

BIRZEIT UNIVERSITY

Birzeit University
Department of Electrical & Computer Engineering
ENCS5343 Computer Vision
Assignment#2

Content-Based Image Retrieval (CBIR) System
Using Color Features and other Techniques

Student: Rama Abdrlrahman

Id : 1191344

Instructur : Dr. Aziz Qaroush

January-2024

Contents

| | |
|---|----|
| 1. Introduction | 3 |
| 1.1 Overview of CBIR | 3 |
| 1.2 Theoretical Background | 3 |
| ○ Color Histograms | 3 |
| ○ Color Moments | 3 |
| ○ CNN feature extraction | 3 |
| 2. System Implementation | 4 |
| 2.1 CBIR System Architecture | 4 |
| 2.2 Programming Tools | 4 |
| 3. Experimental Setup and Results | 4 |
| 3.1 Tasks 1 and Task 2: CBIR System Using Color Features (Experiment with Color Histograms) | 5 |
| ○ Testing for 10 random queries:..... | 5 |
| Testing for specific 10 queries: | 26 |
| 3.2 Task 3: CBIR System Using Color Moments..... | 30 |
| Subsection 3.2.1: Experiment with Equal Weights | 30 |
| Subsection 3.2.2: Experiment with Different Weights | 33 |
| Subsection 3.2.3: Additional Moments..... | 58 |
| 3.3 Task 4: Improvement Using CNN..... | 76 |

1. Introduction

1.1 Overview of CBIR

Content-Based Image Retrieval (CBIR) systems are essential in the field of digital image processing and multimedia databases. CBIR refers to retrieving relevant images from large data sets based on the content of the images themselves, such as colors, shapes, textures, or any other information that can be derived from the image itself. This technology has a wide range of applications, including digital libraries, crime prevention, medical imaging, and even in the entertainment industry for visual effects and graphics.

1.2 Theoretical Background

- **Color Histograms**

A color histogram represents the distribution of colors in an image. It is one of the simplest yet most effective tools for image retrieval. In a CBIR system, color histograms serve as a feature vector to represent the color distribution of an image. The histogram is a count of the number of pixels in an image showing each distinct color. In more technical terms, it represents the frequency of the occurrence of various color bins (ranges of color values) in an image. These histograms are used to compare and retrieve images based on color similarity.

- **Color Moments**

Color moments are statistical measures that summarize the color distribution of an image. The most common moments used are the mean (which gives the average color), standard deviation (which measures the color variability), and skewness (which describes the asymmetry of the color distribution in an image). These moments can be calculated for each color channel and combined to form a feature vector representing the image in a CBIR system.

- **CNN feature extraction**

The VGG16 model is imported from the Keras library and configured to exclude its top classification layer. This allows the model to serve as a powerful feature extractor for images. The extracted features are then used to measure the cosine similarity between query images and the entire dataset, facilitating content-based image retrieval. The system evaluates its performance through metrics such as precision, recall, and F1-score, providing

insights into the effectiveness of the retrieval process. The CNN's ability to capture high-level image features enables the system to discern visual similarities among images in the dataset.

2. System Implementation

2.1 CBIR System Architecture

The architecture of a CBIR system typically consists of several key components:

1. **Image Database:** A collection of images which the system will search to retrieve relevant images.
2. **Feature Extraction Module:** This module processes each image in the database to extract meaningful features such as color histograms or color moments.
3. **Indexing and Storage Module:** Extracted features are indexed and stored for efficient retrieval.
4. **Query tests:** test query images or define features for retrieval. The same feature extraction process is applied to the query image.
5. **Similarity Measurement and Ranking Module:** Computes the similarity between the query features and database image features using distance measures like Euclidean distance. It then ranks the images based on similarity scores.

2.2 Programming Tools

The CBIR system can be implemented using various programming languages and libraries. Python is a popular choice due to its extensive libraries for image processing (such as OpenCV) and data handling (like NumPy and Pandas). In this work I have used the mentioned tools and libraries.

3. Experimental Setup and Results

3.1 Tasks 1 and Task 2: CBIR System Using Color Features (Experiment with Color Histograms)

In this section I will describe the experiments conducted with color histograms using 120, 180, and 360 bins. Including methodologies for calculating precision, recall, F1 score, and computational time. I will also explain the construction of the Receiver Operating Characteristic (ROC) curve and the calculation of the Area Under the Curve (AUC).

Note: In all tests below I am showing small samples of results (images and log output), for the full results you can go to the files or folders mentioned in each part.

- **Testing for 10 random queries:**

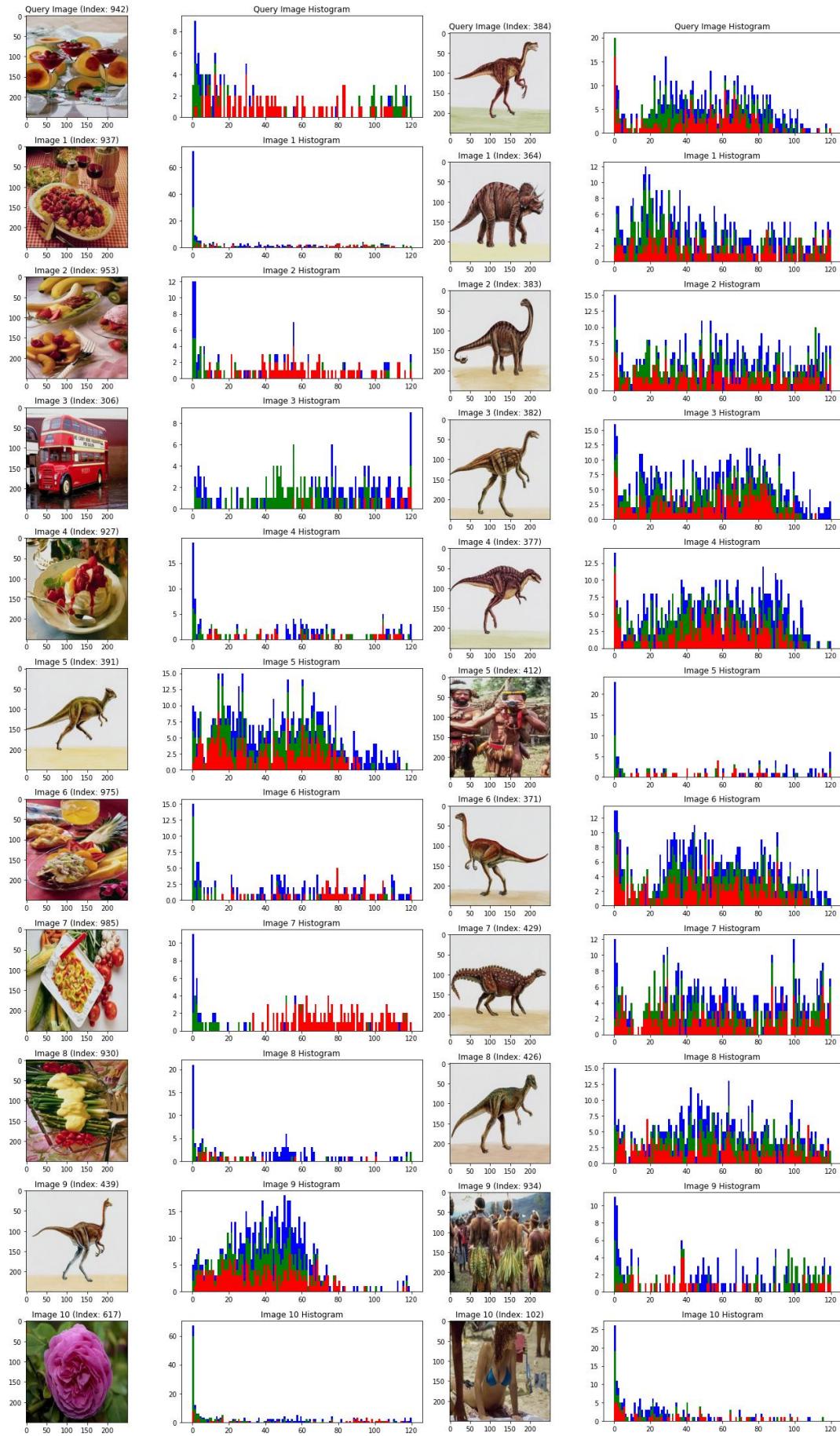
Here 10 random queries where chosen, and same threshold values applied on them. Some results will be shown for some tests, after showing the results, under each test results, we will have a short explanation.

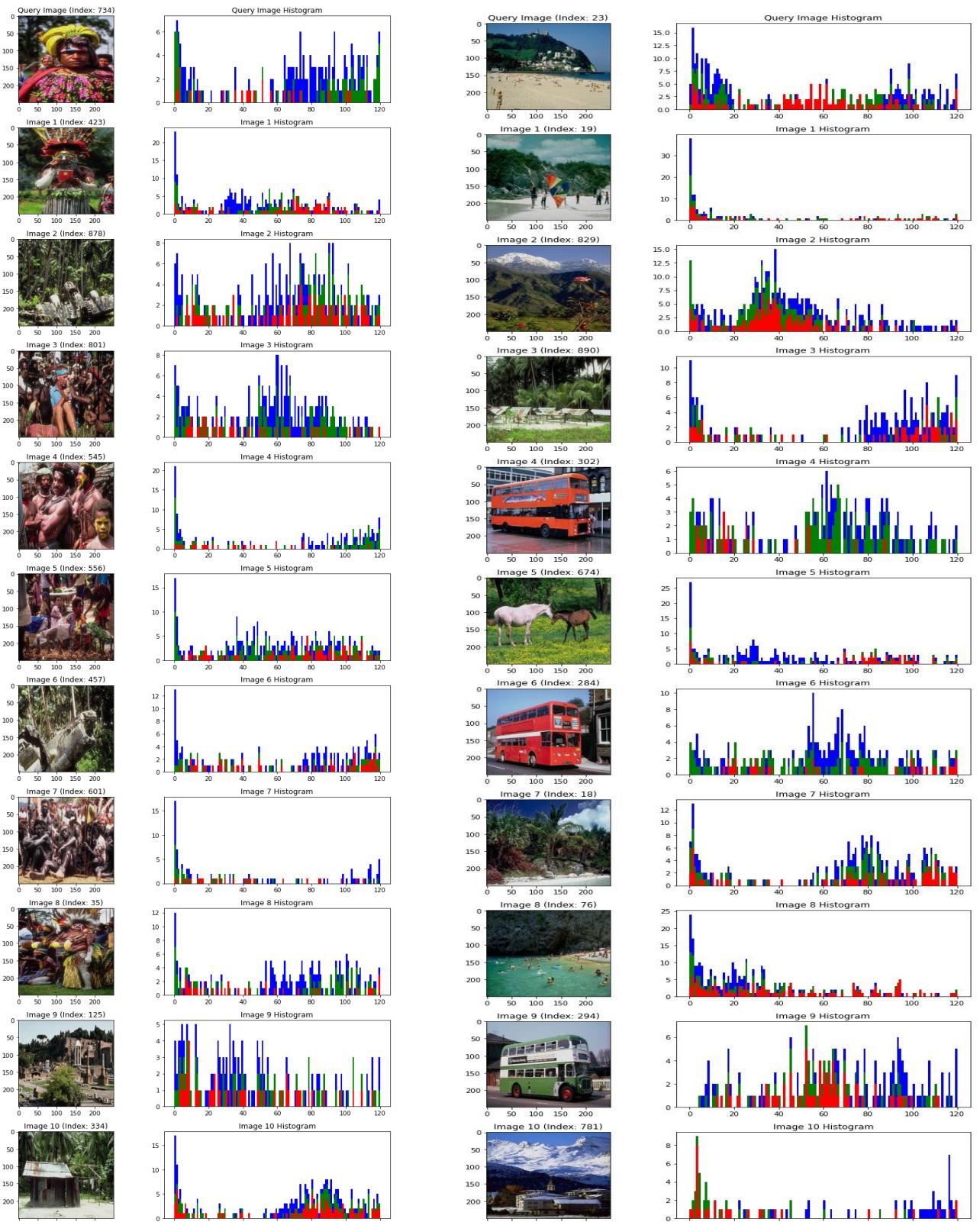
Results Using these thresholds for the 10 queries: (Test 1)

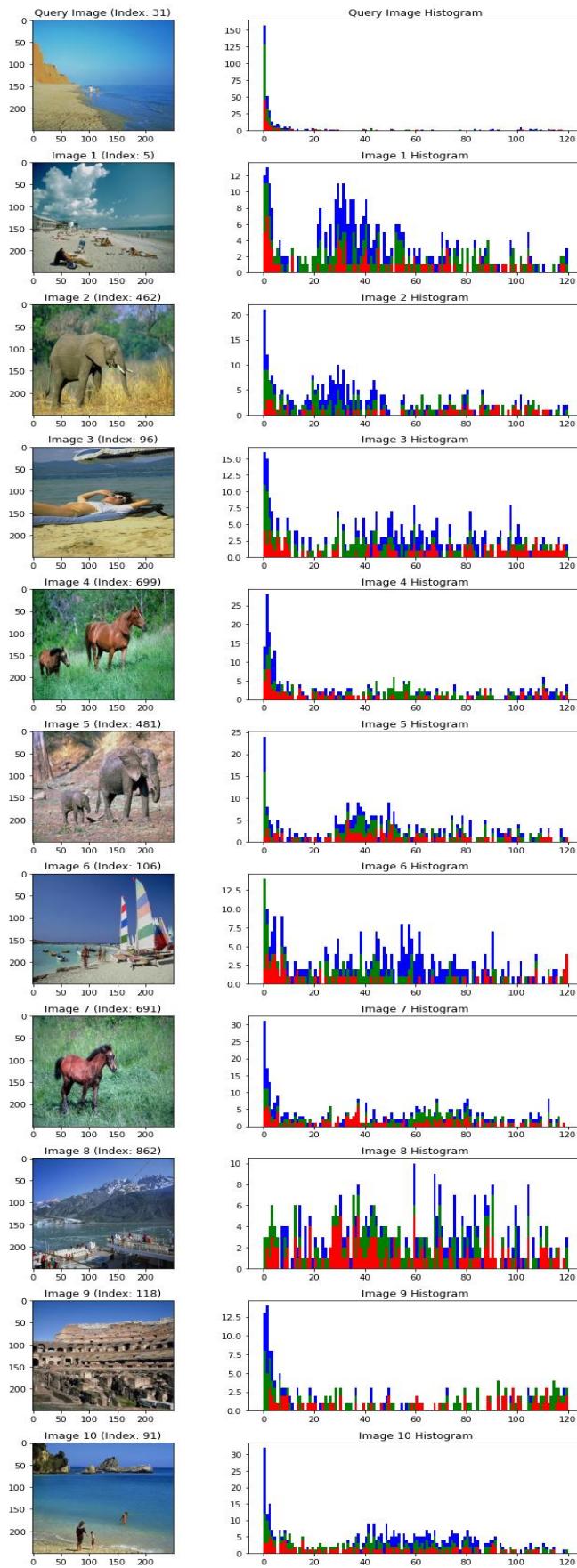
```
quantiles = [ 0.0111, 0.03, 0.05, 0.059, 0.07, 0.09, 0.099, .119, 0.15, 0.188,  
0.19, 0.2, 0.246, 0.3, 0.333, 0.36 ]
```

Results for 120 bins:

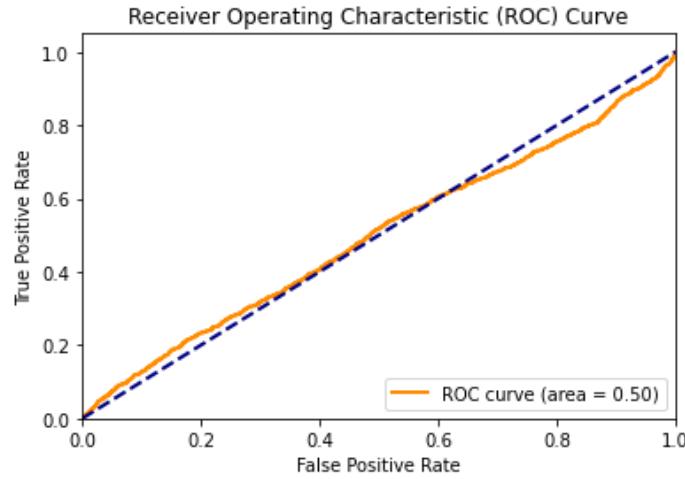
Here I am showing samples, rest of images for this test could be found in folder **RandomQuiriesTest1**.







I have used 10 queries but I can't show them all here since it will take too much space, you can check the folder I mentioned above.



Some results are shown here, the full results including precision, recall, f1score, execution time per query at each threshold for each experiment(bins numbers), and the average values of these measures per query could be found in file **test1.log**

```

532 Time taken for 120 bins: 79.9514 seconds
533
534 --- Average Evaluation Results for 120 Bins for All Queries at Each Threshold
535 Threshold: 0.0524
536 Average Precision: 0.6475
537 Average Recall: 0.0340
538 Average F1 Score: 0.0535
539
540 Threshold: 0.0632
541 Average Precision: 0.6068
542 Average Recall: 0.0678
543 Average F1 Score: 0.0662
544
545 Threshold: 0.0704
546 Average Precision: 0.5537
547 Average Recall: 0.0860
548 Average F1 Score: 0.0946
549
550 Threshold: 0.0731
551 Average Precision: 0.5545
552 Average Recall: 0.0940
553 Average F1 Score: 0.0982
554
555 Threshold: 0.0762
556 Average Precision: 0.5301
557 Average Recall: 0.1030
558 Average F1 Score: 0.1009
559
560 Threshold: 0.0811
561 Average Precision: 0.5080
562 Average Recall: 0.1220
563 Average F1 Score: 0.1054
564
565 Threshold: 0.0831
566 Average Precision: 0.5006
567 Average Recall: 0.1330
568 Average F1 Score: 0.1101
569
570 Threshold: 0.0873
571 Average Precision: 0.4637
572 Average Recall: 0.1660
573 Average F1 Score: 0.1245
574
575 Threshold: 0.0933
576 Average Precision: 0.3660
577 Average Recall: 0.1990
578 Average F1 Score: 0.1270
579
580 Threshold: 0.1001
581 Average Precision: 0.3244
582 Average Recall: 0.2350
583 Average F1 Score: 0.1293
584
585 Threshold: 0.1004
586 Average Precision: 0.2901
587 Average Recall: 0.2370
588 Average F1 Score: 0.1291
589
590 Threshold: 0.1021
591 Average Precision: 0.2914
592 Average Recall: 0.2450
593 Average F1 Score: 0.1288
594

```

```
Threshold: 0.1097
Average Precision: 0.2720
Average Recall: 0.2860
Average F1 Score: 0.1285

Threshold: 0.1185
Average Precision: 0.2233
Average Recall: 0.3340
Average F1 Score: 0.1303

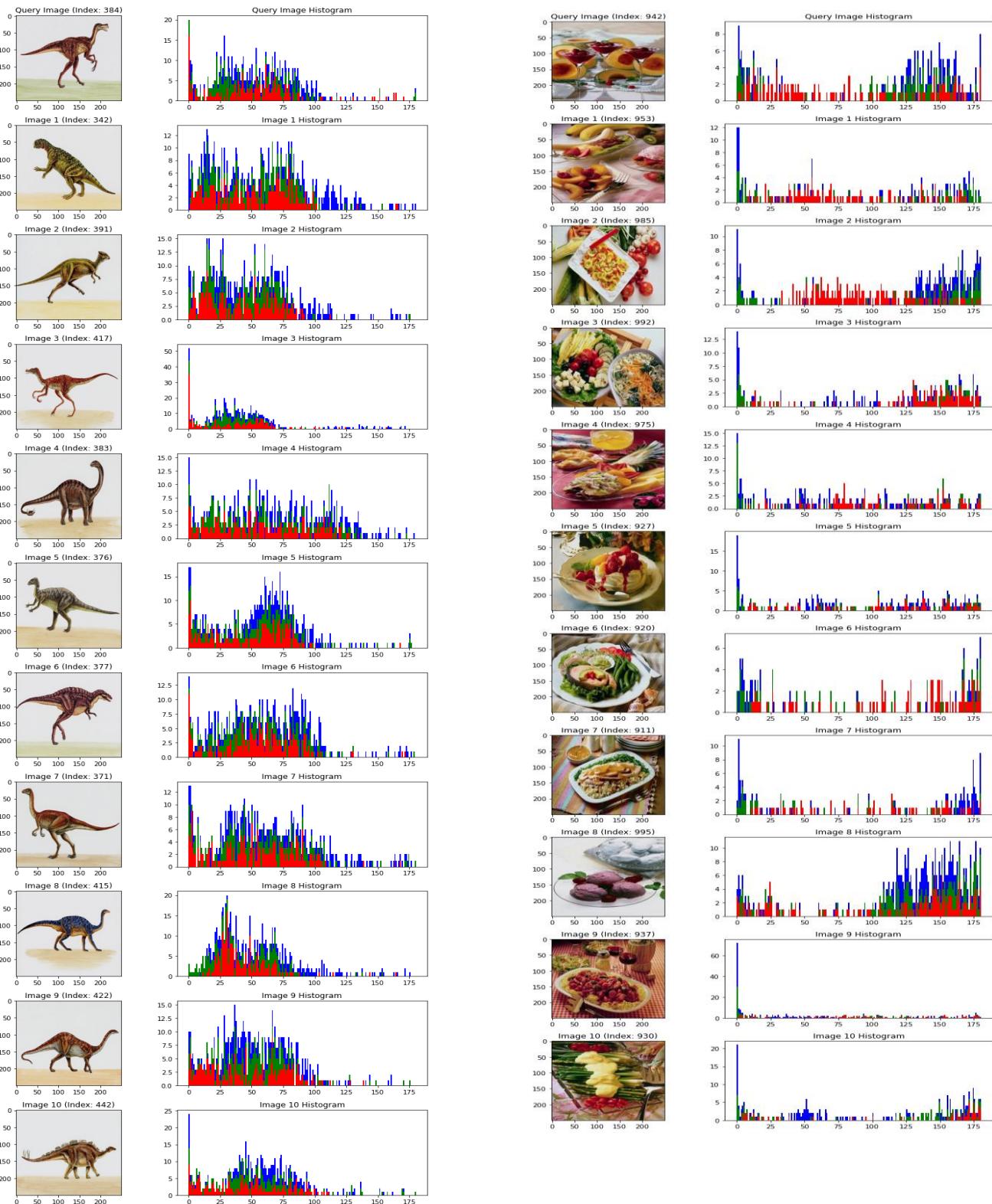
Threshold: 0.1237
Average Precision: 0.2145
Average Recall: 0.3560
Average F1 Score: 0.1293

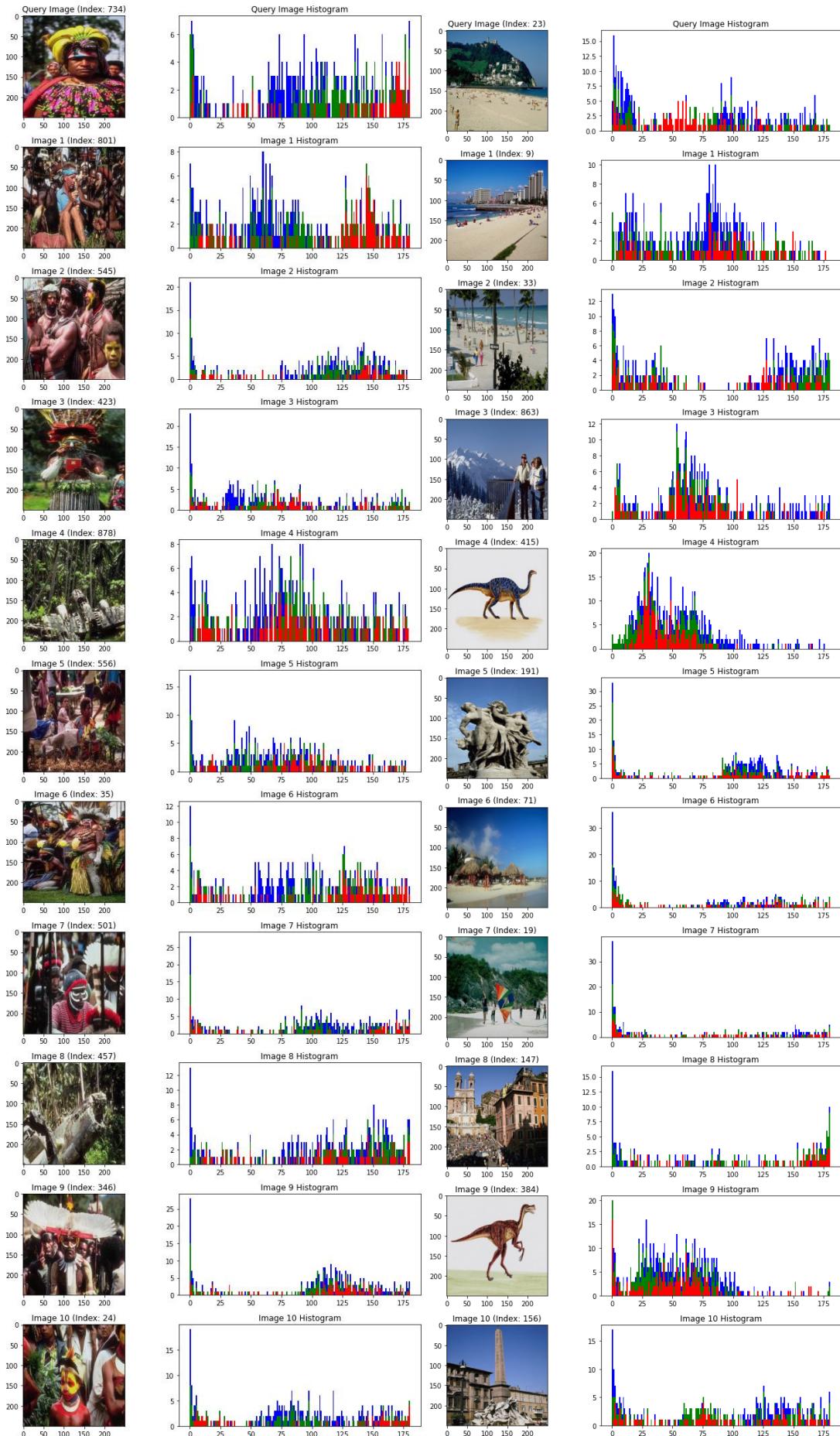
Threshold: 0.1279
Average Precision: 0.2105
Average Recall: 0.3740
Average F1 Score: 0.1291

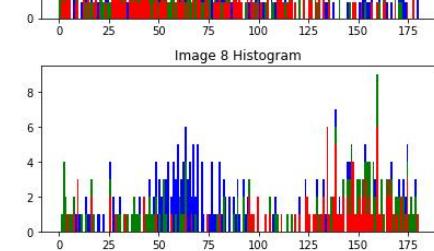
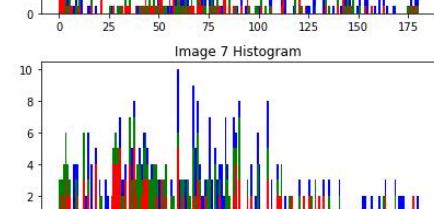
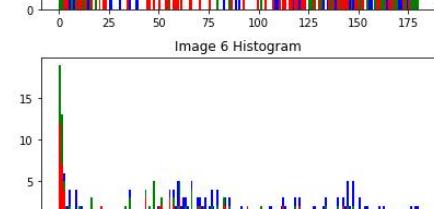
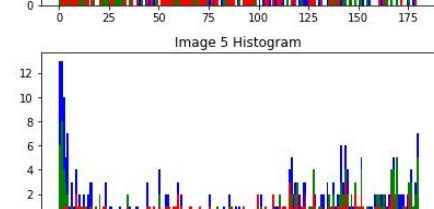
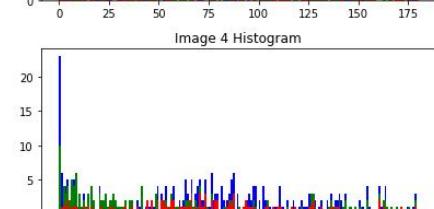
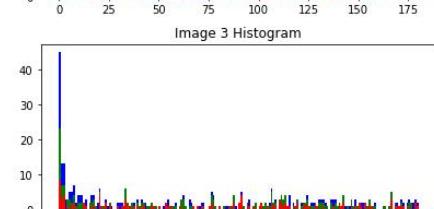
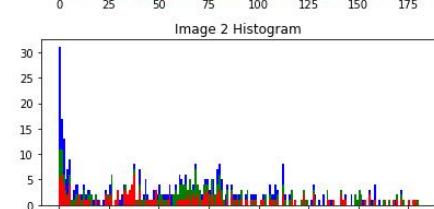
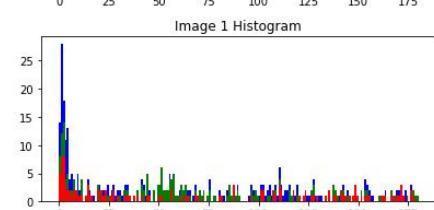
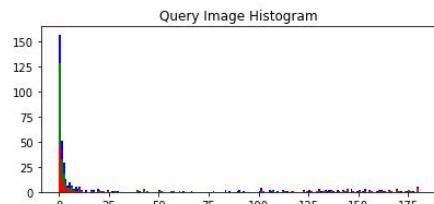
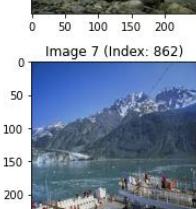
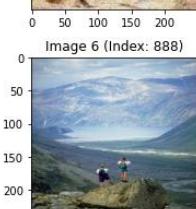
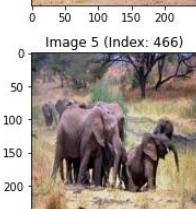
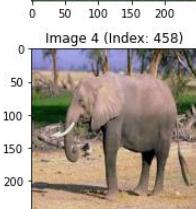
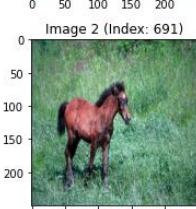
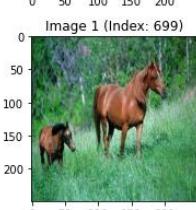
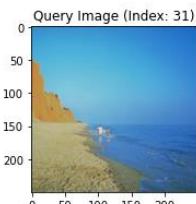
--- Overall Average Evaluation Results for 120 Bins ---
Overall Average Precision: 0.4098
Overall Average Recall: 0.1919
Overall Average F1 Score: 0.1128
```

