

Discussion 9 : Instruction level parallelism in CA-I

Contents

- Motivations
- Techniques exploiting ILP
 - Pipelining
 - Static Multiple Issue: VLIW
 - Dynamic Multiple Issue: Superscalar
 - Speculative Execution
 - Out-of-Order Execution
- Exercises

Why ILP

Recall: The iron law of performance

$$\frac{\text{Time}}{\text{Program}} = \frac{\text{Instructions}}{\text{Program}} \times \frac{\text{Cycles}}{\text{Instructions}} \times \frac{\text{Time}}{\text{Cycle}}$$

Boosting performance without tuning the frequency or rewriting the program.

ILP is to executing multiple instructions in parallel by:

- having multiple datapathes running simultaneously
- utilizing datapath components that are free
- reducing stall incurred data dependencies and controls

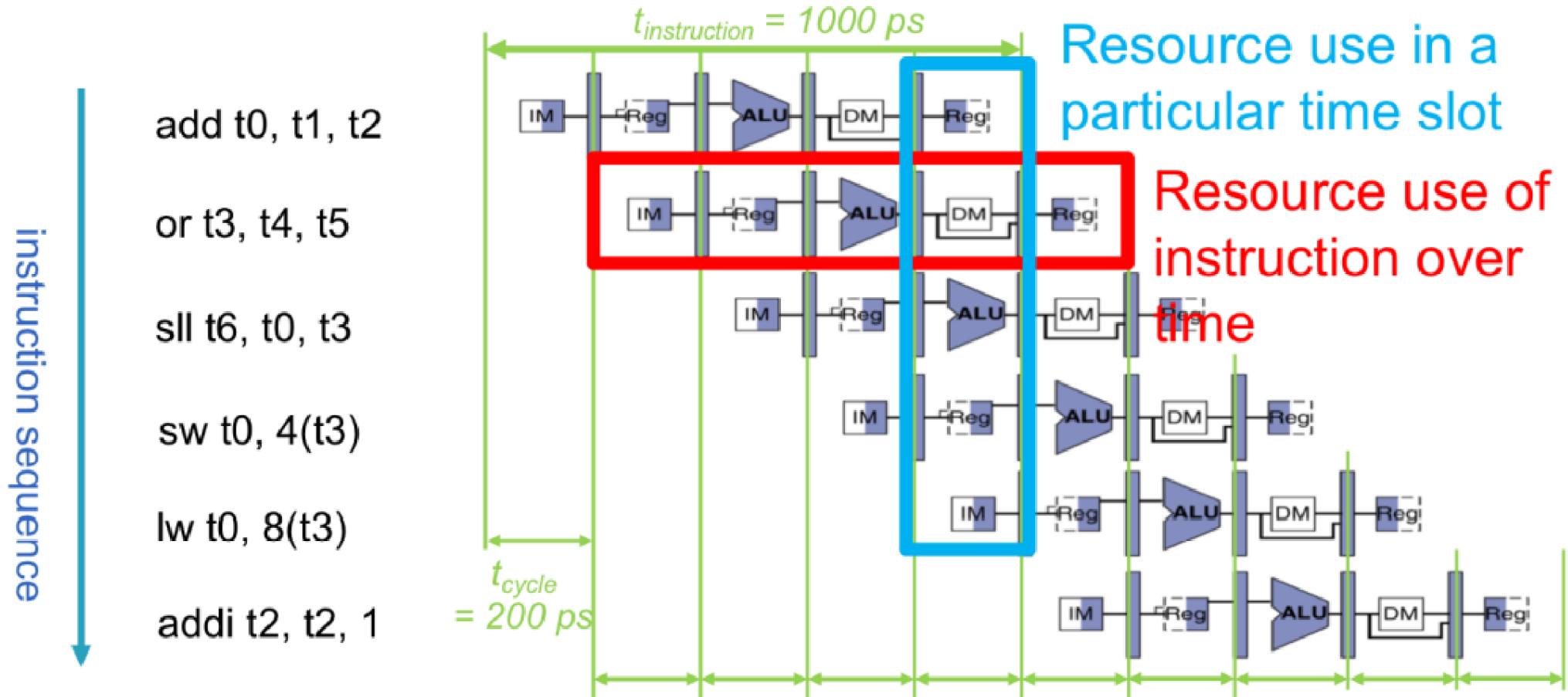
Exploring ILP 1: Pipelining

Recall: five stages of RISC-V datapath

1. **IF**: instruction fetch (Read `InstMem`)
2. **ID**: instruction decode (Read `RegFile`)
3. **EX**: execution (Computation `ALU`)
4. **MEM**: memory access (Read/Write `DataMem`)
5. **WB**: write back (Write `RegFile`)

Time for executing one instruction cannot be shortened because of the *dependencies* between stages.

