



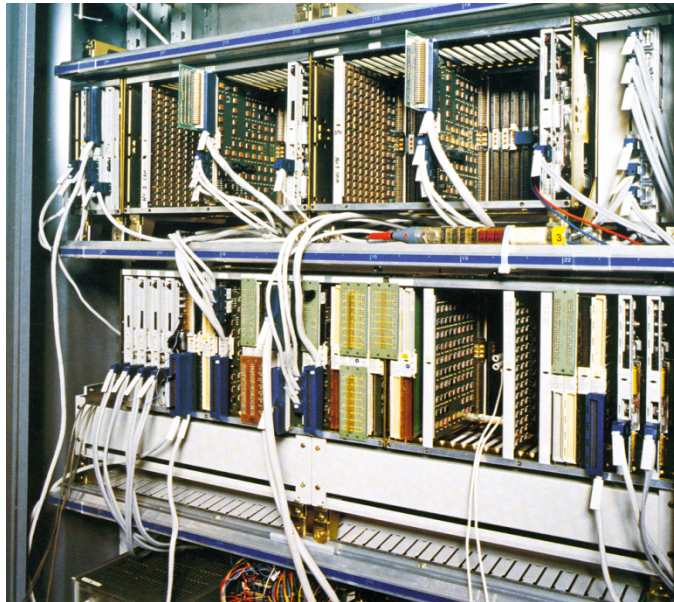
Introduction to Computer Architecture

Erik Hagersten
Uppsala University



UPPSALA
UNIVERSITET

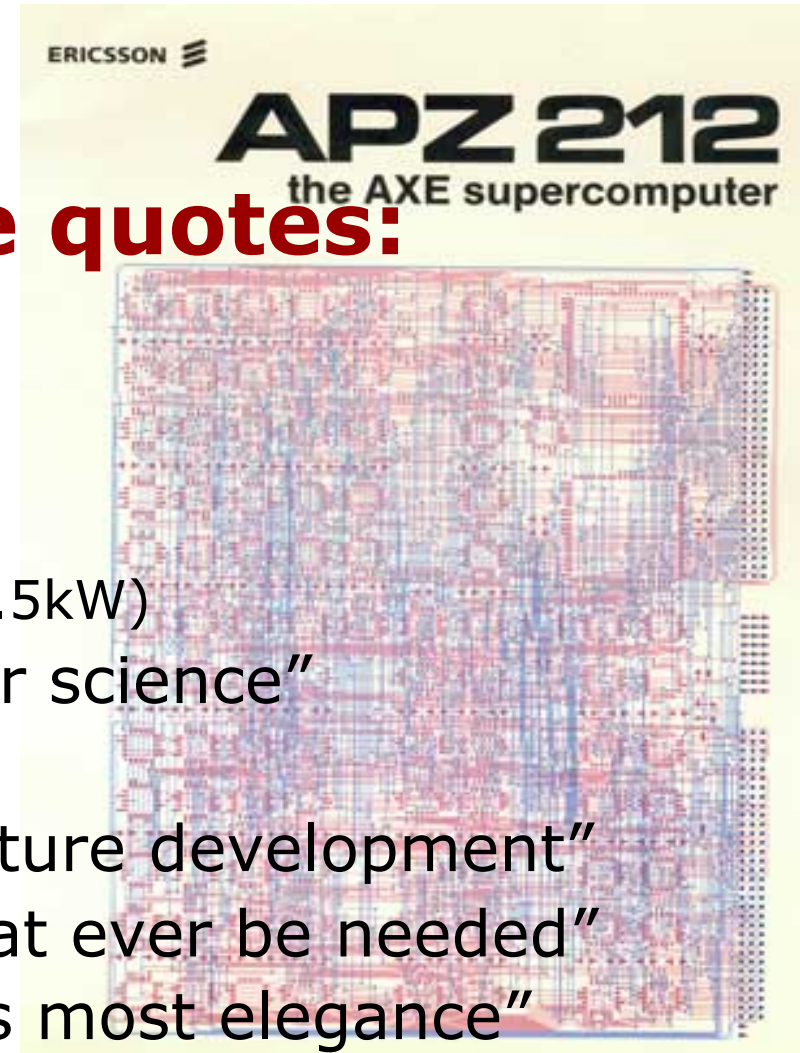
≈35 years ago: APZ 212 @ 5MHz "the Ericsson supercomputer"



PDC
Summer
School
2017

APZ 212

marketing brochure quotes:

