



Simple Instruction Pipelining

Krste Asanovic
Laboratory for Computer Science
M.I.T.

<http://www.csg.lcs.mit.edu/6.823>



Processor Performance Equation

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

- Instructions per program depends on source code, compiler technology, and ISA
- Microcoded DLX from last lecture had cycles per instruction (CPI) of around 7 *minimum*
- Time per cycle for microcoded DLX fixed by microcode cycle time
 - mostly ROM access + next μ PC select logic



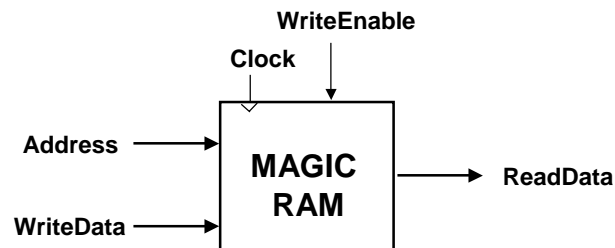
Pipelined DLX

To pipeline DLX:

- First build unpipelined DLX with $CPI=1$
- Next, add pipeline registers to reduce cycle time while maintaining $CPI=1$



A Simple Memory Model



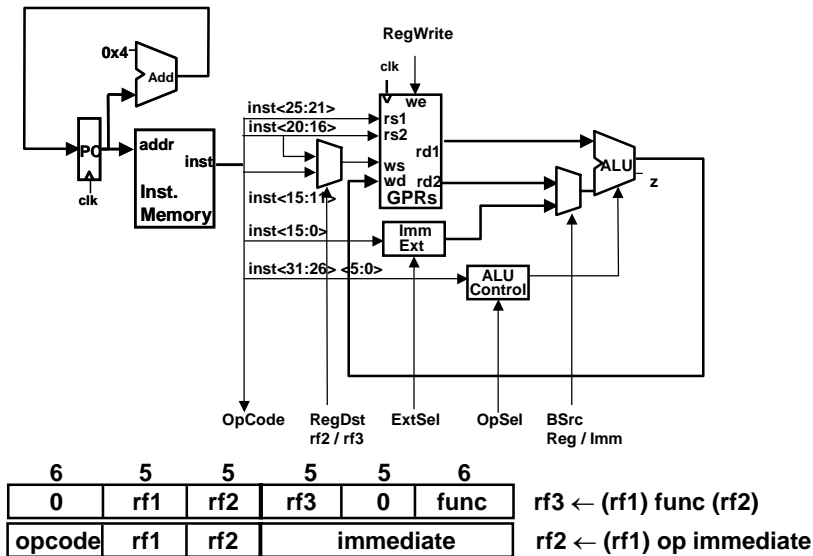
Reads and writes are always completed in one cycle

- a Read can be done any time (i.e. combinational)
- a Write is performed at the rising clock edge if it is enabled

⇒ *the write address and data must be stable at the clock edge*



Datapath for ALU Instructions



Datapath for Memory Instructions

Should program and data memory be separate?

Harvard style: separate

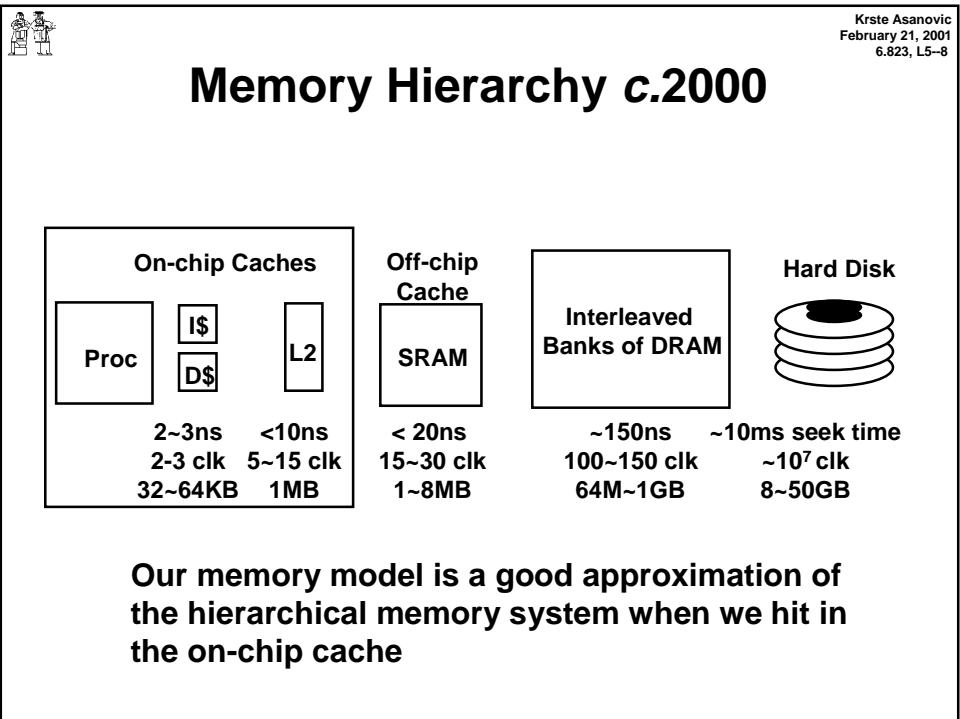
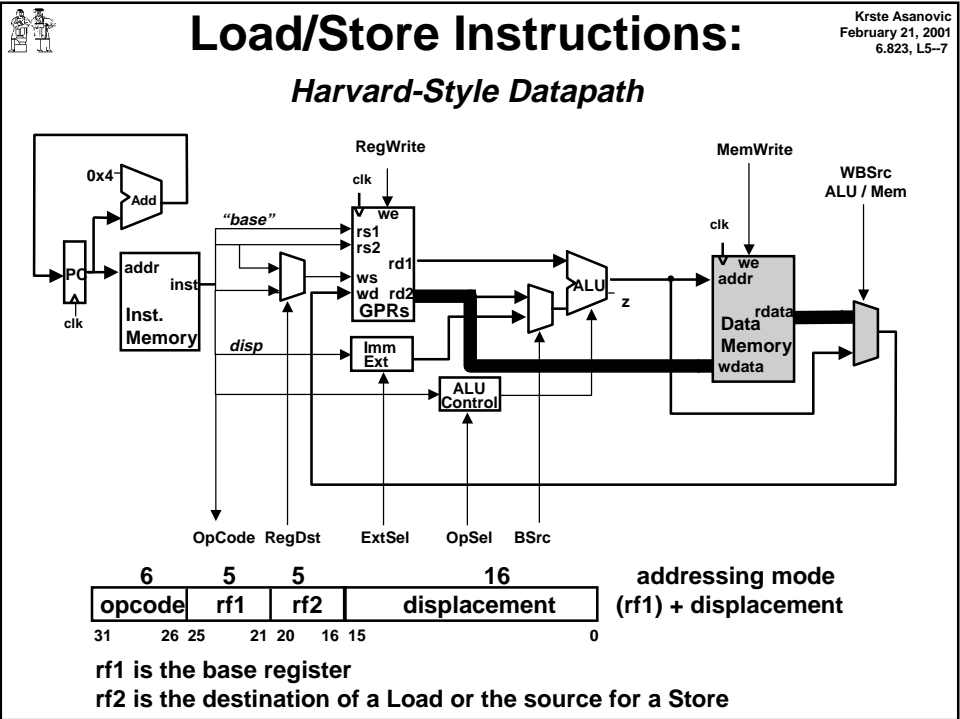
(Aiken and Mark 1 influence)

