

# High level Content-based Image Retrieval

## I. INTRODUCTION

Nowadays with the increase of the computational power, the use of high level techniques to extract the information of images is increasingly common. It makes sense to explore algorithms of how they can be used to create useful and interesting applications.

This report consists in explaining the development of a project of a system to be possible to perform a high level Content-base image retrieval. The project is divided in two parts. The first part is to index a set of images. The second part consists in search in this set by providing an image and the system will retrieve the most similar images. Besides this, it is possible to search by a query of text to try to find images with the text of the query.

## II. ARCHITECTURE

The architecture of the system is divided in two big parts. The first part of the system will index the information that was extracted from the image. The second part is used to search in the indexers and retrieve similar images of the ones that are in the index. Also, there is the possibility to search by a text query and it will search by images with the specific words if detected. In the figure 1 it is possible to see the data flow diagram of the first part of the system. There are three components to extract information from the image. The first one uses mainly object recognition by using YOLO, the second one analyses if the image has text and tries to recognize the words in the image and the last one is dedicated to face recognizing. The humans give a lot of attention to the recognition of faces, so the system also, have that in consideration.

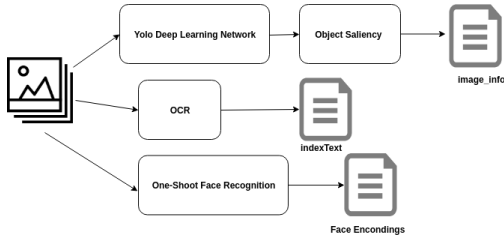


Fig. 1. Data Flow of Index

In the figure 2 is present the data flow that happens when there is a search. When a search is made the input can be a image or text. If it is an image will be recognize the objects of the image using YOLO and also it will try to do the recognition of faces and then retrieving the ranked results of the search. When the search is made by text will only search in the indexer that contains the information of the text in the images.

## III. OBJECT RECOGNITION

The main goal of the system is to retrieval images using high level techniques. To be able to do this is necessary to

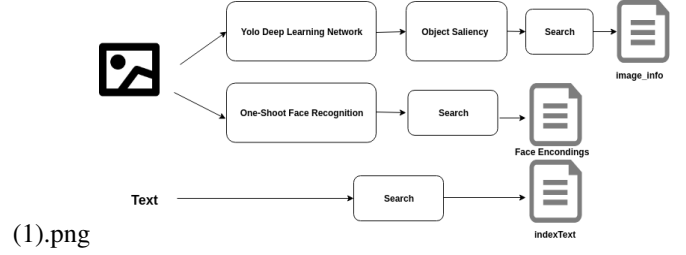


Fig. 2. Search Data Flow

use algorithms for object detection. The algorithm used in this project is YOLO. The goal of YOLO(You only look once) is to detect objects and it uses features learned by a deep convolutional neural network and for that reason the size of the input of the image is not important. The algorithm sends all the image at once through the deep learning network and it divides the image 13x13 grid of cells. The size of the cells depends of the size of the image and all the images were converted to 416\*416,so the cell is 32\*32. Each cell is responsible for predicting a set of boxes in the image. Lastly, for each predicted box there is a threshold for the minimum confidence and each box has a confidence, if the confidence is lower that the threshold the box is discarded. Also, if there are several boxes for the same object they are deleted, this technique is called non-maximum suppression. This algorithm can only detect the classes present in the dataset used to train the model. The used model in this project was trained in the COCO data set, so it has 80 different classes, so it is possible to detect 80 different objects.

