

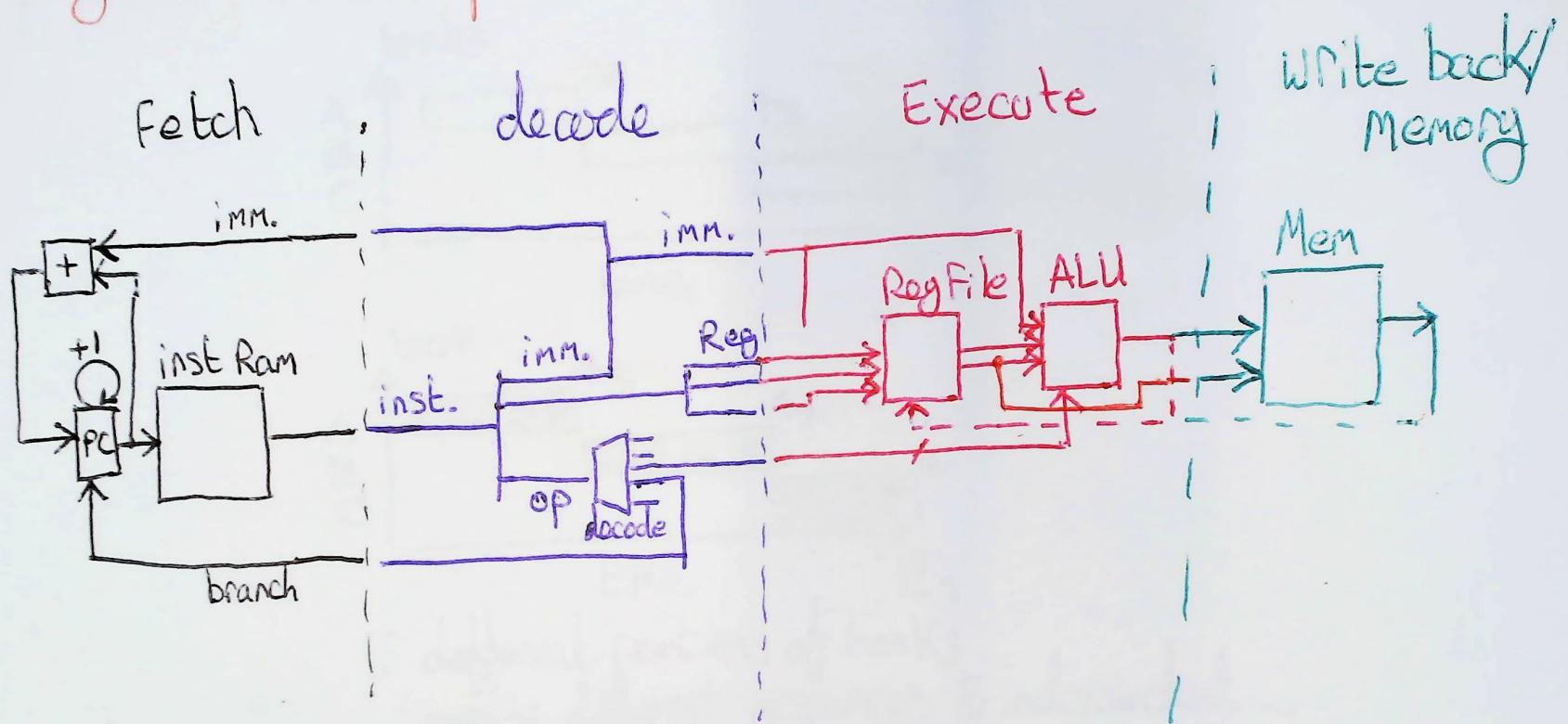
chapter 4: The processor

CPU performance factors for processor

- Instruction Count:
determined by ISA,
compiler, algorithm.

