

CSE4001 Parallel and Distributed Computing
Slot – L11+L12

DUPLICATE IMAGES DETECTION USING PARALLEL PROCESSING

(A PROJECT REPORT)

Submitted by

SHADAB ALAM (18BCE2491)

NAJMUDIN RAIN (18BCE2505)

Submitting To

Prof. Narayanamoorthy M.

Assistant Prof.(Selection Grade)

Computer Science and Engineering

Fall Semester 2020

CONTENTS

Chapter 1: INTRODUCTION

- 1.1 Objective
- 1.2 Motivation
- 1.3 Background
- 1.4 Computing Image Difference

Chapter 2: LITERATURE REVIEW

Chapter 3: IMPLEMENTATION

- 3.1 Existing and Proposed Solutions
- 3.2 Platforms and Tools used
- 3.3 Program Description
- 3.4 Serial Algorithm
- 3.5 Parallel Algorithm
- 3.6 Parallelisation using multiple processors
- 3.7 Parallelisation using threading

Chapter 4: RESULTS AND DISCUSSION

- 4.1 Analysis: Execution Time vs Data Size
- 4.2 Measurement Methodology
- 4.3 Explanation for Behaviour: MultiProcessing in Py
- 4.4 Explanation for Behaviour: MultiThreading in Py
- 4.5 Explanation for Behaviour: Explained with Data
- 4.6 Analysing MultiProcessing v/s MultiThreading v/s Single Core Execution

Chapter 5: CONCLUSION AND FUTURE WORK

REFERENCES

1. INTRODUCTION

1.1 Objective

The main aim for us to make this project is that we want to make a one stop platform to detect duplicate images from a large sample of given images in lesser time. The reason we are implementing it with multiple cores is that we want to compare it to sequential single core execution and write about our findings.

As our usage of mobile grows, the unnecessary duplicate photo and picture files grows in device randomly, ideally every folder of the phone. The duplicate pictures/photos occupy the lots of phone memory and also reduce the operating speed of phone. Manually it is difficult to find and remove them.

We present the duplicate images detection using parallel processing from which you can scan through your entire phone and finds the duplicate image/picture/photo files for you. You can make judgement yourself if you want to delete by selecting them.

At the end we will also analyse the compute time and power consumed of processing on multiple core v/s single core using threads and provide benchmarks for both along with graphical representations.

1.2 Motivation

No doubt we have numerous images on your computer and phones and even the cloud. The trouble with having lots of pictures is that you tend to collect duplicates along the way. We wanted to build a duplicate photo finder listed to get rid of duplicate and similar photos. It would be a wise thing to manage space efficiently.

Dhruv had a recent trip to Pondi, He asks from his friends to upload all clicked images on google drive. Now here comes the tedious and overwhelming problem, he was a bit late in making this scheme and his friends have already exchanged images amongst them on various platforms like WhatsApp, Telegram, Instagram and all the images are scattered and no one ends up with all the pictured. Now he needed to get these images, and thus many duplicate images to be deleted and here he makes a google search: “How to Remove Duplicate Images”.

Finding similarly pictures and duplicate photos can become an overwhelming project. This is where duplicate photo finders come in. Remember that you should always backup all files before doing any deletion, which is a good practice to take regularly, for example to a USB or an external hard drive.

Now the google search brings in some results, where one result is about Duplicate Cleaner. But Duplicate Cleaner free plan has some limitations, limited speed in free mode, so we thought why not to make our own Duplicate Image Cleaner Pro using Parallelism Fundamentals. That must be a long motivation story if think about it that's the story with almost anyone who has gone on a group trip, that's facts.

In conclusion: We end up with idea for our PDC Project

1.3 Background

Let's pretend that we have a huge dataset of stamp images. And we want to take two arbitrary stamp images and compare them to determine if they are identical, or near identical in some way.

In general, we can accomplish this in two ways.

- The first method is to use locality sensitive hashing, which I'll cover in a later blog post.
- The second method is to use algorithms such as Mean Squared Error (MSE) or the Structural Similarity Index (SSIM).

We can use Python to compare two images using Mean Squared Error and Structural Similarity Index.

Let's take a look at the Mean Squared error equation:

$$MSE = \frac{1}{m \ n} \sum_{i=0}^{m-1} \sum_{j=0}^{n-1} [I(i, j) - K(i, j)]^2$$

In order to remedy some of the issues associated with MSE for image comparison, we have the Structural Similarity Index, developed by Wang et al.:

$$SSIM(x, y) = \frac{(2\mu_x\mu_y + c_1)(2\sigma_{xy} + c_2)}{(\mu_x^2 + \mu_y^2 + c_1)(\sigma_x^2 + \sigma_y^2 + c_2)}$$

Using this method, we were able to easily determine if two images were identical or had differences due to slight image manipulations, compression artifacts, or purposeful tampering.

We are going to extend the SSIM approach so that we can visualize the differences between images using OpenCV and Python. Specifically, we'll be drawing bounding boxes around regions in the two input images that differ.

In order to compute the difference between two images we'll be utilizing the Structural Similarity Index, first introduced by Wang et al. in their 2004 paper, Image Quality Assessment: From Error Visibility to Structural Similarity. This method is already implemented in the scikit-image library for image processing.

The trick is to learn how we can determine exactly where, in terms of (x, y)- coordinate location, the image differences are.

To accomplish this, we'll first need to make sure our system has Python, OpenCV, scikit-image, and imutils.

1.4 Computing Image Difference



Fig 1. Manually inspecting the difference between two input images

If you take a second to study the two credit cards, you'll notice that the MasterCard logo is present on the left image but has been Photoshopped out from the right image.

You may have noticed this difference immediately, or it may have taken you a few seconds. Either way, this demonstrates an important aspect of comparing image

differences — sometimes image differences are subtle — so subtle that the naked eye struggles to immediately comprehend the difference.

Computing image difference is important for various reasons one of which is phishing. Attackers can manipulate images ever-so-slightly to trick unsuspecting users who don't validate the URL into thinking they are logging into their banking website — only to later find out that it was a scam.

2. LITERATURE REVIEW

In [1] “Image quality assessment: from error visibility to structural similarity. IEEE transactions on image processing” the authors Wang, Z., Bovik, A. C., Sheikh, H. R., & Simoncelli, E. P. develop a method for assessing perceptual image quality traditionally attempted to quantify the visibility of errors (differences) between a distorted image and a reference image using a variety of known properties of the human visual system. Under the assumption that human visual perception is highly adapted for extracting structural information from a scene, we introduce an alternative complementary framework for quality assessment based on the degradation of structural information. As a specific example of this concept, we develop a Structural Similarity Index and demonstrate its promise through a set of intuitive examples, as well as comparison to both subjective ratings and state-of-the-art objective methods on a database of images compressed with JPEG and JPEG2000

In [3] “Fast image inpainting using similarity of subspace method” the authors Hosoi, T., Kobayashi, K., Ito, K., & Aoki, T. propose a model for Image inpainting is a technique for estimating missing pixel values in an image by using the pixel value information obtained from neighbor pixels of a missing pixel or the prior knowledge derived from learning the object class. In this paper, we propose a fast and accurate image inpainting method using similarity of the subspace. The proposed method generates the subspace from many images related to the object class in the learning step and estimates the missing pixel values of the input image belonging to the same object class so as to maximize the similarity between the input image and the subspace in the inpainting step. Through a set of experiments, we demonstrate that the proposed method exhibits excellent performance in terms of both inpainting accuracy and computation time compared with conventional algorithms.