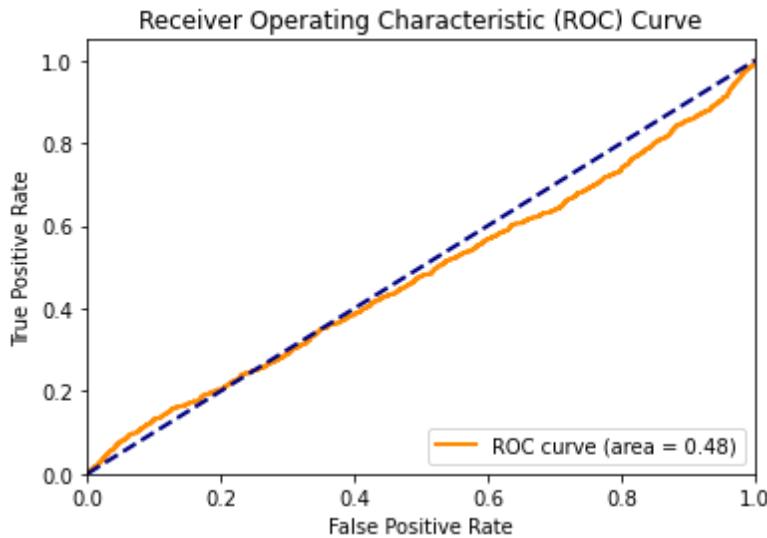
Results for 120 bins:

Here I am showing samples, rest of images for this test could be found in folder [RandomQuiriesTest1](#).









```

1752 Time taken for 180 bins: 100.5498 seconds
1753 --- Average Evaluation Results for 180 Bins for All Queries at Each Threshold -
1754 Threshold: 0.0638
1755 Average Precision: 0.6806
1756 Average Recall: 0.0270
1757 Average F1 Score: 0.0429
1758
1759 Threshold: 0.0765
1760 Average Precision: 0.6280
1761 Average Recall: 0.0570
1762 Average F1 Score: 0.0880
1763
1764 Threshold: 0.0845
1765 Average Precision: 0.5892
1766 Average Recall: 0.0940
1767 Average F1 Score: 0.1118
1768
1769 Threshold: 0.0874
1770 Average Precision: 0.5159
1771 Average Recall: 0.1050
1772 Average F1 Score: 0.1156
1773
1774 Threshold: 0.0906
1775 Average Precision: 0.5145
1776 Average Recall: 0.1200
1777 Average F1 Score: 0.1232
1778
1779 Threshold: 0.0957
1780 Average Precision: 0.4703
1781 Average Recall: 0.1340
1782 Average F1 Score: 0.1211
1783
1784 Threshold: 0.0978
1785 Average Precision: 0.4634
1786 Average Recall: 0.1420
1787 Average F1 Score: 0.1205
1788
1789 Threshold: 0.1023
1790 Average Precision: 0.3704
1791 Average Recall: 0.1570
1792 Average F1 Score: 0.1186
1793
1794 Threshold: 0.1084
1795 Average Precision: 0.3632
1796 Average Recall: 0.1920
1797 Average F1 Score: 0.1293
1798
1799 Threshold: 0.1152
1800 Average Precision: 0.3515
1801 Average Recall: 0.2310
1802 Average F1 Score: 0.1378
1803
1804 Threshold: 0.1156
1805 Average Precision: 0.3499
1806 Average Recall: 0.2330
1807 Average F1 Score: 0.1378
1808
1809 Threshold: 0.1173
1810 Average Precision: 0.3455
1811 Average Recall: 0.2450
1812 Average F1 Score: 0.1391
1813
1814

```

```

Threshold: 0.1249
Average Precision: 0.2643
Average Recall: 0.2880
Average F1 Score: 0.1445

Threshold: 0.1332
Average Precision: 0.2457
Average Recall: 0.3340
Average F1 Score: 0.1473

Threshold: 0.1382
Average Precision: 0.2311
Average Recall: 0.3650
Average F1 Score: 0.1532

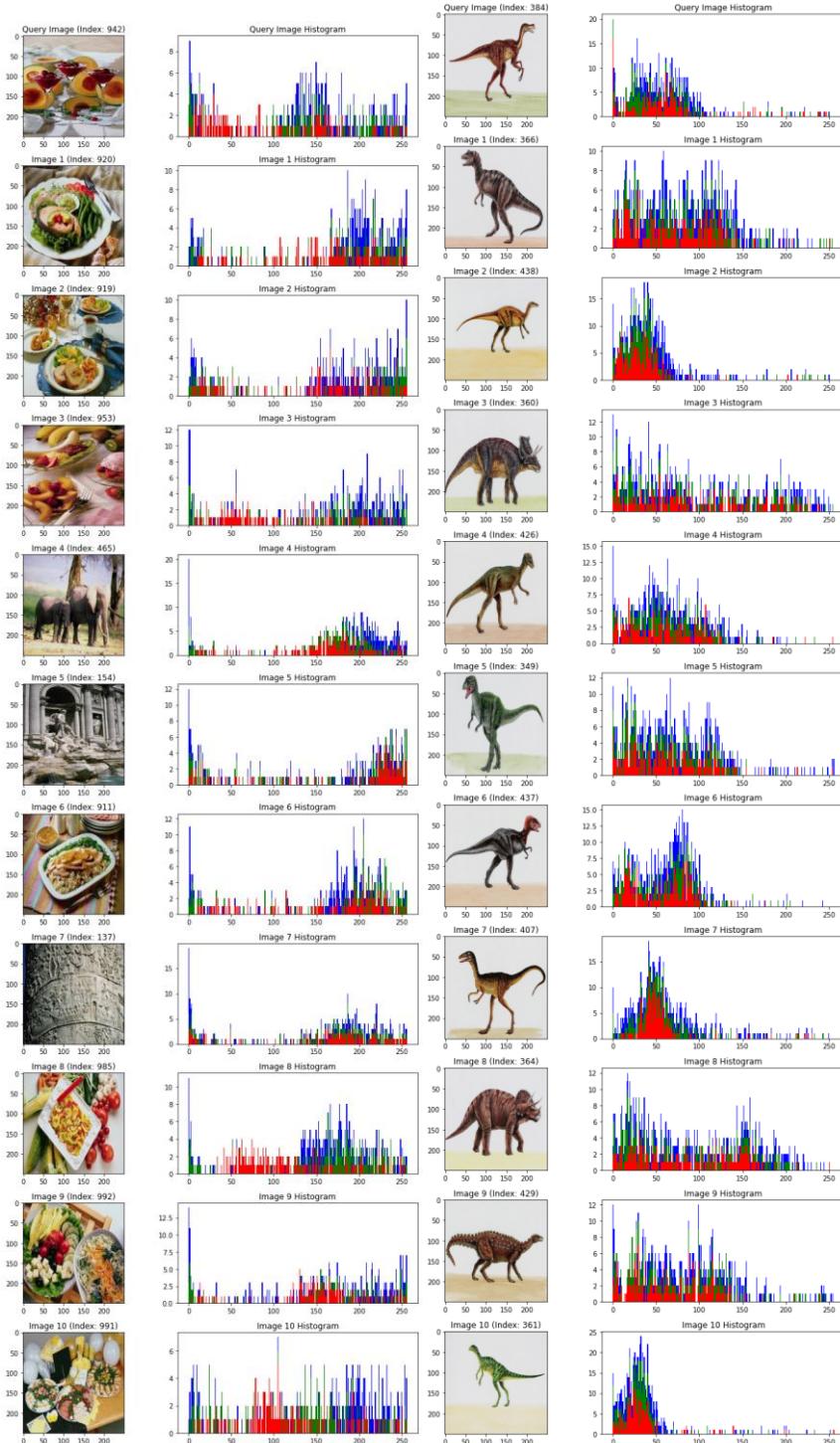
Threshold: 0.1424
Average Precision: 0.2161
Average Recall: 0.3810
Average F1 Score: 0.1516

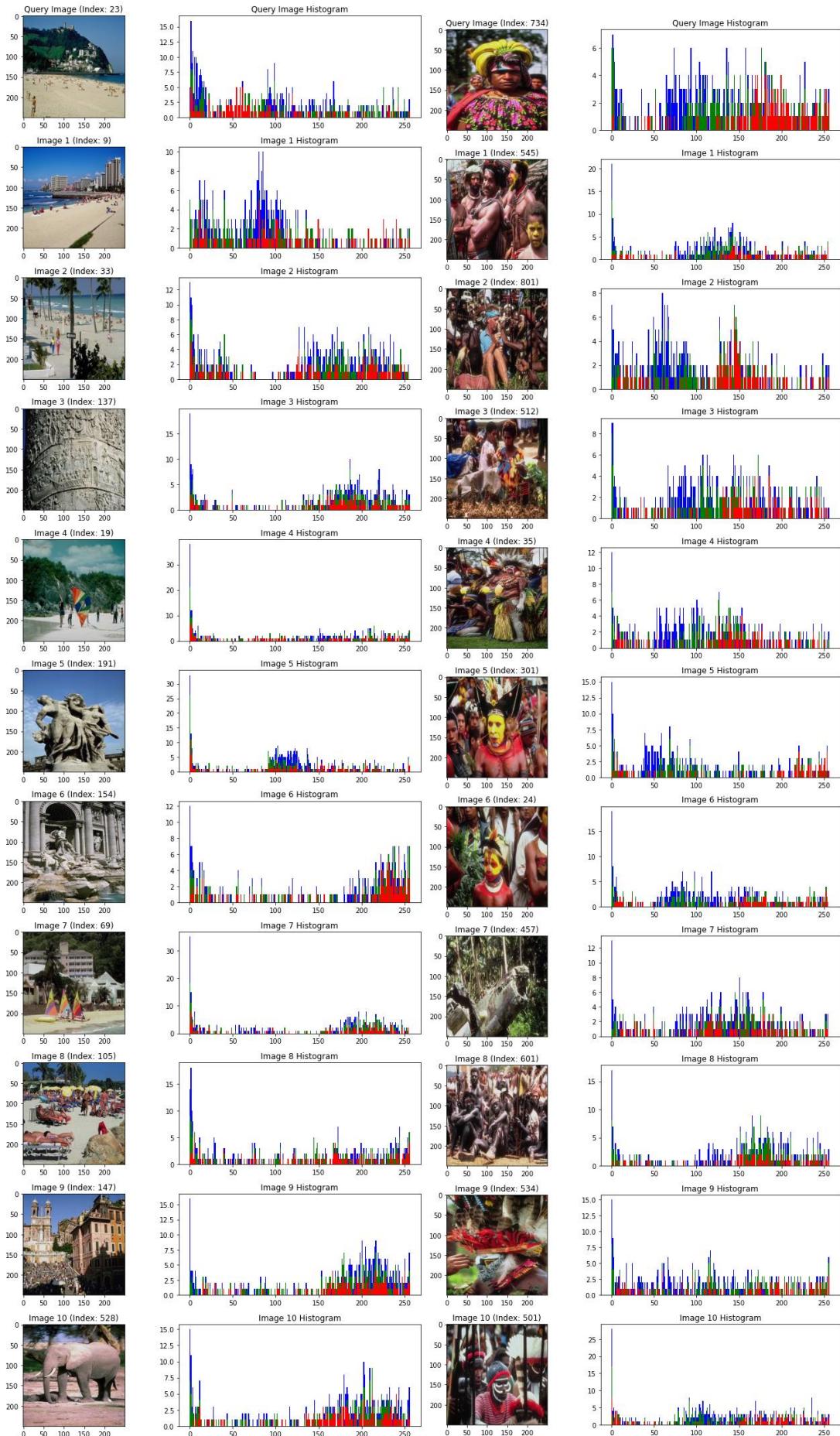
--- Overall Average Evaluation Results for 180 Bins ---
Overall Average Precision: 0.4125
Overall Average Recall: 0.1941
Overall Average F1 Score: 0.1234

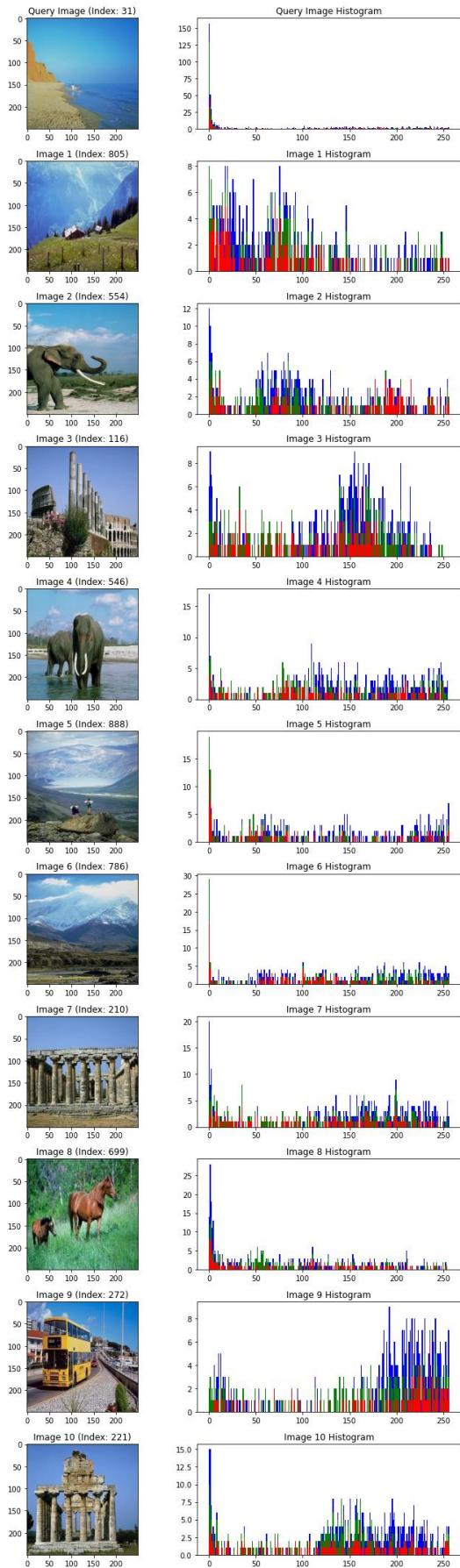
```

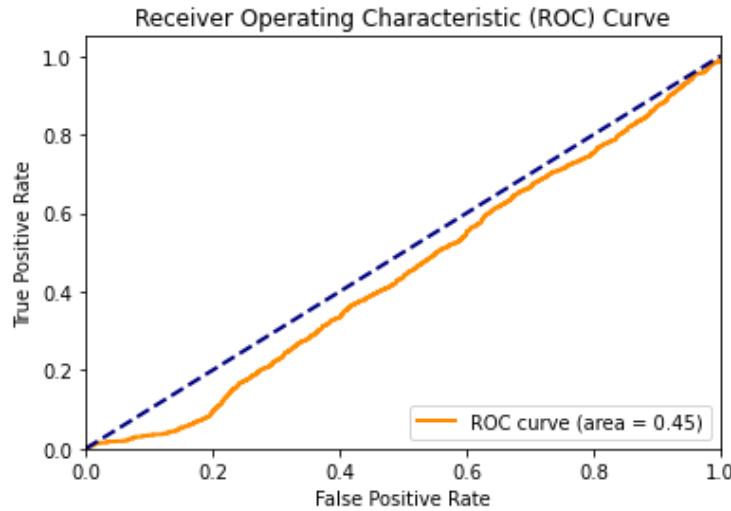
Results for 256 bins:

Here I am showing samples, rest of images for this test could be found in folder [RandomQuiriesTest1](#).









```

678
679
680 Time taken for 256 bins: 130.9678 seconds
681
682 --- Average Evaluation Results for 256 Bins for All Queries at Each Threshold -
683 Threshold: 0.0527
684 Average Precision: 0.6735
685 Average Recall: 0.0250
686 Average F1 Score: 0.0380
687
688 Threshold: 0.0632
689 Average Precision: 0.6235
690 Average Recall: 0.0450
691 Average F1 Score: 0.0548
692
693 Threshold: 0.0698
694 Average Precision: 0.5601
695 Average Recall: 0.0790
696 Average F1 Score: 0.0837
697
698 Threshold: 0.0723
699 Average Precision: 0.5357
700 Average Recall: 0.0890
701 Average F1 Score: 0.0892
702
703 Threshold: 0.0750
704 Average Precision: 0.5197
705 Average Recall: 0.1030
706 Average F1 Score: 0.0920
707
708 Threshold: 0.0793
709 Average Precision: 0.5081
710 Average Recall: 0.1310
711 Average F1 Score: 0.1010
712
713 Threshold: 0.0811
714 Average Precision: 0.5197
715 Average Recall: 0.1460
716 Average F1 Score: 0.1054
717
718 Threshold: 0.0847
719 Average Precision: 0.5052
720 Average Recall: 0.1650
721 Average F1 Score: 0.1061
722
723 Threshold: 0.0897
724 Average Precision: 0.4862
725 Average Recall: 0.1960
726 Average F1 Score: 0.1106
727
728 Threshold: 0.0954
729 Average Precision: 0.4105
730 Average Recall: 0.2260
731 Average F1 Score: 0.1134
732
733 Threshold: 0.0956
734 Average Precision: 0.4196
735 Average Recall: 0.2280
736 Average F1 Score: 0.1138
737
738 Threshold: 0.0970
739 Average Precision: 0.3958
740 Average Recall: 0.2390
741
742 --- Overall Average Evaluation Results for 256 Bins
743 Overall Average Precision: 0.4597
744 Overall Average Recall: 0.1824
745 Overall Average F1 Score: 0.1016
746
747 Completed all experiments.
748

```

Discussing the results above:

- **About the ROC curve and AUC**

120 Bins: The ROC curve has an AUC of 0.50, which suggests that the classifier does no better than random chance in distinguishing between the two classes.

180 Bins: The AUC here is slightly less than before, at 0.48, which indicates a potential decrease in the classifier's performance as the number of bins increases.

256 Bins: The AUC is 0.45, which further suggests a decrease in performance with an increase in the number of bins.

- **About the measures**

As shown on this test, as the number of bins increases from 120 to 256, the overall average precision and recall seem to fluctuate. There isn't a clear trend showing improvement or decline with the increase in bins. However, the precision appears to increase about 0.01 which is not worth to mention, and the overall F1 score appears to slightly decrease.

For 120 Bins the system has average precision and recall rates around 0.40 and an F1 score around 0.11. For 180 Bins precision and recall rates are similar to those with 120 bins, very little increase in precision, the F1 score slightly improves. For 256 Bins the average precision and recall do not improve significantly, very little increase in precision, and the F1 score decreases slightly compared to 180 bins. The time taken for computations increases with the number of bins, which is expected due to the increase in computational complexity.

The results suggest that increasing the number of histogram bins for this color-based retrieval system does not necessarily improve performance. It may slightly degrade the classifier's ability to distinguish between the positive and negative classes, as indicated by the ROC curves' AUC values. Additionally, it seems to increase the computational load without a corresponding increase in accuracy. The performance degrading with more bins has many possible reasons.

dataset issues could be one of the reasons for the poor performance. If the dataset lacks variability in colors between different classes, then a color-based retrieval system might struggle to differentiate between them. There might be significant overlap in color distributions across different classes, making it hard for the histogram method to distinguish between them effectively. Also, if the dataset is too small or does not have enough representative samples for each class, the system may not be able to learn the distinguishing features properly. Besides,

the choice of color space (RGB) could affect performance. Some color spaces might encapsulate the differences between classes better than it. Also, If the dataset is imbalanced, with some classes having significantly more samples than others, it can skew the training process and affect the system's ability to generalize. For some datasets, color may not be a relevant feature for class discrimination. For instance, in datasets where the shape, texture, or context is more important, color histograms might not provide useful information for classification.

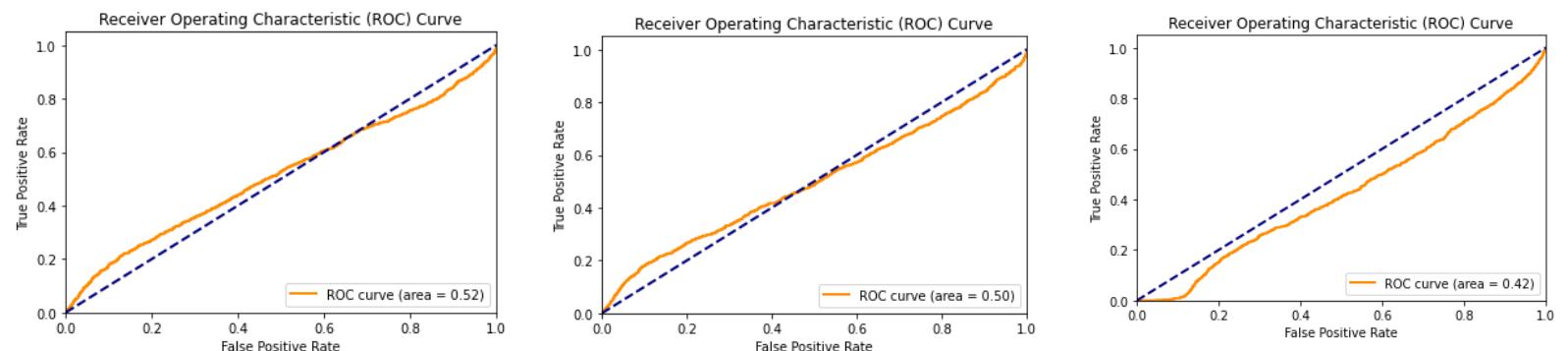
Results Using these thresholds for the 10 queries: (Test 2)

```
quantiles = [ 0.0111, 0.03, 0.05, 0.059, 0.07, 0.09, 0.099, 0.119, 0.15, 0.188, 0.19, 0.2, 0.246, 0.3, 0.333, 0.36, 0.015, 0.02, 0.025, 0.035, 0.04, 0.045, 0.055, 0.06, 0.065, 0.075, 0.08, 0.085, 0.095, 0.1001, 0.105, 0.11, 0.12, 0.125, 0.13, 0.135, 0.14, 0.145, 0.155, 0.16, 0.165, 0.17, 0.175, 0.18, 0.185, 0.195 ]
```

Got the following results:

All images could be found in [RandomQuiriesTest3](#)

For 120 bins, 180 bins, and 256 bins, respectively.



