

# Vision and Cognitive Systems Final Project

Fabio Polito  
230635@studenti.unimore.it

Giordano Costi  
226934@studenti.unimore.it

Stefano Carretti  
227250@studenti.unimore.it

University of Modena and Reggio Emilia

## I. INTRODUCTION

In this paper, we present a method to detect and identify paintings starting from a video taken inside Galleria Estensi, Modena.

Each frame is processed with image processing techniques in order to localize paintings, rectify distortions, and fetch from the database the corresponding work of art.

At the same time, an artificial neural network (YoloV3) detects people, which are localized inside a room of the museum.

## II. RELATED WORKS

To remove camera noise but maintains the contours a Bilateral filter is used. C. Tomasi and R. Manduchi, "Bilateral filtering for gray and color images," Sixth International Conference on Computer Vision (IEEE Cat. No.98CH36271), Bombay, India, 1998, pp. 839-846, doi: 10.1109/ICCV.1998.710815.

OTSU is an automatic thresholding method widely used when the numbers of pixels in each class are close to each other. J. Zhang and J. Hu, "Image Segmentation Based on 2D Otsu Method with Histogram Analysis," 2008 International Conference on Computer Science and Software Engineering, Hubei, 2008, pp. 105-108, doi: 10.1109/CSSE.2008.206.

Opencv Find Counturs function is been used to find counturs in a previously modified frame in order to detect painting Suzuki, S., and Be, K. (1985). Topological structural analysis of digitized binary images by border following. Computer Vision, Graphics, and Image Processing 30, 3246. doi:10.1016/0734-189X(85)90016-7.

The functions approxPolyDP approximate a curve or a polygon with another curve/polygon with less vertices so that the distance between them is less or equal to the specified precision. It uses the Douglas-Peucker algorithm

ORB E. Rublee, V. Rabaud, K. Konolige and G. Bradski, "ORB: An efficient alternative to SIFT or SURF," 2011 International Conference on Computer Vision, Barcelona, 2011, pp. 2564-2571, doi: 10.1109/ICCV.2011.6126544.

YoloV3 is a famous neural network for fast object detection, will be used in this paper for people detection Redmon, Joseph and Ali Farhadi. YOLOv3: An Incremental Improvement. ArXiv abs/1804.02767 (2018): n. pag.

## III. PAINTING DETECTION

All frames are extracted in sequence from the video and processed independently. The followed pipeline is explained afterwards.

### A. Preprocessing

First of all, each frame is converted into black and white and processed with a bilateral filter to remove noise while preserving at the same time the edges.

After that, we apply otsu threshold with the aim to separate the pixels into background and paintings, due to their chromatic difference.

### B. Bounding Box detection

Using the function findContours of OpenCV, we obtain the outlines of the objects in the foreground, among which there will also be the paintings that we are looking for.

We then create the bounding box containing this contours and to eliminate the rectangles identified inside the paintings, we keep only the external ones and eliminate those contained within others.

Among the remaining bounding boxes, those that are not judged as paintings by an SVM model, will get discarded.

### C. SVM

To classify the ROI proposed by the previous pipeline, an SVM model is been trained.

The algorithm takes as input the concatenation of the three histograms, one for each color plane.

The training dataset is composed of 1025 instances taken from bounding boxes derivate from the videos inside the museum and manually labeled. 409 of the samples are labeled as holding a painting and 604 as inaccurate bounding box.

A radial basis function kernel is been exploited for the classification. The model returns False if the rectangle doesn't contain a painting and True if it does.

### D. Precision boosting and paintings segmentation

To obtain a better segmentation within each bounding box we apply a further refinement of the images.

Observing the low light in the videos that were provided to us, the paintings

are significantly darker compared to their frame and do not correspond to the brightness of the paintings in the database. This leads to poor precision in the retrieval step. To cope with this problem, the brightness component is increased for each previously found bounding box.

Afterward, by transforming its format from BGR to HSV and then applying the otsu threshold again, we are able to obtain a more precise distinction between painting and frame/background, which will then be used during the retrieval and rectification steps.

#### IV. PAINTING RETRIEVAL

For the retrieval of the paintings situated in the database we use an approach based on feature detection algorithms. After consulting a paper that performs a comparative analysis between the most well known algorithms [1] to get a general idea of the strengths and weaknesses of the different methods available to us, we carried out some experiments focusing on SIFT, AKAZE and ORB.

The results obtained made us opt for ORB, because overall it gave us more precise results than AKAZE and, unlike SIFT, its free of charge, therefore usable without fees in a possible commercial application.

To save time, the key points of the paintings in the database have been previously calculated and stored using the pickle module [link pickle documentation].

It implements binary protocols for serializing and de-serializing a Python object structure.

This data is loaded only once during the launch of the program. The key points are computed from the bounding boxes detected by the previously described pipeline and afterward, to determine the best matches, the ratio test proposed by D. Lowe in the SIFT paper is performed[Reference to the sift paper].

This measure is obtained by comparing the distance of the closest neighbor to that of the second-closest neighbor.

This measure performs well because correct matches need to have the closest neighbor significantly closer than the closest incorrect match to achieve reliable matching.

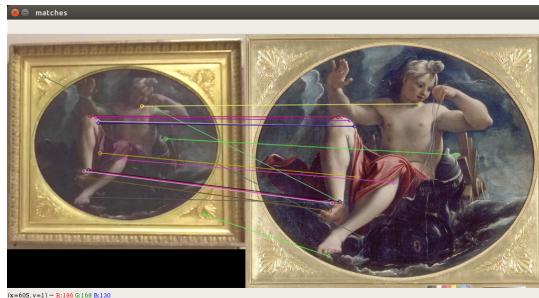


Fig. 1. Orb matches

For false matches, there will likely be a number of other false matches within similar distances due to the high dimensionality of the feature space.

In our implementation we reject all matches in which the distance ratio is greater than 0.75.

This allows us to keep the number of correctly retrieved paintings still high but at the same time to decrease the number of false positive. A ranking with the 5 best matches found is then created and saved to a CSV file to show the results.

To understand if the painting is actually recognized among those in the DB, the average of the key points matched among the best 5 ranked is calculated and, if the first one differs from it for more than significant value, it is considered as correct and shown on the interface.

#### V. PAINTING RECTIFICATION

##### A. Four points transform

On the contour found with the techniques described in **Precision boosting and paintings segmentation** is applied the function approxPolyDP from OpenCV to approximate it to a polygonal curve.

If the shape returned has four vertices we can assume that the process has found a rectangular painting.

Given the four points, we are able to estimate the homography and apply the transformation to rectify the painting.

To calculate the aspect ratio for the projected rectangle we use an implementation based on this paper [] which derives the equations assuming a pinhole camera model.

This pipeline is applied to the paintings that don't get a match on the database.

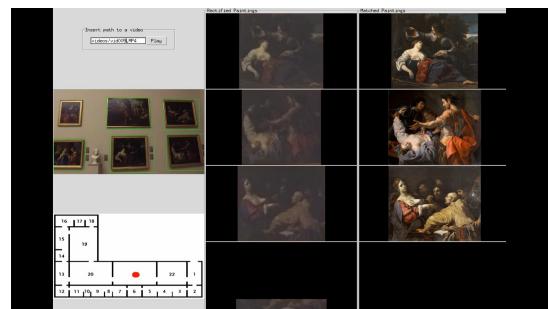


Fig. 2. Painting detection and rectification

##### B. alignImages

Not all the paintings inside the museum have a rectangular shape. This means that the approximation found by approx-PolyDP does not consist of 4 vertices, making the method just explained impractical.

The approach that is used in this case is based on the common keypoints found between the distorted image and the corresponding match in the database.

Through ORB feature matching algorithm the key points are computed from each bounding boxes detected by the previous pipeline and subsequently, to determine the best matches between it and the images of the database, the ratio test is used.

From them we calculate the homography matrix and to avoid mismatches, the RANSAC algorithm is exploited.

To obtain a better result also the inverse warping algorithm is utilized.

#### VI. PEOPLE DETECTION

A neural network (YoloV3) is used for people detection. At inference time each video frame is passed through the network that find all bounding boxes containing one of the objects in our classes list.

The weights for the network are obtained from an already trained network on COCO[ reference to COCO paper], a famous dataset containing 80 different classes.

In our case the only wanted class is the person one, for ease of use the network has not been modified, but from its output will be deleted all the classes with an id different from 0 that is the id of the person class.

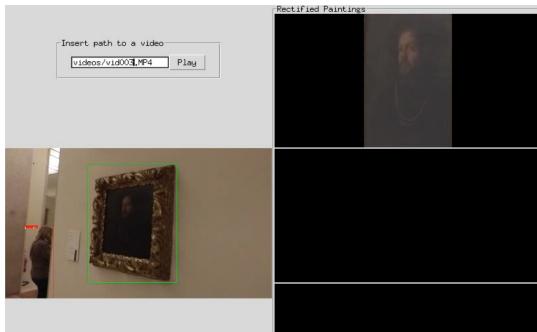


Fig. 3. People detection

A little problem has arisen due to the high number of paintings representing persons, in fact the network detect them also inside the paintings, to prevent these false positive we added a new control on the pixel position in order to cut out all the people detected inside a bounding box previously classified as a painting.



Fig. 4. People detection error



Fig. 5. People detection error fixed

## VII. PEOPLE LOCALIZATION

Starting from a previously detected and retrieved painting from the database, we search the corresponding room in csv file, then a little red dot is printed on the museum image map.

With this method we assume that all the painting detected and the persons are in the same room. Moreover the localization will work only if the painting matched is considered safe that means the painting need to have a good level of matches, otherwise localization is not active.

## VIII. PAINTING REPLACEMENT IN THE 3D MODEL

The pipeline followed to accomplish this task is similar to the one explained before with a little difference. we start directly with the hsv version of the image taken as input (a screenshot from the 3d model); then we apply Otsu threshold, followed by noise removal thanks to opening or closing process, keeping the one that maximizes the number of contours found. At this point we loop over each contour and find an approximation of itself with approxPolyDP; if the approximation has a shape described by 4 vertices, we estimate keypoints with orb and fetch the corresponding image from the database. Once this is done, if the image fetched is judged as a good match, we align that image with respect to the screenshots image plane and superimpose the result on the input image.



Fig. 6. 3dmodel



Fig. 7. 3d model 2

What is worth to mention is that the alignment function that we use in this process is the same of the main pipeline of painting rectification. this function in fact takes as input two images and project the first one on the plane of the second. if the second image is a painting fetched from the database, it means that the first image is rectified; if instead the second image is a painting coming from a frame/screenshot,

it means that we want to distort an image in order to match the perspective of that painting.

## IX. METRICS AND PRECISION

To calculate the precision of the presented method we took almost 10 random frame for 6 different videos and labeled all the paintings by hand using a tool for labeling (normally used for yolo labels). Then with a custom code we calculated precision recall comparing hand made labels with automatically generated ones.

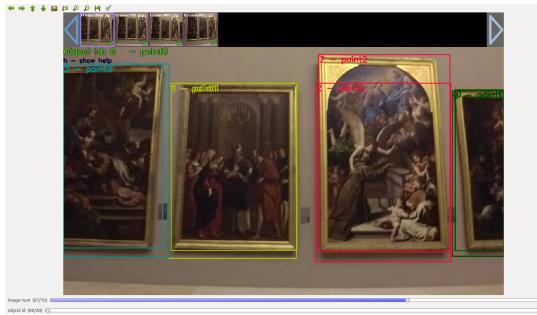


Fig. 8. labeling

TABLE I  
METRICS TABLE

Table Head	Table Column Head		
	Table column subhead	Subhead	Subhead
copy	More table copy <sup>a</sup>		

<sup>a</sup>Sample of a Table footnote.

## DISCUSSIONS

L'approccio seguito si fonda sull'assunzione che nel video in questione sia presente una buona differenza cromatica tra background e paintings. Nel caso questa condizione non sia verificata, come nel caso di una inquadratura più ampia comprendente non solo quadri e parete ma anche corridoio ad esempio, otsu non riesce a dividere correttamente quadri da background e le performance della detection calano. Cronologia degli approcci testati inizialmente si pensato di utilizzare Canny, preceduto da una rimozione iniziale di rumore prodotta dal bilateral filter, per trovare i bordi su cui applicare findContours. I risultati non erano del tutto soddisfacenti in quanto spesso i contorni prodotti da canny non individuavano con precisione il quadro. Abbiamo quindi optato per una soluzione basata sul thresholding dell'immagine RGB, che come già spiegato in precedenza seguita da un raffinamento effettuato a livello HSV. Si è anche testato un metodo basato solo su immagini HSV, che ha rivelato essere molto performante in certe situazioni (in ogni frame vi è una chiara distinzione tra quadro e background) ma meno stabile rispetto al metodo precedentemente esposto.

## REFERENCES

Please number citations consecutively within brackets [1]. The sentence punctuation follows the bracket [2]. Refer simply to the reference number, as in [3]—do not use “Ref. [3]” or “reference [3]” except at the beginning of a sentence: “Reference [3] was the first . . .”

Number footnotes separately in superscripts. Place the actual footnote at the bottom of the column in which it was cited. Do not put footnotes in the abstract or reference list. Use letters for table footnotes.

Unless there are six authors or more give all authors' names; do not use “et al.”. Papers that have not been published, even if they have been submitted for publication, should be cited as “unpublished” [4]. Papers that have been accepted for publication should be cited as “in press” [5]. Capitalize only the first word in a paper title, except for proper nouns and element symbols.

For papers published in translation journals, please give the English citation first, followed by the original foreign-language citation [6].

## REFERENCES

- [1] Tareen, Shaharyar Ahmed Khan, and Zahra Saleem. "A comparative analysis of sift, surf, kaze, akaze, orb, and brisk." 2018 International conference on computing, mathematics and engineering technologies (iCoMET). IEEE, 2018.
- [2] J. Clerk Maxwell, A Treatise on Electricity and Magnetism, 3rd ed., vol. 2. Oxford: Clarendon, 1892, pp.68–73.
- [3] I. S. Jacobs and C. P. Bean, “Fine particles, thin films and exchange anisotropy,” in Magnetism, vol. III, G. T. Rado and H. Suhl, Eds. New York: Academic, 1963, pp. 271–350.
- [4] K. Elissa, “Title of paper if known,” unpublished.
- [5] R. Nicole, “Title of paper with only first word capitalized,” J. Name Stand. Abbrev., in press.
- [6] Y. Yorozu, M. Hirano, K. Oka, and Y. Tagawa, “Electron spectroscopy studies on magneto-optical media and plastic substrate interface,” IEEE Transl. J. Magn. Japan, vol. 2, pp. 740–741, August 1987 [Digests 9th Annual Conf. Magnetics Japan, p. 301, 1982].
- [7] M. Young, The Technical Writer’s Handbook. Mill Valley, CA: University Science, 1989.