However, we can use inactive components to execute the next/previous instruction.



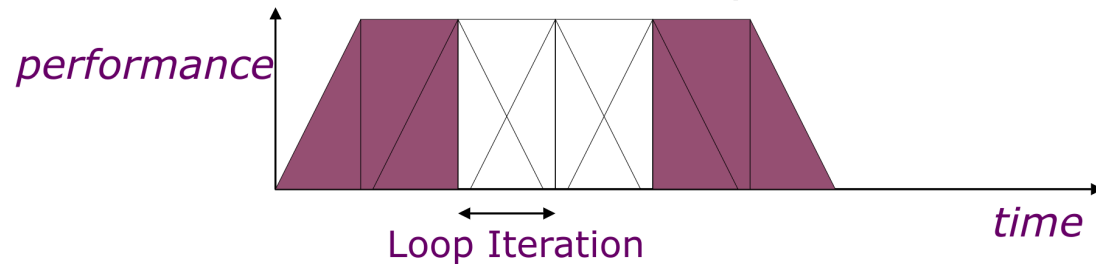
note:

- registers between stages: prevent interference with previous/next instruction.
- Pipelining allows higher clock frequency, but typically increases the latency.

Pipeline: Performance

- without pipeline: One of (IMEM, RegFile, ALU, DMEM) is in use.
- pipelined process: All stages are running, all of (IMEM, RegFile, ALU, DMEM) are in use.

Ideal Performance: n stages \rightarrow $\text{CPI} = 1$, $\text{freq} = n \times \text{freq}_{\text{unpipelined}}$.



1. Starting: Fill pipeline stages with instructions (parallelism increase)
2. Interim: All n stages are running simultaneously (maximum parallelism)
3. Stopping: Stages becomes free (parallelism decrease)

Pipelining: Hazards

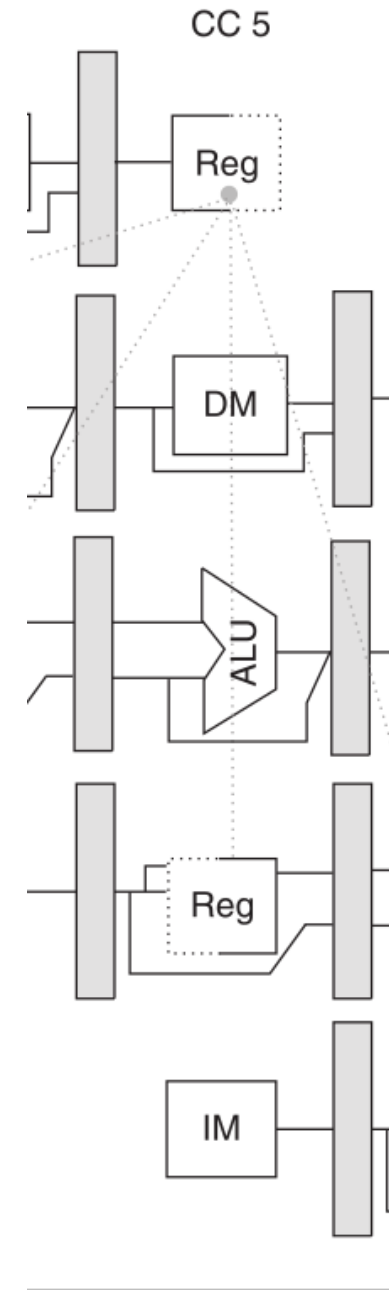
Unable to execute a stage of an instruction due to:

- **Strutruual Hazard:**
The required hardware resources is occupied by other instructions
- **Data Hazard:**
Dependent data not computed and stored yet
- **Control Hazard:**
Jump/Branch about to happen, unable to fetch correct instruction.

