

# ECE 485/585

## Computer Organization and Design

### Lecture 12: Instruction Level Parallelism

Fall 2022

Won-Jae Yi, Ph.D.

Department of Electrical and Computer Engineering  
Illinois Institute of Technology

- Speculation – “Guess”** *an approach where by the compiler or processor guesses the outcome of an instruction to resolve it as a dependence in executing other instructions*
- “Guess” what to do with an instruction
    - Start operation as soon as possible
    - Check whether guess was right
      - If so, complete the operation
      - If not, roll-back and do the right thing
  - Common to static and dynamic multiple issue
  - Examples
    - Speculate on branch outcome
      - Roll back if path taken is different
    - Speculate on load
      - Roll back if location is updated

### Instruction-Level Parallelism (ILP)

- Pipelining: executing multiple instructions in parallel
- To increase ILP
  - ① Deeper pipeline  $5 \rightarrow 10$  stages
    - Less work per stage  $\Rightarrow$  shorter clock cycle
  - ② Multiple-issue processor
    - Replicate pipeline stages  $\Rightarrow$  multiple pipelines
    - Start multiple instructions per clock cycle
      - CPI < 1, so use Instructions Per Cycle (IPC)
      - e.g., 4GHz 4-way multiple-issue
        - 16 BIPS, peak CPI = 0.25, peak IPC = 4
      - But dependencies reduce this in practice

### Multiple Issue Processor

#### Static multiple issue

- Compiler groups instructions to be issued together
- Packages them into “issue slots”
- Compiler detects and avoids hazards

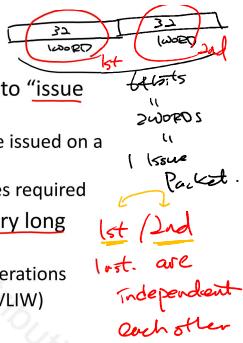
→ look for dependencies

#### Dynamic multiple issue

- CPU examines instruction stream and chooses instructions to issue each cycle
- Compiler can help by reordering instructions
- CPU resolves hazards using advanced techniques at runtime

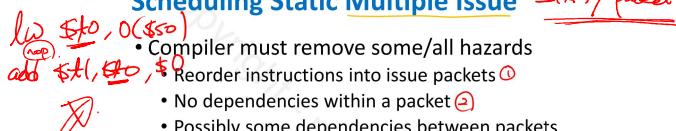
### Static Multiple Issue

- Compiler groups instructions into “issue packets”
  - Group of instructions that can be issued on a single cycle
  - Determined by pipeline resources required
- Think of an issue packet as a very long instruction
  - Specifies multiple concurrent operations  $\Rightarrow$  Very Long Instruction Word (VLIW)



### Scheduling Static Multiple Issue

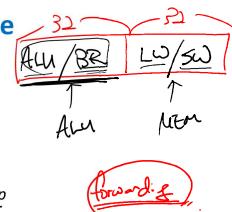
- Compiler must remove some/all hazards
  - Reorder instructions into issue packets ①
    - No dependencies within a packet ②
    - Possibly some dependencies between packets
      - Varies between ISAs; compiler must know!
    - Pad with nop if necessary



### MIPS with Static Dual Issue

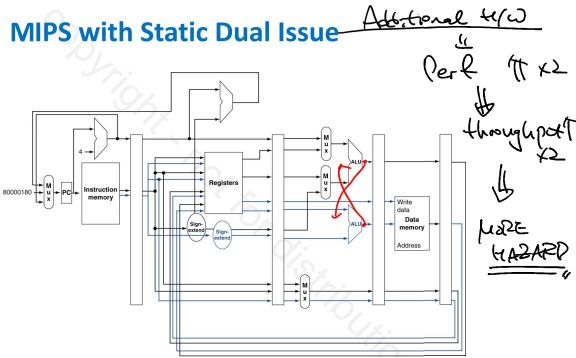
#### Two-issue packets

- One ALU/branch instruction
- One load/store instruction
- 64-bit aligned
  - ALU/branch, then load/store
- Pad an unused instruction with nop



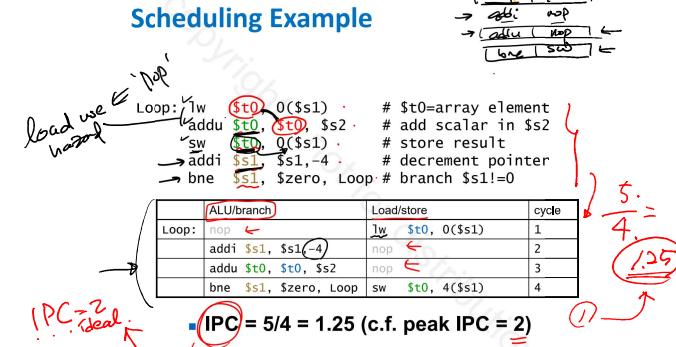
forwarding

Address	Instruction type	Pipeline Stages				
		IF	ID	EX	MEM	WB
n	ALU/branch					
n + 4	Load/store	IF	ID	EX	MEM	WB
n + 8	ALU/branch					
n + 12	Load/store	IF	ID	EX	MEM	WB
n + 16	ALU/branch			IF	ID	EX
n + 20	Load/store			IF	ID	EX



