

DIGITAL DESIGN AND COMPUTER ARCHITECTURE(252-0028-00L), SPRING 2023
OPTIONAL HW 6: SIMD PROCESSING, GPUS, MEMORY ORGANIZATION, AND MEMORY TECHNOLOGY
SOLUTIONS

Instructor: Prof. Onur Mutlu

TAs: Juan Gomez Luna, Mohammad Sadrosadati, Mohammed Alser, Ataberk Olgun, Giray Yaglikci, Can Firtina,
Geraldo De Oliveira Junior, Rahul Bera, Konstantinos Kanellopoulos, Nika Mansouri Ghiasi, Nisa Bostancı,
Rakesh Nadig, Joel Lindegger, İsmail Emir Yüksel, Haocong Luo, Yahya Can Tuğrul, Julien Eudine

Released: Saturday, June 3, 2023

1 Vector Processing I

Consider the following piece of code:

```
for (i = 0; i < 100; i ++)  
    A[i] = ((B[i] * C[i]) + D[i])/2;
```

- (a) Translate this code into assembly language using the following instructions in the ISA (note the number of cycles each instruction takes is shown next to each instruction):

Opcode	Operands	Number of cycles	Description
LEA	Rd, X	1	$Rd \leftarrow \text{address of } X$
LD	Rd, Rs, Rt	11	$Rd \leftarrow \text{MEM}[Rs + Rt]$
ST	Rs, Rt, Ru	11	$\text{MEM}[Rt + Ru] \leftarrow Rs$
MOVI	Rd, imm	1	$Rd \leftarrow \text{imm}$
MUL	Rd, Rs, Rt	6	$Rd \leftarrow Rs \times Rt$
ADD	Rd, Rs, Rt	4	$Rd \leftarrow Rs + Rt$
ADD	Rd, Rs, imm	4	$Rd \leftarrow Rs + \text{imm}$
RSHFA	Rd, Rs, shamt	1	$Rd \leftarrow Rs \ggg \text{shamt}$
BR<CC>	Rs, X	1	Branch to X if Rs satisfies condition code CC

Assume one memory location is required to store each element of the array. Also assume that there are eight register, R0 to R7.

Condition codes are set after the execution of an arithmetic instruction. You can assume typically available condition codes such as zero (EQZ), positive (GTZ), negative (LTZ), non-negative (GEZ), and non-positive (LEZ).

MOVI	R1, 99	// 1 cycle
LEA	R0, A	// 1 cycle
LEA	R2, B	// 1 cycle
LEA	R3, C	// 1 cycle
LEA	R4, D	// 1 cycle
LOOP:		
LD	R5, R2, R1	// 11 cycles
LD	R6, R3, R1	// 11 cycles
MUL	R7, R5, R6	// 6 cycles
LD	R5, R4, R1	// 11 cycles
ADD	R6, R7, R5	// 4 cycles
RSHFA	R7, R6, 1	// 1 cycle
ST	R7, R0, R1	// 11 cycles
ADD	R1, R1, -1	// 4 cycles
BRGEZ	R1 LOOP	// 1 cycle

How many cycles does it take to execute the program?

$$5 + 100 \times 60 = 6005$$

- (b) Now write Cray-like vector assembly code with the minimum number of instructions. Assume that there are eight vector registers and the length of each vector register is 64. Use the following instructions in the vector ISA:

Opcode	Operands	Number of cycles	Description
SET	Vst, imm	1	$Vst \leftarrow \text{imm}$ (Vst: Vector Stride Register)
SET	Vln, imm	1	$Vln \leftarrow \text{imm}$ (Vln: Vector Length Register)
VLD	Vd, X	11, pipelined	$Vd \leftarrow \text{MEM}[\text{address of X}]$
VST	Vs, X	11, pipelined	$\text{MEM}[\text{address of X}] \leftarrow Vs$
VADD	Vd, Vs, Vt	4, pipelined	$Vd \leftarrow Vs + Vt$
VMUL	Vd, Vs, Vt	6, pipelined	$Vd_i \leftarrow Vs_i \times Vt_i$
VRSHFA	Vd, Vs, shamt	1	$Vd_i \leftarrow Vs_i \ggg \text{shamt}$

```

SET      Vln, 50
SET      Vst, 1
VLD      V1, B
VLD      V2, C
VMUL     V4, V1, V2
VLD      V3, D
VADD     V6, V4, V3
VRSHFA   V7, V6, 1
VST      V7, A

VLD      V1, B + 50
VLD      V2, C + 50
VMUL     V4, V1, V2
VLD      V3, D + 50
VADD     V6, V4, V3
VRSHFA   V7, V6, 1
VST      V7, A + 50

```

How many cycles does it take to execute the program on the following processors? Assume that memory is 16-way interleaved.

Vector processor without chaining, 1 port to memory (1 load or store per cycle):

The third load (VLD) can be pipelined with the multiplication (VMUL). However as there is just only one port to memory and no chaining, other operations cannot be pipelined.

Processing the first 50 elements takes 346 cycles as below:

```

| 1 | 1 | 11 | 49 | 11 | 49 | 6 | 49 |
                        | 11 | 49 | 4 | 49 | 1 | 49 | 11 | 49 |

```

Processing the next 50 elements takes 344 cycles as shown below (no need to initialize Vln and Vst as they stay at the same value).

```

| 11 | 49 | 11 | 49 | 6 | 49 |
                        | 11 | 49 | 4 | 49 | 1 | 49 | 11 | 49 |

```

Therefore, the total number of cycles to execute the program is 690 cycles.