Hazard cause pipeline to stall, resulting in $CPI > 1$

Structural hazards

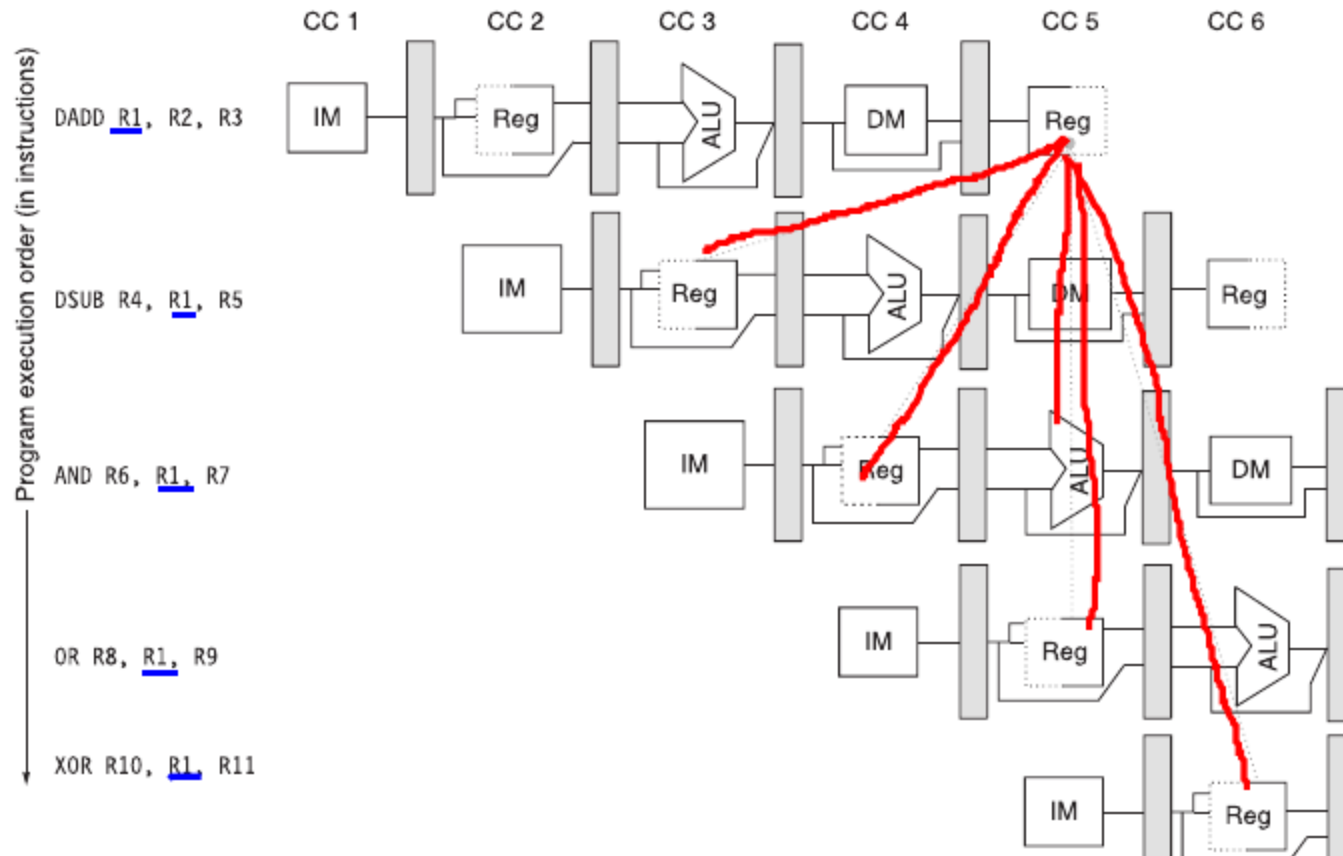
- On RegFile
instruction decode and
register writeback
- On Mem
instruction fetch and memory
read/write
- On ALU / FPU
certain computations requires
more than one cycle to
complete



