# Reverse engineering micro-architectures

Eric Alata    eric.alata@laas.fr

INSA de Toulouse

December 15, 2023

Computer Architecture

Reverse

# Outline

## Computer Architecture

Reverse

# Outline

Computer Architecture

# Computer Architecture

*Computer architecture is the science and art of selecting and interconnecting hardware components to create a computer that meets functional, performance and cost goals.*
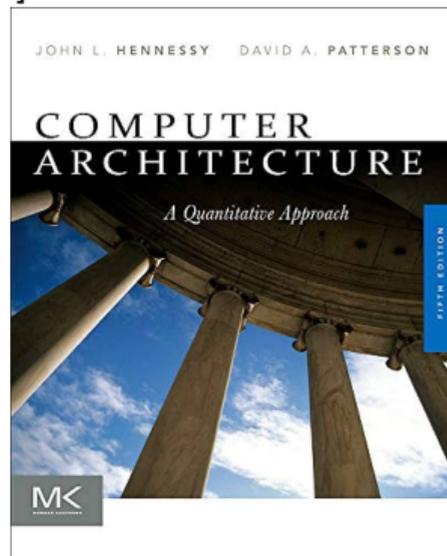
*(Mark Hill [Hil])*

# Textbook

Computer Architecture: *A Quantitative Approach* [HP12]
John Hennessy and David Patterson

- Appendix C
- Chapter 3

## ISA

- sub r0, r1, r2    →    $r0 \leftarrow r1 - r2$
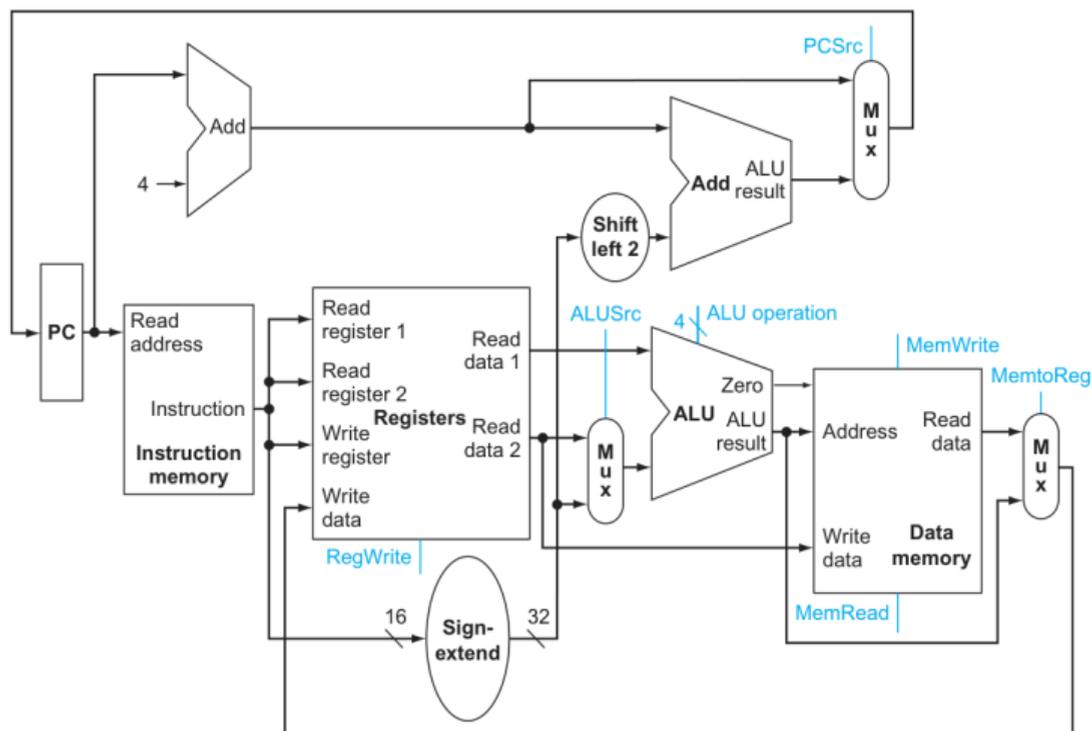
# Simple MIPS architecture [PH13]



