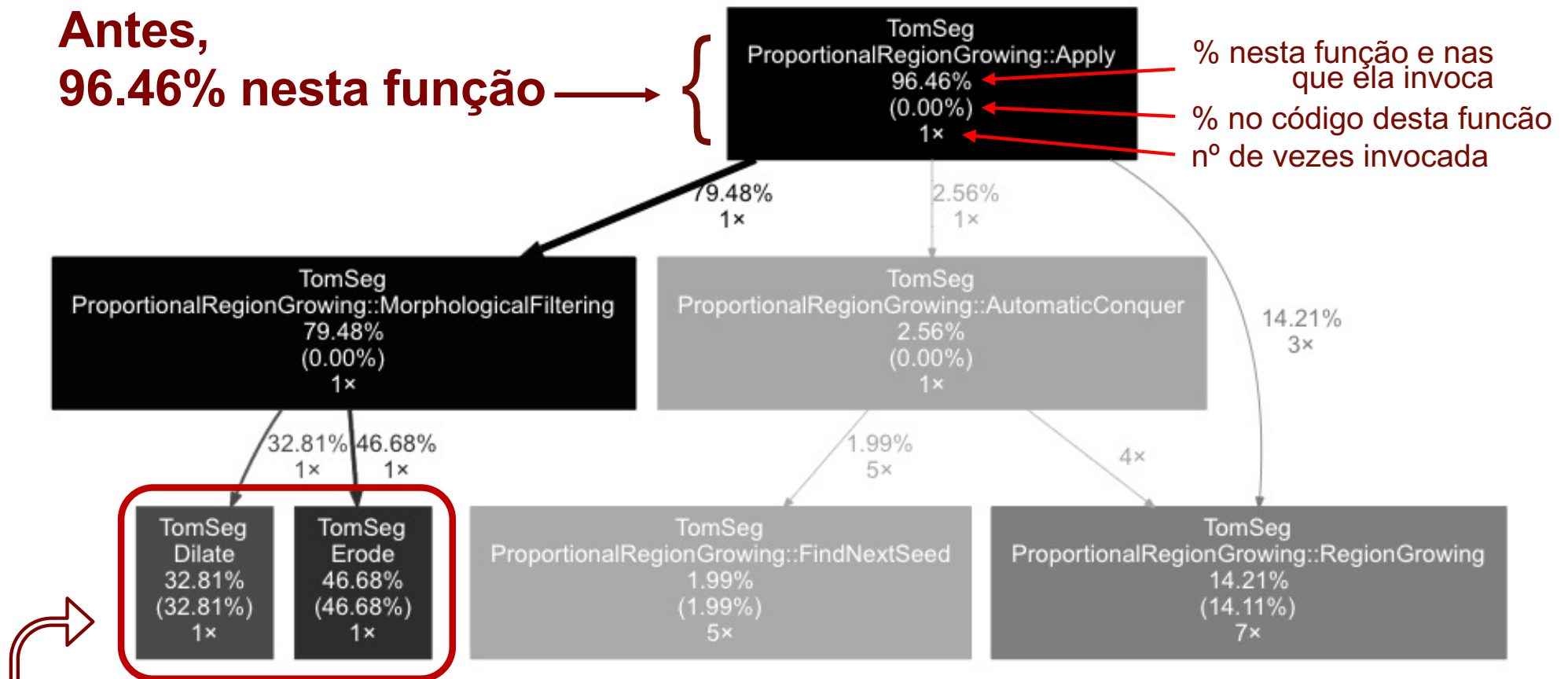- Example: Valgrind (`callgrind`)

Tree view ⟹

**Antes, 96.46% nesta função** →

% nesta função e nas que ela invoca
% no código desta função
nº de vezes invocada

Quase 80% do tempo total é gasto nestas 2 funções da filtragem morfológica!
Conclusão: é aqui que se deve investir para melhorar a *performance* global

Figure 5.7.: Call-graph of the first version of *Propor. Region Growing* (DS3)

Depois, apenas 36.38% na mesma função!