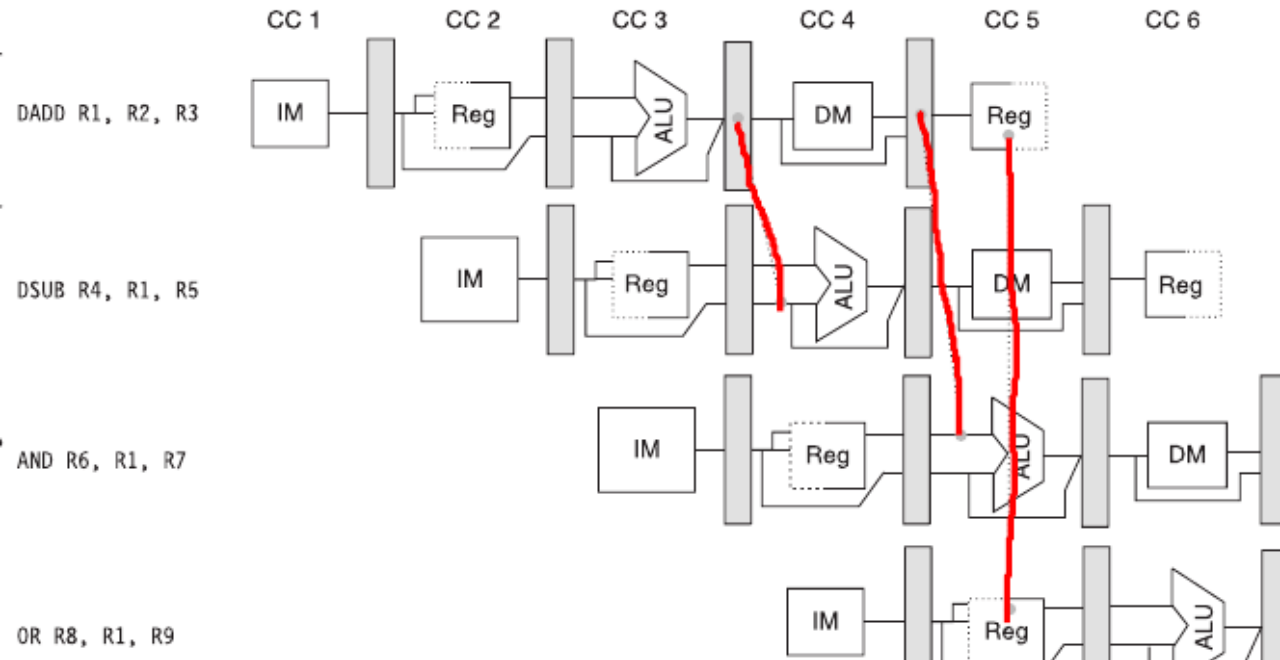
Solving structural hazards

- For `RegFile` : W/R on rise/fall edges respectively.
- For `Mem` : Separation of data memory and instruction memory.
- For `ALU` / `FPU` :
 - Re-ordering the instructions (in compile time | in runtime)
 - Adding more computation resources...

Data hazards



We mainly concern the read after write (RAW) dependency.