Vector processor with chaining, 1 port to memory:

In this case, the first two loads cannot be pipelined as there is only one port to memory and the third load has to wait until the second load has completed. However, the machine supports chaining, so all other operations can be pipelined.

Processing the first 50 elements takes 242 cycles as below:

```

| 1 | 1 | 11 |   49 |   11 |   49 |
                        | 6 |   49 |
                        | 11 |   49 |
                        | 4 |   49 |
                        | 1 |   49 |
                        | 11 |   49 |

```

Processing the next 50 elements takes 240 cycles (same time line as above, but without the first 2 instructions to initialize Vln and Vst).

Therefore, the total number of cycles to execute the program is 482 cycles.

Vector processor with chaining, 2 read ports and 1 write port to memory:

Assuming an in-order pipeline.

The first two loads can also be pipelined as there are two ports to memory. The third load has to wait until the first two loads complete. However, the two loads for the second 50 elements can proceed in parallel with the store.

```

SET      |1|
SET      |1|
VLD      |11| 49 |
VLD      |1|11| 49 |
VMUL      |6| 49 |
VLD      |11| 49 |
VADD      |4| 49 |
VRSHFA    |1| 49 |
VST       |11| 49 |

VLD      |1|11| 49 |
VLD      |11| 49 |
VMUL      |6| 49 |
VLD      |11| 49 |
VADD      |4| 49 |
VRSHFA    |1| 49 |
VST       |11| 49 |

1 + 1 + (11 + 49) + 11 + 4 + 1 + (1 + 11 + 49) + 11 + 4 + 1 + (11 + 49) = 215

```

Therefore, the total number of cycles to execute the program is 215 cycles.

2 Vector Processing II

You are studying a program that runs on a vector computer with the following latencies for various instructions:

- VLD and VST: 50 cycles for each vector element; fully interleaved and pipelined.
- VADD: 4 cycles for each vector element (fully pipelined).
- VMUL: 16 cycles for each vector element (fully pipelined).
- VDIV: 32 cycles for each vector element (fully pipelined).
- VRSHFA: 1 cycle for each vector element (fully pipelined).

Assume that:

- The machine has an in-order pipeline.
 - The machine supports chaining between vector functional units.
 - In order to support 1-cycle memory access after the first element in a vector, the machine interleaves vector elements across memory banks. All vectors are stored in memory with the first element mapped to bank 0, the second element mapped to bank 1, and so on.
 - Each memory bank has an 8 KB row buffer.
 - Vector elements are 64 bits in size.
 - Each memory bank has two ports (so that two loads/stores can be active simultaneously), and there are two load/store functional units available.
- (a) What is the minimum power-of-two number of banks required in order for memory accesses to never stall? (Assume a vector stride of 1.)

64 banks (because memory latency is 50 cycles and the next power of two is 64)

- (b) The machine (with as many banks as you found in part a) executes the following program (assume that the vector stride is set to 1):

```
VLD    V1, A          // V1 ← A
VLD    V2, B          // V2 ← B
VADD   V3, V1, V2     // V3 ← V1 + V2
VMUL   V4, V3, V1     // V4i ← V3i × V1i
VRSHFA V5, V4, 2      // V5i ← V4i >>> 2
```

It takes 111 cycles to execute this program. What is the vector length L (i.e., the number of elements in a vector)?

$L = 40$

```
VLD      |-----50-----|---(L-1)---|
VLD      |1|-----50-----|
VADD                      | -4- |
VMUL                      | -16- |
VRSHFA                    |1|------(L-1)-----|
```

$$1 + 50 + 4 + 16 + 1 + (L - 1) = 71 + L = 111$$

$$\therefore L = 40$$

If the machine did not support chaining (but could still pipeline independent operations), how many cycles would be required to execute the same program?

228 cycles

```
VLD      |-----50-----|---(L-1)---|
VLD      |1|-----50-----|---(L-1)---|
VADD                      | -4- |---(L-1)---|
VMUL                      | -16- |---(L-1)---|
VRSHFA                    |1|---(L-1)---|
```

$$1 + 50 + 4 + 16 + 1 + 4 \times (L - 1) = 68 + 4 \times L = 228 \text{ cycles}$$

- (c) The architect of this machine decides that she needs to cut costs in the machine's memory system. She reduces the number of banks by a factor of 2 from the number of banks you found in part (a) above. Because loads and stores might stall due to bank contention, an *arbiter* is added to each bank so that pending loads from the oldest instruction are serviced first. How many cycles does the program take to execute on the machine with this reduced-cost memory system (but with chaining)?

129 cycles

```

VLD [0]  |----50----|    bank 0 (takes port 0)
...
[31]  |--31--|----50----| bank 31
[32]                |---50---| bank 0 (takes port 0)
...
[39]                |--7--|  bank 7
VLD [0]  |1|----50----|    bank 0 (takes port 1)
...
[31]  |1|--31--|----50----| bank 31
[32]                |---50---| bank 0 (takes port 1)
...
[39]                |--7--|  bank 7
VADD                                |--4--| (tracking last elements)
VMUL                                |--16--|
VRSHFA                                |1|

```

(B[39]: $1 + 50 + 50 + 7$) + 4 + 16 + 1 = 129 cycles

Now, the architect reduces cost further by reducing the number of memory banks (to a lower power of 2). The program executes in 279 cycles. How many banks are in the system?

