

Portfolio Milestone 4

Brady Chin

Colorado State University Global

CSC507-2: Foundations of Operating Systems

Dr. Joseph Issa

February 9th, 2025

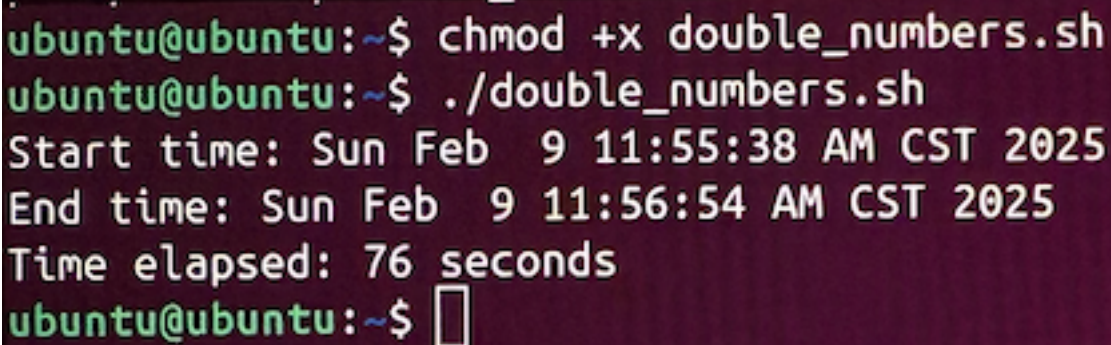
Portfolio Milestone 4

This milestone compares different methods of processing large text files with numerical data using a Bash script and several approaches in Python. By measuring execution times, we can analyze the efficiency of each method and identify how different processing techniques- sequential reading, line-by-line processing, and parallel execution-affect it.

Results

These results show some contrast in run times for methods. Figure 1 shows the execution time of the Bash script. The execution time of a Bash script was 76 seconds.

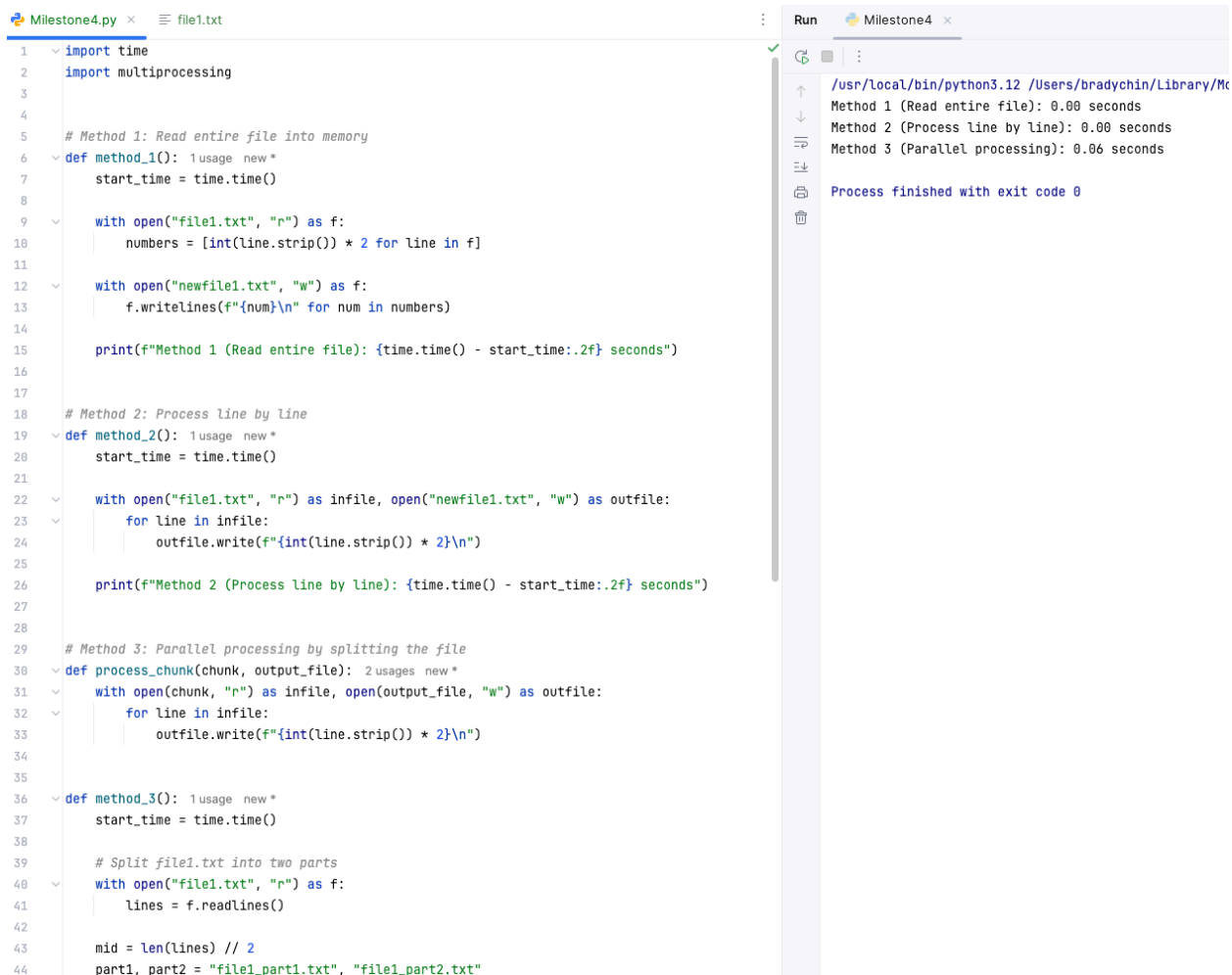
Figure 1: Bash script execution time.