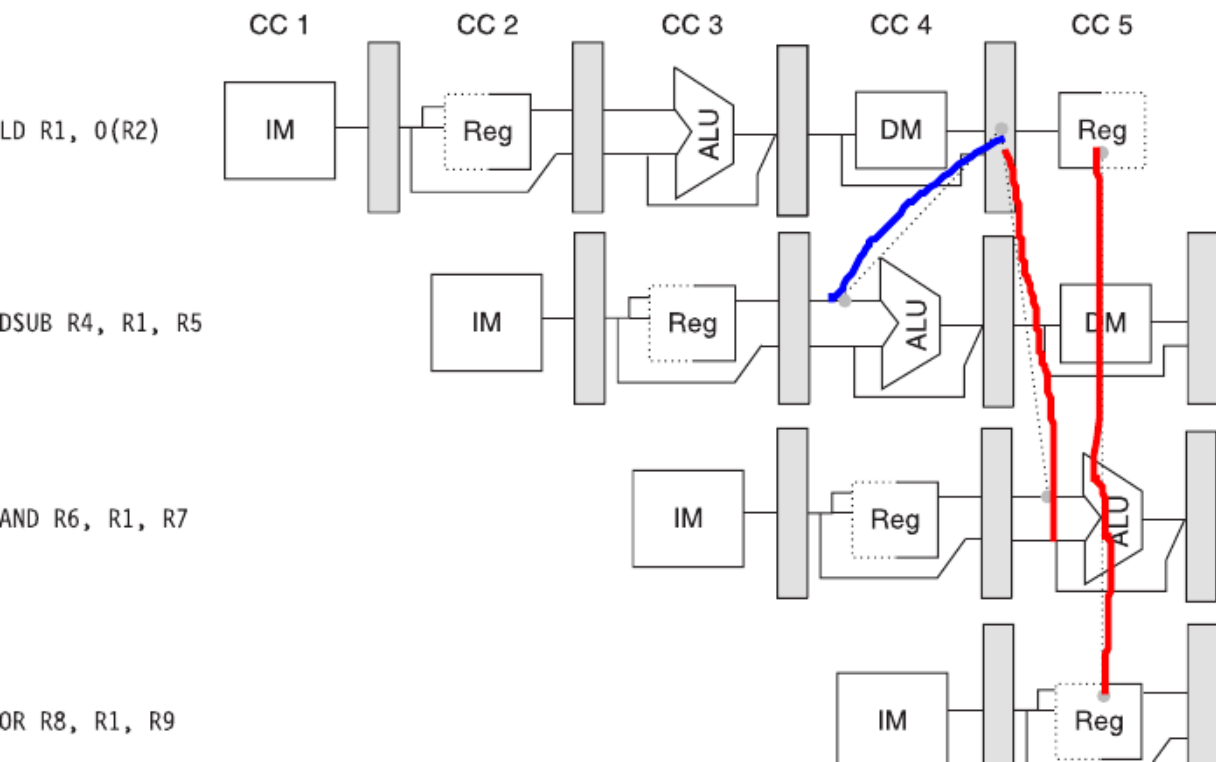


Resolving RAW with forward

Addition result is ready in cycle 3.

1. **sub** : result used in cycle 4
2. **and** : result used in cycle 5