8 banks

```

VLD [0]  |----50----|
...
[8]                |---50---|
...
[16]                |--50--|
...
[24]                |--50--|
...
[32]                |--50--|
...
[39]                |--7--|
VLD [39]                |1|
VADD                                |--4--|
VMUL                                |--16--|
VRSHFA                                |1|

```

$5 \times 50 + 7 + 1 + 4 + 16 + 1 = 279$ cycles \rightarrow 8 banks

- (d) Another architect is now designing the second generation of this vector computer. He wants to build a multicore machine in which 4 vector processors share the same memory system. He scales up the number of banks by 4 in order to match the memory system bandwidth to the new demand. However, when he simulates this new machine design with a separate vector program running on every core, he finds that the average execution time is longer than if each individual program ran on the original single-core system with 1/4 the banks. Why could this be? Provide concrete reason(s).

Row-buffer conflicts (all cores interleave their vectors across all banks).

What change could this architect make to the system in order to alleviate this problem (in less than 20 words), while *only* changing the shared memory hierarchy?

Partition the memory mappings, or using better memory scheduling.

3 SIMD Processing

Suppose we want to design a SIMD engine that can support a vector length of 16. We have two options: a traditional vector processor and a traditional array processor.

(a) Which one is more costly in terms of chip area (circle one)?

1. The traditional vector processor
2. The traditional array processor
3. Neither

Justify your answer.

An array processor requires 16 functional units for an operation whereas a vector processor requires only 1.

(b) Assuming the latency of an addition operation is five cycles in both processors, how long will a VADD (vector add) instruction take in each of the processors (assume that the adder can be fully pipelined and is the same for both processors)?

For a vector length of 1:

The traditional vector processor:

5 cycles

The traditional array processor:

5 cycles

For a vector length of 4:

The traditional vector processor:

8 cycles (5 for the first element to complete, 3 for the remaining 3)

The traditional array processor:

5 cycles

For a vector length of 16:

The traditional vector processor:

20 cycles (5 for the first element to complete, 15 for the remaining 15)

The traditional array processor:

5 cycles

4 GPUs and SIMD I

We define the *SIMD utilization* of a program run on a GPU as the fraction of SIMD lanes that are kept busy with *active threads* during the run of a program. The following code segment is run on a GPU. Each thread executes **a single iteration** of the shown loop. Assume that the data values of the arrays A and B are already in vector registers so there are no loads and stores in this program. (Hint: Notice that there are 2 instructions in each thread.) A warp in the GPU consists of 32 threads, there are 32 SIMD lanes in the GPU. Assume that each instruction takes the same amount of time to execute.

```
for (i = 0; i < N; i++) {  
    if (A[i] % 3 == 0) {    // Instruction 1  
        A[i] = A[i] * B[i]; // Instruction 2  
    }  
}
```

- (a) How many warps does it take to execute this program? Please leave the answer in terms of N .

$$\lceil \left(\frac{N}{32} \right) \rceil$$

- (b) Assume integer arrays A have a repetitive pattern which have 24 ones followed by 8 zeros repetitively and integer arrays B have a different repetitive pattern which have 48 zeros followed by 64 ones. What is the SIMD utilization of this program?

$$0.625$$

$$((24+8 \times 2)/(32 \times 2)) = 40/64 = 0.625$$

- (c) Is it possible for this program to yield a SIMD utilization of 100%? Circle one.

YES

NO

If YES, what should be true about array A for the SIMD utilization to be 100%?

Starting from A[0], consecutive 32 elements of A should be either divisible by 3 or not divisible by 3.

What should be true about array B?

B can be any array of integers.

If NO, explain why not.

(d) Is it possible for this program to yield a SIMD utilization of 56.25%? Circle one.

YES

NO

If YES, what should be true about array A for the SIMD utilization to be 56.25%?

4 out of every 32 elements of A are divisible by 3.

What should be true about array B?

B can be any array of integers.

If NO, explain why not.

(e) Is it possible for this program to yield a SIMD utilization of 50%? Circle one.

YES

NO

If YES, what should be true about array A for the SIMD utilization to be 50%?

What should be true about array B?

If NO, explain why not.

The minimum is when only 1 out of every 32 elements of A is divisible by 3. This yields a 51.5625% usage.

Consider the following example of a program that consists of 3 warps X, Y and Z that are executing the same code segment specified at the top of this question. Assume that the vector below specifies the direction of the branch of each thread within the warp. 1 means the branch in Instruction 1 is resolved to taken and 0 means the branch in Instruction 1 is resolved to not taken.

[illegible]

- [illegible]

There are several answers for this questions but the key is that the taken branch in Z can be combined with either X or Y. However, the taken branch in the first thread of X and Y cannot be merged because they are on the same GPU lane.

- No. Branch divergence happens on the same lane throughout the program resulting in the case where there is no dynamically formed warp.

5 GPUs and SIMD II

We define the *SIMD utilization* of a program run on a GPU as the fraction of SIMD lanes that are kept busy with *active threads* during the run of a program. As we saw in lecture and practice exercises, the SIMD utilization of a program is computed across the *complete run* of the program.

The following code segment is run on a GPU. Each thread executes **a single iteration** of the shown loop. Assume that the data values of the arrays A, B, and C are already in vector registers so there are no loads and stores in this program. (Hint: Notice that there are 6 instructions in each thread.) A warp in the GPU consists of 64 threads, and there are 64 SIMD lanes in the GPU. Please assume that all values in array B have magnitudes less than 10 (i.e., $|B[i]| < 10$, for all i).

```
for (i = 0; i < 1024; i++) {  
    A[i] = B[i] * B[i];  
    if (A[i] > 0) {  
        C[i] = A[i] * B[i];  
        if (C[i] < 0) {  
            A[i] = A[i] + 1;  
        }  
        A[i] = A[i] - 2;  
    }  
}
```

(a) How many warps does it take to execute this program?