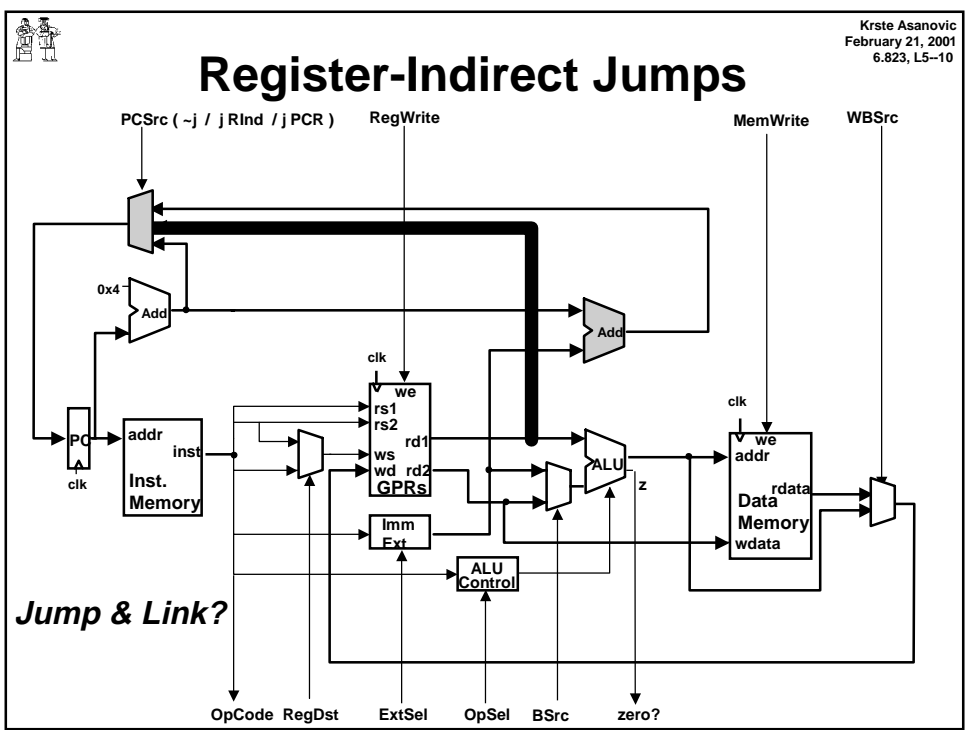
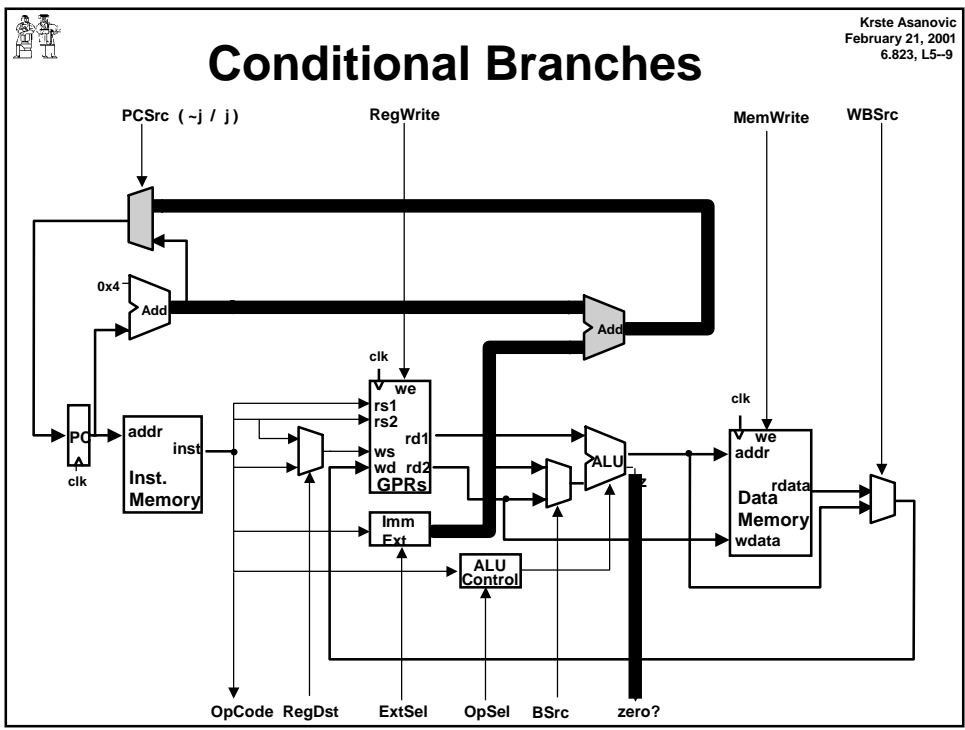
- read-only program memory
 - read/write data memory
- at some level the two memories have to be the same