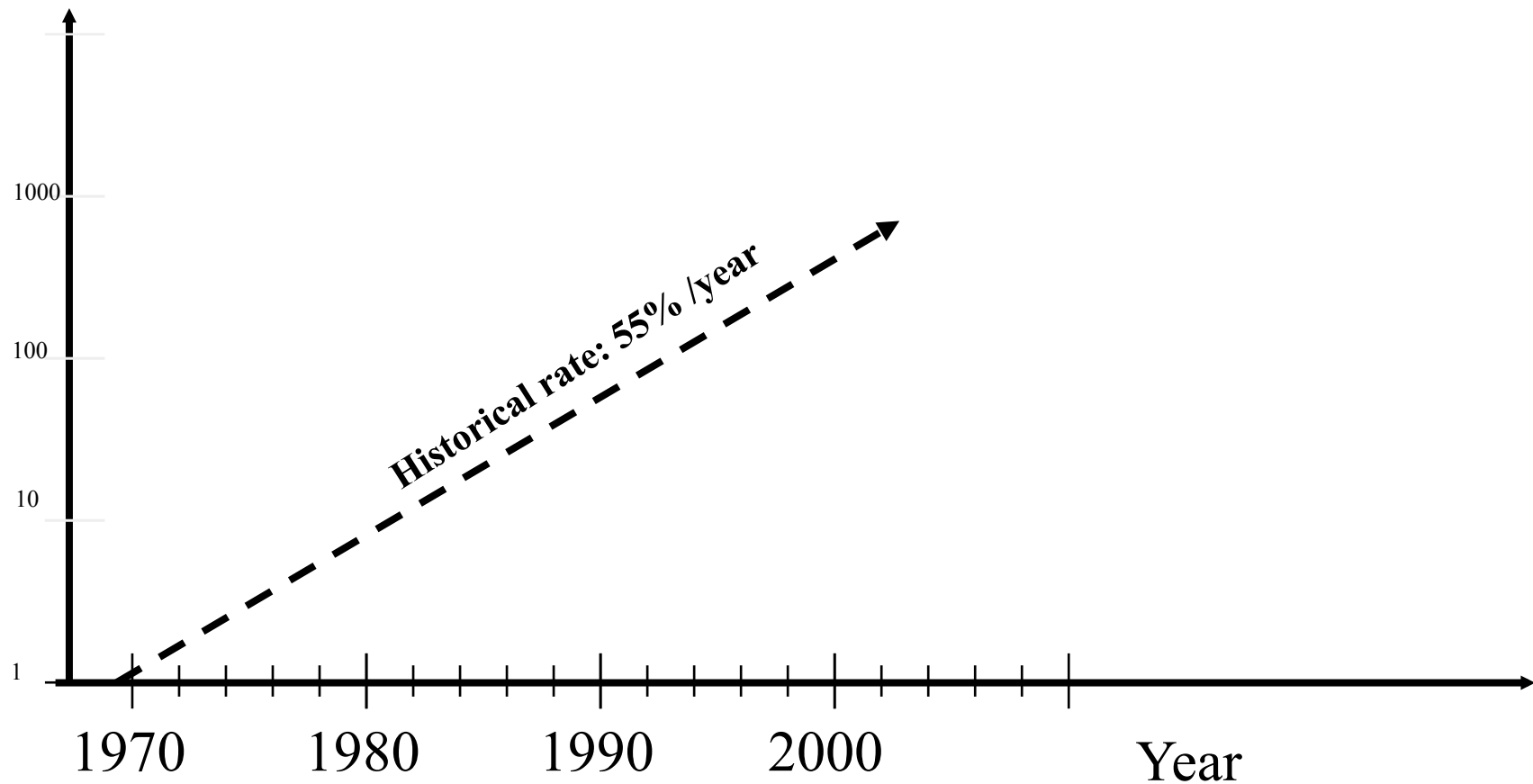


- "Very compact"
 - ✱ 6 times the performance
 - ✱ 1/6:th the size
 - ✱ 1/5 the power consumption (1.5kW)
- "A breakthrough in computer science"
- "Why more CPU power?"
- "All the power needed for future development"
- "...800,000 BHCA, should that ever be needed"
- "SPC computer science at its most elegance"
- "Using 64 kbit memory chips"



CPU Improvements

Relative Performance
[log scale]





How to get efficient architectures...

Very hard today

- ~~Increase clock frequency~~
- Create and explore locality:
 - a) Spatial locality
 - b) Temporal locality
- Create and explore parallelism
 - a) Instruction level parallelism (ILP)
 - b) Thread level parallelism (TLP)
 - c) Memory level parallelism (MLP)
- Speculative execution
 - a) Out-of-order execution
 - b) Branch prediction
 - c) Prefetching



Outline of these lectures

1. Processor implementations
2. Caches and memory system
3. Multiprocessors
4. HW optimizations
5. Multicores
6. *SW optimizations (if there is time...)*



CPU architecture overview

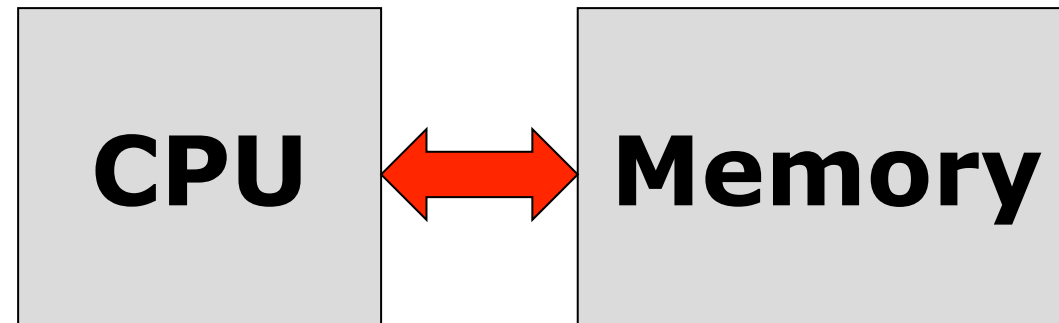
Erik Hagersten
Uppsala University



UPPSALA
UNIVERSITET

How fast is a computer today?

PDC
Summer
School
2017





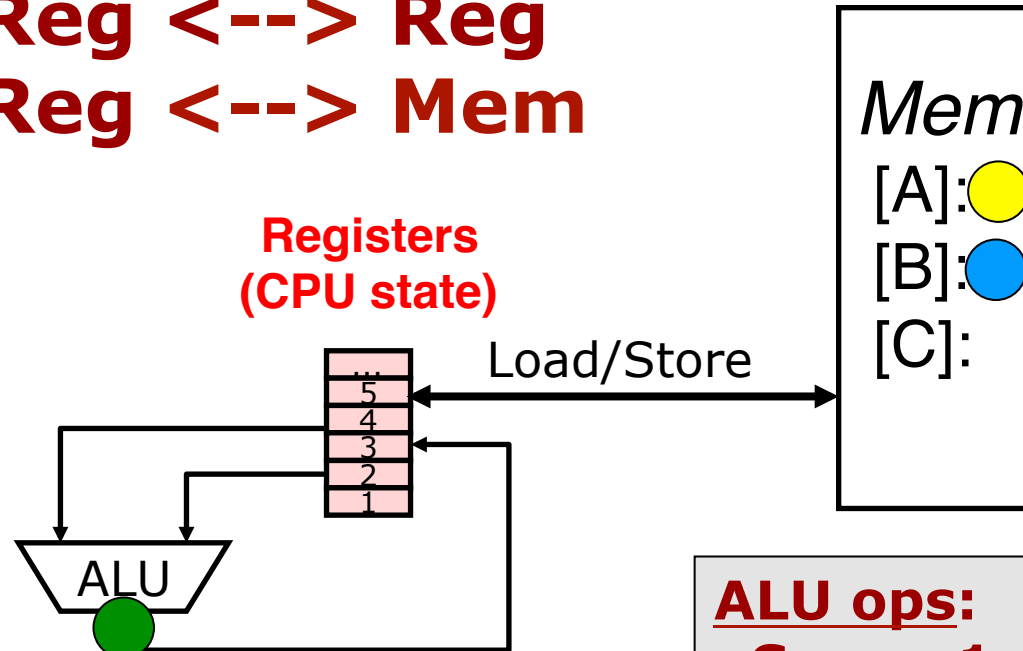
Load/Store architecture (e.g., "RISC")

ALU ops:

Reg <--> Reg

Mem ops:

Reg <--> Mem



Example: $C = A + B$



PC
PC
PC
PC

```
LD R1, [A]
LD R3, [B]
ADD R2, R1, R3
ST R2, [C]
```

ALU ops:

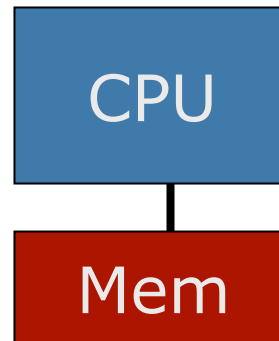
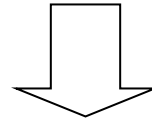
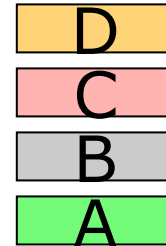
- Source1
- Source2
- Destination

Memory ops



Lifting the CPU hood (simplified...)