16 warps

Number of Warps = $\lceil (\text{Number of threads}) / (\text{Number of threads per warp}) \rceil$

Number of threads = 2^{10} (i.e., one thread per loop iteration).

Number of threads per warp = $64 = 2^6$ (given).

\therefore Number of Warps = $\lceil 2^{10}/2^6 = 2^4 \rceil$

(b) What is the maximum possible SIMD utilization of this program?

100%

- (c) Please describe what needs to be true about array **B** to reach the maximum possible SIMD utilization asked in part (b). (Please cover all cases in your answer)

For every 64 consecutive elements: every value is 0, every value is positive, or every value is negative. Must give all three of these.

- (d) What is the minimum possible SIMD utilization of this program?

132/384

The first two lines must be executed by every thread in a warp (64/64 utilization for each line). The minimum utilization results when a single thread from each warp passes both conditions on lines 2 and 4, and every other thread fails to meet the condition on line 2. The thread per warp that meets both conditions, executes lines 3-6 resulting in a SIMD utilization of 1/64 for each line. The minimum SIMD utilization sums to $(64 \times 2 + 1 \times 4) / (64 \times 6) = 132/384$

- (e) Please describe what needs to be true about array **B** to reach the minimum possible SIMD utilization asked in part (d). (Please cover all cases in your answer)

Exactly 1 of every 64 consecutive elements must be negative. The rest must be zero. This is the only case that this holds true.

6 Memory Organization & Technology

Read the following statements about memory organization & technology. Circle “True” if the statement is true and “False” otherwise. *Note: we will subtract 1 point for each **incorrect** answer and award 0 points for unanswered questions.*

1. A main memory access typically consumes less energy than a register file access.
1. True 2. False
2. Building a larger memory array by increasing the length of the array’s wordlines and bitlines increases the cost (\$) but does not increase the access time of the array.
1. True 2. False
3. Activating a DRAM cell temporarily destroys the value stored in the DRAM cell.
1. True 2. False
4. DRAM cost (\$) per bit is much higher than that of SRAM.
1. True 2. False
5. The memory hierarchy of a typical computer system comprises different memory technologies.
1. True 2. False
6. Recently accessed data should be kept at the bottom-level in the memory hierarchy (e.g., main memory or disk) and not at the top-level (e.g., caches) in the hierarchy.
1. True 2. False
7. A program with no branches has high temporal locality in its instruction memory references.
1. True 2. False
8. A cache that has a block size equal to word size of memory access instructions cannot exploit spatial locality.
1. True 2. False
9. Memory banking enables concurrent access to the memory structure.
1. True 2. False
10. In DRAM, accesses to different rows in one bank can be serviced faster compared to accesses to the same row in one bank.
1. True 2. False
11. PCM is non-volatile, which means PCM retains stored data even when it is powered off.
1. True 2. False
12. If a hypothetical system is not constrained by chip area, memory cost (\$), and energy consumption, DRAM would be the best memory technology to use in that system.
1. True 2. False
13. The entire page table is typically stored in physical memory.
1. True 2. False
14. Virtual-to-physical address translation is on the critical path of a memory access.
1. True 2. False
15. Virtual memory makes programmer’s and microarchitect’s tasks easier.
1. True 2. False

EXTRA EXERCISES FOR PRACTICING

The following exercises are old exam questions that are conceptually similar to the ones above, but with slight alterations. We do not expect or recommend you to solve all of them, unless you think you are struggling with a particular concept, or would like to do practice runs on these old exam questions.

7 Vector Processing (Extra)

Assume a vector processor that implements the following ISA:

Opcode	Operands	Number of cycles	Description
SET	Vst, imm	1	$Vst \leftarrow \text{imm}$ (Vst: Vector Stride Register)
SET	Vln, imm	1	$Vln \leftarrow \text{imm}$ (Vln: Vector Length Register)
VLD	Vd, addr	100, pipelined	$Vd \leftarrow \text{MEM}[\text{addr}]$
VST	Vs, addr	100, pipelined	$\text{MEM}[\text{addr}] \leftarrow Vs$
VADD	Vd, Vs, Vt	5, pipelined	$Vd \leftarrow Vs + Vt$
VMUL	Vd, Vs, Vt	10, pipelined	$Vd_i \leftarrow Vs_i \times Vt_i$
VDIV	Vd, Vs, Vt	20, pipelined	$Vd_i \leftarrow Vs_i / Vt_i$

Assume the following:

- The processor has an in-order pipeline.
 - The size of a vector element is 4 bytes.
 - Vst and Vln are 10-bit registers.
 - The processor *does not* support chaining between vector functional units.
 - The main memory has N banks.
 - Vector elements stored in consecutive memory addresses are interleaved between the memory banks. For example, if a vector element at address A maps to bank B , a vector element at address $A + 4$ maps to bank $(B + 1) \% N$, where $\%$ is the modulo operator and N is the number of banks. N is *not necessarily* a power of two.
 - The memory is byte addressable and the address space is represented using 32 bits.
 - Vector elements are stored in memory in 4-byte-aligned manner.
 - Each memory bank has a 4 KB row buffer.
 - Each memory bank has a single read and a single write port so that a load and a store operation can be performed simultaneously.
 - There are separate functional units for executing VLD and VST instructions.
- (a) What should the minimum value of N be to avoid stalls while executing a VLD or VST instruction, assuming a vector stride of 1? Explain.

100 banks (because the latency of VLD and VST instructions is 100 cycles)

- (b) What should the minimum value of N be to avoid stalls while executing a VLD or VST instruction, assuming a vector stride of 2? Explain.

101 banks.

To avoid stalls, we need to ensure that consecutive vector elements access 100 different banks.

