

CS152 Computer Architecture and
Engineering

VLIW, Vector, and
Multithreaded Machines

Assigned
04/01/2024

Problem Set #4, Version (1.1)

Due April 8
@ 11:59:59PT

<http://inst.eecs.berkeley.edu/~cs152/sp24>

The problem sets are intended to help you learn the material, and we encourage you to collaborate with other students and to ask questions in discussion sections and office hours to understand the problems. However, each student must turn in their own solution to the problems.

The problem sets also provide essential background material for the exam and the midterms. ***This particular problem set will be graded on completion basis.*** However, if you do not work through the problem sets yourself you are unlikely to succeed on the exam or midterms!

We will distribute solutions to the problem set after the deadline to give you feedback.

Assignments must be submitted through [Gradescope](#) by **11:59:59pm PT** on the specified due date. *Box/clearly mark all solutions that don't involve filling in a figure/table. Only boxed/clearly marked solutions and filled in figures/tables will be considered for grading.* See the course website for the policy on [slip days](#) (late submissions).

Name: _____

SID: _____

Collaborators (Name, SID):

Problem 1: Trace Scheduling

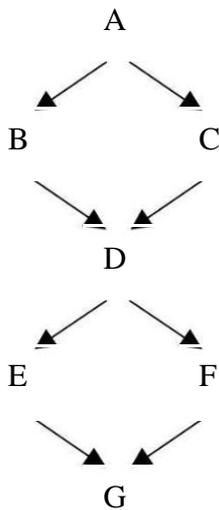
Trace scheduling is a compiler technique that increases ILP by removing control dependencies at compile time, allowing operations following branches to be moved up and speculatively executed in parallel with operations before the branch. It was originally developed for statically scheduled VLIW machines. However, it is a general technique that can be used in different types of machines such as a single-issue RISC-V processor.

Consider the following piece of C code (% is modulus) with basic blocks labeled:

```
A   if (data % 3 == 0)
B     X = V0 / V1;
    else
C     X = V2 / V3;
D   if (data % 4 == 0)
E     Y = V0 * V1;
    else
F     Y = V2 * V3;
G
```

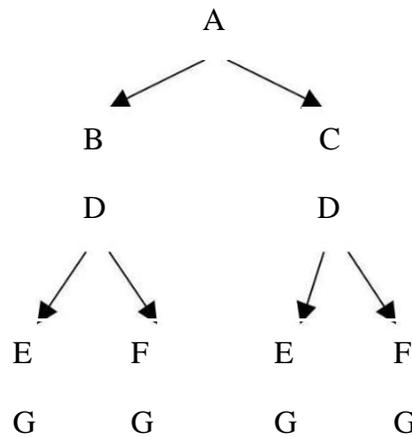
Assume that **data** is a uniformly distributed integer random variable that is set sometime before executing this code.

The program's control flow graph is



Path probabilities:

The decision tree is



A control flow graph and the decision tree both show the possible flow of execution through basic blocks. However, the control flow graph captures the static structure of the program, while the decision tree captures the dynamic execution (history) of the program.

Problem 1.A

On the decision tree, label each path with the probability of traversing that path. For example, the leftmost block will be labeled with the total probability of executing the path ABDEG. (Hint: you might want to write out the cases). Circle the path that is most likely to be executed.

Problem 1.B

Consider the following RISC-V code:

```
A:  lw   x1, data
     remi x2, x1, 3      # x2 <- x1 % 3
     bnez x2, C
B:  div  x3, x4, x5     # X <- V0 / V1
     j   D
C:  div  x3, x6, x7     # X <- V2 / V3
D:  remi x2, x1, 4     # x2 <- x1 % 4
     bnez x2, F
E:  mul  x8, x4, x5    # Y <- V0 * V1
     j   G
F:  mul  x8, x6, x7    # Y <- V2 * V3
G:
```

This code is to be executed on a statically scheduled VLIW machine. It has separate, long latency, **unpipelined** memory, divider, and multiplier units that can be run in parallel. Thus, we must wait for all functional units in a given VLIW instruction to finish before issuing the next instruction.

Note that the instruction to calculate the remainder (REMI) runs on the divider.