In [4] “Image similarity using Fourier transform” the authors Narayanan, S., & Thirivikraman, P. K. use a similarity measure for images based on values from their respective Fourier Transforms is proposed for image registration. The approach uses image content to generate signatures and is not based on image annotation and therefore does not require human assistance. It uses both, the real and complex components of the FFT to compute the final rank for measuring similarity. Any robust approach must accurately represent all objects in an image and depending on the size of the image data set, diverse techniques may need to be followed. This paper discusses implementation of a similarity rating scheme through the Open CV library and introduces a metric for comparison, carried out by considering Intersection bounds of a covariance matrix of two compared images with normalized values of the Magnitude and Phase spectrum. Sample results on a test collection are given along with data using existing methods of image histogram comparison. Results have shown that this method is particularly advantageous in images with varying degrees of lighting.

In [6] “Image Edge Detection Based On Opencv” the authors Xie, G., & Lu, W. introduce a method of image edge detection to determine the exact number of the copper core in the tiny wire based on OpenCV with rich computer vision and image processing algorithms and functions. Firstly, we use high-resolution camera to take picture of the internal structure of the wire. Secondly, we use OpenCV image processing functions to implement image preprocessing. Thirdly we use morphological opening and closing operations to segment image because of their blur image edges. Finally the exact number of copper core can be clearly distinguished through contour tracking. By using of Borland C++ Builder 6.0, experimental results show that OpenCV based image edge detection methods are simple, high code integration, and high image edge positioning accuracy.

In [7] “Parallel astronomical data processing with Python: Recipes for multicore machines” the authors Singh, N., Browne, L. M., & Butler, R. propose parallel processing recipes for multicore machines for astronomical data processing. The target audience is astronomers who use Python as their preferred scripting language and who may be using PyRAF/IRAF for data processing. Three problems of varied complexity were benchmarked on three different types of multicore processors to demonstrate the benefits, in terms of execution time, of parallelizing data processing tasks. The native multiprocessing module available in Python makes it a relatively trivial task to implement the parallel code. We have also compared the three multiprocessing approaches—Pool/Map, Process/Queue and Parallel Python. Our test codes are freely available and can be downloaded from our website.

In [9] “Composable multi-threading for Python libraries.” the author Malakhov, A. tells that python is popular among numeric communities that value it for easy to use number crunching modules like [NumPy], [SciPy], [Dask], [Numba], and many others. These modules often use multi-threading for efficient multi-core parallelism in order to utilize all the available CPU cores. Nevertheless, their threads can interfere with each other leading to overhead and inefficiency if used together in one application. The loss of performance can be prevented if all the multi-threaded parties are coordinated. This paper describes usage of Intel® Threading Building Blocks (Intel® TBB), an open-source cross-platform library for multi-core parallelism [TBB], as the composability layer for Python modules. It helps to unlock additional performance for numeric applications on multi-core systems.

In [10] “PX4: A node-based multithreaded open source robotics framework for deeply embedded platforms” the authors Meier, L., Honegger, D., & Pollefeys, M. present a novel, deeply embedded robotics middleware and programming environment. It uses a multithreaded, publish-subscribe design pattern and provides a Unix-like software interface for micro controller applications. We improve over the state of the art in deeply embedded open source systems by providing a modular and standards-oriented platform. Our system architecture is centered around a publish-subscribe object request broker on top of a POSIX application programming interface. This allows to reuse common Unix knowledge and experience, including a bash-like shell. We demonstrate with a vertical takeoff and landing (VTOL) use case that the system modularity is well suited for novel and experimental vehicle platforms. We also show how the system architecture allows a direct interface to ROS and to run individual processes either as native ROS nodes on Linux or nodes on the micro controller, maximizing interoperability. Our microcontroller-based execution environment has substantially lower latency and better hardware connectivity than a typical Robotics Linux system and is therefore well suited for fast, high rate control tasks.

3. IMPLEMENTATION

3.1 Existing and Proposed Solutions

The trouble with having lots of pictures is that you tend to collect duplicates along the way. There are a lot of duplicate image finder listed to get rid of duplicate and similar photos. It would be a wise thing to manage space efficiently.

Existing Solution:

1. Duplicate Photo Finder
2. Anti-Twin
3. VisiPics
4. Similar Image Search
5. Awesome Duplicate Photo Finder

Proposed Solution:

1. An open source Py package to remove duplicate Images

3.2 Platforms and Tools used

3.2.1 PyCharm

PyCharm is an integrated development environment (IDE) used in computer programming, specifically for the Python language. It is developed by the Czech company JetBrains. It provides code analysis, a graphical debugger, an integrated unit tester, integration with version control systems (VCSes), and supports web development with Django as well as Data Science with Anaconda.

PyCharm is cross-platform, with Windows, macOS and Linux versions. The Community Edition is released under the Apache License, and there is also Professional Edition with extra features – released under a proprietary license.

3.2.2 Visual Studio Code

Visual Studio Code is a free source-code editor made by Microsoft for Windows, Linux and macOS. Features include support for debugging, syntax highlighting, intelligent code completion, snippets, code refactoring, and embedded Git. Users can change the theme, keyboard shortcuts, preferences, and install extensions that add additional functionality. The source code is free and open-source, released under

the permissive MIT License. The compiled binaries are freeware for any use.

3.2.3 Github

GitHub, Inc. is a United States-based global company that provides hosting for software development version control using Git. In 2018, it became a subsidiary of Microsoft for US\$7.5 billion. It offers the distributed version control and source code management (SCM) functionality of Git, plus its own features. It provides access control and several collaboration features such as bug tracking, feature requests, task management, and wikis for every project.

3.2.4 Python3

3.2.5 Python Modules and Libraries

3.2.5.1 OpenCV

OpenCV (Open Source Computer Vision Library) is a library of programming functions mainly aimed at real-time computer vision. Originally developed by Intel, it was later supported by Willow Garage then Itseez (which was later acquired by Intel). The library is cross-platform and free for use under the open-source BSD license.

3.2.5.2 multiprocessing

multiprocessing is a package that supports spawning processes using an API similar to the threading module. The multiprocessing package offers both local and remote concurrency, effectively side-stepping the Global Interpreter Lock by using subprocesses instead of threads. Due to this, the multiprocessing module allows the programmer to fully leverage multiple processors on a given machine. It runs on both Unix and Windows.

The multiprocessing module also introduces APIs which do not have analogs in the threading module. A prime example of this is the Pool object which offers a convenient means of parallelizing the execution of a function across multiple input values, distributing the input data across processes (data parallelism). The following example demonstrates the common practice of defining such functions in a module so that child processes can successfully import that module.

3.2.5.3 threading

This module constructs higher-level threading interfaces on top of the lower level `_thread` module. This module provides low-level primitives for working with multiple threads (also called light-weight processes or tasks) — multiple threads of control sharing their global data space. For synchronization, simple locks (also called mutexes or binary semaphores) are provided. The threading module provides an easier to use and higher-level threading API built on top of this module.

3.3 Program Description

We have created a working Script to detect duplicate images which identifies duplicate images recursively using multiprocessing library of python

In order to compute the difference between two images we'll be utilizing the Structural Similarity Index, first introduced by Wang et al. in their 2004 paper, Image Quality Assessment: From Error Visibility to Structural Similarity. This method is already implemented in the scikit-image library for image processing.

