

Exercise 4.9 a)**ANS**

This code reads four SP floats of 16 bytes

writes two SP floats of 8 bytes for every six FLOPs.

So arithmetic intensity = $6/24$

$$= 0.25 \text{ FLOP per byte of data accessed}$$

Exercise 4.10**ANS**

For the Vector Processor,

Computation Time = Computation Time for scalar execution + Memory Access time

$$= 400 \text{ ms} + (200 + 100) \text{ MB} / (30 \text{ GB/sec})$$

$$= 410 \text{ ms}$$

For the Hybrid Processor,

Computation Time = Computation Time for scalar execution + Memory Access time +
Transfer time between host and local memory + Memory latency

$$= 400 \text{ ms} + (200 + 100) \text{ MB} / (150 \text{ GB/sec}) + (200 + 100) \text{ MB} / (10 \text{ GB/sec}) + 10 \text{ ms}$$

$$= 442 \text{ ms}$$

So the vector processor achieves better performance than the hybrid processor.

Exercise 4.13

a)

ANS

$$\begin{aligned}\text{GFLOPs/sec} &= 1.5 \times 10 \times 0.8 \times 0.85 \times 0.7 \times (32/4) \\ &= 57.12\end{aligned}$$

b)

ANS

1)

$$\begin{aligned}\text{If we increase the lanes by 16, GFLOPs/sec} &= 1.5 \times 10 \times 0.8 \times 0.85 \times 0.7 \times (32/2) \\ &= 114.24\end{aligned}$$

$$\begin{aligned}\text{Speedup} &= 114.24/57.12 \\ &= 2\end{aligned}$$

2)

$$\begin{aligned}\text{If we increase the number of SIMDs to 15, GFLOPs/sec} &= 1.5 \times 15 \times 0.8 \times 0.85 \times 0.7 \times (32/4) \\ &= 85.68\end{aligned}$$

$$\begin{aligned}\text{Speedup} &= 85.68/57.12 \\ &= 1.5\end{aligned}$$

3)

$$\begin{aligned}\text{If we increase the issue rate to 0.95, GFLOPs/sec} &= 1.5 \times 15 \times 0.8 \times 0.95 \times 0.7 \times (32/4) \\ &= 63.84\end{aligned}$$

$$\begin{aligned}\text{Speedup} &= 63.84/57.12 \\ &= 1.12\end{aligned}$$

Exercise 4.16

ANS

The clock rate of a hypothetical GPU is 1.5 GHz, exists 16 SIMD processors, each processor contains 16 single-precision floating point units and off-chip memory bandwidth is 100 GB/sec.

For this GPUs the peak single-precision floating-point throughput is,

$$\begin{aligned} \text{core frequency} \times \text{number of cores} \times \text{number of operations per clock} &= 1.5 \times 16 \times 16 \\ &= 384 \text{ GFLOP/sec} \end{aligned}$$

Assuming the each single precision operation required 4 Byte 2 operands and output one four byte result, sustaining would required the memory bandwidth

$$= 12 \text{ Bytes/ Flop} \times 384 \text{ GFLOPS/sec}$$

$$= 4.608 \text{ TB/s}$$

Throughput is not sustainable because $4.608 \text{ TB/sec} > 100 \text{ GB/sec}$

But still can be achieved in short bursts when using on-chip memory