3. **or** : result fetched in cycle 6

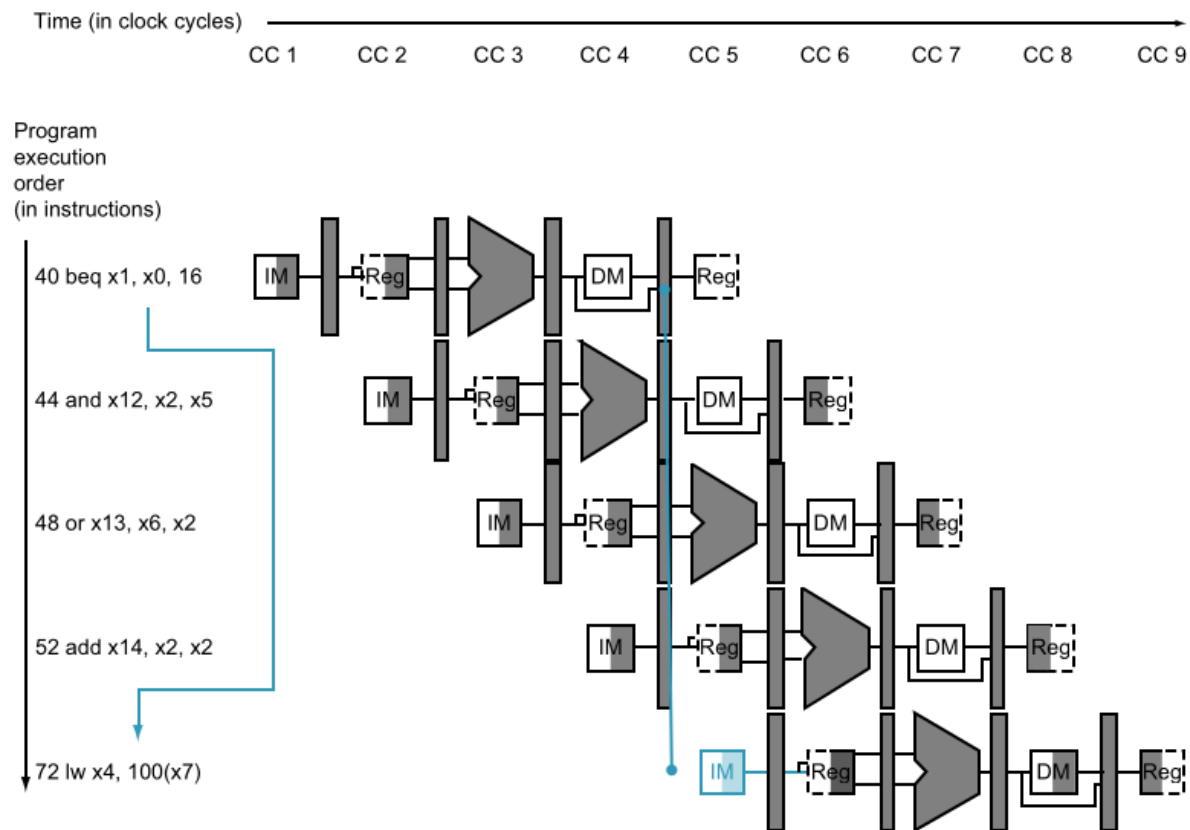
Replace the out-dated value from register with the result from ALU.