The complete Consol output for this test could be found in **Test3.log**:

```
Time taken for 120 bins: 77.9592 seconds

--- Average Evaluation Results for 120 Bins for All Queries
Threshold: 0.0524
Average Precision: 0.7009
Average Recall: 0.0340
Average F1 Score: 0.0581

Threshold: 0.0632
Average Precision: 0.4569
Average Recall: 0.0650
Average F1 Score: 0.0858

Threshold: 0.0704
Average Precision: 0.3808
Average Recall: 0.0990
Average F1 Score: 0.1055

Threshold: 0.0731
Average Precision: 0.3692
Average Recall: 0.1160
Average F1 Score: 0.1120

Threshold: 0.0762
Average Precision: 0.3604
Average Recall: 0.1290
Average F1 Score: 0.1144

Threshold: 0.0811
Average Precision: 0.3579
Average Recall: 0.1630
Average F1 Score: 0.1271

Threshold: 0.0831
Average Precision: 0.3380
Average Recall: 0.1750
Average F1 Score: 0.1293

Threshold: 0.0873
Average Precision: 0.2915
Average Recall: 0.2030
Average F1 Score: 0.1380

Threshold: 0.0933
Average Precision: 0.2377
Average Recall: 0.2290
Average F1 Score: 0.1349

1180 Threshold: 0.1001
1181 Average Precision: 0.2266
1182 Average Recall: 0.2610
1183 Average F1 Score: 0.1391
1184
1185 Threshold: 0.1004
1186 Average Precision: 0.2231
1187 Average Recall: 0.2630
1188 Average F1 Score: 0.1396
1189
1190 Threshold: 0.1021
1191 Average Precision: 0.2210
1192 Average Recall: 0.2670
1193 Average F1 Score: 0.1392
1194
1195 Threshold: 0.1097
1196 Average Precision: 0.1763
1197 Average Recall: 0.2990
1198 Average F1 Score: 0.1388
1199
1200 Threshold: 0.1185
1201 Average Precision: 0.1426
1202 Average Recall: 0.3440
1203 Average F1 Score: 0.1436
1204
1205 Threshold: 0.1237
1206 Average Precision: 0.1489
1207 Average Recall: 0.3810
1208 Average F1 Score: 0.1556
1209
1210 Threshold: 0.1279
1211 Average Precision: 0.1377
1212 Average Recall: 0.3970
1213 Average F1 Score: 0.1534
1214
1215 Threshold: 0.0554
1216 Average Precision: 0.5766
1217 Average Recall: 0.0370
1218 Average F1 Score: 0.0605
1219
1220 Threshold: 0.0584
1221 Average Precision: 0.5269
1222 Average Recall: 0.0470
1223 Average F1 Score: 0.0713
1224
1225 Threshold: 0.0609
1226 Average Precision: 0.4809
1227 Average Recall: 0.0590
1228 Average F1 Score: 0.0815
1229
1230 Threshold: 0.0652
1231 Average Precision: 0.4540
1232 Average Recall: 0.0750
1233 Average F1 Score: 0.0942
1234
1235 Threshold: 0.0671
1236 Average Precision: 0.4400
1237 Average Recall: 0.0830
1238 Average F1 Score: 0.0993
1239
1240 Threshold: 0.0688
1241 Average Precision: 0.4345
1242 Average Recall: 0.0920
1243 Average F1 Score: 0.1029
1244
1245 Threshold: 0.0720
1246 Average Precision: 0.3747
1247 Average Recall: 0.1090
1248 Average F1 Score: 0.1099
1249
1250 Threshold: 0.0734
1251 Average Precision: 0.3675
1252 Average Recall: 0.1170
1253 Average F1 Score: 0.1119
1254
1255 Threshold: 0.0748
1256 Average Precision: 0.3650
1257 Average Recall: 0.1250
1258 Average F1 Score: 0.1146
1259
1260 Threshold: 0.0775
1261 Average Precision: 0.3578
1262 Average Recall: 0.1340
1263 Average F1 Score: 0.1148
1264
1265 Threshold: 0.0787
1266 Average Precision: 0.3586
1267 Average Recall: 0.1460
1268 Average F1 Score: 0.1199
1269
1270 Threshold: 0.0799
1271 Average Precision: 0.3584
1272 Average Recall: 0.1540
1273 Average F1 Score: 0.1227
1274
1275 Threshold: 0.0822
1276 Average Precision: 0.3365
1277 Average Recall: 0.1680
1278 Average F1 Score: 0.1263
1279
1280 Threshold: 0.0834
1281 Average Precision: 0.3303
1282 Average Recall: 0.1770
1283 Average F1 Score: 0.1304
1284
```

```
319 Threshold: 0.0924
320 Average Precision: 0.2537
321 Average Recall: 0.2270
322 Average F1 Score: 0.1365
323
324 Threshold: 0.0942
325 Average Precision: 0.2288
326 Average Recall: 0.2350
327 Average F1 Score: 0.1360
328
329 Threshold: 0.0952
330 Average Precision: 0.2318
331 Average Recall: 0.2370
332 Average F1 Score: 0.1356
333
334 Threshold: 0.0961
335 Average Precision: 0.2348
336 Average Recall: 0.2410
337 Average F1 Score: 0.1362
338
339 Threshold: 0.0969
340 Average Precision: 0.2297
341 Average Recall: 0.2440
342 Average F1 Score: 0.1353
343
344 Threshold: 0.0978
345 Average Precision: 0.2285
346 Average Recall: 0.2490
347 Average F1 Score: 0.1364
348
349 Threshold: 0.0987
350 Average Precision: 0.2282
351 Average Recall: 0.2550
352 Average F1 Score: 0.1395
353
354 Threshold: 0.0995
355 Average Precision: 0.2246
356 Average Recall: 0.2570
357 Average F1 Score: 0.1379
358
359 Threshold: 0.1013
360 Average Precision: 0.2236
361 Average Recall: 0.2650
362 Average F1 Score: 0.1404
363
364
365 --- Overall Average Evaluation Results for 120 Bins ---
366 Overall Average Precision: 0.3189
367 Overall Average Recall: 0.1870
368 Overall Average F1 Score: 0.1231
369
```

```
Time taken for 180 bins: 103.2280 seconds
```

```
4
5
6 --- Overall Average Evaluation Results for 180 Bins ---
7 Overall Average Precision: 0.3450
8 Overall Average Recall: 0.1753
9 Overall Average F1 Score: 0.1144
0
```

```
1 Time taken for 256 bins: 127.4501 seconds
2
```

```
5
6 --- Overall Average Evaluation Results for 256 Bins ---
7 Overall Average Precision: 0.3685
8 Overall Average Recall: 0.2018
9 Overall Average F1 Score: 0.1403
0
```

Discussing the results above:

- **ROC and AUC**

For 120 bins, the AUC is 0.52, which is slightly better than random guessing but still indicates poor discriminative ability.

For 180 bins, the AUC drops to 0.50, which suggests no better than random chance.

For 256 bins, the AUC is 0.42, which is even worse than random chance, indicating that the model is performing poorly at this bin size.

- **Evaluation Metrics:**

With 120 bins, the overall average precision is around 0.3189, recall is 0.1870, and F1 score is 0.1231.

With 180 bins, the overall average precision decreases slightly to 0.3450, recall is 0.1753, and F1 score is 0.1144.

With 256 bins, there's an increase in overall average precision to 0.3685, recall also increases to 0.2018, and F1 score goes up to 0.1403.

when looking at the overall average precision, recall, and F1 score, there is a slight improvement as the number of bins increases to 256. This could indicate that for this particular dataset and retrieval task, having more bins allows for a more detailed color histogram, which may capture more nuances in the color distributions of the images.

However, an AUC below 0.5 for the 256 bins model suggests that despite the apparent improvements in precision and recall, the model may be making decisions that are systematically incorrect, or the method used to construct the ROC curve may have some issues. The time taken for computations increases with the number of bins, which is expected.

While the precision and recall seem to improve with more bins, the ROC AUC values suggest that the system is not effectively discriminating between the classes, especially at

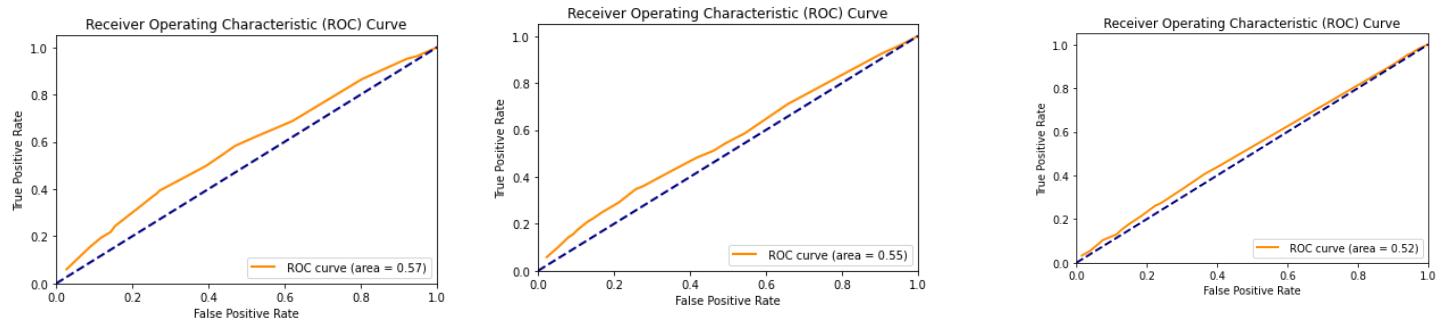
256 bins. These conflicting signals imply that while the system might be retrieving more relevant items, it is also making many incorrect classifications. This scenario may warrant further investigation into the dataset, and sure the feature extraction method itself, since it is very basic simple model.

Results Using these thresholds for the 10 queries: (Test 3)

For the following thresholds:

quantiles = [0.0111,0.00111 , 0.03,0.05,0.059,0.07, 0.09,0.099,,119,0.15,0.188, 0.19,0.2,0.246,0.3,0.333,0.36, 0.4,0.5,0.8,.88,.85,.66,,9,0.77,0.67,.73,.95, .999,1]

AUC values: Decreasing from 0.57 for 120 bins , to .55 for 180 bins, to 0.52 for 256 bins, as shown below.



The overall average measures:

```

24
25
26 --- Overall Average Evaluation Results for 120 Bins ---
27 Overall Average Precision: 0.2379
28 Overall Average Recall: 0.4960
29 Overall Average F1 Score: 0.1441
30

```

```

5
6 --- Overall Average Evaluation Results for 180 Bins ---
7 Overall Average Precision: 0.2100
8 Overall Average Recall: 0.5389
9 Overall Average F1 Score: 0.1648
10

```

```

85
86 --- Overall Average Evaluation Results for 256 Bins ---
87 Overall Average Precision: 0.2617
88 Overall Average Recall: 0.5390
89 Overall Average F1 Score: 0.1631
90

```

AUC of 0.57 for Bins 120: An AUC value of 0.57 is generally considered slightly better than random guessing. However, it's not a very strong performance. In practical terms, this indicates that the model's ability to distinguish between classes is only marginally better than chance.

AUC of 0.55 for Bins 180: Here, the data is divided into 180 bins, and the model's AUC drops slightly to 0.55. This decrease, albeit small, suggests that increasing the bins does not improve the model's performance; in fact, it slightly worsens it.

AUC of 0.52 for Bins 256: With the data divided into 256 bins, the AUC further decreases to 0.52. This is a very modest improvement over random guessing and suggests that the model is not effectively learning the distinctions between classes at this level of data segmentation.

Results Using these thresholds for the 10 queries: (Test 4)

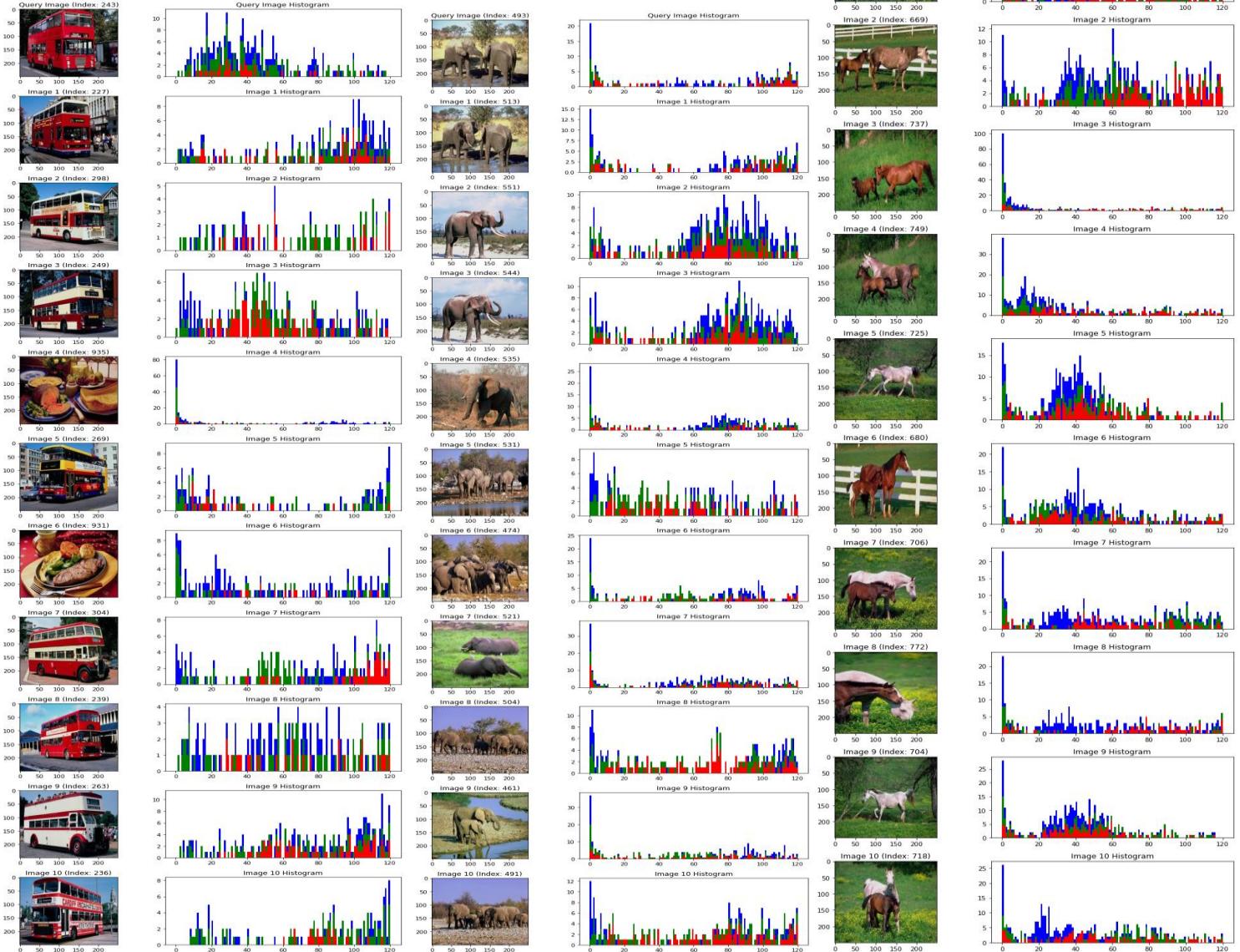
I would mention that at some point, for another test, Test4, I found that AUC increase with increasing bins but still very low increasing (From 0.39 , to 0.41, to 0.51), Which indicate that increasing number of bins has a good impact in performance.

You can find the run results of this test on [Test4.log](#)

Testing for specific 10 queries:

Here 10 specific queries chosen , each was tested on specific thresholds values suits the ranked distances to it .

Some samples images: I couldn't add more images, since it will take the full space.



Here I am showing overall results of `test5.log`, in `test5.log` and `Task1and2-specific query.log` you can find the testing results with different thresholds for these specific 10 queries, both gave near results to each other's.

```
Time taken for 120 bins: 81.5891 seconds
```

```
0
1 --- Overall Average Evaluation Results for 120 Bins ---
2 Overall Average Precision: 0.2537
3 Overall Average Recall: 0.4427
4 Overall Average F1 Score: 0.2372
5
```

```
1
2 Time taken for 180 bins: 101.6776 seconds
3
```

```
--- Overall Average Evaluation Results for 180 Bins ---
Overall Average Precision: 0.4138
Overall Average Recall: 0.1114
Overall Average F1 Score: 0.1120
```

```
Time taken for 256 bins: 128.7877 seconds
```

```
4031 --- Overall Average Evaluation Results for 256 Bins ---
4032 Overall Average Precision: 0.4920
4033 Overall Average Recall: 0.0868
4034 Overall Average F1 Score: 0.0882
4035
```

Here, because I chosen thresholds per query , I got better results for 120 bins as shown , but as increasing number of bins , precision increase but recall decrease , and overall performance decrease (f1 score decreased).

Some good results on some thresholds:

```
Threshold: 0.1625
avg-Precision: 0.20
avg-Recall: 0.27
avg-F1 Score: 0.23
Execution Time: 0.4638 seconds

Threshold: 0.1828
avg-Precision: 0.17
avg-Recall: 0.30
avg-F1 Score: 0.22
Execution Time: 0.4638 seconds

Threshold: 0.2051
avg-Precision: 0.17
avg-Recall: 0.39
avg-F1 Score: 0.24
Execution Time: 0.4648 seconds

Threshold: 0.2063
avg-Precision: 0.17
avg-Recall: 0.39
avg-F1 Score: 0.23
Execution Time: 0.4648 seconds
```

```
Threshold: 0.1156
avg-Precision: 0.26
avg-Recall: 0.21
avg-F1 Score: 0.23
Execution Time: 0.5651 seconds

Threshold: 0.1252
avg-Precision: 0.22
avg-Recall: 0.23
avg-F1 Score: 0.23
Execution Time: 0.5651 seconds

Threshold: 0.1412
avg-Precision: 0.19
avg-Recall: 0.28
avg-F1 Score: 0.23
Execution Time: 0.5651 seconds

Threshold: 0.1479
avg-Precision: 0.19
avg-Recall: 0.31
avg-F1 Score: 0.23
Execution Time: 0.5651 seconds

Threshold: 0.2063
avg-Precision: 0.16
avg-Recall: 0.37
avg-F1 Score: 0.23
Execution Time: 0.5651 seconds

Threshold: 0.1828
avg-Precision: 0.15
avg-Recall: 0.40
avg-F1 Score: 0.22
Execution Time: 0.5651 seconds

Threshold: 0.2051
avg-Precision: 0.14
avg-Recall: 0.41
avg-F1 Score: 0.21
Execution Time: 0.5651 seconds

Threshold: 0.2063
avg-Precision: 0.14
avg-Recall: 0.41
avg-F1 Score: 0.21
Execution Time: 0.5666 seconds

Threshold: 0.2117
avg-Precision: 0.14
avg-Recall: 0.42
avg-F1 Score: 0.21
Execution Time: 0.5666 seconds
```

Conclusion for this part:

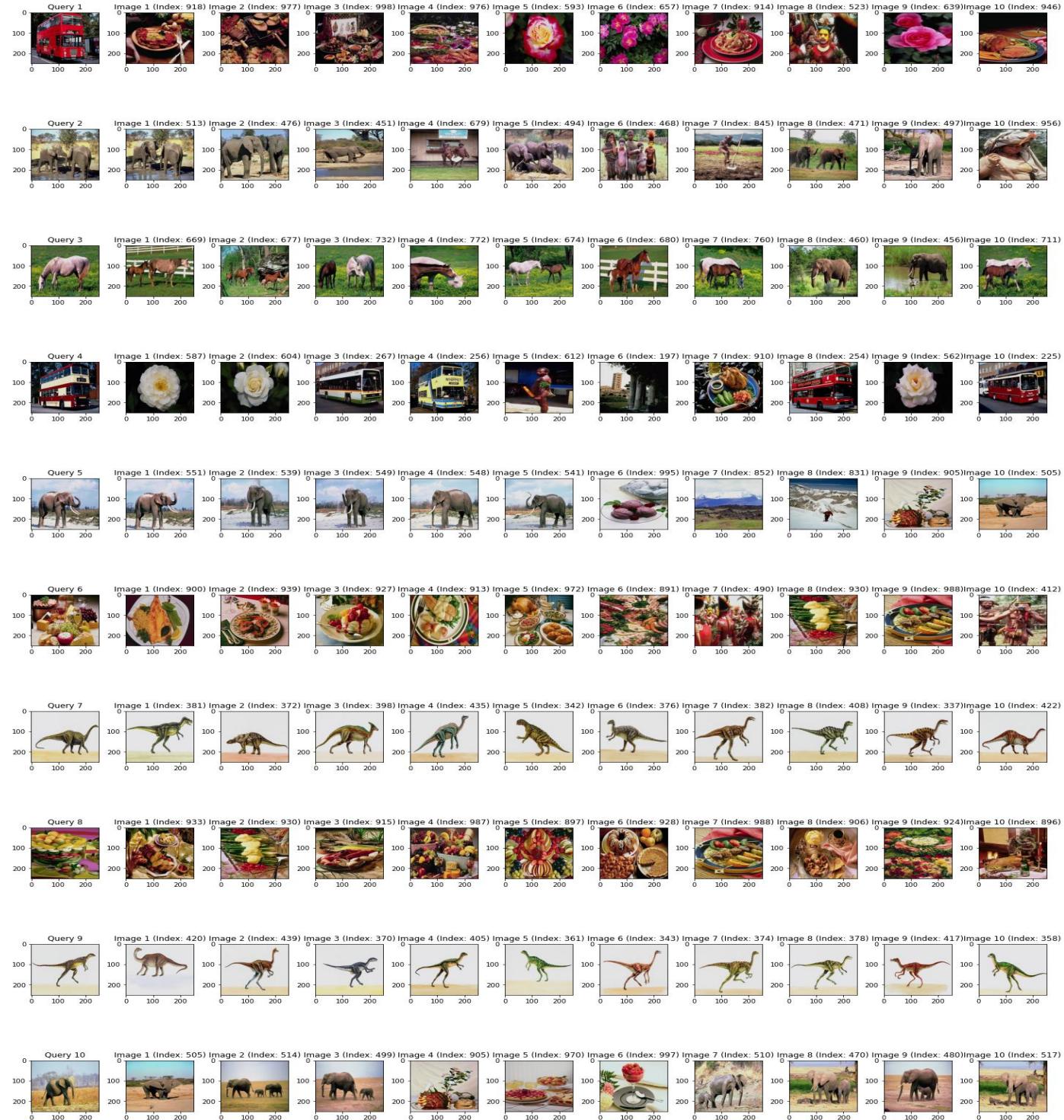
In good datasets higher bin counts may lead to better discrimination between relevant and irrelevant images for the task and potentially improve the AUC and ROC performance. In this work, the evaluation of the model's performance with varying bin sizes and thresholds reveals nuanced findings. The ROC and AUC analysis indicate that, overall, the model struggles to effectively discriminate between classes. While the precision, recall, and F1 score show a slight improvement as the number of bins increases to 256, the conflicting signals from the AUC values suggest potential issues in constructing the ROC curve. The reasons of this decrease in AUC are explained before during discussion results.

The results from different tests, such as Test 4, demonstrate varying impacts on performance with an increasing number of bins. In some cases, the AUC increases marginally, indicating a potential positive influence on the model's discriminatory ability. However, this improvement is modest and may not be consistent across different datasets or retrieval tasks.

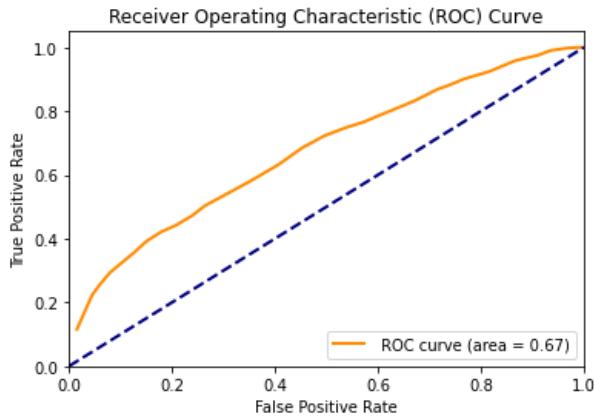
The choice of thresholds per query in Test 5 leads to improved results with 120 bins, but as the number of bins increases, precision improves while recall decreases, resulting in an overall decrease in performance (F1 score). This suggests a delicate balance between granularity and computation time, with 120 bins providing a medium level of granularity, 180 bins offering higher resolution, and 256 bins providing the highest level of granularity.