To be able to develop the object detection using YOLO were followed mainly two different sites [1] and [2] The first thing done to use this algorithm was download the weighs of the pre-trained newtork's weighs and the configuration file of the network. The next step was configure the parameters to run the algorithm. The value of the confidence threshold used is 0.5 and for the non-maximum suppression is 0.4 and the the file with the names of the classes was loaded and also the network. After this with the path of the folder with images to index received by argument, it is possible to iterate through all the images of the folder and create a blob for each one and each blob is passed through the network to calculate the output vector. 3 This output vector has the probabilities of each class labels and also the information of the bounding boxes. There is a function to process this vector. [3] The process consists in iterate though the vector and save the results in different lists. So, in the end it is calculated lists to have the classes found in the image, the confidence of each class and the information of the bounding boxes. The last step is to do the Non Maximum Suppression, so it will be calculated the bounding boxes that have more than the confidence threshold and the bounding

box with overlap that are surrounding the same object will be deleted. This function does more things after this, however will not be explained in this section, because after the last step the object recognition using YOLO is totally done.

??

$$Y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix}$$

Fig. 3. Output vector

#### IV. VISUAL SALIENCY

With the approach described in III we have only the object count for each image, with no consideration of what objects are more relevant to it. Our approach is based on [9], where a way to compute visual saliency that is similar to our perception of it is proposed. Specifically, we use [4] from the reference image in Figure 4 we compute the saliency map seen in Figure 5. the shown saliency map is a matrix of number in the  $[0, 1]$  range representing the visual saliency. Looking at Figure 6, where the classes and the bounding boxes were already predicted, we can see that the handbag on the right side of the image doesn't contribute much to the picture but according to our previous model it mattered just as much as the fire extinguisher. Figure 7 shows the bounding boxes overlapped on the saliency map. To estimate the relevance of an object to the picture we calculate what percentage of the salience it's bounding box occupies. The exact way this is used later for similarity is described in IX

#### V. ONE-SHOT FACE RECOGNITION

The human visual system is highly sensitive to both detecting and recognizing faces. As such we take that into account in our retrieval system. The one-shot learning problem is relevant as we may only have one image of each person in our



Fig. 5. Saliency map

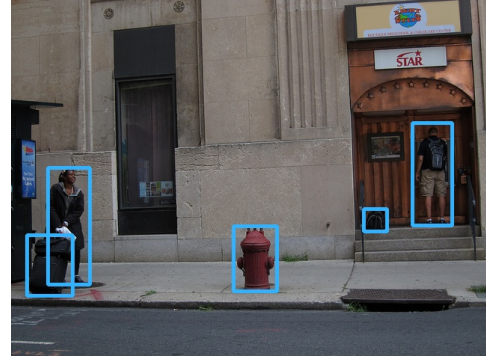


Fig. 6. Reference image with bounding boxes for predicted classes

database. This is achieved in our implementation with Siamese Networks, the goal of Siamese Networks is not to classify images but instead to learn the similarity of inputs. Figure 8 illustrates this. After the network is trained, embeddings can then be calculated, they are an high level representation of the image that are especially useful for face differentiation. In our model specifically our embeddings are composed of 120 units. To do all of this we use [5] module with trained weights, already fine tuned for this specific problem. Face recognition is done whenever a person is detected in the image and the embeddings are stored. When a new image query is given, it's embeddings are compared with existing ones and a distance measure is calculated, in our case the Euclidean distance. One interesting thing about this distance is we can



Fig. 4. Reference image

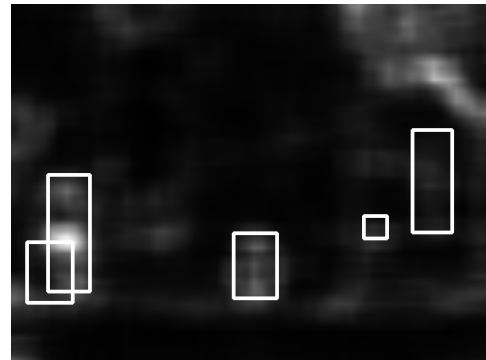


Fig. 7. Saliency map with bounding boxes for predicted classes

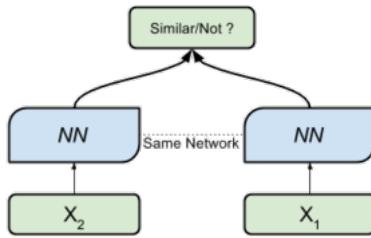


Fig. 8. Siamese Network model

not only retrieve a specific person but also give results ordered by face similarity.

