

# COMP4211 Seminar

## Intro to Instruction-Level Parallelism

04S1 Week 02  
Oliver Diessel

## Overview

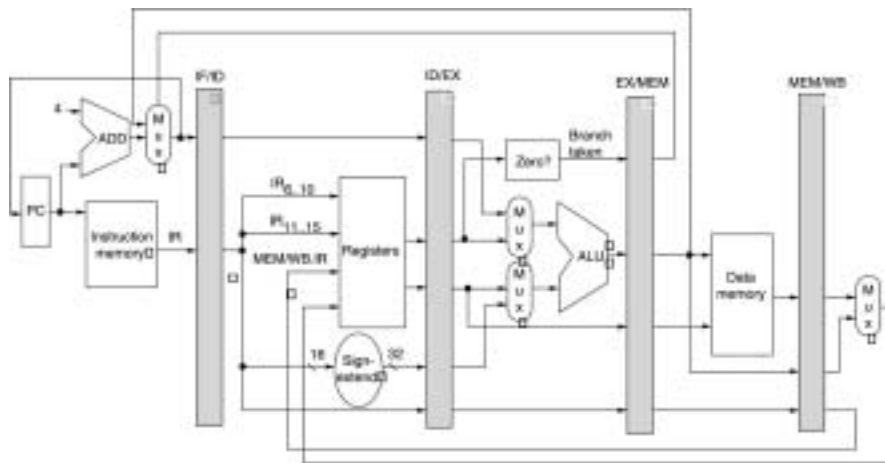
- Review pipelining from COMP3211
- Hennessy & Patterson, Chapters 3 & 4

04S1

COMP4211 Seminar

W02S2

## Five stage statically scheduled pipeline



© 2003 Elsevier Science (USA). All rights reserved.

04S1

COMP4211 Seminar

W02S3

## Pipeline characteristics

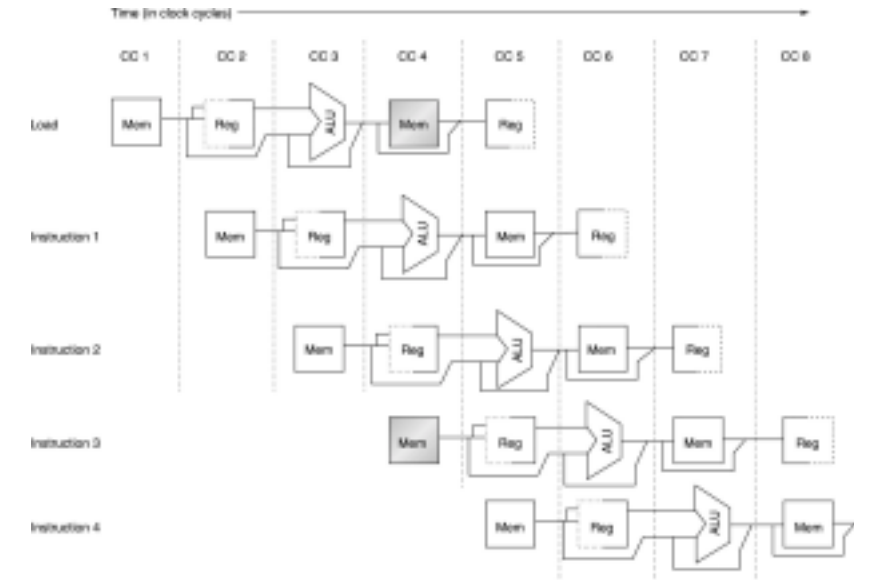
- Parallelism
- 1 instruction issued per cycle
- Reduced performance due to hazards:
  - Structural
    - E.g. single memory
  - Data
    - Use forwarding/stall
  - Control
    - Cope with hardware and software techniques

