
Content-based Image Search

Parallelization using the OpenMP library

Group members:

Faridreza Momtaz Zandi

Mahya EhsaniMehr

Ali Mojahed

Spring 2023

::General description::

In this project, which was written in Python, we wanted to be able to find similar images using an image query by having a dataset of random and unrelated images. For this, there are several important factors that we must apply to the images and use each of them to finally make the best choice for the final images. These factors allow us to remove noises, obtain a histogram, and apply appropriate features to obtain a feature vector and finally find the similarity between two images by finding the Euclidean distance.

The first problem we encountered was not the possibility of using OpenMP commands in Python, and only one library supported it called Numba, which was very troublesome and could only work in very special and isolated conditions, by checking libraries like pypm, we concluded that Python is not a suitable language for working with OpenMP and it is better to use languages that support this library natively, so we decided to rewrite the program in our own C++ language.

::Code Review::

Together, we will review the program code and the use of each section:

```
1  #include <iostream>
2  #include <fstream>
3  #include <vector>
4  #include <dirent.h>
5  #include <opencv2/opencv.hpp>
6  #include <chrono>
7  #include <omp.h>
```

The code includes several libraries for I/O operations, image processing, scheduling, and parallelization using OpenMP.

```
11 // Function to resize an image to a given size
12 cv::Mat resizeImage(const cv::Mat& image, int width, int height)
13 {
14     cv::Mat resizedImage;
15     cv::resize(image, resizedImage, cv::Size(width, height));
16     return resizedImage;
17 }
```

resizeImage: resizes an image to the specified width and length.

```
19 // Function to compute the mean value of an image
20 double computeMean(const cv::Mat& image)
21 {
22     cv::Scalar meanVal = cv::mean(image);
23     return meanVal[0];
24 }
25
```

computeMean: calculates the mean value of an image.

```

26 // Function to compute the median value of an image
27 double computeMedian(const cv::Mat& image)
28 {
29     cv::Mat sortedImage;
30     cv::sort(image, sortedImage, cv::SORT_EVERY_COLUMN | cv::SORT_ASCENDING);
31     int totalPixels = image.rows * image.cols;
32
33     double medianValue;
34     if (totalPixels % 2 == 0)
35     {
36         int index1 = (totalPixels / 2) - 1;
37         int index2 = index1 + 1;
38         medianValue = (sortedImage.at<uchar>(index1) + sortedImage.at<uchar>(index2)) / 2.0;
39     }
40     else
41     {
42         int index = (totalPixels / 2);
43         medianValue = sortedImage.at<uchar>(index);
44     }
45
46     return medianValue;
47 }
48

```

computeMedian: calculates the median value of an image.

```

49 // Function to compute the histogram of an image
50 cv::Mat computeHistogram(const cv::Mat& image)
51 {
52     int numBins = 256; // Number of bins for the histogram
53     int histSize[] = { numBins };
54     float range[] = { 0, 256 };
55     const float* ranges[] = { range };
56     int channels[] = { 0 }; // Compute histogram only for the first channel (grayscale image)
57
58     cv::Mat histogram;
59     cv::calcHist(&image, 1, channels, cv::Mat(), histogram, 1, histSize, ranges);
60
61     return histogram;
62 }
63

```

computeHistogram: calculates the histogram of an image.

```

64 // Function to compute the standard deviation of an image
65 double computeStandardDeviation(const cv::Mat& image)
66 {
67     double mean = computeMean(image);
68     cv::Scalar meanSquaredDiff = cv::mean((image - mean).mul(image - mean));
69     return std::sqrt(meanSquaredDiff[0]);
70 }
71

```

computeStandardDeviation: Calculates the standard deviation of an image.

```

72 // Function to compute the Hu Moments of an image
73 cv::Mat computeHuMoments(const cv::Mat& image)
74 {
75     cv::Moments moments = cv::moments(image);
76     cv::Mat huMoments;
77     cv::HuMoments(moments, huMoments);
78     return huMoments;
79 }
80