## VI. TRANSFER LEARNING

In this section will be explaining the training process to train the model to detect trees. This part was quite difficult because there isn't a lot of information and it is necessary to install the drivers of the graphic card from Nvidia, the cuda and Cudnn. Furthermore, there is a lot of different configurations to train the models and it is not possible to know the output without testing. Most of the times that it is not a problem, however train a model takes many hours. Another problem is that the most promising configurations can't run in our laptop because the graphic card doesn't have enough memory. Despite, all this problems and difficulties a model was trained and the process will explain.

The first step necessary to train data is to have the data and this data has to be labeled. The images of trees with the labels were download from [6] and this data set from Google has around of 660 different classes and all the images are labeled. To download the images it was used a script created by [7] and after 1 hour it was stopped the script and the total of images download were 362. To train the model was used the Darknet. Darknet is an implementation in C and CUDA to support training data using the GPU. To train is necessary to create some configuration files. Two of them have the absolute path for the images to be used to train and the other one with images to be used to test. To create this two files it was used a script from the last source and it will choose 10 % of the images to be used to test and the rest of the images to train. Then, some changes were made in some files in the code C to be able to save more frequently the intermediate weights. The next step was to download the pre-trained model from darknet. The model used was darknet53.conv.74. The use of this pre-trained model will accelerate the process of training even though the large data set may not contain the object you are trying to detect. The last two files that were create was a file to give the specifications and the location of needed files to darknet and the other one with the name of the class that will be trained. The last step before starting the train is to configure the architecture that will be used. The base file used was the yolov3-tiny.cfg because of the constrains in the hardware. In this file was set that it will be used 24 batches and 8 subdivisions. This means that in each iteration will be used 24 images to update the parameters of neural network.

```
1099: 3.637688, 4.712093 avg, 0.001000 rate, 1.561469 seconds, 26376 Images
Loaded: 0.000838 seconds
Region 16 Avg IOU: 0.663184, Class: 0.955552, Obj: 0.010632, No Obj: 0.003994, .SR: 1.000000, .7SR: 0.250000, count: 4
Region 23 Avg IOU: 0.456580, Class: 0.970622, Obj: 0.000761, No Obj: 0.000725, .SR: 0.500000, .7SR: 0.000000, count: 2
Region 16 Avg IOU: 0.647252, Class: 0.939695, Obj: 0.025384, No Obj: 0.004038, .SR: 0.875000, .7SR: 0.250000, count: 8
Region 23 Avg IOU: 0.470347, Class: 0.935039, Obj: 0.018396, No Obj: 0.001019, .SR: 0.527778, .7SR: 0.055556, count: 30
Region 16 Avg IOU: 0.570514, Class: 0.959161, Obj: 0.014028, No Obj: 0.003502, .SR: 0.600000, .7SR: 0.000000, count: 5
Region 23 Avg IOU: .nan, Class: .nan, Obj: .nan, No Obj: 0.000597, .SR: .nan, .7SR: .nan, count: 0
Region 16 Avg IOU: 0.650923, Class: 0.939702, Obj: 0.012624, No Obj: 0.003711, .SR: 0.750000, .7SR: 0.500000, count: 4
Region 23 Avg IOU: 0.530774, Class: 0.925903, Obj: 0.020019, No Obj: 0.000803, .SR: 0.333333, .7SR: 0.000000, count: 3
Region 16 Avg IOU: 0.491248, Class: 0.959115, Obj: 0.018525, No Obj: 0.004123, .SR: 0.200000, .7SR: 0.200000, count: 5
Region 23 Avg IOU: 0.521246, Class: 0.933184, Obj: 0.017941, No Obj: 0.000740, .SR: 0.571429, .7SR: 0.000000, count: 7
Region 16 Avg IOU: 0.736763, Class: 0.960526, Obj: 0.014658, No Obj: 0.003595, .SR: 1.000000, .7SR: 0.000000, count: 5
Region 23 Avg IOU: .nan, Class: .nan, Obj: .nan, No Obj: 0.000402, .SR: .nan, .7SR: .nan, count: 0
Region 16 Avg IOU: 0.609830, Class: 0.954589, Obj: 0.020199, No Obj: 0.003872, .SR: 0.888889, .7SR: 0.000000, count: 9
Region 23 Avg IOU: 0.649878, Class: 0.940039, Obj: 0.004760, No Obj: 0.000680, .SR: 0.666667, .7SR: 0.333333, count: 3
Region 16 Avg IOU: 0.589012, Class: 0.958722, Obj: 0.020184, No Obj: 0.003538, .SR: 0.769231, .7SR: 0.000000, count: 13
Region 23 Avg IOU: 0.520708, Class: 0.948516, Obj: 0.013034, No Obj: 0.000987, .SR: 0.600000, .7SR: 0.000000, count: 5
```