04S1

COMP4211 Seminar

W02S4

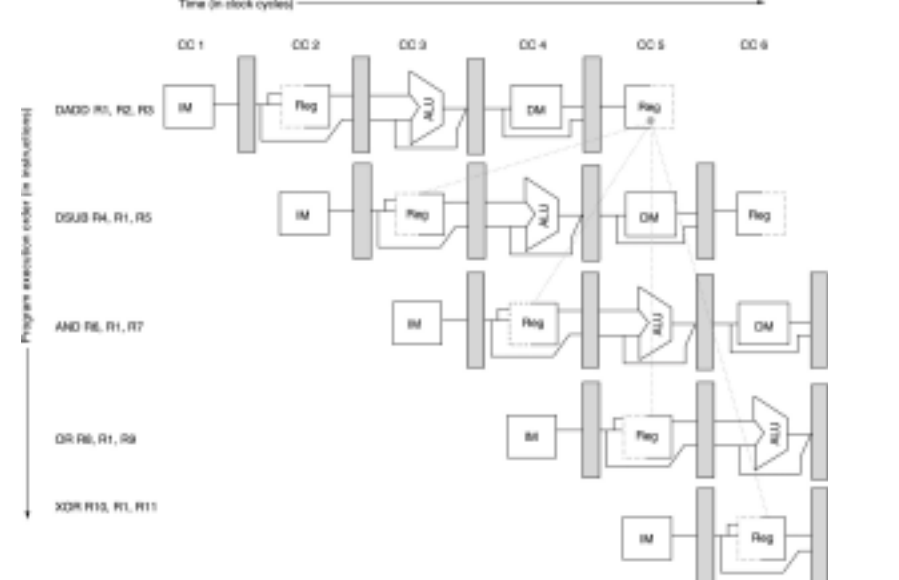
# Structural hazard example



04

S5

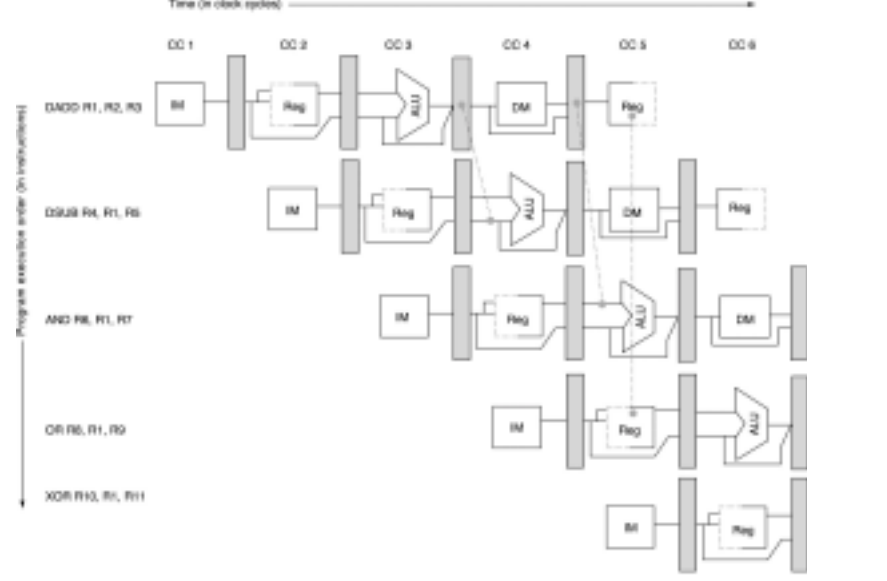
# Data hazard examples



C

V02S6

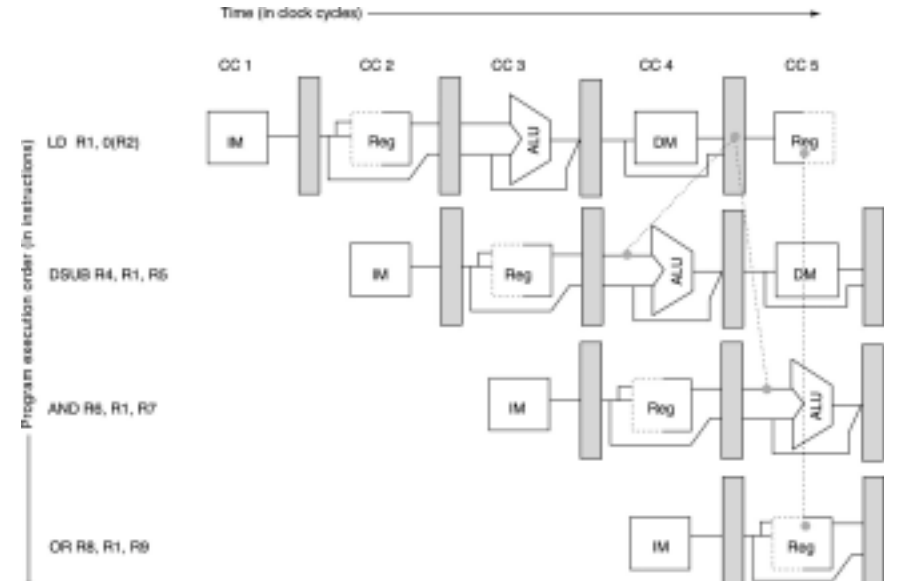
# Data hazard remedy ñ forwarding



C

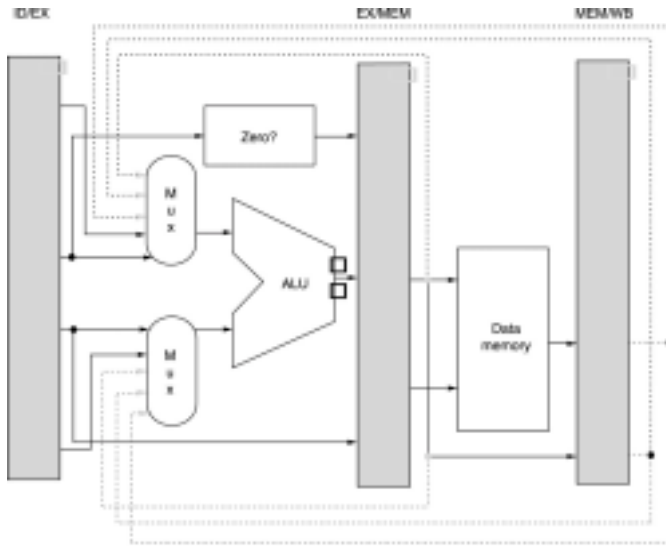
W02S7

# Data hazard needing stall



C

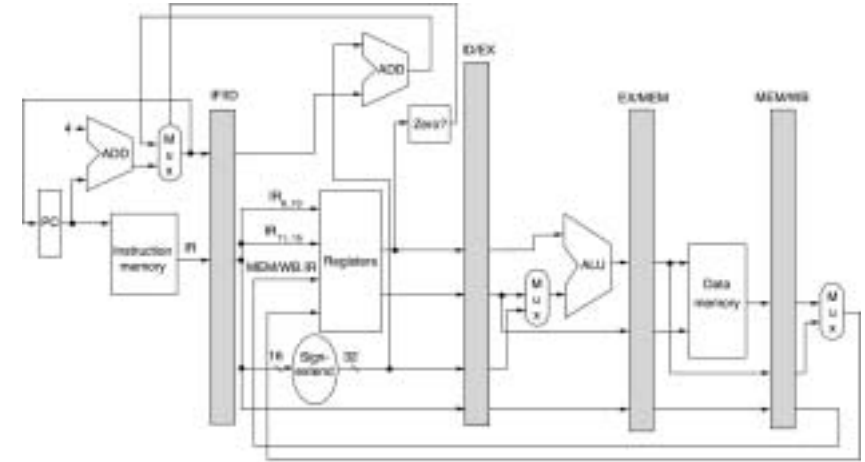
## Forwarding hardware



04S1

W02S9

