

NYU Tandon School of Engineering

Fall 2022, ECE 6913

Homework Assignment 2

Instructor: Azeez Bhavnagarwala, email: ajb20@nyu.edu

Course Assistants

Varadraj Kakodkar (vns2008), Kartikay Kaushik (kk4332), Siddhanth Iyer (si2152), Swarnashri Chandrashekar (sc8781), Karan Sheth (kk4332), Haotian Zheng (hz2687), Haoren Zhang (kk4332), Varun Kumar (vs2411)

Homework Assignment 2 [released Tuesday September 20th 2022] [due Friday September 30th by 11:59PM]

You *are allowed* to discuss HW assignments with anyone. You are *not allowed* to share your solutions with other colleagues in the class. Please feel free to reach out to the Course Assistants or the Instructor during office hours or by appointment if you need any help with the HW. Please enter your responses in this Word document after you download it from NYU Classes. *Please use the Brightspace portal to upload your completed HW.*

1. After graduating, you are asked to become the lead computer designer at Hyper Computers, Inc. Your study of usage of high-level language constructs suggests that procedure calls are one of the most expensive operations. You have invented a scheme that reduces the loads and stores normally associated with procedure calls and returns. The first thing you do is run some experiments with and without this optimization. Your experiments use the same state-of-the-art optimizing compiler that will be used with either version of the computer. These experiments reveal the following information:

- The clock rate of the unoptimized version is 5% higher.
- 30% of the instructions in the unoptimized version are loads or stores.
- The optimized version executes 2/3 as many loads and stores as the unoptimized version. For all other instructions the dynamic counts are unchanged.
- All instructions (including load and store) take one clock cycle.

Which is faster? Justify your decision quantitatively.

We can know that:

Total Execution time = Number of Instruction * CPI * Cycle Time = Number of Instruction * CPI / Clock Rate

According to the problem:

1. Unoptimized Clock Rate = 1.05 Optimized Clock Rate, Optimized Clock Rate = 95.24% Unoptimized Clock Rate

2. Optimized Instruction count = $(30\% * 2/3 + 70\%)$ Unoptimized Instruction count = 90% Unoptimized Instruction count

3. The optimized execution time = Optimized Instruction count * CPI / Optimized Clock Rate = $(90\% \text{ Unoptimized Instruction count}) * \text{CPI} / (95.24\% \text{ Unoptimized Clock Rate}) = 94.5\%$ (Unoptimized Instruction count * CPI / Unoptimized Clock Rate) = 94.5% Unoptimized execution time

Therefore, we can know the optimized version is faster since its execution time is 5.5% less than the unoptimized version.

2. General-purpose processes are optimized for general-purpose computing. That is, they are optimized for behavior that is generally found across a large number of applications. However, once the domain is restricted somewhat, the behavior that is found across a large number of the target applications may be different from general-purpose applications. One such application is deep learning or neural networks. Deep learning can be applied to many different applications, but the fundamental building block of inference—using the learned information to make decisions—is the same across them all. Inference operations are largely parallel, so they are currently performed on graphics processing units, which are specialized more toward this type of computation, and not to inference in particular. In a quest for more performance per watt, Google has created a custom chip using tensor processing units to accelerate inference operations in deep learning. This approach can be used for speech recognition and image recognition, for example. This problem explores the trade-offs between this process, a general-purpose processor (Haswell E5-2699 v3) and a GPU (NVIDIA K80), in terms of performance and cooling. If heat is not removed from the computer efficiently, the fans will blow hot air back onto the computer, not cold air. Note: The differences are more than processor—on-chip memory and DRAM also come into play. Therefore statistics are at a system level, not a chip level.

a. If Google's data center spends 70% of its time on workload A and 30% of its time on workload B when running GPUs, what is the speedup of the TPU system over the GPU system?

b. Google's data center spends 70% of its time on workload A and 30% of its time on workload B when running GPUs, what percentage of Max IPS does it achieve for each of the three systems?

c. Building on (b), assuming that the power scales linearly from idle to busy power as IPS grows from 0% to 100%, what is the performance per watt of the TPU system over the GPU system?

d. If another data center spends 40% of its time on workload A, 10% of its time on workload B, and 50% of its time on workload C, what are the speedups of the GPU and TPU systems over the general-purpose system?

e. A cooling door for a rack cost \$4000 and dissipates 14 kW (into the room; additional cost is required to get it out of the room). How many Haswell-, NVIDIA-, or Tensor-based servers can you cool with one cooling door, assuming TDP in Figures 1.27 and 1.28?

f. Typical server farms can dissipate a maximum of 200 W per square foot. Given that a server rack requires 11 square feet (including front and back clearance), how many servers from part (e) can be placed on a single rack, and how many cooling doors are required?

