

3. NOVEMBER 2023

PROBLEM SET 1

Computing at Scale in Machine Learning: Distributed
computing and algorithmic approaches
WS_23/24

Submitted by:

Group ABA consisting of:

Kalson, Akshay

Biswal, Anuj

Hakari, Bhavanishankar Ramesh

Problem 2

2a)

In the provided code in the file pi-montecarlo.py, we look at the code step by step and see whether the program is executed sequentially or there's some level of parallelisation in the code, Here's a function-wise report on whether the processing within each function is sequential:

-sample_pi(n):

```
1  import argparse # See https://docs.python.org/3/library/argparse.html
2  import random
3  from math import pi
4
5  def sample_pi(n):
6      """ Perform n steps of Monte Carlo simulation for estimating Pi/4.
7          Returns the number of successes. """
8      s = 0
9      for i in range(n):
10         x = random.random()
11         y = random.random()
12         if x**2 + y**2 <= 1.0:
13             s += 1
14     return s
```

Function Purpose: This function performs a Monte Carlo simulation for estimating $\pi/4$ by generating random points and counting how many fall within a quarter of the unit circle.

Sequential Processing: The code inside this function is entirely sequential. It generates random points one at a time within a loop and checks whether each point is inside the quarter of the unit circle. This process is inherently sequential.

-compute_pi(args):

```

17 def compute_pi(args):
18     random.seed(1)
19
20     n_total = args.steps
21     s_total = sample_pi(n_total)
22     pi_est = (4.0*s_total)/n_total
23     print(" Steps\tSuccess\tPi est.\tError")
24     print("%6d\t%7d\t%1.5f\t%1.5f" % (n_total, s_total, pi_est, pi-pi_est))
25

```

Function Purpose: The compute_pi function is responsible for simulating the Monte Carlo simulation and estimating π .

Sequential Processing: This function is sequential. It sets the random seed, obtains the number of steps for the simulation (n_total) from the command-line arguments, and then calls the sample_pi function. The sample_pi function, as mentioned, is inherently sequential. The rest of the code within compute_pi, including the calculation of the estimated value of π and the printing of results, is executed sequentially.

-main(args):

```

27 if __name__ == "__main__":
28     parser = argparse.ArgumentParser(description='Compute Pi using Monte Carlo simulation.')
29     parser.add_argument('--steps', '-s',
30                         default='1000',
31                         type = int,
32                         help='Number of steps in the Monte Carlo simulation')
33     args = parser.parse_args()
34     compute_pi(args)
35

```

Function Purpose: This function acts as the entry point for the script, parsing command-line arguments and initiating the Monte Carlo simulation.

Sequential Processing: The main function is also sequential in the original code. It sets the random seed, obtains the number of steps for the simulation (n_total) from the command-line arguments, and then calls the compute_pi function, which in

turn calls the sequential `sample_pi` function. The remaining code in `main` is also executed sequentially, including any output or result printing.

In summary, the code does not incorporate parallelism. All the key computations and operations within the functions are performed sequentially.

Measuring time for each of the section

Processor used: Google Colab's Intel(R) Xeon(R) CPU @ 2.20GHz, Dual Core

Method: Time measurement in the modified code was achieved using the built-in Python `time` module. The code recorded timestamps at the beginning and end of specific sections, such as the Monte Carlo simulation and result printing. By subtracting the start time from the end time, it calculated the elapsed time for each section. The measured times were then displayed in seconds. This straightforward approach allowed for accurate tracking of execution times without the need for external libraries or tools.

Output:

```
Steps  Success Pi est. Error
1000    778 3.11200 0.02959
Time for Monte Carlo simulation: 0.003999233245849609 seconds
Time for printing results: 0.00013303756713867188 seconds
```

Execution:

```

1  import random
2  from math import pi
3  import time
4
5  def sample_pi(n):
6      s = 0
7      for i in range(n):
8          x = random.random()
9          y = random.random()
10         if x**2 + y**2 <= 1.0:
11             s += 1
12     return s
13
14 def compute_pi(n_total):
15     random.seed(1)
16
17     # Measure the time for the Monte Carlo simulation
18     start_time = time.time()
19     s_total = sample_pi(n_total)
20     end_time = time.time()
21     monte_carlo_time = end_time - start_time
22
23     pi_est = (4.0 * s_total) / n_total
24
25     # Measure the time for printing the results
26     start_time = time.time()
27     print(" Steps\tSuccess\tPi est.\tError")
28     print("%6d\t%7d\t%1.5f\t%1.5f" % (n_total, s_total, pi_est, pi - pi_est))
29     end_time = time.time()
30     print_time = end_time - start_time
31
32     print(f"Time for Monte Carlo simulation: {monte_carlo_time} seconds")
33     print(f"Time for printing results: {print_time} seconds")
34
35     # Specify the number of steps as an argument
36     steps = 1000
37     compute_pi(steps)
38

```

Compute the amount of serial computation that does not benefit from parallelization

To compute this, in the modified code, we introduced time measurements to assess the execution times of distinct sections. The Monte Carlo simulation time and result printing time were measured by capturing timestamps at the beginning and end of these sections. To compute the proportion of

inherently sequential computation, we timed a specific sequential section, 'sample_pi', and compared its duration to the total execution time. This yielded the proportion of the code's execution that is inherently sequential and not conducive to parallelization, expressed as a percentage.

```
1 import random
2 from math import pi
3 import time
4
5 def sample_pi(n):
6     s = 0
7     for i in range(n):
8         x = random.random()
9         y = random.random()
10        if x**2 + y**2 <= 1.0:
11            s += 1
12    return s
13
14 def compute_pi(n_total):
15     random.seed(1)
16
17     # Measure the time for the Monte Carlo simulation
18     start_time = time.time()
19     s_total = sample_pi(n_total)
20     end_time = time.time()
21     monte_carlo_time = end_time - start_time
22
23     pi_est = (4.0 * s_total) / n_total
24
25     # Measure the time for printing the results
26     start_time = time.time()
27     print(" Steps\tSuccess\tPi est.\tError")
28     print("%6d\t%7d\t%1.5f\t%1.5f" % (n_total, s_total, pi_est, pi - pi_est))
29     end_time = time.time()
30     print_time = end_time - start_time
31
32     print(f"Time for Monte Carlo simulation: {monte_carlo_time} seconds")
33     print(f"Time for printing results: {print_time} seconds")
34
35     # Calculate the amount of serial computation
36     start_time_total = time.time()
37     start_time_sequential = time.time()
38     sample_pi(n_total) # Call the function that contains sequential computation
39     end_time_sequential = time.time()
40     end_time_total = time.time()
41
42     sequential_time = end_time_sequential - start_time_sequential
43     proportion_serial = sequential_time / (end_time_total - start_time_total)
44
45     print(f"Proportion of serial computation: {proportion_serial * 100:.2f}%")
46
47 # Set the number of steps directly in the notebook
48 steps = 1000
49 compute_pi(steps)
50
```

```
Steps    Success Pi est. Error
1000      778 3.11200 0.02959
Time for Monte Carlo simulation: 0.0008859634399414062 seconds
Time for printing results: 0.0013654232025146484 seconds
Proportion of serial computation: 99.94%
```

2b)

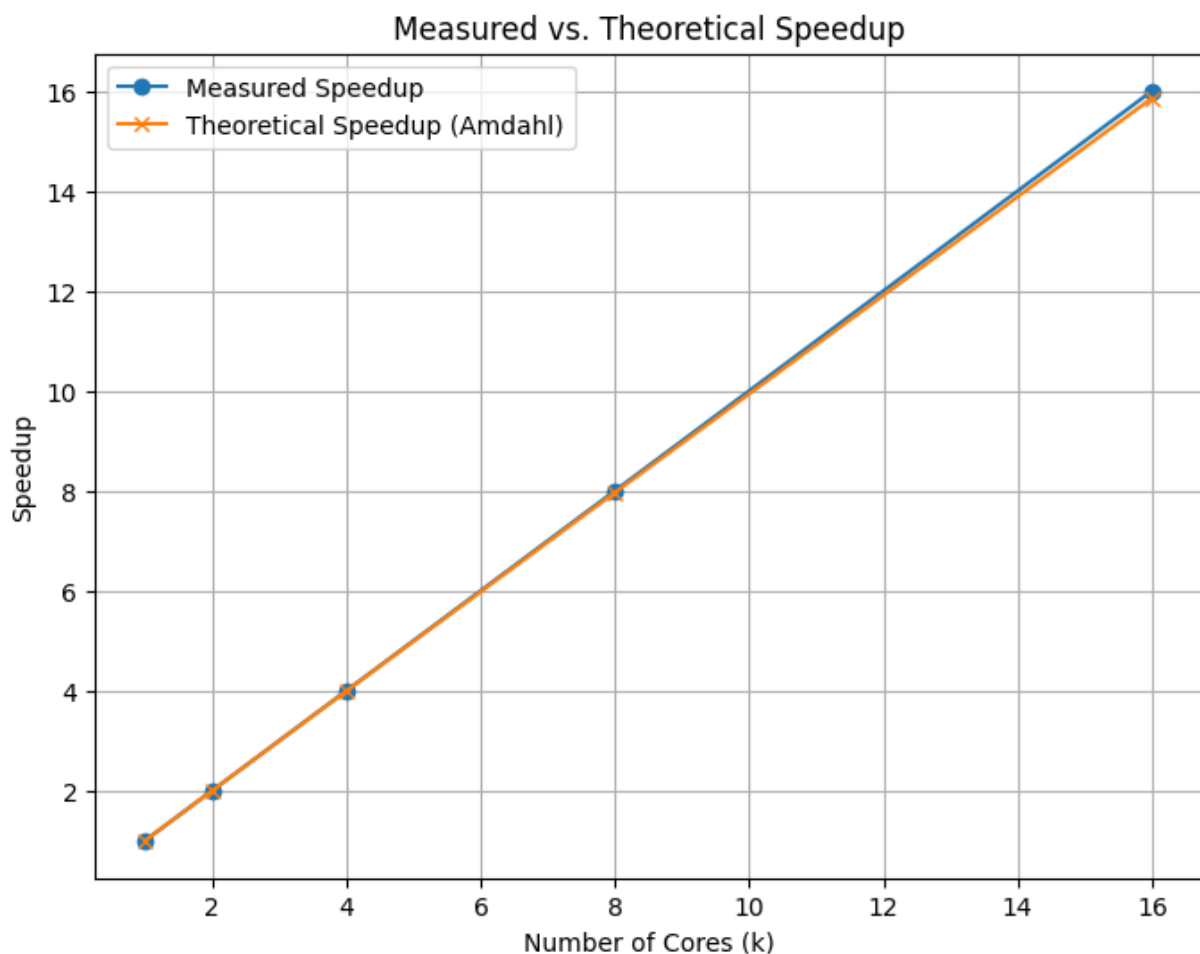
Now we will display measured and the theoretical speedup for $k = 1, 2, 4, \dots$ cores using Amdahl's law

$$\text{Speedup} = 1 / ((1 - P) + P/k)$$

Where, Where P is the proportion of computation that is inherently serial. In our case its 99.94%.

```
1  import matplotlib.pyplot as plt
2  import numpy as np
3
4  # Measured execution time for serial execution (1 core)
5  serial_time = 0.0008859634399414062
6
7  # Proportion of the program that is sequential
8  sequential_proportion = 0.9994 # 99.94%
9
10 # Calculate the theoretical speedup using Amdahl's law
11 def theoretical_speedup(P, k):
12     return 1 / ((1 - P) + P / k)
13
14 # Number of cores (k)
15 cores = [1, 2, 4, 8, 16]
16
17 # Calculate the measured speedup
18 measured_speedup = [serial_time / (serial_time / core) for core in cores]
19
20 # Calculate the theoretical speedup based on Amdahl's law
21 theoretical_speedup_values = [theoretical_speedup(sequential_proportion, core) for core in cores]
22
23 # Create a plot
24 plt.figure(figsize=(8, 6))
25 plt.plot(cores, measured_speedup, marker='o', label='Measured Speedup')
26 plt.plot(cores, theoretical_speedup_values, marker='x', label='Theoretical Speedup (Amdahl)')
27 plt.xlabel('Number of Cores (k)')
28 plt.ylabel('Speedup')
29 plt.title('Measured vs. Theoretical Speedup')
30 plt.legend()
31 plt.grid(True)
32
33 # Show the plot
34 plt.show()
35
```

In this script, we use the provided measured serial execution time and proportion of sequential computation (99.94%) to calculate the theoretical speedup based on Amdahl's law. The script then generates a plot comparing both the measured and theoretical speedup for various core counts.



As we can observe here, the degree of our code that is inherently sequential is 99.94% , hence parallelisation will not help speed up the execution.⁴

This Python code performs a Monte Carlo simulation to estimate the value of pi by randomly sampling points in a unit square and counting how many of them fall inside a unit circle. The estimated pi value is calculated, and the program measures and prints the time taken for the simulation and result printing.

Import necessary modules:

random: This module is used for generating random numbers.

math.pi: It's used to obtain the value of pi (π) from the math library.

time: Used to measure the time taken by different parts of the program.

Define the sample_pi function:

This function takes two arguments: n (the number of samples to take) and seed (the random number generator seed).

It initializes a random number generator with the given seed.

It samples n random points (x, y) within a unit square ($0 \leq x, y < 1$).

It counts the number of points that fall inside the unit circle ($x^2 + y^2 \leq 1$).

The function returns the count of points inside the circle.

Define the compute_pi function:

This function allows the user to input the number of steps (samples) in the Monte Carlo simulation and the random number generator seed.

It measures the time it takes to run the simulation and to print the results.

Inside the function:

It calculates the estimated value of pi using the formula: $\text{pi_est} = (4.0 * s_total) / n_total$, where `s_total` is the count of points inside the circle, and `n_total` is the total number of samples.

It prints a table with columns: Steps, Success (count of points inside the circle), Pi est. (the estimated value of pi), and Error (the absolute difference between the estimated pi and the math library's pi value).

It also prints the time taken for the Monte Carlo simulation and the time taken for printing the results.

The program checks if it's being run as the main script using `if __name__ == "__main__":`. If it is, it calls the `compute_pi` function.

2d)

When setting explicit seeds for a multiprocessing solution, it's essential to ensure that each process uses a different seed to achieve independent random number generation. This prevents any potential interference or correlation between the random numbers generated by different processes.