Assume that the load takes x cycles, the divider takes y cycles, and the multiplier takes z cycles. Approximately how many cycles does this code take *in the best case*, *in the worst case*, and *on average*? (Ignore the latency of integer instructions such as `bnez` and `j`.)

Problem 1.C

With trace scheduling, we can obtain the following code:

```
ACF: ld    x1, data
      div  x3, x6, x7    # X  <- V2 / V3
      mul  x8, x6, x7    # Y  <- V2 * V3
A:    remi x2, x1, 3     # x2 <- x1 % 3
      bnez x2, D
B:    div  x3, x4, x5    # X  <- V0 / V1
D:    remi x2, x1, 4     # x2 <- x1 % 4
      bnez x2, G
E:    mul  x8, x4, x5    # Y  <- V0 * V1
G:
```

We optimize only for the most common path, but the other paths are still correct.

- i) Approximately how many cycles does the new code take *in the best case, in the worst case* and *on average*?
- ii) If memory takes the most cycles (i.e. $x > y, z$), is this code faster *in the best case, in the worst case* and *on average* than the code in Problem 1.B?
- iii) If either of the functional units has the longest latency (i.e. $y, z > x$), is this code faster *in the best case, in the worst case* and *on average* than the code in Problem 1.B?

You might find the *max* function useful.

The following hint might help simplify the comparison with Problem 1.B.

$$\forall a, b \geq 0, \quad \frac{1}{4}a + \frac{5}{12}b \leq \frac{2}{3} \max(a, b)$$

Problem 2: VLIW machines

In this problem, we consider the execution of a code segment on a VLIW processor. The code we consider is the SAXPY kernel, which scales a single-precision vector X by a single-precision constant A , adding this quantity to a single-precision vector Y .

```
for(i = 0; i < N; i++) {  
    Y[i] = Y[i] + A*X[i];  
}
```

```
loop:  1.flw   f1, 0(x1)      # f1 = X[i]  
       2.fmul  f2, f0, f1   # A * X[i]  
       3.flw   f3, 0(x2)    # f3 = Y[i]  
       4.fadd  f4, f2, f3   # f4 = Y[i] + A*X[i]  
       5.fsw   f4, 0(x2)    # Y[i] = f4  
       6.addi  x1, x1, 4     # bump pointer  
       7.addi  x2, x2, 4     # bump pointer  
       8.bne  x1, x3, loop  # loop
```

Now we have a VLIW machine with seven execution units:

- two ALU units, latency one cycle, also used for branch operations
- three memory units, latency three cycles, fully pipelined, each unit can perform either a store or a load
- two FPU units, latency four cycles, fully pipelined, one unit can perform **fadd** operations, the other **fmul** operations.

Below is the format of a VLIW instruction:

Int Op 1	Int Op 2	Mem Op 1	Mem Op 2	Mem Op 3	FP Add	FP Mul
----------	----------	----------	----------	----------	--------	--------

Our machine has no interlocks. The result of an operation is written to the register file immediately after it has gone through the corresponding execution unit: one cycle after issue for ALU operations, three cycles for memory operations and four cycles for FPU operations. The old values can be read from the registers until they have been overwritten.

When writing code for this machine, you may:

- 1) Assume the arrays are long (> 32 elements)
- 2) Assume the arrays have an even number of elements
- 3) Reorder instructions and change immediates as long as the code is functionally correct

Do not count floating point memory operations towards FLOPS in this problem.

Problem 2.A: No Code Optimization

Schedule instructions for the VLIW machine in Table P4.2-1 without loop unrolling and software pipelining. What is the throughput of the loop in the code in floating point operations per cycle (FLOPS/cycle)?

Problem 2.B: Loop Unrolling

Schedule instructions for the VLIW machine in Table P4.2-2 only with loop unrolling. Write the assembly code by unrolling the loop once. What is the throughput of the loop in the code in floating point operations per cycle (FLOPS/cycle)? What is the speedup over Problem 2.A?

Problem 2.C: Software Pipelining

Schedule instructions for the VLIW machine in Table P4.2-3 only with software pipelining. Include the prologue and the epilogue in Table P4.2-3. What is the throughput of the loop in the code in floating point operations per cycle (FLOPS/cycle)? What is the speedup over Problem 2.A?