Fig. 9. Print of one iteration when training



Fig. 10. Prediction using trained model

The subdivision will divide this iteration in 8, this is use full because sometimes the GPU is not enough capable to handle all the images at once, so by creating subdivisions there isn't that problem. Also, another parameters were changed to try to have a better output in the lowest amount of time possible, like for example the learning rate. The last step is to start the training and wait several hours. A rule that sometimes is done is to wait for the stabilization of the loss. The figure 9 is a print of one iteration of the process when we training our model. In the figure 13 is the command used to train the mode and in the figure 12 is the command to make a prediction using the model. After more or less 10 hours of training the model it is possible to see some results. In the figures 10 and 11 are some examples using the test set.

9  
10  
11  
12  
13

To put the detection of the trees in our project basically every image is passed trough the trained network to detect trees. The result final of object detection in our image is the union of the results of this network with the YOLO network. The trained network is only capable of detect trees, so there is no intersection of results.

## VII. OCR

To increase the capabilities of the system was decided to add OCR to our system. OCR stands for optical character



Fig. 11. Prediction 2 using trained model

```

rute@rute-X550JK:~/Documents/cadelras/vc/progs/darknet$ ./darknet detector test obj.data cfg/yolov3-tiny.cfg backup/yolov3-tiny_5000.weights data/treed.jpg
layer      filters      size  output              BFLOPs
 0 conv     16  3 x 3 / 1    416 x 416 x 3    -> 416 x 416 x 16  0.150 BFLOPs
 1 max       2  2 x 2 / 1    416 x 416 x 16    -> 208 x 208 x 16  0.399 BFLOPs
 2 conv     32  3 x 3 / 1    208 x 208 x 16    -> 208 x 208 x 32  0.399 BFLOPs
 3 max       2  2 x 2 / 1    208 x 208 x 32    -> 104 x 104 x 32  0.399 BFLOPs
 4 conv     64  3 x 3 / 1    104 x 104 x 32    -> 104 x 104 x 64  0.399 BFLOPs
 5 max       2  2 x 2 / 1    104 x 104 x 64    -> 52 x 52 x 64  0.399 BFLOPs
 6 conv    128  3 x 3 / 1     52 x 52 x 64    -> 52 x 52 x 128  0.399 BFLOPs
 7 max       2  2 x 2 / 1     52 x 52 x 128    -> 26 x 26 x 128  0.399 BFLOPs
 8 conv    256  3 x 3 / 1     26 x 26 x 128    -> 26 x 26 x 256  0.399 BFLOPs
 9 max       2  2 x 2 / 1     26 x 26 x 256    -> 13 x 13 x 256  0.399 BFLOPs
10 conv   512  3 x 3 / 1     13 x 13 x 256    -> 13 x 13 x 512  1.595 BFLOPs
11 max       2  2 x 2 / 1     13 x 13 x 512    -> 13 x 13 x 512  1.595 BFLOPs
12 conv  1024  3 x 3 / 1     13 x 13 x 512    -> 13 x 13 x1024  1.595 BFLOPs
13 conv   256  1 x 1 / 1     13 x 13 x1024    -> 13 x 13 x 256  0.089 BFLOPs
14 conv   128  3 x 3 / 1     13 x 13 x 256    -> 13 x 13 x 128  0.399 BFLOPs
15 conv    18  1 x 1 / 1     13 x 13 x 128    -> 13 x 13 x 18  0.003 BFLOPs
16 yolo
17 route 13
18 conv   128  1 x 1 / 1     13 x 13 x 128    -> 13 x 13 x 128  0.011 BFLOPs
19 upsample 2x 13 x 13 x 128 -> 26 x 26 x 128
20 route 18 8
21 conv   256  3 x 3 / 1     26 x 26 x 128    -> 26 x 26 x 256  1.196 BFLOPs
22 conv    18  1 x 1 / 1     26 x 26 x 256    -> 26 x 26 x 18  0.006 BFLOPs
23 yolo
loading weights from backup/yolov3-tiny_5000.weights...done!
data/treed.jpg Predicted in 0.007339 seconds.
tree: 69%