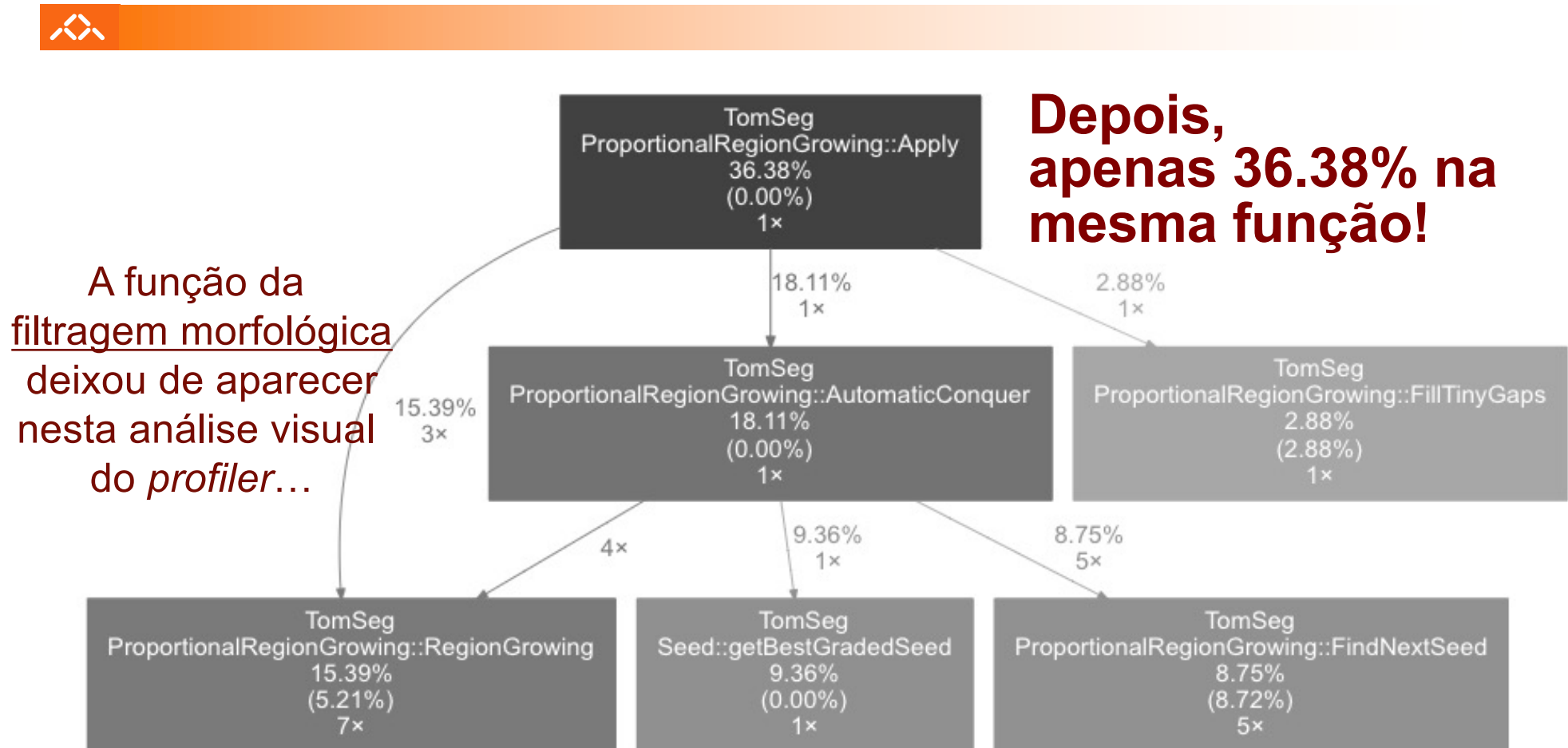A função da filtragem morfológica deixou de aparecer nesta análise visual do *profiler*…

Figure 5.9.: Call-graph from the last version of *Propor. Region Growing* (DS3)