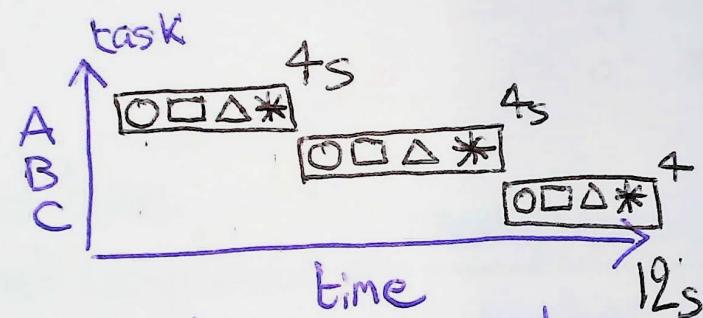
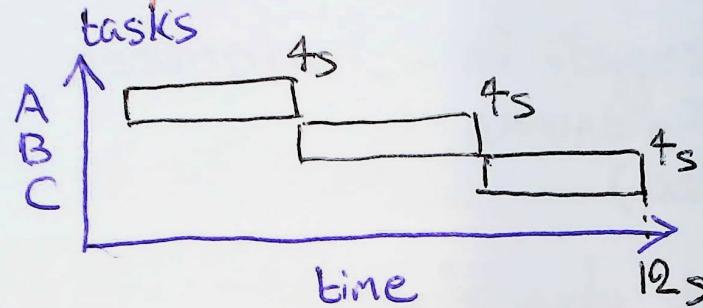
- clocks per instruction (CPI)
or Instructions per clock (IPC)
determined design of architecture
and Technology (nm)

Logic critical path



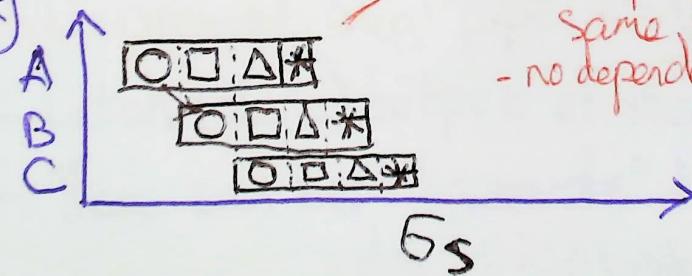
critical path: longest possible path
through the gates of a circuit.

Pipelining concept



If different portions of tasks
require different resources & independent -

pipelining



→ each portion takes about
same time.
- no dependencies between portions

How to calculate speedup of pipelining

assumption: - all stages take same amount of time.

(balanced pipeline)

- "steady state" speedup is of interest.

$$\text{Pipelined time of task completion} = \frac{\text{non-pipelined time of task}}{\text{number of pipeline stages}}$$

$$\underline{\text{ideal speedup}} = \frac{\text{non-pipelined time}}{\text{pipelined time}} = \frac{\# \text{ pipeline stages}}{\text{stages}}$$

Things that prevent ideal pipelining

- Latch (store) output of each stage has overhead
- There can be "hazards" that prevent ideal pipelining.

in processor easy to avoid hazard with "noop".

