



# Big Data Analytics 2024 Project 1

Oboni Anower

13923163

June 2024

---

## Introduction

This Big Data Analytics project focuses on developing efficient programming strategies for downloading videos from YouTube to use them for further analysis. The goal is to use parallel processing techniques such as Threading, Multiprocessing, and Asynchronous programming to perform simultaneous tasks by taking advantage of modern computational resources. The programming is done in Python, as it offers an extensive libraries and frameworks for parallel processing and multimedia analysis.

This report will discuss about the parallel programming using Threads to download videos, compare it's performance with a serial execution, and discuss space-time complexities for both. It will also discuss about parallel programming libraries used to complete five video analysis tasks. One of such tasks are Audio extraction from the videos, which will be used to compare the execution time taken for different processing methods. Once “*main.py*” is run, all the functionalities will be executed one after another, in order of their call (see Appendix 1).

## Folder Structure

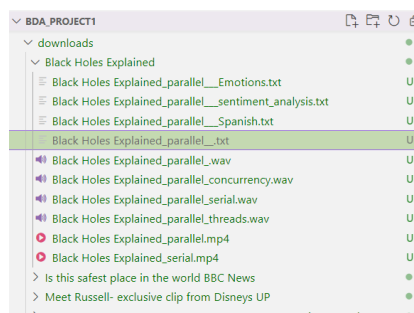


Image 1: dependant files in folder

For every YouTube Video, a folder is generated and named after the video title. The downloaded videos are saved in their respective folders. Once the `main()` function is run successfully, each folder will contain 11 files. Every task, starting from serial test download to subtasks will be downloaded to these folders.

GitHub Repo: [https://github.com/0anower/BDA\\_Project1.git](https://github.com/0anower/BDA_Project1.git)

## Libraries

```
pip install pytube moviepy speechrecognition spacy nltk NRClex Path
TextBlob googletrans==4.0.0-rc1
```

# YouTube Video Downloads

## Video URLs

“*video\_urls.txt*” text file generated from a list containing 12 URLs.

## Serial and Parallel Downloads

While both serial and parallel downloads are run using the same download function `video_download()`, each video in serial programming is downloaded sequentially and only one video is downloaded at a time. Although the program is much simpler, the process tends to be slower, especially if dealing with large volume of data.

Since downloading videos are I/O bound task from a shared resource (list of URLs in same text file), the download speed may often slow down due to the network’s latency or low disc-space. However, the process can be sped up with the help of parallel programming like threads. A thread used over multiprocessing because instead of utilising multiple processors, threads share their memory address for a single process, significantly reducing the overheads.

```
THE NUN II | OFFICIAL TRAILER Download completed
Download started at: 02:07:25 .....Phil Lempert's 2 minute Speech Demo
Phil Lempert's 2 minute Speech Demo Download completed
Download started at: 02:07:28 .....Self Motivation | Brendan Clark | TEDxYouth@BarnstableHS
Self Motivation | Brendan Clark | TEDxYouth@BarnstableHS Download completed
Download started at: 02:07:31 .....Is this safest place in the world? BBC News
Is this safest place in the world? BBC News Download completed
Serial download finished in 34.80964469999992 seconds
```

*Image 2: Serial video download*

```
Download started at: 02:07:47 .....PERCENTAGE INCREASES USING A MULTIPLIER | *2 minute maths*
Download started at: 02:07:49 .....THE NUN II | OFFICIAL TRAILER
Download started at: 02:07:51 .....Phil Lempert's 2 minute Speech Demo
Black Holes Explained Download completed
THE NUN II | OFFICIAL TRAILER Download completed
PERCENTAGE INCREASES USING A MULTIPLIER | *2 minute maths* Download completed
Download started at: 02:07:53 .....Is this safest place in the world? BBC News
Phil Lempert's 2 minute Speech Demo Download completed
Download started at: 02:07:55 .....Self Motivation | Brendan Clark | TEDxYouth@BarnstableHS
The 'Busy' Life of the Sloth | BBC Earth Download completed
Self Motivation | Brendan Clark | TEDxYouth@BarnstableHS Download completed
Is this safest place in the world? BBC News Download completed
Parallel download finished in 24.178622099999984 seconds
```