```

computeHuMoments: Computes the Hu moments of an image.

```

81 // Function to compute the Euclidean distance between two feature vectors
82 double computeDistance(const cv::Mat& queryFeatureVector, const cv::Mat& datasetFeatureVector)
83 {
84     cv::Mat diff;
85     cv::absdiff(queryFeatureVector, datasetFeatureVector, diff);
86     cv::Mat squaredDiff = diff.mul(diff);
87     cv::Scalar sum = cv::sum(squaredDiff);
88     return std::sqrt(sum[0]);
89 }

```

computeDistance: calculates the Euclidean distance between two feature vectors.

```

91  std::vector<std::string> getFilesInDirectory(const std::string& dirPath)
92  {
93      std::vector<std::string> fileNames;
94      DIR* directory;
95      struct dirent* entry;
96
97      directory = opendir(dirPath.c_str());
98      if (directory != nullptr)
99      {
100         while ((entry = readdir(directory)) != nullptr)
101         {
102             std::string fileName = entry->d_name;
103             if (fileName != "." && fileName != "..")
104             {
105                 fileNames.push_back(dirPath + fileName);
106             }
107         }
108         closedir(directory);
109     }
110
111     return fileNames;
112 }
113

```

The `getFilesInDirectory` function takes a directory path as input and returns a vector of file names in that directory. It uses the `dirent.h` library to read the directory contents.

```

114 int main()
115 {
116     auto start = high_resolution_clock::now();
117     std::string datasetPath = "./dataset/images/"; // Path to the dataset images
118     std::string queryImagePath = "000000124442.jpg"; // Path to the query image
119
120     std::vector<std::string> imagePaths = getFilesInDirectory(datasetPath);
121
122     // Load and resize the query image
123     cv::Mat queryImage = cv::imread(queryImagePath, cv::IMREAD_GRAYSCALE);
124     queryImage = resizeImage(queryImage, 500, 500);
125
126     // Apply noise reduction techniques to the query image
127     cv::Mat queryImageSmallBlur;
128     cv::Mat queryImageLargeBlur;
129     cv::Mat queryImageGaussianBlur;
130
131     cv::blur(queryImage, queryImageSmallBlur, cv::Size(3, 3));
132     cv::blur(queryImage, queryImageLargeBlur, cv::Size(9, 9));
133     cv::GaussianBlur(queryImage, queryImageGaussianBlur, cv::Size(9, 9), 0);
134
135     // Compute features for the query image
136     double queryMean = computeMean(queryImage);
137     double queryMedian = computeMedian(queryImage);
138     double queryStdDev = computeStandardDeviation(queryImage);
139     cv::Mat queryHuMoments = computeHuMoments(queryImage);
140     cv::Mat queryHistogram = computeHistogram(queryImage);
141
142     // Compute features for all dataset images
143     std::vector<std::string> bestImagePaths;
144     std::vector<double> bestImageDistances(20, std::numeric_limits<double>::max());
145

```

The main function is the entry point of the program. Performs the following steps:

- Initializes the start time to measure the duration of the run.
- Defines the path of dataset images and query image.
- It calls the function `getFilesInDirectory` to get the paths of all the images in the dataset.
- Loads the query image and resizes it.
- It applies noise reduction techniques (small blurs, large blurs and Gaussian blurs) to the query image.
- Calculates various features (mean, median, standard deviation, Ho-moments and histogram) for the query image.