```

Fig. 12. Commando to make a prediction using the trained model

recognition, so after doing the object recognizing using YOLO, the detection of text and the recognition will be performed. The web page follow to help us develop this part was [8]. The Tesseract 4 was installed to be able to do this and this is the software that will recognize the characters.

The pipeline of OCR VII starts by passing the image through a deep learning model. This model is the EAST deep learning text detector. The output will be the bounding box coordinates of the detected text. This output will be used in the Tesseract to recognize the characters. In our project the last step of the pipeline is not used, because the recognized text is only saved in the inverted index created to text and the last step consists in show to the user the output of both algorithms.

To handle the recognizing of the text was created a file `text_recognizer.py` and from the `indexer.py` it is called the function `verifyText()` from that file created and the argument is the image and the return is a list with the words detected in the

```

rute@rute-X550JK:~/Documents/cadelras/vc/progs/darknet$ ./darknet detector train
obj.data cfg/yolov3-tiny.cfg backup/yolov3-tiny_5000.weights > /home/rute/Docum
ents/cadelras/vc/vc1819-76505-76532/project/proj1/training/log.txt

```

Fig. 13. Commando to train the model

given image. This function is called after the object detection is done using Yolo. After receiving the return of the function will be verified if it was detected text in the image and if it is true all the special characters will be removed of the words and it will be applied stemming. The stemming consists in normalizing the worlds for example if it is found the word fairly, it will be converted to fair. This stemming used is only prepare to support English words and the library used was 'nltk'. After this the last step is to count how many times each word appears in the image to save also in the inverted index that information.

## VIII. INDEXING

As show in II The indexing is divided in 3 parts: the ocr, the facial recognition and the object detection indexes. Stored in `face_encodings`, `imgs_info` and `index text`.

### A. Face Encodings

The Face Encodings file contains a dictionary of image name to embedding list with the embeddings of all the faces found in the image

### B. Images Info

The Image Info file contains a dictionary, where each image name is a key and points to an inner dictionary where we have the classes found in the images as keys along with the information about their count in the image and their saliency score

### C. Index Text

The index text is a dict with words as keys, pointing to the images they appear in, and in what number

## IX. PUTTING IT ALL TOGETHER

To weight all of these factors together we set guidelines we should follow

- Having the same classes as those in the query image is more important than having the same number of instances of a specific class
- The more instances of an object, the less important its difference to the query image becomes
- If a person is recognized their images should be given top priority, ranked then according to other factors