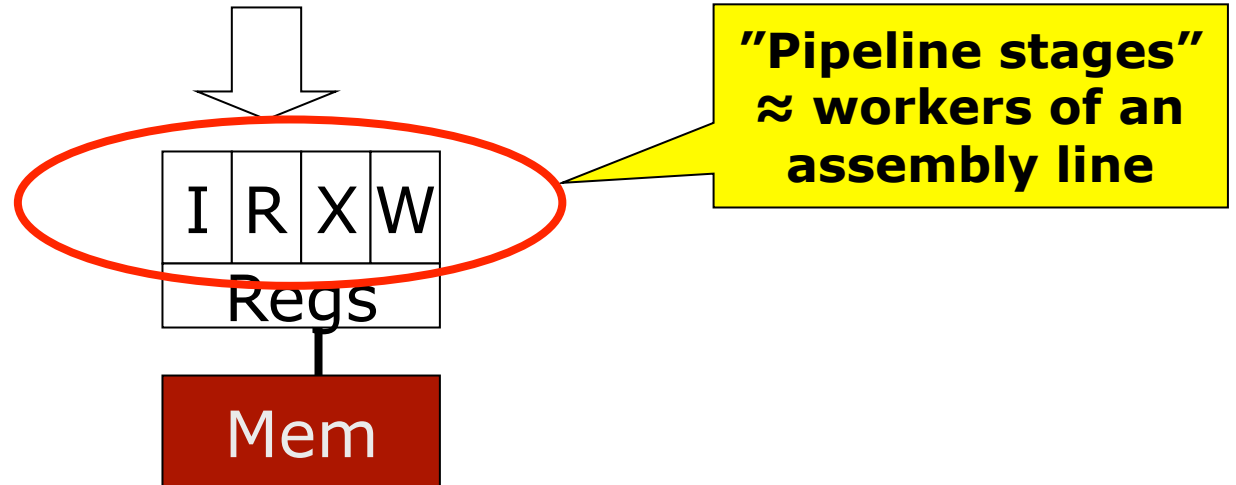
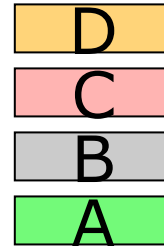
Instructions:





Pipeline

Instructions:





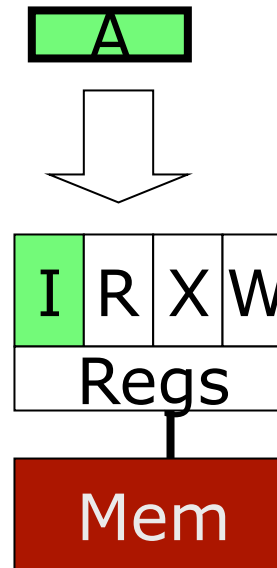
UPPSALA
UNIVERSITET



PDC
Summer
School
2017

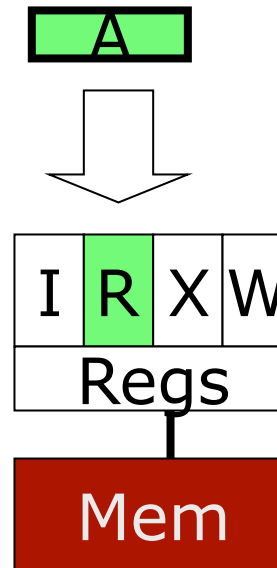


Pipeline



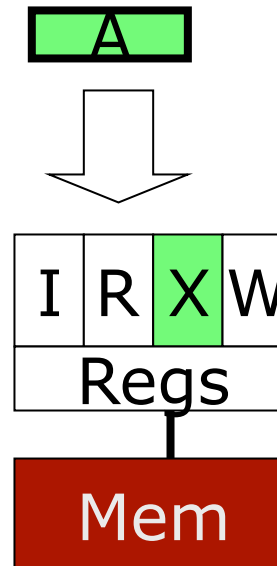


Pipeline



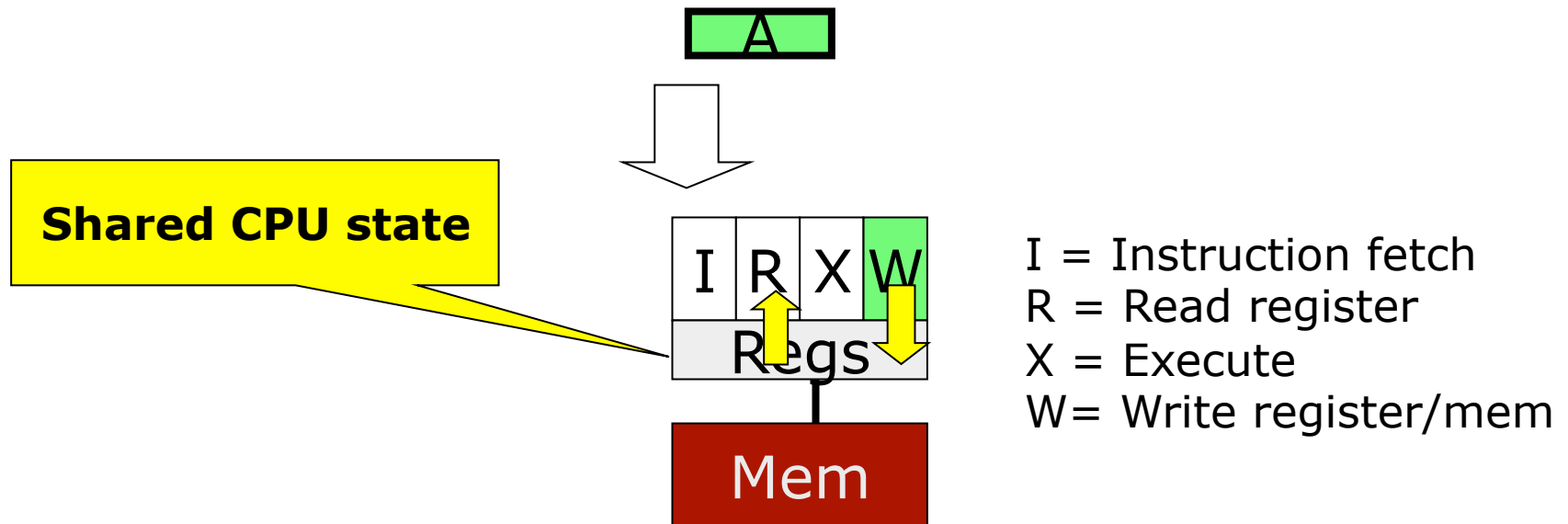


Pipeline





Pipeline:

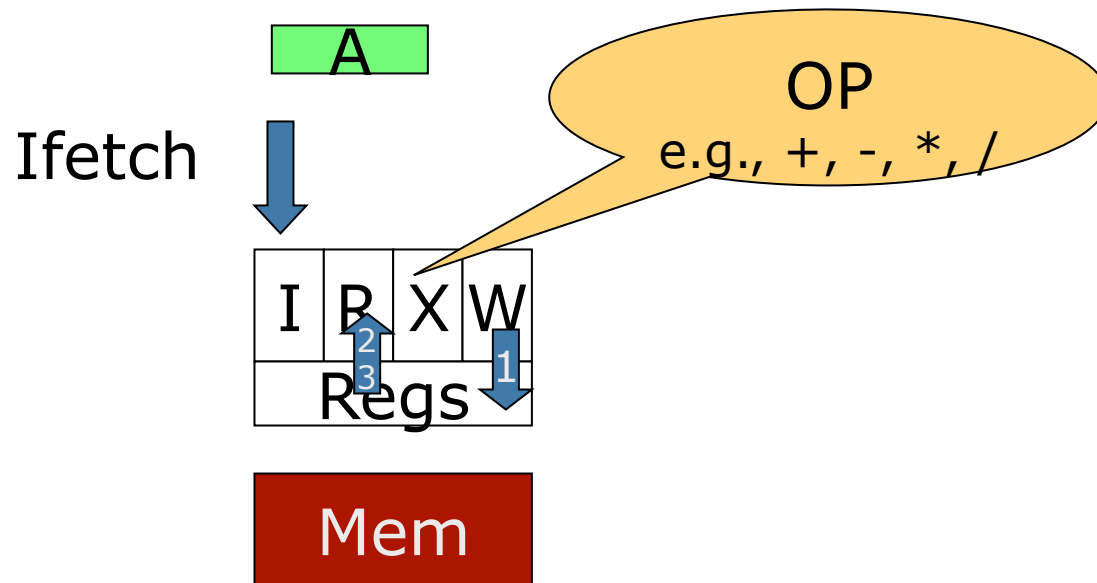




Register Operations [ALU operation]

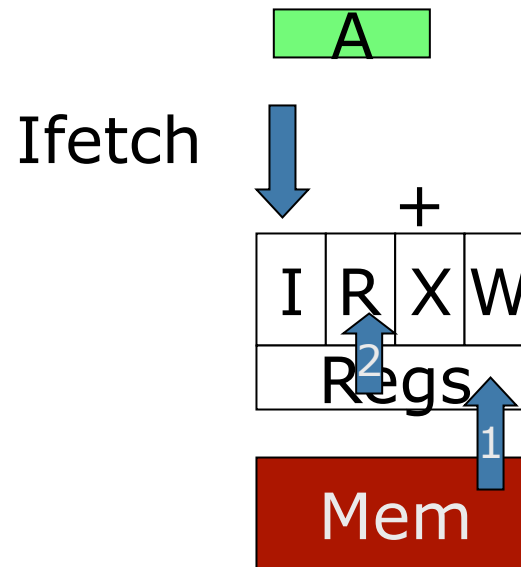
ADD R1, R2, R3

(a.k.a. $R1 := R2 + R3$)



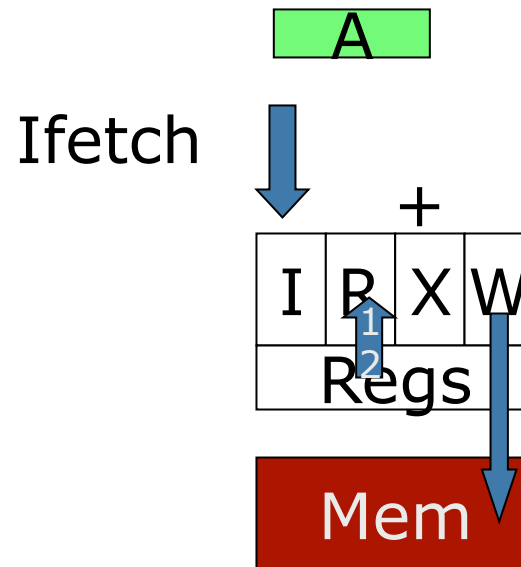
Load Operation:

LD R1, mem[cnst+R2]



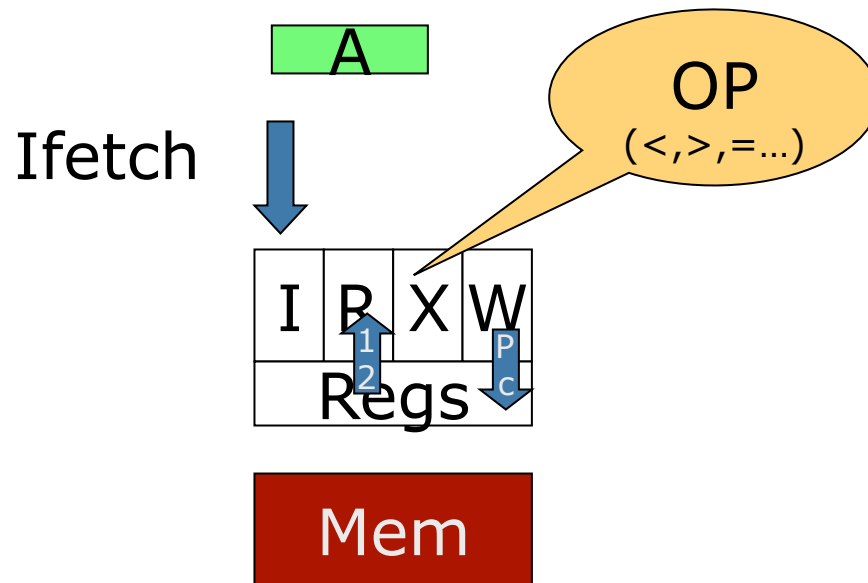
Store Operation:

ST R1, mem[cnst+R2]



Branch Operations:

if (R1 < Const) GOTO mem[R2]

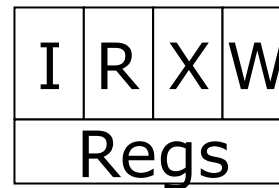
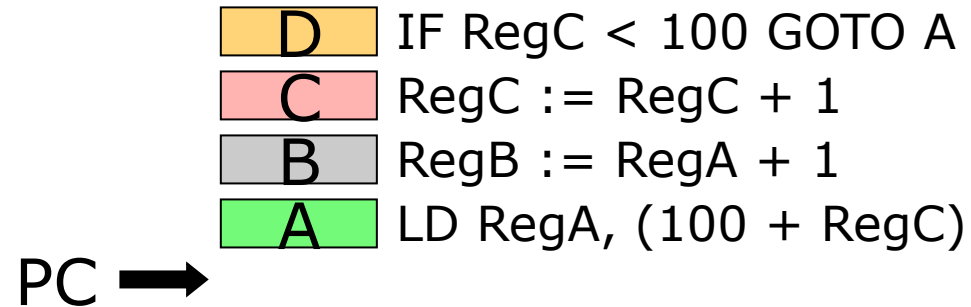


PC = Program Counter.