With a vector stride of 2, consecutive elements of a vector will map to every other bank. For example, if the first element maps to bank 0, the next element will map to bank 2, and so on.

With 100 banks, the 51st element of a vector will map to bank $100\%100 = 0$, conflicting with the first element of the vector.

However, with 101 banks, the 51st element will map to bank 1, which was skipped by the previous vector elements.

Let's assume there are 102 elements in a vector and the first elements accesses bank 0. The 101 banks will be accessed in the following order:

$(0, 2, \dots, 100, 102, 104, \dots, 200, 202)\%101 = (0, 2, \dots, 100, 1, 3, \dots, 99, 0)$

We can see that none of the elements conflict in the DRAM banks. Note that, when the last vector elements accesses bank 0, the bank is already available for a new access because the 100 cycle latency of accessing the first element is overlapped by accessing the other 101 elements.

- (c) Assume:

- A machine that has a memory with as many banks as you found in part (a).
- The vector stride is set to 1.
- The value of the vector length is set to M (but we do *not* know M)

The machine executes the following program:

```
VLD  V1, A           // V1 ← MEM[A]
VLD  V2, (A + 32768)  // V2 ← MEM[A + 32768]
VADD V3, V1, V1       // V3 ← V1, V1
VMUL V4, V2, V3       // V4i ← V2i, V3i
VST  V4, (A + 32768*2) // MEM[A + 32768*2] ← V4
```

It takes 4,306 cycles to execute the above program. What is M ? Explain.

$M = 1000$

```

VLD  |--100--|--(M-1)--|
VLD  |--100--|--(M-1)--|
VADD  |--5--|--(M-1)--|
VMUL  |--10--|--(M-1)--|
VST   |--100--|--(M-1)--|

```

$(M + 100 - 1) + 100 + (M - 1) + 10 + (M - 1) + 100 + (M - 1) = 306 + 4 * M = 4306 \rightarrow M = 1000$
elements

- (d) If we modify the vector processor to *support chaining*, how many cycles would be required to execute the same program in part (c)? Explain.

2,308 cycles

Let L the vector length, i.e., the number of elements of a vector.

```

VLD  |--100--|--(L-1)--|
VLD  |--100--|--(L-1)--|
VADD  |--1--|--5--|--(L-1)--| (this is delayed because the processor
                                executes the instructions in order)
VMUL  |--10--|--(L-1)--|
VST   |--100--|--(L-1)--|

```

$100 + (L - 1) + 100 + 10 + 100 + (L - 1) = 310 + 2 \times 1000 - 2 = 2308$ cycles

8 Vector Processing (Extra)

A vector processor includes the following instructions in its ISA:

Opcode	Operands	Latency (cycles)	Description
VLD	$V_i, \#Address$	60, fully interleaved and pipelined	$V_i \leftarrow Mem[Address]$
VST	$V_i, \#Address$	60, fully interleaved and pipelined	$Mem[Address] \leftarrow V_i$
VADD	V_i, V_j, V_k	8, fully pipelined	$V_i \leftarrow V_j + V_k$
VMUL	V_i, V_j, V_k	16, fully pipelined	$V_i \leftarrow V_j * V_k$

Assume that:

- The machine has an in-order pipeline.
- Each vector element is 64 bits in size.
- In order to support 1-element per cycle memory throughput for vector elements, after the first element in a vector, the machine interleaves vector elements across memory banks. All vectors are stored in memory with the first element mapped to bank 0, the second element mapped to bank 1, etc.
- Memory accesses within a vectorized memory request must be issued in order.
- Each memory bank has two ports (so that two loads/stores can be active simultaneously), and there are two load/store functional units available.

Answer the following questions about this vector computer.

- (a) The number of memory banks in this vector processor is a power of two. What should the minimum number of banks be to avoid stalls while executing a VLD or VST instruction, assuming a vector stride of 1? Explain.

64 banks, because the memory latency is 60 cycles and the next power of two is 64.

- (b) Translate the following piece of code, called SAXPY, into vector code that executes in the minimum possible number of cycles on this vector processor. Constant **alpha** is stored in vector register **V1**. (Note: Assume vector length register and vector stride register are already introduced to appropriate values.)

```
for (i = 0; i < VLEN; i++){
    temp = alpha * X[i];
    Z[i] = temp + Y[i];
}
```

```
VLD V2, X
VLD V3, Y
VMUL V4, V1, V2
VADD V5, V4, V3
VST Z, V5
```

- (c) Consider a machine with as many banks as you found in part (a) executes the program you wrote in part (b). Assume that the vector stride is set to 1. The machine does *not* support chaining between vector functional units. If the machine takes 340 cycles to finish the program, what is the vector length? Show your work.

50.

Explanation.

```

VLD    |---60---|--(VLEN-1)--|
VLD    |1|---60---|--(VLEN-1)--|
VMUL           |--16--|--(VLEN-1)--|
VADD                                |-8-|--(VLEN-1)--|
VST                                     |---60---|--(VLEN-1)--|

```

$$60 + 16 + 8 + 60 + (VLEN-1) \times 4 = 340$$

$$VLEN = 50$$

- (d) The architect of the machine improves the design to support chaining. For the VLEN you found in the previous part, what is the total number of cycles for the same program? Show your work.

193.

Explanation.

```

VLD    |---60---|--(VLEN-1)--|
VLD    |1|---60---|--(VLEN-1)--|
VMUL           |--16--|--(VLEN-1)--|
VADD                                |-8-|--(VLEN-1)--|
VST                                     |---60---|--(VLEN-1)--|

```

$$60 + 16 + 8 + 60 + (VLEN-1) = 193$$

- (e) Now the architect decides that she needs to cut costs in the machine's memory system by reducing the number of banks. Because loads and stores might stall due to bank contention, an arbiter is added to each bank so that pending loads from the oldest instruction are serviced first. She observes the program now takes 481 cycles (with chaining). What is the new number of banks? Show your work. (Note: Recall the number of banks is a power of two.)

