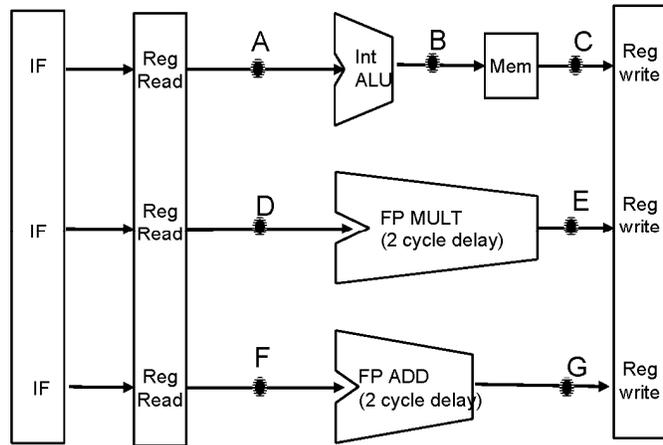


Question 3 (15 points) Consider the VLIW architecture shown below which consists of three units; an ALU unit (1 cycle delay), a FP add unit (2 cycle delay) and a FP multiply unit (2 cycle delay). Assume that the architecture executes the shown MIPS instruction set (L.D, S.D, ADD.D, MULT.D, L, S, ADD), with each VLIW instruction composed of up to three MIPS instructions, one for each unit. The points marked A, B, ... , G in the figure are used to identify forwarding paths. For example, $C \rightarrow D$ denotes a forwarding path from point C to point D. The usefulness of the path $C \rightarrow D$ can be demonstrated by the example shown in the following table where it is used to forward F0 from the “LD F0, 100(R1)” instruction to the “MULT.D F1, F0, F2” instruction:

Example for the use of the forwarding path from C → D			
	Int ALU	FP MULT	FP ADD
Cycle 1	L.D F0, 100(R1)		
Cycle 2			
Cycle 3		MULT.D F1, F0, F2	

- L.D F1, I(R1) // load to FP reg.
- S.D F1, I(R1) // store to FP reg.
- ADD.D F1, F2, F3 // FP add
- MULT.D F1, F2, F3 // FP multiply
- L R1, I(R2) // load to int reg.
- S R1, I(R2) // store to int reg.
- ADD R1, R2, R3 // Integer add

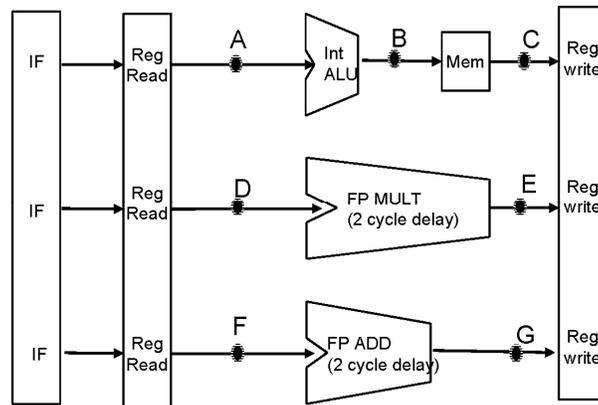


Complete the following tables to demonstrate an example of a use case for the forwarding path indicated in each table. You may use any of the above instructions in any order and with any registers. If a path does not have any use given the above instructions (should not be implemented) indicate so by writing “not useful” next to the table.

Example for the use of the forwarding path from E → F			
	Int ALU	FP MULT	FP ADD
Cycle 1			
Cycle 2			
Cycle 3			

Example for the use of the forwarding path from E → A			
	Int ALU	FP MULT	FP ADD
Cycle 1			
Cycle 2			
Cycle 3			

L.D F1, I(R1) // load to FP reg.
 S.D F1, I(R1) // store to FP reg.
 ADD.D F1, F2, F3 // FP add
 MULT.D F1, F2, F3 // FP multiply
 L R1, I(R2) // load to int reg.
 S R1, I(R2) // store to int reg.
 ADD R1, R2, R3 // Integer add



Example for the use of the forwarding path from E → B			
	Int ALU	FP MULT	FP ADD
Cycle 1			
Cycle 2			
Cycle 3			

Example for the use of the forwarding path from C → A			
	Int ALU	FP MULT	FP ADD
Cycle 1			
Cycle 2			
Cycle 3			

Example for the use of the forwarding path from C → B			
	Int ALU	FP MULT	FP ADD
Cycle 1			
Cycle 2			
Cycle 3			