## Control hazard ñ hardware amelioration



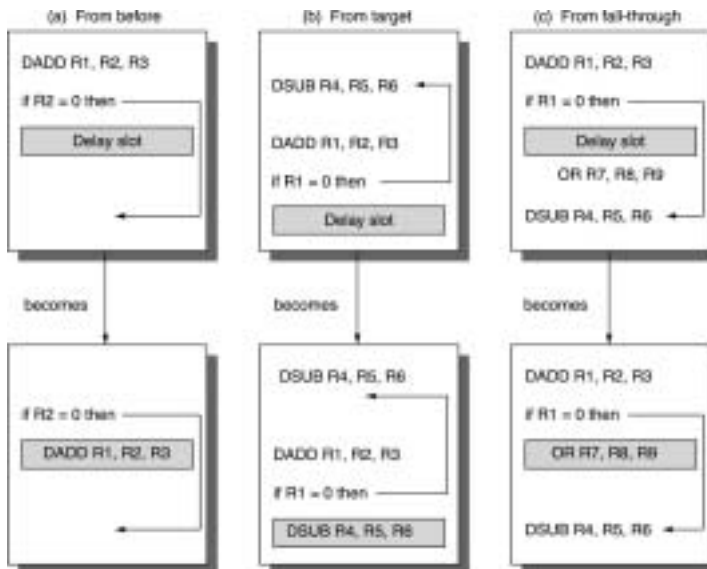
© 2003 Elsevier Science (USA). All rights reserved.

04S1

COMP4211 Seminar

W02S10

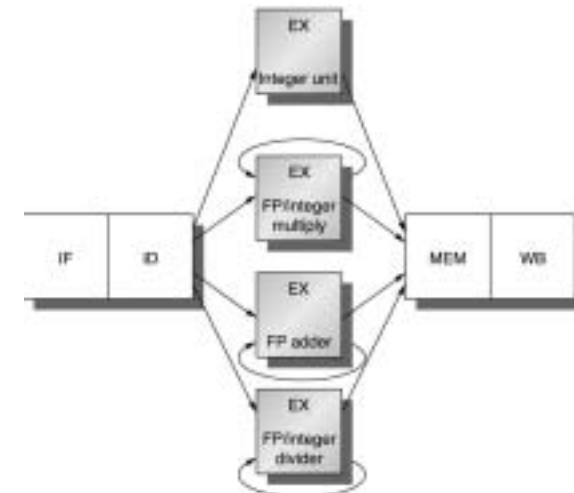
## Control stalls ñ software amelioration



04S1

W02S11

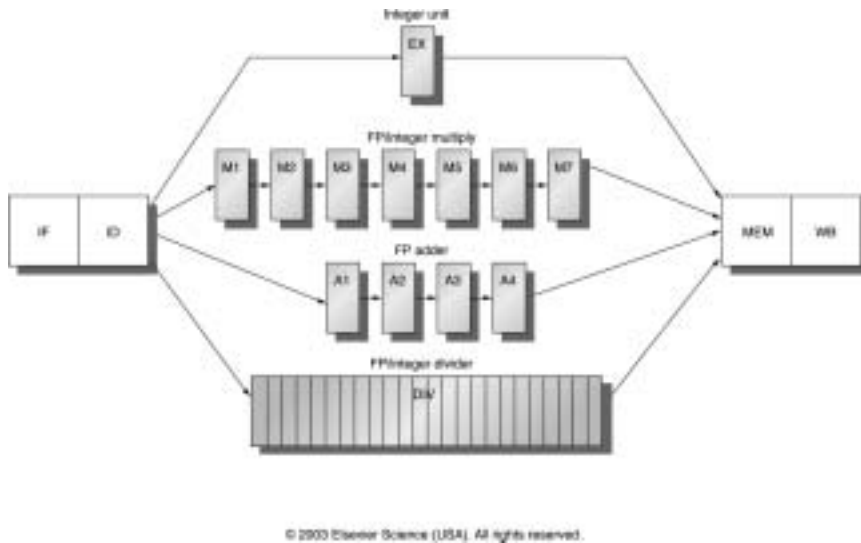
## Extending MIPS to handle FP operations



© 2003 Elsevier Science (USA). All rights reserved.

04S1

W02S12



04S1

COMP4211 Seminar

W02S13

## Instruction-level parallelism

- ï Pipelining commonly used since 1985 to overlap the execution & improve performance ñsince instructions evaluated in parallel, known as *instruction-level parallelism (ILP)*
- ï Here we look at extending pipelining ideas by increasing the amount of parallelism exploited among instructions
- ï Start by looking at limitation imposed by data & control hazards, then look at increasing the ability of the processor to exploit parallelism