16.

Explanation.

```

VLD
[0] |---60---| bank0
...
[15] |-15-|---60---| bank15
[16] |---60---| bank0
...
[31] |-15-|---60---| bank15
[32] |---60---| bank0
...
[47] |-15-|---60---| bank15
[48] |---60---| bank0
[49] |1|---60---| bank1
VLD
[0] |1|---60---| bank0 (port1)
...
[15] |1|-15-|---60---| bank15 (port1)
[16] |---60---| bank0 (port1)
...
[31] |-15-|---60---| bank15 (port1)
[32] |---60---| bank0 (port1)
...
[47] |-15-|---60---| bank15 (port1)
[48] |---60---| bank0 (port1)
[49] |1|---60---| bank1 (port1)
VMUL
[0] |--16--|
...
[15] |-15-|--16--|
[16] |--16--|
...
[31] |-15-|--16--|
[32] |--16--|
...
[47] |-15-|--16--|
[48] |--16--|
[49] |1|--16--|
VADD
[0] |-8-|
...
[15] |-15-|-8-|
[16] |-8-|
...
[31] |-15-|-8-|
[32] |-8-|
...
[47] |-15-|-8-|
[48] |-8-|
[49] |1|-8-|
VST
[0] |---60---| bank 0 (it cannot start until port 0 is free)
...
[15] |-15-|---60---| bank 15
[16] |---60---| bank 0
...
[31] |-15-|---60---| bank 15
[32] |---60---| bank 0
...
[47] |-15-|---60---| bank 15
[48] |---60---| bank 0
[49] |1|---60---| bank 1

```

9 GPUs and SIMD (Extra)

We define the *SIMD utilization* of a program that runs on a GPU as the fraction of SIMD lanes that are kept busy with *active threads* during the run of the program. As we saw in lecture and practice exercises, the SIMD utilization of a program is computed across the *complete run* of the program.

The following code segment is run on a GPU. Each thread executes **a single iteration** of the shown loop. Assume that the data values of the arrays A and B are already in vector registers, so there are no loads and stores in this program. (Hint: Notice that there are 3 instructions in each iteration.) A warp in the GPU consists of 32 threads, and there are 32 SIMD lanes in the GPU.

```
for (i = 0; i < 1025; i++) {  
    if (A[i] < 33) {           // Instruction 1  
        B[i] = A[i] << 1;     // Instruction 2  
    }  
    else {  
        B[i] = A[i] >> 1;     // Instruction 3  
    }  
}
```

Please answer the following six questions.

- (a) How many warps does it take to execute this program?

33 warps

Number of Warps = $\lceil (\text{Number of threads}) / (\text{Number of threads per warp}) \rceil$

Number of threads = $2^{10} + 1$ (i.e., one thread per loop iteration).

Number of threads per warp = $32 = 2^5$ (given).

\therefore Number of Warps = $\lceil \frac{2^{10}+1}{2^5} \rceil = 2^5 + 1 = 33$

- (b) What is the *maximum* possible SIMD utilization of this program? (Hint: The warp scheduler does *not* issue instructions when *no* threads are active).

$\frac{1025}{1056}$

Even though all active threads in a warp follow the same execution path, the last warp will only have one active thread.

- (c) Please describe what needs to be true about array **A** to reach the maximum possible SIMD utilization asked in part (b). (Please cover all cases in your answer.)

For every 32 consecutive elements of **A**, every element should be lower than 33 (**if**), or greater than or equal to 33 (**else**). (NOTE: The solution is correct if both cases are given.)

- (d) What is the *minimum* possible SIMD utilization of this program?

$$\frac{1025}{1568}$$

Instruction 1 is executed by every active thread ($\frac{1025}{1056}$ utilization). Then, part of the threads in each warp executes **Instruction 2** and the other part executes **Instruction 3**. We consider that **Instruction 2** is executed by α threads in each warp (except the last warp), where $0 < \alpha \leq 32$, and **Instruction 3** is executed by the remaining $32 - \alpha$ threads. The only active thread in the last warp executes either **Instruction 2** or **Instruction 3**. The other instruction is not issued for this warp.

The minimum SIMD utilization sums to $\frac{1025 + \alpha \times 32 + (32 - \alpha) \times 32 + 1}{1056 + 1024 + 1024 + 32} = \frac{1025}{1568}$.

- (e) Please describe what needs to be true about array **A** to reach the minimum possible SIMD utilization asked in part (d). (Please cover all cases in your answer.)

For every 32 consecutive elements of **A**, part of the elements should be lower than 33 (**if**), and the other part should be greater than or equal to 33 (**else**).

- (f) What is the SIMD utilization of this program if $A[i] = i$? Show your work.

$$\frac{1025}{1072}$$

Instruction 1 is executed by every active thread ($\frac{1025}{1056}$ utilization).

Instruction 2 is executed by the first 33 threads, i.e., all threads in the first warp and one thread in the second warp.

Instruction 3 is executed by the remaining active threads.

The SIMD utilization sums to $\frac{1025+32+1+31+960+1}{1056+32+32+32+960+32} = \frac{2050}{2144} = \frac{1025}{1072}$.

10 GPUs and SIMD (Extra)

We define the *SIMD utilization* of a program run on a GPU as the fraction of SIMD lanes that are kept busy with *active threads* during the run of a program.