Pipeline Hazards

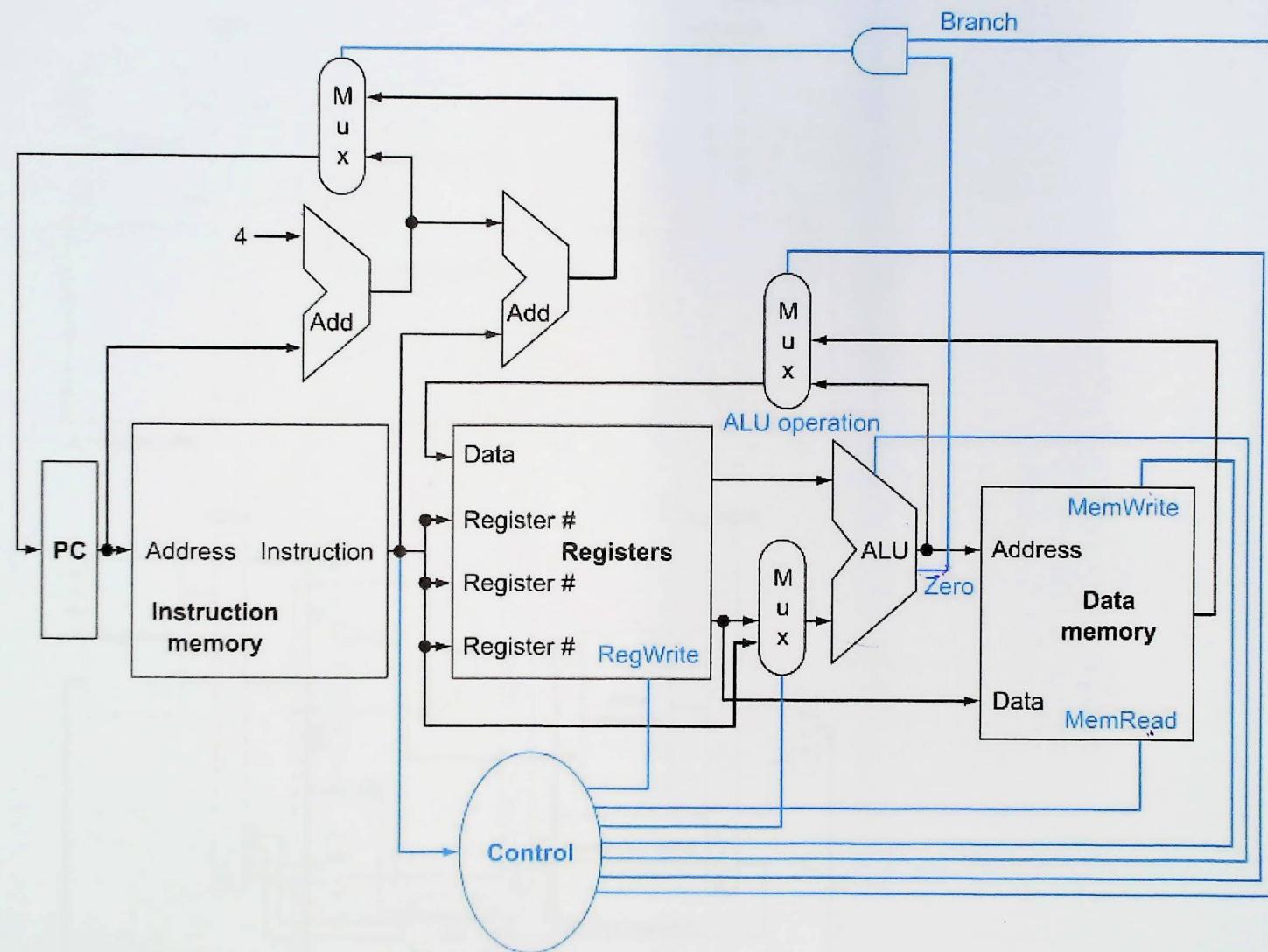
(dependencies that prevent ideal pipelining)

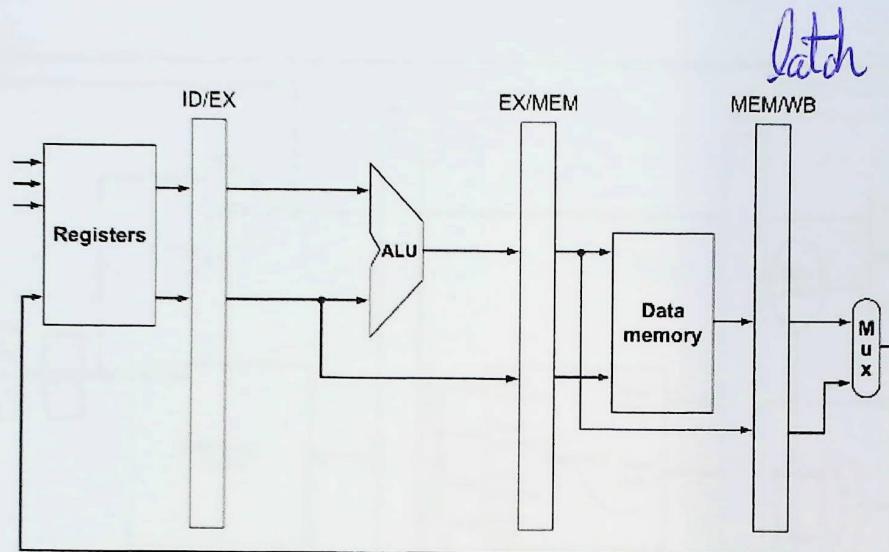
3 main type of hazards:

- structural: when different stages of pipeline access same hardware resource

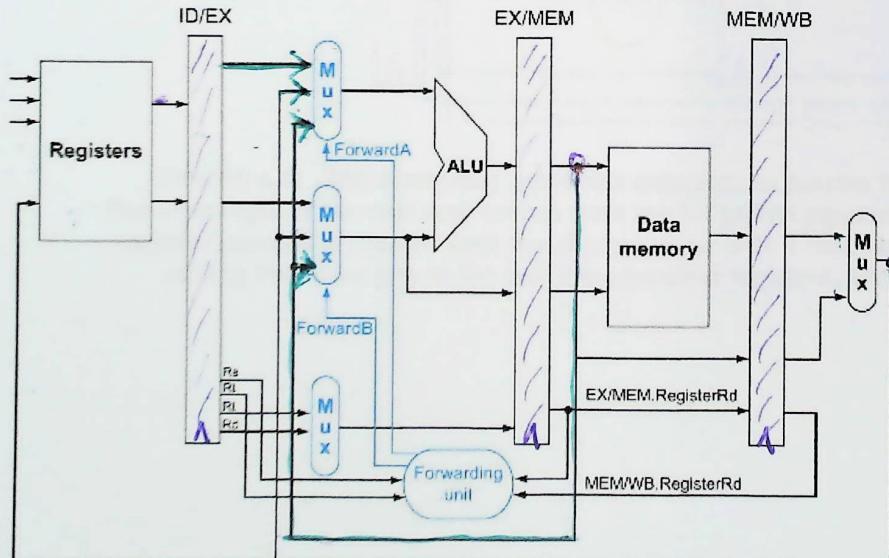
- data: when result of a prior instruction is need to execute next instruction.

- control: when the next instruction to fetch depends on a prior instruction.