A terminal window with a dark purple background and light green text. The text shows the execution of a Bash script named 'double_numbers.sh'. The first line shows the user making the script executable with 'chmod +x double_numbers.sh'. The second line shows the user running the script with './double_numbers.sh'. The output displays the start time as 'Sun Feb 9 11:55:38 AM CST 2025', the end time as 'Sun Feb 9 11:56:54 AM CST 2025', and the time elapsed as '76 seconds'. The prompt returns to the user's shell.

```
ubuntu@ubuntu:~$ chmod +x double_numbers.sh
ubuntu@ubuntu:~$ ./double_numbers.sh
Start time: Sun Feb 9 11:55:38 AM CST 2025
End time: Sun Feb 9 11:56:54 AM CST 2025
Time elapsed: 76 seconds
ubuntu@ubuntu:~$
```

Figure 2 shows the script and the execution times for all three of the Python methods. Methods 1 and 2, which read the whole file into memory and process line by line, respectively, finished in 0.00 seconds. This is indicative of how efficient Python is at handling such operations, further aided by operating system-level optimizations such as disk caching. Surprisingly, Method 3, parallel processing through chunking, was a bit slower at 0.06 seconds.

Figure 2: Python script and execution times.



The screenshot shows a Python IDE with a script named 'Milestone4.py' and its execution results. The script defines three methods for processing a file: Method 1 (Read entire file), Method 2 (Process line by line), and Method 3 (Parallel processing by splitting the file). The execution results show that Method 1 took 0.00 seconds, Method 2 took 0.00 seconds, and Method 3 took 0.06 seconds. The process finished with exit code 0.

```
1 import time
2 import multiprocessing
3
4
5 # Method 1: Read entire file into memory
6 def method_1(): 1 usage new *
7     start_time = time.time()
8
9     with open("file1.txt", "r") as f:
10         numbers = [int(line.strip()) * 2 for line in f]
11
12     with open("newfile1.txt", "w") as f:
13         f.writelines(f"{num}\n" for num in numbers)
14
15     print(f"Method 1 (Read entire file): {time.time() - start_time:.2f} seconds")
16
17
18 # Method 2: Process Line by Line
19 def method_2(): 1 usage new *
20     start_time = time.time()
21
22     with open("file1.txt", "r") as infile, open("newfile1.txt", "w") as outfile:
23         for line in infile:
24             outfile.write(f"{int(line.strip()) * 2}\n")
25
26     print(f"Method 2 (Process line by line): {time.time() - start_time:.2f} seconds")
27
28
29 # Method 3: Parallel processing by splitting the file
30 def process_chunk(chunk, output_file): 2 usages new *
31     with open(chunk, "r") as infile, open(output_file, "w") as outfile:
32         for line in infile:
33             outfile.write(f"{int(line.strip()) * 2}\n")
34
35
36 def method_3(): 1 usage new *
37     start_time = time.time()
38
39     # Split file1.txt into two parts
40     with open("file1.txt", "r") as f:
41         lines = f.readlines()
42
43     mid = len(lines) // 2
44     part1, part2 = "file1_part1.txt", "file1_part2.txt"
```

Run Milestone4

/usr/local/bin/python3.12 /Users/bradychin/Library/Mc

Method 1 (Read entire file): 0.00 seconds
Method 2 (Process line by line): 0.00 seconds
Method 3 (Parallel processing): 0.06 seconds

Process finished with exit code 0

Surprises

Surprisingly, the above results show that Method 3 was not the fastest of the three Python approaches. Intuitively, it seems that parallelizing the splitting process should pay off for larger files, but in practice the single-threaded approaches were fast enough that their simplicity won out over the more complex multi-processing version. This suggests that Python's built-in file handling is already highly optimized for such a simple task as reading and modifying numbers in a file, and parallel processing at this scale is quite unnecessary.

Another aspect that was surprising was the extreme difference between Bash and Python. The Bash script took more than a minute since it reads line by line, modifies, and

appends in a new file, which implies high I/O overhead (Yesmin, F., 2018). While Python handles the operations of a file much better, especially reading great bulks of data. This experiment really reinforces that for tasks involving heavy file manipulation, Python is the far superior choice due to its speed and built-in optimizations.

References

Yesmin, F., (2018) *How to read file line by line in Bash script*. Linux Hint.

https://linuxhint.com/read_file_line_by_line_bash/