The trick is to learn how we can determine exactly where, in terms of (x, y)- coordinate location, the image differences are.

To accomplish this, we'll first need to make sure our system has Python, OpenCV, scikit-image, and imutils.

3.4 Serial Algorithm

Search for Images Recursively from a given starting directory

Compare two images by their RGB values and find image similarity using cv2 library sequentially

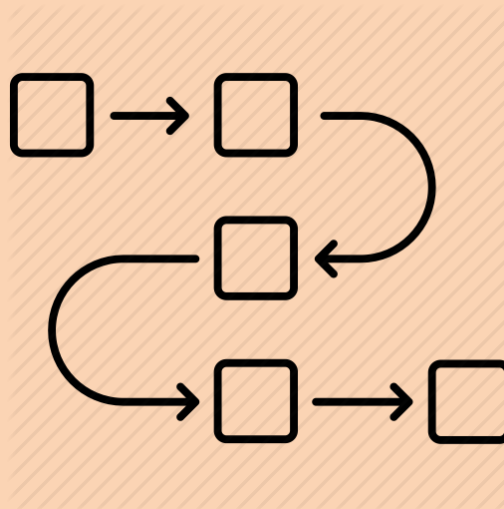


Fig 2.1 Sequential Algorithm

In computer science, a sequential algorithm or serial algorithm is an algorithm that is executed sequentially – once through, from start to finish, without other processing executing – as opposed to concurrently or in parallel. The term is primarily used to contrast with concurrent algorithm or parallel algorithm; most standard computer algorithms are sequential algorithms, and not specifically identified as such, as sequentialness is a background assumption. Concurrency and parallelism are in general distinct concepts, but they often overlap – many distributed algorithms are both concurrent and parallel – and thus "sequential" is used to contrast with both, without distinguishing which one. If these need to be distinguished, the opposing pairs sequential/concurrent and serial/parallel may be used.

"Sequential algorithm" may also refer specifically to an algorithm for decoding a convolutional code.

3.5 Parallel Algorithm

Here we use the divide and conquer method.

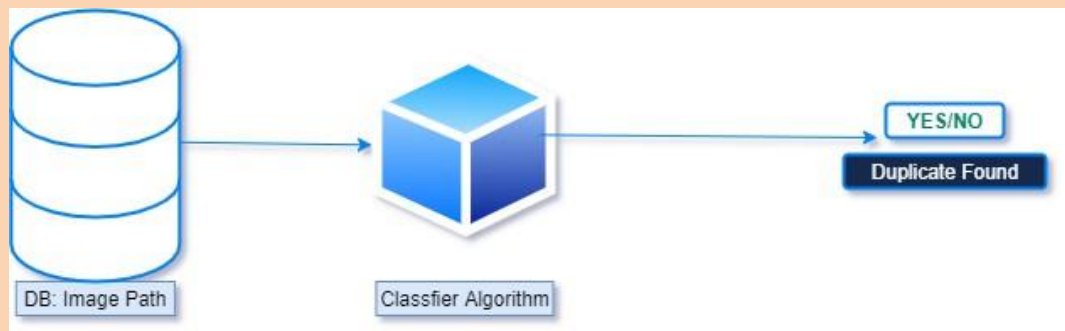


Fig 3.1 Sequential Execution

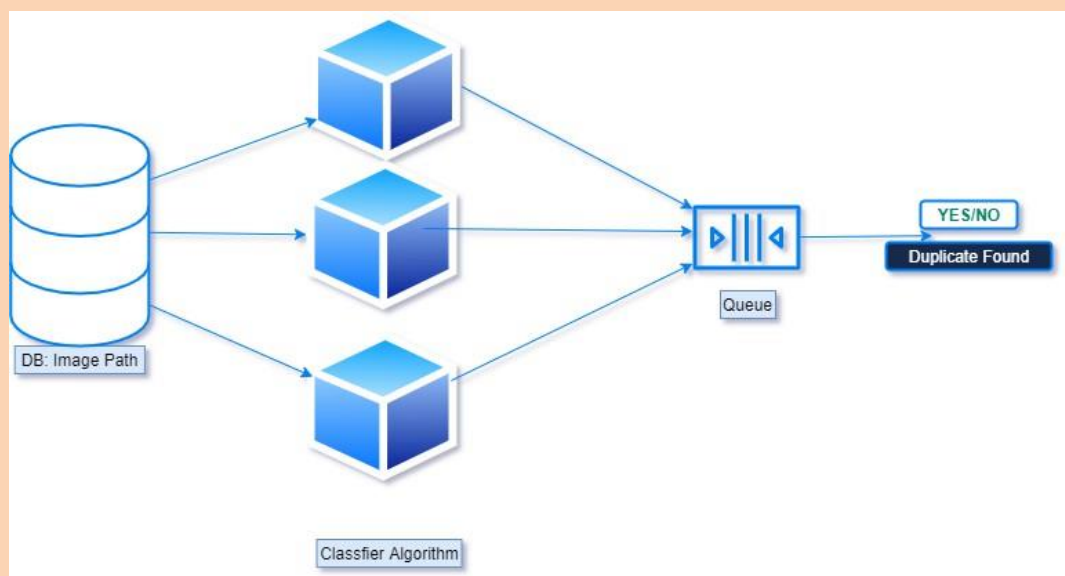


Fig 3.2 Parallel Execution

Divide and Conquer

This technique can be divided into the following three parts:

1. **Divide:** This involves dividing the problem into some sub problem.
2. **Conquer:** Sub problem by calling recursively until sub problem solved.
3. **Combine:** The Sub problem Solved so that we will get find problem solution.

□ Divide And Conquer algorithm:

```

DAC (a, i, j)
{
    if (small(a, i, j))
        return (Solution(a, i, j))
    else
        m = divide(a, i, j)           // f1(n)
        b = DAC(a, i, mid)           // T(n/2)
        c = DAC(a, mid+1, j)         // T(n/2)
        d = combine(b, c)             // f2(n)
    return (d)
}

```

□ Recurrence Relation for DAC algorithm: This is recurrence relation for above program.

$O(1)$ if n is small
 $T(n) = f_1(n) + 2T(n/2) + f_2(n)$

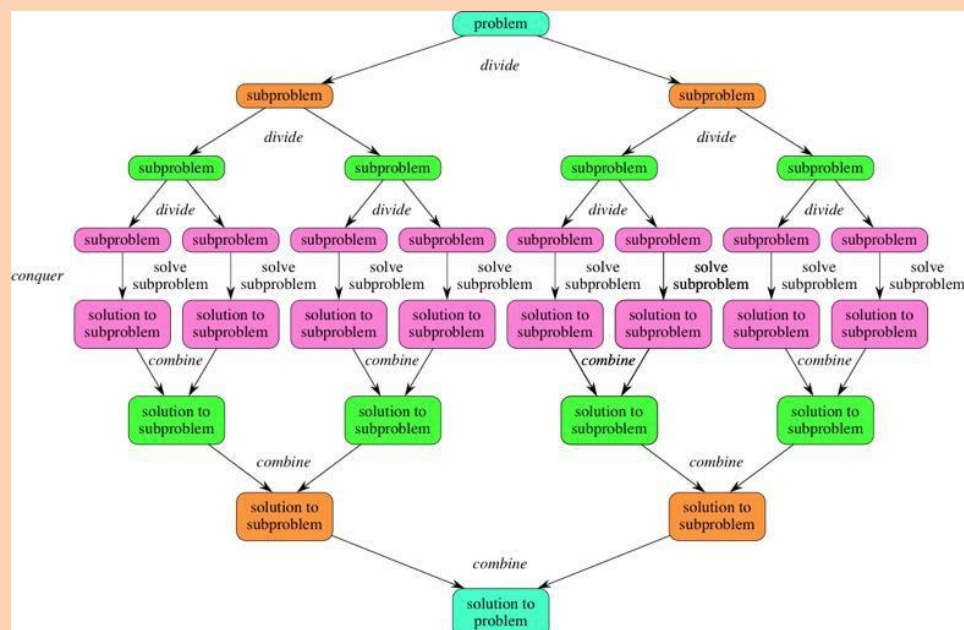


Fig 4 Divide and Conquer Algorithm

3.6 Parallelisation using multiple processors

Search for Images Recursively from a given starting directory

Compare two images by their RGB values and find image similarity using cv2 library by spawning processes using multiprocessing library API