**FIGURE 4.11   The simple datapath for the core MIPS architecture combines the elements required by different instruction classes.** The components come from Figures 4.6, 4.9, and 4.10. This datapath can execute the basic instructions (load-store word, ALU operations, and branches) in a single clock cycle. Just one additional multiplexor is needed to integrate branches. The support for jumps will be added later.
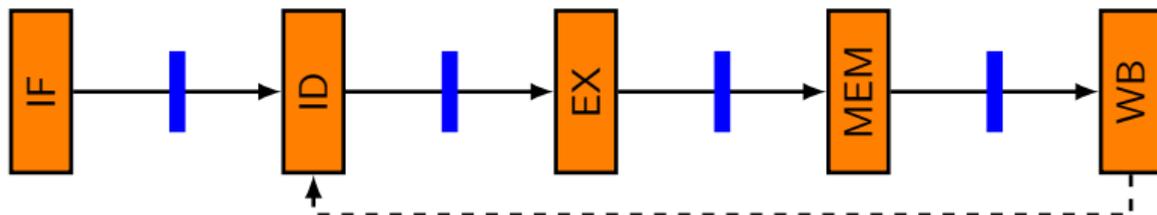
# Pipeline

Classical 5 stages

- Instruction fetch – IF
- Instruction decode – ID
- Execute – EX
- Memory access – MEM
- Write back – WB

Pipeline Hazards

- Data hazards
- Structural hazards
- Control hazards
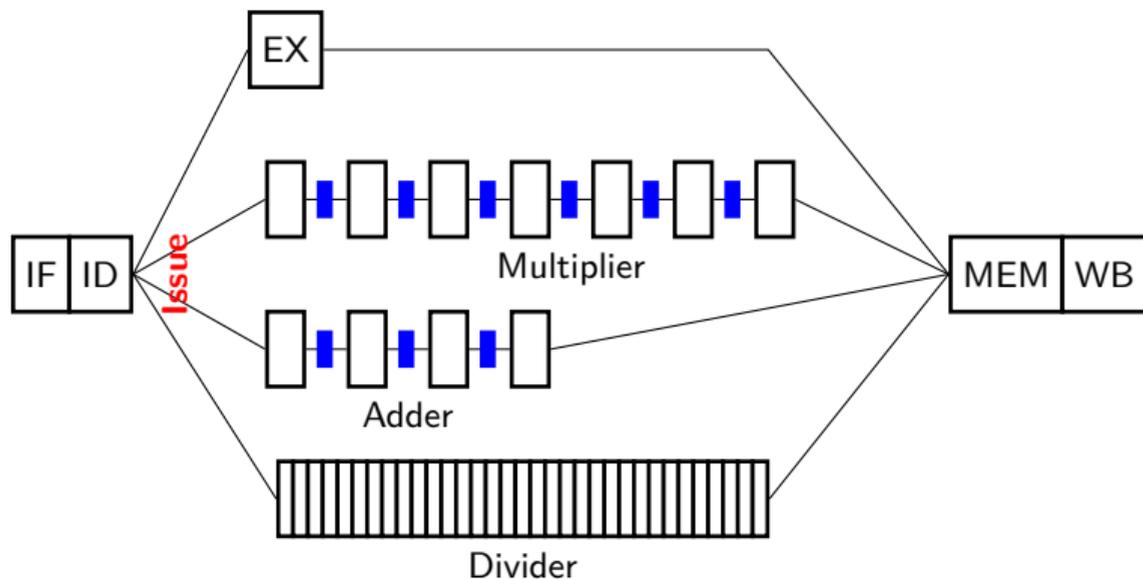- ⇒ Forwarding
- ⇒ Bubble/stall

## Example 1

- Registers: `ri = i + 10`
- `add r0, r1, r2`
- `add r3, r4, r5`
- `add r6, r7, r8`
- Structural hazard (on registers ID and WB)

# Example 2

- Registers: `ri = i + 10`
- `add r0, r1, r2`
- `add r3, r0, r0`

# Multicycle pipeline



| Latency 3<br>(w/ forwarding) | add r0, r1, r2 | IF | ID | A1 | A2 | A3 | A4 | MEM | WB | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | sub r3, r0, r5 | | IF | ID | ID | ID | ID | A1 | A2 | A3 | A4 | MEM | WB |

| Initiation interval 1 | add r0, r1, r2 | IF | ID | A1 | A2 | A3 | A4 | MEM | WB | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | add r3, r4, r5 | | IF | ID | A1 | A2 | A3 | A4 | MEM | WB | |
| | add r6, r7, r8 | | | IF | ID | A1 | A2 | A3 | A4 | MEM | WB |

# Example 3

- Registers: ri = i + 10
- add r0, r1, r2
- div r3, r4, r5    div by 0?
- add r6, r7, r8
- add r9, r10, r11
- add r12, r13, r14

# Multicyle pipeline

- Multiple paths with various latencies
- Multiple instructions may write to the register file at the same time!
- ⇒ Multiple ports & stall
- Write back stage may proceed instruction out-of-order
- ⇒ Check registers usage of issued instructions & stall if needed
  (set in ID, cleared in WB)
- Write-after-write, Read-after-write & write-after-read...