System	Chip	TDP	Idle power	Busy power
General-purpose	Haswell E5-2699 v3	504 W	159 W	455 W
Graphics processor	NVIDIA K80	1838 W	357 W	991 W
Custom ASIC	TPU	861 W	290 W	384 W

Figure 1.27 Hardware characteristics for general-purpose processor, graphical processing unit-based or custom ASIC-based system, including measured power

System	Chip	Throughput			% Max IPS		
		A	B	C	A	B	C
General-purpose	Haswell E5-2699 v3	5482	13,194	12,000	42%	100%	90%
Graphics processor	NVIDIA K80	13,461	36,465	15,000	37%	100%	40%
Custom ASIC	TPU	225,000	280,000	2000	80%	100%	1%

Figure 1.28 Performance characteristics for general-purpose processor, graphical processing unit-based or custom ASIC-based system on two neural-net workloads

a.

For workload A, the speedup of TPU over GPU is:

$$225000 / 13461 = 16.715$$

For workload B, the speedup of TPU over GPU is:

$$280000 / 36465 = 7.679$$

The total speedup of TPU is:

$$1 / (70\% / 16.715 + 30\% / 7.679) = 12.354$$

b.

$$\text{Haswell E5-2699 v3: } 70\% * 42\% + 30\% * 100\% = 59.4\%$$

$$\text{NVIDIA K80: } 70\% * 37\% + 30\% * 100\% = 55.9\%$$

$$\text{TPU: } 70\% * 80\% + 30\% * 100\% = 86\%$$

c.

Haswell E5-2699 v3: $159 + (455 - 159) * 59.4\% = 334.824 \text{ W}$

NVIDIA K80: $357 + (991 - 357) * 55.9\% = 711.406 \text{ W}$

TPU: $290 + (384 - 290) * 86\% = 370.84 \text{ W}$

The performance per watt of TPU over GPU is:

$$(86\% / 370.84 \text{ W}) / (55.9\% / 711.406 \text{ W}) = 2.951$$

d.

For workload A, the speedup of GPU over General-purpose is:

$$13461 / 5482 = 2.455$$

For workload A, the speedup of TPU over General-purpose is:

$$225000 / 5482 = 41.043$$

For workload B, the speedup of GPU over General-purpose is:

$$36465 / 13194 = 2.764$$

For workload B, the speedup of TPU over General-purpose is:

$$280000 / 13194 = 21.222$$

For workload C, the speedup of GPU over General-purpose is:

$$15000 / 12000 = 1.25$$

For workload C, the speedup of TPU over General-purpose is:

$$2000 / 12000 = 0.167$$

Therefore, the total speedup over general-purpose is

$$\text{NVIDIA K80: } 1 / (40\% / 2.455 + 10\% / 2.764 + 50\% / 1.25) = 1.669$$

$$\text{TPU: } 1 / (40\% / 41.043 + 10\% / 21.222 + 50\% / 0.167) = 0.332$$

e.

Number of Haswell E5-2699 v3: $14\text{kW} / 504 = 27.778$, floor to 27

Number of NVIDIA K80: $14\text{kW} / 1838 = 7.617$, floor to 7

Number of TPU: $14\text{kW} / 861 = 16.260$, floor to 16

f.

According to the problem, the dissipation of a single rack is:

$$11 * 200 = 2.2 \text{ kW}$$

Number of Haswell E5-2699 v3: $2.2\text{kW} / 504 = 4.365$, floor to 4

Number of NVIDIA K80: $2.2\text{kW} / 1838 = 1.197$, floor to 1

Number of TPU: $2.2\text{kW} / 861 = 2.555$, floor to 2

Besides, $14\text{kW} > 2.2\text{kW}$, so 1 cooling door is enough.