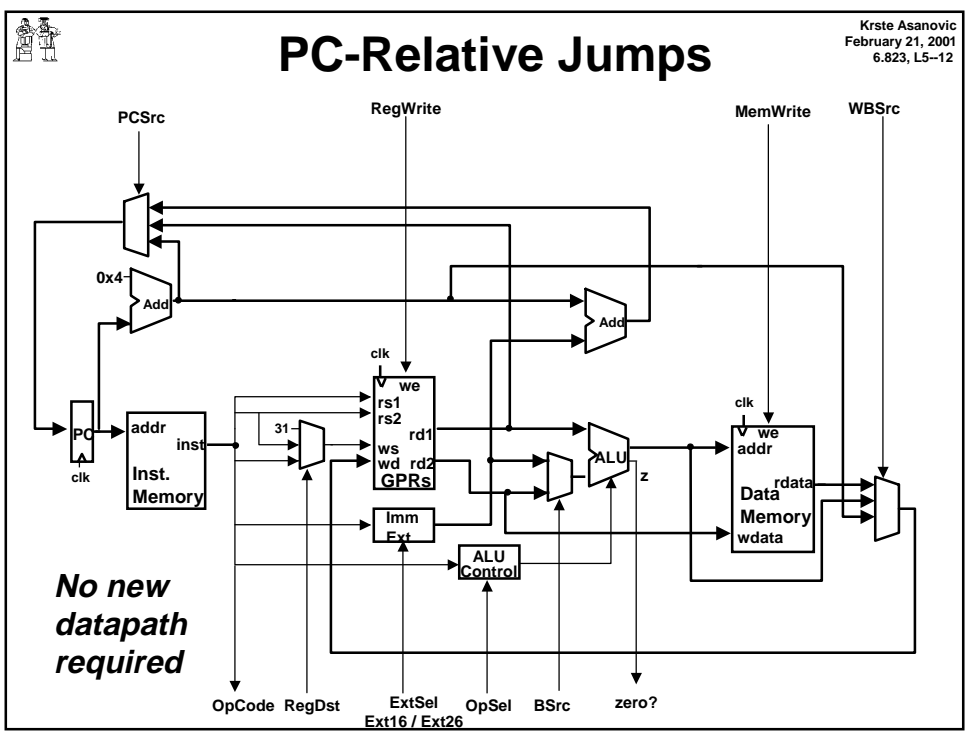
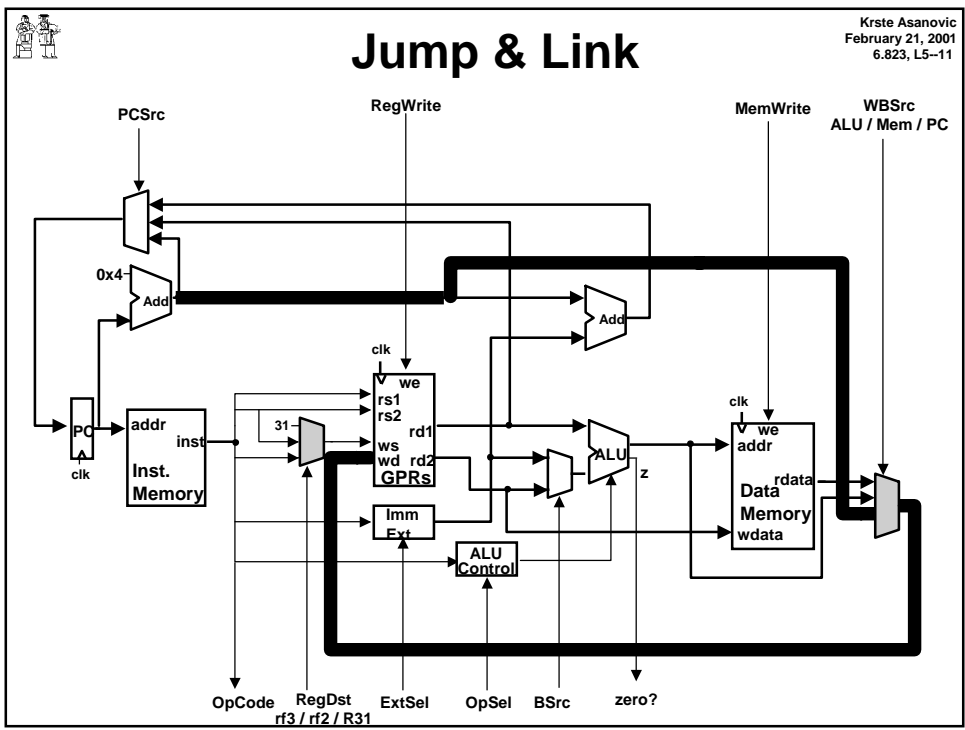
Princeton style: the same

