

CSC420 Project Report

Hirotaka Ishihara Yichun Zhang

December 2019

1 Introduction

In this project, we are going to perform analysis on videos, in particular, news broadcast. There are various tasks associated with the video analysis, such as shot detection, feature detection and recognition, and feature tracking between frames. Video processing has been a popular area in computer vision for a long time. In the project, we are going to employ novel ideas from the past few years to complete our analysis.

2 Shot Detection

Shot detection is crucial in video processing, it is a useful helper for other computer vision operations like face detection in consecutive frames[10]. I will perform three shot detection algorithms in 2.1 on the given video clips [5].

2.1 Method

The first two algorithms: Average Intensity Measurement 2.1.2 and Color Histogram Comparison 2.1.3, are very widely used in shot detection. They in general compares the color distribution of two images, which is robust to camera motion. The third method 2.1.4 is a relatively novel measurement, instead of collecting pixel color information, it focuses on object features.



Figure 1: case where traditional color histogram doesn't work in shot detection

2.1.1 Adaptive Thresholding

Before describing the details of the methods, I would talk about out threshold policy. In this project, the threshold policy we choose is called Adaptive Thresholding [14]. There are two hyperparameters in this thresholding, W denotes the window size and by default is an odd integer. T_c is a constant. For each dissimilarity score at time slot t , the threshold is

$$m_t = \max(\mu_{left} + T_c \sqrt{\sigma_{left}}, \mu_{right} + T_c \sqrt{\sigma_{right}})$$

$\mu_{left}, \sigma_{left}$ is calculated from the dissimilarity values within interval $[d - W/2, d]$, and the similar manner for $\mu_{right}, \sigma_{right}$.

2.1.2 Mean Intensity Measurement (MIM)

This method is described briefly in [13], which computes the average intensity of each image channel, and sum the difference over all channels with the values of the next frame.

2.1.3 Color Histogram Comparison (CHC)

The traditional color histogram described in [13] takes intensity values at gray-scale, and calcutes the entire image intensities into one histogram. While we think this may not be ro-



Figure 2: the red box moves between frames

bust, Figure 1 is the case. If we take color histogram of these two images, the difference would zero. Our method slightly changes the traditional one. We first take color histogram of all three channels. Instead of taking the entire image size as the input, we make a small patch with size 16 to gather color information.

2.1.4 Edge Change Ratio (ECR)

The edge change ratio [11] is a method focuses more on the feature information within the image. Figure 2 shows a case that a box moves between frames but within a certain boundary. ECR tracks how edge features move between frames. This needs pre-processing on two frames.

Implementation

- We first calculate edge features (our implement applies canny edge with $\sigma = 5$) of both images, denoting both the image I_t, I_{t+1} , and processed canny edges C_t, C_{t+1} . Here we used OpenCV Canny Edge Detection functons.[2]
- Then apply dilation on both edge feature D_t^-, D_{t+1}^- , then we calculate the inverse

$$D_t = 1 - D_t^-$$

$$D_{t+1} = 1 - D_{t+1}^-$$

- Take binary operations between C_t, D_{t+1} and C_{t+1}, D_t . A binary operation is defined as $binary(A, B) = A * B$ when A, B

only have values zero and one.

$$EC_{out} = \frac{binary(C_t, D_{t+1})}{C_t}$$

$$EC_{in} = \frac{binary(C_{t+1}, D_t)}{C_t + 1}$$

- The final edge change ratio is

$$ECR = max(EC_{in}, EC_{out})$$

2.2 Evaluation

[R:Recall, P:Precision, F1: F Score]
[C:Correct, M:Missed, F:False]

Clip	C	M	F	R	P	F1
#clip 1	1	0	2	1.00	0.33	0.50
#clip 2	8	0	2	1.00	0.80	0.89
#clip 3	4	2	9	0.67	0.40	0.50

Table 1: metrics of *Mean Intensity Measurement* over three clips

Clip	C	M	F	R	P	F1
#clip 1	1	0	2	1.00	0.33	0.50
#clip 2	8	0	7	1.00	0.53	0.70
#clip 3	4	2	8	0.67	0.33	0.44

Table 2: metrics of *Color Histogram Comparison* over three clips

Clip	C	M	F	R	P	F1
#clip 1	1	0	0	1.00	1.00	1.00
#clip 2	6	2	0	0.75	1.00	0.86
#clip 3	5	1	3	0.83	0.62	0.71

Table 3: metrics of *Edge Change Ratio* over three clips

The graph in Figure 3 4 5 shows the performance of three measurements respectively. The yellow vertical spans are the true video shot detections, where those thin lines indicate the hard cuts and wide spans indicate dissolve transitions. We also have the following tables 1 2 3 show the metrics of the three algorithm. The mathematical expressions for Recall, Precision and F measure (F1) [9] are

$$Recall = \frac{Correct}{Correct + Missed}$$

$$Precision = \frac{Correct}{Correct + False}$$

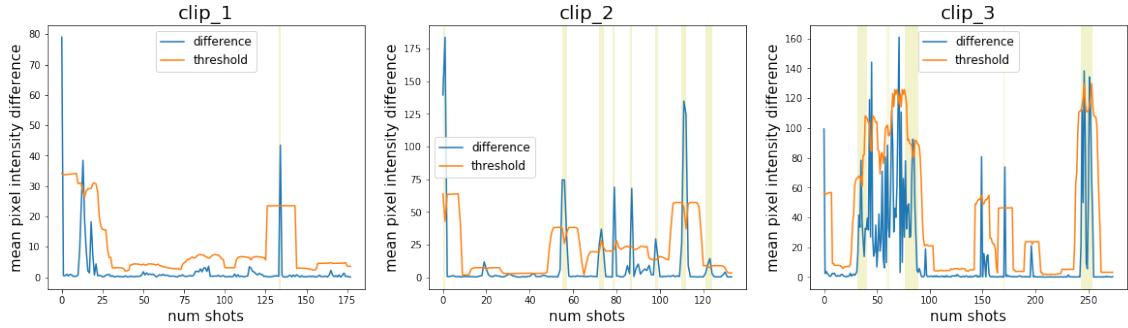


Figure 3: performance of Mean Intensity Measurement

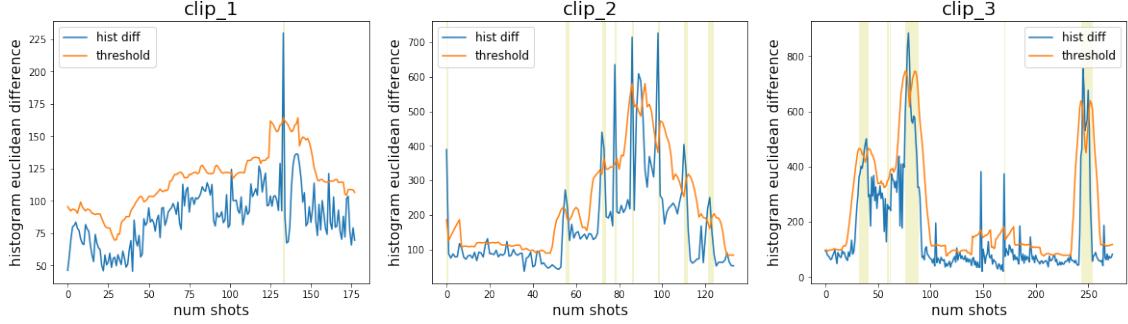


Figure 4: performance of Color Histogram Comparison

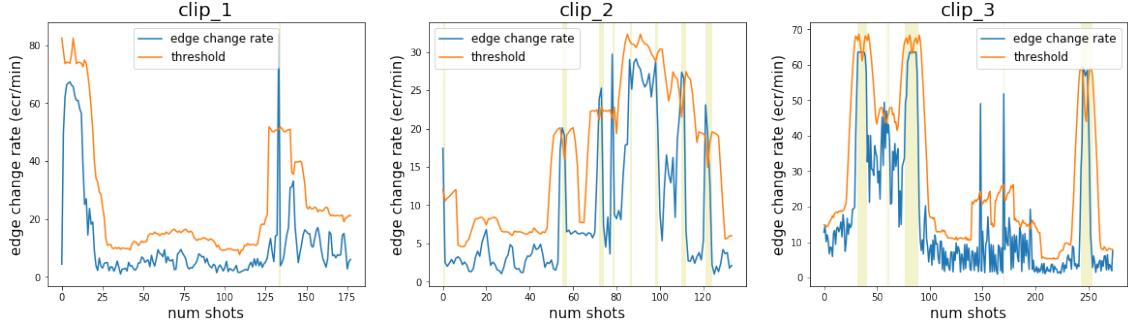


Figure 5: performance of Edge Change Ratio

$$F1 = \frac{2 \times Recall \times Precision}{Recall + Precision}$$

Intensity based algorithms tend to have a higher Recall and low Precision, since the given video clips[5] generally have frequently intensity changes, the algorithms capture both correct shot boundaries and intensity noise. On the other hand, the feature based algorithm (ECR) generally out performs the other two algorithms. Observe that the worst performance of the three clips[5] is the clip 3. This is because a rapid changing comic scenes (Marvel comic), the ECR is not robust to video clips that have low sampling rate, since in this case the

object has a large movement gap which causes the False detections.

3 Feature Detection

Feature detection field has been rapidly improved as the Deep Learning concept came out. Some deep learning algorithms now have extraordinary performances. However, there are still some widely used techniques that are non deep learning based which also have good performances. For face detection, we are going to use Haar Cascade Classifier[12], and for logo detection, we adopt a method that learnt from



Figure 6: missed face detection

lecture materials.

3.1 Face Detection

We use Haar Cascade Classifier[12] as our face detector. The algorithm first collects a large number of features, we need to figure out which features among all of these best determine a face. And this is done by using Adaboost, a collection of weak learner. All the features are stacked together with different amount, from less to more, grouped into the classifier. The algorithm is also fast at real time because of the Cascade concept, it immediately discard a sliding window in detecting when it fails at any position inside the classifier. Though in terms of the performance, CNN is better, but it is more efficient in training and has lighter model. So for simple task such as face detection, Haar Cascade has always been a good choice.

3.1.1 Face Detection Results

From Table 4, we can see that the Haar Cascade Classifier has better Precision than Recall. This is due to the Cascade property discussed previously. In contrast, the performance in the first two clips is better than the last one, especially for Recall measurement. This is because the angle that people facing at in clips 3 is not towards the camera like in Fig 6. So the Haar feature filters are not able to apply on these case variations. Even though in general people has a common front face feature, the style on each left

or right may vary. Unlike convolutional neural networks, which is of high freedom to perform optimal kernels (as weights to train), Haar features are manually determined, thus with a tighter limitation.

[R:Recall, P:Precision, F1: F Score]
[C:Correct, M:Missed, F:False]

Clip	C	M	F	R	P	F1
#clip 1	253	45	10	0.85	0.96	0.90
#clip 2	373	73	4	0.84	0.99	0.91
#clip 3	329	165	12	0.67	0.96	0.78

Table 4: metrics of *Haar Cascade Classifier* over three clips

3.2 Logo Detection

The method for Logo detection we used is applying the knowledge we learnt from class: Scale Invariant Feature Transform (SIFT) [7], match keypoints with RANSAC [4], and apply homography, mapping the Logo to the location in the given video clips [5].

The advantage of these sequence of operations over machine learning is that it only requires the minimum amount of data (1 is enough), and fast to apply, without the training process. While the drawback is also inevitable, which the SIFT [7] algorithm is very sensitive to noise. In addition, the algorithm is not prepared for multiple appearances of the same Logo. As shown in Fig 7, the Logo of NBC news is half-transparent, so the detection result is not optimal. While in a relatively stable, less noise environment Fig 8, the detection is more accurate.

By using this method, though we are unable to improve Recall, the it is able to improve Precision. This is done by pruning bad mapping points, such us negative width and height, or points that are mapped out of the image, we abandon the boundary box in these case.



Figure 7: false Logo detection & positive face detection



Figure 8: positive Logo detection & face detection

3.2.1 Logo Detection Evaluation

Logo	[R:Recall, P:Precision, F1: F Score]			[C:Correct, M:Missed, F:False]		
	C	M	F	R	P	F1
nbc news	92	76	10	0.55	0.90	0.68
nbc excl	96	10	9	0.91	0.91	0.91
the voice	57	0	0	1.0	1.0	1.0
clevver	77	36	0	0.68	1.0	0.81
flick	22	8	0	0.73	1.0	0.85

Table 5: metrics of *Logo Detection* over all Logos

The outcome varies between Logo types.

Like the one showed in Fig 8, the type of Logo with solid background (i.e Logo has high opacity) tend to be robust in detection. In this method, the transparency of the Logo introduces noise in position detection.

4 Face Tracking

Perform face tracking on video or real time camera helps gain both spatial and temporal information. The most difficult part lies on the assignment problem. Most algorithm perform assignment on consecutive frames. We have come up a simple method that performs tracking even between non-continues frames and shots.

4.1 Algorithm

Algorithm 1 Greedy Detection Tracking

```

1: initiate a MAX_HEAP
2: for object in 1st frame do
3:   MAX_SCORE = 1
4:   SECOND_MAX = 0
5: for object in 2nd frame do
6:   calculate similarity score
7:   if score > MAX_SCORE then
8:     SECOND_MAX = MAX_SCORE
9:     MAX_SCORE = score
10:  end if
11: end for
12: if RATIO > THRESHOLD then
13:   prune the pair
14: else
15:   push pair to heap
16: end if
17: end for
18: while MAX_HEAP not empty do
19:   obj1, obj2 = MAX_HEAP.pop
20:   if obj2 already belongs to a tracking then
21:     continue
22:   end if
23:   if neither obj1 nor obj2 belongs to a tracking then
24:     initiate new tracking with obj1 and obj2
25:   end if
26:   if obj1 belongs to a track then
27:     add obj2 to the same track
28:   end if
29: end while
30: LEFT_OVER = [object in 1st frame not get paired]
31: carry LEFT_OVER to the next round

```

The algorithm starts from the very first two consecutive frames, the process is similar to find the matching in SIFT matching[7].

4.2 Similarity Function for tracking

The essential parts of performing an accurate object tracking are detection and similarity measurement. In this subsection, we talk about how we set up the similarity function. In order to determine if two objects are very similar, we focus on both color distribution and feature information. To tackle color problem, we chose the color histogram difference measurement, over all three channels, this give a detailed color difference of two images. For feature comparison, we chose Histograms of Oriented Gradients (HOG).

To compute HOG, we first calculate both the

magnitude and the direction of the gradients:

$$|\nabla I(x, y)| = \sqrt{I_x^2 + I_y^2}$$

$$\theta(x, y) = \arctan(I_y/I_x)$$

We take a 2×2 cell with each cell size of 8×8 pixel blocks sliding over the image, quantizing the directions into 9 bins with 20 degree increment. The gradient with larger magnitude votes into the bin. Then normalize 4 of such blocks together (2×2). So, in total, our similarity function S with Color Histogram difference CH_D and Histogram of Oriented Gradients difference HOG_D is:

$$S(I_1, I_2) = \frac{1}{HOG_D(I_1, I_2) + CH_D(I_1, I_2) + 1}$$

4.3 Performance

The algorithm works better in situation where number of people in the same scenes is small. Our algorithm can tackle discontinuous shot tracking (e.g identify and track people in 1st shot scene and 3rd scene). For frames that involves many people, it usually swaps the tracking sequence, like Fig 9 10

5 Gender Classification

Automated gender recognition is important in many applications such as biometric and human computer interaction. As a result, classifying gender from face images becomes an important research area. We tried two methods to perform this task.

5.1 ResNet18

Deeper neural networks turn out to be more powerful theoretically, but in practice they suffer from the degradation problem and are harder



Figure 9: false Logo detection & positive face detection



Figure 10: positive Logo detection & face detection

to train due to vanishing gradient. Residual blocks are proposed to solve this problem by using shortcut connections.

We originally trained a residual network—Resnet18. This is an end-to-end approach, that we feed the model with an image and would directly obtain the class of it. Data are separated into training data and validation data in order to see the performance. Our neural net was trained with an Adam optimizer [6] and cross entropy loss. To fasten the training process, the model was trained with a minibatch of size 32. However, after 25 epochs, the training accuracy approached 100%, but validation accuracy is only around 80%, even after adding one dropout layer and applying weight decay. This implies that our model faces the problem of overfitting.

We thought this problem can be alleviated by employing more regularization techniques. However, training the neural net took a large amount of time. In my experiment, training 30 epochs requires more than 2 hours. Due to the time concerns, we didn't do lots of experiments

on it.

For further improvement, experiments such like cropping part of the image before feeding it to the neural net are considered. Also, representing the image in lower dimensional space may help since it decreases the number of parameters in the model.

5.2 Histogram of Oriented Gradient Feature

Another method we tried is simply extracting features from images and then feeding the features into a classifier. In real world, people may find that outline of faces contributes a lot to gender recognition from facial image. In order to represent the edge information, we choose the Histogram of Oriented Gradient[3] as the descriptor.

According to the paper [3], classic HOG features are obtained by normalizing the image, computing the gradient of each pixel, creating a histogram of edge directions and doing block normalization. One of the main contributions of the paper is that they use dense grid of cells, allowing overlapping both horizontally and vertically when sliding the windows. Also, for better accuracy, they apply local contrast normalization to each block before sending them to the classifier.

In my implementation, each individual histogram has 20 signed bins, with an interval length of 18.

Since people could tell the gender of a person from a black-white face image, I transformed the image to gray scale before extracting the HOG features. The only preprocessing we perform is image resizing. Observe that the resolution of image can affect the extracted edge feature. Images with low resolution contain mainly

the information of position and coarse edges while large scale images include more of fine edges. Once a face is detected, it is cropped and resized to four scales— $16 \times 16, 32 \times 32, 64 \times 64, 128 \times 128$. Then the feature is extracted using a window of size $2 \times 2, 4 \times 4, 8 \times 8, 16 \times 16$ respectively, either with or without overlapping. After obtaining the input data, we fit a svm.SVC classifier to it and perform classification. [8].

5.2.1 Result

After comparing the accuracy scores on test data, I chose the classifier trained by 64×64 images with overlapped windows. Most faces in first two clips are classified correctly. Good results are in Fig 11, Fig 12. There are a few



Figure 11: Correct classifications from clip1.

failures shown in Fig 13, Fig 14. The top example in 13 could be caused by the change of il-



Figure 12: Correct classifications from clip2.



Figure 13: Bad results from clip1.

lumination within the face. This may give us a wrong edge which would influence the result. The potential reason for the other two failures could be the vague outlines.

But when applied the trained classifier to clip3, we obtained low precision with all classifiers. I found the outlines of faces in clip3 are not as clear as those in clip1 and clip2. This could be one possible reason for the failure, since the classifier is trained on edge features. Another possible reason is that faces in the training data set are mostly front faces, but in clip3, the man in the last frame bows his head. In order to improve the generalization of the classifier, we could include more images with various head poses in the training set.

5.2.2 Future Work

- As previously mentioned, the classifier may fail when faces are not perfectly front or edges are not clear. It could also fail when there exists a partial mask in the face. One improvement I could make is to gather



Figure 14: Bad result from clip2.

training images with different head poses and augment the input by smoothing the faces or applying masks on randomly chosen areas. Rotating or translating the training data may also help.

- The method has not yet been tested on other datasets. So, future work includes testing it on larger test sets.
- Alexandre and Lu (2010) [1] shows that decision fusion improves the result significantly. Instead of training a single classifier, this paper proposed to train various classifiers using different image scales or different features, predict the gender separately and make final decision via major-

ity vote. I think above approach can be included for the future improvement since given a specific image, some features may be easy to obtain while others are not. We're interested of extracting geometric face features such as the relative position of eyes, nose and lip, as well as texture features shown in the paper.

- Feature extraction process consumes a lot of time. This could be shorten by code optimization

6 Conclusion

In this project, we found that the traditional SVM classifier with facial feature performs well when the face edges are clear but with a large possibility to fail when the illumination varies. Also, the resolution of the face image would affect the classification results.

7 Contribution

Hirotaka Ishihara

- 3 shot detection
- face detection
- logo detection
- face tracking and similarity function
- all code above

Yichun Zhang

- gender classification

References

- [1] Luis A Alexandre. “Gender recognition: A multiscale decision fusion approach”. In: *Pattern recognition letters* 31.11 (2010), pp. 1422–1427.
- [2] G. Bradski. “The OpenCV Library”. In: *Dr. Dobb’s Journal of Software Tools* (2000).
- [3] Navneet Dalal and Bill Triggs. “Histograms of oriented gradients for human detection”. In: 2005.
- [4] Martin A. Fischler and Robert C. Bolles. “Random Sample Consensus: A Paradigm for Model Fitting with Applications to Image Analysis and Automated Cartography”. In: *Commun. ACM* 24.6 (June 1981), pp. 381–395. ISSN: 0001-0782. DOI: [10.1145/358669.358692](https://doi.acm.org/10.1145/358669.358692). URL: <http://doi.acm.org/10.1145/358669.358692>.
- [5] Gary B. Huang et al. “Learning to Align from Scratch”. In: *Proceedings of the 25th International Conference on Neural Information Processing Systems - Volume 1*. NIPS’12. Lake Tahoe, Nevada: Curran Associates Inc., 2012, pp. 764–772. URL: <http://dl.acm.org/citation.cfm?id=2999134.2999220>.
- [6] Diederik P Kingma and Jimmy Ba. “Adam: A method for stochastic optimization”. In: *arXiv preprint arXiv:1412.6980* (2014).
- [7] David G Lowe. “Object Recognition from Local Scale-Invariant Features”. In: () .
- [8] F. Pedregosa et al. “Scikit-learn: Machine Learning in Python”. In: *Journal of Machine Learning Research* 12 (2011), pp. 2825–2830.

- [9] C. J. Van Rijsbergen. *Information Retrieval*. 2nd. Newton, MA, USA: Butterworth-Heinemann, 1979. ISBN: 0408709294.
- [10] Makarand Tapaswi, Martin Bauml, and Rainer Stiefelhagen. “Knock! Knock! Who is it? probabilistic person identification in TV-series”. In: *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition* i (2012), pp. 2658–2665. ISSN: 10636919. DOI: [10 . 1109 / CVPR . 2012 . 6247986](https://doi.org/10.1109/CVPR.2012.6247986).
- [11] Aidan Totterdell. “An Algorithm for Detecting and Classifying Scene Breaks in MPEG Video Bit Streams “”. In: *Dublin City University* September (1998). DOI: [10 . 1 . 1 . 74 . 5363](https://doi.org/10.1.1.74.5363).
- [12] P. Viola and M. Jones. “Haar-like”. In: *Cvpr* 1 (2001), pp. I-511-I-518. ISSN: 1063-6919. DOI: [10 . 1109 / CVPR . 2001 . 990517](https://doi.org/10.1109/CVPR.2001.990517). URL: [http : / / ieeexplore . ieee . org / document/990517/](http://ieeexplore.ieee.org/document/990517/).
- [13] Y. Yusoff, W. Christmas, and J. Kittler. “A study on automatic shot change detection”. In: *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)* 1425 (1998), pp. 177–189. ISSN: 16113349. DOI: [10 . 1007 / 3 – 540 – 64594 – 2{_\}94](https://doi.org/10.1007/3-540-64594-2{_\}94).
- [14] Y Yusoff, W Christmas, and J Kittler. “Video Shot Cut Detection Using Adaptive Thresholding”. In: January 2000 (2014).

Appendix

Decemnber 2019

1 Brief Pipeline

Pipeline Shot Detection

```
import numpy as np
import matplotlib.pyplot as plt
import cv2
import os
import sys
sys.path.insert(0, os.path.abspath('../'))

from utils import *
from ShotDetection import *

data = load_data()
shotDetector = ShotDetector(data)

clip_1 images importing...
clip_2 images importing...
clip_3 images importing...
Done !
```

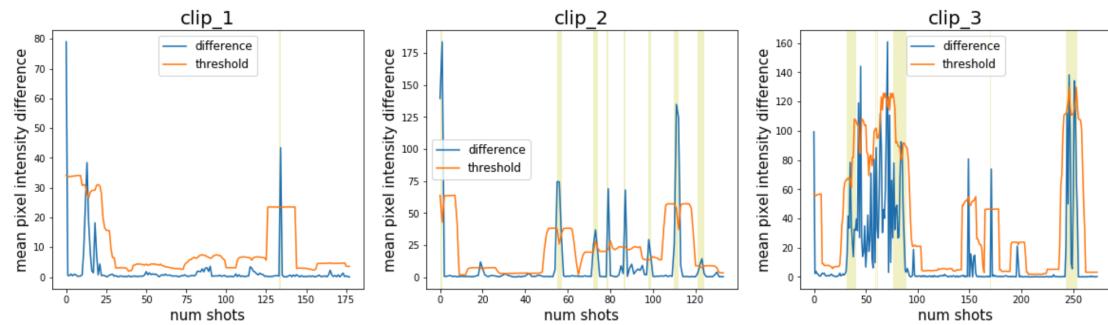
Shot Detection

Mean Pixel Intensity Measurement

```
shotDetector.show_detection('mean_pixel_intensity')

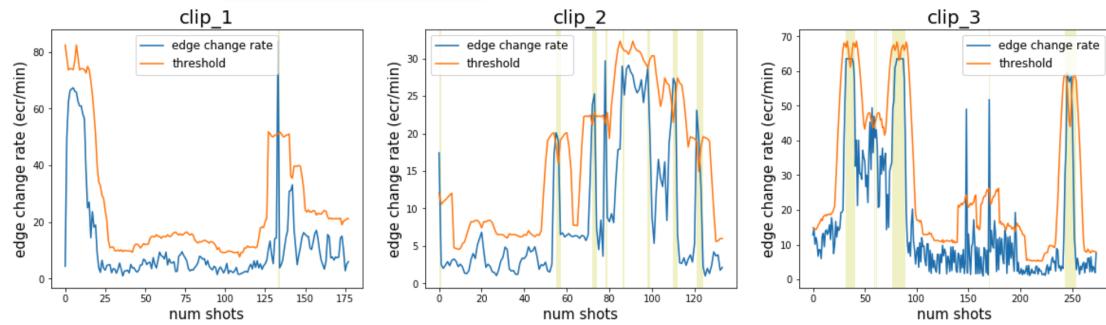
method used: mean_pixel_intensity

[Clip_1]  Recall: 1.00  Precision: 0.33  Combine: 0.50  Correct: 1  Missed: 0  False detect: 2
[Clip_2]  Recall: 1.00  Precision: 0.80  Combine: 0.89  Correct: 8  Missed: 0  False detect: 2
[Clip_3]  Recall: 0.67  Precision: 0.40  Combine: 0.50  Correct: 4  Missed: 2  False detect: 6
```



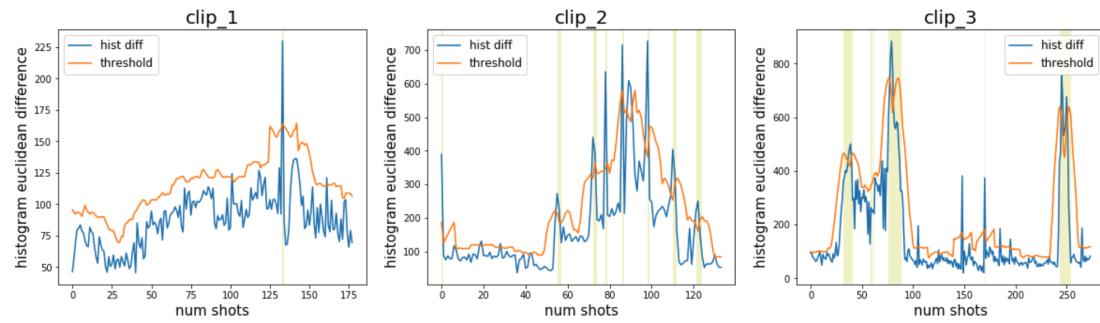
Edge Change Rate

```
shotDetector.show_detection('edge_change_rate', True)  
method used: edge_change_rate  
  
time used: 0.010544061660766602  
[Clip_1] Recall: 1.00 Precision: 1.00 Combine: 1.00 Correct: 1 Missed: 0 False detect: 0  
time used: 0.0033278465270996094  
[Clip_2] Recall: 0.75 Precision: 1.00 Combine: 0.86 Correct: 6 Missed: 2 False detect: 0  
time used: 0.002835988998413086  
[Clip_3] Recall: 0.83 click to scroll output; double click to hide e: 0.71 Correct: 5 Missed: 1 False detect: 3
```



Color Histogram

```
: shotDetector.show_detection('color_histogram_difference', True)  
method used: color_histogram_difference  
  
time used: 0.013776063919067383  
[Clip_1] Recall: 1.00 Precision: 0.33 Combine: 0.50 Correct: 1 Missed: 0 False detect: 2  
time used: 0.0026578903198242188  
[Clip_2] Recall: 1.00 Precision: 0.53 Combine: 0.70 Correct: 8 Missed: 0 False detect: 7  
time used: 0.0028378963470458984  
[Clip_3] Recall: 0.67 Precision: 0.33 Combine: 0.44 Correct: 4 Missed: 2 False detect: 8
```



```

: import os
import sys
import cv2
import matplotlib.pyplot as plt

from FaceDetector import *
from LogoDetector import *
from GenderDetector import *
from DetectorWraper import *
from utils import *
from config import *
from moviepy.editor import *

data = load_data()

clip_1 images importing...
clip_2 images importing...
clip_3 images importing...
Done !

```

1. load gender detector model ¶

```

: gender = GenderDetector()
0.8571428571428571

```

2. Initiate all detectors need for a clip

```

: detectorList = [FaceDetector(),
#                 LogoDetector('NBC'),
#                 LogoDetector('NBC exclusive'),
#                 LogoDetector('voice'),
#                 LogoDetector('clevver news'),
#                 LogoDetector('super woman'),
#                 LogoDetector('flick'),
#                 LogoDetector('marvel'),
]

detector = DetectorWraper(detectorList, CLIP_3, gender)
# image = detector.visualize_detection(image)

: detector.tracking(data["X"][2], 0.5, relabel=True)

```

1.1 Python Files

DetectorWraper.py

```

1
2 import cv2
3 import numpy as np
4 import heapq
5 from skimage.feature import hog
6 from config import *
7 from utils import *
8
9
10 class DetectorWraper(object):
11     """
12         A class that contains all the detectors need for a clip detection
13     """
14     def __init__(self, detectors, clip, genderDetect=None):
15         self.clip = clip
16         self.genderDetect = genderDetect
17         self.detectors = []
18         for detector in detectors:

```

```

19     self.detectors.append(detector)
20 self.num_save = 0
21 # later greedy algorithm for tracking
22 self.num_detect = 0
23 self.detection_frames = {}
24 self.detection_index= {}
25 self.track_frame = None
26
27 def apply_detection(self, image):
28 """
29     Return a list of tuple: List[(text, position)]
30 """
31 pos_list = []
32 for c in self.detectors:
33     pos_list += c.get_position(image)
34 return pos_list
35
36 def visualize_detection(self, image):
37 """
38     Return a image with visualized locations detected
39 """
40 H, W, _ = image.shape
41 pos_list = self.apply_detection(image)
42 detections = {}
43 hasDetection = False
44 for i, L in enumerate(pos_list):
45     text, coordinates = L[0], L[1]
46     COLOR = COLORS[text]
47     for x, y, w, h in coordinates:
48         # prune bad homography points
49         if x < 0 or y < 0 or x + w > W or \
50             y + h > H or w <= 1 or h <= 1:
51             continue
52         # add the detection to the dict for tracking
53         detections[self.num_detect] = (x, y, w, h)
54         self.detection_index[self.num_detect] = (x, y, w, h, self.num_save, text)
55         self.num_detect += 1
56         hasDetection = True
57         # if the detection is human
58         if text == 'face':
59             gender = self.genderDetect.classify(image[y:y+h, x:x+w, :])
60             gender = 'female' if gender[0] < 0.5 else 'male'
61             cv2.putText(image, gender, (x + w // 2, y + h + 5),
62                         cv2.FONT_HERSHEY_SIMPLEX, 0.6, COLOR, 2, cv2.LINE_AA)
63
64             image = cv2.rectangle(image, (x, y), (x + w, y + h), COLOR, 2)
65             cv2.putText(image, text, (x, y - 5),
66                         cv2.FONT_HERSHEY_SIMPLEX, 0.6, COLOR, 2, cv2.LINE_AA)
67         if hasDetection:
68             self.detection_frames[self.num_save] = detections
69             self.num_save +=1
70         return image

```

```

71
72     def save_detection(self, image):
73         """
74             Save the visualized image
75         """
76
77         img = self.visualize_detection(image)
78         img = cv2.cvtColor(img, cv2.COLOR_RGB2BGR)
79         cv2.imwrite(f'{SAVE_PATH}{self.clip}{self.num_save}.jpg', img)
80
81
82     def tracking(self, frameList, thres, relabel=False):
83         """
84             Perform detection tracking.
85             Input: the list the video frames frameList
86             Output: Example: 1. Dict('track 1':[obj_1, obj_3], 'track 2': [obj_2] ...)
87         """
88
89         if relabel:
90             print('detecting objects ...')
91             for frame in frameList:
92                 self._get_positions(frame)
93
94             track_obj = {}
95             track_frame = {}
96             num_track = 0
97             left_over = []
98             detected_frames = list(self.detection_frames.keys())
99
100            print('start tracking ...')
101            for i, frame in enumerate(detected_frames[:-1]):
102                # extract frame image
103                img1 = frameList[detected_frames[i]]
104                img2 = frameList[detected_frames[i + 1]]
105
106                # frame detections
107                frame_1 = self.detection_frames[detected_frames[i]]
108                frame_2 = self.detection_frames[detected_frames[i + 1]]
109
110                # detection id
111                detects_1 = list(frame_1.keys())
112                detects_2 = list(frame_2.keys())
113
114
115                # ===== frame wise compare stage (consecutive) =====
116                # init a heap for later ranking
117                HEAP = []
118
119                # get only human left over
120                only_human = []
121                detects_1 += left_over
122                for obj_1 in detects_1:
123                    max_score = -9999
124                    next_score = -9999
125                    pair = None
126                    ratio = None
127
128                    x1, y1, w1, h1, f, text1, _ = self.detection_index[obj_1]
129
130                    # ponly track human
131                    if text1 != 'face' and text1 != 'super woman':
132                        continue

```

```

123     only_human.append(obj_1)
124     for obj_2 in detects_2:
125         x2, y2, w2, h2, _, text2, _ = self.detection_index[obj_2]
126         # only track human
127         if text2 != 'face' and text2 != 'super woman':
128             continue
129
130         patch_1 = img1[y1 : y1 + h1, x1 : x1 + w1, :]
131         patch_2 = img2[y2 : y2 + h2, x2 : x2 + w2, :]
132         score = self.similarity_score(patch_1, patch_2)
133
134         if score > max_score:
135             next_score = max_score
136             max_score = score
137             pair = obj_1, f, obj_2
138             ratio = next_score / max_score
139             # calculate min over max ratio
140             if ratio is None or ratio > thres:
141                 continue
142             else: heapq.heappush(HEAP, (-max_score, (pair[0], pair[1], pair[2])))
143
144             # ===== track adding stage =====
145             repeated_obj1 = []
146             repeated_obj2 = []
147             for _ in range(len(HEAP)):
148                 track = heapq.heappop(HEAP)
149                 obj_1, frame_idx, obj_2 = track[1]
150                 # if already paired, skip
151                 if obj_2 in repeated_obj2:
152                     continue
153
154                 repeated_obj1.append(obj_1)
155                 repeated_obj2.append(obj_2)
156
157                 isPreviousTracked = False
158                 for track_id in track_obj.keys():
159                     if obj_1 in track_obj[track_id]:
160                         track_obj[track_id].append(obj_2)
161                         x, y, w, h, f, t, _ = self.detection_index[obj_2]
162                         self.detection_index[obj_2] = (x, y, w, h, f, t, track_id)
163                         isPreviousTracked = True
164                         break
165                 if isPreviousTracked:
166                     continue
167                 # if no previous track, add new one
168                 track_obj[num_track] = [obj_1, obj_2]
169                 x1, y1, w1, h1, f1, t1, _ = self.detection_index[obj_1]
170                 x2, y2, w2, h2, f2, t2, _ = self.detection_index[obj_2]
171                 self.detection_index[obj_1] = (x1, y1, w1, h1, f1, t1, num_track)
172                 self.detection_index[obj_2] = (x2, y2, w2, h2, f2, t2, num_track)
173                 num_track += 1
174

```

```

175     left_over = [obj for obj in only_human if obj not in repeated_obj]
176     print(len(left_over))
177
178     self.track_obj = track_obj
179     print('saving annotations ...')
180     self._annotate_images(frameList)
181     print('Done !')
182
183 def _get_positions(self, image):
184     """
185     SHOULD NOT BE CALL EXPLICITY, HIDDEN FUNCTION
186     return detection position without visualize positon
187     """
188
189     H, W, _ = image.shape
190     pos_list = self.apply_detection(image)
191     detections = {}
192     hasDetection = False
193
194     for i, L in enumerate(pos_list):
195         text, coordinates = L[0], L[1]
196         for x, y, w, h in coordinates:
197             if x < 0 or y < 0 or x + w > W or \
198                 y + h > H or w <= 1 or h <= 1:
199                 continue
200             # add the detection to the dict for tracking
201             if text == 'face' or text == 'super woman':
202                 self.detection_index[self.num_detect] = (x, y, w, h, self.num_save, text, -1)
203             else:
204                 self.detection_index[self.num_detect] = (x, y, w, h, self.num_save, text, -2)
205             detections[self.num_detect] = (x, y, w, h)
206             self.num_detect += 1
207             hasDetection = True
208
209     if hasDetection:
210         self.detection_frames[self.num_save] = detections
211     self.num_save +=1
212
213 def _annotate_images(self, frameList):
214     """
215     SHOULD NOT BE CALL EXPLICITY, HIDDEN FUNCTION
216     annotate positions in each frames with given frames
217     """
218
219     image_array = frameList
220     for i, image in enumerate(image_array):
221         if i in list(self.detection_frames.keys()):
222             for obj in list(self.detection_frames[i].keys()):
223                 x, y, w, h, frame, text, track_id = self.detection_index[obj]
224                 COLOR = COLORS[text]
225                 # if the detection is human
226                 if text == 'face':
227                     text = text + " id:{}".format(track_id)
228                     # predict
229                     gender = self.genderDetect.classify(image[y:y+h, x:x+w, :])
230                     gender = 'female' if gender[0] < 0.5 else 'male'

```

```

227         cv2.putText(image, gender, (x + w // 2 - 10, y + h + 15),
228                     cv2.FONT_HERSHEY_SIMPLEX, 0.6, COLOR, 2, cv2.LINE_AA)
229
230     image_array[i] = cv2.rectangle(image_array[i], (x, y), (x + w, y + h), COLOR,
231                                    2)
232     cv2.putText(image_array[i], text, (x, y - 5),
233                 cv2.FONT_HERSHEY_SIMPLEX, 0.6, COLOR, 2, cv2.LINE_AA)
234
235     cv2.imwrite(f'{SAVE_PATH}{self.clip}{i}.jpg', image_array[i])
236
237 def similarity_score(self, img1, img2):
238     """
239     Calculate the similarity score used for detection tracking
240     """
241
242     # resize into the same shape first
243     if img1.shape != img2.shape:
244         v, h = max(img1.shape[0], img2.shape[0]), max(img1.shape[1], img2.shape[1])
245         dim = (h, v)
246         h_scale = min(img1.shape[1], img2.shape[1]) / h
247         v_scale = min(img1.shape[0], img2.shape[0]) / v
248         img1 = cv2.resize(img1, dim, interpolation=cv2.INTER_AREA)
249         img2 = cv2.resize(img2, dim, interpolation=cv2.INTER_AREA)
250
251     # histogram
252     diff = 0
253     for c in range(3):
254         hist1 = cv2.calcHist([img1], [c], None, [256], [0, 256])
255         hist2 = cv2.calcHist([img2], [c], None, [256], [0, 256])
256         diff += np.linalg.norm(hist1 - hist2)
257
258     # HoG
259     fd1, _ = hog(img1, orientations=8, pixels_per_cell=(16, 16),
260                  cells_per_block=(1, 1), visualize=True, multichannel=True)
261     fd2, _ = hog(img2, orientations=8, pixels_per_cell=(16, 16),
262                  cells_per_block=(1, 1), visualize=True, multichannel=True)
263
264     # Combine both
265     dist = np.linalg.norm(fd1 - fd2)
266     aim = mean_pixel_intensity_diff(img1, img2)
267     score = 1 / (dist + diff + aim + 1)
268
269     return score

```

GenderDetector.py

```
1 import glob
2 from skimage import transform
3 import numpy as np
4 from sklearn.model_selection import train_test_split
5 from sklearn.svm import SVC
6 import cv2
7
8 from config import *
9
10 class GenderDetector(object):
11     """
12         A class detect gender given detected faces
13     """
14     def __init__(self, img_size=64, window_size=(8, 8), overlap=2, intrvlInDegree=18):
15         self.img_size = img_size
16         self.window_size = window_size
17         self.overlap = overlap
18         self.intrvlInDegree = intrvlInDegree
19         self.data_0 = self._get_train_data(female_path)
20         self.data_1 = self._get_train_data(male_path)
21         self.classifier = self._get_shape_classifier()
22
23
24     def classify(self, face):
25         """
26             classify female or male given croped face
27         """
28         face = self._resize_image(face)
29         feature = self._get_shape_feature(face)
30         feature = np.array(feature)
31         feature = feature.reshape((1, feature.shape[0]*feature.shape[1]))
32         result = self.classifier.predict(feature)
33
34         return result
35
36
37     def _detect_face(self, img, min_nghb=5):
38         """
39             Detect face, used to training
40         """
41
42         face_cascade = cv2.CascadeClassifier(train_svm)
43         img_gray = cv2.cvtColor(img, cv2.COLOR_RGB2GRAY)
44         faces = face_cascade.detectMultiScale(img_gray, minNeighbors=min_nghb, minSize
45 = (30,30))
46         crops = [] #store in a list, can think of converting to an array if vectorization
47         helps
48         for (x,y,w,h) in faces:
49             crops.append(img_gray[y:y+h, x:x+w])
50
51         return crops
52
53     def _resize_image(self, img):
54         return transform.resize(img, (self.img_size, self.img_size), mode='reflect',
```

```

    anti_aliasing=True)

50
51     def _get_train_data(self, path):
52         img_paths = glob.glob(pathname=path)
53         imgs = [cv2.cvtColor(cv2.imread(img_path), cv2.COLOR_BGR2RGB) for img_path in
54         img_paths]
55         data = []
56         for image in imgs:
57             crops = self._detect_face(image)
58             data.extend([self._resize_image(crop) for crop in crops])
59         return data

60
61     def _get_shape_feature(self, img):
62         if not 360%self.intrvlInDegree == 0:
63             raise RuntimeError
64         elif img.shape[0]%self.window_size[0]!=0 or img.shape[1]%self.window_size[1]!=0:
65             raise RuntimeError
66         else:
67             # Compute image gradient using Sobel filter
68             sobelx = cv2.Sobel(img, cv2.CV_64F, 1, 0, ksize=3)
69             sobely = cv2.Sobel(img, cv2.CV_64F, 0, 1, ksize=3)

70             # Get the gradient magnitude and orientation
71             magnitudes = np.round(np.hypot(sobelx, sobely))
72             angles = cv2.phase(sobelx, sobely, angleInDegrees=True)

73             # Divide angle by the length of interval(e.g. 18), take the floor of that
74             # number, then we get class for each gradient
75             categories = np.floor_divide(angles, self.intrvlInDegree)
76             # print("Category\n", categories.shape, categories)
77             # Compute total number of windows
78             num_window = (self.overlap*img.shape[0]//self.window_size[0])*(self.overlap*
img.shape[1]//self.window_size[1])

79             # Create the storage space for the image feature
80             feature = np.zeros((num_window, 360//self.intrvlInDegree))

81
82             window_idx = 0
83             for i in range(0, img.shape[0], self.window_size[0]//self.overlap):
84                 for j in range(0, img.shape[1], self.window_size[1]//self.overlap):
85                     # print("window", window_idx, "\n")
86                     magn_window = magnitudes[i:i+self.window_size[0], j:j+self.
window_size[1]].flatten()
87                     category_window = categories[i:i+self.window_size[0], j:j+self.
window_size[1]].flatten()

88                     for idx_bin in range(360//self.intrvlInDegree):
89                         # Sum the maginitude where the angle is in the current interval
90                         range(category)
91                         indices = np.argwhere(category_window==idx_bin)
92                         feature[window_idx, idx_bin] = np.sum(magn_window[indices])
93                         # print("Current feature window: \n", feature[window_idx, :])

```

```

95     #           Increment the index of window
96     window_idx += 1
97
98     return feature
99
100
101    def _get_shape_classifier(self):
102        features_0 = []
103        features_1 = []
104
105        # Extract shape features for both classes
106        for img in self.data_0:
107            feature = self._get_shape_feature(img)
108            features_0.append(feature.flatten())
109
110        for img in self.data_1:
111            feature = self._get_shape_feature(img)
112            features_1.append(feature.flatten())
113
114        # Turn it from list to numpy array
115        features_0 = np.stack(features_0)
116        features_1 = np.stack(features_1)
117
118
119        # Label the features
120        features_10 = np.column_stack((features_0, np.zeros(features_0.shape[0])))
121        features_11 = np.column_stack((features_1, np.ones(features_1.shape[0])))
122
123        # Combine two classes together
124        features_labeled = np.vstack((features_10, features_11))
125
126        # Shuffle the data
127        features_labeled = np.random.permutation(features_labeled)
128
129
130        # Split it to X and y
131        X = features_labeled[:, :features_labeled.shape[1]-1]
132        y = features_labeled[:, -1]
133
134
135        X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.2)
136
137
138        svcclassifier = SVC(kernel='linear')
139        svcclassifier.fit(X_train, y_train)
140
141        # Compute the accuracy using cross validation
142        print(svcclassifier.score(X_test,y_test))
143
144        return svcclassifier

```

LogoDetector.py

```

1 import cv2
2 import itertools
3 from config import *
4 from utils import *
5
6 COLOR = COLORS['face']
7 sift = cv2.xfeatures2d.SIFT_create()
8
9
10 class LogoDetector(object):
11     """
12         A class performs Logo detection given video clips
13     """
14     def __init__(self, logName):

```

```

15     self.match = cv2.BFMatcher(cv2.NORM_L1, crossCheck=False)
16     template = cv2.imread(LOGO[logName])
17     self.template = cv2.cvtColor(template, cv2.COLOR_RGB2GRAY)
18     self.logName = logName
19
20     def get_position(self, image):
21         """
22             Return position list
23         """
24
25         image = cv2.cvtColor(image, cv2.COLOR_RGB2GRAY)
26         # do rancan and homography
27         pos = self._homography(image)
28         if pos is None:
29             return []
30         top_l, bottom_l, bottom_r, top_r = pos
31         x, y = top_l[0], top_l[1]
32         w, h = top_r[0] - top_l[0], bottom_l[1] - top_l[1]
33         return [[self.logName, [(x, y, w, h)]]]
34
35     def _visualize_detection(self, pos_list, image):
36         """
37             Visualize detection
38         """
39
40         for i, L in enumerate(pos_list):
41             text, coordinates = L[0], L[1]
42             for x, y, w, h in coordinates:
43                 image = cv2.rectangle(image, (x, y), (x + w, y + h), COLOR, 2)
44                 cv2.putText(image, text, (x, y - 5),
45                             cv2.FONT_HERSHEY_SIMPLEX, 0.6, COLOR, 1, cv2.LINE_AA)
46         return image
47
48     def _homography(self, image):
49         """
50             Compute homography from image to the given video clip
51         """
52
53         # get matching points
54         kp1, des1 = sift.detectAndCompute(image, None)
55         kp2, des2 = sift.detectAndCompute(self.template, None)
56         bf = cv2.BFMatcher()
57         matches = bf.knnMatch(des1, des2, k=2)
58         good = []
59
60         # prune the matching
61         for m,n in matches:
62             if m.distance < 0.5 * n.distance:
63                 good.append([m])
64
65         # if no good matching, return
66         if len(good) < 4:
67             return
68
69         # extract location info from matching point
70         kp1_index = [m[0].queryIdx for m in good]
71         kp2_index = [m[0].trainIdx for m in good]
72         kp1 = np.array([k.pt for k in np.array(kp1)[kp1_index]])

```

```

67     kp2 = np.array([k.pt for k in np.array(kp2)[kp2_index]])
68     # calculate homography using RANSAC
69     H, _ = cv2.findHomography(kp2, kp1, cv2.RANSAC, 0.1)
70     if H is None:
71         return
72     # get location of the self.template in the image
73     height, width = self.template.shape
74     top_left = H.dot(np.array([[0, 0, 1]]).reshape(3, 1)).flatten()
75     bottom_left = H.dot(np.array([[0, height, 1]]).reshape(3, 1)).flatten()
76     bottom_right = H.dot(np.array([[width, height, 1]]).reshape(3, 1)).flatten()
77     top_right = H.dot(np.array([[width, 0, 1]]).reshape(3, 1)).flatten()
78     # rescale the points
79     top_left = top_left[:2] / top_left[-1]
80     bottom_left = bottom_left[:2] / bottom_left[-1]
81     bottom_right = bottom_right[:2] / bottom_right[-1]
82     top_right = top_right[:2] / top_right[-1]
83     # int type
84     top_left = top_left.astype(int)
85     bottom_left = bottom_left.astype(int)
86     bottom_right = bottom_right.astype(int)
87     top_right = top_right.astype(int)
88
89     return top_left, bottom_left, bottom_right, top_right

```

FaceDetector.py

```

1 import cv2
2 import numpy as np
3 from config import *
4 from utils import *
5
6 COLOR = COLORS['face']
7
8
9 class FaceDetector(object):
10 """
11 A class performs face detect in video clips
12 """
13 def __init__(self, List=[face1, face2, profile1, profile2]):
14     self._classifiers = []
15     for xml in List:
16         self._classifiers.append(cv2.CascadeClassifier(xml))
17     self.category = 'face'
18
19 def get_position(self, image):
20 """
21 Return list of positions that detected faces
22 """
23     image = cv2.cvtColor(image, cv2.COLOR_RGB2GRAY)
24     pos = []
25     for c in self._classifiers:
26         p = c.detectMultiScale(image, 1.3, 3)
27         if len(p) != 0:
28             pos += p.tolist()

```

```

29     pos_list = self._prune_overlap(pos)
30     return [['face', pos_list]]
31
32 def _visualize_detection(self, pos_list, image):
33     """
34     Visualize the face detection
35     """
36     for i, L in enumerate(pos_list):
37         text, coordinates = L[0], L[1]
38         for x, y, w, h in coordinates:
39             image = cv2.rectangle(image, (x, y), (x + w, y + h), COLOR, 2)
40             cv2.putText(image, text, (x, y - 5),
41                         cv2.FONT_HERSHEY_SIMPLEX, 0.7, COLOR, 2, cv2.LINE_AA)
42     return image
43
44 def _prune_overlap(self, List):
45     """
46     prune face detection with overly larged overlap
47     """
48     result = []
49     for pos in List:
50         if len(result) == 0:
51             result.append(pos)
52             continue
53         include = True
54         for old_pos in result:
55             overlap = self._overlap(pos, old_pos)
56             if overlap > 0.2:
57                 include = False
58         if include:
59             result.append(pos)
60     return result
61
62 def _overlap(self, pos1, pos2):
63     """
64     Calculate IoU region given two boxes
65     """
66     x1, y1, w1, h1 = pos1
67     x2, y2, w2, h2 = pos2
68     # calculate the intersection
69     top = max(x1, x2), min(y1 + h1, y2 + h2)
70     bottom = min(x1 + w1, x2 + w2), max(y1 + h1, y2 + h2)
71     interset = max(0, min(y1 + w1, y2 + w2) - max(y1, y2)) * \
72                 max(0, min(x1 + w1, x2 + w2) - max(x1, x2))
73     # calculate the union + intersection
74     union1 = w1 * h1
75     union2 = w2 * h2
76
77     return interset / float(union1 + union2 - interset)

```

ShotDetector.py

```

1 import matplotlib.pyplot as plt
2 import numpy as np

```

```

3 import time
4 from tempfile import TemporaryFile
5 from utils import *
6 from config import *
7
8 MEAN_PIXEL_INTENSITY = 'mean_pixel_intensity'
9 ECR = 'edge_change_rate'
10 HIST_DIFF = 'color_histogram_difference'
11
12
13 class ShotDetector(object):
14     """
15     A class detect shot transitions given video clips
16     """
17
18     def __init__(self, data):
19         self.data = data
20         self.img_array1 = data['X'][0]
21         self.img_array2 = data['X'][1]
22         self.img_array3 = data['X'][2]
23
24     def show_detection(self, method, *arg):
25         """
26             Create plot of shot detection and the threshold given 'method'
27         """
28
29         if method == MEAN_PIXEL_INTENSITY:
30             print(f'method used: {MEAN_PIXEL_INTENSITY}\n')
31             if len(arg) > 0:
32                 self._mean_p_int(window_size=arg[0], T=arg[1])
33             else: self._mean_p_int()
34
35         elif method == ECR:
36             print(f'method used: {ECR}\n')
37             if len(arg) > 0:
38                 self._edge_change_rate(load=arg[0])
39             else: self._edge_change_rate()
40
41         elif method == 'color_histogram_difference':
42             print(f'method used: {HIST_DIFF}\n')
43             if len(arg) > 0:
44                 self._hist_diff(load=arg[0])
45             else: self._hist_diff()
46
47         else:
48             raise Exception(f"Method <{method}> valid")
49
50     # ===== hidden function =====
51
52     def _mean_p_int(self, window_size=15, T=4):
53         """
54             Apply mean intensity measurement
55         """
56
57         window_sizes, Ts = [19, 15, 15], [5, 2, 7]
58         fig = plt.figure(figsize=(20, 5))
59         for clip in range(3):
60             window_size, T, diff = window_sizes[clip], Ts[clip], []

```

```

55     img_array = self.data['X'][clip]
56     # clip 1
57     for i in range(len(img_array) - 1):
58         diff.append(mean_pixel_intensity_diff(img_array[i - 1], img_array[i]))
59     frames, thres = adaptive_threshold(diff, window_size=window_size, T=T)
60     # plot clip
61     plt.subplot(1, 3, clip + 1)
62     plt.plot(np.arange(len(diff)), diff, label='difference')
63     plt.plot(np.arange(len(thres)), thres, label='threshold')
64     for shots in self.data['Y'][clip].tolist():
65         plt.axvspan(shots[0], shots[1], facecolor='y', alpha=0.2)
66     plt.xlabel('num shots', size=15)
67     plt.ylabel('mean pixel intensity difference', size=15)
68     plt.legend(fontsize=12)
69     plt.title(f'clip_{clip + 1}', size=20)
70     # brief evaluate
71     Recall, Precision, Combine, C, M, F = Evaluate(frames, self.data['Y'][clip])
72     print(f"[Clip_{clip + 1}]  Recall: {Recall:.2f}" + \
73           f"  Precision: {Precision:.2f}  Combine: {Combine:.2f}  Correct: {C}" \
74           + \
75           f"  Missed: {M}  False detect: {F}")
76     plt.show()
77
78 def _edge_change_rate(self, load=True):
79     """
80     Apply Edge Change Ratio
81     """
82     window_sizes, Ts = [15, 13, 19], [6, 3, 3]
83     fig = plt.figure(figsize=(20, 5))
84     for clip in range(3):
85         window_size, T, ecr = window_sizes[clip], Ts[clip], []
86         img_array = self.data['X'][clip]
87         # clip 1
88         start = time.time()
89         if load:
90             try:
91                 ecr = np.load(f'{CACHE}ecr_ecr_{clip}.npy')
92                 frames = np.load(f'{CACHE}frames_ecr_{clip}.npy')
93                 thres = np.load(f'{CACHE}thres_ecr_{clip}.npy')
94             except:
95                 for i in range(1, len(img_array)):
96                     ecr.append(EdgeChangeRate(img_array[i - 1], img_array[i], sigma=3))
97                     # apply threshold
98                     ecr = ecr / np.min(ecr)
99                     frames, thres = adaptive_threshold(ecr, window_size=window_size, T=T)
100                    np.save(f'{CACHE}ecr_ecr_{clip}.npy', ecr)
101                    np.save(f'{CACHE}frames_ecr_{clip}.npy', frames)
102                    np.save(f'{CACHE}thres_ecr_{clip}.npy', thres)
103            else:
104                for i in range(1, len(img_array)):
105                    ecr.append(EdgeChangeRate(img_array[i - 1], img_array[i], sigma=3))
106                    # apply threshold

```

```

106         ecr = ecr / np.min(ecr)
107         frames, thres = adaptive_threshold(ecr, window_size=window_size, T=T)
108         np.save(f'{CACHE}ecr_ecr_{clip}.npy', ecr)
109         np.save(f'{CACHE}frames_ecr_{clip}.npy', frames)
110         np.save(f'{CACHE}thres_ecr_{clip}.npy', thres)
111
112
113     print(f"time used: {time.time() - start}")
114
115     # plot clip
116     plt.subplot(1, 3, clip + 1)
117     plt.plot(np.arange(len(ecr)), ecr, label='edge change rate')
118     plt.plot(np.arange(len(thres)), thres, label='threshold')
119     for shots in self.data['Y'][clip].tolist():
120         plt.axvspan(shots[0], shots[1], facecolor='y', alpha=0.2)
121     plt.xlabel('num shots', size=15)
122     plt.ylabel('edge change rate (ecr/min)', size=15)
123     plt.legend(fontsize=12)
124     plt.title(f'clip_{clip + 1}', size=20)
125
126     # brief evaluate
127     Recall, Precision, Combine, C, M, F = Evaluate(frames, self.data['Y'][clip])
128     print(f"[Clip_{clip + 1}]    Recall: {Recall:.2f}" + \
129           f"    Precision: {Precision:.2f}    Combine: {Combine:.2f}    Correct: {C}" + \
130           f"    Missed: {M}    False detect: {F}")
131     plt.show()
132
133
134     def _hist_diff(self, load=True):
135         """
136             Apply histogram difference
137         """
138         window_sizes, Ts = [19, 13, 19], [7, 5, 6]
139         fig = plt.figure(figsize=(20, 5))
140         for clip in range(3):
141             window_size, T, hist_diff = window_sizes[clip], Ts[clip], []
142             img_array = self.data['X'][clip]
143             # clip 1
144             start = time.time()
145             if load:
146                 try:
147                     hist_diff = np.load(f'{CACHE}diff_hist_{clip}.npy')
148                     frames = np.load(f'{CACHE}frames_hist_{clip}.npy')
149                     thres = np.load(f'{CACHE}thres_hist_{clip}.npy')
150                 except:
151                     for i in range(1, len(img_array)):
152                         hist_diff.append(color_hist_diff(img_array[i - 1], img_array[i],
153 32))
154
155             # apply threshold
156             frames, thres = adaptive_threshold(hist_diff, window_size=window_size,
157 , T=T)
158             np.save(f'{CACHE}diff_hist_{clip}.npy', hist_diff)
159             np.save(f'{CACHE}frames_hist_{clip}.npy', frames)
160             np.save(f'{CACHE}thres_hist_{clip}.npy', thres)
161         else:

```

```

156         for i in range(1, len(img_array)):
157             hist_diff.append(color_hist_diff(img_array[i - 1], img_array[i], 32))
158             # apply threshold
159             frames, thres = adaptive_threshold(hist_diff, window_size=window_size, T=T)
160             np.save(f'{CACHE}diff_hist_{clip}.npy', hist_diff)
161             np.save(f'{CACHE}frames_hist_{clip}.npy', frames)
162             np.save(f'{CACHE}thres_hist_{clip}.npy', thres)
163
164             print(f"time used: {time.time() - start}")
165             # plot clip
166             plt.subplot(1, 3, clip + 1)
167             plt.plot(np.arange(len(hist_diff)), hist_diff, label='hist diff')
168             plt.plot(np.arange(len(thres)), thres, label='threshold')
169             for shots in self.data['Y'][clip].tolist():
170                 plt.axvspan(shots[0], shots[1], facecolor='y', alpha=0.2)
171             plt.xlabel('num shots', size=15)
172             plt.ylabel('histogram euclidean difference', size=15)
173             plt.legend(fontsize=12)
174             plt.title(f'clip_{clip + 1}', size=20)
175             # brief evaluate
176             Recall, Precision, Combine, C, M, F = Evaluate(frames, self.data['Y'][clip])
177             print(f"[Clip_{clip + 1}]  Recall: {Recall:.2f}" + \
178                 f"  Precision: {Precision:.2f}  Combine: {Combine:.2f}  Correct: {C}" + \
179                 f"  Missed: {M}  False detect: {F}")
180             plt.show()

```

utils.py

```

1 import numpy as np
2 import numba
3 from numba import jit
4 import cv2
5 import glob
6 from moviepy.editor import *
7 import matplotlib.pyplot as plt
8 from scipy import ndimage
9 from skimage import color, feature, transform
10
11 data = './data/'
12
13 @numba.jit(forceobj=True)
14 def load_data():
15     """
16     load all three clips to a dict variable
17     """
18     # ===== first clip =====
19     print('clip_1 images importing...')
20     fileName1 = sorted(glob.glob(data + 'clip_1/*.jpg'))
21     img1 = [cv2.imread(img) for img in fileName1]
22     labell = np.zeros((len(img1), 1))
23     shot1 = np.array([(133, 134)])
24     # ===== second clip =====
25     print('clip_2 images importing...')

```

```

26     fileName2 = sorted(glob.glob(data + 'clip_2/*.jpg'))
27     img2 = [cv2.imread(img) for img in fileName2]
28     label2 = np.zeros((len(img2), 1))
29     shot2 = np.array([(0, 1), (55, 57), (72, 74), (78, 79),
30                       (86, 87), (98, 99), (110, 112), (121, 124)])
31     # ===== third clip =====
32     print('clip_3 images importing...')
33     fileName3 = sorted(glob.glob(data + 'clip_3/*.jpg'))
34     img3 = [cv2.imread(img) for img in fileName3]
35     label3 = np.zeros((len(img3), 1))
36     shot3 = np.array([(32, 41), (59, 60), (61, 62),
37                       (76, 89), (170, 171), (243, 254)])
38     print('Done !')
39     test = {'X': [img1, img2, img3], 'Y': [shot1, shot2, shot3]}
40     return test
41
42 def image_to_video(path):
43     images = sorted(glob.glob(path + '*.jpg'))
44     clips = [ImageClip(m).set_duration(2) for m in images]
45     concat_clip = concatenate_videoclips(clips, method="compose")
46     concat_clip.write_videofile(path + "test.mp4", fps=3)
47
48 def color_hist_diff(image1, image2, patch_size=16):
49     """
50     calculate the color histogram difference of two images
51     """
52     # get dimesion and patch
53     vertical, horizontal, Z = image1.shape
54     v_patch = vertical // 16
55     h_patch = horizontal // 16
56     # calculate difference
57     diff = 0
58     for z in range(Z):
59         img1, img2 = image1[:, :, z], image2[:, :, z]
60         for i in range(0, img1.shape[0] - v_patch + 1, patch_size):
61             for j in range(0, img1.shape[1] - h_patch + 1, patch_size):
62                 patch1 = img1[i:i + v_patch, j:j + h_patch].flatten()
63                 patch2 = img2[i:i + v_patch, j:j + h_patch].flatten()
64                 hist1 = np.histogram(patch1, bins=np.arange(257), density=True)
65                 hist2 = np.histogram(patch2, bins=np.arange(257), density=True)
66                 diff += np.linalg.norm(hist1[0] - hist2[0])
67     return diff
68
69 def mean_pixel_intensity_diff(img1, img2):
70     diff = 0
71     for z in range(img1.shape[-1]):
72         diff += np.abs(np.mean(img1[:, :, z]) - np.mean(img2[:, :, z]))
73     return diff
74
75 # ===== calculate Edge Change Ratio =====
76
77 @numba.jit(forceobj=True)

```

```

78 def EdgeChangeRate(img1, img2, iteration=1, sigma=5):
79     # convert to gray
80     gray1 = color.rgb2gray(img1)
81     gray2 = color.rgb2gray(img2)
82     # get white background and edge
83     black1 = feature.canny(gray1, sigma=sigma).astype("uint8") # background: 0 edge: 1
84     black2 = feature.canny(gray2, sigma=sigma).astype("uint8")
85     # count number of edge pixel
86     E1 = max(np.sum(black1 == 1), 0.1)
87     E2 = max(np.sum(black2 == 1), 0.1)
88     # dilate both image
89     kernel = np.ones((3, 3)).astype("uint8")
90     dilate1 = cv2.dilate(black1, kernel).astype("uint8")
91     dilate2 = cv2.dilate(black2, kernel).astype("uint8")
92     # combine
93     imgIn = black1 * dilate2
94     imgOut = black2 * dilate1
95     # count edge change pixel
96     C1 = np.sum(imgIn == 1)
97     C2 = np.sum(imgOut == 1)

98
99     return max(1 - min(C1 / E2, C2 / E1), 0)

100

101
102 # ===== get result =====
103
104 def adaptive_threshold(dissimilarity, window_size=9, T=5):
105     # Dugad el. al Model
106     if type(dissimilarity) != np.ndarray:
107         dissimilarity = np.array(dissimilarity)
108     # split into two windows
109     w = window_size // 2
110     dissimilarity = np.pad(dissimilarity, pad_width=w, mode='constant')
111     thresholds = []
112     for i in range(w, len(dissimilarity) - w):
113         left = dissimilarity[i - w : i]
114         right = dissimilarity[i : i + w]
115         threshold = max(np.mean(left) + T * np.sqrt(np.std(left)),
116                         np.mean(right) + T * np.sqrt(np.std(right)))
117         thresholds.append(threshold)

118     shot_index = np.arange(len(dissimilarity) - window_size + 1)[dissimilarity[w : -w] >
119                         thresholds]
120     return shot_index, thresholds

121
122
123 # ===== Evaluation =====
124
125 def Evaluate(pred, target):
126     # C: correct detection
127     # M: missed detection
128     # F: false detection

```

```

129 C, M, F = 0, 0, 0
130 # pred may have multiple frames detection for dissolve
131 total_correct = 0
132 for shot in target:
133     # count for dissolve
134     if type(shot) not in [int, float]:
135         # output correctly predict one of the dissolve frames
136         detected = sum([output in range(shot[0], shot[1] + 1) for output in pred])
137         C += 1 if detected > 0 else 0
138         M += 0 if detected != 0 else 1
139         total_correct += detected
140     # hard cut case
141     else:
142         C += 1 if shot in pred else 0
143         M += 0 if shot in pred else 1
144         total_correct += 1 if shot in pred else 0
145     # false detection
146 F = len(pred) - total_correct
147 # Recall
148 Recall = C / (C + M)
149 Precision = C / (C + F) if (C + F) != 0 else 0
150 Combine = 2 * Recall * Precision / (Recall + Precision) if (Recall + Precision) != 0
151 else 0
152 return Recall, Precision, Combine, C, M, F

```

config.py

```

1 """
2 Configurations for the entire project
3 """
4
5 # color setting
6 COLORS = {'face': (255, 255, 0),
7            'NBC': (0, 255, 0),
8            'NBC exclusive': (0, 255, 255),
9            'voice': (255, 0, 0),
10           'flick': (255, 0, 255),
11           'clevver news': (255, 255, 255),
12           'super woman': (0, 0, 255),
13           'marvel': (230, 9, 237)
14          }
15 # train data
16 female_path = './data/train_data/female/*.jpg'
17 male_path = './data/train_data/male/*.jpg'
18 # xml file for haar cascade
19 face1 = './xml/haarcascade_frontalface_default.xml'
20 face2 = './xml/mallick_haarcascade_frontalface_default.xml'
21 profile1 = './xml/haarcascade_profileface.xml'
22 profile2 = './xml/mallick_haarcascade_profileface.xml'
23 train_svm = './xml/haarcascade_frontalface_alt_tree.xml'
24 # logo template file path
25 LOGO = {'NBC': './Logo/nbc.png',
26          'NBC exclusive': './Logo/nbc_exclusive.png',
27          'voice': './Logo/the_voice.png',

```

```
28     'flick': './Logo/flick.png',
29     'twitter': './Logo/twitter.jpg',
30     'clevver news': './Logo/clevver_news.png',
31     'super woman': './Logo/white_super_woman.png',
32     'marvel': './Logo/marvel.png',
33   }
34 # save labeled image
35 SAVE_PATH = './save/'
36 CACHE = './cache/'
37 CLIP_1 = 'clip_1/'
38 CLIP_2 = 'clip_2/'
39 CLIP_3 = 'clip_3/'
```