- Launches containers to store the best image paths and their corresponding distances.

```
146     omp_set_num_threads(3); // Set the number of threads to 3
147
148     #pragma omp parallel for
149     for (int i = 0; i < imagePaths.size(); ++i)
150     {
151         // Load and resize the dataset image
152         cv::Mat datasetImage = cv::imread(imagePaths[i], cv::IMREAD_GRAYSCALE);
153         datasetImage = resizeImage(datasetImage, 500, 500);
154
155         // Apply noise reduction techniques to the dataset image
156         cv::Mat datasetImageSmallBlur;
157         cv::Mat datasetImageLargeBlur;
158         cv::Mat datasetImageGaussianBlur;
159
160         cv::blur(datasetImage, datasetImageSmallBlur, cv::Size(3, 3));
161         cv::blur(datasetImage, datasetImageLargeBlur, cv::Size(9, 9));
162         cv::GaussianBlur(datasetImage, datasetImageGaussianBlur, cv::Size(9, 9), 0);
163
164         // Compute features for the dataset image
165         double datasetMean = computeMean(datasetImage);
166         double datasetMedian = computeMedian(datasetImage);
167         double datasetStdDev = computeStandardDeviation(datasetImage);
168         cv::Mat datasetHuMoments = computeHuMoments(datasetImage);
169         cv::Mat datasetHistogram = computeHistogram(datasetImage);
170
171         // Compute similarity using Euclidean distance
172         double distance = computeDistance(cv::Mat(1, 1, CV_64F, &datasetMean), cv::Mat(1, 1, CV_64F, &queryMean))
173             + computeDistance(cv::Mat(1, 1, CV_64F, &datasetMedian), cv::Mat(1, 1, CV_64F, &queryMedian))
174             + computeDistance(cv::Mat(1, 1, CV_64F, &datasetStdDev), cv::Mat(1, 1, CV_64F, &queryStdDev))
175             + computeDistance(datasetHuMoments, queryHuMoments)
176             + computeDistance(datasetHistogram, queryHistogram);
177     }
```

It sets the number of threads to 3 using the OpenMP `omp_set_num_threads` function.

It uses OpenMP parallelization with a parallel `omp pragma` for instructions to process dataset images concurrently.

In the parallel region, each dataset image is loaded, resized, and processed similarly to the request image.

The similarity between the query image and each dataset image is calculated using the Euclidean distance formula.

The best image paths and distances are updated based on the calculated similarity values.


```

177
178 // Update the best image paths if the current distance is smaller than the current best distances
179 #pragma omp critical
180 {
181     for (int j = 0; j < bestImageDistances.size(); ++j)
182     {
183         if (distance < bestImageDistances[j])
184         {
185             bestImagePaths.insert(bestImagePaths.begin() + j, imagePath[i]);
186             bestImagePaths.resize(20);
187             bestImageDistances.insert(bestImageDistances.begin() + j, distance);
188             bestImageDistances.resize(20);
189             break;
190         }
191     }
192 }
193 }
194
195 // Save the names of the best 20 images to a file
196 std::ofstream outputFile("best_images.txt");
197 if (outputFile.is_open())
198 {
199     for (const auto& imagePath : bestImagePaths)
200     {
201         outputFile << imagePath << std::endl;
202     }
203     outputFile.close();
204 }
205 else
206 {
207     std::cerr << "Failed to open the output file." << std::endl;
208     return 1;
209 }
210 auto stop = high_resolution_clock::now();
211 auto duration = duration_cast<microseconds>(stop - start);
212 std::cout << "\nTime taken by function: " << duration.count() << " microseconds";
213
214 return 0;
215 }

```

After the parallel region, the names of the top 20 images are stored in a text file named "best_images.txt".

If the file cannot be opened, an error message is printed.

The end time is recorded and the running time is calculated.

The execution time is printed on the console.

::presentation of results::

First, we executed the above program with a smaller dataset and in four different forms, one to four fragments, and the results of this execution are shown in the diagram below.