(von Neumann's influence)

- A Load or Store instruction requires accessing the memory more than once during its execution









Single-Cycle Hardwired Control

We will assume

- clock period is sufficiently long for all of the following steps to be “completed”:

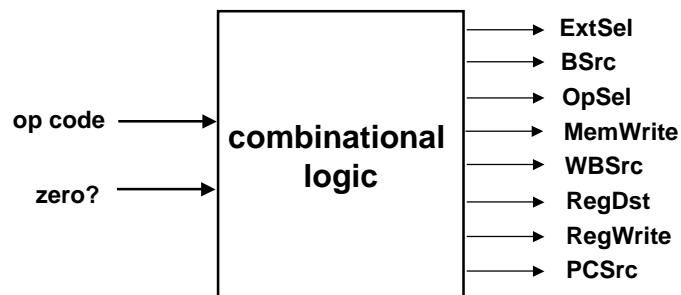
1. instruction fetch
2. decode and register fetch
3. ALU operation
4. data fetch if required
5. register write-back setup time

$$\Rightarrow t_C > t_{IFetch} + t_{RFetch} + t_{ALU} + t_{DMem} + t_{RWB}$$

- At the rising edge of the following clock, the PC, the register file and the memory are updated

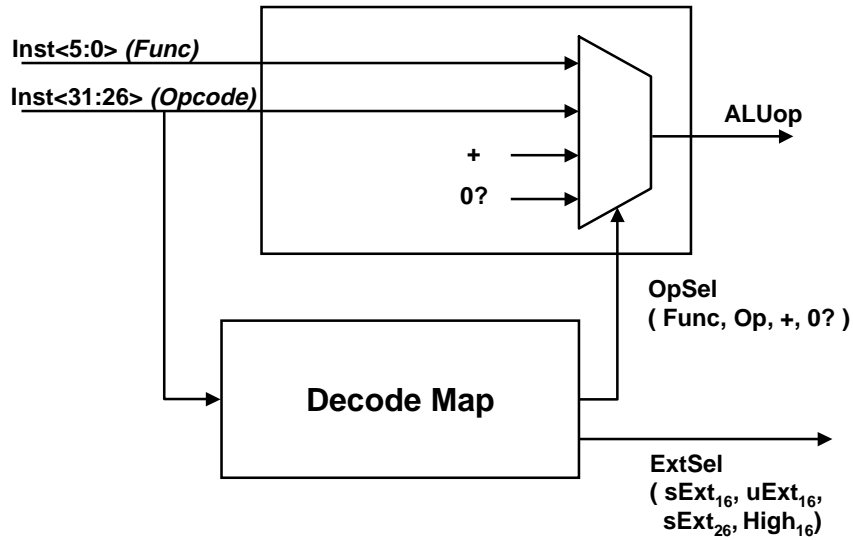


Hardwired Control is pure Combinational Logic

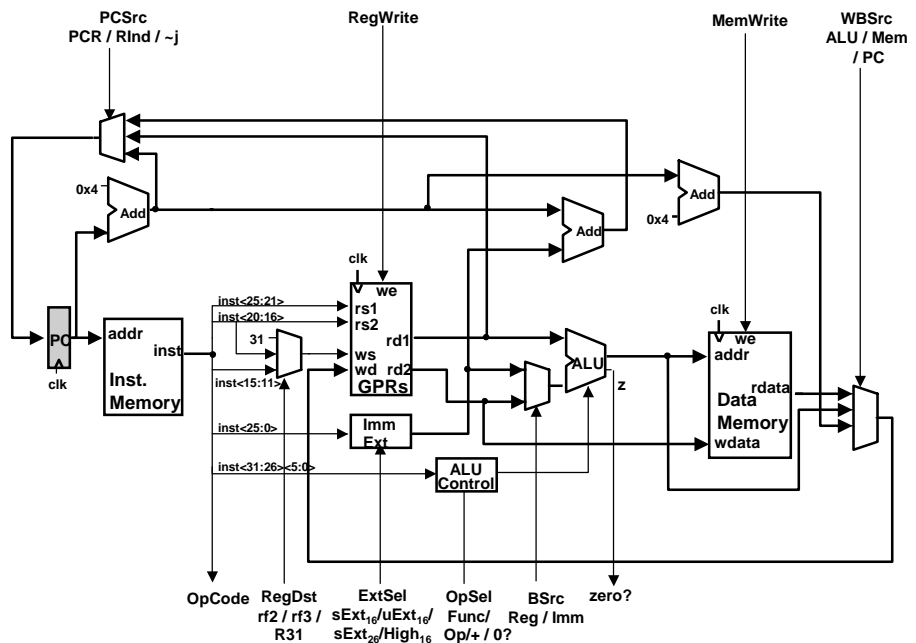




ALU Control & Immediate Extension



Hardwired Control *worksheet*





Hardwired Control Table

	Ext Sel	B Src	Op Sel	Mem Write	Reg Write	WB Src	Reg Dst	PC Src
ALU ALUu								
ALUi ALUui								
LW SW								
BEQZ _{taken} BEQZ _{-taken}								
J JAL								
JR JALR								