Contexts and start methods

Depending on the platform, multiprocessing supports three ways to start a process. These start methods are

- **spawn**

The parent process starts a fresh python interpreter process. The child process will only inherit those resources necessary to run the process objects `run()` method. In particular, unnecessary file descriptors and handles from the parent process will not be inherited. Starting a process using this method is rather slow compared to using `fork` or `forkserver`.

- **fork**

The parent process uses `os.fork()` to fork the Python interpreter. The child process, when it begins, is effectively identical to the parent process. All resources of the parent are inherited by the child process. Note that safely forking a multithreaded process is problematic.

- **forkserver**

When the program starts and selects the `forkserver` start method, a server process is started. From then on, whenever a new process is needed, the parent process connects to the server and requests that it fork a new process. The fork server process is single threaded so it is safe for it to use `os.fork()`. No unnecessary resources are inherited.

3.7 Parallelisation using threading

Search for Images Recursively from a given starting directory

Compare two images by their RGB values and find image similarity using cv2 library by passing each comparison to a new thread

Python on the CPython interpreter does not support true multi-core execution via multithreading. However, Python DOES have a Threading library. So what is the benefit of using the library if we (supposedly) cannot make use of multiple cores?

Many programs, particularly those relating to network programming or data input/output (I/O) are often network-bound or I/O bound. This means that the Python interpreter is awaiting the result of a function call that is manipulating data from a "remote" source such as a network address or hard disk. Such access is far slower than reading from local memory or a CPU-cache.

Hence, one means of speeding up such code if many data sources are being accessed is to generate a thread for each data item needing to be accessed.

For example, consider a Python code that is scraping many web URLs. Given that each URL will have an associated download time well in excess of the CPU processing capability of the computer, a single-threaded implementation will be significantly I/O bound.

By adding a new thread for each download resource, the code can download multiple data sources in parallel and combine the results at the end of every download. This means that each subsequent download is not waiting on the download of earlier web pages. In this case the program is now bound by the bandwidth limitations of the client/server(s) instead.

However, many financial applications ARE CPU-bound since they are highly numerically intensive. They often involve large-scale numerical linear algebra solutions or random statistical draws, such as in Monte Carlo simulations. Thus as far as Python and the GIL are concerned, there is no benefit to using the Python Threading library for such tasks.

Python offers four possible ways to handle that. First, you can execute functions in parallel using the multiprocessing module. Second, an alternative to processes are threads. Technically, these are lightweight processes, and are outside the scope of this article. For further reading you may have a look at the Python threading module. Third, you can call external programs using the `system()` method of the `os` module, or methods provided by the `subprocess` module, and collect the results afterwards.

The multiprocessing module covers a nice selection of methods to handle the parallel execution of routines. This includes processes, pools of agents, queues, and pipes.

3.8 Proposed Algorithm

□ To find duplicate images

```
def check_duplicate(self, image_):
    """Prints Duplicate Image Name"""
    try:
        image_to_compare = cv2.imread(image_)
        if original.shape == image_to_compare.shape:

            difference =
cv2.subtract(original, image_to_compare)
            b, g, r = cv2.split(difference)

            if cv2.countNonZero(b) == 0 and cv2.countNonZero(g)
== 0 and cv2.countNonZero(r) == 0:
                print('Duplicate Found')
```

□ To generate data for thread processing

```
class DuplicateFinder(threading.Thread):

    def _init_(self, queue, queue2):
        """Initialize the thread"""
        threading.Thread._init_(self)
        self.lock = threading.Lock()
        self.queue = queue
        self.queue2 = queue2

    def run(self):
        # global image_count
        """Run the Thread"""

        while True:
            image_path = self.queue.get()

            self.check_duplicate(image_path)

            # send a signal to the queue that the job is done
            self.queue.task_done()

    def check_duplicate(self, image_):
        global image_count
        """Prints Duplicate Image Name"""
        try:
            image_to_compare = cv2.imread(image_)
            if original.shape == image_to_compare.shape:

                difference = cv2.subtract(original, image_to_compare)
                b, g, r = cv2.split(difference)

                if cv2.countNonZero(b) == 0 and cv2.countNonZero(g) == 0 and cv2.countNonZero(r) == 0:
                    # print('Duplicate Found')
                    sift = cv2.xfeatures2d.SIFT_create()
                    kp_1, desc_1 = sift.detectAndCompute(original, None)
                    kp_2, desc_2 = sift.detectAndCompute(image_to_compare, None)

                    index_params = dict(algorithm=0, trees=5)
                    search_params = dict()
```

```

        flann = cv2.FlannBasedMatcher(index_params,
search_params)

        matches = flann.knnMatch(desc_1, desc_2, k=2)

        good_points = []
        for m, n in matches:
            if m.distance < 0.6*n.distance:
                good_points.append(m)

        # Define how similar they are
        number_keypoints = 0
        if len(kp_1) <= len(kp_2):
            number_keypoints = len(kp_1)
        else:
            number_keypoints = len(kp_2)

        self.queue2.put([image_count, round(time.time() -
start_time, 5)]) # we can skip image count but that's not fair

        self.lock.acquire()
        try:
            image_count += 1
        finally:
            self.lock.release()
    except Exception as e:
        pass

def main(image_dir):
    """
    Run the program
    """

    queue = Queue()
    queue2 = Queue() # queue to take plotting data

    # create a thread pool and give them a queue
    for i in range(5000): # I don't want my virtual memory to get
bunked up
        t = DuplicateFinder(queue, queue2)
        t.setDaemon(True)
        t.start()

    #give the queue some data | consumer and producer code huh
    for image_ in image_dir:
        queue.put(image_)

    # wait for the queue to finish
    queue.join()
    queue2.join()

    print(list(queue2.queue))