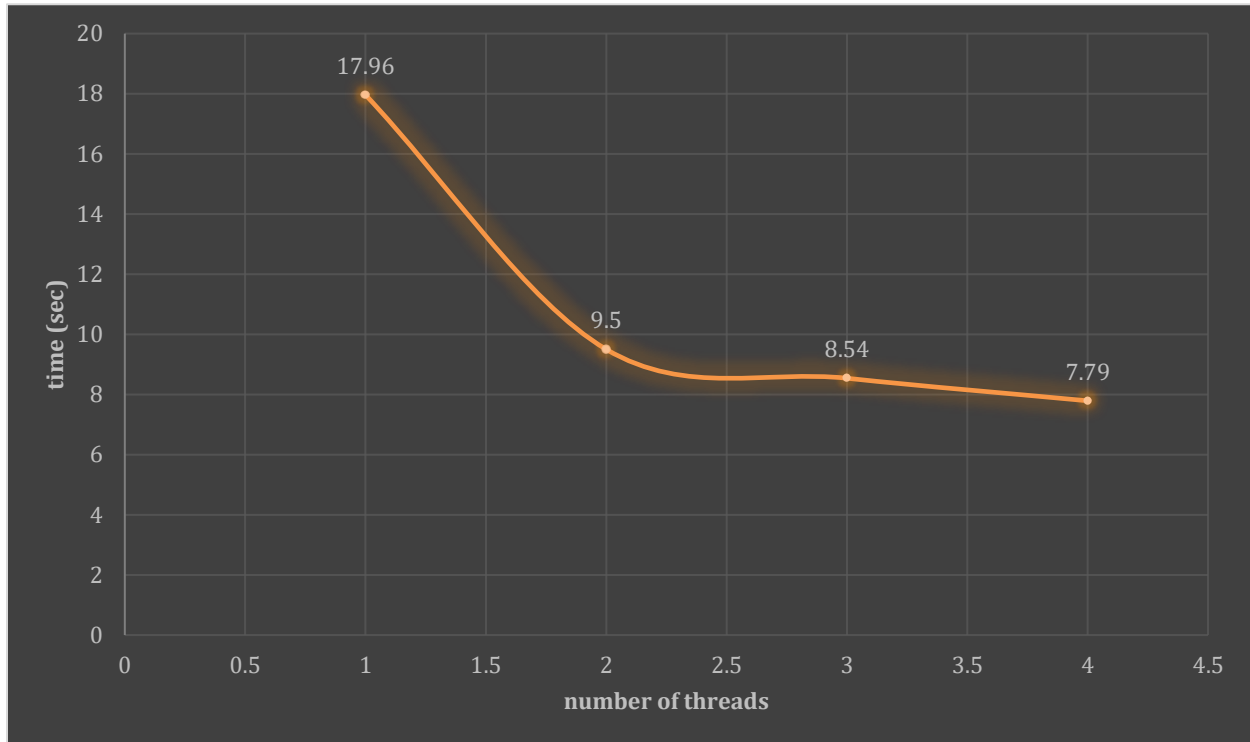


chart 1: executing code on 1215 images with 1 to 4 threads

Now we run the program with a larger dataset and in the form of one to eight chunks this time and record all the states, as we can see in the diagram below:

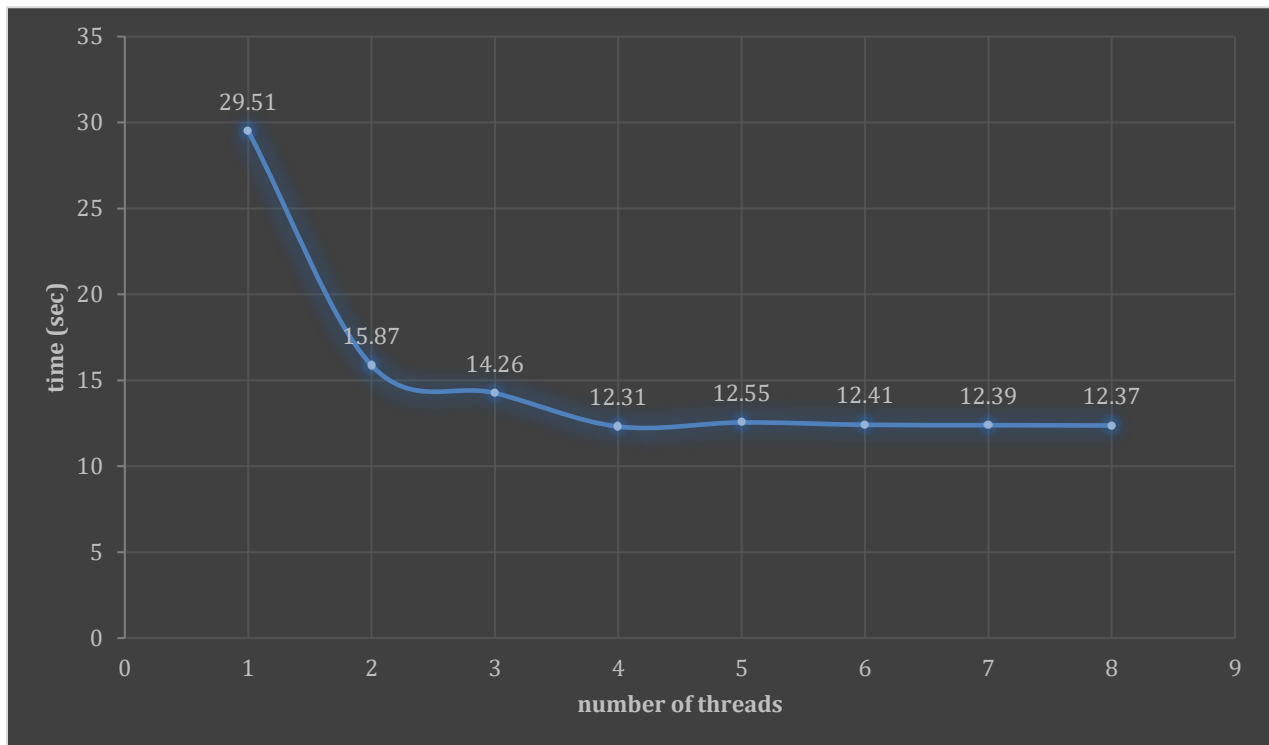
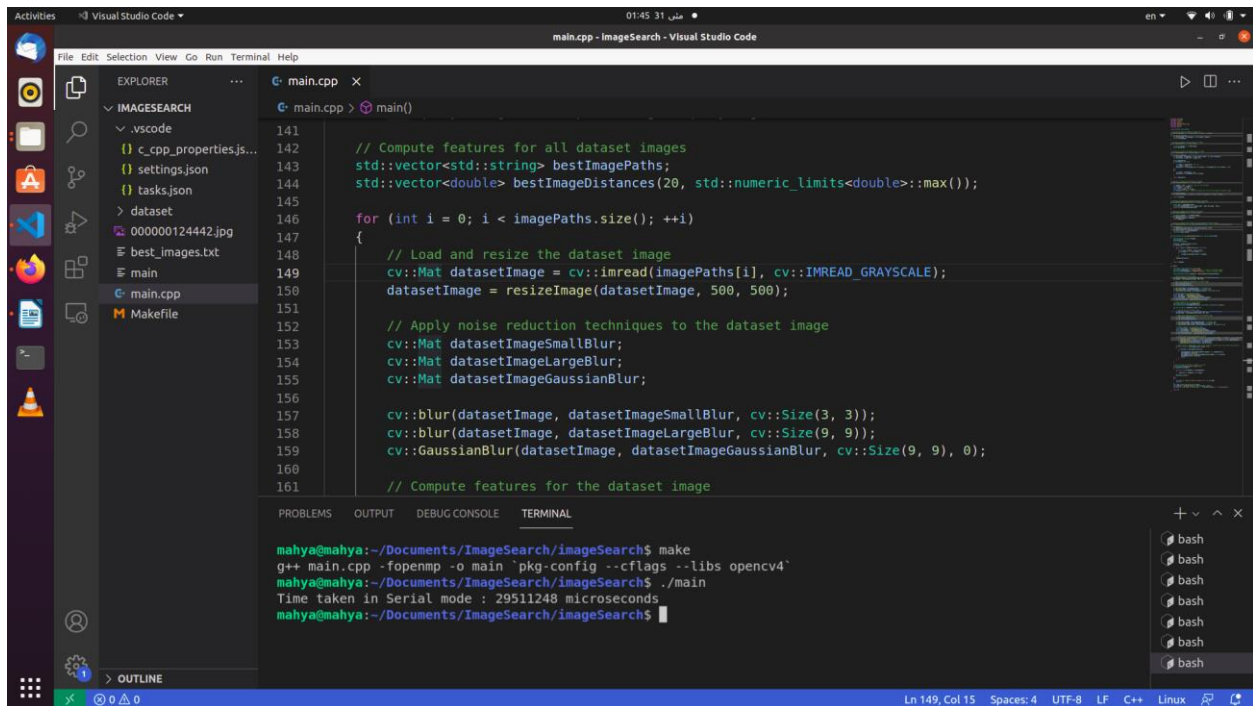