BSrc = Reg / Imm WBSrc = ALU / Mem / PC RegDst = rf2 / rf3 / R31
 PCSrc = PCR/RInd / ~j



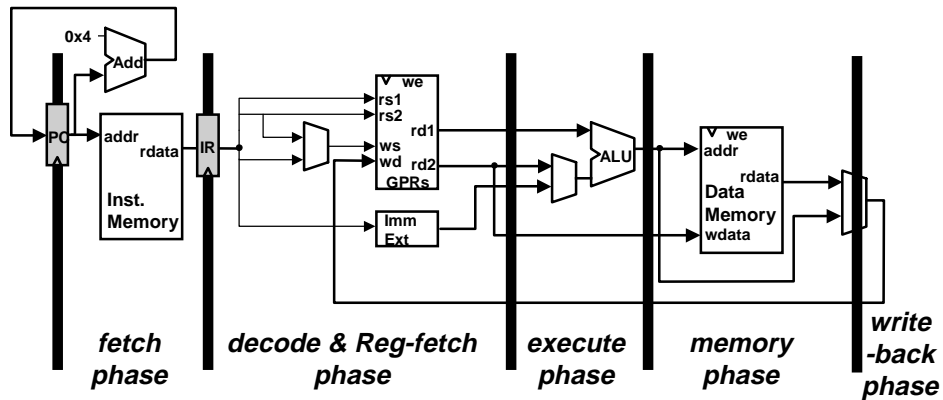
Hardwired Control Table: *Harvard DLX*

	Ext Sel	B Src	Op Sel	Mem Write	Reg Write	WB Src	Reg Dst	PC Src
ALU ALUu	* *	Reg Reg	Func Func	no no	yes yes	ALU ALU	rf3 rf3	~j ~j
ALUi ALUui	sExt ₁₆ uExt ₁₆	Imm Imm	Op Op	no no	yes yes	ALU ALU	rf2 rf2	~j ~j
LW SW	sExt ₁₆ sExt ₁₆	Imm Imm	+ +	no yes	yes no	Mem *	rf2 *	~j ~j
BEQZ _{zero?=1} BEQZ _{zero?=0}	sExt ₁₆ sExt ₁₆	* *	0? 0?	no no	no no	* *	* *	PCR ~j
J JAL	sExt ₂₆ sExt ₂₆	* *	* *	no no	no yes	* PC	* R31	PCR PCR
JR JALR	* *	* *	* *	no no	no yes	* PC	* R31	RInd RInd

BSrc = Reg / Imm WBSrc = ALU / Mem / PC RegDst = rf2 / rf3 / R31
 PCSrc1 = j / ~j PCSrc2 = PCR / RInd



Pipelined DLX Datapath



Clock period can be reduced by dividing the execution of an instruction into multiple cycles

$$t_C > \max \{t_{IM}, t_{RF}, t_{ALU}, t_{DM}, t_{RW}\} = t_{DM} \text{ (probably)}$$

However, CPI will increase unless instructions are pipelined



How to divide the datapath into stages

Suppose memory is significantly slower than other stages. In particular, suppose

$$t_{IM} = t_{DM} = 10 \text{ units}$$

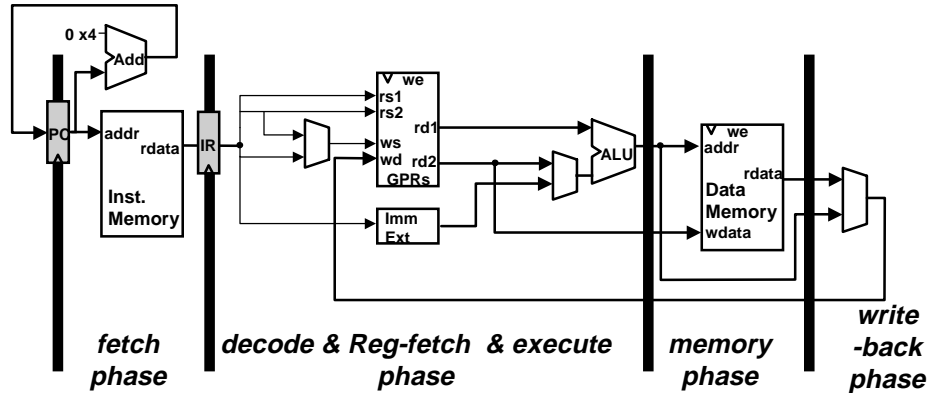
$$t_{ALU} = 5 \text{ units}$$

$$t_{RF} = t_{RW} = 1 \text{ unit}$$

Since the slowest stage determines the clock, it may be possible to combine some stages without any loss of performance