3.2 Task 3: CBIR System Using Color Moments

Subsection 3.2.1: Experiment with Equal Weights



ROC and AUC:



Some measures result at many thresholds for different queries:

```
At Threshold: 0.1818, Precision: 0.17, Recall: 0.04, F1 Score: 0.05
Average values for all thresholds for Query Index: 959
avg-Precision: 0.29
avg-Recall: 0.26
avg-F1 Score: 0.23
Query Index: 440
```

```
At Threshold: 0.1818, Precision: 0.40, Recall: 0.41, F1 Score: 0.40
Average values for all thresholds for Query Index: 440
avg-Precision: 0.31
avg-Recall: 0.28
avg-F1 Score: 0.25
Query Index: 979
```

```
Average values for all thresholds for Query Index: 979
avg-Precision: 0.31
avg-Recall: 0.29
avg-F1 Score: 0.26
Query Index: 373
```

```
At Threshold: 0.1818, Precision: 0.56, Recall: 0.53, F1 Score: 0.53
Average values for all thresholds for Query Index: 373
avg-Precision: 0.34
avg-Recall: 0.30
avg-F1 Score: 0.29
Query Index: 472
```

```
At Threshold: 0.1818, Precision: 0.70, Recall: 0.70, F1 Score: 0.70
Average values for all thresholds for Query Index: 472
avg-Precision: 0.34
avg-Recall: 0.28
avg-F1 Score: 0.27
```

```

5 Threshold: 0.1818
3 Average Precision: 0.2632
9 Average Recall: 0.2000
3 Average F1 Score: 0.2273
1
2 Threshold: 0.1907
3 Average Precision: 0.2796
4 Average Recall: 0.2600
5 Average F1 Score: 0.2694
6
7 Threshold: 0.2029
3 Average Precision: 0.2881
9 Average Recall: 0.3400
3 Average F1 Score: 0.3119
1
2 Threshold: 0.2155
3 Average Precision: 0.2676
4 Average Recall: 0.3800
5 Average F1 Score: 0.3140
6
7 Threshold: 0.2060
3 Average Precision: 0.2960
9 Average Recall: 0.3700
3 Average F1 Score: 0.3289
1
2 Threshold: 0.2187
3 Average Precision: 0.2550
4 Average Recall: 0.3800
5 Average F1 Score: 0.3052
6

```

```

--- Overall Average Evaluation Results ---
Overall Average Precision: 0.2953
Overall Average Recall: 0.1414
Overall Average F1 Score: 0.1604

```

```
Total Execution Time: 12.0261 seconds
```

Full run results could be found on [Task3part1.log](#), the results above show better f1scores compared with the ones we got on the task2 when we used the same queries and same thresholds.

They got a ROC AUC score of 0.67. An AUC of 0.67 suggests that the model has moderate discriminatory ability. As adjusting the threshold can impact precision, recall, and the F1 score, the images above show good thresholds since they give good trade off.

Subsection 3.2.2: Experiment with Different Weights

1. Emphasizing Mean (Average Color)

- Weights: [3, 1, 1] * 3 (3 times the weight for the mean of each color channel)

This implies that the mean of each color channel is given three times the weight compared to the other color moments (variance and skewness). Such a weighting scheme suggests a hypothesis that the average color of an image is more significant for retrieval purposes than its color distribution and asymmetry.

Auc measure was 0.66



```
--- Overall Average Evaluation Results ---  
Overall Average Precision: 0.2898  
Overall Average Recall: 0.0953  
Overall Average F1 Score: 0.1274  
  
Total Execution Time: 11.0965 seconds  
Completed all experiments
```

```
946  
947     Threshold: 0.1799  
948     Average Precision: 0.2500  
949     Average Recall: 0.2300  
950     Average F1 Score: 0.2396  
951  
952     Threshold: 0.1929  
953     Average Precision: 0.2759  
954     Average Recall: 0.3200  
955     Average F1 Score: 0.2963  
956  
957     Threshold: 0.1830  
958     Average Precision: 0.2500  
959     Average Recall: 0.2400  
960     Average F1 Score: 0.2449  
961  
962     Threshold: 0.1964  
963     Average Precision: 0.2720  
964     Average Recall: 0.3400  
965     Average F1 Score: 0.3022  
966
```

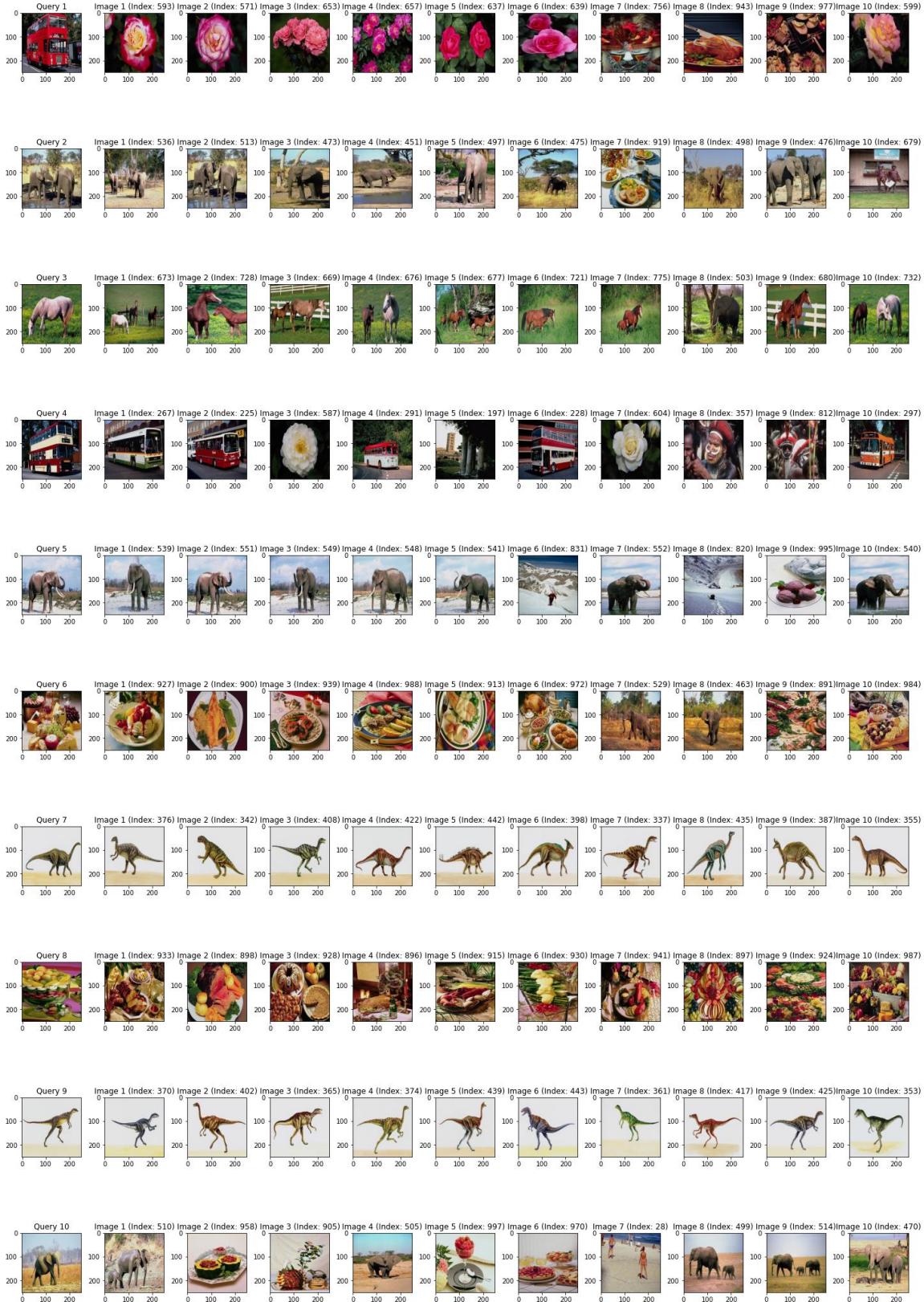
```
At threshold: 0.1964, precision: 0.30, recall: 0.37, f1 score: 0.34  
Average values for all thresholds for Query Index: 373  
avg-Precision: 0.34  
avg-Recall: 0.32  
avg-F1 Score: 0.28  
Query Index: 472
```

```
Average values for all thresholds for Query Index: 472  
avg-Precision: 0.33  
avg-Recall: 0.30  
avg-F1 Score: 0.26
```

```
9 Average values for all thresholds for Query Index: 979  
9 avg-Precision: 0.31  
L avg-Recall: 0.30  
9 avg-F1 Score: 0.25
```

- o Increasing the weight of mean more : weights = [10, 1, 1]*3

Here, the mean is given even more prominence, suggesting that the average color is thought to be far more critical compared to other weights.



Auc measure was 0.66, and the other measures as shown :

```
68 --- Overall Average Evaluation Results ---
69 Overall Average Precision: 0.2806
70 Overall Average Recall: 0.0910
71 Overall Average F1 Score: 0.1219
72
73 Total Execution Time: 10.0233 seconds
74 Completed all experiments
```

```
Threshold: 0.1771
Average Precision: 0.2500
Average Recall: 0.2200
Average F1 Score: 0.2340

Threshold: 0.1902
Average Precision: 0.2696
Average Recall: 0.3100
Average F1 Score: 0.2884

Threshold: 0.1803
Average Precision: 0.2581
Average Recall: 0.2400
Average F1 Score: 0.2487

Threshold: 0.1938
Average Precision: 0.2705
Average Recall: 0.3300
Average F1 Score: 0.2973
```

```
At threshold: 0.1938, precision: 0.37, recall: 0.37, f1 score: 0.33
Average values for all thresholds for Query Index: 373
avg-Precision: 0.33
avg-Recall: 0.31
avg-F1 Score: 0.28
```

```
At threshold: 0.1938, precision: 0.37, recall: 0.37, f1 score: 0.33
Average values for all thresholds for Query Index: 979
avg-Precision: 0.31
avg-Recall: 0.30
avg-F1 Score: 0.25
```

Comparing Results:

By increasing the weight of the mean color, one might anticipate an improvement in the precision, recall, and F1 score if the mean color is indeed a more critical factor for the retrieval process, as the system would become more attuned to the average color differences between images. However, overemphasizing a single aspect can lead to the

neglect of other significant features like texture or shape, which are partly reflected by the variance and skewness of colors. This could result in suboptimal performance. In this case, there is no decline in performance, but rather a lack of enhancement; the outcomes remain consistent. For instance, with equal weights, the AUC was 0.67, and upon emphasizing the mean, it stands at 0.66. As for the F1 score, which was initially 0.1604, it slightly decreased to 0.1274 after tripling the weight of the mean, and with a tenfold increase in the mean's weight, there was virtually no change, dropping to just 0.1219. Therefore, in this scenario, giving additional weight to the mean color alone does not yield a significant benefit.

The AUC (Area Under the Curve) for both the $[3, 1, 1] * 3$ and $[10, 1, 1] * 3$ weighted scenarios is 0.66, which is the same. This indicates that there is no significant difference in the discriminative power of the classifier when changing the weights from a 3-times emphasis to a 10-times emphasis on the mean color. This could be interpreted in several ways:

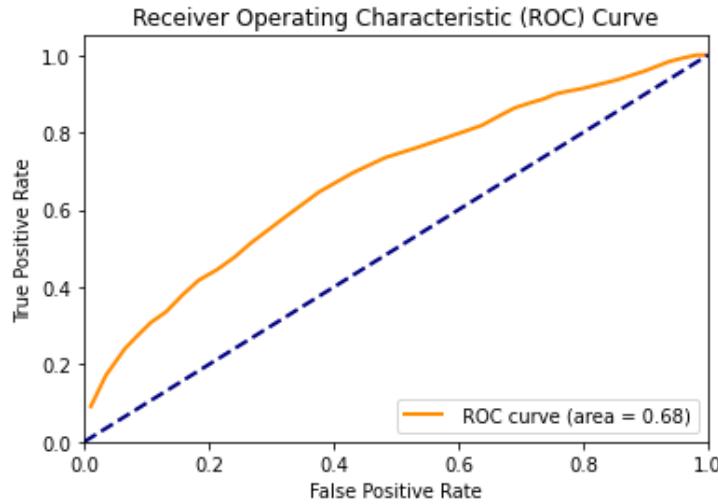
- The mean color is not the only critical feature: It's possible that the mean color alone isn't sufficiently discriminative, and by focusing too much on it, the system is not improving.
- Diminishing Returns: There might be a point beyond which giving more weight to the mean color doesn't contribute to a better separation of classes because the system needs a balanced consideration of all color moments.
- Dataset Characteristics: The specific dataset used for evaluation might not have significant variations in mean color across different classes, making it a less effective feature for discrimination.
- System Limitations: The CBIR system itself might have limitations, such as in the way it uses the weighted features, that prevent it from benefiting from the changed weights.

2. Emphasizing Contrast (Standard Deviation)

- Weights: [1, 3, 1] * 3 (3 times the weight for the standard deviation of each color channel)



Auc measure was **0.68**, and the other measures as shown:



```
Threshold: 0.2248
Average Precision: 0.3171
Average Recall: 0.3900
Average F1 Score: 0.3498

Threshold: 0.2325
Average Precision: 0.3134
Average Recall: 0.4200
Average F1 Score: 0.3590

Threshold: 0.2403
Average Precision: 0.3025
Average Recall: 0.4900
Average F1 Score: 0.3740

Threshold: 0.2479
Average Precision: 0.3027
Average Recall: 0.5600
Average F1 Score: 0.3930

Threshold: 0.2595
Average Precision: 0.2791
Average Recall: 0.6000
Average F1 Score: 0.3810

Threshold: 0.2750
Average Precision: 0.2679
Average Recall: 0.7100
Average F1 Score: 0.3890

Threshold: 0.2908
Average Precision: 0.2542
Average Recall: 0.7600
Average F1 Score: 0.3810

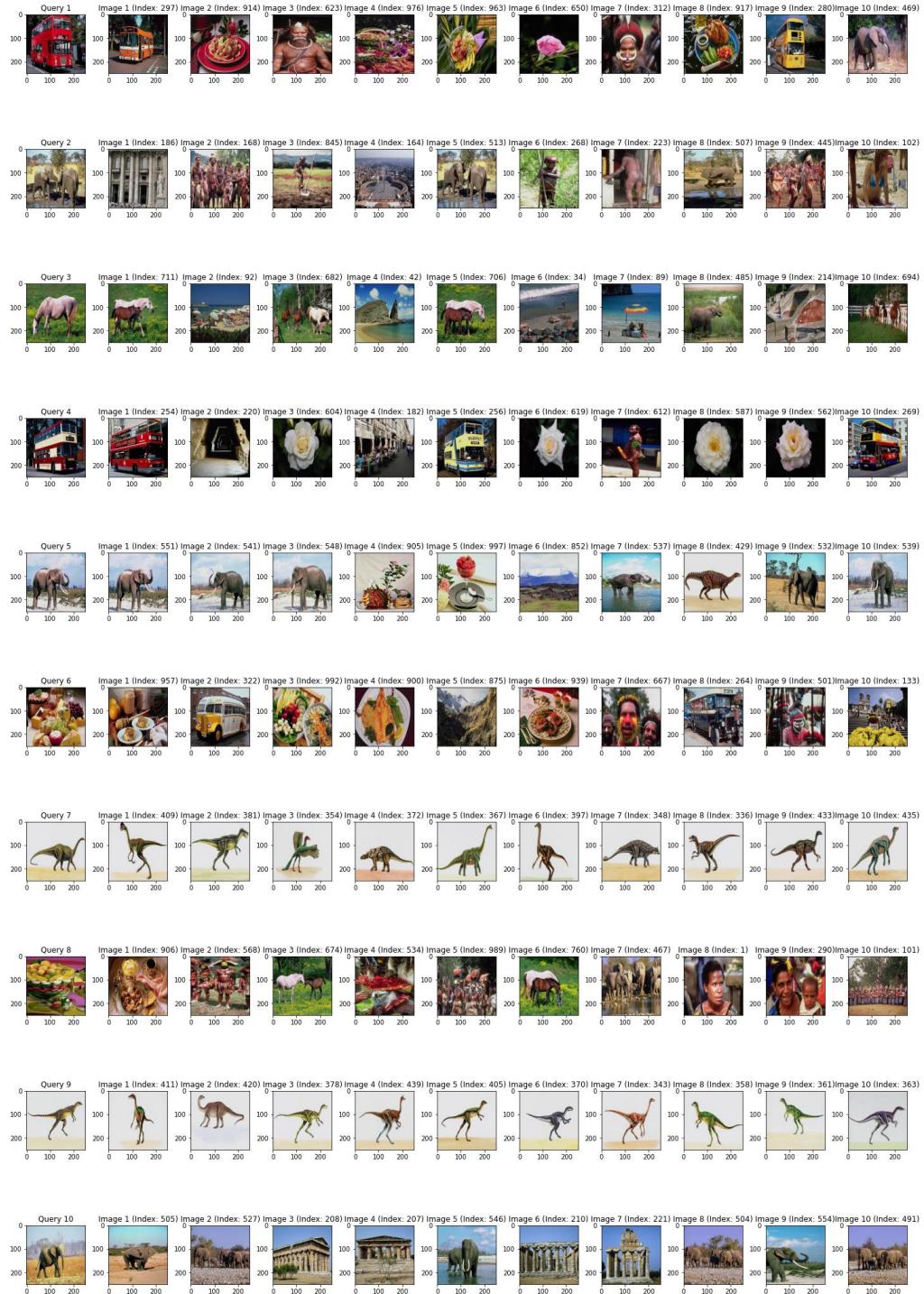
Threshold: 0.2789
Average Precision: 0.2667
Average Recall: 0.7200
Average F1 Score: 0.3892

Threshold: 0.2948
Average Precision: 0.2516
Average Recall: 0.7800
Average F1 Score: 0.3805

--- Overall Average Evaluation Results ---
Overall Average Precision: 0.3055
Overall Average Recall: 0.2102
Overall Average F1 Score: 0.1877

Total Execution Time: 10.3223 seconds
```

- Weights: $[1, 10, 1] * 3$ (3 times the weight for the standard deviation of each color channel), f1 score increased from the value when had equal weight



```

--- Overall Average Evaluation Results ---
Overall Average Precision: 0.2190
Overall Average Recall: 0.3058
Overall Average F1 Score: 0.2183

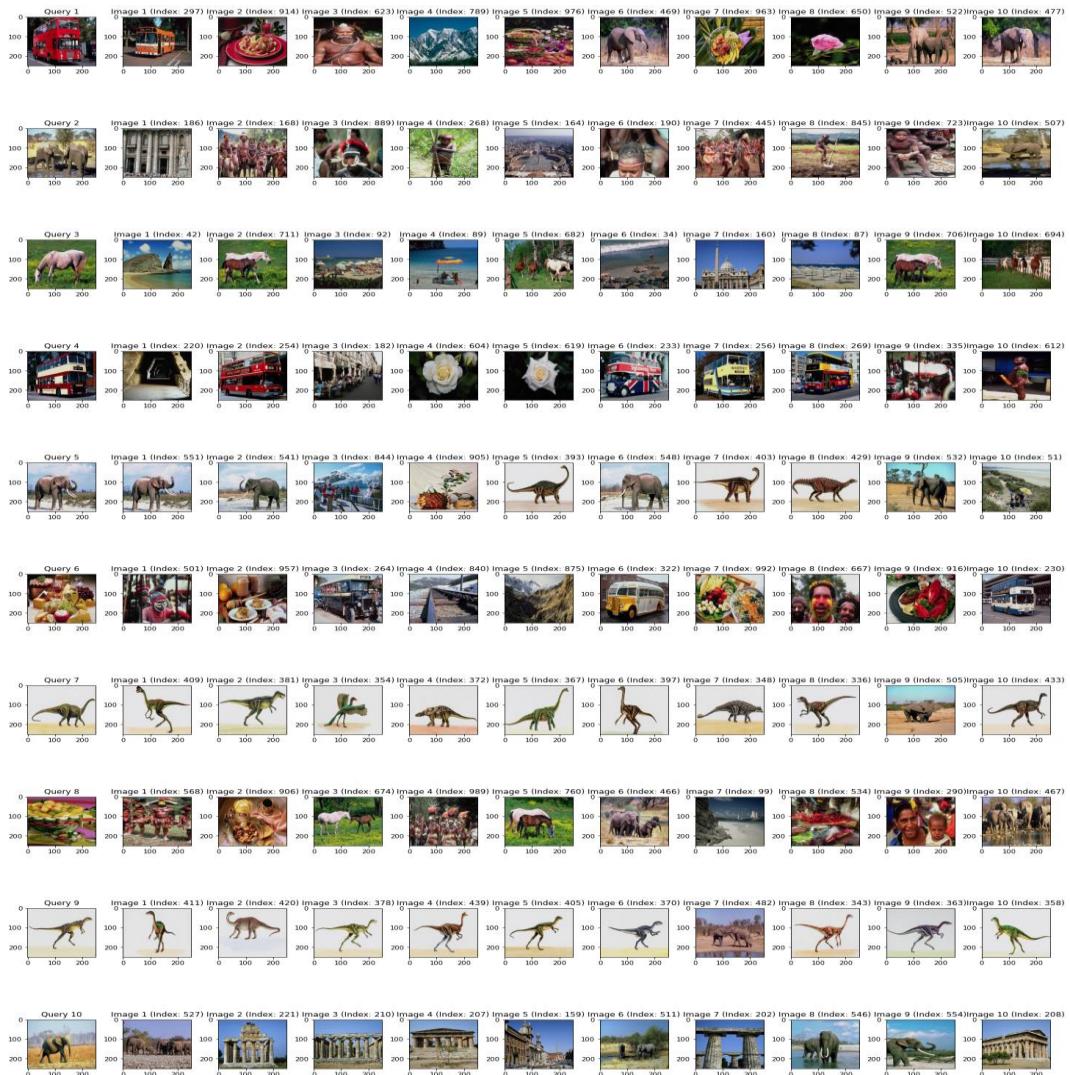
Total Execution Time: 10.3502 seconds
Completed all experiments.

```

AUC was 0.64, f1 score increased more.

- Weights: [1, 20, 1] * 3 (3 times the weight for the standard deviation of each color channel)

AUC was 0.62, f1 score remained as for weight 10 .



```
--- Overall Average Evaluation Results ---  
Overall Average Precision: 0.1916  
Overall Average Recall: 0.2965  
Overall Average F1 Score: 0.2073
```

Comparing Results:

When the Content-Based Image Retrieval (CBIR) system was configured with equal weights across all color channels, the resulting Area Under the Curve (AUC) was 0.67, with an F1 score of 0.1604. Adjusting the system to emphasize contrast by increasing the weight of the standard deviation of each color channel to three times (with weights [1, 3, 1] * 3) resulted in a slight improvement in the AUC to 0.68. This suggests that contrast may play a more critical role in the retrieval process than previously assumed.

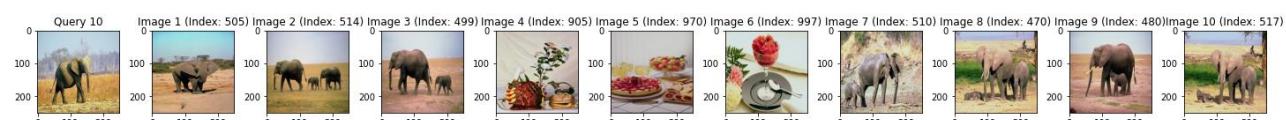
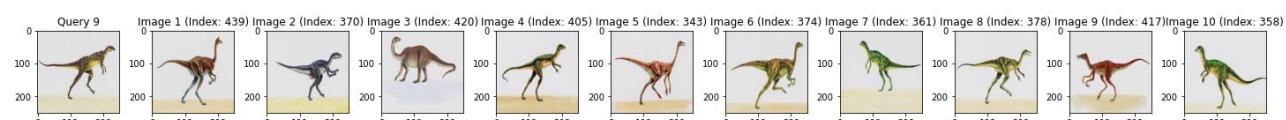
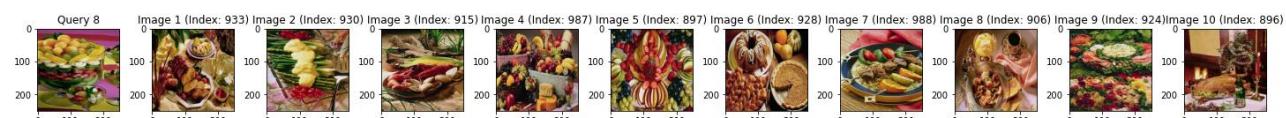
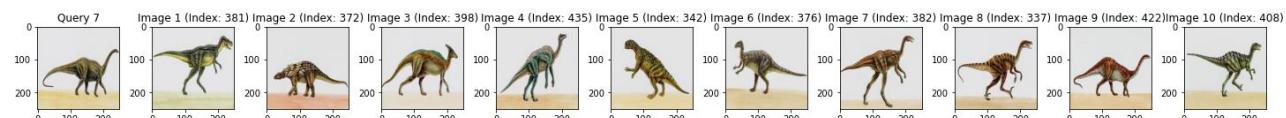
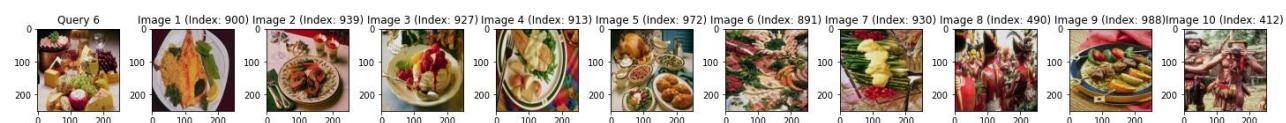
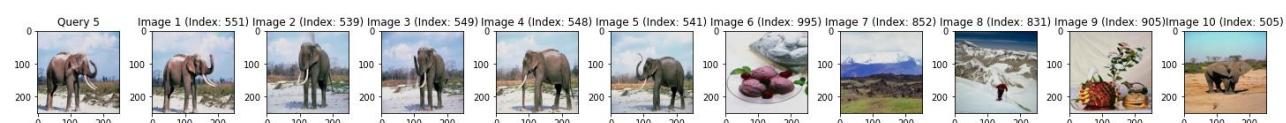
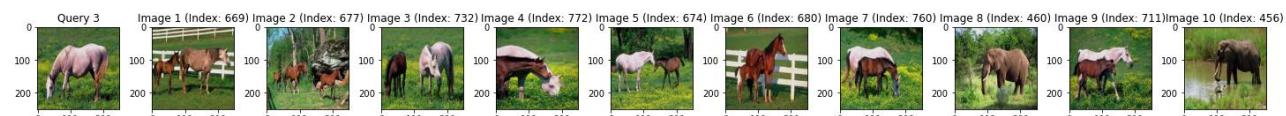
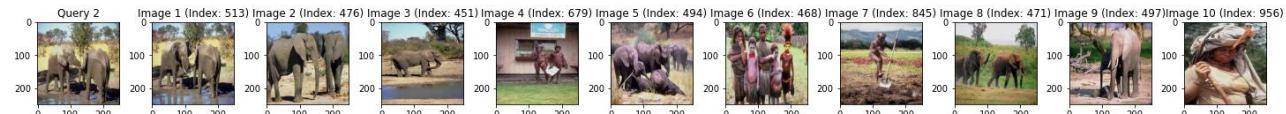
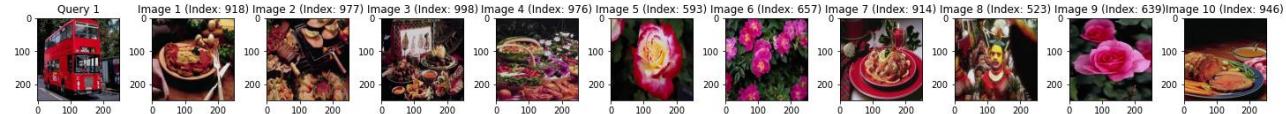
Further emphasizing contrast by increasing the weight to ten times the other factors (with weights [1, 10, 1] * 3) led to a decrease in AUC to 0.64. However, it's notable that the F1 score increased compared to the scenario with equal weights, indicating that while the model's overall discriminative power as measured by AUC decreased, its harmonic balance of precision and recall improved.