Problem 2.D: Loop Unrolling + Software Pipelining

Schedule instructions for the VLIW machine in Table P4.2-4 with both loop unrolling and software pipelining. Unroll the loop once as in Problem 2.B. Include the prologue and the epilogue in Table P4.2-4. What is the throughput of the loop in the code in floating point operations per cycle (FLOPS/cycle)? What is the speedup over Problem 2.A?

Problem 3: Vector Machines

In this problem, we analyze the performance of vector machines. We start with a baseline vector processor with the following features:

- **32 elements per vector register**
- **8 lanes**
- **One ALU per lane: 1 cycle latency**
- **One load/store unit per lane: 4 cycle latency, fully pipelined**
- **No dead time**
- **No support for chaining**
- **Scalar instructions execute on a separate 5-stage pipeline**

To simplify the analysis, we assume a magic memory system with no bank conflicts and no cache misses.

We consider execution of the following loop:

```
C code  
for (i = 0; i < N; i++) {  
    C[i] = A[i] + B[i] - 1;  
}
```

```
loop:  1.LV      V1, (x1)      # load A  
       2.LV      V2, (x2)      # load B  
       3.ADDV    V3, V1, V2    # A + B  
       4.SUBVS   V4, V3, x4    # subtract x4 = 1  
       5.SV      V4, (x3)      # store C  
       6.ADDI    x1, x1, 128    # bump pointer  
       7.ADDI    x2, x2, 128    # bump pointer  
       8.ADDI    x3, x3, 128    # bump pointer  
       9.SUBI    x5, x5, 32     # i++ (x5 = N)  
      10.BNEZ   x5, loop      # loop
```

Problem 3.A: Simple Vector Processor

Complete the pipeline diagram in Table P4.4-1 of the baseline vector processor running the given code. The following **supplementary information** explains the diagram:

Scalar instructions execute in 5 cycles: fetch (**F**), decode (**D**), execute (**X**), memory (**M**), and writeback (**W**). A vector instruction is also fetched (**F**) and decoded (**D**). Then, it stalls (—) until its required vector functional unit is available. With no chaining, a dependent vector instruction stalls until the previous instruction finishes writing back all of its elements. A vector instruction is pipelined across all the lanes in parallel. For each element, the operands are read (**R**) from the vector register file, the operation executes on the load/store unit (**M**) or the ALU (**X**), and the result is written back (**W**) to the vector register file. A stalled vector instruction does not block a scalar instruction from executing.

Inst #	cycle																																										
	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40			
1	F	D	R	M1	M2	M3	M4	W																																			
1				R	M1	M2	M3	M4	W																																		
1					R	M1	M2	M3	M4	W																																	
1						R	M1	M2	M3	M4	W																																
2		F	D	-	-	-	R	M1	M2	M3	M4	W																															
2								R	M1	M2	M3	M4	W																														
2									R	M1	M2	M3	M4	W																													
2										R	M1	M2	M3	M4	W																												
3			F	D	-	-	-	-	-	-	-	-	-	-	-	R	X1	W																									
3																	R	X1	W																								
3																		R	X1	W																							
3																			R	X1	W																						
4			F	D	-																																						
4																																											
4																																											
4																																											
5				F	D																																						
5																																											
5																																											
5																																											
6					F	D	X	M	W																																		
7						F	D	X	M	W																																	
8							F	D	X	M	W																																
9								F	D	X	M	W																															
10									F	D	X	M	W																														
1										F	D	-																															
1																																											
1																																											
1																																											

Table P4.3-1: Vector Pipeline Diagram (8 Lanes without Chaining)

Problem 3.B: Hardware Optimization (Chaining)

In this question, we analyze the performance benefits of chaining and additional lanes. Vector chaining is done through the register file and an element can be read (**R**) on the same cycle in which it is written back (**W**), or it can be read on any later cycle (the chaining is *flexible*). For this question, we always assume 32 elements per vector register, so there are 4 elements per lane with 8 lanes, and 1 element per lane with 32 lanes.