chart 2: executing code on 2015 images with 1 to 8 threads

Now, in the following images, we can see some examples of program execution in different conditions:



```
141
142
143 // Compute features for all dataset images
144 std::vector<std::string> bestImagePaths;
145 std::vector<double> bestImageDistances(20, std::numeric_limits<double>::max());
146
147 for (int i = 0; i < imagePaths.size(); ++i)
148 {
149     // Load and resize the dataset image
150     cv::Mat datasetImage = cv::imread(imagePaths[i], cv::IMREAD_GRAYSCALE);
151     datasetImage = resizeImage(datasetImage, 500, 500);
152
153     // Apply noise reduction techniques to the dataset image
154     cv::Mat datasetImageSmallBlur;
155     cv::Mat datasetImageLargeBlur;
156     cv::Mat datasetImageGaussianBlur;
157
158     cv::blur(datasetImage, datasetImageSmallBlur, cv::Size(3, 3));
159     cv::blur(datasetImage, datasetImageLargeBlur, cv::Size(9, 9));
160     cv::GaussianBlur(datasetImage, datasetImageGaussianBlur, cv::Size(9, 9), 0);
161
162     // Compute features for the dataset image
```

```
mahya@mahya:~/Documents/ImageSearch/imageSearch$ make
g++ main.cpp -fopenmp -o main `pkg-config --cflags --libs opencv4`
mahya@mahya:~/Documents/ImageSearch/imageSearch$ ./main
Time taken in Serial mode : 29511248 microseconds
mahya@mahya:~/Documents/ImageSearch/imageSearch$
```

Figure1 execution in serial mode

```

main.cpp
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
// Compute features for all dataset images
std::vector<std::string> bestImagePaths;
std::vector<double> bestImageDistances(20, std::numeric_limits<double>::max());

omp_set_num_threads(2);

#pragma omp parallel for
for (int i = 0; i < imagePaths.size(); ++i)
{
    // Load and resize the dataset image
    cv::Mat datasetImage = cv::imread(imagePaths[i], cv::IMREAD_GRAYSCALE);
    datasetImage = resizeImage(datasetImage, 500, 500);

    // Apply noise reduction techniques to the dataset image
    cv::Mat datasetImageSmallBlur;
    cv::Mat datasetImageLargeBlur;
    cv::Mat datasetImageGaussianBlur;
}

mahya@mahya:~/Documents/ImageSearch/imageSearch$ make
g++ main.cpp -fopenmp -o main `pkg-config --cflags --libs opencv4`
mahya@mahya:~/Documents/ImageSearch/imageSearch$ ./main
Time taken in Parallel mode : 15877104 microseconds
mahya@mahya:~/Documents/ImageSearch/imageSearch$

```

Figure 2 execution in parallel and two thread mode

```

main.cpp
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
// Compute features for all dataset images
std::vector<std::string> bestImagePaths;
std::vector<double> bestImageDistances(20, std::numeric_limits<double>::max());

omp_set_num_threads(4);

#pragma omp parallel for
for (int i = 0; i < imagePaths.size(); ++i)
{
    // Load and resize the dataset image
    cv::Mat datasetImage = cv::imread(imagePaths[i], cv::IMREAD_GRAYSCALE);
    datasetImage = resizeImage(datasetImage, 500, 500);

    // Apply noise reduction techniques to the dataset image
    cv::Mat datasetImageSmallBlur;
    cv::Mat datasetImageLargeBlur;
    cv::Mat datasetImageGaussianBlur;
}

mahya@mahya:~/Documents/ImageSearch/imageSearch$ make
g++ main.cpp -fopenmp -o main `pkg-config --cflags --libs opencv4`
mahya@mahya:~/Documents/ImageSearch/imageSearch$ ./main
Time taken in Parallel mode : 12317483 microseconds
mahya@mahya:~/Documents/ImageSearch/imageSearch$

```

Figure 3 execution in parallel and four thread which is the best

::Hardware specifications::

The machine on which this code was executed and the results obtained on Ubuntu 20.04 operating system and hardware are as follows:

Machine model: Lenovo IdeaPad 500-ISK

Processor: Intel core™ i5-6200U CPU @ 2.30GHz 2.4 GHz

Physical cores: 2 pcs

Supported chips: 4 pcs

RAM memory: 8 GB