Unavoidable stall

- `ld x1, 0(x2)` : result of `ld` ready in cycle 4
- `sub x4, x1, x5` : value of `x1` used in cycle 4

Mitigation: Reorder the code, inserting an independent instruction after `ld`.



Control hazards

PC changes on rising edge of cycle 5, if branch taken.

Unable to determine the correct instruction to fetch in cycle 2, 3, 4.

Speculative execution

Idea: *continue execution, undo the wrong actions if necessary*

- eager execution: *executing both if branch and else branch*
- branch prediction: *guess the result of branch comparison*
 - Nothing wrong on correct prediction
 - Flush pipeline on wrong prediction
(branch compare result available in cycle 3, 2 instructions affected)

Performance tip for C/C++ programmer:

- short-circuit logical operators `&&` and `||` involves branch prediction.
- bitwise operators `&` and `|` are eagerly evaluated.
- See also: CppCon 2021 branchless programming [video slides](#)

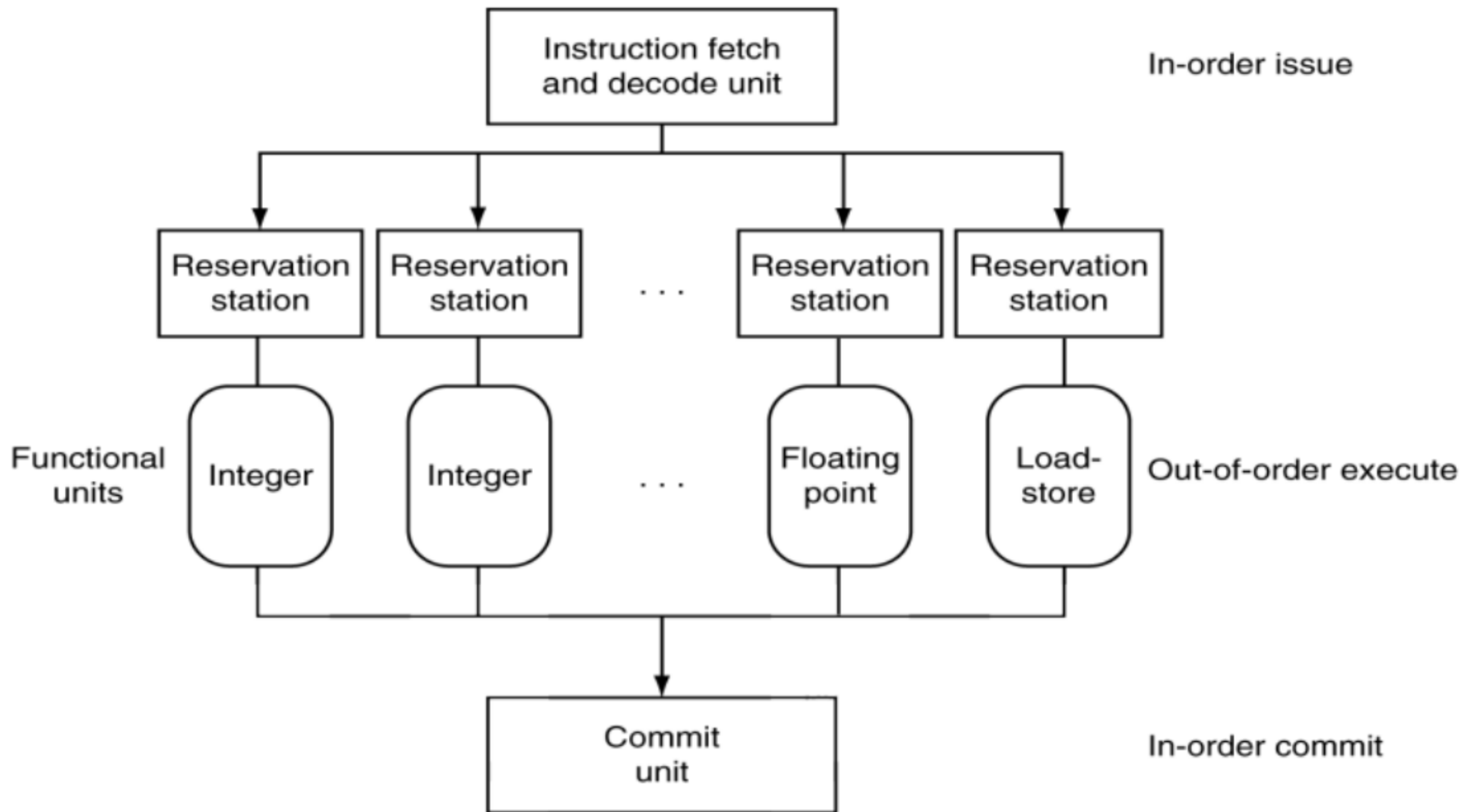
Exploiting ILP 2: Multiple Issue

What if we pack multiple ALUs, FPUs and Load/Store units into the processor?

- Single Issue: fetching and start executing a **single** instruction per cycle.
Not utilizing the extra computation power.
Ideal performance: $CPI = 1$
- Multiple Issue: fetching and starting executing **multiple** instructions per cycle.
Exploiting the additional hardware resources.
Ideal performance: $CPI < 1$

Typically paired with the following techniques:

- Out of order execution: help mitigating data dependencies (RAW, WAR, WAW)
- Speculative execution: help reduce impact of branches