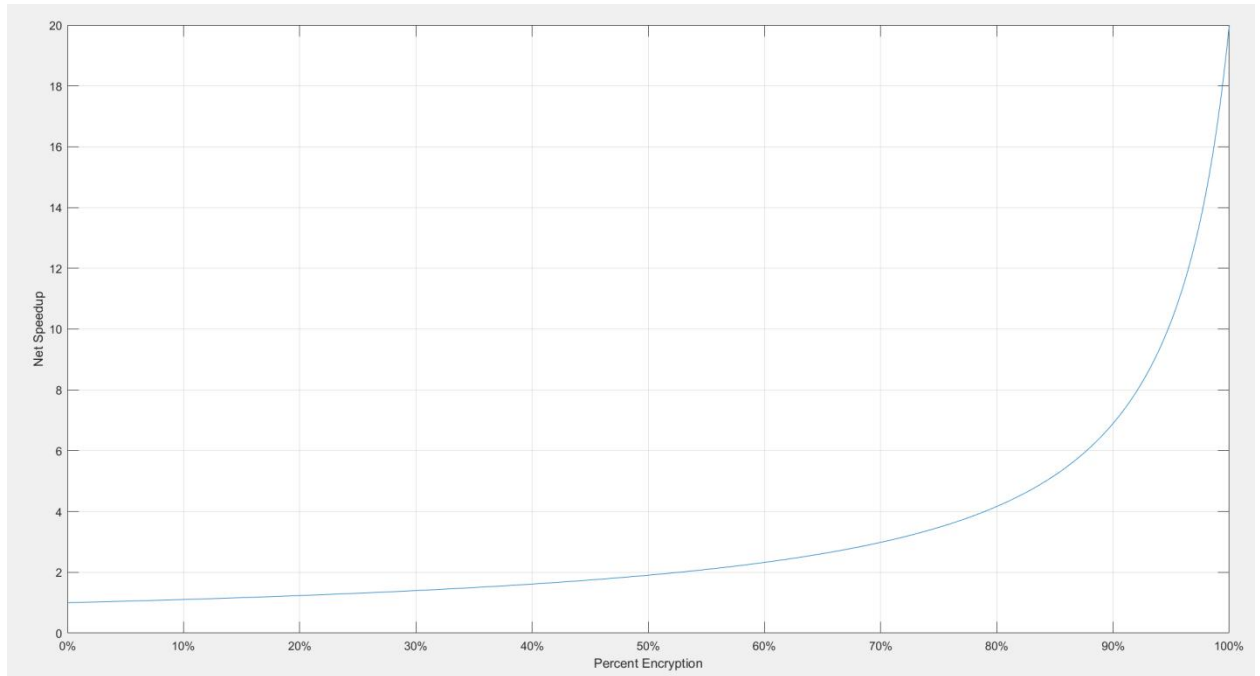
3. In this exercise, assume that we are considering enhancing a quad-core machine by adding encryption hardware to it. When computing encryption operations, it is 20 times faster than the normal mode of execution. We will define percentage of encryption as the percentage of time in the original execution that is spent performing encryption operations. The specialized hardware increases power consumption by 2%.

a. Draw a graph that plots the speedup as a percentage of the computation spent performing encryption. Label the y-axis “Net speedup” and label the x-axis “Percent encryption.”

Assume $x\%$ of work is encrypted, then the speed up is:

$$1 / (x\% / 20 + (1 - x\%))$$

The plot is:



b. With what percentage of encryption will adding encryption hardware result in a speedup of 2?

When the speedup is 2, we can have:

$$1 / (x\% / 20 + (1 - x\%)) = 2$$

We can solve that $x = 52.631$

Thus, 52.631% of encryption will result in a speedup of 2.

c. What percentage of time in the new execution will be spent on encryption operations if a speedup of 2 is achieved?

According to problem b., 52.631% of encryption will result in a speedup of 2. The percentage of time in the new execution will be spent on encryption is:

$$(52.631\% * 1 / 20) / ((52.631\% * 1 / 20) + (1 - 52.631\%)) = 5.263\%$$

4. Assume that we make an enhancement to a computer that improves some mode of execution by a factor of 10. Enhanced mode is used 50% of the time, measured as a percentage of the execution time when the enhanced mode is in use. Recall that Amdahl's Law depends on the fraction of the original, unenhanced execution time that could make use of enhanced mode. Thus, we cannot directly use this 50% measurement to compute speedup with Amdahl's Law.

a. What is the speedup we have obtained from fast mode?

According to the problem, now we have 50% time executed in unenhanced mode and 50% time executed in enhanced mode. If the enhanced part of work is executed in unenhanced mode, it will require 500% time to finish.

Thus, the speed up is:

$$(50\% + 500\%) / (50\% + 50\%) = 5.5$$

b. What percentage of the original execution time has been converted to fast mode?

Assume x% of original execution time has been converted to fast mode:

$$\text{The overall speedup} = 1 / (x\% / 10 + (1 - x\%)) = 5.5$$

Thus, we can solve that x% = 90.91%

90.91% of the original execution time has been converted to fast mode.

5. When parallelizing an application, the ideal speedup is speeding up by the number of processors. This is limited by two things: percentage of the application that can be parallelized and the cost of communication. Amdahl's Law takes into account the former but not the latter.

a. What is the speedup with N processors if 80% of the application is parallelizable, ignoring the cost of communication?

Assume the original execution time is T.

The speedup is:

$$T / (20\% * T + 80\% * T / N) = 1 / (0.2 + 0.8 / N)$$

b. What is the speedup with eight processors if, for every processor added, the communication overhead is 0.5% of the original execution time.

The speedup with communication time is:

$$T / (20\% * T + 80\% * T / 8 + 0.5\% * T * 8) = 1 / 0.34 = 2.941$$

c. What is the speedup with eight processors if, for every time the number of processors is doubled, the communication overhead is increased by 0.5% of the original execution time?

Assume the original communication time is 0. After 1 processor is added, the communication time is still 0. If 2 processors are added, the communication time is 0.5%.

The speedup with increased communication time is:

$$T / (20\% * T + 80\% * T / 8 + 0.5\% * 3 * T) = 1 / 0.46 = 3.175$$

d. What is the speedup with N processors if, for every time the number of processors is doubled, the communication overhead is increased by 0.5% of the original execution time?

Assume the original communication time is 0. After 1 processor is added, the communication time is 0.5% of execution time.

The speedup with increased communication time is:

$$T / (20\% * T + 80\% * T / N + 0.5\% * \text{floor}(\log_2(N)) * T)$$

e. Write the general equation that solves this question: What is the number of processors with the highest speedup in an application in which P% of the original execution time is parallelizable, and, for every time the number of processors is doubled, the communication is increased by 0.5% of the original execution time?

Here assume the percentage P% is a constant number.

The speedup is:

$$T / ((1 - P\%) * T + P\% * T / N + 0.5\% * \text{floor}(\log_2(N)) * T) = 1 / ((1 - P\%) + P\% / N + 0.5\% * \log_2(N))$$

The optimal number of processors should make the partial derivative of speedup equals to 0, which is:

$$\frac{d}{dN} \left(\frac{1}{(1 - P\%) + \log_2 N * 0.5\% + \frac{P\%}{N}} \right) = 0$$

Solve this we can have:

$$N = P * 2 \ln 2$$

6. Your company has just bought a new 22-core processor, and you have been tasked with optimizing your software for this processor. You will run four applications on this system, but the resource requirements are not equal. Assume the system and application characteristics listed in Table 1.1 below (from textbook)

Table 1.1 Four applications

Application	A	B	C	D
% resources needed	41	27	18	14
% parallelizable	50	80	60	90

The percentage of resources of assuming they are all run in serial. Assume that when you parallelize a portion of the program by X, the speedup for that portion is X.

a. How much speedup would result from running application A on the entire 22-core processor, as compared to running it serially?

50% of A can be parallelizable, so the speed up of that portion is 50.

The total speedup of A is:

$$1 / (50\% / 50 + 50\%) = 1.961$$

b. How much speedup would result from running application D on the entire 22-core processor, as compared to running it serially?

90% of D can be parallelizable, so the speed up of that portion is 90.