A special register pointing to the next instruction to execute

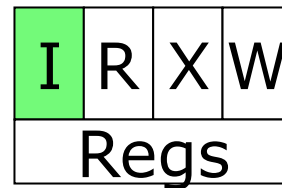
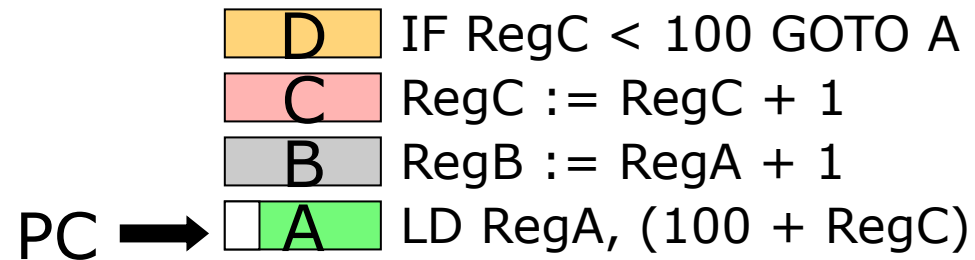


Initially

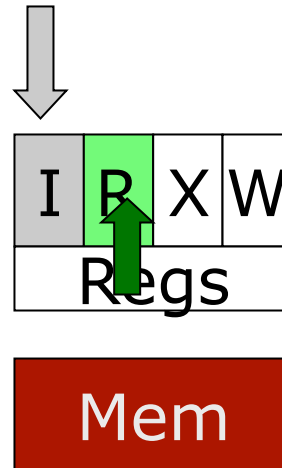
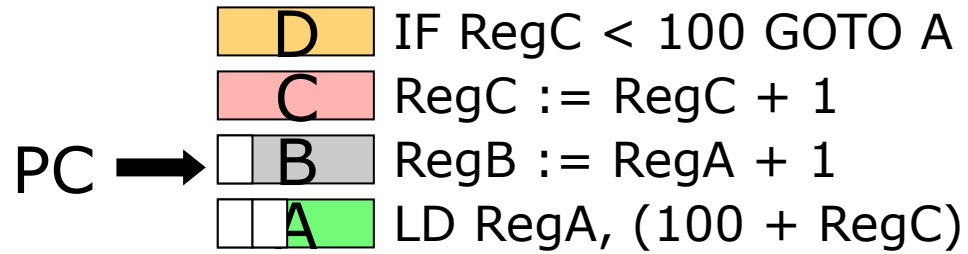




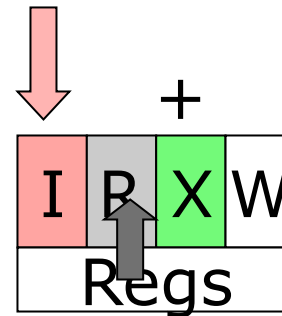
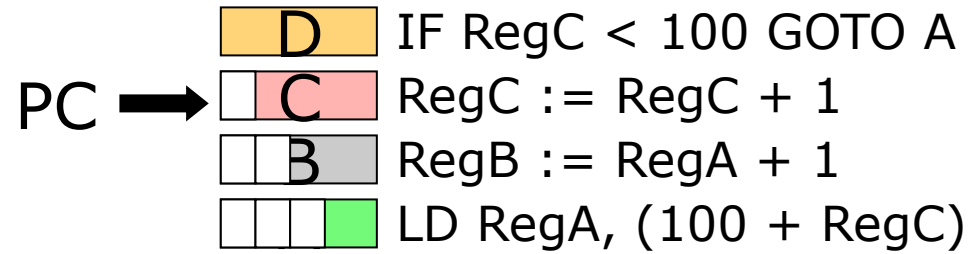
Cycle 1



Cycle 2



Cycle 3



Cycle 4

PC →

	D
--	---

 IF RegC < 100 GOTO A

	C
--	---

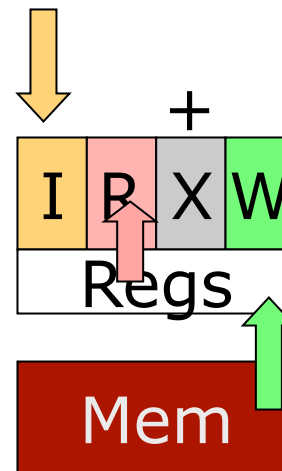
 RegC := RegC + 1

		X
--	--	---

 RegB := RegA + 1

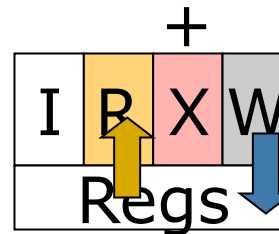
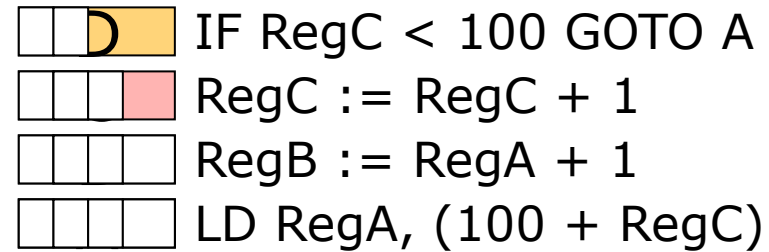
--	--	--	--

 LD RegA, (100 + RegC)



Cycle 5

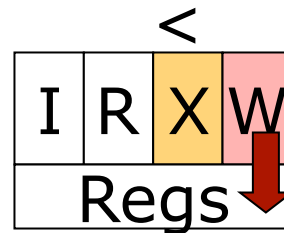
PC →



Cycle 6

PC →

□	□	□	■	IF RegC < 100 GOTO A
□	□	□	□	RegC := RegC + 1
□	□	□	□	RegB := RegA + 1
□	□	□	□	LD RegA, (100 + RegC)



Cycle 7

PC →

--	--	--	--

 IF RegC < 100 GOTO A

--	--	--	--

 RegC := RegC + 1

--	--	--	--

 RegB := RegA + 1

A:

--	--	--	--

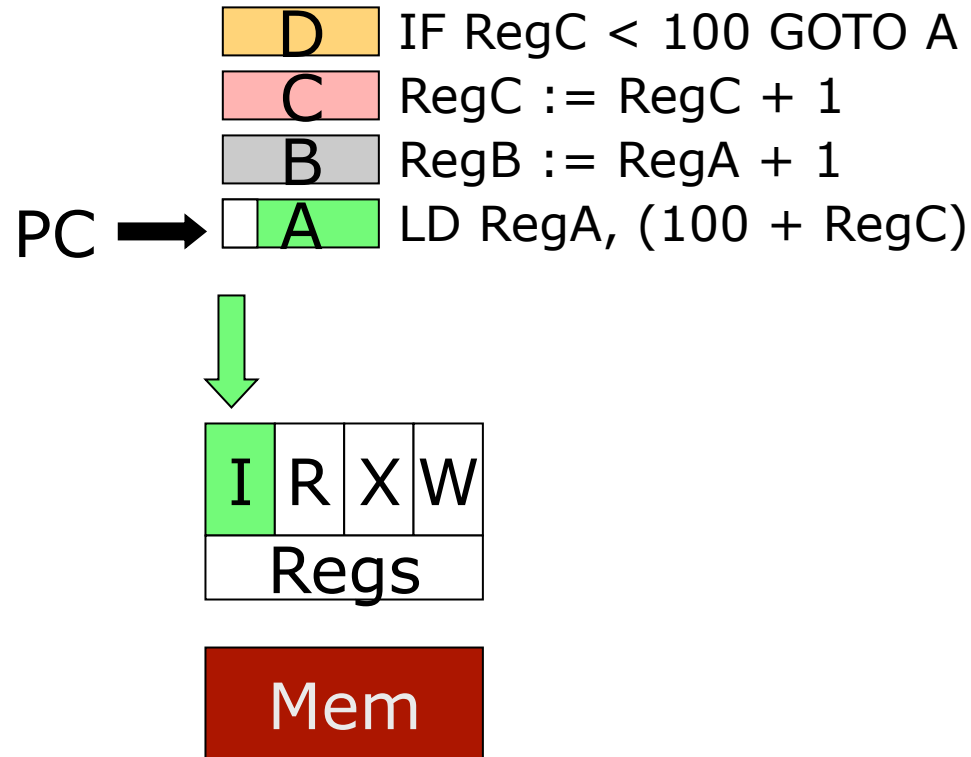
 LD RegA, (100 + RegC)

I	R	X	W
Regs			

Branch:
Write addr. of "A" to PC

Mem

Cycle 8

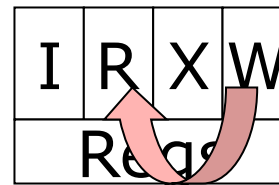


**First example of "Instruction Level Parallelism" (ILP)
-- working on several instructions simultaneously.
Improves "throughput" of instructions, not latency**

Data dependency ☹️

Previous execution example wrong!

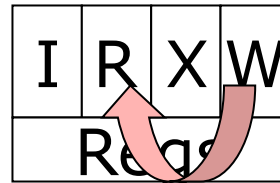
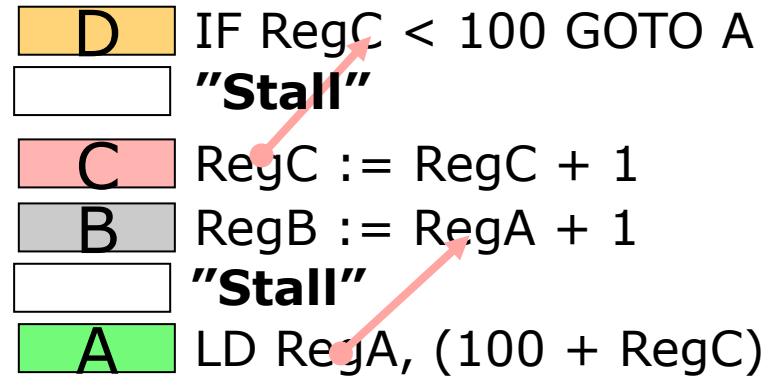
D	IF RegC < 100 GOTO A
C	RegC := RegC + 1
B	RegB := RegA + 1
A	LD RegA, (100 + RegC)



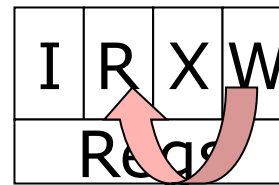
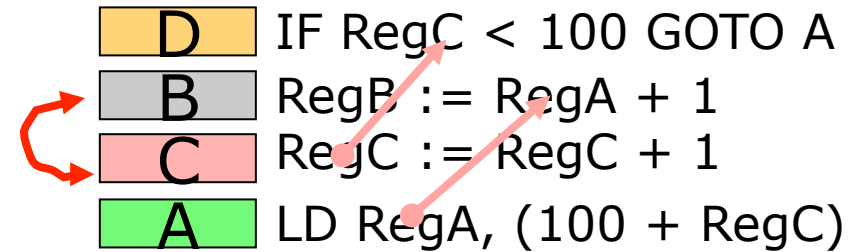
Mem