## Data Hazards

- Read-after-write (RAW) – True dependence
  $i_1$: add r0, r1, r2
  $i_2$: sub r3, r0, r4
  Swap instructions $\Rightarrow$ sub uses wrong value for r0
  Stall pipeline ($+$ Bypassing)

- Write-after-read (WAR) – Anti dependence
  $i_1$: add r0, r1, r2
  $i_2$: sub r1, r3, r4
  Swap instructions $\Rightarrow$ add uses wrong value for r1
  Not in in-order pipeline ; possible in out-of-order execution

- Write-after-write (WAW) – Output dependence
  $i_1$: add r0, r1, r2
  $i_2$: sub r0, r3, r4
  Swap instructions $\Rightarrow$ wrong value stored in r0

- Read-after-read?

## Structural Hazards

- Resource used twice in one cycle
- For instance, memory (instruction & data)
- Stall one instruction, usually the younger

## Control Hazards

- Must fetch next instruction before the branch outcome is known
- Strategy: always *not-taken* and flush if needed
- Strategy: branch prediction strategy
- Branch Target Buffer

# Advanced strategies

- Multiple issue: two or more instructions per stage at once
- Superscalar architectures

# In-order pipeline to Out-of-order execution

- Passing functionality to reduce structural hazard
- Use of an instruction buffer
- Dispatch: allocate en entry in the instruction buffer
- Issue: send an instruction from the instruction buffer to the execution unit

# Dynamic scheduling

- How many cycles for the execution of the following program?
  $i_1$: div r0, r1, r2
  $i_2$: mul r3, r0, r5
  $i_3$: add r6, r7, r8
- Instruction $i_3$ does not depend on $i_1$ and $i_2$
- It could be executed in parallel...
- Multiple strategies (mainly scoreboard and Tomasulo)
- Focus on Tomasulo algorithm [Tom67]