The following code segment is run on a GPU. Each thread executes **a single iteration** of the shown loop. Assume that the data values of the arrays A, B, and C are already in vector registers so there are no loads and stores in this program. (Hint: Notice that there are 6 instructions in each thread.) A warp in the GPU consists of 64 threads, and there are 64 SIMD lanes in the GPU.

```
for (i = 0; i < 4096; i++) {  
    if (B[i] < 8888) {  
        A[i] = A[i] * C[i];  
        A[i] = A[i] + B[i];  
        C[i] = B[i] + 1;  
    }  
    if (B[i] > 8888) {  
        A[i] = A[i] * B[i];  
    }  
}
```

- (a) How many warps does it take to execute this program?

Warps = (Number of threads) / (Number of threads per warp)
Number of threads = 2^{12} (i.e., one thread per loop iteration).
Number of threads per warp = $64 = 2^6$ (given).
Warps = $2^{12}/2^6 = 2^6$

- (b) When we measure the SIMD utilization for this program with one input set, we find that it is 134/320. What can you say about arrays A, B, and C? Be precise (Hint: Look at the “if” branch).

A:

Nothing

B:

2 in every 64 of B's elements are less than 8888, the rest are 8888

C:

Nothing

(c) Is it possible for this program to yield a SIMD utilization of 100% (circle one)?

YES

NO

If YES, what should be true about arrays A, B, C for the SIMD utilization to be 100%? Be precise. If NO, explain why not.

Yes. All consecutive 64 elements of B should be either:

- (1) All of B's elements are equal to 8888, or
- (2) All of B's elements are less than 8888, or
- (3) All of B's elements are greater than 8888.

(d) What is the lowest SIMD utilization that this program can yield? Explain.

132/384. 1 in every 64 of B's elements are greater than 8888, and 1 in every 64 of B's elements are less than 8888, and the rest of the elements are 8888.

11 GPUs and SIMD (Extra)

We define the *SIMD utilization* of a program that runs on a GPU as the fraction of SIMD lanes that are kept busy with *active threads* during the run of a program. As we saw in lecture and practice exercises, the SIMD utilization of a program is computed across the *complete run* of the program.

The following code segment is run on a GPU. A warp in the GPU consists of 32 threads, and there are 32 SIMD lanes in the GPU. Each thread executes **a single iteration** of the shown loop. Assume that the data values of the arrays A and B are already in vector registers so there are no loads and stores in this program. Both A and B are arrays of integers. (Hint: Notice that there are 5 instructions in each iteration.)

```
for (i = 0; i < 1024; i++) {  
    A[i] = 2 * B[i];           // Instruction 1  
    if (A[i / 32] == 0) {      // Instruction 2  
        int a = 2 * A[i];      // Instruction 3  
        if (a == i) {          // Instruction 4  
            A[i] += 1;         // Instruction 5  
        }  
    }  
}
```

Please answer the following questions.

- (a) How many warps does it take to execute this program?

32 warps.

Explanation:

The number of warps is calculated as:

$$\#Warp s = \lceil \frac{\#Total_threads}{\#Warp_size} \rceil,$$

where

$$\#Total_threads = 1024 = 2^{10} \text{ (i.e., one thread per loop iteration),}$$

and

$$\#Warp_size = 32 = 2^5 \text{ (given).}$$

Thus, the number of warps needed to run this program is:

$$\#Warp s = \lceil \frac{2^{10}}{2^5} \rceil = 2^5 = 32.$$

- (b) What needs to be true about array B to achieve 100% utilization? Show your work. (Hint: The warp scheduler does not issue instructions where no threads are active).

Elements $B[x]$, with x between 0 and 31, should be different from 0.

Explanation:

Threads with an ID i enter the first `if` statement if the element of B with index $\lfloor \frac{i}{32} \rfloor$ (i.e., the first 32 elements of B) is equal to 0. If that is *not* the case, all threads execute only **Instruction 1** and **Instruction 2**. Thus, the utilization is 100%.

- (c) What is the minimum possible SIMD utilization of this program?

$$\text{Minimum } SIMD_utilization = \frac{4096+32}{4096+1024} = \frac{4128}{5120}.$$

- (d) What needs to be true about array B to achieve the minimum possible SIMD utilization? Show your work. (Please cover all cases in your answer.)

The minimum SIMD utilization is achieved when there is one (and only one) thread per warp that executes **Instruction 5**. Only one thread every four threads can potentially execute **Instruction 5**, since index **i** needs to be a multiple of 4 for a thread to enter the second **if** statement. Thus, we know about array B:

- B[x], with x between 1 and 31, is equal to 0.
- For $i \geq 32$ and i a multiple of 4, only 1 out of 8 values of B fulfills $4 * B[i] = i$.

- (e) If B[0] = 0, what is the maximum possible SIMD utilization of this program?

$$\text{Maximum } SIMD_utilization = \frac{2048+32+32+8}{2048+32+32+32} = \frac{2120}{2144}.$$

- (f) What needs to be true about array `B` to achieve the maximum possible SIMD utilization, if `B[0] = 0`? Show your work. (Please cover all cases in your answer.)

The maximum SIMD utilization is achieved when only the first warp enters the first `if` statement. In this case, the maximum number of threads that can execute `Instruction 5` is 8. This can only happen if:

- `B[0] = 0`;
- `B[x]`, with `x` between 1 and 31, is different from 0.
- `B[4 * x]`, with `x` between 1 and 7, is equal to `x`.

12 GPUs and SIMD (Extra)

We define the *SIMD utilization* of a program that runs on a GPU as the fraction of SIMD lanes that are kept busy with *active threads* during the run of a program.