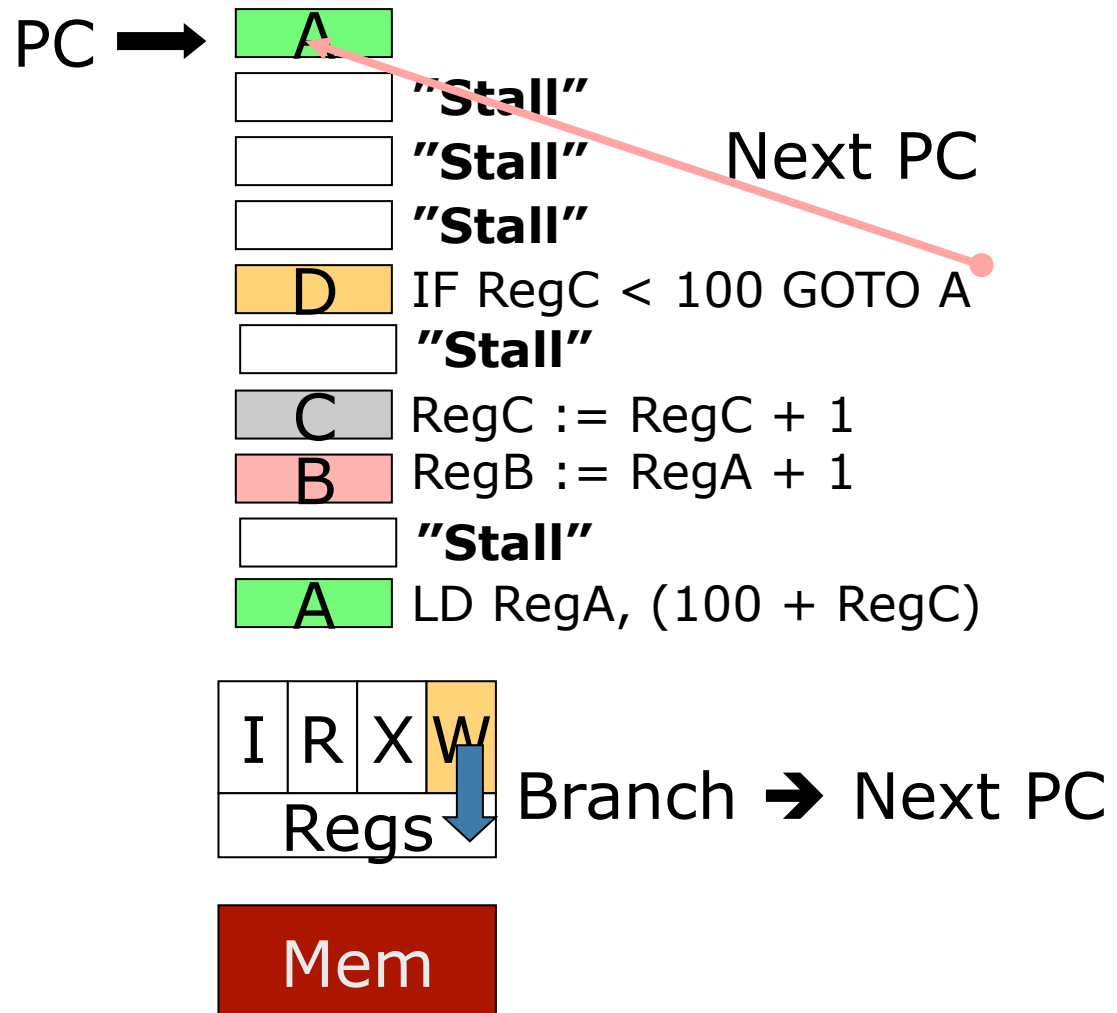
Data dependency fix 1: pipeline stalls



Data dependency fix 2: Compiler instruction scheduling



Branch delays ☹️



9 cycles per iteration of 4 instructions ☹️
 Need longer "basic blocks" with independent instr.

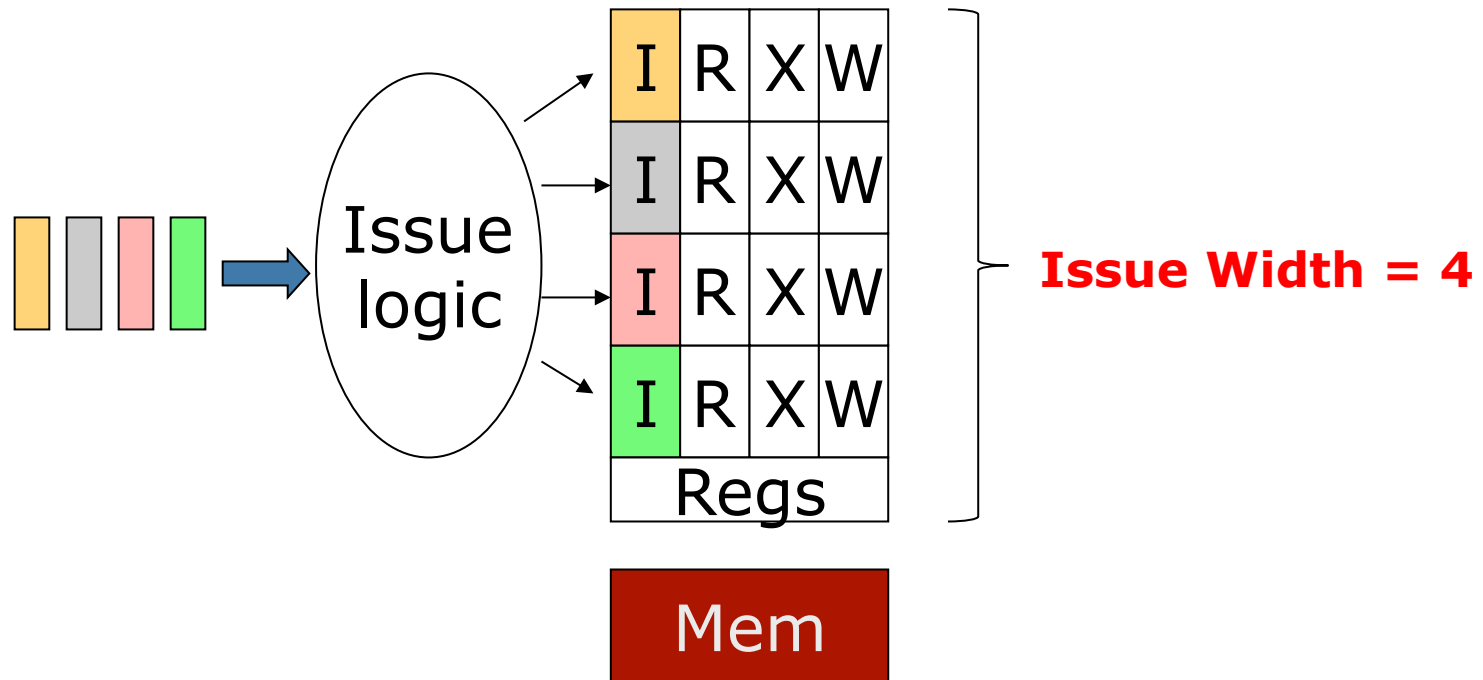
Pipeline Challenges

- Balance the pipeline stages
- Setup and hold time overhead
- Minimize pipeline stalls
- Predict and perform speculative work
- Undo speculative work



It is actually a lot worse!