The total speedup of D is:

$$1 / (90\% / 90 + 10\%) = 9.091$$

c. Given that application A requires 41% of the resources, if we statically assign it 41% of the cores, what is the overall speedup if A is run parallelized but everything else is run serially?

$$41\% * 22 = 9.02, 9 \text{ cores are statically assigned}$$

The overall speedup is:

$$1 / (41\% * (50\% / 9 + 50\%) + 27\% + 18\% + 14\%) = 1.223$$

d. What is the overall speedup if all four applications are statically assigned some of the cores, relative to their percentage of resource needs, and all run parallelized?

$$41\% * 22 = 9.02, 9 \text{ cores are statically assigned to A}$$

$$27\% * 22 = 5.94, 6 \text{ cores are statically assigned to B}$$

$$18\% * 22 = 3.96, 4 \text{ cores are statically assigned to C}$$

$14\% * 22 = 3.08$, 3 cores are statically assigned to D

The overall speedup is:

$$1 / (41\% * (50\% / 9 + 50\%) + 27\% * (80\% / 6 + 20\%) + 18\% * (60\% / 4 + 40\%) + 14\% * (90\% / 3 + 10\%)) = 2.115$$

e. Given acceleration through parallelization, what new percentage of the resources are the applications receiving, considering only active time on their statically-assigned cores?

$$\text{Resources for A: } 41\% * (50\% / 9 + 50\%) = 22.78\%$$

$$\text{Resources for B: } 27\% * (80\% / 6 + 20\%) = 9\%$$

$$\text{Resources for C: } 18\% * (60\% / 4 + 40\%) = 9.9\%$$

$$\text{Resources for D: } 14\% * (90\% / 3 + 10\%) = 5.6\%$$

7. When making changes to optimize part of a processor, it is often the case that speeding up one type of instruction comes at the cost of slowing down something else. For example, if we put in a complicated fast floating-point unit, that takes space, and something might have to be moved farther away from the middle to accommodate it, adding an extra cycle in delay to reach that unit. The basic Amdahl's Law equation does not take into account this trade-off.

a. If the new fast floating-point unit speeds up floating-point operations by, on average, 2x, and floating-point operations take 20% of the original program's execution time, what is the overall speedup (ignoring the penalty to any other instructions)?

The overall speedup is :

$$1 / (20\% / 2 + 80\%) = 1.111$$

b. Now assume that speeding up the floating-point unit slowed down data cache accesses, resulting in a 1.5x slowdown (or 2/3 speedup). Data cache accesses consume 10% of the execution time. What is the overall speedup now?

The overall speedup is :

$$1 / (20\% / 2 + 10\% / (2 / 3) + 70\%) = 1.0526$$

c. After implementing the new floating-point operations, what percentage of execution time is spent on floating-point operations? What percentage is spent on data cache accesses?

The percentage of execution time spent on floating-point operations is:

$$(20\% / 2) / (20\% / 2 + 10\% / (2 / 3) + 70\%) = 10.53\%$$

The percentage of execution time spent on data cache access is:

$$(10\% / (2 / 3)) / (20\% / 2 + 10\% / (2 / 3) + 70\%) = 15.79\%$$

8. When making changes to optimize part of a processor, it is often the case that speeding up one type of instruction comes at the cost of slowing down something else. For example, if we put in a complicated fast floating-point unit, that takes space, and something might have to be moved farther away from the middle to accommodate it, adding an extra cycle in delay to reach that unit. The basic Amdahl's Law equation does not take into account this trade-off.

a. If the new fast floating-point unit speeds up floating-point operations by, on average, 2x, and floating-point operations take 20% of the original program's execution time, what is the overall speedup (ignoring the penalty to any other instructions)?

The overall speedup is :

$$1 / (20\% / 2 + 80\%) = 1.111$$

b. Now assume that speeding up the floating-point unit slowed down data cache accesses, resulting in a 1.5x slowdown (or 2/3 speedup). Data cache accesses consume 10% of the execution time. What is the overall speedup now?

The overall speedup is :

$$1 / (20\% / 2 + 10\% / (2 / 3) + 70\%) = 1.0526$$

c. After implementing the new floating-point operations, what percentage of execution time is spent on floating-point operations? What percentage is spent on data cache accesses?

The percentage of execution time spent on floating-point operations is:

$$(20\% / 2) / (20\% / 2 + 10\% / (2 / 3) + 70\%) = 10.53\%$$

The percentage of execution time spent on data cache access is:

$$(10\% / (2 / 3)) / (20\% / 2 + 10\% / (2 / 3) + 70\%) = 15.79\%$$