Pushing the emphasis on contrast even further to twenty times the other factors (with weights [1, 20, 1] * 3) saw the AUC decrease further to 0.62. Despite this continued decline in AUC, the F1 score remained at the level it was when the weight was ten, indicating a plateau in the benefits gained from further increasing the weight of contrast in terms of F1 score.

3. Emphasizing Color Symmetry (Skewness)

- Weights: [1, 1, 3] * 3 (3 times the weight for skewness of each color channel)

Auc was 0.67, f1 score 0.146



```

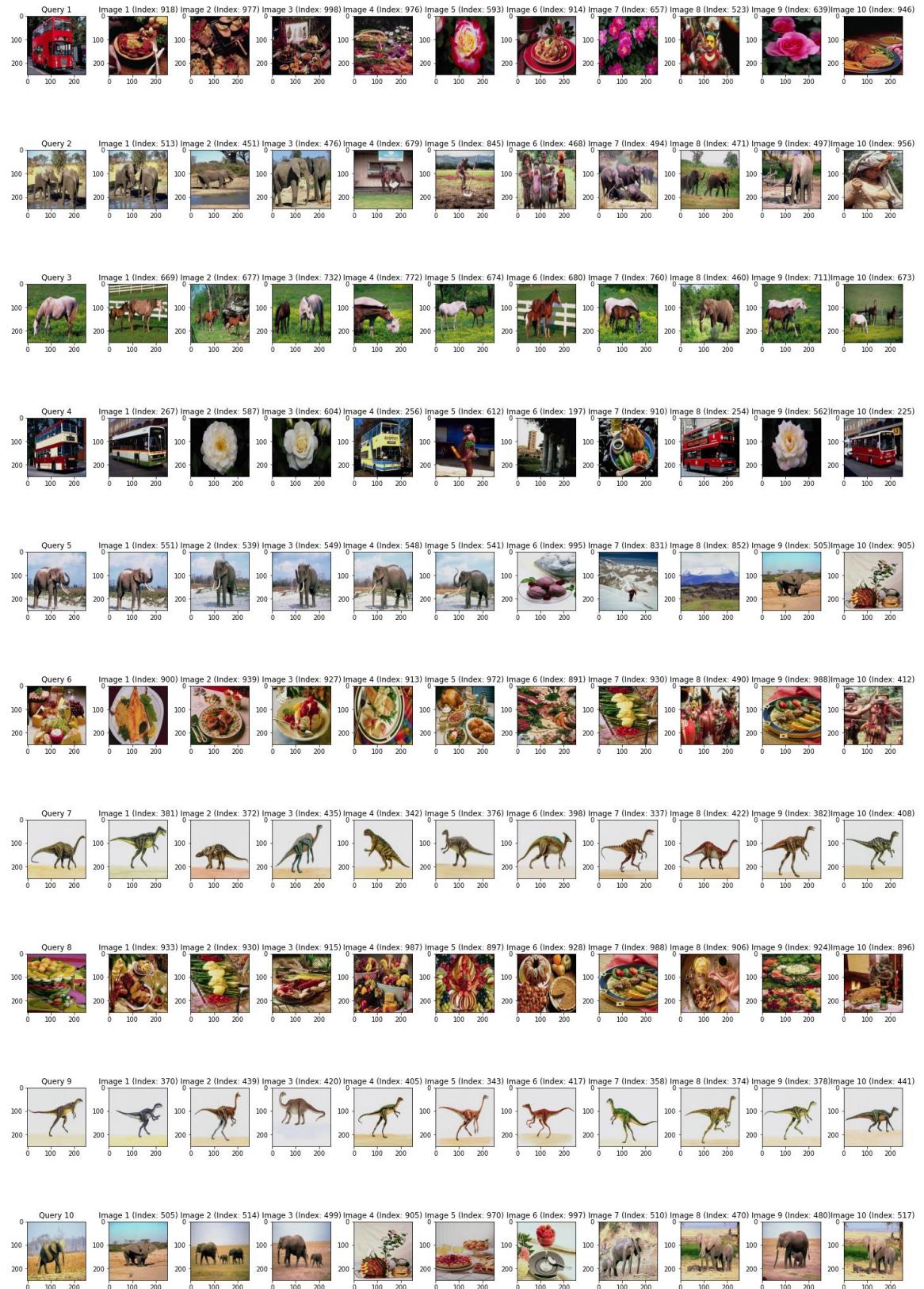
67 --- Overall Average Evaluation Results ---
68 Overall Average Precision: 0.2997
69 Overall Average Recall: 0.1157
70 Overall Average F1 Score: 0.1460
71
72
73 Total Execution Time: 10.7500 seconds
74
```

```

+1
42 Threshold: 0.1903
43 Average Precision: 0.2717
44 Average Recall: 0.2500
45 Average F1 Score: 0.2604
46
47 Threshold: 0.2025
48 Average Precision: 0.2881
49 Average Recall: 0.3400
50 Average F1 Score: 0.3119
51
52 Threshold: 0.2151
53 Average Precision: 0.2676
54 Average Recall: 0.3800
55 Average F1 Score: 0.3140
56
57 Threshold: 0.2056
58 Average Precision: 0.2927
59 Average Recall: 0.3600
60 Average F1 Score: 0.3229
61
62 Threshold: 0.2183
63 Average Precision: 0.2550
64 Average Recall: 0.3800
65 Average F1 Score: 0.3052
66
```

- Weights: [1, 1, 6] * 3 (3 times the weight for skewness of each color channel)

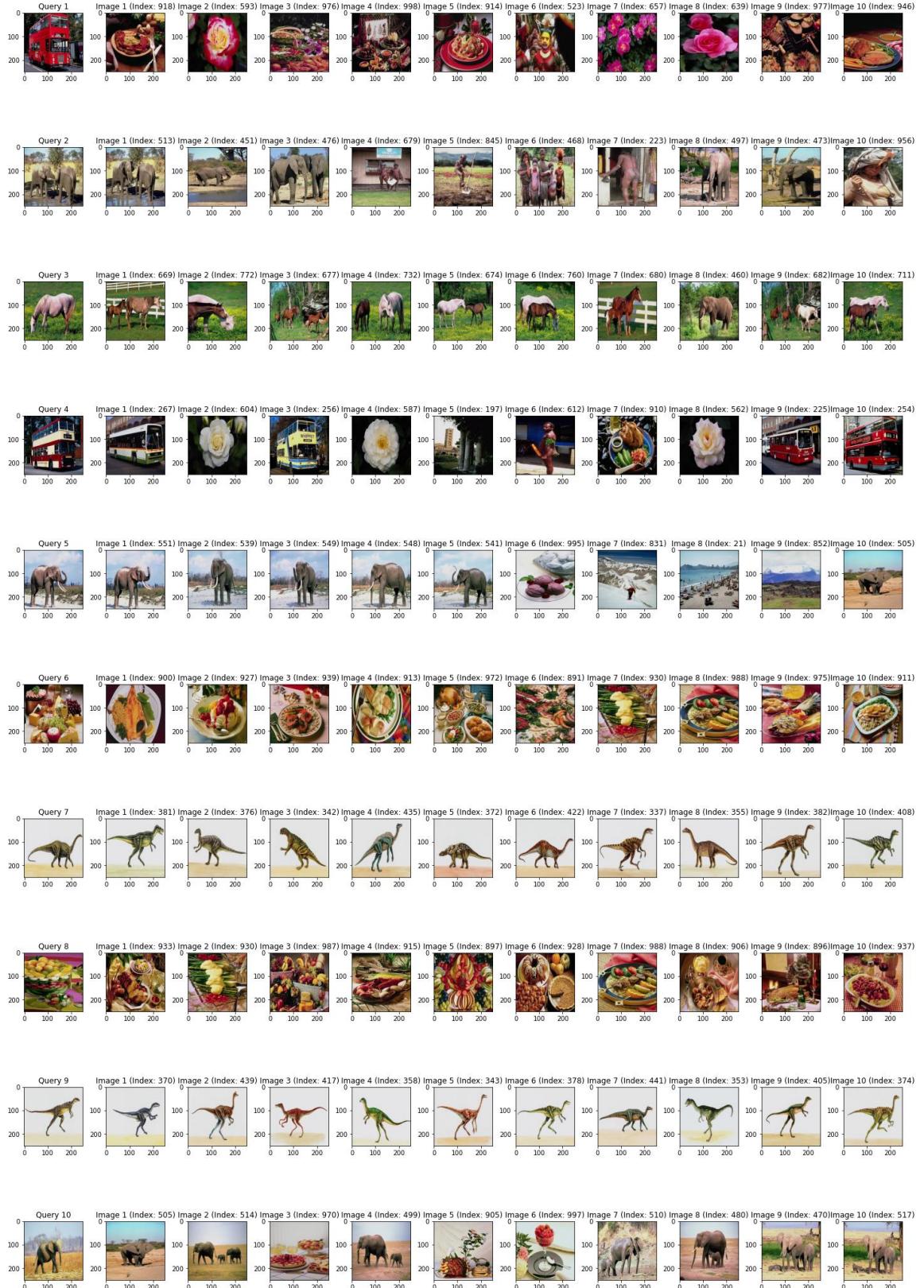
AUC was 0.67, f1 score 0.1463



```
1
2     Threshold: 0.1891
3     Average Precision: 0.2717
4     Average Recall: 0.2500
5     Average F1 Score: 0.2604
6
7     Threshold: 0.2011
8     Average Precision: 0.2906
9     Average Recall: 0.3400
10    Average F1 Score: 0.3134
11
12    Threshold: 0.2136
13    Average Precision: 0.2676
14    Average Recall: 0.3800
15    Average F1 Score: 0.3140
16
17    Threshold: 0.2042
18    Average Precision: 0.2951
19    Average Recall: 0.3600
20    Average F1 Score: 0.3243
21
22    Threshold: 0.2168
23    Average Precision: 0.2603
24    Average Recall: 0.3800
25    Average F1 Score: 0.3089
26
27
28    --- Overall Average Evaluation Results ---
29    Overall Average Precision: 0.3023
30    Overall Average Recall: 0.1157
31    Overall Average F1 Score: 0.1463
32
33    Total Execution Time: 10.1494 seconds
34    Completed all experiments
```

- Weights: [1, 1, 15] * 3 (3 times the weight for skewness of each color channel)

AUC was 0.67 , f1 score 0.1381



```

Threshold: 0.1778
Average Precision: 0.2697
Average Recall: 0.2400
Average F1 Score: 0.2540

Threshold: 0.1889
Average Precision: 0.3019
Average Recall: 0.3200
Average F1 Score: 0.3107

Threshold: 0.2006
Average Precision: 0.2741
Average Recall: 0.3700
Average F1 Score: 0.3149

Threshold: 0.1918
Average Precision: 0.2957
Average Recall: 0.3400
Average F1 Score: 0.3163

Threshold: 0.2035
Average Precision: 0.2734
Average Recall: 0.3800
Average F1 Score: 0.3180

--- Overall Average Evaluation Results ---
Overall Average Precision: 0.3033
Overall Average Recall: 0.1080
Overall Average F1 Score: 0.1381

Total Execution Time: 9.8142 seconds
Completed all experiments

```

Comparing Results:

When the Content-Based Image Retrieval (CBIR) system utilized equal weights for all color channels, it achieved an Area Under the Curve (AUC) of 0.67 and an F1 score of 0.1604. Subsequently, the system was adjusted to emphasize color symmetry by tripling the weight for the skewness of each color channel (with weights [1, 1, 3] * 3). This alteration maintained the AUC at 0.67, but the F1 score saw a slight decrease to 0.146.

Increasing the emphasis on color symmetry further, by assigning six times the weight to skewness (with weights [1, 1, 6] * 3), resulted in the AUC remaining unchanged at 0.67, and the F1 score experienced a negligible increase to 0.1463. This minute change suggests that the impact of skewness on the retrieval performance was minimal.

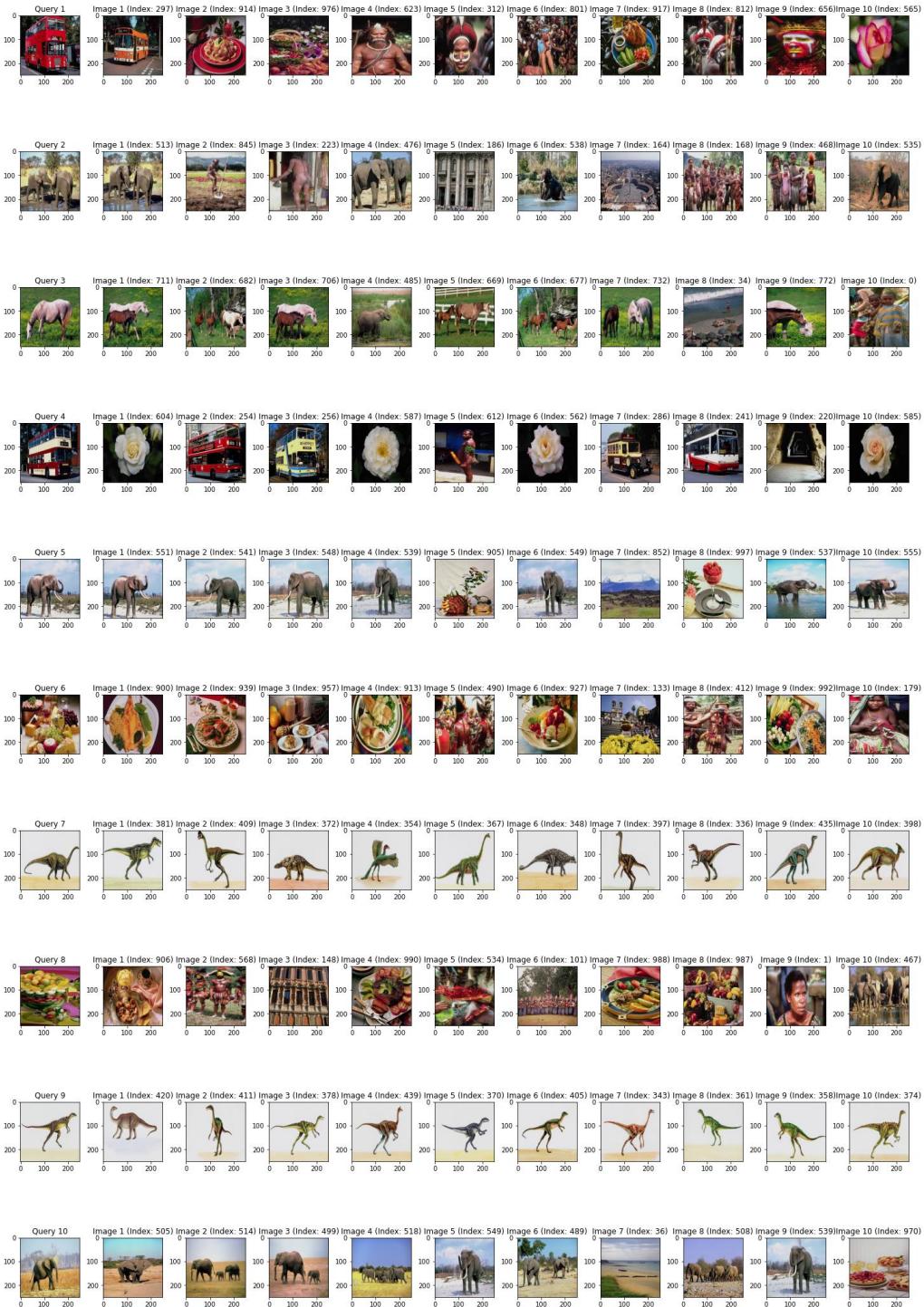
An even more pronounced weight adjustment was tested by assigning fifteen times the weight to skewness (with weights [1, 1, 15] * 3). Despite this significant emphasis, the AUC persisted at 0.67, while the F1 score dipped slightly to 0.1381.

These experiments indicate that while skewness is an aspect of color distribution in images, its influence on the retrieval performance of the CBIR system is limited. Even with increased weighting, the system's ability to discriminate between relevant and non-relevant images, as measured by the AUC, remained constant, and the precision-recall balance, reflected by the F1 score, showed a slight overall decline. The results suggest there is a threshold beyond which additional emphasis on color symmetry does not translate into improved performance, highlighting the necessity for a balanced approach to feature weighting in image retrieval systems.

4. *Custom Scenario*

- **weights** = [2, 8, 2] *3, giving more weight for standard deviation since it gave the best improvement.

Auc was 0.67 , 0.1962 , at some thresholds it reached 0.34 and 0.42



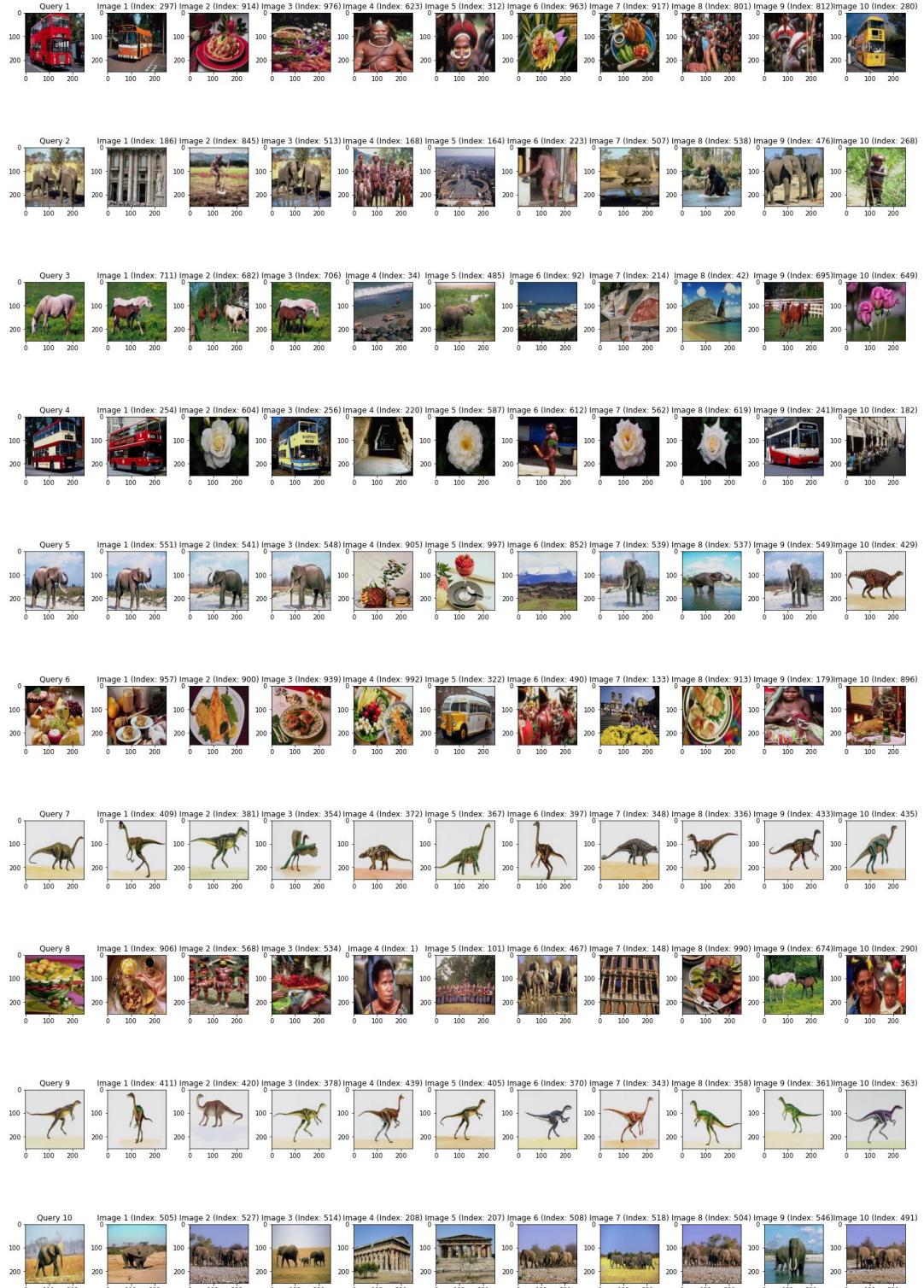
```

11
12 Threshold: 0.2260
13 Average Precision: 0.3187
14 Average Recall: 0.5100
15 Average F1 Score: 0.3923
16
17 Threshold: 0.2341
18 Average Precision: 0.3000
19 Average Recall: 0.5400
20 Average F1 Score: 0.3857
21
22 Threshold: 0.2422
23 Average Precision: 0.2935
24 Average Recall: 0.5900
25 Average F1 Score: 0.3920
26
27 Threshold: 0.2501
28 Average Precision: 0.2800
29 Average Recall: 0.6300
30 Average F1 Score: 0.3877
31
32 Threshold: 0.2620
33 Average Precision: 0.2586
34 Average Recall: 0.6800
35 Average F1 Score: 0.3747
36
37 Threshold: 0.2780
38 Average Precision: 0.2500
39 Average Recall: 0.7500
40 Average F1 Score: 0.3750
41
42 Threshold: 0.2943
43 Average Precision: 0.2185
44 Average Recall: 0.7800
45 Average F1 Score: 0.3414
46
47 Threshold: 0.2821
48 Average Precision: 0.2413
49 Average Recall: 0.7600
50 Average F1 Score: 0.3663
51
52 Threshold: 0.2984
53 Average Precision: 0.2174
54 Average Recall: 0.8000
55 Average F1 Score: 0.3419
56
57 --- Overall Average Evaluation Results -
58 Overall Average Precision: 0.2802
59 Overall Average Recall: 0.2423
60 Overall Average F1 Score: 0.1962
61
62 Total Execution Time: 9.8468 seconds

```

- weights = [3, 20, 2]*3

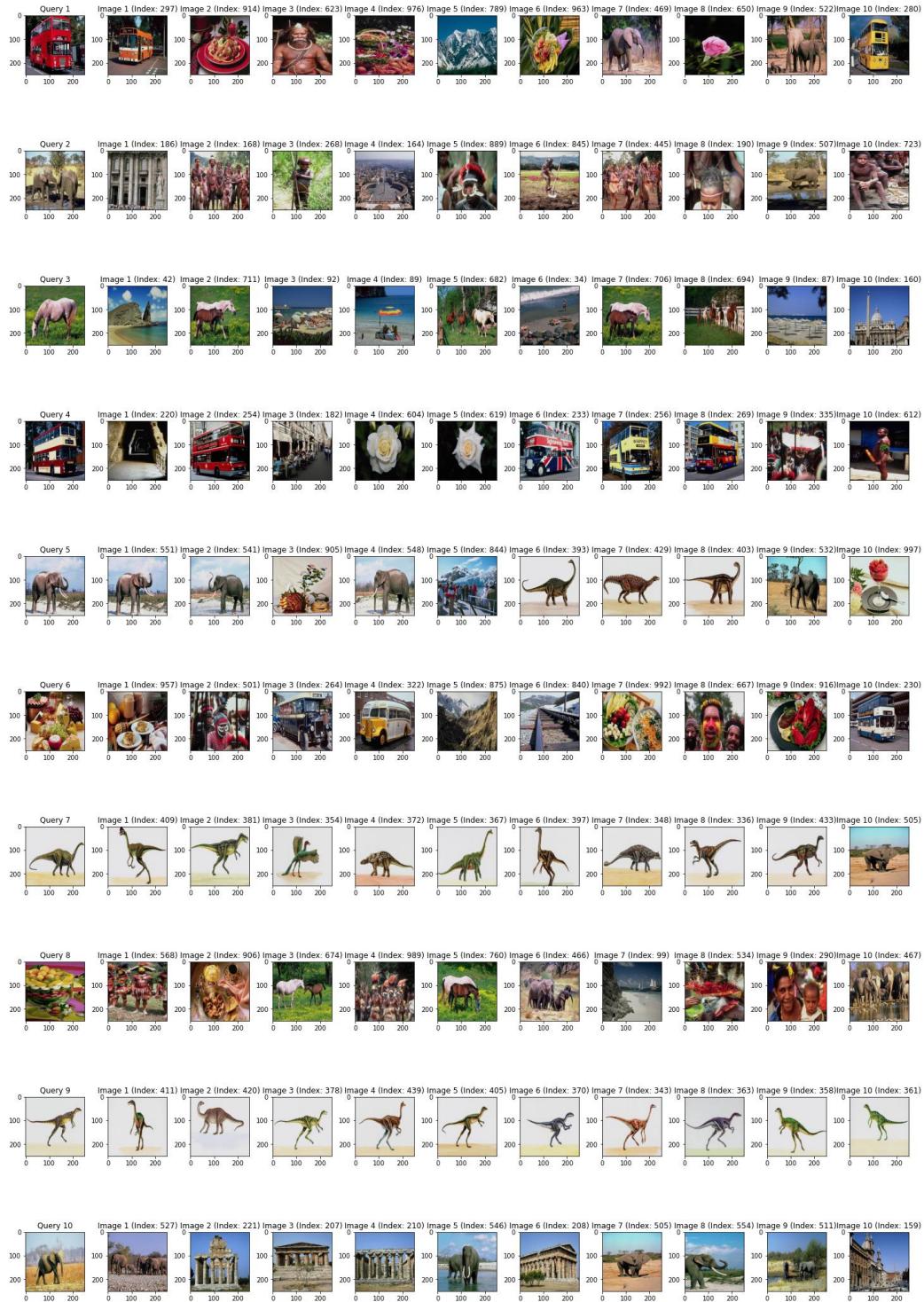
Auc gave 0.66 , f1 score was 0.2103



```
21
22 Threshold: 0.2133
23 Average Precision: 0.2605
24 Average Recall: 0.5600
25 Average F1 Score: 0.3556
26
27 Threshold: 0.2212
28 Average Precision: 0.2625
29 Average Recall: 0.6300
30 Average F1 Score: 0.3706
31
32 Threshold: 0.2291
33 Average Precision: 0.2538
34 Average Recall: 0.6600
35 Average F1 Score: 0.3667
36
37 Threshold: 0.2368
38 Average Precision: 0.2464
39 Average Recall: 0.6900
40 Average F1 Score: 0.3632
41
42 Threshold: 0.2483
43 Average Precision: 0.2271
44 Average Recall: 0.7200
45 Average F1 Score: 0.3453
46
47 Threshold: 0.2636
48 Average Precision: 0.2094
49 Average Recall: 0.7600
50 Average F1 Score: 0.3283
51
52 Threshold: 0.2785
53 Average Precision: 0.1971
54 Average Recall: 0.8100
55 Average F1 Score: 0.3170
56
57 Threshold: 0.2673
58 Average Precision: 0.2042
59 Average Recall: 0.7700
60 Average F1 Score: 0.3229
61
62 Threshold: 0.2823
63 Average Precision: 0.1920
64 Average Recall: 0.8200
65 Average F1 Score: 0.3112
66
67
68 --- Overall Average Evaluation Results ---
69 Overall Average Precision: 0.2353
70 Overall Average Recall: 0.2880
71 Overall Average F1 Score: 0.2103
72
73 Total Execution Time: 10.5262 seconds
Completed all experiments
```

- weights = [3, 50, 2]*3 , increasing more!