The weighting scheme we came up with is divided into 3 groups:

### A. Object weights

Object weight take into account the number of items in the query and reference image ( $y$  and  $\hat{y}$ ) and the saliency score, noted as  $rel(x)$ . A table with the object weights can be seen in IX-A



TABLE I  
OBJECT WEIGHT CALCULATION

| $y$   | $\hat{y}$ | cost  |
|-------|-----------|---|
| 0     | 0         | 0   |
| 0     | $> 1$     | $\log(1 + \hat{y}) * e^{rel(\hat{y})}$                    |
| $> 1$ | 0         | $1 + \log(y) * e^{rel(y)}$                                |
| $> 1$ | $> 1$     | $ \log(y) - \log(\hat{y})  * e^{ rel(y) - rel(\hat{y}) }$ |

TABLE II  
FACE DISCOUNT CALCULATION

|               |                                    |
|---------------|------------------------------------|
| $d \leq 6$    | -10                                |
| $0.6 < d < 1$ | $\min(\log((d - 0.6) * 2.5), -10)$ |
| $d > 1$       | 0                                  |

### B. Face discounts

When faces are detected, the cost of the picture is brought down, extremely so if it's a match ( $d \leq 0.6$ ). The face discount calculations can be seen in IX-B and is based on the distance between embeddings,  $d$ .

### C. OCR Weight

To perform the raking when it is made a search by text it is created a dictionary to save the results. The key of the dictionary is the path of the images and the value the score. The score of each image is calculated by summing the times of the word of the query appears on the image. In the end the results are normalized by dividing the score by the sum of all the scores.

## X. RESULTS

Here we show that the more advanced features of our retrieval system work

### A. Recognizing faces

For Figure 14 we got results X-A X-A X-A

### B. Recognizing text

A search for "Nokia" gets us the results seen in Figure X-B

### C. Recognizing objects

With 15, we got the following top 3 results



Fig. 14. Face Query



Fig. 15. Object Query





## XI. USAGE

To index a set of images it is necessary to run the indexing.py file with the flat -path and to the front of the flag is place the path.

If in the same directory of the file it is the directory with the images and the name of that directory is album to run the file it is done the next command.

See the following command :

```
$ python3 indexing.py --path album/
```

To search it is similar, however the name of the file is search.py and in from of the argument path it is placed the directory of the image.

See the following command :

```
$ python3 search.py --path image.jpg
```

If it is to search by text the argument is -text and at front is placed the query.

See the following command :

```
$ python3 search.py --text "Selling_car"
```

## REFERENCES

- [1] <https://blog.paperspace.com/how-to-implement-a-yolo-object-detector-in-pytorch/>.
- [2] <https://www.learnopencv.com/deep-learning-based-object-detection-using-yolov3-with-opencv-python-c/>.
- [3] <https://medium.com/machine-learning-bites/deeplearning-series-object-detection-and-localization-yolo-algorithm-r-cnn-71d4dfd07d5f>.
- [4] URL: <https://docs.opencv.org/3.0-beta/modules/saliency/doc/saliency.html>.
- [5] URL: [https://github.com/ageitgey/face\\_recognition](https://github.com/ageitgey/face_recognition).
- [6] <https://storage.googleapis.com/openimages/web/visualizer/index.html?set=train&c=%2Fm%2F0mkg>.
- [7] <https://www.learnopencv.com/training-yolov3-deep-learning-based-custom-object-detector/>.
- [8] <https://www.pyimagesearch.com/2018/09/17/opencv-ocr-and-text-recognition-with-tesseract/>.
- [9] Xiaodi Hou and Liqing Zhang. “Saliency detection: A spectral residual approach”. In: *Computer Vision and Pattern Recognition, 2007. CVPR’07. IEEE Conference on*. IEEE. 2007, pp. 1–8.