04S1

COMP4211 Seminar

W02S14

## Two main approaches

- ï Two largely separable approaches to exploiting ILP:
  - ñ Dynamic techniques (Ch. 3) depend upon hardware to locate parallelism
  - ñ Static techniques (Ch. 4) rely much more on software
- ï Practical implementations typically involve a mix or some crossover of these approaches
- ï Dynamic, hardware-intensive approaches dominate the desktop and server markets; examples include Pentium, Power PC, and Alpha
- ï Static, compiler-intensive approaches have seen broader adoption in the embedded market, except, for example, IA-64 and Itanium

04S1

COMP4211 Seminar

W02S15

- ï What are the features of programs & processors that limit the amount of parallelism that can be exploited among instructions?
- ï How are programs mapped to hardware?
- ï Will a program property limit performance? If so, when?
- ï Recall

$$\text{Pipeline CPI} = \text{Ideal CPI} + \text{Structural stalls} + \text{Data hazard stalls} + \text{Control stalls}$$

- ï In order to increase *instructions/cycle (IPC)* we need to pay increasing attention to dealing with stalls

04S1

COMP4211 Seminar

W02S16

- The amount of parallelism available within a *basic block* ñ a straight-line code sequence with no branches in or out except to the entry and from the exit ñ is quite small
- Typical dynamic branch frequency is often between 15% and 25% ñ between 4 and 7 instructions execute between branch pairs ñ these instructions are likely to depend upon each other, and thus the overlap we can exploit within a basic block is typically less than the average block size
- To obtain substantial performance enhancements, we must exploit ILP across multiple basic blocks

## Loop-level parallelism

- *Loop-level parallelism* increases the amount of parallelism available among iterations of a loop
- In the code:
 

```
for (i=1; i<1000; i++)
    x[i] = x[i] + y[i];
```

 every iteration can overlap with any other, although within a loop iteration there is little or no opportunity for overlap
- We will examine techniques for unrolling loops to convert the loop-level parallelism to ILP
- Another approach to exploiting loop-level parallelism is to use vector instructions. While processors that exploit ILP have almost totally replaced vector processors, vector instruction sets may see a renaissance for use in graphics, DSP, and multimedia

## Data dependences and hazards

- In order to exploit ILP we must determine which instructions can be executed in parallel
- *Parallel* instructions can execute simultaneously without causing stalls assuming there are no structural hazards (sufficient resources)
- *Dependent* instructions are not parallel and must be executed in order, although they may be often be partially overlapped
- Three types of dependences exist:
  - ñ *Data* dependences
  - ñ *Name* dependences
  - ñ *Control* dependences

## Data dependence

- Instruction *j* is data dependent on instruction *i* if:
  - ñ *i* produces a result that may be used by *j*, or
  - ñ *j* is data dependent on instruction *k*, and *k* is data dependent on instruction *i*.
- Example
 

```
Loop:   L.D   F0,0(R1)    ;F0 = array element
        ADD.D F4,F0,F2    ;add scalar in F2
        S.D   F4,0(R1)    ;store result
        ADDUI R1,R1,#-8   ;dec pointer 8 bytes
        BNE  R1,R2,LOOP   ;branch R1!=R2
```

 has data dependences on consecutive pairs of instructions

- ï Processors with pipeline interlocks will detect a hazard and stall if such instructions are scheduled simultaneously
- ï Compilers for processors without interlocks that rely on compiler scheduling cannot schedule dependent instructions to allow complete overlap
- ï Dependences are properties of *programs* ñ whether a dependence results in an actual hazard being detected and whether that hazard causes a stall is a property of the *pipeline organization*
- ï A dependence
  1. Indicates the possibility of a hazard;
  2. Determines the order in which results must be calculated; and
  3. Sets an upper bound on how much parallelism can possibly be exploited

04S1

COMP4211 Seminar

W02S21

- ï Since data dependences limit the amount of ILP we can exploit, we focus on how to overcome those limitations
- ï A dependence can be overcome by
  1. Maintaining the dependence but avoiding a hazard, and
  2. Eliminating the dependence by transforming the code.
- ï Here we primarily consider hardware techniques for scheduling the code dynamically as it is executed