The following code segments are run on a GPU. We assume that (1) **A** resides in memory and is shared by all threads, (2) **s** resides in a register and is private to each thread, and (3) the code segments are correct (i.e., do not think about any correctness issues when answering this question).

A warp in the GPU consists of 32 threads, and there are 32 SIMD lanes in the GPU. Each thread executes **a single iteration** of the outermost loop (with index **i**). Assume that the data values of the array **A** are already in vector registers so there are no memory loads and stores in this program. (Hint: Notice that there are 4 instructions in each iteration of the outermost loop of both code segments.)

```
s = 1;
for (i = 0; i < 1024; i++) {
    for (j = 0; j < 10; j++) { // Inst. 1
        if (i % (2 * s) == 0) // Inst. 2
            A[i] += A[i + 1]; // Inst. 3
        s = s << 1;           // Inst. 4
    }
}
```

Code Segment 1

```
s = 512;
for (i = 0; i < 1024; i++) {
    for (j = 0; j < 10; j++) { // Inst. 1
        if (i < s)             // Inst. 2
            A[i] += A[i + s]; // Inst. 3
        s = s >> 1;           // Inst. 4
    }
}
```

Code Segment 2

Please answer the following questions.

- (a) How many warps does it take to execute these code segments?

32 warps.

Explanation:

The number of warps is calculated as:

$$\#Warp s = \lceil \frac{\#Total_threads}{\#Warp_size} \rceil,$$

where

$$\#Total_threads = 1024 = 2^{10} \text{ (i.e., one thread per loop iteration),}$$

and

$$\#Warp_size = 32 = 2^5 \text{ (given).}$$

Thus, the number of warps needed to run this program is:

$$\#Warp s = \lceil \frac{2^{10}}{2^5} \rceil = 2^5 = 32.$$

- (b) What is the SIMD utilization of the first iteration of the inner loop ($j = 0$) for Code Segment 1? Show your work. (Hint: The warp scheduler does *not* issue instructions when no thread is active).

The utilization of the first iteration ($j = 0$) of Code Segment 1 is $\frac{7}{8}$.

Explanation:

Instructions 1, 2, and 4 are executed by all threads in Code Segment 1.

In Code Segment 1, $s = 1$ during the first iteration. Thus, only even numbered threads fulfill the predicate of the `if` statement, and only half of the threads of each warp execute Instruction 3.

Code Segment 1, $j = 0$: $SIMD_utilization = \frac{1024+1024+512+1024}{1024+1024+1024+1024} = \frac{7}{8}$.

- (c) What is the SIMD utilization of the first iteration of the inner loop ($j = 0$) for Code Segment 2? Show your work. (Hint: The warp scheduler does *not* issue instructions when no thread is active).

The utilization of the first iteration ($j = 0$) of Code Segment 2 is 100%.

Explanation:

Instructions 1, 2, and 4 are executed by all threads in Code Segment 2.

In Code Segment 2, $s = 512$ during the first iteration. Thus, only threads with $i < 512$ fulfill the predicate of the `if` statement, and all threads of only half of the warps execute Instruction 3.

Code Segment 2, $j = 0$: $SIMD_utilization = \frac{1024+1024+512+1024}{1024+1024+512+1024} = \frac{7}{7} = 100\%$.

- (d) What is the SIMD utilization of any iteration of the inner loop ($0 \leq j < 10$) for Code Segment 1? Show your work. (Hint: Derive an analytical expression, which may be piecewise).

As mentioned in part b, Instructions 1, 2, and 4 are executed by all threads.

In Code Segment 1, with $0 \leq j < 5$, all 32 warps are active, but the number of active threads per warp divides by half in each iteration. With $5 \leq j < 10$, only one thread per warp is active, and the number of active warps divides by half in each iteration. As a result:

Code Segment 1, iteration j :

$$SIMD_utilization = \begin{cases} \frac{3072+2^{(9-j)}}{4096}, & \text{if } 0 \leq j < 5 \\ \frac{3072+2^{(9-j)}}{3072+32 \cdot 2^{(9-j)}}, & \text{if } 5 \leq j < 10 \end{cases} \quad (1)$$

- (e) What is the SIMD utilization of any iteration of the inner loop ($0 \leq j < 10$) for Code Segment 2? Show your work. (Hint: Derive an analytical expression, which may be piecewise).

As mentioned in part b, Instructions 1, 2, and 4 are executed by all threads.

In Code Segment 2, with $0 \leq j < 5$, all 32 threads per warp are active, but the number of active warps divides by half in each iteration. With $5 \leq j < 10$, only one warp is active, and the number of active threads divides by half in each iteration. As a result:

Code Segment 2, iteration j :

$$SIMD_utilization = \begin{cases} \frac{3072+32 \cdot 2^{(4-j)}}{3072+32 \cdot 2^{(4-j)}} = 100\%, & \text{if } 0 \leq j < 5 \\ \frac{3072+2^{(9-j)}}{3072+32}, & \text{if } 5 \leq j < 10 \end{cases} \quad (2)$$

- (f) Is there any iteration ($0 \leq j < 10$) where both code segments have the same utilization? Explain your reasoning.

Yes, with $j = 9$ only one thread of only one warp is active, since only one thread (out of 1024) is needed to perform the last addition.

- (g) Which code is expected to run faster on a GPU? Explain your reasoning.

Code Segment 2 is faster because it has less intra-warp divergence, and thus higher SIMD utilization. In each iteration (except the last one), the number of warps that Code Segment 2 schedules is smaller than the number of warps that Code Segment 1 schedules. This results in fewer execution cycles.