*Image 3: Parallel Video Download*

Image 1 shows the time taken to execute the same `video_download()` to download each video serially, and concurrently using a thread pool. A thread pool is collection of pre-installed threads that are inactive but on standby. Thus, instead of creating a new thread every time for every download, the thread pool can be used to execute of large number of threads concurrently. It is also faster compared to serial execution. A *semaphore*=5 is used for parallel download, to ensure that no more than 5 threads can access the downloading function at the same time.



```

download_log.txt
1 "Thread_ID": ThreadPoolExecutor-0_0, "TimeStamp": 02:34:18, 24 June 2024, "URL": "https://www.youtube.com/watch?v=HisYsqszq0&ab_channel=Daretodo.Motivation", "Download": True
2 "Thread_ID": ThreadPoolExecutor-0_2, "TimeStamp": 02:34:19, 24 June 2024, "URL": "https://www.youtube.com/watch?v=TK2lvByDN_g", "Download": True
3 "Thread_ID": ThreadPoolExecutor-0_1, "TimeStamp": 02:34:20, 24 June 2024, "URL": "https://www.youtube.com/watch?v=nmPMKcYEkm&ab_channel=WaltDisneyStudios", "Download": True
4 "Thread_ID": ThreadPoolExecutor-0_3, "TimeStamp": 02:34:20, 24 June 2024, "URL": "https://www.youtube.com/watch?v=4x5kjvpQZ-4", "Download": True
5 "Thread_ID": ThreadPoolExecutor-0_4, "TimeStamp": 02:34:20, 24 June 2024, "URL": "https://www.youtube.com/watch?v=t5khm-VjEu4", "Download": True
6 "Thread_ID": ThreadPoolExecutor-0_0, "TimeStamp": 02:34:26, 24 June 2024, "URL": "https://www.youtube.com/watch?v=eyIvcxeMEs8", "Download": True
7 "Thread_ID": ThreadPoolExecutor-0_1, "TimeStamp": 02:34:26, 24 June 2024, "URL": "https://www.youtube.com/watch?v=AMvOp5noJ58", "Download": True
8 "Thread_ID": ThreadPoolExecutor-0_4, "TimeStamp": 02:34:28, 24 June 2024, "URL": "https://www.youtube.com/watch?v=QF-oyCwaArU&ab_channel=WarnerBros.Pictures", "Download": True
9 "Thread_ID": ThreadPoolExecutor-0_4, "TimeStamp": 02:34:28, 24 June 2024, "URL": "https://www.youtube.com/watch?v=C8Im8MC6Q0Q&ab_channel=SupermarketGuru", "Download": True
10 "Thread_ID": ThreadPoolExecutor-0_2, "TimeStamp": 02:34:30, 24 June 2024, "URL": "https://www.youtube.com/watch?v=ndMKTNsR9K&ab_channel=BBCEarth", "Download": True
11 "Thread_ID": ThreadPoolExecutor-0_0, "TimeStamp": 02:34:31, 24 June 2024, "URL": "https://www.youtube.com/watch?v=rLXclBFdVvE&ab_channel=TEDxTalks", "Download": True
12 "Thread_ID": ThreadPoolExecutor-0_1, "TimeStamp": 02:34:32, 24 June 2024, "URL": "https://www.youtube.com/watch?v=0aqiPyTjv8&ab_channel=BBCNews", "Download": True

```

Image 4: download\_log.txt

Although this project requires only 10-15 downloads, a concurrent thread pool is chosen because it reuses a fixed number of threads from the pool, in this case at most 5, which reduces overhead. This is because when a working thread in the pool completes its execution, that same thread can do another available task in the same process. Additionally, when one thread is terminated from the pool, a replacement will complete the following task. Image 3 is “download\_log.txt”, generated after the download completion, which shows that five threads work together concurrently. A *mutex* is used in to generate this log, to prevent multiple threads trying to log at the same time. Instead, it allows only a single thread to write at a time to maintain consistency.

## Time and Space Complexities

### Serial Execution



```

video_download.py
67 # Test SERIAL video download from txt file
68 def serial_download(video_list, download_path):
69     print(f' Serial Download Started...')
70     start=time.perf_counter()
71
72     # extract youtube URLs to download each serially
73     with open(video_list, 'r') as readurlfile:
74         url_list = [url.strip() for url in readurlfile.readlines()]
75         for url in url_list:
76             video_download(url, download_path, download_mode='serial')
77
78     end=time.perf_counter()
79     print(f' Serial download finished in {end-start} seconds\n')

```

Image 5: Serial Programming

**Time Complexity:** YouTube videos are downloaded using `video_download()` which takes constant time  $O(1)$  for each video. The Since there are 12 videos in the List, the time complexity to loop through each video to download all 12 will be  $O(n)$ .