# Tomasulo algorithm

Key elements

- Register renaming
- Instruction queue
- Load/store buffer
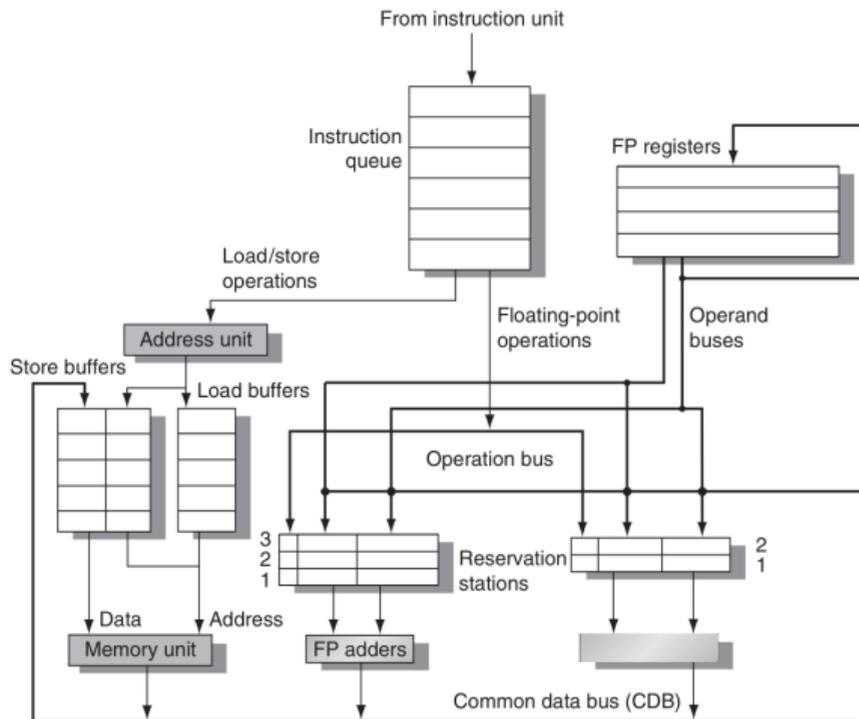- Reservation stations
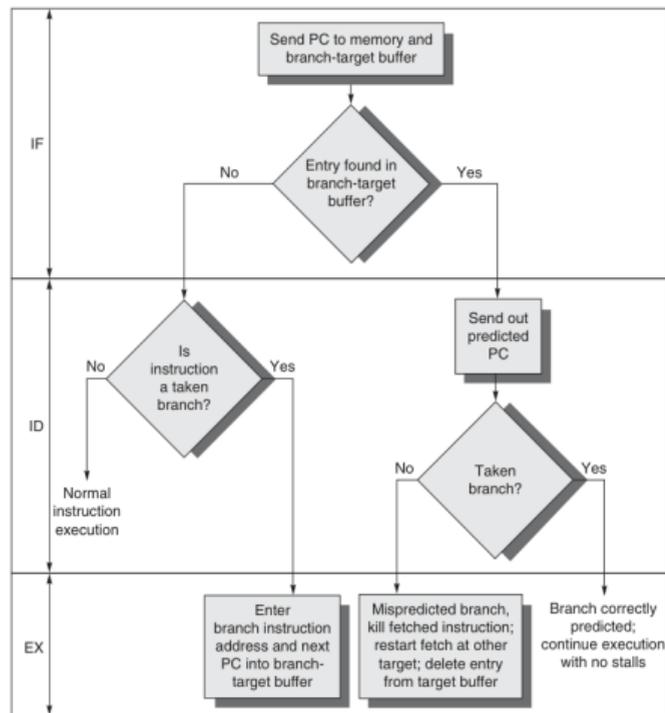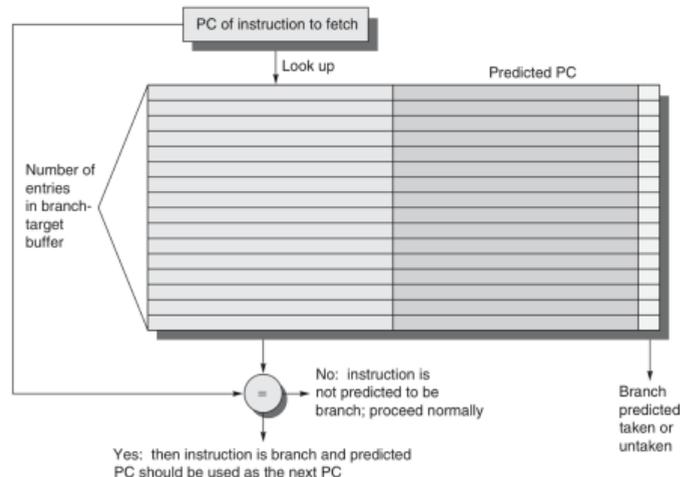- Multiple units
- Common data bus



**Figure 3.6 The basic structure of a MIPS floating-point unit using Tomasulo's algorithm.** Instructions are sent

# Speculation [McF93]

- An instruction is (control-)dependant on some set of branches
  ```
  if (a)   ins1; // this instruction depends on a
  if (b)   ins2; // this instruction depends on b but not on a
  ```
- Branch prediction predicts the execution as accurate as possible
- Branch prediction could be wrong (mis-prediction)
- Speculation means to execution instructions according to predictions
- Speculative execution recovery   wrong prediction → roll-back execution
- **Remaining traces of misprediction?**
- Useful if the condition is computer after the target address (classical 5 stages?)

# Speculation [McF93]

- Types of dynamic branch predictors
  - 1-bit Branch-Prediction Buffer
  - 2-bit Branch-Prediction Buffer
  - Correlating Branch Prediction Buffer
  - Tournament Branch Predictor
  - **Branch Target Buffer**
  - Return Address Predictors
  - Integrated Instruction Fetch Units



**Figure 3.21  A branch-target buffer.** The PC of the instruction being fetched is matched against a set of instruction



**3.22  The steps involved in handling an instruction with a branch-target buffer.**

# Tomasulo algorithm

Major issue: precise exception & speculation

- Tomasulo: in-order issue, out-of-order execution and out-of-order completion
- Consider the following example
  $i_1$: div r0, r1, r2
  $i_2$: mul r3, r0, r5
  $i_3$: add r6, r7, r8
- Suppose an Hardware interruption occurs 10 cycles after $i_1$ is issued
- Instructions $i_1$ and $i_2$ have not been committed
- Instruction $i_3$ has already been committed!
- What about specution? mis-prediction as an exception?
- ⇒ Use a reorder-buffer out-of-order execution & in-order commit