```

Code- Detect duplicate multi core

```
import cv2
import numpy as np
import glob
import time
import multiprocessing
import threading

"""
Multiprocessing
"""

image_name = 'cat_d1.jpeg' # original image path goes here
original = cv2.imread(image_name)

def find_duplicates(image_):
    duplicates = ''
    try:
        image_to_compare = cv2.imread(image_)
        if original.shape == image_to_compare.shape:

            difference = cv2.subtract(original, image_to_compare)
            b, g, r = cv2.split(difference)

            if cv2.countNonZero(b) == 0 and cv2.countNonZero(g) == 0 and
cv2.countNonZero(r) == 0:
                duplicates = image_

            sift = cv2.xfeatures2d.SIFT_create()
            kp_1, desc_1 = sift.detectAndCompute(original, None)
            kp_2, desc_2 = sift.detectAndCompute(image_to_compare, No
ne)

            index_params = dict(algorithm=0, trees=5)
            search_params = dict()
            flann = cv2.FlannBasedMatcher(index_params, search_params
)

            matches = flann.knnMatch(desc_1, desc_2, k=2)

            good_points = []
            for m, n in matches:
                if m.distance < 0.6*n.distance:
```

```

        good_points.append(m)

        # Define how similar they are
        number_keypoints = 0
        if len(kp_1) <= len(kp_2):
            number_keypoints = len(kp_1)
        else:
            number_keypoints = len(kp_2)

        return duplicates

    except Exception as e:
        pass

if __name__ == '__main__':

    def filebrowser(ext='', directory=''):
        """
        returns files with an extension
        """
        return [f for f in glob.glob(f"{directory}**/*{ext}", recursive=True)

]

    image_dir = filebrowser(ext='.jpeg', directory='') # directory='' ==> sea
rch images from current to inner sub-directories
    image_dir += filebrowser(ext='.jpg', directory='')
    ## print(image_dir)

    # # 1) Check if 2 images are equals | parallel on cores
    start_time = time.time()
    pool = multiprocessing.Pool() # Equal to number of cores | octa for this
pc

    inputs = image_dir
    outputs_async = pool.map_async(find_duplicates, inputs)
    print(outputs_async.get())

    print("--- %s seconds ---" % (time.time() - start_time))
    print('Program Executed Completely')

```

Output-

```
[Running] python -
u "d:\VScode Python\Duplicate_detection_script\tempCodeRunnerFile.py"
[None, None, None, None, None, None, None, None, None, None, None, None, None,
, None, None, None, None, None, None, None, None, None, None, None, None, None, Non
e, None, None, None, None, None, None, None, None, None, None, None, None, No
ne, None, None, None, None, None, None, None, None, None, None, None, None, N
one, None, None, None, None, None, None, None, None, None, None, None, None, None,
None, None, None, None, None, None, None]
--- 1.2918365001678467 seconds ---
Program Executed Completely
```

2-Detect Duplicate Single Core-

```
import cv2
import numpy as np
import glob
import time
import multiprocessing
import threading

"""
Sequentiall on single core
"""

def filebrowser(ext='', directory=''):
    """
    returns files with an extension
    """
    return [f for f in glob.glob(f"{directory}**/*{ext}", recursive=True)]

image_dir = filebrowser(ext='.jpeg', directory='') # directory='' ==> search
images from current to inner sub-directories
image_dir += filebrowser(ext='.jpg', directory='')

## print(image_dir)

image_name = 'cat_d1.jpeg' # original image path goes here
original = cv2.imread(image_name)

start_time = time.time()

for image_ in image_dir:
    try:
```

```

image_to_compare = cv2.imread(image_)
if original.shape == image_to_compare.shape:

    difference = cv2.subtract(original, image_to_compare)
    b, g, r = cv2.split(difference)

    if cv2.countNonZero(b) == 0 and cv2.countNonZero(g) == 0 and cv2.
countNonZero(r) == 0:
        print(f'Duplicates Found: {image_name} is Duplicate of {image
_}')

    sift = cv2.xfeatures2d.SIFT_create()
    kp_1, desc_1 = sift.detectAndCompute(original, None)
    kp_2, desc_2 = sift.detectAndCompute(image_to_compare, None)

    index_params = dict(algorithm=0, trees=5)
    search_params = dict()
    flann = cv2.FlannBasedMatcher(index_params, search_params)

    matches = flann.knnMatch(desc_1, desc_2, k=2)

    good_points = []
    for m, n in matches:
        if m.distance < 0.6*n.distance:
            good_points.append(m)

    # Define how similar they are
    number_keypoints = 0
    if len(kp_1) <= len(kp_2):
        number_keypoints = len(kp_1)
    else:
        number_keypoints = len(kp_2)
except Exception as e:
    pass

print("--- %s seconds ---" % (time.time() - start_time))
print('Program Executed Completely')

```

Output-

```
[Running] python -
u "d:\VScode Python\Duplicate_detection_script\tempCodeRunnerFile.py"
Duplicates Found: cat_d1.jpeg is Duplicate of cat_d1.jpeg
Duplicates Found: cat_d1.jpeg is Duplicate of data_generating_script\cat_d1.j
peg
Duplicates Found: cat_d1.jpeg is Duplicate of data_generating_script\images_i
n\cat_d1 - Copy.jpeg
Duplicates Found: cat_d1.jpeg is Duplicate of data_generating_script\images_i
n\cat_d1.jpeg
Duplicates Found: cat_d1.jpeg is Duplicate of Duplicate_detection_script\imag
es_in\cat_d1 - Copy.jpeg
Duplicates Found: cat_d1.jpeg is Duplicate of Duplicate_detection_script\imag
es_in\cat_d1.jpeg
Duplicates Found: cat_d1.jpeg is Duplicate of images_in\cat_d1.jpeg
--- 2.298617124557495 seconds ---
Program Executed Completely
```

3-Dectect Duplicate.py

```
import cv2
import numpy as np
import glob
import time
import multiprocessing
import threading

"""
Sequentiall on single core
"""

# def filebrowser(ext='', directory=''):
#     """
#     returns files with an extension
#     """
#     return [f for f in glob.glob(f"{directory}**/*{ext}", recursive=True)]

# image_dir = filebrowser(ext='.jpeg', directory='') # directory='' ==> searc
h images from current to inner sub-directories
# image_dir += filebrowser(ext='.jpg', directory='')

# ## print(image_dir)

# image_name = 'cat_d1.jpeg' # origional image path goes here
# original = cv2.imread(image_name)
```



```

# start_time = time.time()

# for image_ in image_dir:
#     try:
#         image_to_compare = cv2.imread(image_)
#         if original.shape == image_to_compare.shape:

#             difference = cv2.subtract(original, image_to_compare)
#             b, g, r = cv2.split(difference)

#             if cv2.countNonZero(b) == 0 and cv2.countNonZero(g) == 0 and cv
2.countNonZero(r) == 0:
#                 print(f'Duplicates Found: {image_name} is Duplicate of {ima
ge_}')

#                 sift = cv2.xfeatures2d.SIFT_create()
#                 kp_1, desc_1 = sift.detectAndCompute(original, None)
#                 kp_2, desc_2 = sift.detectAndCompute(image_to_compare, None
)

#                 index_params = dict(algorithm=0, trees=5)
#                 search_params = dict()
#                 flann = cv2.FlannBasedMatcher(index_params, search_params)

#                 matches = flann.knnMatch(desc_1, desc_2, k=2)

#                 good_points = []
#                 for m, n in matches:
#                     if m.distance < 0.6*n.distance:
#                         good_points.append(m)

#                 # Define how similar they are
#                 number_keypoints = 0
#                 if len(kp_1) <= len(kp_2):
#                     number_keypoints = len(kp_1)
#                 else:
#                     number_keypoints = len(kp_2)
#             except Exception as e:
#                 pass

# print("--- %s seconds ---" % (time.time() - start_time))
# print('Program Executed Completely')