**Space Complexity:** To download each of the 12 videos in the list, so the list will take space  $O(n)$ . Although 12 videos will be downloaded, only 1 video is being downloaded at a time using single thread in a single process. Thus, the space complexity of download remains constant  $O(1)$ . Thus, the total space complexity can be  $O(n) + O(1) = O(n)$

## Parallel Execution

A screenshot of a code editor showing a Python script named 'video\_download.py'. The script defines a function 'parallel\_download' that takes 'video\_list' and 'download\_path' as arguments. It reads URLs from 'video\_urls.txt', sets a semaphore for 5 concurrent downloads, and uses 'ThreadPoolExecutor' to submit download tasks. The code is numbered from 84 to 188. A blue box highlights the lines where tasks are submitted to the executor.

```
84 def parallel_download(video_list, download_path):
85     """performs concurrent execution with threads
86     to download videos from video_urls.txt"""
87
88     # Loads youtube URLs before downloading
89     with open(video_list, 'r') as readurlfile:
90         url_list = [url.strip() for url in readurlfile.readlines()]
91
92     # 5 concurrent download using semaphore
93     max_download=5
94     p_semaphore = threading.BoundedSemaphore(max_download)
95     start=time.perf_counter()
96     print(f' PARALLEL Download Started...')
97
98     # submits video_download func to execute concurrently using thread-pool
99     with ThreadPoolExecutor(max_workers=max_download) as executor:
100         parallel_downloads = [executor.submit(video_download, url, download_path, p_semaphore, download_mode='parallel')
101                               for url in url_list]
102
103     # main thread waits for all submitted downloads to be completed before moving on
104     for download in parallel_downloads:
105         download.result()
106
107     end=time.perf_counter()
108     print(f' Parallel download finished in {end-start} seconds\n')
```

Image 6: Parallel Programming using ThreadPoolExecutor

**Time Complexity:** Thread Pool Executor in the script can handle up to 5 downloads at time. Assuming constant time for each download,  $n=12$  tasks are distributed into groups of 5. Since the download is concurrent, in each group it will take maximum  $t$  time to download 5 videos in parallel, considering the time to be linear. So, the time complexity will be  $O(n/5) = O(n)$

**Space Complexity:** `parallel_downloads` is a list that keeps a reference to the submitted Future objects. Up to 12 Future objects will be created to store in the list because there are 12 videos URLs need to be processed. This makes a space complexity linear, of  $O(12) = O(n)$ .

## Analysis Subtasks

Among the recommended libraries, `multiprocessing` library was implemented as the choice of audio extraction and transcription because these involve decoding videos, extracting audio streams, and processing them, which makes them computationally intensive CPU bound tasks. This is because they distribute the workload across multiple CPU cores.

However, for sentiment analysis, translation, and emotion extraction, the decision was to use `threads`, as they perform fast and use less computational resources compared to multiprocessing. Creating and managing threads also cause less overhead than processes. Additionally, the communication between threads is faster compared to the communication between processes. When executions using Threads and Multiprocessing are compared, it is seen that threads usually tend to perform faster than processes.

For text translation, Google's translation API interface from `googletrans` library is used which can autodetect languages and translate as desired.

Finally, the comparison between different processes such as serial, multiprocessing, threading, and concurrency are experimented for the audio extraction functionality.

## Comparison

```
compare_processes.txt M X
compare_processes.txt
1 Audio Extraction Using Multiprocessing : Extraction completed in 31.457961699999994 seconds
2 Audio Extraction Serially : Extraction completed in 32.127706100000007 seconds
3 Audio Extraction Using Threads : Extraction completed in 33.2074737000000264 seconds
4 Audio Extraction Using Concurrency : Extraction completed in 22.1718985000000225 seconds
```

*Image 7: Time duration for each process of Audio Extraction*

Image 5 shows a text file “*compare\_processes.txt*” which is generated from **task 5** at the end of each audio extraction process, showing the time taken for each of the execution. It shows that concurrency using process pool is the fastest, followed by multiprocessor execution. In this case, threads took the longest because they tend to slow down CPU-bound programs. Concurrent ProcessPoolExecutor takes the shortest time because a pool of processes asynchronously executes functions and call them independently using multiple cores.

# Appendix

## Appendix 1: Folder Structure once the main() function is run successfully.