Auc was 62 , f1 score was 0.2107



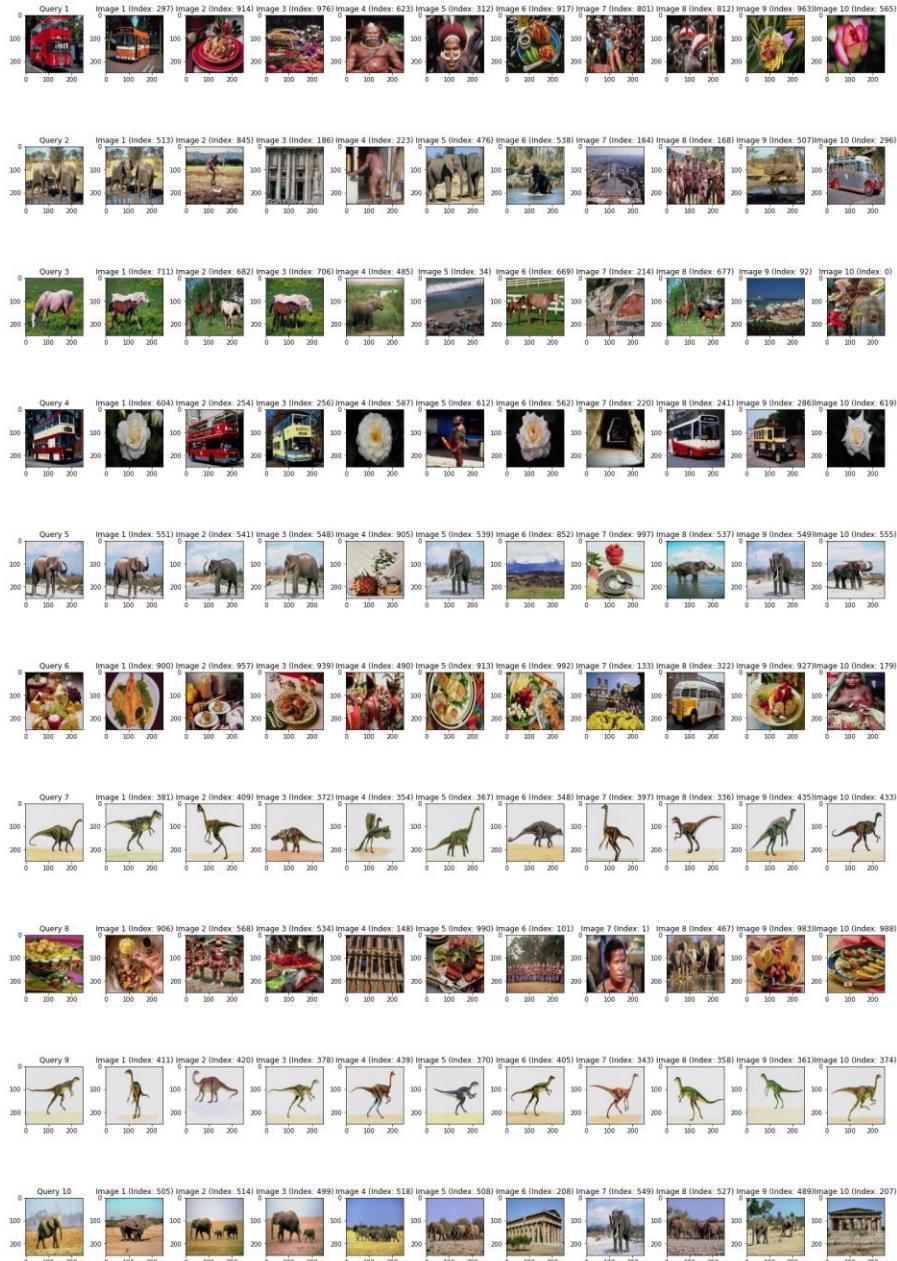
```

.067 --- Overall Average Evaluation Results ---
.068 Overall Average Precision: 0.1970
.069 Overall Average Recall: 0.3000
.070 Overall Average F1 Score: 0.2107
.071
.072
.073 Total Execution Time: 10.2847 seconds
.074 Completed all experiments

```

- weights = [10, 50, 20]*3

Auc was 67, f1 score was 0.1999



```

1022 Threshold: 0.2217
1023 Average Precision: 0.2865
1024 Average Recall: 0.5300
1025 Average F1 Score: 0.3719
1026
1027 Threshold: 0.2297
1028 Average Precision: 0.2823
1029 Average Recall: 0.5900
1030 Average F1 Score: 0.3819
1031
1032 Threshold: 0.2378
1033 Average Precision: 0.2826
1034 Average Recall: 0.6500
1035 Average F1 Score: 0.3939
1036
1037 Threshold: 0.2458
1038 Average Precision: 0.2588
1039 Average Recall: 0.6600
1040 Average F1 Score: 0.3718
1041
1042 Threshold: 0.2578
1043 Average Precision: 0.2509
1044 Average Recall: 0.7200
1045 Average F1 Score: 0.3721
1046
1047 Threshold: 0.2736
1048 Average Precision: 0.2269
1049 Average Recall: 0.7600
1050 Average F1 Score: 0.3494
1051
1052 Threshold: 0.2896
1053 Average Precision: 0.2109
1054 Average Recall: 0.8100
1055 Average F1 Score: 0.3347
1056
1057 Threshold: 0.2776
1058 Average Precision: 0.2197
1059 Average Recall: 0.7600
1060 Average F1 Score: 0.3408
1061
1062 Threshold: 0.2936
1063 Average Precision: 0.2071
1064 Average Recall: 0.8200
1065 Average F1 Score: 0.3306
1066
1067
1068 --- Overall Average Evaluation Results ---
1069 Overall Average Precision: 0.2560
1070 Overall Average Recall: 0.2633
1071 Overall Average F1 Score: 0.1999
1072
1073 Total Execution Time: 9.9895 seconds

```

Results

In the first scenario, weights were adjusted to [2, 8, 2] * 3, giving more weight to the standard deviation, which previously showed promising improvement. The AUC remained steady

comparing:

custom

at 0.67, but the F1 score improved to 0.1962, indicating a better balance between precision and recall. Notably, at certain thresholds, the F1 score reached higher peaks of 0.34 and 0.42, suggesting that at specific operational points, the system demonstrated considerably improved performance.

The next configuration, with weights [3, 20, 2] * 3, yielded a slight decrease in the AUC to 0.66 but an enhanced F1 score of 0.2103. This further increase in the weight of standard deviation continued to show an overall positive trend in retrieval quality as assessed by the F1 score.

Pushing the standard deviation weight even more to [3, 50, 2] * 3, the AUC dropped to 0.62, but the F1 score saw a marginal improvement to 0.2107. This suggests that an extreme emphasis on standard deviation may start to adversely affect the system's discriminatory power (as AUC decreases) while still marginally improving the F1 score.

The final scenario tested a significant alteration in the weights to [10, 50, 20] * 3. Here, the AUC returned to 0.67, and the F1 score was 0.1999. This indicates that a substantial increase in the weights of all color channels provided a balance that brought the AUC back to the level of the baseline, with an improved F1 score, although not the highest observed.

conclusion of this part:

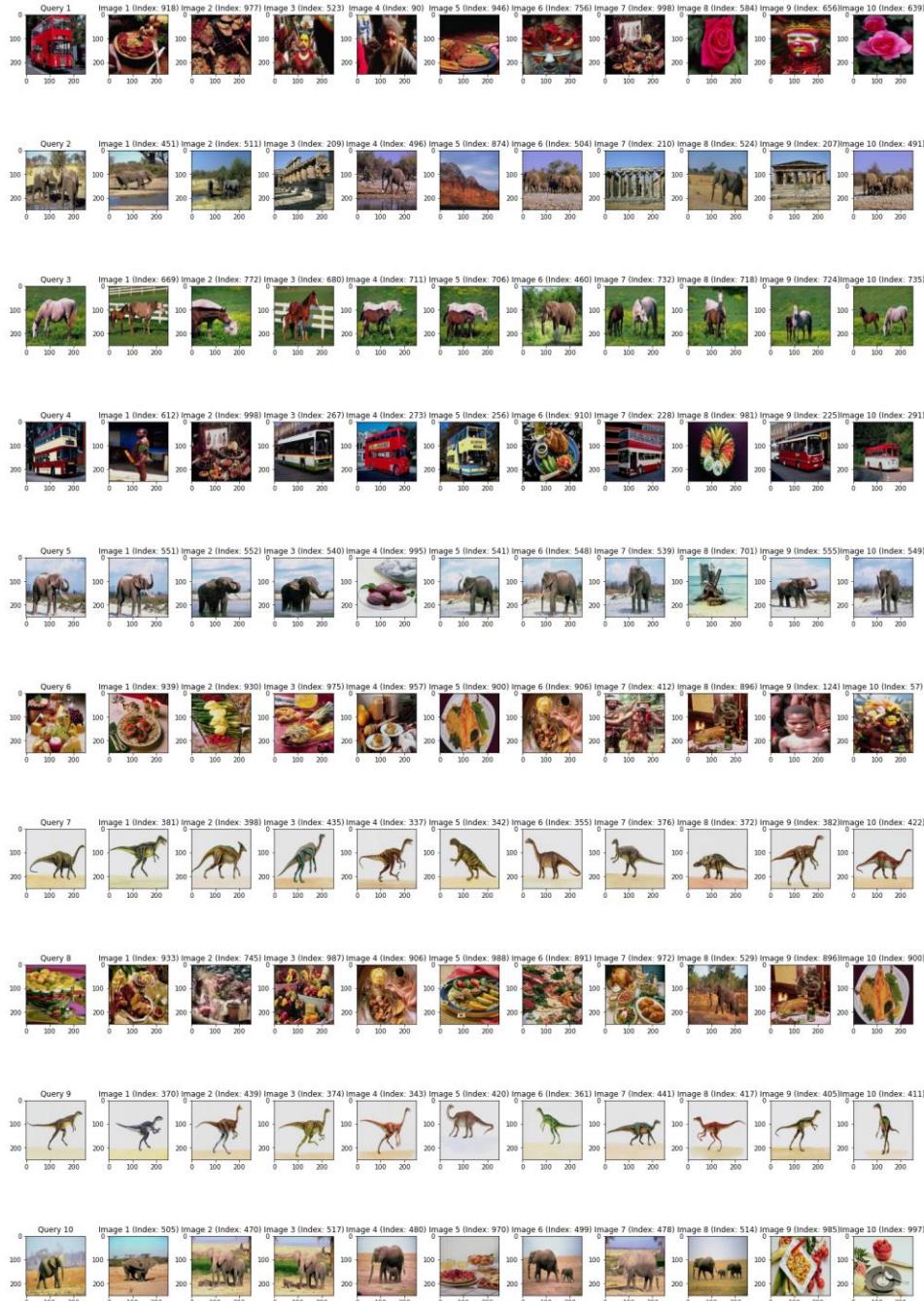
These scenarios reveal a complex relationship between the weights of the color channels and the performance of the CBIR system. While certain weight distributions can enhance the F1 score, suggesting an improved precision-recall balance, they may not necessarily increase the AUC. Also, while standard deviation seems to have a positive impact on the precision and recall balance of the system, it does not unilaterally improve all aspects of the system's performance. This suggests that while standard deviation is an important feature, the system also relies on a balanced weighting of other features to perform optimally across all metrics.

Subsection 3.2.3: Additional Moments

1) Balanced Approach

- Weights: [1, 1, 1, 1, 1, 1] * 3 (Equal weight for all color moments)

Auc was 0.68 , f1 score 0.2327



```

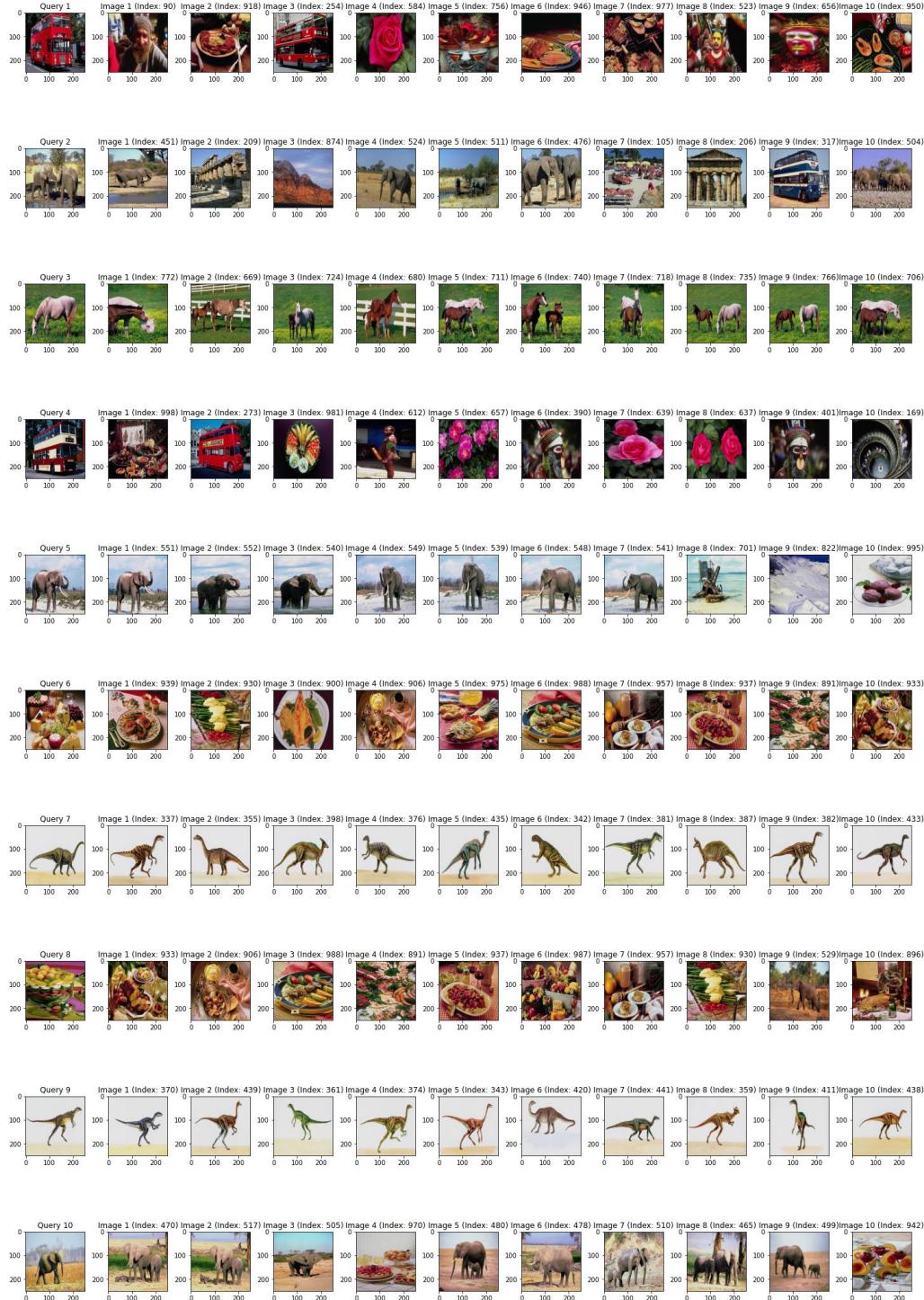
1 At Threshold: 0.5528, Precision: 0.15, Recall: 0.84, F1 Score: 0.2
2 Average values for all thresholds for Query Index: 979
3 avg-Precision: 0.29
4 avg-Recall: 0.31
5 avg-F1 Score: 0.23
6
7 Average values for all thresholds for Query Index: 373
8 avg-Precision: 0.32
9 avg-Recall: 0.34
10 avg-F1 Score: 0.27
11 Query Index: 472
12
13
14
15
16
17
18
19
20
21
22 Threshold: 0.2443
23 Average Precision: 0.2933
24 Average Recall: 0.6100
25 Average F1 Score: 0.3961
26
27 Threshold: 0.2543
28 Average Precision: 0.2889
29 Average Recall: 0.6500
30 Average F1 Score: 0.4000
31
32 Threshold: 0.2643
33 Average Precision: 0.2936
34 Average Recall: 0.6900
35 Average F1 Score: 0.4119
36
37 Threshold: 0.2742
38 Average Precision: 0.2893
39 Average Recall: 0.7000
40 Average F1 Score: 0.4094
41
42 Threshold: 0.2891
43 Average Precision: 0.2812
44 Average Recall: 0.7200
45 Average F1 Score: 0.4045
46
47 Threshold: 0.3086
48 Average Precision: 0.2637
49 Average Recall: 0.7200
50 Average F1 Score: 0.3861
51
52 Threshold: 0.3279
53 Average Precision: 0.2559
54 Average Recall: 0.7600
55 Average F1 Score: 0.3829
56
57 Threshold: 0.3136
58 Average Precision: 0.2616
59 Average Recall: 0.7300
60 Average F1 Score: 0.3852
61
62 Threshold: 0.3326
63 Average Precision: 0.2476
64 Average Recall: 0.7600
65 Average F1 Score: 0.3735
66
67
68 --- Overall Average Evaluation Results ---
69 Overall Average Precision: 0.3132
70 Overall Average Recall: 0.2765
71 Overall Average F1 Score: 0.2327
72
73 Total Execution Time: 10.7751 seconds

```

2) Emphasizing Median Color Value

- Weights: [1, 1, 1, 3, 1, 1] * 3 (3 times the weight for the median of each color channel)

Auc was 0.66, f1 score 0.1850



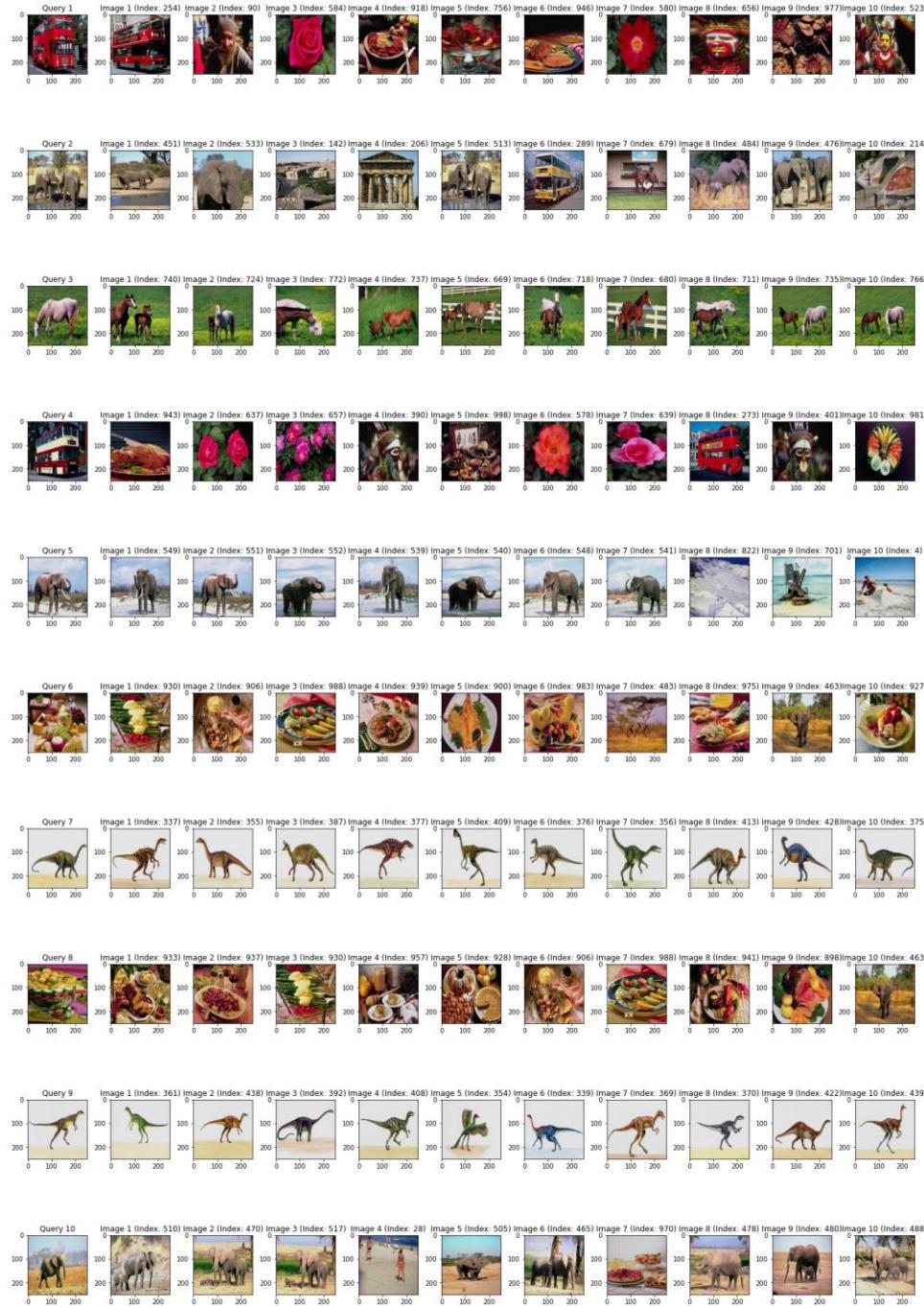
```
225  
226  
227 Threshold: 0.1981  
228 Average Precision: 0.2660  
229 Average Recall: 0.2500  
230 Average F1 Score: 0.2577  
231  
232 Threshold: 0.2049  
233 Average Precision: 0.2475  
234 Average Recall: 0.2500  
235 Average F1 Score: 0.2488  
236  
237 Threshold: 0.2117  
238 Average Precision: 0.2455  
239 Average Recall: 0.2700  
240 Average F1 Score: 0.2571  
241  
242 Threshold: 0.2218  
243 Average Precision: 0.2377  
244 Average Recall: 0.2900  
245 Average F1 Score: 0.2613  
246  
247 Threshold: 0.2356  
248 Average Precision: 0.2083|  
249 Average Recall: 0.3000  
250 Average F1 Score: 0.2459  
251  
252 Threshold: 0.2497  
253 Average Precision: 0.1902  
254 Average Recall: 0.3100  
255 Average F1 Score: 0.2357  
256  
257 Threshold: 0.2391  
258 Average Precision: 0.2081  
259 Average Recall: 0.3100  
260 Average F1 Score: 0.2490  
261  
262 Threshold: 0.2532  
263 Average Precision: 0.1845  
264 Average Recall: 0.3100  
265 Average F1 Score: 0.2313  
266  
267  
268 --- Overall Average Evaluation Results ---  
269 Overall Average Precision: 0.3408  
270 Overall Average Recall: 0.1505  
271 Overall Average F1 Score: 0.1850  
272  
273 Total Execution Time: 9.5296 seconds  
274 Completed all experiments
```

```
Average values for all thresholds for Query Index: 373  
avg-Precision: 0.34  
avg-Recall: 0.33  
avg-F1 Score: 0.28  
Query_Index: 472
```

```
At Threshold: 0.2552, Precision: 0.15, Recall: 0.75, F1: 0.11  
Average values for all thresholds for Query Index: 979  
avg-Precision: 0.31  
avg-Recall: 0.30  
avg-F1 Score: 0.24
```

- Weights: [1, 1, 1, 10, 1, 1] * 3 (3 times the weight for the median of each color channel)

Auc was 0.65, f1 score 0.1643



```
Threshold: 0.1545
Average Precision: 0.3137
Average Recall: 0.1600
Average F1 Score: 0.2119

Threshold: 0.1612
Average Precision: 0.3220
Average Recall: 0.1900
Average F1 Score: 0.2390

Threshold: 0.1679
Average Precision: 0.2923
Average Recall: 0.1900
Average F1 Score: 0.2303

Threshold: 0.1747
Average Precision: 0.2632
Average Recall: 0.2000
Average F1 Score: 0.2273

Threshold: 0.1850
Average Precision: 0.2553
Average Recall: 0.2400
Average F1 Score: 0.2474

Threshold: 0.1987
Average Precision: 0.2250
Average Recall: 0.2700
Average F1 Score: 0.2455

Threshold: 0.2132
Average Precision: 0.2101
Average Recall: 0.2900
Average F1 Score: 0.2437

Threshold: 0.2023
Average Precision: 0.2295
Average Recall: 0.2800
Average F1 Score: 0.2523

Threshold: 0.2169
Average Precision: 0.2083
Average Recall: 0.3000
Average F1 Score: 0.2459

--- Overall Average Evaluation Results ---
Overall Average Precision: 0.3812
Overall Average Recall: 0.1217
Overall Average F1 Score: 0.1643

Total Execution Time: 10.1377 seconds
Completed all experiments
```