```

```

"""
Multiprocessing
"""

# image_name = 'c2.jpeg' # original image path goes here
# original = cv2.imread(image_name)

# def find_duplicates(image_):
#     duplicates = ''
#     try:
#         image_to_compare = cv2.imread(image_)
#         if original.shape == image_to_compare.shape:

#             difference = cv2.subtract(original, image_to_compare)
#             b, g, r = cv2.split(difference)

#             if cv2.countNonZero(b) == 0 and cv2.countNonZero(g) == 0 and
# cv2.countNonZero(r) == 0:
#                 duplicates = image_

#                 sift = cv2.xfeatures2d.SIFT_create()
#                 kp_1, desc_1 = sift.detectAndCompute(original, None)
#                 kp_2, desc_2 = sift.detectAndCompute(image_to_compare,
None)

#                 index_params = dict(algorithm=0, trees=5)
#                 search_params = dict()
#                 flann = cv2.FlannBasedMatcher(index_params, search_params)

#                 matches = flann.knnMatch(desc_1, desc_2, k=2)

#                 good_points = []
#                 for m, n in matches:
#                     if m.distance < 0.6*n.distance:
#                         good_points.append(m)

#                 # Define how similar they are
#                 number_keypoints = 0
#                 if len(kp_1) <= len(kp_2):
#                     number_keypoints = len(kp_1)
#                 else:

```

```

#             number_keypoints = len(kp_2)

#             return duplicates

#         except Exception as e:
#             pass

# if __name__ == '__main__':

#     def filebrowser(ext='', directory=''):
#         """
#         returns files with an extension
#         """
#         return [f for f in glob.glob(f"{directory}**/*{ext}", recursive=True)]

#     image_dir = filebrowser(ext='.jpeg', directory='') # directory='' ==> search images from current to inner sub-directories
#     image_dir += filebrowser(ext='.jpg', directory='')
#     ## print(image_dir)

#     # # 1) Check if 2 images are equals | parallel on cores
#     start_time = time.time()
#     pool = multiprocessing.Pool() # Equal to number of cores | octa for this pc

#     inputs = image_dir
#     outputs_async = pool.map_async(find_duplicates, inputs)
#     print(outputs_async.get())

#     print("--- %s seconds ---" % (time.time() - start_time))
#     print('Program Executed Completely')

"""
Threading
"""

image_name = 'cat_d1.jpeg' # original image path goes here
original = cv2.imread(image_name)

```

```

def find_duplicates(image_=''):
    duplicates = ''
    try:
        image_to_compare = cv2.imread(image_)
        if original.shape == image_to_compare.shape:

            difference = cv2.subtract(original, image_to_compare)
            b, g, r = cv2.split(difference)

            if cv2.countNonZero(b) == 0 and cv2.countNonZero(g) == 0 and
cv2.countNonZero(r) == 0:
                print(f'Duplicates Found: {image_name} is Duplicate of {i
mage_}')

                sift = cv2.xfeatures2d.SIFT_create()
                kp_1, desc_1 = sift.detectAndCompute(original, None)
                kp_2, desc_2 = sift.detectAndCompute(image_to_compare, No
ne)

                index_params = dict(algorithm=0, trees=5)
                search_params = dict()
                flann = cv2.FlannBasedMatcher(index_params, search_params
)

                matches = flann.knnMatch(desc_1, desc_2, k=2)

                good_points = []
                for m, n in matches:
                    if m.distance < 0.6*n.distance:
                        good_points.append(m)

                # Define how similar they are
                number_keypoints = 0
                if len(kp_1) <= len(kp_2):
                    number_keypoints = len(kp_1)
                else:
                    number_keypoints = len(kp_2)

    except Exception as e:
        pass

if __name__ == '__main__':

```

```

def filebrowser(ext='', directory=''):
    """
    returns files with an extension
    """
    return [f for f in glob.glob(f"{directory}**/*{ext}", recursive=True)
]

image_dir = filebrowser(ext='.jpeg', directory='') # directory='' ==> sea
rch images from current to inner sub-directories
image_dir += filebrowser(ext='.jpg', directory='')

start_time = time.time()

for image_ in image_dir:
    t = threading.Thread(target=find_duplicates, args=(image_,))
    t.daemon = False
    t.start()

print("--- %s seconds ---" % (time.time() - start_time))
print('Program Executed Completely')

```

OUTPUT-

```

[Running] python -
u "d:\VScode Python\Duplicate_detection_script\detect_duplicate.py"
--- 0.03650474548339844 seconds ---
Program Executed Completely
Duplicates Found: cat_d1.jpeg is Duplicate of data_generating_script\images_i
n\cat_d1 - Copy.jpeg
Duplicates Found: cat_d1.jpeg is Duplicate of data_generating_script\images_i
n\cat_d1.jpeg
Duplicates Found: cat_d1.jpeg is Duplicate of cat_d1.jpeg
Duplicates Found: cat_d1.jpeg is Duplicate of Duplicate_detection_script\imag
es_in\cat_d1 - Copy.jpeg
Duplicates Found: cat_d1.jpeg is Duplicate of Duplicate_detection_script\imag
es_in\cat_d1.jpeg
Duplicates Found: cat_d1.jpeg is Duplicate of data_generating_script\cat_d1.j
peg
Duplicates Found: cat_d1.jpeg is Duplicate of images_in\cat_d1.jpeg

```

4. RESULTS AND DISCUSSION

4.1 Analysis: Execution Time vs Data Size

Processing Type -->	Single Core & Single Threaded	Octa Core & Single Threaded	Single Core & Multiple Threaded
6 Images or 4 MB	0.2550806999206543 seconds	0.5232882499694824 seconds	0.0 seconds
184 images or 64 MB	4.634257793426514 seconds	2.0565671920776367 seconds	1.0612797737121582 seconds
1472 images or 520 MB	38.83846831321716 seconds	11.805004835128784 seconds	10.827009201049805 seconds

Fig 5. Analysis: Execution Time vs Data Size

Note: Threads are not synchronized, after sync time taken for 520MB is 13.5 sec using 100 Threads

4.2 Measurement Methodology

OS: Windows 10

Value Calculation: Mean of 3 Executions

Data Used: Consistent with all Scripts

Data Selection: No Preference

Error Prone areas and why:

- ☐ Image Processing TimeSheet: It's possible that processor or thread worked on bigger images first

Note: We have taken extreme measures to ensure that workload on all 3 executions is similar

4.3 Explanation for Behaviour: MultiProcessing in Py

- ☐ Multiprocessing adds CPUs to increase computing power.
- ☐ Multiple processes are executed concurrently.
- ☐ Creation of a process is time-consuming and resource intensive.
- ☐ Multiprocessing can be symmetric or asymmetric.
- ☐ The multiprocessing library in Python uses separate memory space, multiple CPU cores, bypasses GIL limitations in CPython, child processes are killable (ex. function calls in program) and is much easier to use.
- ☐ Some caveats of the module are a larger memory footprint and IPC's a little more complicated with more overhead.

4.4 Explanation for Behaviour: MultiThreading in Py

- ☐ Multithreading creates multiple threads of a single process to increase computing power.
- ☐ Multiple threads of a single process are executed concurrently.
- ☐ Creation of a thread is economical in both sense time and resource.
- ☐ The multithreading library is lightweight, shares memory, responsible for responsive UI and is used well for I/O bound applications.
- ☐ The module isn't killable and is subject to the GIL.
- ☐ Multiple threads live in the same process in the same space, each thread will do a specific task, have its own code, own stack memory, instruction pointer, and share heap memory.
- ☐ If a thread has a memory leak it can damage the other threads and parent process.