## Hazards in the Dual-Issue MIPS

- More instructions executing in parallel
- ① EX data hazard**
  - Forwarding avoided stalls with single-issue
  - Now can't use ALU result in load/store in same packet
    - add \$t0, \$s0, \$s1 → load \$s2, 0(\$s1)
    - load \$s2, 0(\$s1) → add \$t0, \$s0, \$s1
    - Split into two packets, effectively a stall
- ② Load-use hazard**
  - Still one cycle use latency, but now two instructions
  - More aggressive scheduling required



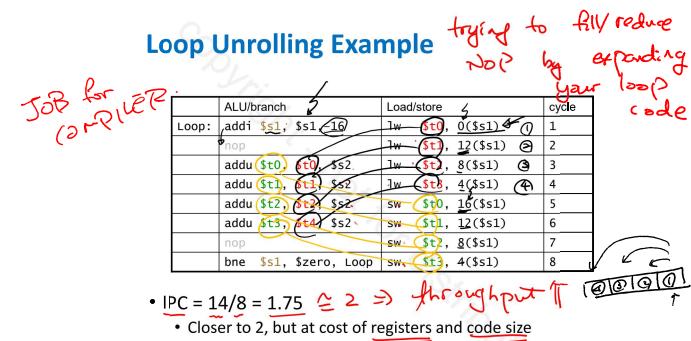
## Loop Unrolling

- Replicate loop body to expose more parallelism
  - Reduces loop-control overhead
- Use different registers per replication
  - Called "register renaming"
  - Avoid loop-carried "anti-dependencies"
    - Store followed by a load of the same register
    - a.k.a. "name dependence"
      - Reuse of a register name

Loop Unrolling  
→ A technique to get more throughput from loops that access arrays in which multiple copies of the loop body are made and instructions from different iterations are scheduled together

for ( $i=0$ ;  $i < \text{cnt}$ ;  $i++$ ) → for ( $i=0$ ;  $i < \text{cnt}$ ;  $i++$ )  
 $a[i] = b[i] + c;$   
 $a[i+1] = b[i+1] + c;$   
 $a[i+2] = b[i+2] + c;$   
 $a[i+3] = b[i+3] + c;$   
 $a[i+4] = b[i+4] + c;$

13 14



## Dynamic Multiple Issue

- "Superscalar" processors
- CPU decides whether to issue 0, 1, 2, ... each cycle
  - Avoiding structural and data hazards
- Avoids the need for compiler scheduling
  - Though it may still help
  - Code semantics ensured by the CPU

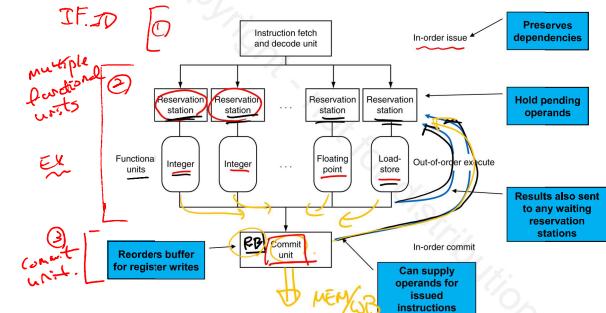
## Dynamic Pipeline Scheduling

- Allow the CPU to execute instructions out of order to avoid stalls
  - But commit result to registers in order
- Example
 

lw \$t0, 20(\$s2)  
 addu \$t1, \$t0, \$t2  
 sub \$t3, \$s4, \$t3  
 slli \$t5, \$s4, 20

Can start sub while addu is waiting for lw

## Dynamically Scheduled CPU



Reservation Station : a buffer within a functional unit that holds the operand and the operation

Commit Unit : the unit in dynamic/out-of-order execution pipeline that decides when it's safe to release the results of an operation to programmer-visible registers & memory.

Reorder Buffer : a buffer that holds results in a dynamically scheduled processor until it's safe to store the results to registers & memory.

19

## Register Renaming

- Reservation stations and reorder buffer effectively provide register renaming
- On instruction issue, instruction copied to reservation station
  - If operand is available in register file or reorder buffer
    - Copied to reservation station
    - No longer required in the register; can be overwritten
  - If operand is not yet available
    - It will be provided to the reservation station by a function unit
    - Register update may not be required

20

## Speculation

- Predict branch and continue issuing
  - Don't commit until branch outcome determined
- Load speculation
  - Avoid load and cache miss delay
    - Predict the effective address
    - Predict loaded value
    - Load before completing outstanding stores
    - Bypass stored values to load unit
  - Don't commit load until speculation cleared

21

## Why Do Dynamic Scheduling?

- Why not just let the compiler schedule code?
- Not all stalls are predictable
  - e.g., cache misses
- Can't always schedule around branches
  - Branch outcome is dynamically determined
- Different implementations of an ISA have different latencies and hazards

22

## Does Multiple Issue Work?

### The BIG Picture

- Yes, but not as much as we'd like
- Programs have real dependencies that limit ILP
- Some dependencies are hard to eliminate
  - e.g., pointer aliasing
- Some parallelism is hard to expose
  - Limited window size during instruction issue
- Memory delays and limited bandwidth
  - Hard to keep pipelines full
- Speculation can help if done well

23