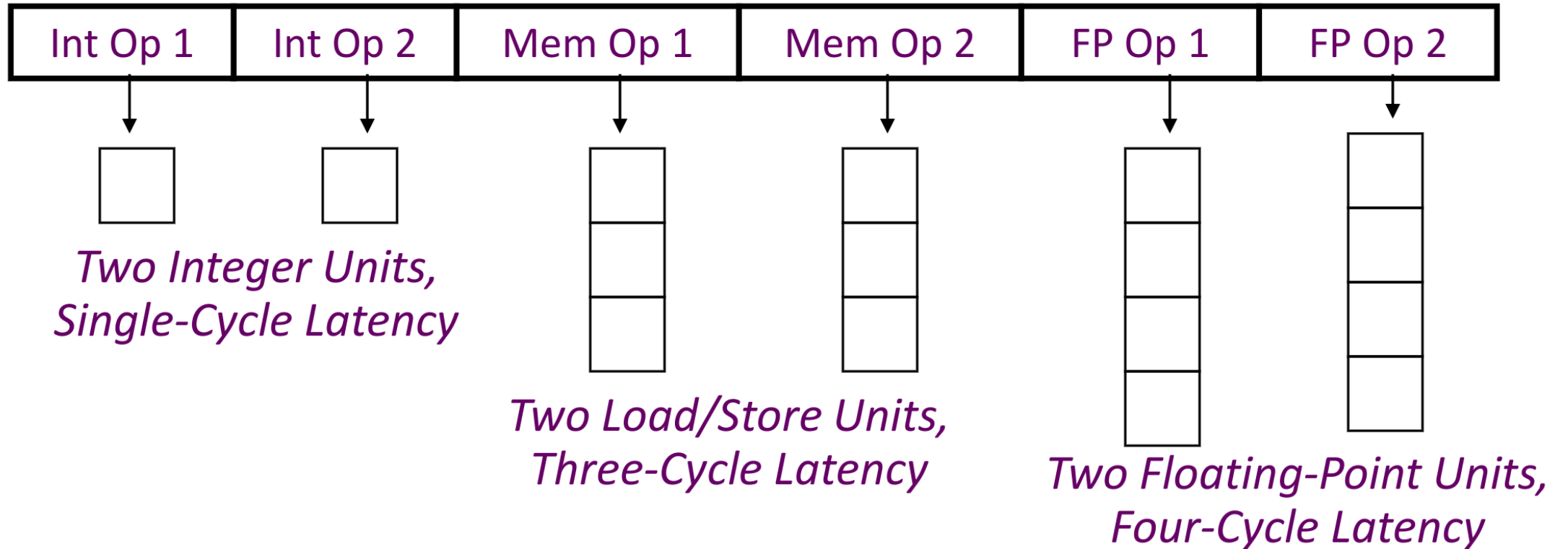


Common name	Issue structure	Hazard detection	Scheduling	Distinguishing characteristic	Examples
Superscalar (static)	Dynamic	Hardware	Static	In-order execution	Mostly in the embedded space: MIPS and ARM, including the Cortex-A53
Superscalar (dynamic)	Dynamic	Hardware	Dynamic	Some out-of-order execution, but no speculation	None at the present
Superscalar (speculative)	Dynamic	Hardware	Dynamic with speculation	Out-of-order execution with speculation	Intel Core i3, i5, i7; AMD Phenom; IBM Power 7
VLIW/LIW	Static	Primarily software	Static	All hazards determined and indicated by compiler (often implicitly)	Most examples are in signal processing, such as the TI C6x
EPIC	Primarily static	Primarily software	Mostly static	All hazards determined and indicated explicitly by the compiler	Itanium

Figure 3.19 The five primary approaches in use for multiple-issue processors and the primary characteristics that distinguish them. This chapter has focused on the hardware-intensive techniques, which are all some form of superscalar. Appendix H focuses on compiler-based approaches. The EPIC approach, as embodied in the IA-64 architecture, extends many of the concepts of the early VLIW approaches, providing a blend of static and dynamic approaches.

VLIW: Very Long Instruction Word

Instead of using a simply 32-bit instruction word to encode one operation, VLIW packs multiple operators into one instruction word (generally longer).



Static issuing and static scheduling.

Compiler responsible for identifying possible parallelism and grouping OPs

Superscalar

To fetch, issue for execution and finish more than one instruction per cycle.

t=0	t=1	t=2	t=3	t=4	t=5	t=6	t=7	t=8
F 1	F 3	F 5	F 7	F 9				
F 2	F 4	F 6	F 8	F 10				
	D 1	D 3	D 5	D 7	D 9			
	D 2	D 4	D 6	D 8	D 10			
		X 1	X 3	X 5	X 7	X 9		
		X 2	X 4	X 6	X 8	X 10		
			M 1	M 3	M 5	M 7	M 9	
			M 2	M 4	M 6	M 8	M 10	
				W 1	W 3	W 5	W 7	W 9
				W 2	W 4	W 6	W 8	W 10

Example of a 2-issue, 5-stage pipelined processor executing 10 instructions.

Hardware responsible for dependencies analysis and dynamic issuing.

Additional Exercises

- previous exams of CS110@ShanghaiTech
- previous exams of CS61c@Berkeley
- previous exams of CS211@ShanghaiTech
- previous exams of CS152@Berkeley

Reference & Acknowledgement

- Slides from **CS211@ShanghaiTech** by Prof. Chundong Wang
- Slides from **CS152@UC Berkeley** by Krste Asanovic
- Chapter 4 of **Computer Organization and Design 2nd edition**
- Chapter 3 of **Computer Architecture: A Quantitative Approach**