4.5 Explanation for Behaviour: Explained with Data

- For Smaller Sample Thread works better but on creation of too many threads then comes context switch in play which acts as a barrier
- But Creating a Thread is a cheap process when creating limited number of threads as they all are sharing same resources and hence quite faster than Multiple Processing at that time.
- Multiple Processors have their own resource bucket, they are GIL safe and can work in parallel because of this reason with increase in data they worked better than threading

Image Number	Single Core	Multiple Core	Multiple Threads (100)
1	0.16798	0.10107	0.09719
2	0.32265	0.19383	0.12229
3	0.46663	0.31848	0.13236
4	0.58964	0.47482	0.13739
5	0.76592	0.65649	0.13739

4.6 Analysing MultiProcessing v/s MultiThreading v/s Single Core Execution



Fig 6.1 Image vs Time graph for single core execution

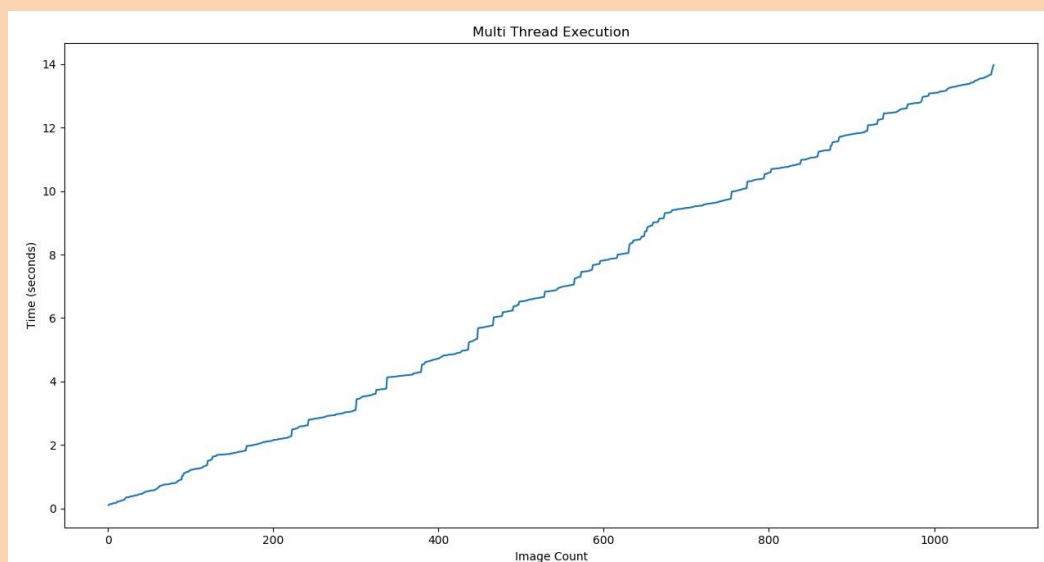


Fig 6.2 Image vs Time graph for multi thread execution

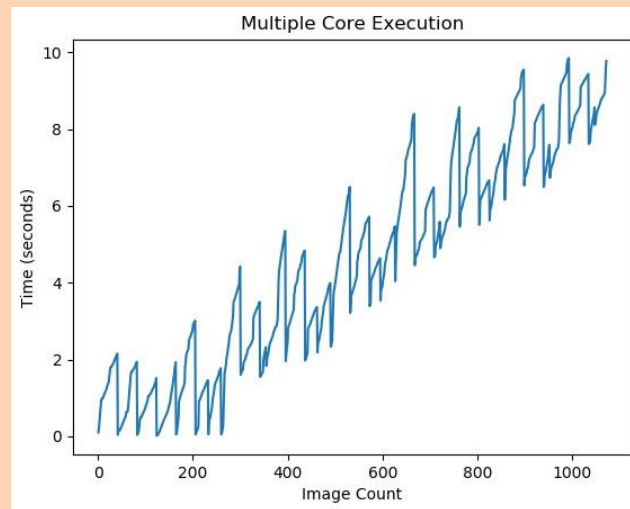


Fig 6.3 Image vs Time graph for multi core execution

Explanation

Here we are having this pattern because each core had an equal amount of work divided, now they all merging their chunked worked output together which is creating this wave-like pattern. It is observable this line graph can be broken down into chunks of the group of 8 cores generating output together for their particular work batch.



Fig 6.4 Merged graph for multi core execution

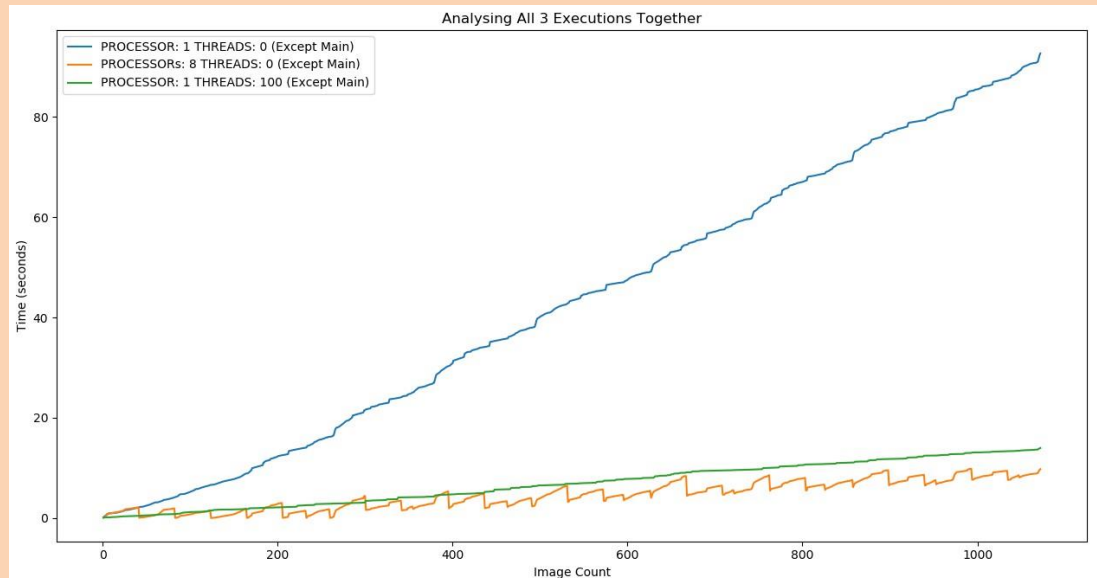


Fig 6.5 Analysis of all 3 executions

Explanation

We can clearly observe here that the Time Taken for Single Core, Single Threaded application increases exponentially with Time. But for MultiProcessing and MultiThreading, it increases with a little slope value. An addition to the observation is the fact that MultiProcessing Library is beating the thread as data is increased this is because MultiProcessing is resource extensive and so takes higher time to start it's processing but then it processes the data like a charm whereas Thread have interrupts which are continuously degrading the performance of the application and it is not able to perform any better.