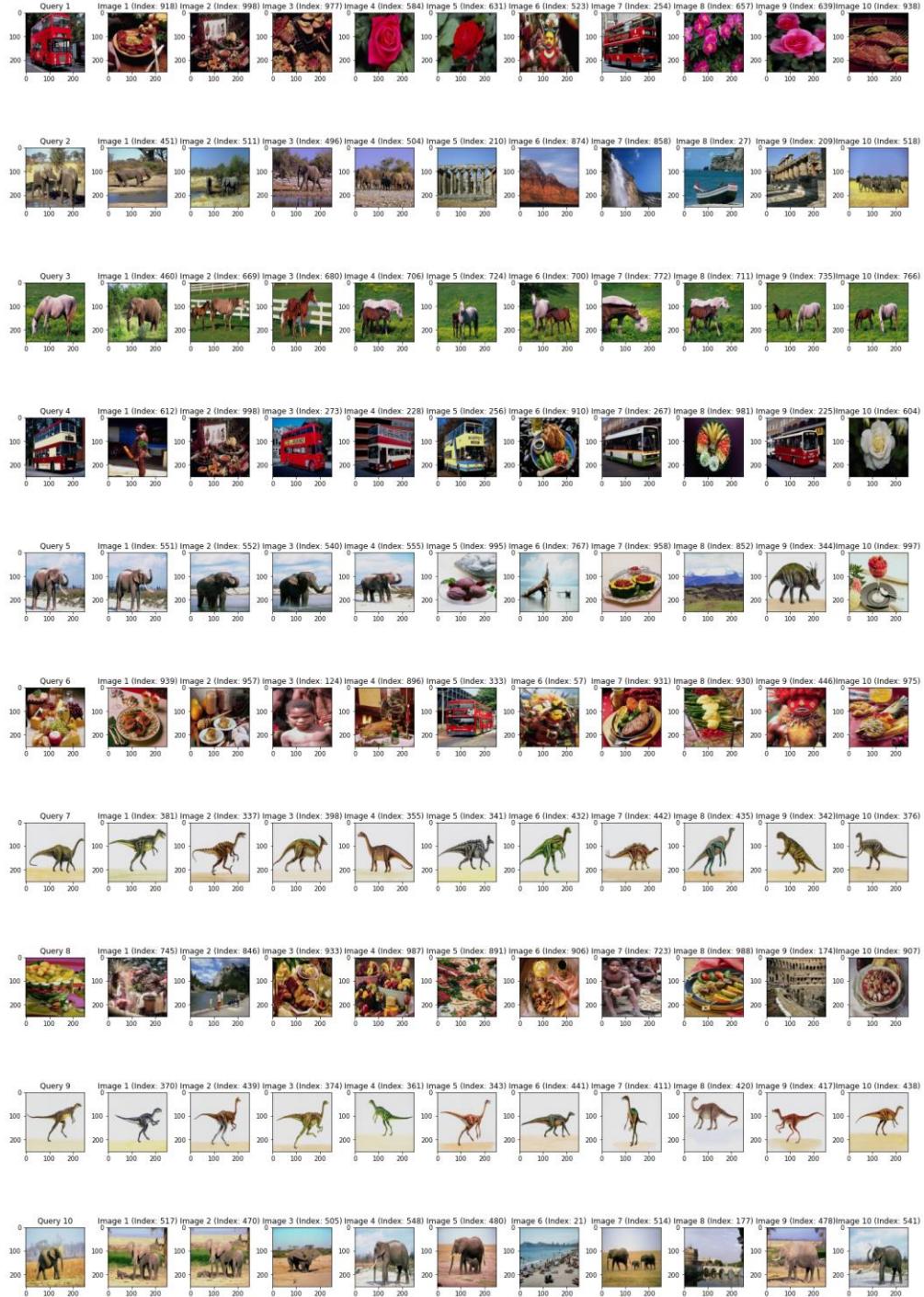
```
At threshold: 0.2169, precision: 0.50, recall: 0.59, f1 score: 0.49
Average values for all thresholds for Query Index: 373
avg-Precision: 0.37
avg-Recall: 0.32
avg-F1 Score: 0.28
Query Index: 472
```

```
At Threshold: 0.2169, Precision: 0.14, Recall: 0.70, F1 Score: 0.14
Average values for all thresholds for Query Index: 979
avg-Precision: 0.35
avg-Recall: 0.29
avg-F1 Score: 0.24
Query Index: 472
```

3) Emphasizing Most Common Color (Mode)

- Weights: [1, 1, 1, 1, 3, 1] * 3 (3 times the weight for the mode of each color channel)

Auc was 67, f1 score 0.2597



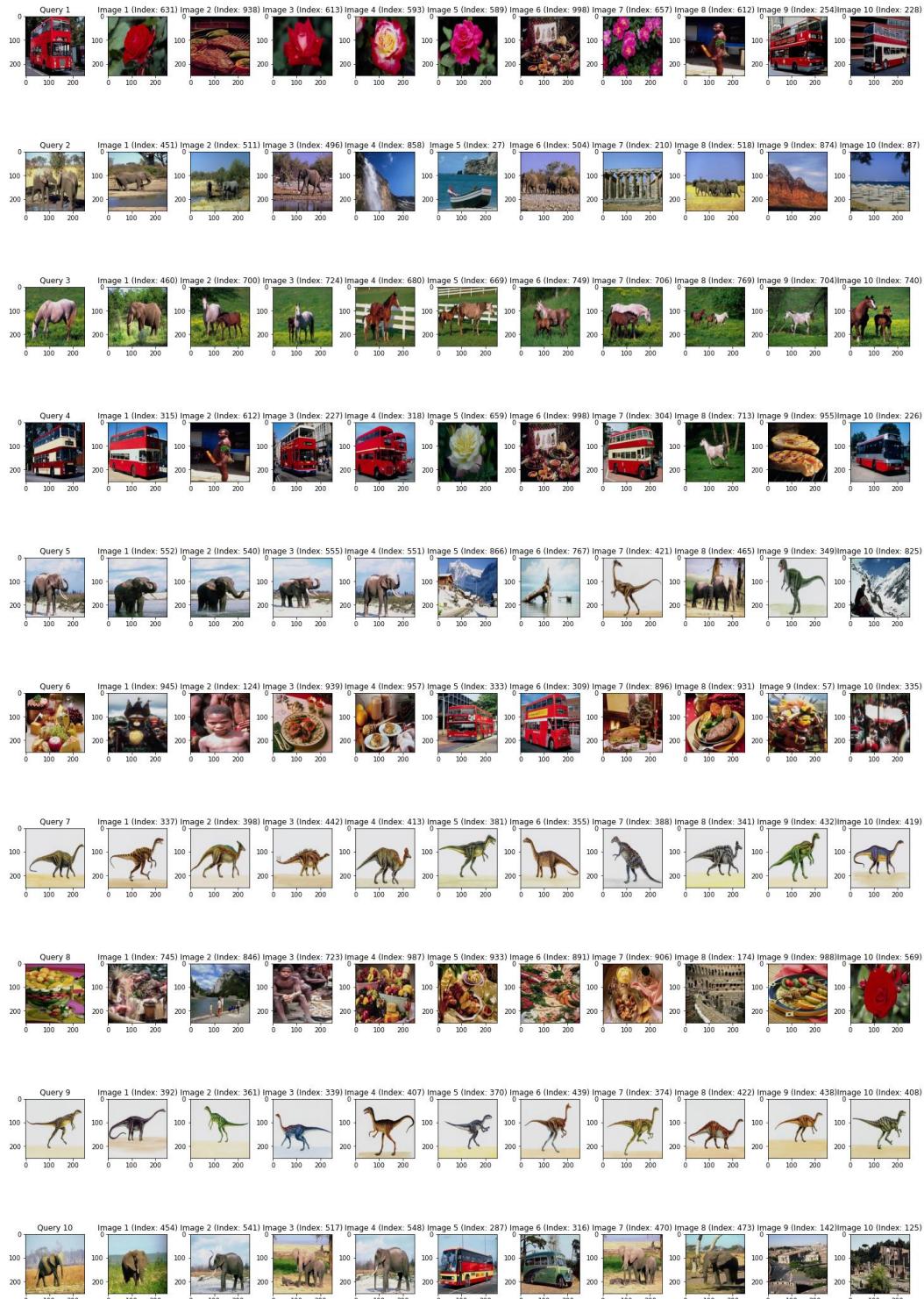
```

22 Threshold: 0.2744
23 Average Precision: 0.2863
24 Average Recall: 0.6700
25 Average F1 Score: 0.4012
26
27 Threshold: 0.2898
28 Average Precision: 0.2827
29 Average Recall: 0.6700
30 Average F1 Score: 0.3976
31
32 Threshold: 0.3050
33 Average Precision: 0.2757
34 Average Recall: 0.6700
35 Average F1 Score: 0.3907
36
37 Threshold: 0.3199
38 Average Precision: 0.2742
39 Average Recall: 0.6800
40 Average F1 Score: 0.3908
41
42 Threshold: 0.3424
43 Average Precision: 0.2672
44 Average Recall: 0.7000
45 Average F1 Score: 0.3867
46
47 Threshold: 0.3703
48 Average Precision: 0.2552
49 Average Recall: 0.7300
50 Average F1 Score: 0.3782
51
52 Threshold: 0.3972
53 Average Precision: 0.2468
54 Average Recall: 0.7600
55 Average F1 Score: 0.3725
56
57 Threshold: 0.3768
58 Average Precision: 0.2561
59 Average Recall: 0.7400
60 Average F1 Score: 0.3805
61
62 Threshold: 0.4047
63 Average Precision: 0.2460
64 Average Recall: 0.7700
65 Average F1 Score: 0.3729
66
67 --- Overall Average Evaluation Results ---
68 Overall Average Precision: 0.2980
69 Overall Average Recall: 0.3590
70 Overall Average F1 Score: 0.2597
71
72 Total Execution Time: 10.1787 seconds
Completed all measurements

```

- Weights: [1, 1, 1, 1, 15, 1] * 3 (3 times the weight for the mode of each color channel)

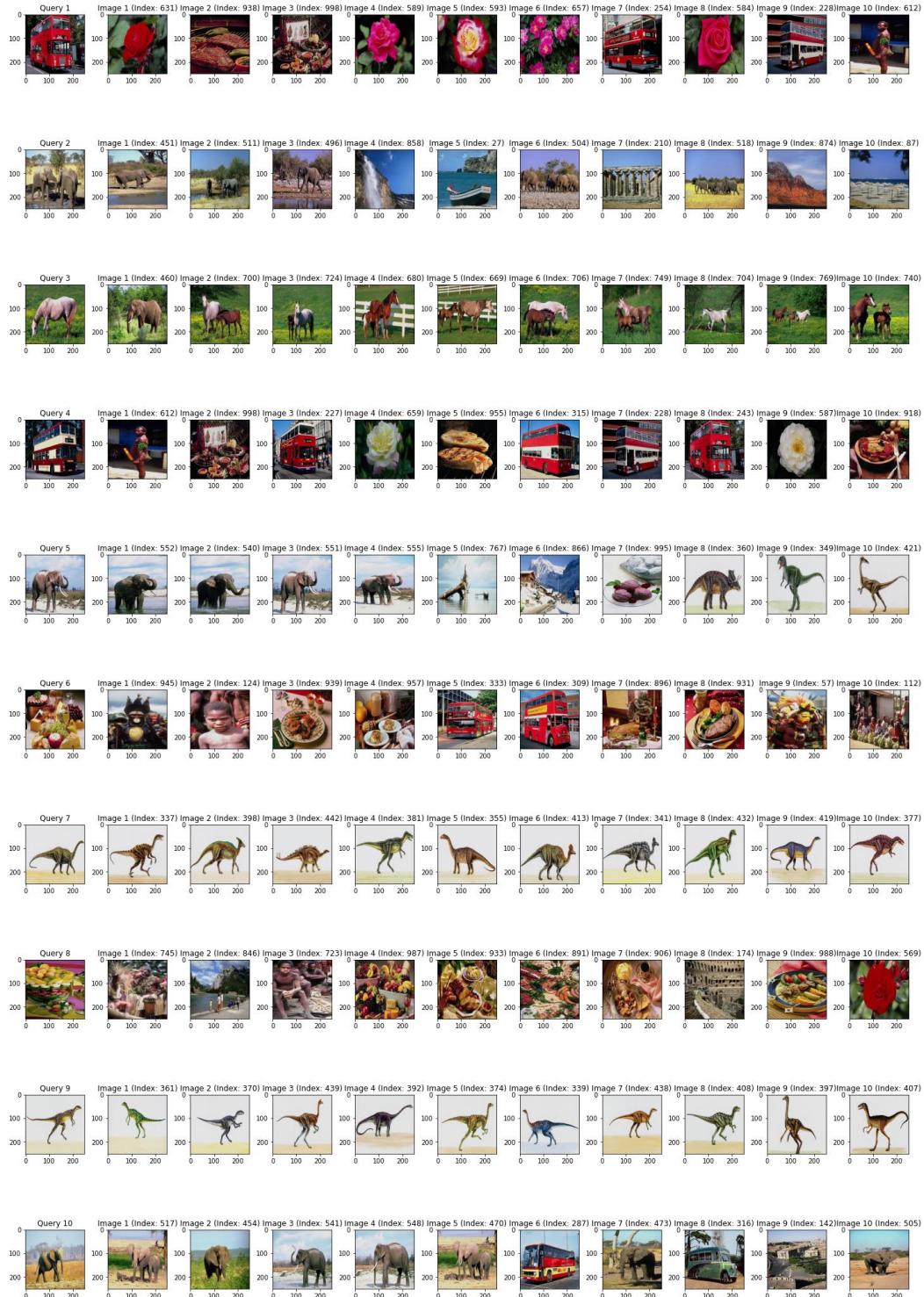
Auc was 66, f1 score 0.2319



```
.1
.2 Threshold: 0.2790
.3 Average Precision: 0.2863
.4 Average Recall: 0.6700
.5 Average F1 Score: 0.4012
.6
.7 Threshold: 0.2960
.8 Average Precision: 0.2827
.9 Average Recall: 0.6700
.0 Average F1 Score: 0.3976
.1
.2 Threshold: 0.3129
.3 Average Precision: 0.2769
.4 Average Recall: 0.6700
.5 Average F1 Score: 0.3918
.6
.7 Threshold: 0.3296
.8 Average Precision: 0.2731
.9 Average Recall: 0.6800
.0 Average F1 Score: 0.3897
.1
.2 Threshold: 0.3537
.3 Average Precision: 0.2700
.4 Average Recall: 0.7100
.5 Average F1 Score: 0.3912
.6
.7 Threshold: 0.3826
.8 Average Precision: 0.2578
.9 Average Recall: 0.7400
.0 Average F1 Score: 0.3824
.1
.2 Threshold: 0.4108
.3 Average Precision: 0.2468
.4 Average Recall: 0.7600
.5 Average F1 Score: 0.3725
.6
.7 Threshold: 0.3896
.8 Average Precision: 0.2577
.9 Average Recall: 0.7500
.0 Average F1 Score: 0.3836
.1
.2 Threshold: 0.4182
.3 Average Precision: 0.2420
.4 Average Recall: 0.7600
.5 Average F1 Score: 0.3671
.6
.7
.8 --- Overall Average Evaluation Results ---
.9 Overall Average Precision: 0.2923
.0 Overall Average Recall: 0.3448
.1 Overall Average F1 Score: 0.2319
.2
.3 Total Execution Time: 9.9486 seconds
Completed 11 experiments
```

- Weights: [1, 1, 1, 1, 10, 1] * 3 (3 times the weight for the mode of each color channel)

Auc was 66, f1 score 0.2331



```

--- Overall Average Evaluation Results ---
Overall Average Precision: 0.2838
Overall Average Recall: 0.3435
Overall Average F1 Score: 0.2331

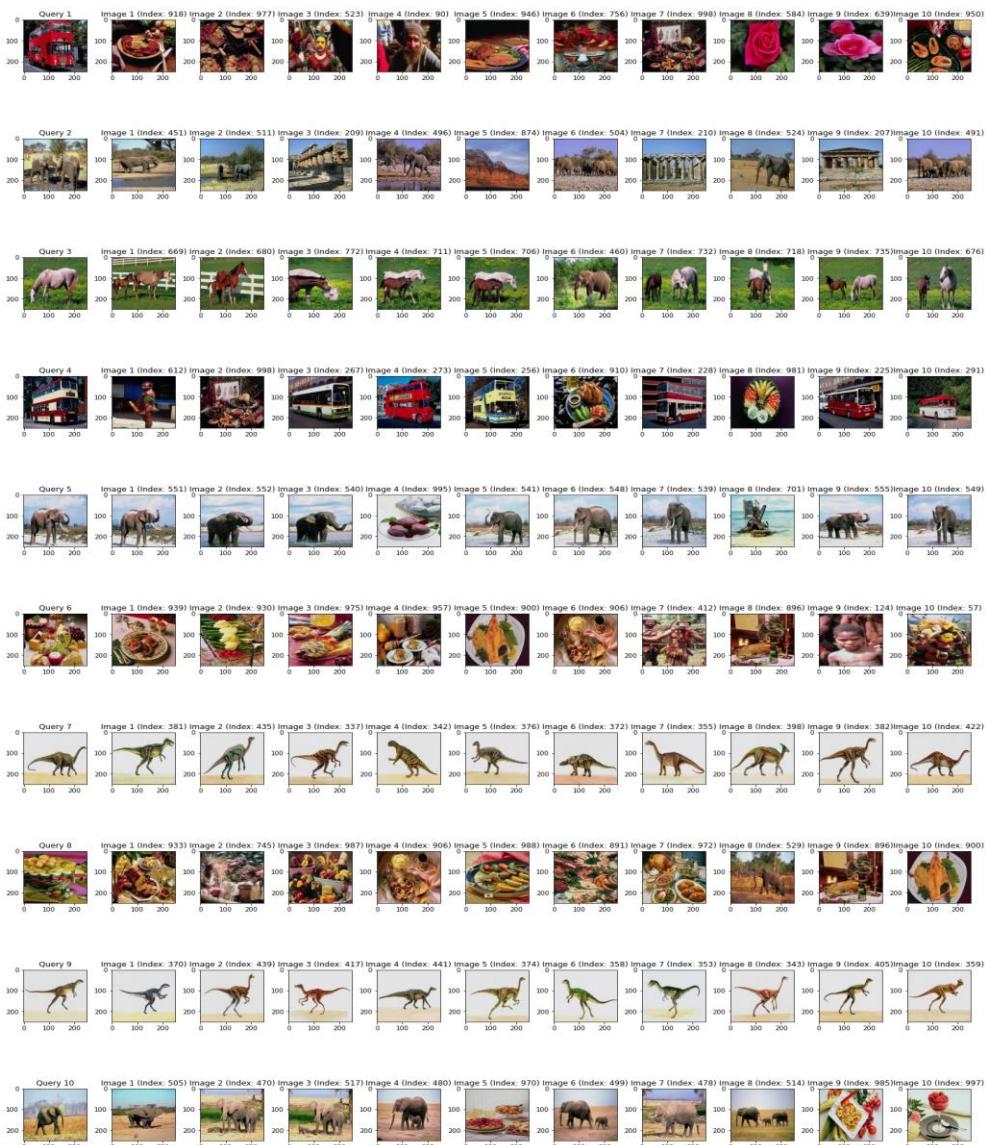
Total Execution Time: 9.9612 seconds

```

4) Emphasizing Color Distribution Shape (Kurtosis)

- Weights: [1, 1, 1, 1, 1, 3] * 3 (3 times the weight for the kurtosis of each color channel)

AUC was 0.68 , f1 score 0.2332



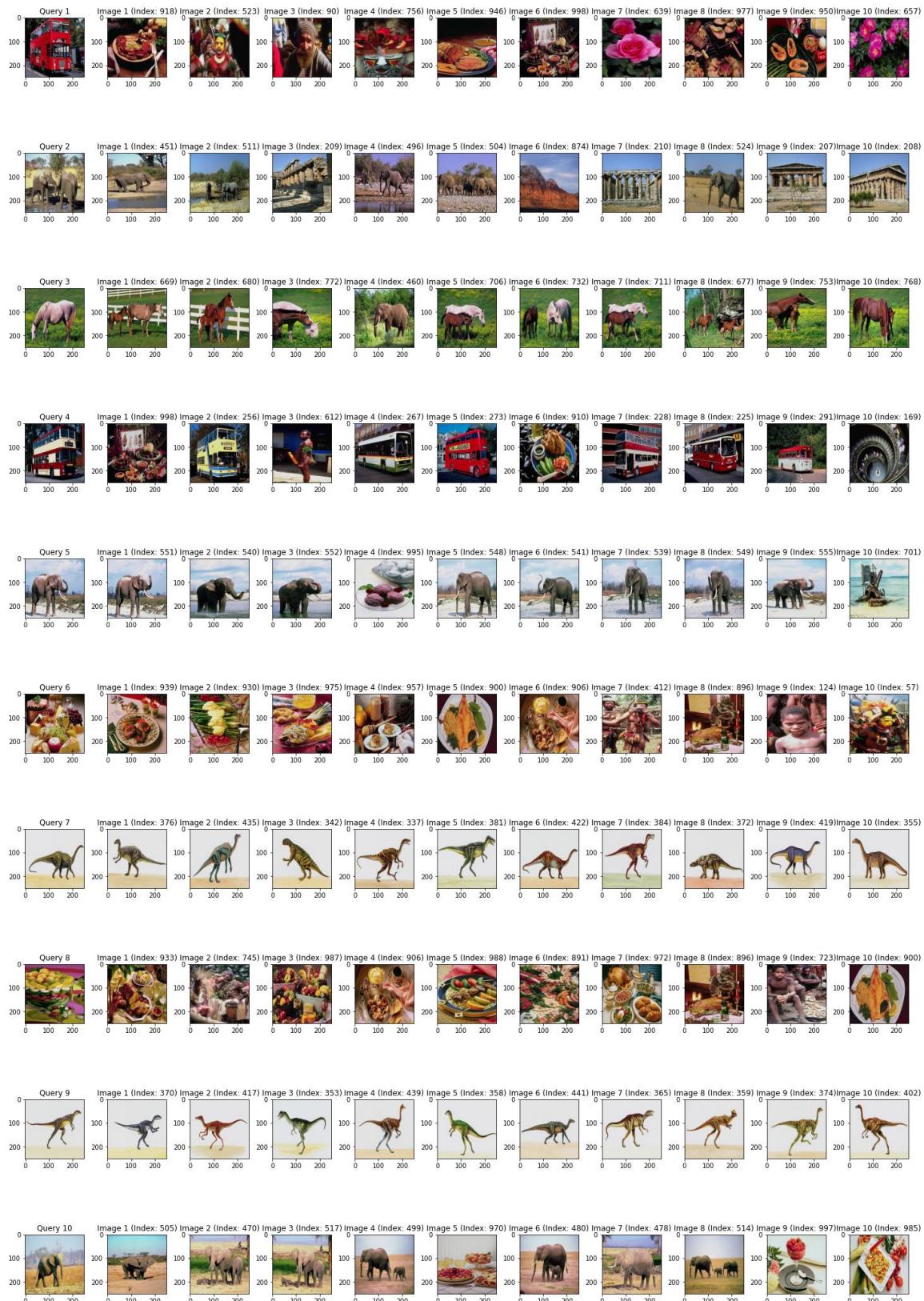
```

1
2 Threshold: 0.2442
3 Average Precision: 0.2913
4 Average Recall: 0.6000
5 Average F1 Score: 0.3922
6
7 Threshold: 0.2542
8 Average Precision: 0.2928
9 Average Recall: 0.6500
0 Average F1 Score: 0.4037
1
2 Threshold: 0.2639
3 Average Precision: 0.2918
4 Average Recall: 0.6800
5 Average F1 Score: 0.4084
6
7 Threshold: 0.2737
8 Average Precision: 0.2881
9 Average Recall: 0.7000
0 Average F1 Score: 0.4082
1
2 Threshold: 0.2885
3 Average Precision: 0.2802
4 Average Recall: 0.7200
5 Average F1 Score: 0.4034
6
7 Threshold: 0.3078
8 Average Precision: 0.2628
9 Average Recall: 0.7200
0 Average F1 Score: 0.3850
1
2 Threshold: 0.3270
3 Average Precision: 0.2559
4 Average Recall: 0.7600
5 Average F1 Score: 0.3829
6
7 Threshold: 0.3128
8 Average Precision: 0.2616
9 Average Recall: 0.7300
0 Average F1 Score: 0.3852
1
2 Threshold: 0.3316
3 Average Precision: 0.2476
4 Average Recall: 0.7600
5 Average F1 Score: 0.3735
6
7
8 --- Overall Average Evaluation Results ---
9 Overall Average Precision: 0.3137
0 Overall Average Recall: 0.2767
1 Overall Average F1 Score: 0.2332
2
3 Total Execution Time: 10.3334 seconds
Completed all experiments.