# Outline

Computer Architecture

Reverse

# Outline

Reverse

# Reverse engineering micro-architecture

- Why *to reverse*?
  - Understand the behavior of the component
  - Produce more efficient software
  - Identify undocumented instructions and behavior
  - Identify weakness (vulnerabilities?)
- What *to reverse*?
  - Almost all elements of the micro-architecture
  - Is the element available in the micro-architecture?
  - Where is this element in the datapath?
  - What are the main dimensions of this element?
- How *to reverse*?
  - Difficulty: we don't have a direct acces to the micro-architecture!
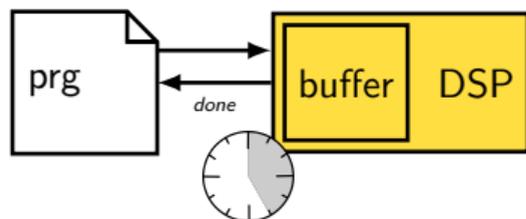  - We only have acces to ISA
  - Need for a methodology

# Example

- Target: a DSP used to process data in batch
- The internal memory is finite: batch buffer
- A data can be send to the DSP using the `add <data>` instruction
- Data are internally processed using the `compute` instruction
  And the buffer is flushed
- On `add`, if the buffer is full then `compute` is executed
- Processing time of a data is constant
  $\Delta(process\ a\ data) = 1$   cc         (cc: clock cycle)
- Processing time of the buffer also depends on an overhead
  $\Delta(\texttt{compute}) = \#data\ in\ buffer\ \times \Delta(process\ a\ data) + 128$   cc
- Processing time of `add`?
- Element under reverse: batch buffer → its size (number of data slots)

<p style="text-align:center">How to guess the size of the batch buffer?</p>

# Example

- $\Delta(\text{process a data}) = 1$   cc
- $\Delta(\texttt{compute}) = \#\text{data in buffer} + 128$   cc
- $\Delta(\texttt{add}) = \begin{cases} \#\text{buffer size} + 1 + 128 & \text{cc,} \quad \text{if the buffer is full} \\ 1 & \text{cc,} \qquad\qquad\qquad \text{otherwise} \end{cases}$
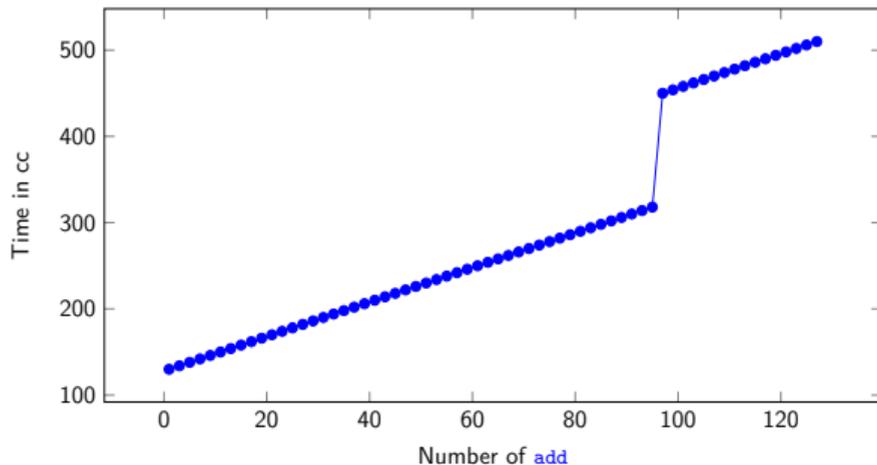


```
1   times = []
2   for n in range(1, 128, 2):
3       execute(["compute"])
4       time = execute((["add"] * n) + ["compute"])
5       times.append((n, time))
6   print("n t")
7   for (n, t) in times: print(n, t)
```
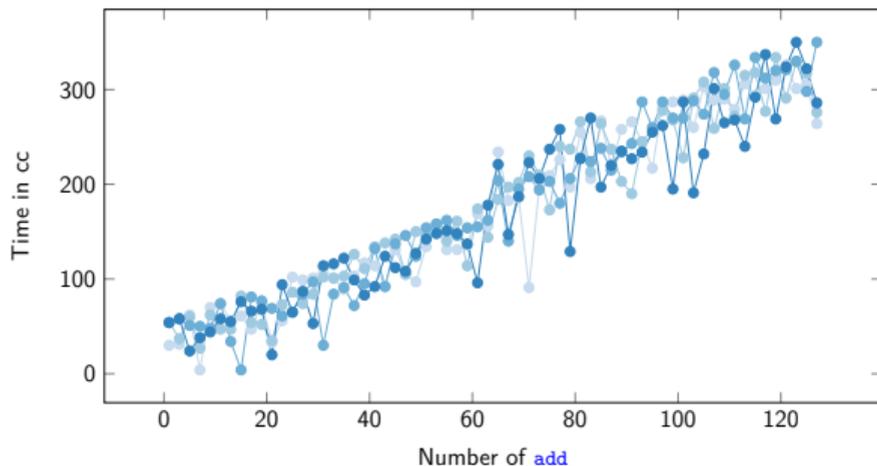
# Example