04S1

COMP4211 Seminar

W02S22

- ï Data values may flow from instruction to instruction
  - ñ via registers (in which case dependence detection is reasonably straightforward since register names are fixed in the instructions, although intervening branches may cause correctness concerns), or
  - ñ via memory (in which case dependences are more difficult to detect because of aliasing i.e.  $100(R4) = 20(R6)$  and effective addresses such as  $20(R6)$  may change from one execution of an instruction to the next)
- ï We will examine hardware for detecting data dependences that involve memory locations and will see the limitations of these techniques
- ï Compiler techniques for detecting dependences (Ch. 4) may be examined later

04S1

COMP4211 Seminar

W02S23

## Name dependences

- ï A *name dependence* occurs when two instructions use the same register or memory location, called a *name* but there is no flow of data between the instructions associated with that name
- ï Two types, defined between an instruction *i* that *precedes* instruction *j*:
  1. An *antidependence* occurs when *j* writes to a name that *i* reads ñ the original ordering must be preserved
    - e.g.    ADD    R1, R3, R4
    - LD     R4, 0(R0)
  2. An *output dependence* occurs when *i* and *j* write to the same name ñ also requires order to be preserved
    - e.g.    ADD    R4, R3, R1
    - LD     R4, 0(R0)

04S1

COMP4211 Seminar

W02S24

- ï Since a name dependence is not a true (data) dependence, instructions involved in a name dependence can be executed simultaneously or be reordered if the name (register or memory location) is changed to avoid the conflict
- ï Renaming is more easily done for registers, and it can be done either statically by the compiler, or dynamically by hardware

## Data hazards

- ï A hazard is created whenever there is a dependence between instructions and they are close enough that the overlap caused by pipelining or reordering would change the order of access to the operand involved in the dependence
- ï We must then preserve *program order* i.e., the order instructions would execute in if executed sequentially
- ï Our goal is to exploit parallelism by preserving program order *only where it affects the outcome of the program*
- ï Detecting & avoiding hazards ensures the necessary program order is preserved

- ï Three types of data hazards depending upon the order of read and write accesses
- ï By convention, hazards are named by the ordering that must be preserved
- ï Consider two instructions  $i$  and  $j$  with  $i$  occurring before  $j$  in program order. The possible hazards are:
  - ñ **RAW (read after write)** ñ  $j$  tries to read a source before  $i$  writes it. This hazard is most common and corresponds to true data dependence
    - e.g. `LD R4, 0(R0)`
    - `ADD R1, R3, R4`

- ñ **WAW (write after write)** ñ  $j$  tries to write an operand before it is written by  $i$ . This hazard corresponds to output dependence. They occur in pipelines that write in more than one stage or allow instructions to proceed when previous ones are stalled. It is not present in the simple statically scheduled 5 stage pipeline that only writes in the WB stage.
- ñ **WAR (write after read)** ñ  $j$  tries to write a destination before it has been read by  $i$ . Arises from antidependence. Cannot occur in static issue pipelines when reads are earlier in the pipeline than writes. Can occur when some instruction writes early in the pipeline and another reads late, or when instructions can be reordered
- ï Note that **RAR (read after read)** is not a hazard

## Control dependence

- A control dependence determines the ordering of an instruction  $i$  with respect to a branch instruction so that  $i$  is executed in correct program order, and only when it should be. see later

e.g.

```
if p1 {
    s1;
};
if p2 {
    s2;
}
```

## Overcoming data hazards with dynamic scheduling

- Simple statically scheduled pipelines fetch instructions and issue them unless stalled due to some data dependence that cannot be hidden by forwarding
- Once stalled, no further instructions are fetched or issued until the dependence is cleared
- We now begin to explore *dynamic scheduling* in which the hardware rearranges instruction execution to reduce stalls while maintaining data flow and exception behaviour
- This technique allows us to handle dependences that are unknown at compile time (e.g. a memory reference) and allows code that was compiled with one pipeline in mind to be efficiently executed on a different pipeline
- Unfortunately, the benefits of dynamic scheduling are gained at the cost of a significant increase in hardware complexity

- A dynamically scheduled processor attempts to avoid stalls in the presence of dependences.
- In contrast, static pipeline scheduling by the compiler (Ch. 4) tries to minimize stalls by separating dependent instructions to avoid hazards