```

- Weights: [1, 1, 1, 1, 1, 10] * 3 (3 times the weight for the kurtosis of each color channel)

AUC was 0.68 , f1 score 0.2336



```

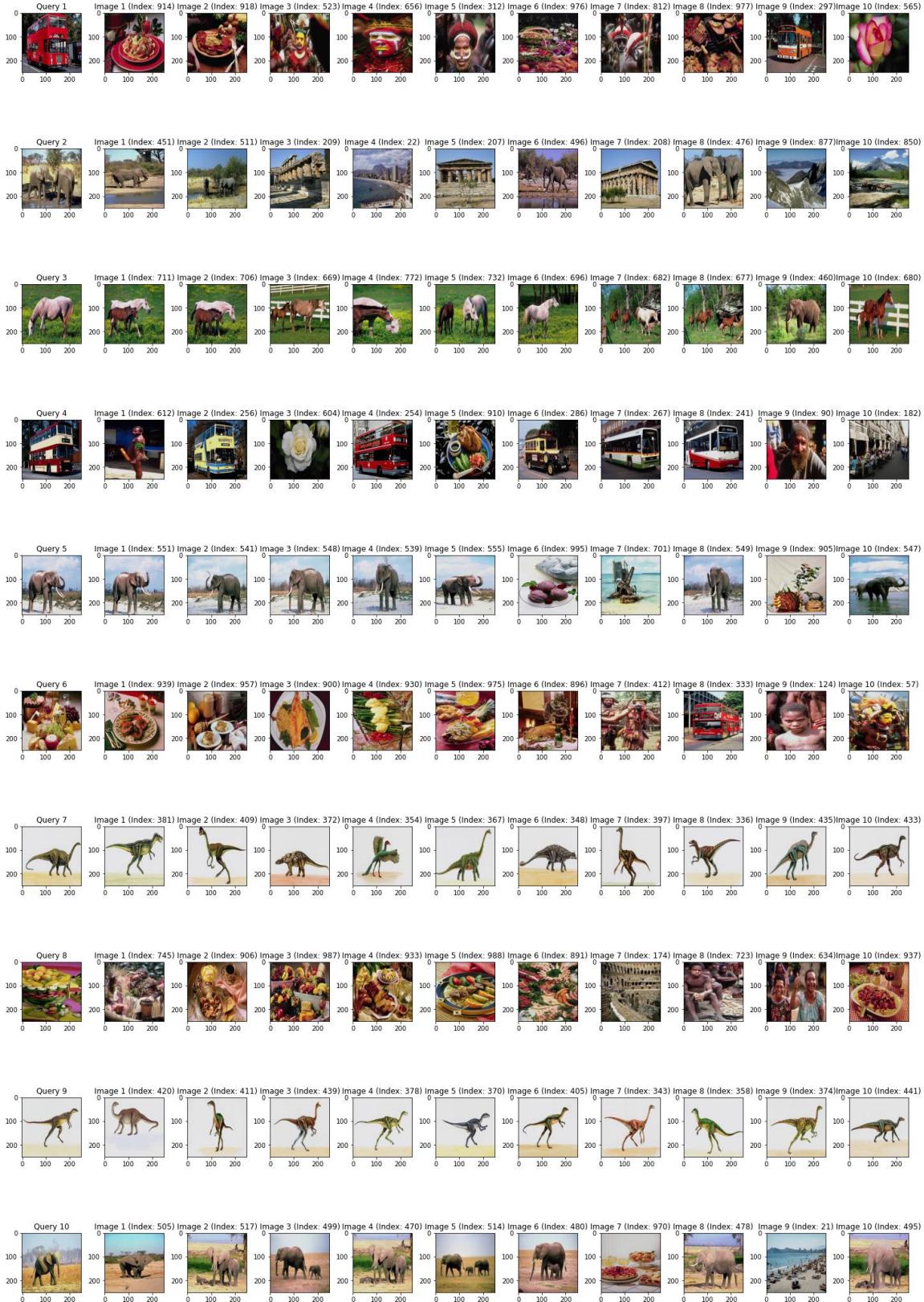
1 Threshold: 0.2293
2 Average Precision: 0.2828
3 Average Recall: 0.5600
4 Average F1 Score: 0.3758
5
6 Threshold: 0.2379
7 Average Precision: 0.2823
8 Average Recall: 0.5900
9 Average F1 Score: 0.3819
10
11 Threshold: 0.2464
12 Average Precision: 0.2928
13 Average Recall: 0.6500
14 Average F1 Score: 0.4037
15
16 Threshold: 0.2549
17 Average Precision: 0.2851
18 Average Recall: 0.6700
19 Average F1 Score: 0.4000
20
21 Threshold: 0.2676
22 Average Precision: 0.2688
23 Average Recall: 0.6800
24 Average F1 Score: 0.3853
25
26 Threshold: 0.2842
27 Average Precision: 0.2647
28 Average Recall: 0.7200
29 Average F1 Score: 0.3871
30
31 Threshold: 0.3006
32 Average Precision: 0.2542
33 Average Recall: 0.7500
34 Average F1 Score: 0.3797
35
36 Threshold: 0.2883
37 Average Precision: 0.2635
38 Average Recall: 0.7300
39 Average F1 Score: 0.3873
40
41 Threshold: 0.3046
42 Average Precision: 0.2492
43 Average Recall: 0.7500
44 Average F1 Score: 0.3741
45
46
47 --- Overall Average Evaluation Results ---
48 Overall Average Precision: 0.3209
49 Overall Average Recall: 0.2750
50 Overall Average F1 Score: 0.2336
51
52 Total Execution Time: 10.2251 seconds
53 Completed all experiments

```

5) Additional scenario

- weights=[2, 10, 2, 3, 3, 6]*3

AUC was 0.69, f1 score 0.2897



```
--- Overall Average Evaluation Results ---
Overall Average Precision: 0.3268
Overall Average Recall: 0.3768
Overall Average F1 Score: 0.2897
```

```
Total Execution Time: 10.6207 seconds
```

```
5 Threshold: 0.2720
6 Average Precision: 0.2982
7 Average Recall: 0.6800
8 Average F1 Score: 0.4146
9
10 Threshold: 0.2819
11 Average Precision: 0.2975
12 Average Recall: 0.7200
13 Average F1 Score: 0.4211
14
15 Threshold: 0.2916
16 Average Precision: 0.2903
17 Average Recall: 0.7200
18 Average F1 Score: 0.4138
19
20 Threshold: 0.3013
21 Average Precision: 0.2780
22 Average Recall: 0.7200
23 Average F1 Score: 0.4011
24
25 Threshold: 0.3108
26 Average Precision: 0.2657
27 Average Recall: 0.7200
28 Average F1 Score: 0.3881
29
30 Threshold: 0.3247
31 Average Precision: 0.2598
32 Average Recall: 0.7300
33 Average F1 Score: 0.3832
34
35 Threshold: 0.3420
36 Average Precision: 0.2607
37 Average Recall: 0.7900
38 Average F1 Score: 0.3921
39
40 Threshold: 0.3584
41 Average Precision: 0.2455
42 Average Recall: 0.8100
43 Average F1 Score: 0.3767
44
45 Threshold: 0.3461
46 Average Precision: 0.2581
47 Average Recall: 0.8000
48 Average F1 Score: 0.3902
49
50 Threshold: 0.3626
51 Average Precision: 0.2426
52 Average Recall: 0.8200
53 Average F1 Score: 0.3744
54
```

```
At Threshold: 0.3626, Precision: 0.36, Recall: 0.59, F1 Score: 0.3744
Average values for all thresholds for Query Index: 373
avg-Precision: 0.28
avg-Recall: 0.35
avg-F1 Score: 0.27
```

Discussion of 3.3:

Previous Results with Mean, Standard Deviation, and Skewness:

- The baseline with equal weights (mean, standard deviation, skewness) achieved an AUC of 0.67 and an F1 score of 0.1604.
- Adjusting weights to favor standard deviation showed some improvements in the F1 score but mixed results for AUC.

Results with Additional Moments:

- Balanced Approach with All Moments:

Weights: [1, 1, 1, 1, 1, 1] * 3

AUC: 0.68, F1 score: 0.2327

Including all moments with equal weighting resulted in a slight improvement in AUC and a significant increase in the F1 score compared to the baseline, indicating that these additional moments contribute positively to the system's performance.

Emphasizing Median Color Value:

Slight emphasis ([1, 1, 1, 3, 1, 1] * 3): AUC decreased to 0.66, F1 score to 0.1850.

Strong emphasis ([1, 1, 1, 10, 1, 1] * 3): AUC further decreased to 0.65, F1 score to 0.1643.

It appears that emphasizing the median color value diminishes the system's effectiveness, as both the AUC and F1 score decreased compared to the balanced approach.

Emphasizing Most Common Color (Mode):

Moderate emphasis ([1, 1, 1, 1, 3, 1] * 3): AUC was 0.67, F1 score 0.2597.

Strong emphasis ([1, 1, 1, 1, 15, 1] * 3): AUC decreased to 0.66, F1 score to 0.2319.

Different strong emphasis ([1, 1, 1, 1, 10, 1] * 3): AUC remained at 0.66, F1 score slightly higher at 0.2331.

These results suggest that while the mode is an important feature, overemphasis beyond a certain threshold does not necessarily yield better retrieval performance.

Emphasizing Color Distribution Shape (Kurtosis):

Moderate emphasis ($[1, 1, 1, 1, 1, 3] * 3$): AUC was 0.68, F1 score 0.2332.

Strong emphasis ($[1, 1, 1, 1, 1, 10] * 3$): AUC was consistent at 0.68, F1 score marginally increased to 0.2336.

Emphasizing kurtosis slightly improves the F1 score without affecting the AUC, indicating that kurtosis contributes to a balanced precision-recall trade-off.

Additional Custom Scenario:

Weights: $[2, 10, 2, 3, 3, 6] * 3$

AUC was 0.69, F1 score 0.2897

This custom scenario, which gives varying emphasis across different moments, resulted in the best AUC and F1 score among all the configurations tested, suggesting that a more nuanced combination of all moments leads to the most effective performance.

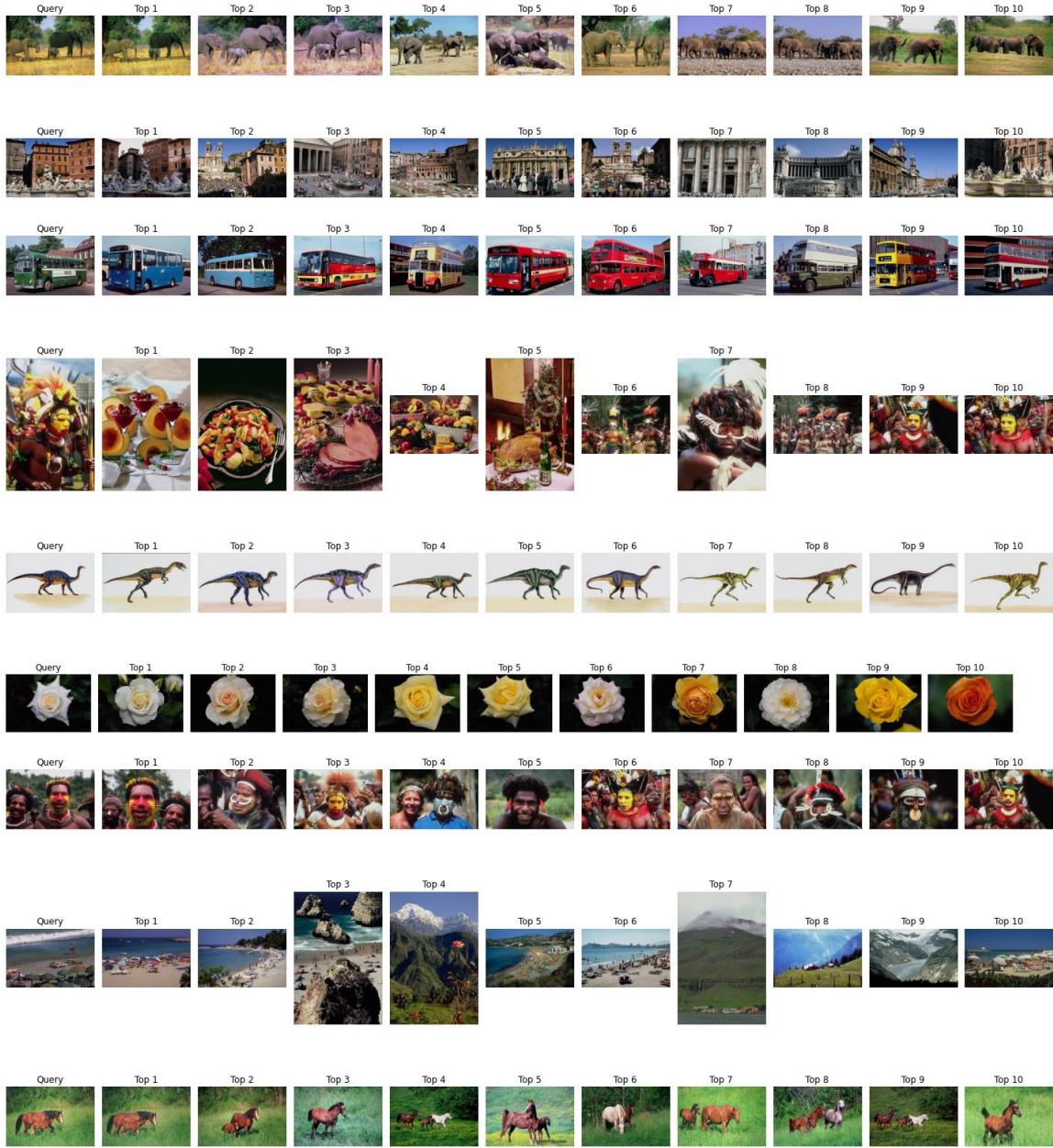
Comment on task 2 & 3 :

A CBIR system using color moments seems to be more sensitive to threshold adjustments, and fine-tuning this parameter can lead to significant changes in performance. The histogram-based approach seems to be more sensitive to the number of bins, which affects its ability to discriminate between different colors in the images.

3.3 Task 4: Improvement Using CNN

Each query with its top 10 similar images:

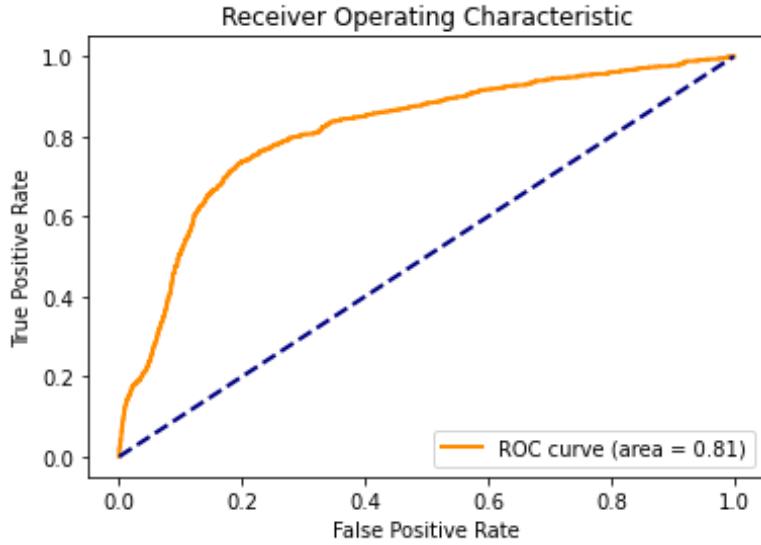




Appling measures for the top 100 similar images to the query.

```
N_QUERIES = 10
TOP_K = 100
TIME_RUNNING = 0.01 - 0.04
```

AUC was 0.81 as shown below:



The AUC of 0.81 indicates a good level of separability; the model is able to distinguish between the positive class (similar images) and the negative class (dissimilar images) with a reasonable degree of accuracy.

Measures for each query and average measures for all queries :

```

Query 1: Precision = 0.5500, Recall = 0.5500, F1-score = 0.5500
1/1 [=====] - 0s 117ms/step
Query 2: Precision = 0.4000, Recall = 0.4000, F1-score = 0.4000
1/1 [=====] - 0s 113ms/step
Query 3: Precision = 0.5000, Recall = 0.5000, F1-score = 0.5000
1/1 [=====] - 0s 112ms/step
Query 4: Precision = 0.3800, Recall = 0.3800, F1-score = 0.3800
1/1 [=====] - 0s 111ms/step
Query 5: Precision = 0.6800, Recall = 0.6800, F1-score = 0.6800
1/1 [=====] - 0s 119ms/step
Query 6: Precision = 0.3800, Recall = 0.3800, F1-score = 0.3800
1/1 [=====] - 0s 113ms/step
Query 7: Precision = 0.3800, Recall = 0.3800, F1-score = 0.3800
1/1 [=====] - 0s 114ms/step
Query 8: Precision = 0.5800, Recall = 0.5800, F1-score = 0.5800
1/1 [=====] - 0s 131ms/step
Query 9: Precision = 0.3000, Recall = 0.3000, F1-score = 0.3000
1/1 [=====] - 0s 110ms/step
Query 10: Precision = 0.4500, Recall = 0.4500, F1-score = 0.4500
Average values for all thresholds for Query Index: 9

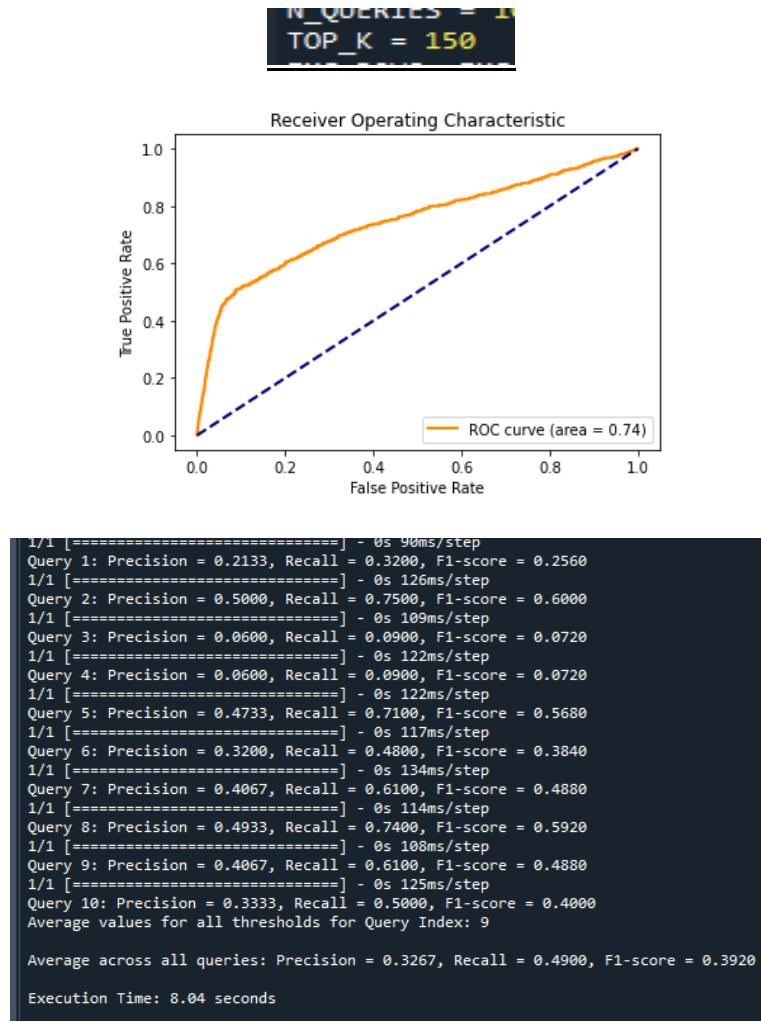
Average across all queries: Precision = 0.4600, Recall = 0.4600, F1-score = 0.4600

```

Execution Time: 7.88 seconds

The average precision across all queries is 0.4600, which suggests that when the model retrieves a set of images, about 46% of the retrieved images are relevant. The average recall (also 0.4600) indicates that the model retrieves about 46% of all relevant images in the dataset. The average F1-score, which balances precision and recall, is also 0.4600, showing that there is a balance between precision and recall in this model's performance.

Increasing top-K affects AUC and the whole measures decrease as shown:



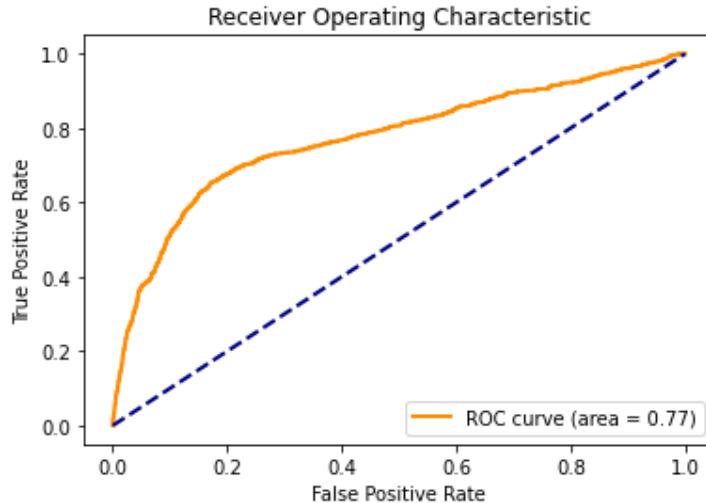
TOP_K = 150, AUC = 0.74:

- The AUC is slightly lower compared to TOP_K = 100, suggesting that increasing the number of top results considered introduces more false positives, thus reducing the model's ability to distinguish between similar and dissimilar images effectively.
- The precision drops to 0.3267, indicating that the inclusion of more images in the top results leads to more false positives.
- The recall increases to 0.4900, suggesting that while more relevant images are retrieved, the precision suffers.

- The F1-score is 0.3920, which is lower than the TOP_K = 100 scenario, reflecting the trade-off between precision and recall.

Decreasing top-k less than 100:

```
N_QUERIES = 10
TOP_K = 77
```



```
1/1 [=====] - 0s 154ms/step
Query 1: Precision = 0.8701, Recall = 0.6700, F1-score = 0.7571
1/1 [=====] - 0s 104ms/step
Query 2: Precision = 0.0649, Recall = 0.0500, F1-score = 0.0565
1/1 [=====] - 0s 110ms/step
Query 3: Precision = 0.4156, Recall = 0.3200, F1-score = 0.3616
1/1 [=====] - 0s 110ms/step
Query 4: Precision = 0.5455, Recall = 0.4200, F1-score = 0.4746
1/1 [=====] - 0s 120ms/step
Query 5: Precision = 0.2727, Recall = 0.2100, F1-score = 0.2373
1/1 [=====] - 0s 110ms/step
Query 6: Precision = 0.6234, Recall = 0.4800, F1-score = 0.5424
1/1 [=====] - 0s 109ms/step
Query 7: Precision = 0.5974, Recall = 0.4600, F1-score = 0.5198
1/1 [=====] - 0s 113ms/step
Query 8: Precision = 0.7403, Recall = 0.5700, F1-score = 0.6441
1/1 [=====] - 0s 113ms/step
Query 9: Precision = 0.5584, Recall = 0.4300, F1-score = 0.4859
1/1 [=====] - 0s 110ms/step
Query 10: Precision = 0.4935, Recall = 0.3800, F1-score = 0.4294
Average values for all thresholds for Query Index: 9

Average across all queries: Precision = 0.5182, Recall = 0.3990, F1-score = 0.4508

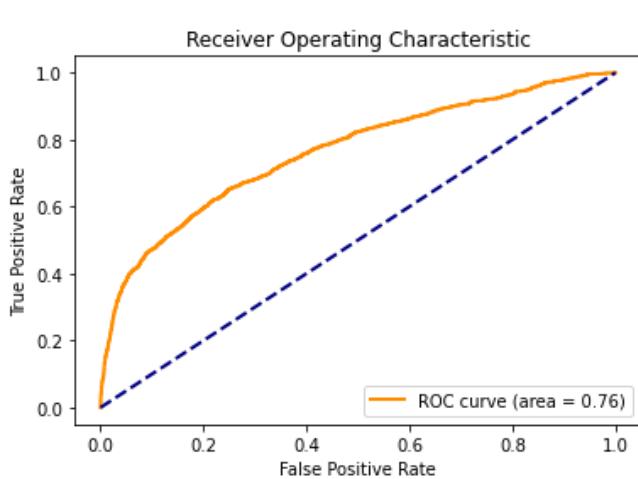
Execution Time: 8.64 seconds
```

TOP_K = 77, AUC = 0.77:

- This AUC is lower than the best performing TOP_K = 100 but higher than the TOP_K = 150, suggesting that there might be an optimal range around TOP_K = 100 where the model performs better.

- The average precision of 0.5182 is an improvement over both the TOP_K = 100 and 150 scenarios, suggesting that a smaller set of results might be more precise.
- The average recall drops to 0.3990, indicating fewer relevant images are retrieved compared to the TOP_K = 150 scenario.
- The F1-score of 0.4508 is slightly lower than the TOP_K = 100 scenario, which suggests a better balance between precision and recall was achieved at TOP_K = 100.

And when top-k is 50 , measures decreased:



```

Query 1: Precision = 0.3000, Recall = 0.1500, F1-score = 0.2000
1/1 [=====] - 0s 121ms/step
Query 2: Precision = 0.5000, Recall = 0.2500, F1-score = 0.3333
1/1 [=====] - 0s 120ms/step
Query 3: Precision = 0.3400, Recall = 0.1700, F1-score = 0.2267
1/1 [=====] - 0s 113ms/step
Query 4: Precision = 0.5000, Recall = 0.2500, F1-score = 0.3333
1/1 [=====] - 0s 124ms/step
Query 5: Precision = 0.5800, Recall = 0.2900, F1-score = 0.3867
1/1 [=====] - 0s 116ms/step
Query 6: Precision = 0.6200, Recall = 0.3100, F1-score = 0.4133
1/1 [=====] - 0s 113ms/step
Query 7: Precision = 0.7400, Recall = 0.3700, F1-score = 0.4933
1/1 [=====] - 0s 113ms/step
Query 8: Precision = 0.7400, Recall = 0.3700, F1-score = 0.4933
1/1 [=====] - 0s 114ms/step
Query 9: Precision = 0.3600, Recall = 0.1800, F1-score = 0.2400
1/1 [=====] - 0s 113ms/step
Query 10: Precision = 0.5600, Recall = 0.2800, F1-score = 0.3733
Average values for all thresholds for Query Index: 9

Average across all queries: Precision = 0.5240, Recall = 0.2620, F1-score = 0.3493
  
```

TOP K = 50, AUC = 0.76:

- The AUC is almost similar to TOP_K = 77, indicating consistent model performance in distinguishing between similar and dissimilar images.
- The precision increases to 0.5240, which is the highest among all scenarios, suggesting that a smaller set of top results yields more accurate predictions.
- However, the recall drops to 0.2620, which is the lowest among the scenarios, indicating that many relevant images are being missed.
- The F1-score of 0.3493 reflects the trade-off, showing that while the precision is high, the recall suffers significantly, resulting in a lower overall F1-score.