Minimizing Critical Path



$$t_C > \max \{t_{IM}, t_{RF} + t_{ALU}, t_{DM}, t_{RW}\}$$

Write-back stage takes much less time than other stages.
Suppose we combined it with the memory phase

⇒ increase the critical path by 10%



Maximum Speedup by Pipelining

For the 4-stage pipeline, given

$t_{IM} = t_{DM} = 10$ units, $t_{ALU} = 5$ units, $t_{RF} = t_{RW} = 1$ unit
 t_C could be reduced from 27 units to 10 units
 ⇒ speedup = 2.7

However, if $t_{IM} = t_{DM} = t_{ALU} = t_{RF} = t_{RW} = 5$ units

The same 4-stage pipeline can reduce t_C from 25 units to 10 units
 ⇒ speedup = 2.5

But, since $t_{IM} = t_{DM} = t_{ALU} = t_{RF} = t_{RW}$, it is possible to achieve higher speedup with more stages in the pipeline.

A 5-stage pipeline can reduce t_C from 25 units to 5 units
 ⇒ speedup = 5



Technology Assumptions

We will assume

- A small amount of very fast memory (caches) backed up by a large, slower memory
- Fast ALU (at least for integers)
- Multiported Register files (slower!).

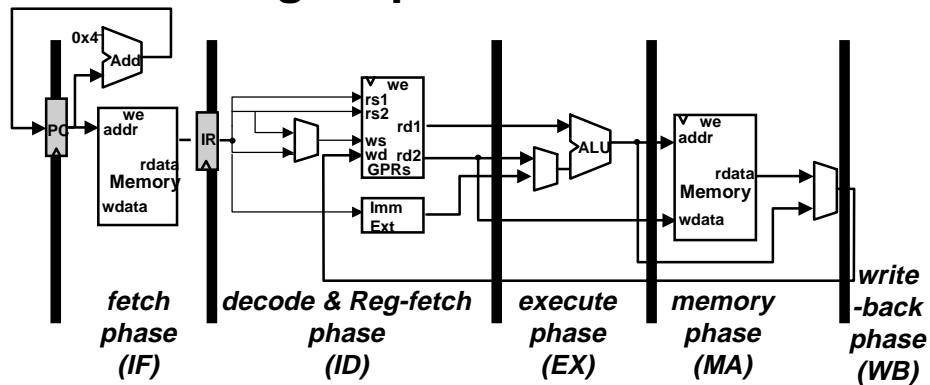
It makes the following timing assumption valid

$$t_{IM} \approx t_{RF} \approx t_{ALU} \approx t_{DM} \approx t_{RW}$$

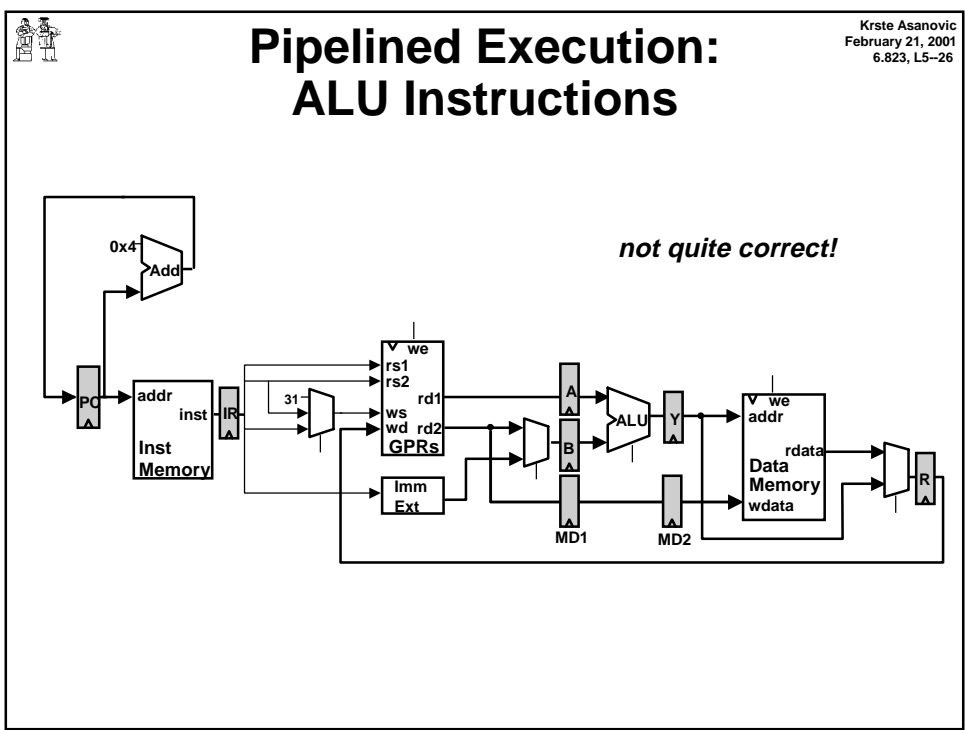
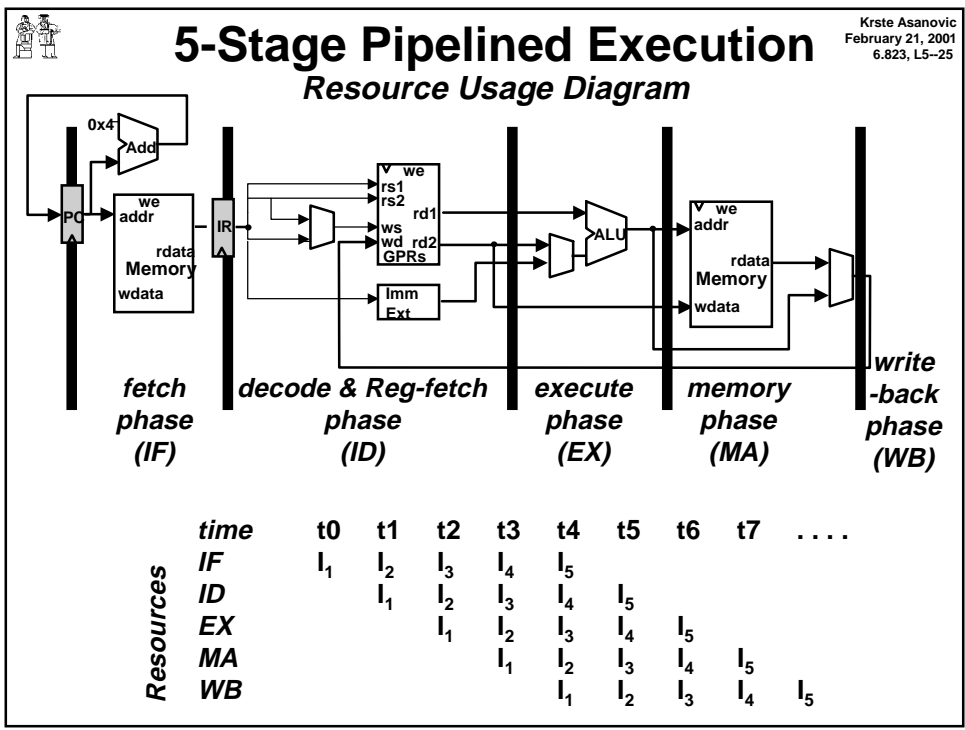
A 5-stage pipelined Harvard-style architecture will be the focus of our detailed design