a. No forwarding



b. With forwarding

bypassing

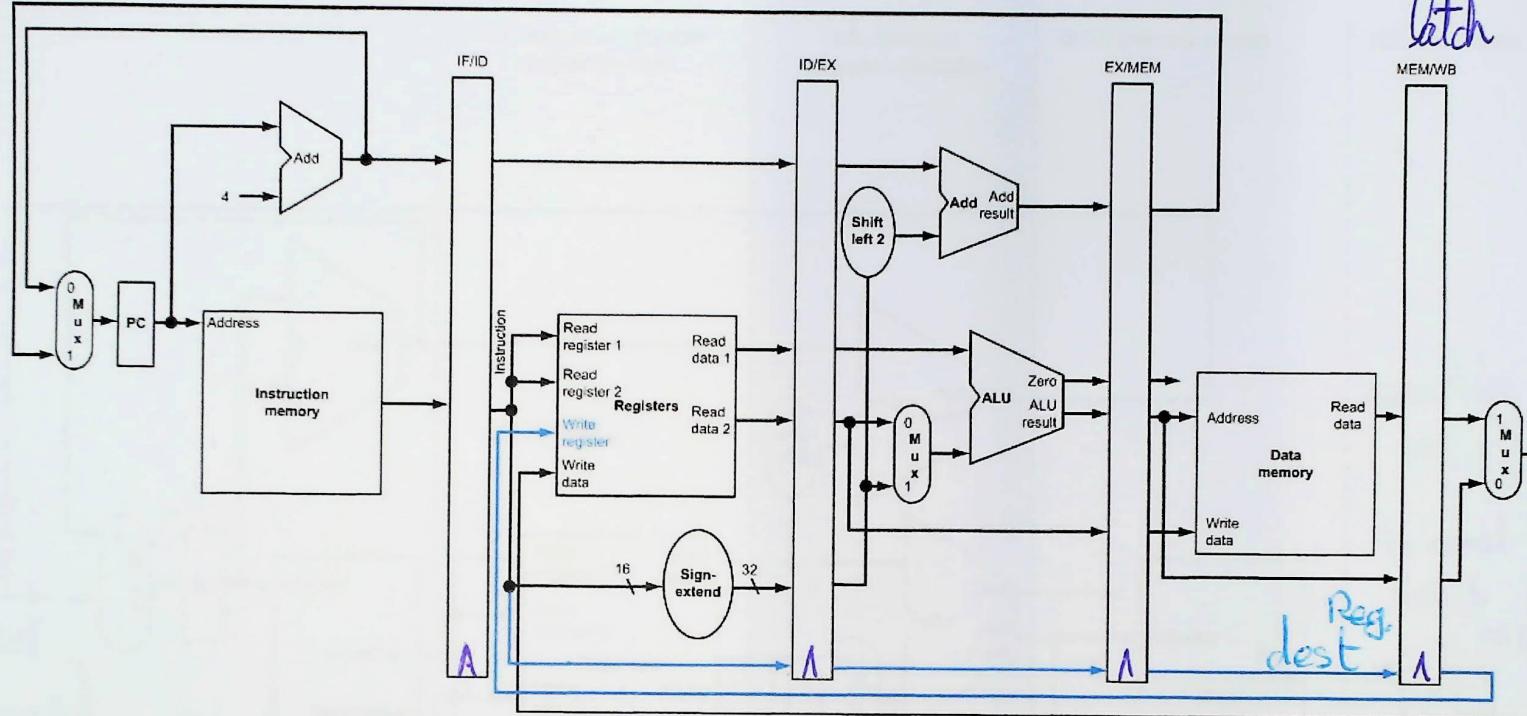
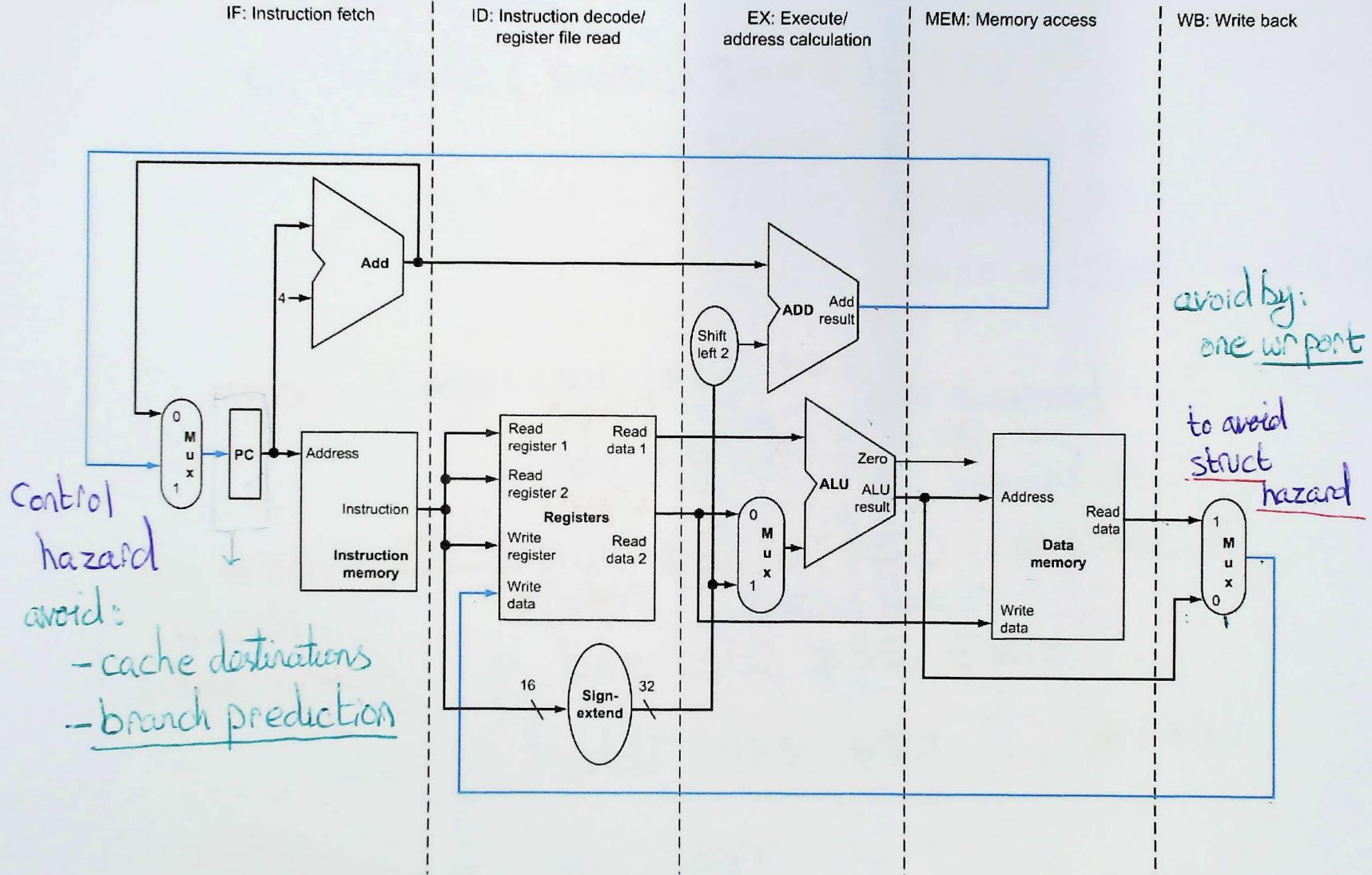