Example for the use of the forwarding path from B → A			
	Int ALU	FP MULT	FP ADD
Cycle 1			
Cycle 2			
Cycle 3			

Example for the use of the forwarding path from B → D			
	Int ALU	FP MULT	FP ADD
Cycle 1			
Cycle 2			
Cycle 3			

Question 4 (10 points): Consider a superscalar architecture with two pipelined units, one for load/store and one for ALU instructions. The following two tables indicate the order of execution of two threads, A and B, and the latencies mandated by dependences between instructions. For example, A2 and A3 should execute after A1 and there should be at least four cycles between the execution of A3 and A5 and at least three cycles between A4 and A5. In other words, the tables indicate the schedule if the instructions of each of the threads are executed with no multithreading.

Note that an instruction that executes on one pipeline (for example A1 on the load/store pipeline) cannot execute on the other pipeline. Also, you cannot issue the instructions out of order.

time	Load/store pipeline	ALU pipeline
t	A1	
t+1	A3	A2
t+2		A4
t+3		
t+4		
t+5		
t+6	A5	
t+7	A6	A7

time	Load/store pipeline	ALU pipeline
t	B1	B2
t+1		B3
t+2	B4	B5
t+3		
t+4	B6	
t+5		B7
t+6	B8	
t+7		

Show the execution schedule for the two threads on the two pipelines assuming:

(a) Fine grain multithreading

(b) Simultaneous multithreading
(with priority always given to thread A)

time	Load/store pipeline	ALU pipeline
t		
t+1		
t+2		
t+3		
t+4		
t+5		
t+6		
t+7		
t+8		
t+9		
t+10		
t+11		
t+12		
t+13		
t+14		

time	Load/store pipeline	ALU pipeline
t		
t+1		
t+2		
t+3		
t+4		
t+5		
t+6		
t+7		
t+8		
t+9		
t+10		
t+11		
t+12		
t+13		
t+14		

Question 5 (15 points): Consider a system with a two level cache where 20% of the memory references are “stores” and 80% are “loads”. Out of all the memory references made by the CPU, 94% are found in the L1 cache and 50% of the L1 cache misses are found in L2. Assume that it takes one cycle to access L1, 10 cycles to access L2 (either to write a word or a block or to read a word or a block) and 100 cycles to access memory (either to write or read a block). Assume also that “write allocate” is used and that there is a very large and fast "write buffer" between L2 and memory while there is no “write buffers” between L1 and L2 (the CPU stalls until the writing is complete onto L2).

(a) Compute the average memory access time if L2 uses "write back" while L1 uses "write through".

(i) the average number of cycles for a load is

(ii) the average number of cycles for a store is

(iii) the average memory access time is

(b) Compute the average memory access time if both L1 and L2 use "write back" assuming that, on average, 40% of the evicted blocks are dirty.

Question 7 (15 points): Consider the multiplication of two 100×100 matrices on a system with a 64KB fully associative cache and 32B blocks (4 words, each 8 bytes). Answer the following assuming row-wise allocation and noting that $100 \times 100 \times 8 \approx 80\text{KB}$ is needed to store a 100×100 matrix.

(a) When executing the following code

```
for (i = 0 ; i < 100 ; i++)
  for (j = 0 ; j < 100 ; j++)
  {
    r = 0;
    for (k = 0; k < 100; k++)
      r = r + A[i][k]*B[k][j];
    C[i][j] = r;
  };
```

the average miss rate for accesses to the elements of A is _____

the average miss rate for accesses to the elements of B is _____

(b) When executing the following code, which computes two rows of C in each iteration of the outer loop,

```
for (i = 0 ; i < 100 ; i=i+2)
  for (j = 0 ; j < 100 ; j++)
  {
    r1 = 0 ;
    r2 = 0 ;
    for (k = 0; k < 100; k++)
      {
        r1 = r1 + A[i][k]*B[k][j];
        r2 = r2 + A[i+1][k]*B[k][j];
      }
    C[i][j] = r1;
    C[i+1][j] = r2;
  };
```

the average miss rate for accesses to the elements of A is _____

the average miss rate for accesses to the elements of B is _____

(c) It is possible to generalize the code in (b) so that m rows of C are computed in each iteration of the outer loop, where m is one of 4, 5, 10, 20, 25 or 50. What value of m would you use to obtain the lowest average cache miss rate of A and B (you do not have to write the code for that generalization – also you may ignore the miss rate for C)? – You need to give a good reason for your answer.