Modern CPUs: "superscalars" with ≈ 4 parallel pipelines

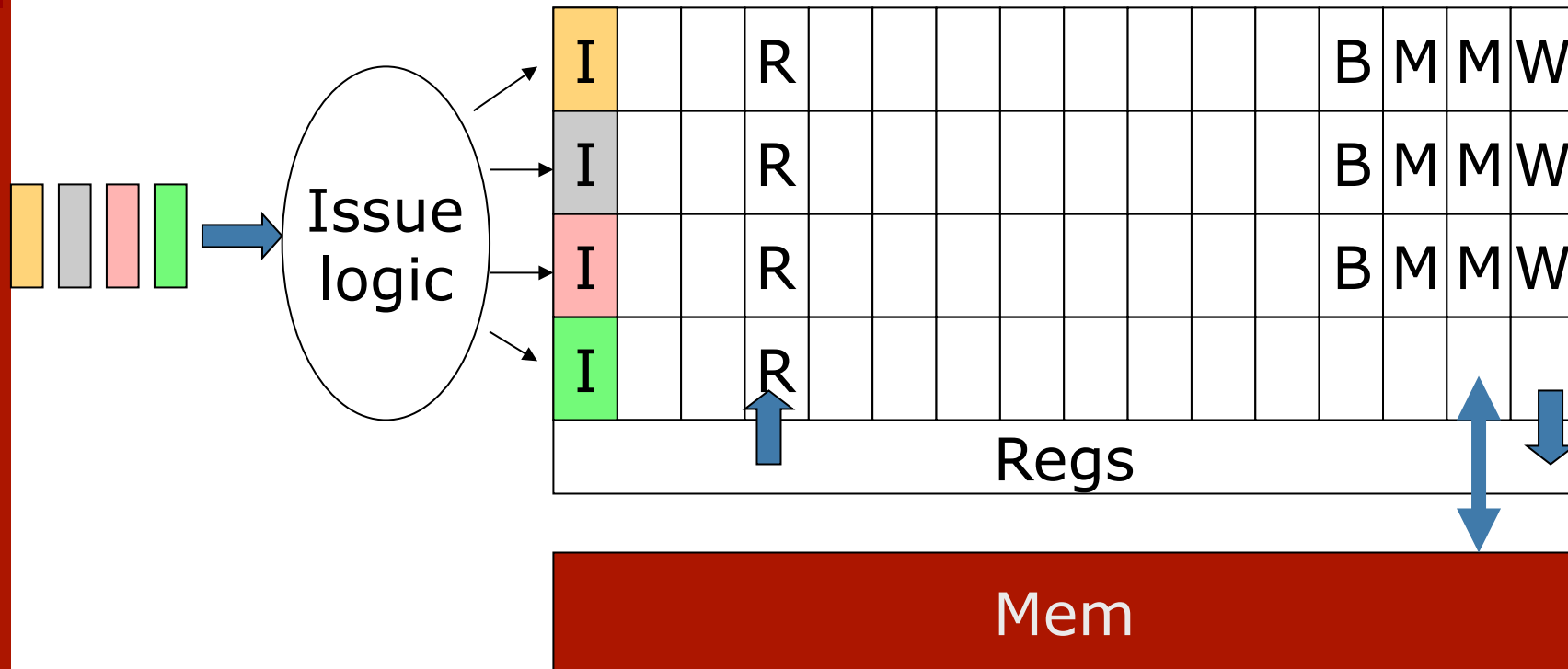


- +More ILP \rightarrow Higher throughput
- More complicated architecture
- Branch delay more expensive (more instr. wasted)
- Harder to find "enough" independent instr. (need 4 instr. between write and use)



It is actually a lot worse!

Modern CPUs: $\sim 10-20$ stages/pipe

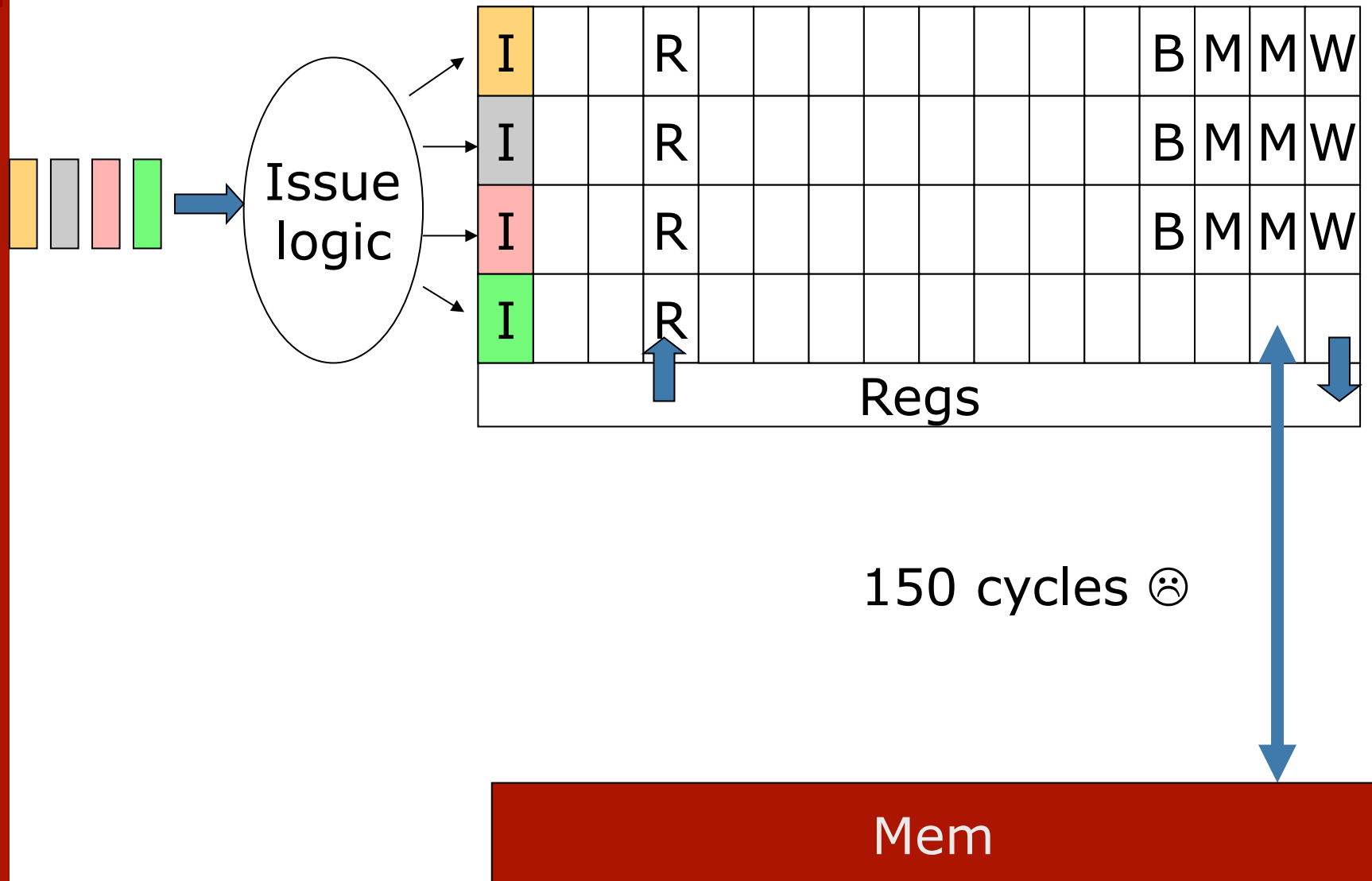


- +Explores even more ILP
- +Shorter cycletime (higher frequency)
- Branch delay even more expensive
- Even harder to find "enough" independent instr.



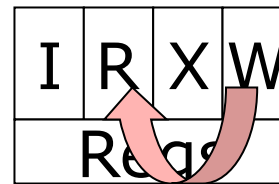
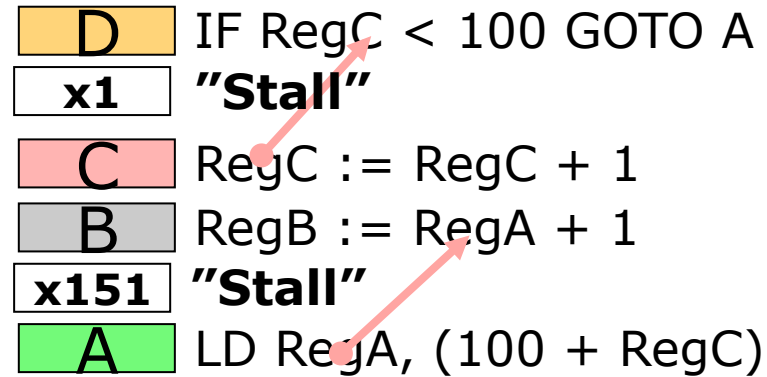
It is actually a lot worse!

DRAM access: ~150 CPU cycles





Pipeline delays gets worse





Fix: Use a cache