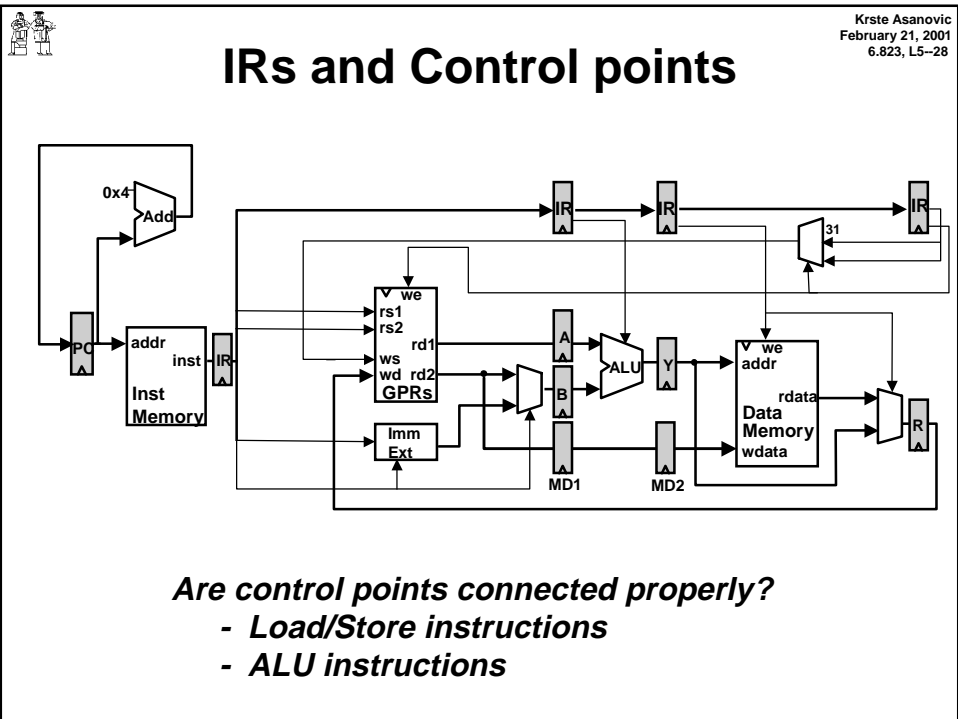
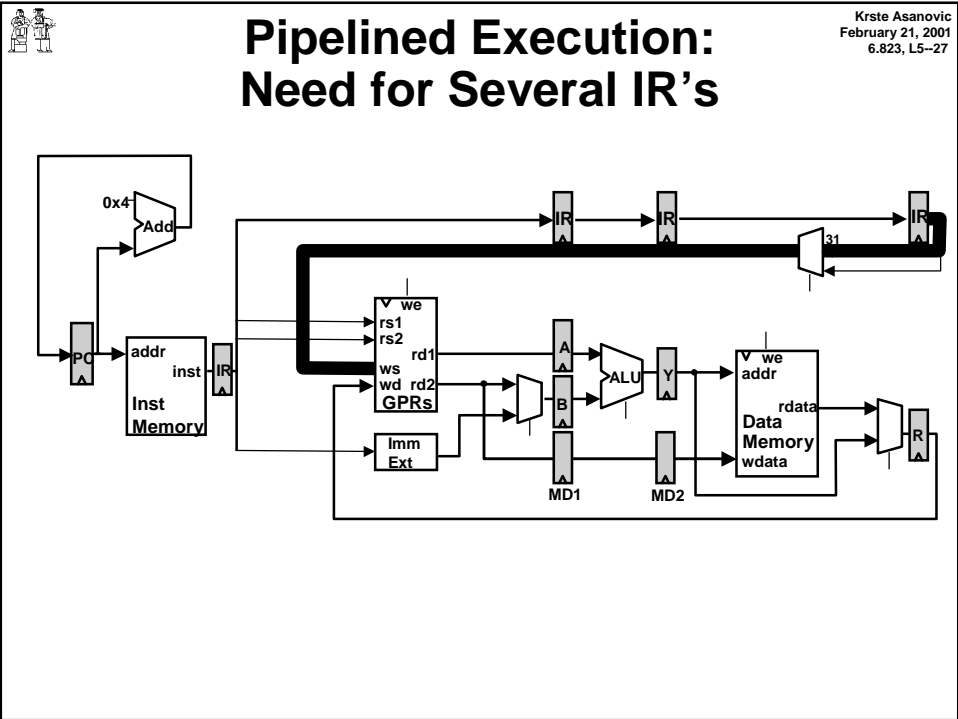