- $\Delta(\textit{process a data}) = 1$   cc
- $\Delta(\texttt{compute}) = \#\textit{data in buffer} + 128$   cc
- $\Delta(\texttt{add}) = \begin{cases} \#\textit{buffer size} + 1 + 128 & \text{cc,} \quad \text{if the buffer is full} \\ 1 & \text{cc,} \qquad\qquad\qquad \text{otherwise} \end{cases}$



What is the size of the buffer?

# Example

- Real hardware components are more complex
- Branch prediction unit, caches, etc.
- These different units can speed up or slow down treatments
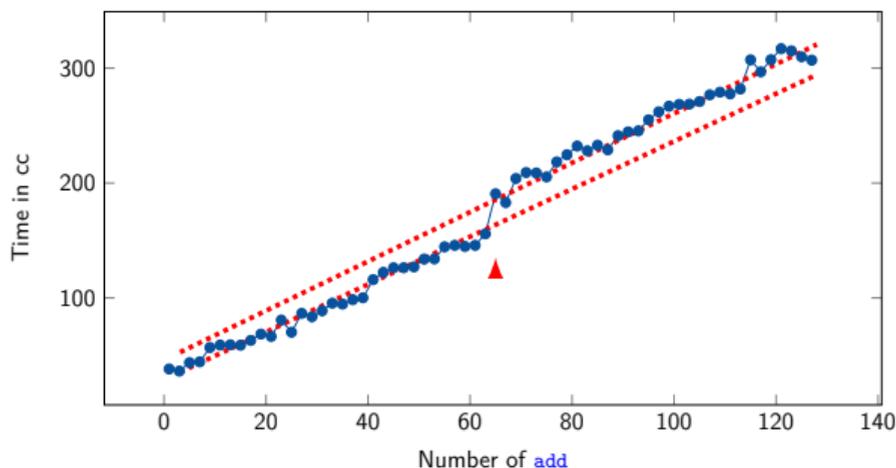- Makes timing very noisy



What is the size of the buffer?

# Example

```
1    times = [(n, []) for n in range(1, 128, 2)]
2    for k in range(nb_experiences):
3      for i in range(len(times)):
4        execute(["compute"])
5        time = execute((["add"] * times[i][0]) + ["compute"])
6        times[i][1].append(time)
7    print("n " + " ".join(["t%d" % i for i in range(nb_experiences)]) + " tm")
8    for (n, ts) in times: print(n, " ".join(map(str, ts)), np.mean(ts))
```



What is the size of the buffer?

# Specific topics

- Superscalar architecture
- Scoreboard

# References I

📄 Cs/ece 752: Advanced computer architecture i,
https://pages.cs.wisc.edu/~markhill/cs752/Fall1999, Accessed:
2023-06-27.

📄 John L. Hennessy and David A. Patterson, Computer architecture: A quantitative
approach, 5 ed., Morgan Kaufmann, Amsterdam, 2012.

📄 Scott McFarling, Combining branch predictors, Tech. report, 1993.

📄 David A. Patterson and John L. Hennessy, Computer organization and design, fifth
edition: The hardware/software interface, 5th ed., Morgan Kaufmann Publishers
Inc., San Francisco, CA, USA, 2013.

📄 R. M. Tomasulo, An efficient algorithm for exploiting multiple arithmetic units,
IBM Journal of Research and Development **11** (1967), no. 1, 25–33.