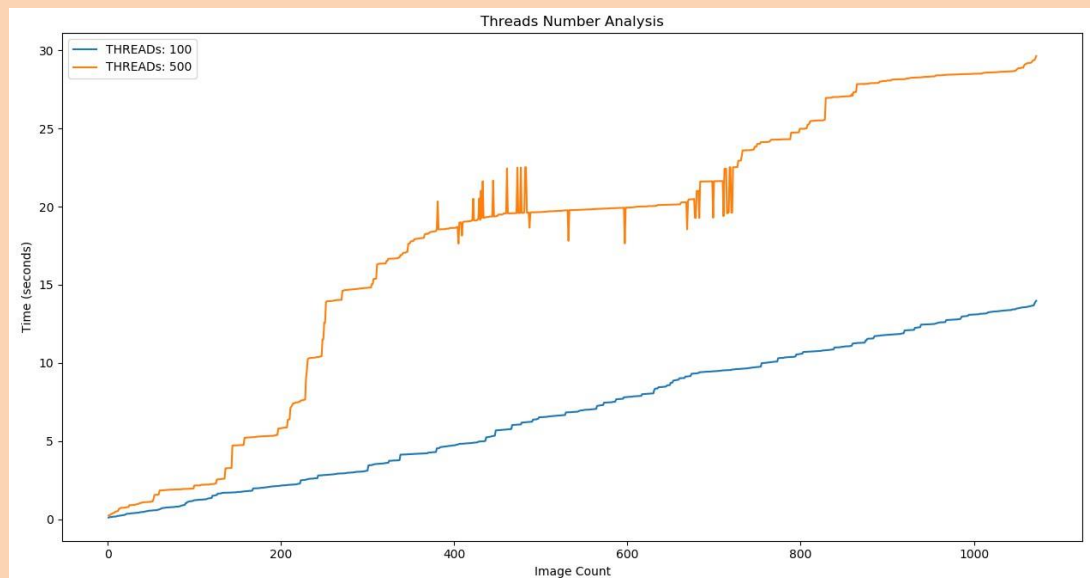


Fig 6.6 How many threads is too many?

Explanation

When we increase the number of threads then Time Taken to process the complete shall have decreased but what we observe is that because Interrupts things get worsen and it instead takes more time to process the complete data. When we further increased the number of threads to 5000 then the output had only 586/1072 outputs generated and else went missing, one possible explanation for this may be that older threads were killed because of the heavy load to save the program from closing abruptly.

So we conclude from this that higher number of threads add to the overhead of thread creation and the interrupts degrade the performance.

Some Strangely Interesting Finding:

- Synchronized Thread Script with **5K threads**
 - Time Taken: **10.31 sec**
 - Time Taken to process 1st Image: **2.16 sec**
 - Strange finding: **Only 586/1072 were processed**
- Similar Happened when executed with 500 Threads but only a few times
 - Rise in Number of Interrupts, OS kicks in

We found that:

- CPU Utilization becomes 100% for both MultiProcessing and MultiThreading
- CPU utilization remains lower and around 50% which is normal for Single Processing and Single Threaded application

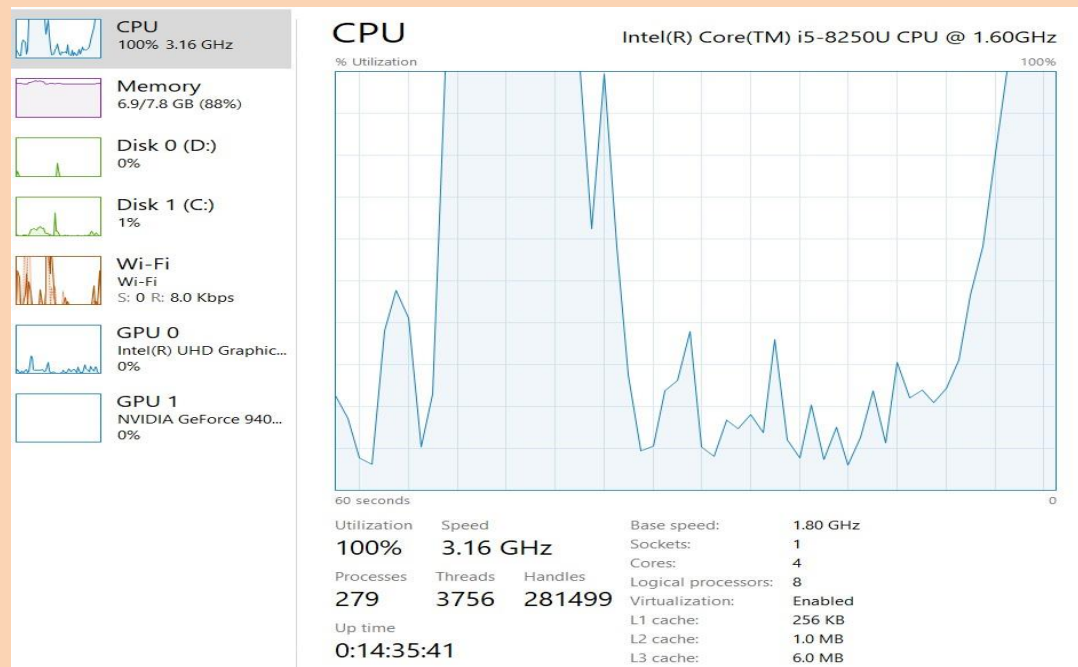


Fig 7. CPU Utilization

5. CONCLUSION AND FUTURE WORK

Use of Threading is enough and more suited in Applications like GUI and Networking as it's less resource intensive, Both Multi Processing and Multi-Threading have similar performance in I/O bound operations.

Since Most online servers provide with single core only so Threaded apps are more used but for Research purposes we have need more computation power and parallel execution with less complication and there comes Multiple Processors in to rescue.

Future Scope:

- ☐ Added functionality of removing the images as soon as we find it.
- ☐ Adding a GUI
- ☐ Web Application/Mobile Application
- ☐ Scaling the algorithm with time and making it more feasible for even larger dataset.

REFERENCES

- [1] Wang, Z., Bovik, A. C., Sheikh, H. R., & Simoncelli, E. P. (2004). Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4), 600-612.
- [2] García, G. B., Suarez, O. D., Aranda, J. L. E., Tercero, J. S., Gracia, I. S., & Enano, N. V. (2015). *Learning image processing with OpenCV*. Packt Publishing Ltd.
- [3] Hosoi, T., Kobayashi, K., Ito, K., & Aoki, T. (2011, September). Fast image inpainting using similarity of subspace method. In *2011 18th IEEE International Conference on Image Processing* (pp. 3441-3444). IEEE.
- [4] Narayanan, S., & Thirivikraman, P. K. (2015). Image similarity using fourier transform. *Journal Impact Factor*, 6(2), 29-37.
- [5] Bradski, G., & Kaehler, A. (2008). *Learning OpenCV: Computer vision with the OpenCV library*. " O'Reilly Media, Inc."
- [6] Xie, G., & Lu, W. (2013). Image Edge Detection Based On Opencv. *International Journal of Electronics and Electrical Engineering*, 1(2), 104-6.
- [7] Singh, N., Browne, L. M., & Butler, R. (2013). Parallel astronomical data processing with Python: Recipes for multicore machines. *Astronomy and Computing*, 2, 1-10.
- [8] Saha, A. (2012). Parallel programming in C and Python. *Linux journal*, 2012(217), 4.
- [9] Malakhov, A. (2016, July). Composable multi-threading for Python libraries. In *Proceedings of the Python in Science Conference*.
- [10] Meier, L., Honegger, D., & Pollefeys, M. (2015, May). PX4: A node-based multithreaded open source robotics framework for deeply embedded platforms. In *2015 IEEE international conference on robotics and automation (ICRA)* (pp. 6235-6240). IEEE.