5-Stage Pipelined Execution



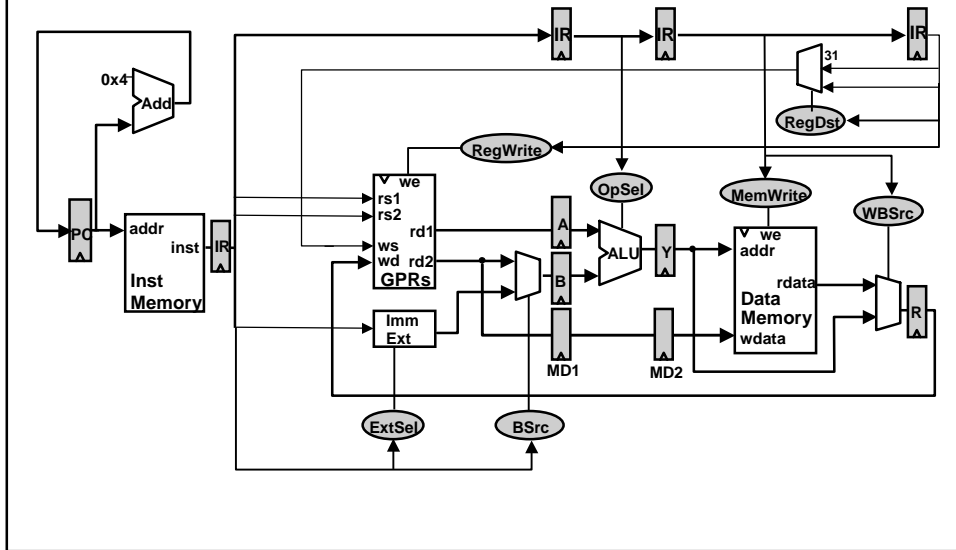
time	t0	t1	t2	t3	t4	t5	t6	t7
instruction1	IF ₁	ID ₁	EX ₁	MA ₁	WB ₁				
instruction2		IF ₂	ID ₂	EX ₂	MA ₂	WB ₂			
instruction3			IF ₃	ID ₃	EX ₃	MA ₃	WB ₃		
instruction4				IF ₄	ID ₄	EX ₄	MA ₄	WB ₄	
instruction5					IF ₅	ID ₅	EX ₅	MA ₅	WB ₅



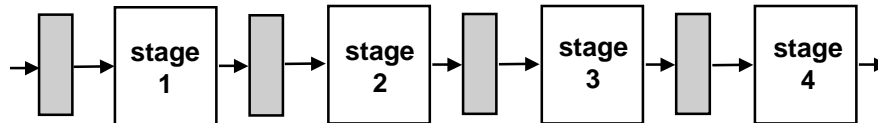




Pipelined DLX Datapath *without jumps*



An Ideal Pipeline



- All objects go through the same stages
 - No sharing of resources between any two stages
 - Propagation delay through all pipeline stages is equal
 - The scheduling of an object entering the pipeline is not affected by the objects in other stages
- These conditions generally hold for industrial assembly lines. An instruction pipeline, however, cannot satisfy the last condition. Why?*

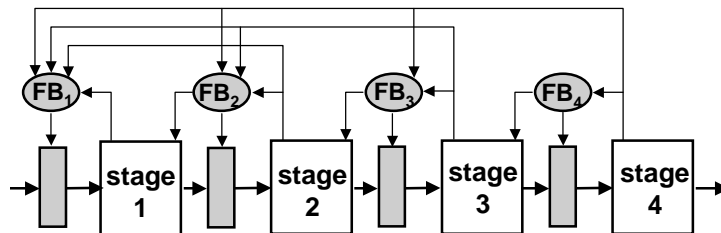


How Instructions can Interact with each other in a pipeline

- An instruction in the pipeline may need a resource being used by another instruction in the pipeline
structural hazard
- An instruction may produce data that is needed by a later instruction
data hazard
- In the extreme case, an instruction may determine the next instruction to be executed
control hazard (branches, interrupts,...)



Feedback to Resolve Hazards



Controlling pipeline in this manner works provided *the instruction at stage $i+1$ can complete without any interference from instructions in stages 1 to i* (otherwise deadlocks may occur)

Feedback to previous stages is used to *stall or kill instructions*