To analyze performance, we calculate the total number of cycles per vector loop iteration by summing the number of cycles between the issuing of successive vector instructions. For example, in Question 3.A, Inst #1(LV) begins execution in cycle 3, Inst #2(LV) in cycle 7 and Inst #3(ADDV) in cycle 16. Therefore, there are 4 cycles between Inst #1 and Inst #2 and 9 cycles between Inst #2 and Inst #3.

First, fill in Table P4.3-2 for 8 lanes with chaining, Table P4.3-3 for 16 lanes with chaining, and Table P4.3-4 for 32 lanes with chaining.

Note that, with 8 lanes and chaining, Inst #4(SUBVS) cannot issue 2 cycles after Inst #3(ADDV) because there is only one ALU per lane. This would be possible in the absence of a structural hazard.

Also, complete the following table. The first row corresponds to the baseline 8-lane vector processor with no chaining. The second row adds flexible chaining to the baseline processor, and the last two rows increase the number of lanes from 8 to 32.

Vector Processor Configuration	Number of cycles between successive vector instructions					Total cycles per vector loop iteration
	#1(LV) #2(LV)	#2(LV) #3(ADDV)	#3(ADDV) #4(SUBVS)	#4(SUBVS) #5(SV)	#5(SV) #1(LV)	
8 lanes, no chaining	4	9	6	6	4	29
8 lanes, chaining						
16 lanes, chaining						
32 lanes, chaining						

Inst #	Cycle																																										
	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40			
1	F	D	R	M1	M2	M3	M4	W																																			
1				R	M1	M2	M3	M4	W																																		
1					R	M1	M2	M3	M4	W																																	
1						R	M1	M2	M3	M4	W																																
2		F	D	-	-	-	R	M1	M2	M3	M4	W																															
2								R	M1	M2	M3	M4	W																														
2									R	M1	M2	M3	M4	W																													
2										R	M1	M2	M3	M4	W																												
3			F	D	-	-	-	-	-	-	-	R	X1	W																													
3													R	X1	W																												
3														R	X1	W																											
3															R	X1	W																										
4				F	D	-																																					
4																																											
4																																											
4																																											
5					F	D																																					
5																																											
5																																											
5																																											
6						F	D	X	M	W																																	
7							F	D	X	M	W																																
8								F	D	X	M	W																															
9									F	D	X	M	W																														
10										F	D	X	M	W																													
1											F	D	-																														
1												F	D	-																													
1													F	D	-																												
1														F	D	-																											

Table P4.3-2: 8 Lanes with Chaining

Inst	Cycle																																							
------	-------	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--

#	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	
1	F	D	R	M1	M2	M3	M4	W																																	
1				R	M1	M2	M3	M4	W																																
2		F	D	-																																					
2																																									
3			F	D	-																																				
3																																									
4				F	D	-																																			
4																																									
5					F	D	-																																		
5																																									
6						F	D	X	M	W																															
7							F	D	X	M	W																														
8								F	D	X	M	W																													
9									F	D	X	M	W																												
10										F	D	X	M	W																											
1											F	D	-																												
1																																									

Table P4.3-3: 16 Lanes with Chaining

Inst #	Cycle																																											
	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40				
1	F	D	R	M1	M2	M3	M4	W																																				
2		F	D																																									
3			F	D	-																																							
4				F	D	-																																						
5					F	D	-																																					
6						F	D	X	M	W																																		
7							F	D	X	M	W																																	
8								F	D	X	M	W																																
9									F	D	X	M	W																															
10										F	D	X	M	W																														
1											F	D	-																															

Table P4.3-3: 32 Lanes with Chaining

Problem 4: Multithreading

This problem evaluates the effectiveness of multithreading using a simple database benchmark. The benchmark searches for an entry in a linked list built from the following structure, which contains a key, a pointer to the next node in the linked list, and a pointer to the data entry.

```
struct node {
    int key;
    struct node *next;
    struct data *ptr;
}
```

The following RISC-V code shows the core of the benchmark, which traverses the linked list and finds an entry with a particular key.

```
loop:  LW    x3, 0(x1)           # load a key
       LW    x4, 4(x1)         # load the next pointer
       SEQ   x3, x3, x2        # set x3 if x3 == x2
       BNEZ  x3, end           # found the entry
       ADD   x1, x0, x4        # check the next node
       BNEZ  x4, loop
end:
```

We run this benchmark on a single-issue in-order processor. The processor can fetch and issue (dispatch) one instruction per cycle. If an instruction cannot be issued due to a data dependency, the processor stalls. Integer instructions have single-cycle latency; the result can be used in the next cycle. For example, if SEQ is executed in cycle 1, BNEZ can be executed in cycle 2. We also assume that the processor has a perfect branch predictor with no penalty for both taken and not-taken branches.