FIGURE 4.41 The corrected pipelined datapath to handle the load instruction properly.

The write register number now comes from the MEM/WB pipeline register along with the data. The register number is passed from the ID pipe stage until it reaches the MEM/WB pipeline register, adding five more bits to the last three pipeline registers. This new path is shown in color.



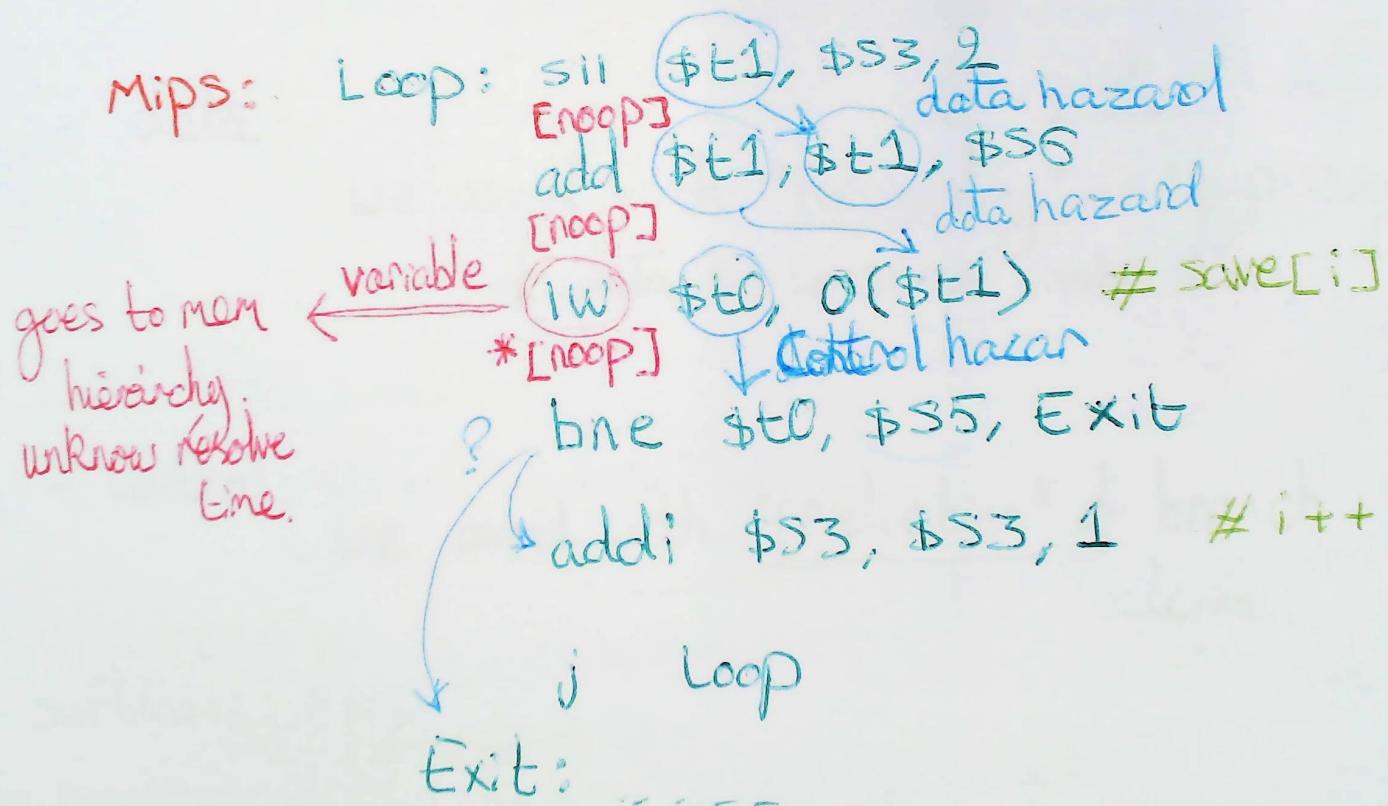
Example of pipeline hazards code

C: while (save[i] == k) i++;

assump.: i is in \$53

k is in \$55

base of save is in \$56



How we avoid different pipeline hazards [in Hardware]

- structural :

We avoid by having single point of access to resource.

"writeback"
→ wb phase

- data :

we avoid with "Forwarding" or "bypassing" of data to consumer instruction.

→ a quick path data connection

- Control:

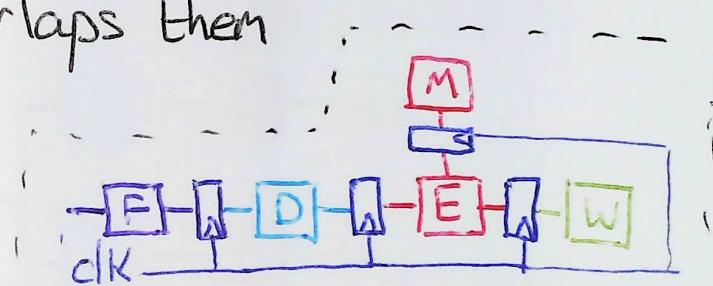
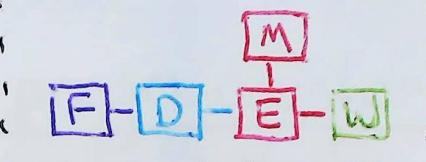
we avoid with "prediction" of branch outcome.

Software: noop

Different Types of parallelism

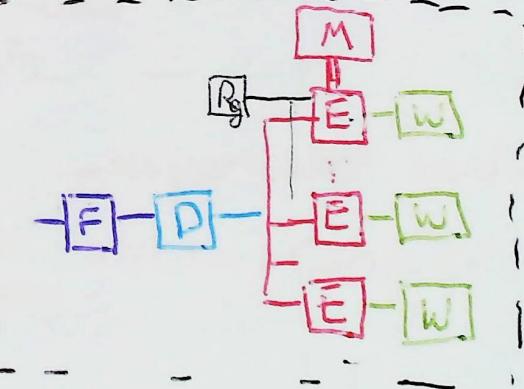
1) sub-instruction parallelism

- pipelining
- break instructions into portions and overlaps them



2) data-level parallelism

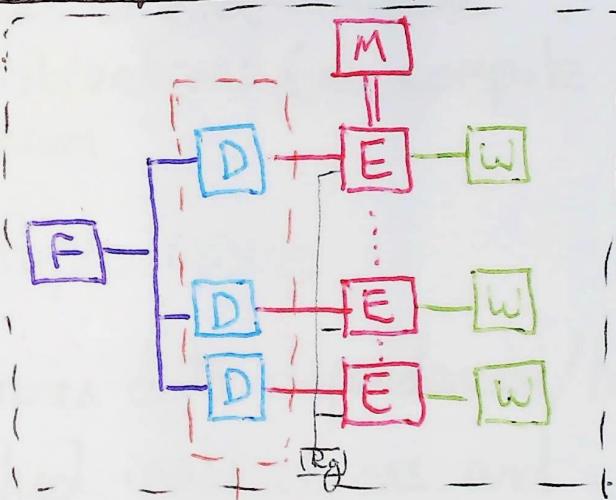
- one instruction works on multiple data.
- SIMD (single instruction multiple data)
- common in media/GPU processors



3) Instruction Level parallelism

- single fetch. single program counter.
- But possibly can fetch multiple inst. per ck.
- Can execute multiple inst. per clock.

Forms {
1- superscalar *
2- Very
Long Inst. Width:
(VLIW)



pull out → - superscalar: hardware.
independent - VLIW: compiler.
instructions

3.a) ILP: static multiple issue

- pack multiple independent instructions together, and fetched together
- Very Long Instruction width architecture (VLIW)
- Compiler need to find independent instructions. (at compile time).
e.g. Intel Itanium

3.b) ILP: dynamic multiple issue

- hardware examines data dependencies/hazards between fetched instructions. and clears them for execute.
- Superscalar

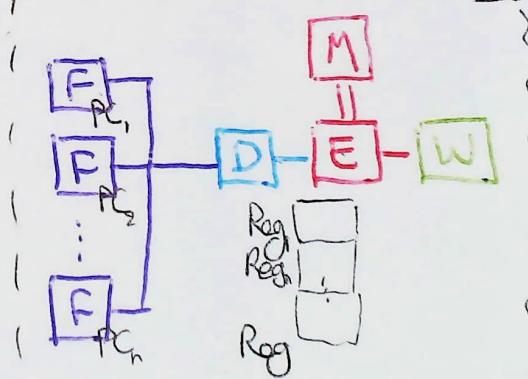
4) Thread Level parallelism

Multithreading:

- managed by operating system (OS).
- OS stores and retrieves the program counter and reg file of a program, to switch it in/out of execution.

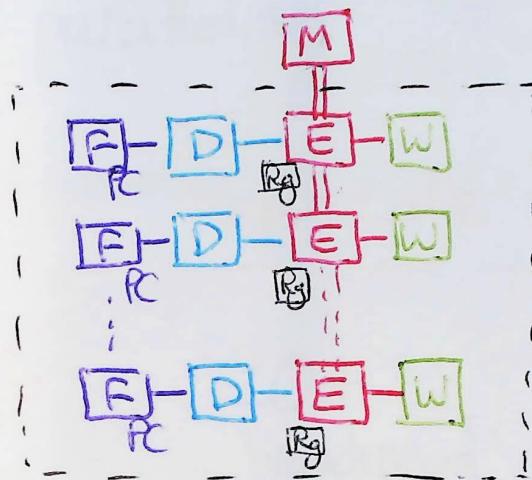
Simultaneous multithreading

- keep multiple program counter and regfile.
- dynamically pick instructions from threads to keep processor busy.
- Intel: "hyperthreading"



5) multicore - Full thread-level parallelism

- stamp out copies of your "core" processor.
- important: how to hook up all the cores to a unified memory.



Contemporary processors and forms of parallelism

	AMD Ryzen 7000	Intel Raptor Lake	NVidia 4090 GPU
Pipelining (sub-instruction parallel)	20 branch pred.	20+ branchpred	deep
data level	special SIMD for Media	special SIMD for Media	SIMD 32 default
Instruction level	superscalar 5X	superscalar 5X	None
Thread level	2 threads per core	1 thread per core 2 threads per core	32 threads per core
Core level	16 core	15 "E" + 8 "P" 2 threads / thread	Very long! {