Problem 4.A

Assume that our system does not have a cache. Each memory operation directly accesses main memory which has a latency of 50 cycles. The load/store unit is fully pipelined, and non-blocking. After the processor issues a memory operation, it can continue executing instructions until it reaches an instruction that is dependent on an outstanding memory operation.

How many cycles does it take to execute one iteration of the loop in steady state? End Cycle refers to the cycle when an instruction's result can be used.

Instruction	Start Cycle	End Cycle
LW x3, 0(x1)		
LW x4, 4(x1)		
SEQ x3, x3, x2		
BNEZ x3, End		
ADD x1, x0, x4		
BNEZ x1, Loop		

Problem 4.B

Now we add zero-overhead multithreading to our pipeline. A processor executes multiple threads, each of which performs an independent search. Hardware mechanisms schedule a thread to execute each cycle.

In our first implementation, the processor switches to a different thread every cycle using fixed round robin scheduling (similar to CDC 6600 PPUs). Each of the N threads executes one instruction every N cycles. What is the **minimum** number of threads that we need to fully utilize the processor, i.e., execute one instruction per cycle?

Problem 4.C

How does multithreading affect throughput (number of keys the processor can find within a given time) and latency (time processor takes to find an entry with a specific key)? Assume the processor switches to a different thread every cycle and is fully utilized. Check the correct boxes.

	Throughput	Latency
Better		
Same		
Worse		

Problem 4.D

We change the processor to only switch to a different thread when an instruction cannot execute due to data dependency. What is the minimum number of threads to fully utilize the processor in steady state? Note that the processor issues instructions in-order in each thread.

Problem 5: Multithreading

Consider a single-issue in-order multithreading processor that is similar to the one described in Problem 4.

Each cycle, the processor can fetch and issue one instruction that performs any of the following operations:

- **load/store, 13-cycle latency (fully pipelined)**
- **integer add, 1-cycle latency**
- **floating-point add, 6-cycle latency (fully pipelined)**
- **branch, no delay slots, 1-cycle latency**

The processor **does not have a cache**. Each memory operation directly accesses main memory. If an instruction cannot be issued due to a data dependency, the processor stalls. We also assume that the processor has a perfect branch predictor with no penalty for both taken and not-taken branches.

Your job is to analyze the processor utilizations for the following two thread-switching implementations:

Fixed Switching: the processor switches to a different thread every cycle using fixed round robin scheduling. Each of the N threads executes an instruction every N cycles.

Data-dependent Switching: the processor only switches to a different thread when an instruction cannot execute due to a data dependency.

Each thread executes the following RISC-V code:

```
loop:  LD      f2, 0(x1)      # load data into f2
      ADDI   x1, x1, 4      # bump pointer
      FADD  f3, f3, f2      # f3 = f3 + f2
      BNE   f2, f4, loop    # continue if f2 != f4
```

Problem 5.A

What is the minimum number of threads that we need to fully utilize the processor for each implementation in steady state? Briefly explain your reasoning.

Fixed Switching: _____ **Thread(s)**

Data-dependent Switching: _____ **Thread(s)**

Problem 5.B

What is the minimum number of threads that we need to fully utilize the processor in steady state for each implementation if we change the **load/store latency to 1-cycle (but keep the 6-cycle floating-point add)**? Briefly explain your reasoning.

Fixed Switching: _____ **Thread(s)**

Data-dependent Switching: _____ **Thread(s)**

Problem 5.C

Consider a **Simultaneous Multithreading (SMT)** machine with limited hardware resources. **Circle** the following hardware constraints that can limit the total number of threads that the machine can support. For the item(s) that you circle, **briefly describe** the minimum requirement to support N threads.

(A) Number of Functional Unit:

(B) Number of Physical Registers:

(C) Data Cache Size:

(D) Data Cache Associativity: