



UNIVERSITY OF BORDEAUX 1

INTERNSHIP REPORT
MASTER OF SOFTWARE ENGINEERING (2009-2011)

Semantic Saliency of Video Content: Application to Quality Assessment and Efficient Video Retrieval

Author:
Ho Tien Lam

Supervisor:
Prof. Jenny
BENOIS-PINEAU

October 31, 2011

Abstract

Enjoying video on digital television as well as Internet has become a significant demand at the present time. People demand not only a fast system response time but also a high video quality. However, through the transmission networks and processing systems, some amounts of errors may be introduced in the video signal, so video quality assessing is an important problem. Many researchers proposed a lot of methods for evaluating video quality. Researchers at LaBRI have proposed to assess video quality taking into account visual saliency. Visual saliency of an item is the state or quality by which it stands out relative to its neighbors. Spatial and temporal saliency were used for defining an objective quality assessment metric.

The aim of my internship was to help in studying if a semantic saliency, based on faces, could improve this metric. In this internship, I have studied how to create semantic saliency maps, if human faces on video frame attract viewers' attention, and how to use a face detection algorithm for automatic semantic map construction.

Contents

Acknowledgments	4
Introduction	9
1 Context	10
1.1 PUF	10
1.2 LaBRI	11
1.3 Video Analysis and Indexing Group	11
2 Visual Saliency Detection	13
2.1 BBBox Annotation Software	13
2.1.1 Introduction of Qt Framework	14
2.1.2 Functionality of BBBox software	14
2.1.3 Missing Features of BBBox software	15
2.1.4 Architecture of BBBox software	16
2.1.5 Modification of BBBox software	18
2.2 Ground-truth for Semantic Saliency	22
2.2.1 Introduction of OpenCV library	22
2.2.2 Annotate Sources of Videos	23
2.2.3 Generate Gaussian Map	23
2.3 Statistics of Faces in Saliency Map	28
2.3.1 Computational Methods	28
2.3.2 Statistical Results	30
2.4 Automatic Face Detection	33
2.4.1 Introduction of Viola–Jones algorithm	33
2.4.2 Gaussian Map Construction	34
2.4.3 Statistical Results with Automatic Detection	34
3 Conclusion and Future Work	38
3.1 Conclusion	38
3.2 Future Work	39

List of Figures

1	Process Overview	6
2	Example of Bounding boxes and Gaussian map	7
3	Example of Error map and Ground-truth saliency map	8
2.1	Screenshot of BBBox	14
2.2	Class diagram of BBBox	16
2.3	UML diagram of Observer pattern	17
2.4	UML diagram of Observer pattern used in BBBox	18
2.5	Flowchart for solving the first problem	19
2.6	Flowchart for solving the second problem	21
2.7	1D Gaussian function	24
2.8	2D Gaussian function	24
2.9	Processing steps for creating Gaussian Map	25
2.10	Example of overlapped BBox	29
2.11	Statistics of faces in saliency maps	31
2.12	Statistics of faces in error maps	31
2.13	Statistics of faces in combination of saliency maps and error maps	32
2.14	Cascade Architecture	34
2.15	Statistics of faces (automatic detection) in saliency maps	35
2.16	Statistics of faces (automatic detection) in error maps	36
2.17	Statistics of faces (automatic detection) in combination of saliency maps and error maps	37

List of Tables

2.1	Description of some important classes	26
2.2	Percentage of variation in mean values	36

Acknowledgments

Firstly, I would like to thank Professor Jenny Benois-Pineau, head of Video Analysis and Indexing research group, for giving me an opportunity to do this internship at the LaBRI.

I would like to thank Boris Mansencal and Hugo Boujut who help me in understanding and proceeding all the tasks. They have reviewed and given me many useful comments on my report.

I also thank other members of Video Analysis and Indexing research group, Laëtitia Letoupin and Svebor Karaman, who help me in solving some technical issues.

Finally, I want to express my gratitude to Professor Anne Dicky, who is very kind to me. She was one of my teachers when I studied the master course at PUF.

Introduction

According to physiological and psychological evidences, we know that human beings do not pay equal attention to all exposed visual information, but only focus on certain areas known as focus of attention (FOA) or saliency regions [6]. The LaBRI's researchers have proposed a novel objective quality assessment metric using spatio-temporal saliency map [3]. In the literature, the saliency of the visual scene is characterized by two saliency maps called *spatial* and *temporal* saliency maps. The spatial saliency map is based on color or contrast, for example. The temporal saliency map is computed with the residual motion in the visual scene with regard to global and camera motion. And the spatio-temporal saliency map is the result of the fusion step.

As a part in the large project of Video analysis and indexing group, the aim of my work is to generate semantic saliency map. These results will be integrated into the existing method [3] to assess video quality. Figure 1 gives an overview to the whole process. To enhance this method, the semantic saliency map is generated in two ways, manual way and automatic way. The details of manual way is described in section 2.1 and 2.2. The results of this task will be used in automatic way which is described in 2.4. Section 2.3 will show a statistics of faces in saliency maps.

Below we introduce some important terms used in this report.

Saliency Map is used for representing the saliency areas in an image and guiding the selection of attended locations, based on the spatial distribution of saliency [4].

Gaussian Map is a gray-scale saliency map which assigns a saliency value to each pixel. This value is calculated by using the formula of Gaussian function. This function will be introduced in section 2.2.3.

Error Map is a gray-scale image. The white areas in this map represent the error locations in a distorted image. This map is produced by video decoder in processing systems.

Bounding Box is a rectangle which surrounds an object of interest. In context of this study, we are only concerned by bounding boxes around human faces. The purpose of these bounding boxes is to represent semantic saliency regions in an image. The bounding boxes may be generated manually by using BBBox Annotation Software, introduced in section 2.1, or automatically by face detection algorithm, introduced in section 2.4.

Figure 2 illustrates an example of bounding boxes and the corresponding Gaussian map, and figure 3 gives an example of error map and ground-truth saliency map produced by eye-tracker. In Gaussian map, the white areas show the salient locations in the image. The pixel value becomes whiter when the point goes closer to the center of bounding box. This kind of map and the saliency map from eye-tracking device show visual areas which actually attract viewers' attention. However, these maps are generated manually by human or supporting device. So it is not effective for automatic video quality assessment. The purpose of using these maps is to compare with the quality assessment metric using automatically produced saliency maps.

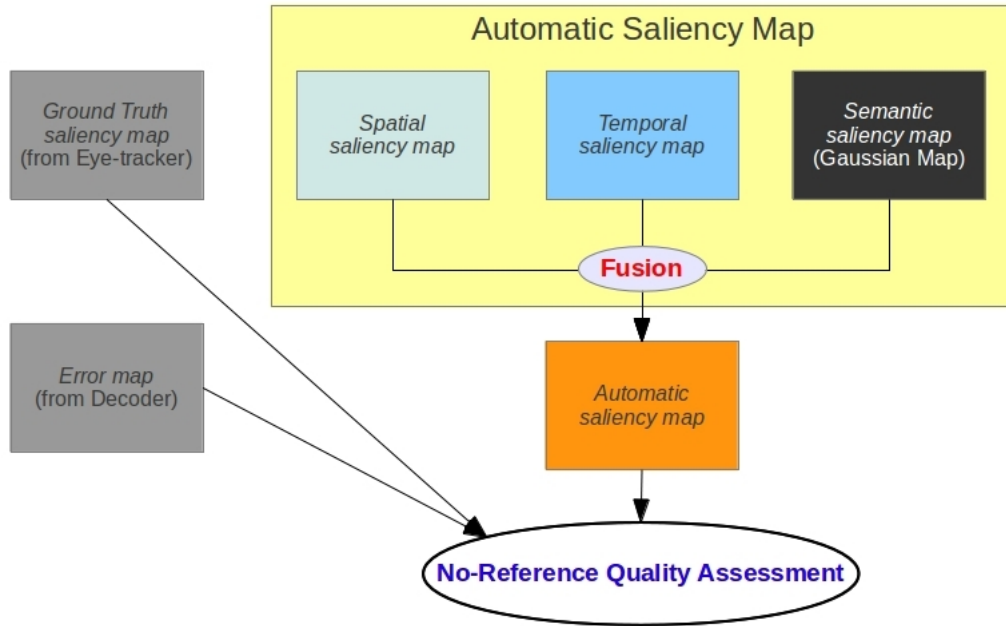
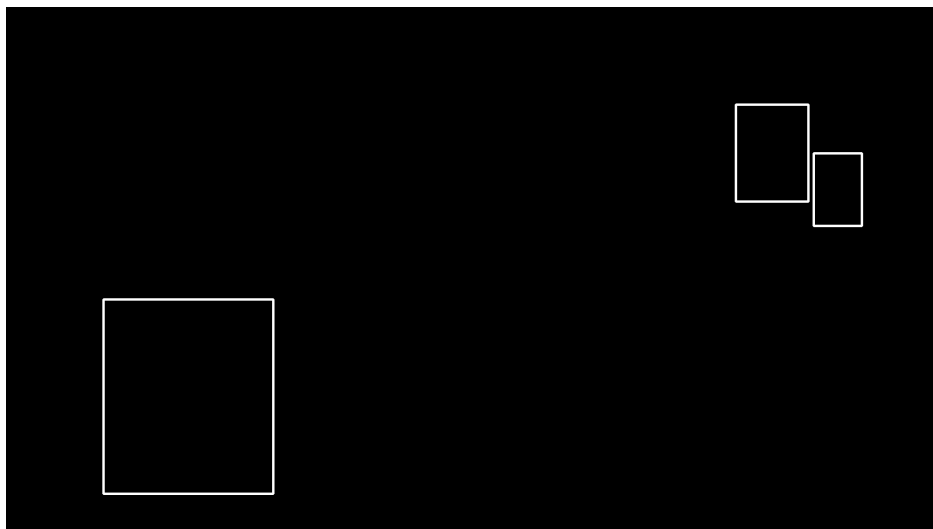
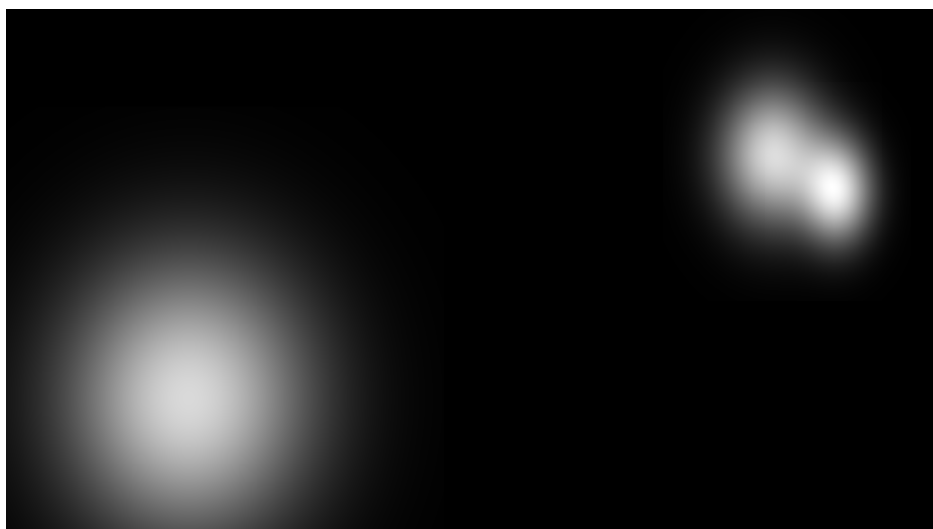


Figure 1: Process Overview



(a) Bounding boxes



(b) Gaussian map

Figure 2: Example of Bounding boxes and Gaussian map



(a) Error map



(b) Ground-truth saliency map

Figure 3: Example of Error map and Ground-truth saliency map

We will introduce some basic concepts of video quality assessment¹.

Subjective video quality assessment The idea of this assessment is: video sequences are shown to the group of viewers and then their opinion is recorded and averaged to evaluate the quality of each video sequence.

Objective video quality assessment are mathematical models that approximate results of subjective quality assessment, but are based on criteria and metrics that can be measured objectively and automatically evaluated by a computer program.

Objective methods are classified based on the availability of the original video signal, which is considered to be of high quality (generally not compressed). Therefore, they can be classified as Full Reference Methods (FR), Reduced Reference Methods (RR) and No-Reference Methods (NR) [3].

FR Methods compute the quality difference by comparing every pixel in each image of the distorted video to its corresponding pixel in the original video.

RR Methods extract some features of both videos and compare them to give a quality score. They are used when all the original video is not available, e.g. in a transmission with a limited bandwidth.

NR Methods try to assess the quality of a distorted video without any reference to the original video.

¹http://en.wikipedia.org/wiki/Video_quality

Chapter 1

Context

This report presents the work that I have done during the internship at the LaBRI. The internship which is a part of Master of Informatics program lasts six months from July 2011 to December 2011. In this section, PUF, LaBRI, and Video Analysis and Indexing group are introduced.

1.1 PUF

This is the place where I studied the master course. PUF¹ (Pôle Universitaire Français) is a center of French universities in Vietnam. PUF in Hanoi and Ho Chi Minh City were established by the cooperation between the French and Vietnamese governments. After five years of educating, PUF has proved the excellent education quality to many generations of students. Graduated from the best universities of France, many students have the opportunities to find good jobs or to continue studying at higher degrees.

There are many academic programs from bachelor to master degree as follows:

- Bachelor in Economics and Management
- Bachelor of Informatics
- Master of Biotechnology
- Master of Business and Economics
- Master of Informatics
- Master of Science In Applied Mathematics

¹<http://www.pufhcm.edu.vn>

1.2 LaBRI

The LaBRI² (Laboratoire Bordelais de Recherche en Informatique) is a French research laboratory in field of computer science. It is associated with the CNRS (UMR 5800), the University of Bordeaux 1, the IPB and the University of Bordeaux 2. It has been the partner of INRIA since 2002. In March 2011, it had a total of more than 300 members including 103 teaching/research staff, 36 research staff, 22 administrative and technical staff and more than 146 doctoral students and post-docs.

The members of the laboratory are divided in six departments:

- Combinatorics and Algorithms
- Image and Sound
- Languages, Systems and Networks
- Formal Methods
- Models and Algorithms for Bio-informatics and Data Visualization
- Supports and Algorithms for High Performance Numerical Applications

1.3 Video Analysis and Indexing Group

I worked as a trainee in Image and Sound team. There are four research groups inside Image and Sound team:

- Structuring and analysis of images
- 3D modeling, visualization and interaction
- Video analysis and indexing
- Modeling of sound and music

My internship took place in the Video Analysis and Indexing group³ (for short AIV group). The Video analysis and indexing subject deals with three major topics. The first topic concerns the segmentation of animated images. The second topic concerns macro-segmentation of video documents into sequences and scenes with the double objective of indexing and navigation.

²<http://www.labri.fr>

³<http://www.labri.fr/projet/AIV/projets.php>

The third topic concerns cross-media segmentation, that is the joint analysis of video and audio flow in digital audio-visual content.

My subject relates to the topic of video quality assessment. The traditional ways of evaluating video quality are based on the metrics like signal-to-noise ratio (SNR) and peak signal-to-noise ratio (PSNR). However, these measures do not match well with subjective perception of human beings. Therefore, many researchers have attempted to integrate human vision system into quality metric designing.

Chapter 2

Visual Saliency Detection

The goal of my training placement is to create “semantic saliency” maps of HD video which express the presence of a semantic object–“face”. We thus present related tasks below:

1. Adapt a software of face annotation for annotating HD videos.
2. Annotate the ground truth: trace the bounding boxes of faces and store their coordinates observed in video frames by human.
3. For manually annotated bounding boxes of faces, we compute the saliency maps with the help of Gaussian filtering.
4. Detect faces automatically by using Viola-Jones face detector of OpenCV library.

In the following sub-chapter, we will thus describe our contribution into adaptation of face-annotation GUI.

2.1 BBBox Annotation Software

To annotate videos, we used the BBBox software developed in LaBRI on the basis of Viola–Jones object detection algorithm. This application is written in C++, and it uses Qt¹ framework to build graphical user interface (GUI). The aim of this task is to annotate high-definition videos. In this task, we have to work with Qt framework which is described in section 2.1.1.

¹<http://qt.nokia.com>

2.1.1 Introduction of Qt Framework

This section introduces Qt Framework, a free and open source software. Qt is a cross-platform application framework that is widely used for developing software with a GUI and non-GUI programs. At first, Qt is produced by Trolltech, an Norwegian company, in May 1995. Currently, it is developed by Nokia's Qt Development Frameworks division after Nokia's acquisition of Trolltech. Since Trolltech's birth, Qt has become a product used by thousands of customers and hundreds of thousands of open source developers all around the world. The community licenses Qt under both open source licenses (LGPL and GPL), as well as a commercial license.

Using Qt Framework, developers can build C++ applications that run natively on many platforms such as Windows, Linux/Unix, Mac OS X, and embedded Linux, without making source code changes. For more information, please refer to the book *"C++ GUI Programming with Qt 4"*[1].

2.1.2 Functionality of BBBox software

BBBox is used for annotating the human faces in each frame of a video. The figure 2.1 shows the screenshot of this software.

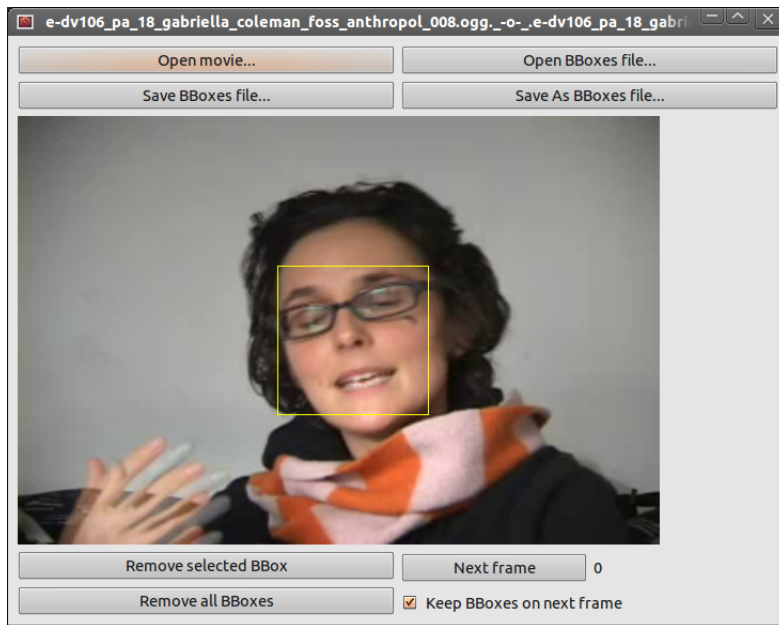


Figure 2.1: Screenshot of BBBox

On GUI of BBBox, there are some buttons with the following functionality:

- Button "Open movie...": select a video to annotate
- Button "Open BBoxes file...": open annotating result
- Button "Save BBoxes file...": save the current annotating result
- Button "Save As BBoxes file...": save as another annotating file
- Button "Remove selected BBox": remove the desired bounding box in current frame
- Button "Remove all BBoxes": remove all bounding boxes in current frame
- Button "Next frame": move forward one frame

To annotate a video, at first user opens it. The frame number, which begins at 0, will be displayed next to the "Next frame" button. Then by dragging mouse, user can draw the bounding boxes which covers human faces. It is also possible to move the bounding boxes by clicking into the area of one box and dragging it to new coordinate. In one frame, user can also draw as many boxes as they want.

After annotating some frames of the video, user can save the current result to a file. We call this file BBoxes file later on. To review the annotating result or continue with annotation, user can open the previous BBoxes file by using the "Open BBoxes file..." button. BBoxes file is just a text file which contains one line per processed frame. Each line of text has the format as follows.

```
frameNb    nbBBoxes    xMin    yMin    width    height
```

The above information describes the frame number containing bounding boxes (*frameNb*), the number of boxes (*nbBBoxes*), the position of top-left corner (*xMin*, *yMin*), and the dimension of box (*width*, *height*) respectively.

2.1.3 Missing Features of BBBox software

Suppose that user open a full-HD video with display resolution (1920 x 1080) bigger than that of desktop screen – what will happen? The problem with the actual version of BBBox is that it cannot adjust the window size to fit

board and other events from the window system, and paints a representation of itself on the screen. A widget that is not embedded in a parent widget is called a window. In Qt, `QMainWindow` and the various subclasses of `QDialog` are the most common window types. An application's main window is created by subclassing `QMainWindow`.

The `QLabel` class is used for displaying a text or an image. To display an image, we use the `QPixmap` class as a paint device. In definition of `QLabel` class, there is a property, `pixmap`, which holds the label's pixmap. This property is modified through the `setPixmap` method.

The `ImageArea` class is both a subclass of `QLabel` and a child widget of `MainWindow`. This widget will receive mouse and resize events from user. The way which `ImageArea` object handles events is similar to that of Java Swing². Whenever a user clicks a button or adjusts window size, many objects in the application may need to react to the change. For example, when a user clicks the "Next frame" button, the `ImageArea` object will have to update the new frame on its graphical representation. To implement this, the BBBox software uses the Observer pattern³. The structure of this pattern is illustrated in figure 2.3 taken from Wikipedia. Figure 2.4 shows how this pattern is used in this software.

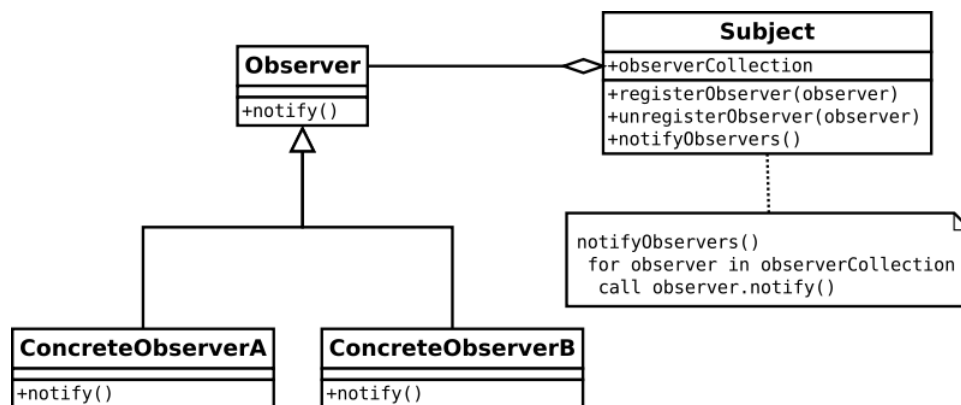


Figure 2.3: UML diagram of Observer pattern

²<http://docs.oracle.com/javase/tutorial/ui/overview/intro.html>

³http://en.wikipedia.org/wiki/Observer_pattern

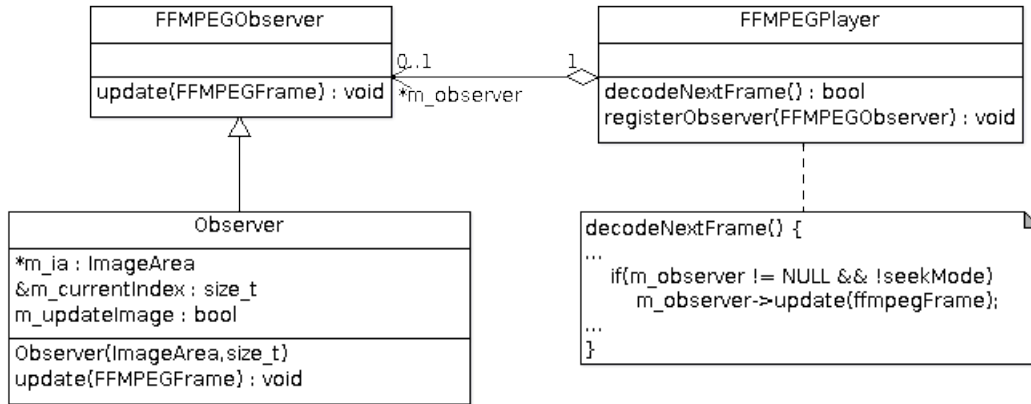


Figure 2.4: UML diagram of Observer pattern used in BBBox

Besides the classes listed in figure 2.2, this software also uses some classes in another library called *libVideoFeatureBridge*. This library is developed by members in the AIV group, Hugo Boujut and Boris Mansencal. Some classes of *libVideoFeatureBridge* used in this software are **FFMPEGPlayer**, **FFMPEGObserver**, and **FFMPEGFrame**.

The word “AABBox” stands for Axis-Aligned Bounding Box. AABBox always has the rectangle shape; however, BBox (Bounding Box) may have the parallelogram shape. That is the difference between these two kinds of bounding boxes. Later on we only use the term “BBox”, “bounding box”, or “box” to mention Axis-Aligned Bounding Box. The **AABBox** class defines the attributes of a bounding box and some utility methods.

In summary, the **MainWindow** class takes charge of building user interface, receiving events, and event-handling. The **ImageArea** class is in charge of displaying each frame on main window and dealing with events. There may be zero or non-zero human face(s) on each frame. The **ImageArea** object, therefore, holds a list of **AABBox** objects.

2.1.5 Modification of BBBox software

As mentioned in section 2.1.3, there are two problems with the actual version of BBBox. To solve the first problem, the flowchart in figure 2.5 is carried out. This is also the flowchart of the `update(FFMPEGFrame)` method in **Observer** class.

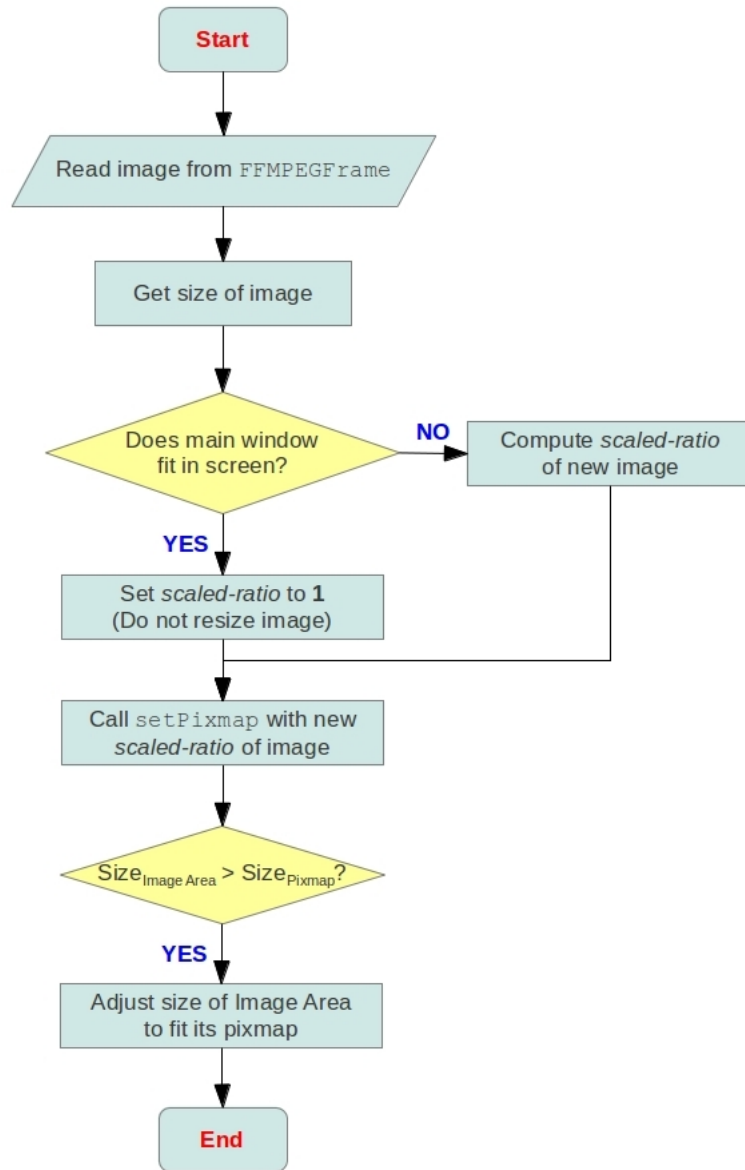


Figure 2.5: Flowchart for solving the first problem

In order to check if main window fits in desktop screen, we need two input parameters. They are dimensions of desktop screen and image which will be shown on `ImageArea`. One difficulty that we had in solving this issue is that we have to calculate main window size before calling `setPixmap()` method on `ImageArea`. Suppose that we do in inverse way, i.e. this method is called before the computation of main window size. What will happen in this case? This will make main window bigger than screen. We have overcome this ob-

stacle by computing approximately the main window size. This computation is done with dimension of displayed image, button height, height of title bar, spacing and margin values of layout.

The flowchart in figure 2.6 describes the steps for solving second problem. It is also the flowchart of the `resizeEvent(QResizeEvent)` method in `ImageArea` class. In this method, we need to scale not only the `pixmap` of `ImageArea` but also the bounding boxes displayed on this `pixmap`. The reason is that there are two different coordinates of a bounding box. One coordinate is used for displaying on the scaled `pixmap`, and the other is used for storing in `BBoxes` file. The latter one is the `BBox` coordinate on original video scene (`pixmap`). The task of transforming `BBox` coordinate between original size and displaying size is also used in the accessor methods of `ImageArea` class, `getBBoxes()` and `setBBoxes()`. The `MainWindow` class gets and sets `BBox` coordinates through these methods. Thus, the displaying size of `BBox` takes effect only in the `ImageArea` class, and its original size takes effect in the `MainWindow` class.

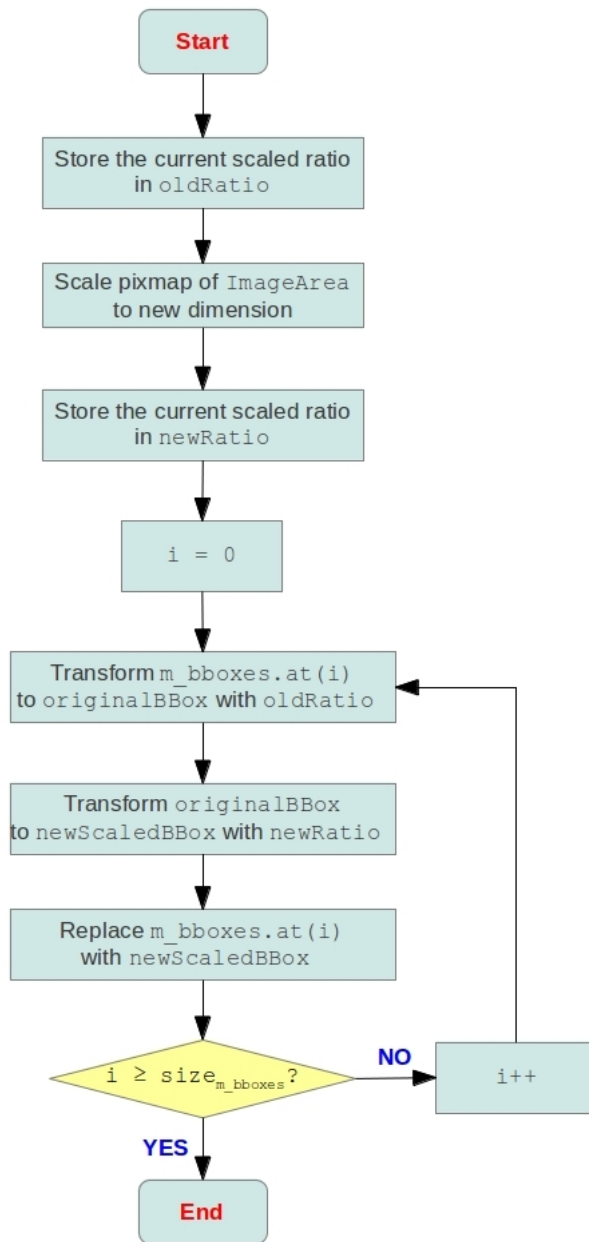


Figure 2.6: Flowchart for solving the second problem

2.2 Ground-truth for Semantic Saliency

Gaussian Map, one kind of saliency map, is generated by using the results of annotating videos and Gaussian function. For this task and the task described in section 2.3, we need to use OpenCV library in processing images. The following part provides some information on this library.

2.2.1 Introduction of OpenCV library

OpenCV⁴ (Open Source Computer Vision) is an open-source library that includes hundreds of computer vision algorithms. OpenCV has many modules. Each module includes several shared or static libraries. Below we introduce some main modules in this library.

- **core:** a compact module defining basic data structures, including the multi-dimensional array **Mat** and basic functions used by all other modules.
- **imgproc:** an image processing module that includes linear and non-linear image filtering, geometrical image transformations, color space conversion, histograms, etc.
- **video:** a video analysis module that includes motion estimation, background subtraction, and object tracking algorithms.
- **calib3d:** basic multiple-view geometry algorithms, single and stereo camera calibration, object pose estimation, stereo correspondence algorithms, and elements of 3D reconstruction.
- **features2d:** salient feature detectors, descriptors, and descriptor matchers.
- **objdetect:** detection of objects and instances of the predefined classes (for example, faces, eyes, mugs, people, cars, etc.).
- **highgui:** an easy-to-use interface to video capturing, image and video codecs, as well as simple UI capabilities.

⁴<http://opencv.willowgarage.com/wiki/>

2.2.2 Annotate Sources of Videos

The aim of this task is to simulate a perfect face detection, and the result of annotating video is one BBoxes file per video. With this result, we can construct perfect semantic saliency maps.

There are several sources of videos. Each one contains many videos of same content, but the quality of videos is different. These videos contain the different amount of errors in frames. We annotate them by using the modified version of BBBox software.

2.2.3 Generate Gaussian Map

Gaussian functions⁵ are widely used in statistics where they describe the normal distributions, in signal processing where they serve to define Gaussian filters, in image processing where two-dimensional Gaussians are used for Gaussian blurs. In this task, Gaussian functions are used for generating saliency map. The definition of this function will be presented in the following section.

Introduction of Gaussian function

In mathematics, a Gaussian function is a function of the form:

$$f(x) = ae^{-\frac{(x-b)^2}{2c^2}} \quad (2.1)$$

for some real constants $a, b, c > 0$, and $e \approx 2.718281828$ (Euler's number).

The graph of a Gaussian has a shape of symmetric "bell curve" that quickly falls off towards plus/minus infinity. The parameter a is the height of the curve's peak, b is the position of the center of the peak, and c controls the width of the "bell".

For creating Gaussian map, we used a particular form of two-dimensional Gaussian function:

$$f(x, y) = Ae^{-\left(\frac{(x-x_0)^2}{2\sigma_x^2} + \frac{(y-y_0)^2}{2\sigma_y^2}\right)} \quad (2.2)$$

In above formula, the coefficient A is the amplitude, x_0, y_0 is the center and σ_x, σ_y are the x and y dimension of the blob. The figure 2.7 and figure 2.8, taken from Wikipedia, give examples of one-dimensional Gaussian functions and two-dimensional Gaussian function.

⁵http://en.wikipedia.org/wiki/Gaussian_function

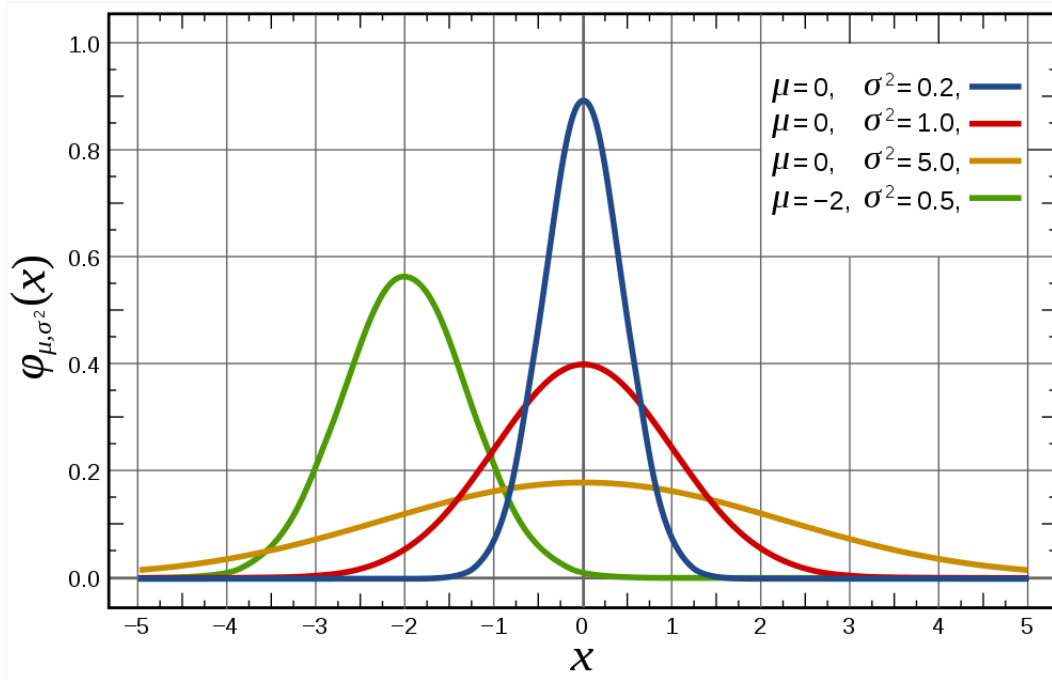


Figure 2.7: 1D Gaussian function

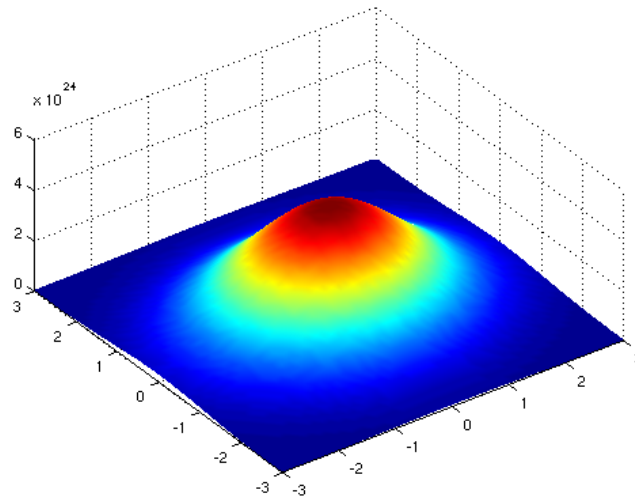


Figure 2.8: 2D Gaussian function

Program Description

The program for generating Gaussian Map uses the result of section 2.2.2 and OpenCV library. The figure 2.9 summarizes the processing steps.



Figure 2.9: Processing steps for creating Gaussian Map

In this program, we created the following classes.

Class name	Description
AnnotatedFrame	Stores a frame number and all coordinates of bounding boxes on that frame
BBoxFileReader	Opens a BBoxes file and creates an object of AnnotatedFrame class from information in each text line

GaussianMap	Computes pixel values by using formula of 2D Gaussian function and outputs a gray-scale image of Gaussian map
-------------	---

Table 2.1: Description of some important classes

For the task of constructing Gaussian map and outputting a corresponding image, we used the matrix structure. The class `Mat` represents a 2D numerical array that can act as a matrix. There are many different ways to create `Mat` object. Here is the way we used inside the `GaussianMap` class:

```

void
GaussianMap::set(const AABBoxCollection &bboxes, int width,
                int height)
{
    if (!bboxes.empty()) {
        // Create two Mat objects
        // m_gaussianMat: outputs image of Gaussian map
        // m_floatMat: stores pixel values in float number
        m_gaussianMat = Mat(height, width, CV_8UC1);
        m_floatMat = Mat::zeros(height, width, CV_32FC1);
        for (int i = 0; i < bboxes.size(); i++) {
            AABBox bbox = bboxes[i];
            if (bbox.width() < m_bboxMinSize) {
                bbox.set(bbox.xMin(), bbox.yMin(), m_bboxMinSize,
                        bbox.height());
            }
            if (bbox.height() < m_bboxMinSize) {
                bbox.set(bbox.xMin(), bbox.yMin(), bbox.width(),
                        m_bboxMinSize);
            }
            setPixelValue(bbox, width, height);
        }
        // Get max of m_floatMat
        float maxMat = 0;
        for(int row = 0; row < m_floatMat.rows; row++)
            for(int col = 0; col < m_floatMat.cols; col++)
                if (maxMat < m_floatMat.at<ELT_TYPE>(row, col))
                    maxMat = m_floatMat.at<ELT_TYPE>(row, col);
        // Normalize
        assert(maxMat != 0);
        m_floatMat.convertTo(m_gaussianMat, CV_8UC1, 255/maxMat);
    }
}

```

Listing 2.1: A method used for constructing Gaussian map

After creating `Mat` objects, we scan the list of bounding boxes, `bboxes`. This list is stored in `AnnotatedFrame` object. For each box, we modify its width and height if they are smaller than a predefined value, `m_bboxMinSize`. Then we call the method `setPixelValue()` with box coordinate and image dimension (`width`, `height`). This method is used for computing all pixel values inside the area of bounding box. The pixel value is calculated on the basis of formula 2.2 with $A = 255$. Then we find the maximum of `m_floatMat` matrix to convert into a matrix of `unsigned char` type (`m_gaussianMat`).

2.3 Statistics of Faces in Saliency Map

Statistics on faces in saliency maps will let us know whether human faces attract viewers' attention. Or how many errors appearing on faces and where the errors happen on frames would actually attract viewers' attention? In [6], the authors assume that an error that appears on a saliency region is much more annoying than an error appearing in an inconspicuous area.

Input data that we use in this task includes:

1. **BBoxes files**: are the result of annotating videos through the BBox software. See section 2.2.2 for more details.
2. **Saliency maps**: are created by using eye-tracking device. They will let us know the areas on one video frame to which viewers pay more attention.
3. **Error maps**: are extracted from video decoder in processing systems. These maps will show the distorted areas on each frame.

By using these kinds of map and faces' locations in BBoxes files, we compute one value per video according to the methods introduced in next section. After that, we plot all these numbers on one chart for comparison.

2.3.1 Computational Methods

There are three different computation methods. You can see the new terms, *Overlapped BBox* and *Overlapped pixel*, in these methods.

Overlapped pixel is the common pixel of both saliency area (white area) and bounding box. Thus, it has positive (> 0) value.

Overlapped BBox is the bounding box which overlaps with saliency area.

To check if one BBox overlaps with saliency area, we read all pixel values of saliency map within area of BBox. If there exists one pixel having positive value, that BBox is counted as a overlapped BBox. In figure 2.10, the yellow BBox on the left does not overlap the saliency area, and the green one on the right is an overlapped BBox.

The calculation of three methods is presented in following sections. For all methods, we compute firstly value f_i for the i -th frame in the video. Then we calculate the mean value of all previously-computed values. The

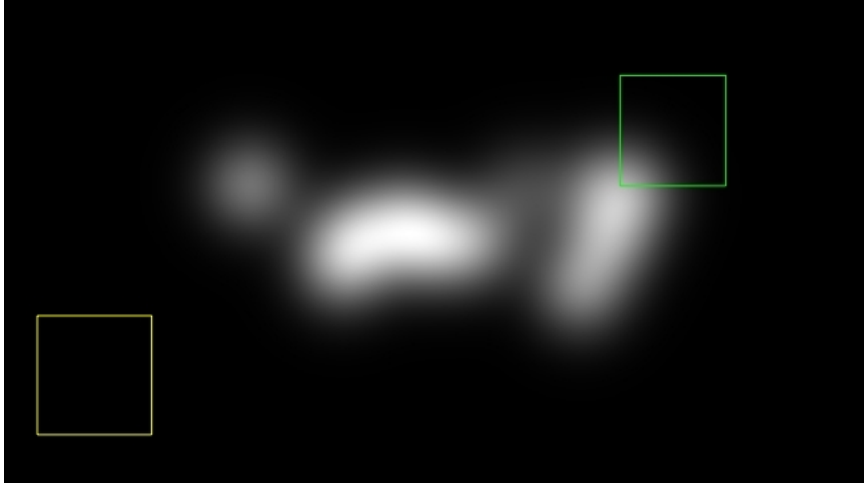


Figure 2.10: Example of overlapped BBox

only difference between these methods is the calculation formula of f_i . The mean value is computed in following equation.

$$f = \frac{\sum_{i=0}^N f_i}{N} \quad (2.3)$$

One remark is that the N number in all equations denotes *the number of processed frames*, not the number of frames contained in a video. This means that we do not process all frames, but we process the frames containing faces. To easily remember, we call method 1, 2, 3 by shortened names: *COB* (Count Overlapped BBoxes), *COP* (Count Overlapped Pixels), and *FOB* (Find Overlapped BBoxes) respectively. These names represent the main task in each computational procedure.

Method COB

This method computes the value f_i for frame i as follows.

$$f_i = \frac{\text{Number of overlapped BBoxes}}{\text{Number of BBoxes on frame}}, i = 0..N \quad (2.4)$$

Method COP

This method computes the value f_i for frame i as follows.

$$f_i = \frac{\text{Total percentage of overlapped pixels}}{\text{Number of BBoxes on frame}}, i = 0..N \quad (2.5)$$

In above formula, we compute percentage p_i of overlapped pixels for one BBox. Then we get the numerator in equation 2.5 by calculating the sum of all p_i values.

$$p_i = \frac{\text{Number of overlapped pixels in one BBox}}{\text{Area}_{BBox}} \quad (2.6)$$

and

$$\text{Area}_{BBox} = \text{width}_{BBox} \cdot \text{height}_{BBox}$$

Method FOB

This method computes the value f_i for frame i as follows.

$$f_i = \begin{cases} 1 & , \text{ if there exists one overlapped BBox} \\ 0 & , \text{ otherwise} \end{cases} \quad (2.7)$$

2.3.2 Statistical Results

This section illustrates computed results in graphs, and the Gnuplot⁶ tool is used for plotting. As mentioned above, we calculate one value per video by using the methods introduced in previous section. There are many sources of videos. We choose one of sources which contains nine videos. The best-quality video is video 9. However, the remaining videos have worse quality, and the amount of noise appearing on one frame in every video is different. After the computation step is done, we plot all these numbers on one diagram for comparison. The statistical results in this part will be used for later comparison in section 2.4.3.

Firstly, through saliency maps, we examine the percentage of faces appearing in the saliency areas. Three methods are used for computing mean values.

Observing the graph in figure 2.11, we see that method COB and FOB have produced the same values. However, the values produced by method COP are slightly different because we compute the percentage of overlapped pixels for one bounding box. With this computational way, we know exactly how many percent of faces overlap with saliency areas. The mean values range from 0.954194 to 1 (100%). This result shows that viewers pay more attention to faces on frame .

⁶<http://www.gnuplot.info/>

Secondly, through error maps, we examine the percentage of faces appearing in the distorted areas. Only method COB and COP are used for computing mean values. These values are plotted on the graph in figure 2.12.

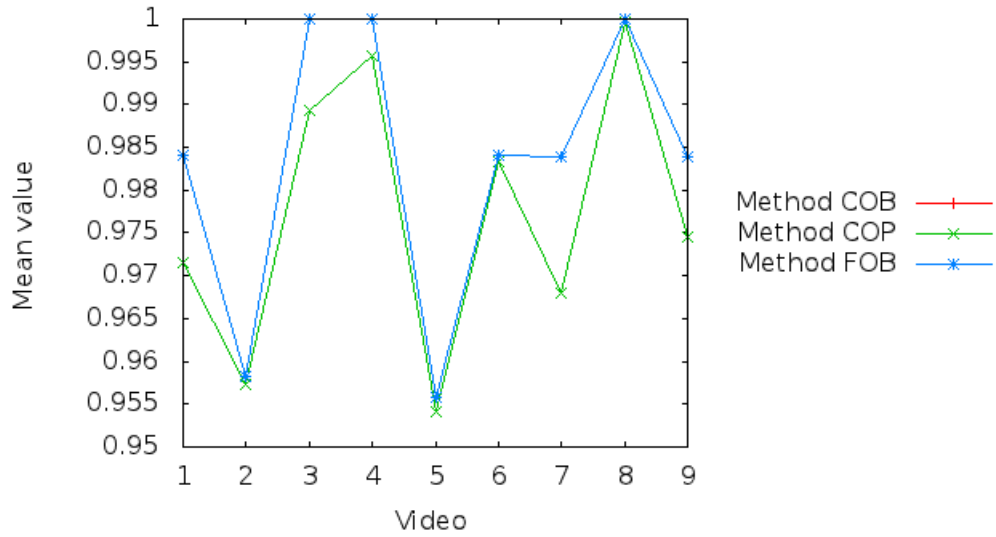


Figure 2.11: Statistics of faces in saliency maps

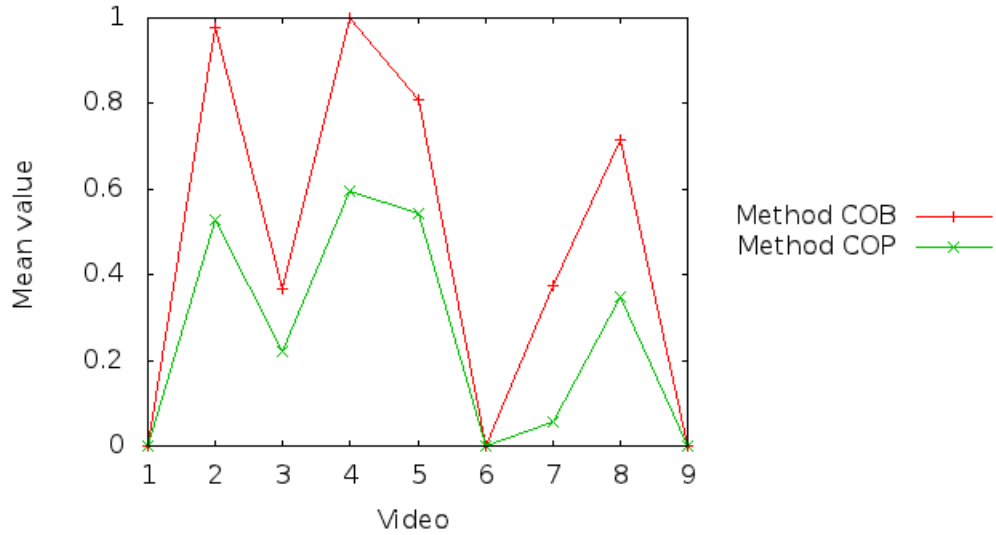


Figure 2.12: Statistics of faces in error maps

Finally, we build a combination of saliency map and error map. This job is performed by calling the `mul` method on `Mat` object as in following code snippet.

```
// Get list of BBoxes on AnnotatedFrame object
AABBoxCollection bboxes = aframe.getBBoxes();
// Performs an element-wise multiplication of
// the two matrices (errorMap & eyeTrackerMap)
Mat A = errorMap.mul(eyeTrackerMap);
// Compute value for one frame
float v = compute(A, bboxes, method);
```

Listing 2.2: A method used for constructing Gaussian map

Using these combined map, we examine how many errors appearing on faces and the percentage of errors would attract viewers' attention. Method COB and COP are used for computing mean values, and these values are plotted on the graph in figure 2.13. The videos having high mean values (video 2, 4, 5, and 8) will make viewers more annoyed than other videos.

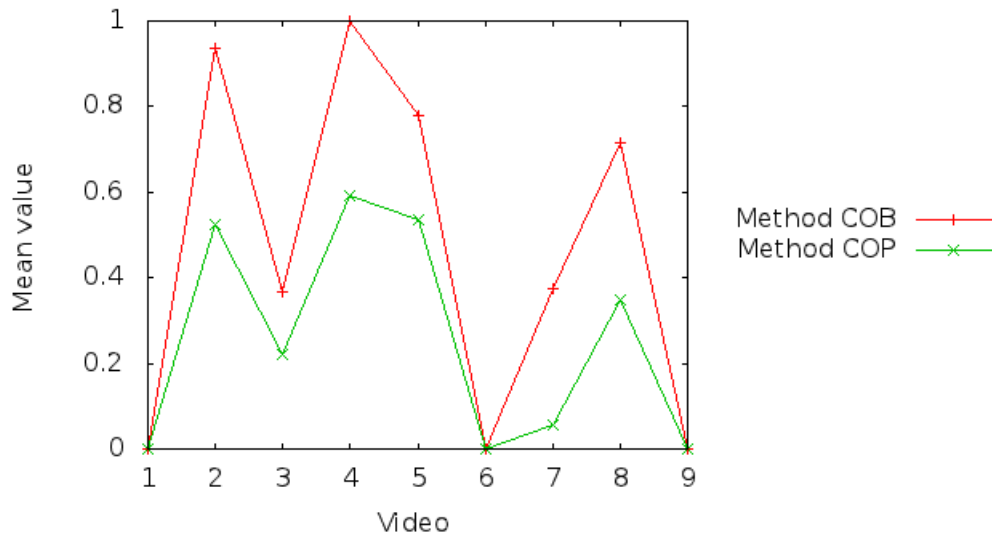


Figure 2.13: Statistics of faces in combination of saliency maps and error maps

2.4 Automatic Face Detection

The purpose of this task is to add automatic face detection for the automatic semantic map construction. We give a brief introduction of Viola–Jones object detection algorithm in below section. This algorithm is implemented in an example program⁷ coming with OpenCV library.

2.4.1 Introduction of Viola–Jones algorithm

This algorithm was proposed by Paul Viola and Michael Jones [5]. The proposed object detection framework provides competitive object detection rates in real-time. There are three main contributions of this framework. We will introduce these ideas briefly below.

The first contribution is a new image representation called an integral image that allows for very fast feature evaluation. Their object detection framework classifies images based on the value of simple features. One critical motivations for using features rather than the pixels directly is that the feature-based system operates much faster than a pixel-based system. With the use of the integral image, rectangular features can be evaluated in constant time.

The second contribution is a method for constructing a classifier by selecting a small number of important features using AdaBoost⁸. In order to select the features and train the classifier, Viola–Jones use a modified version of the AdaBoost algorithm developed by Freund and Schapire in 1995.

The third major contribution is a method for combining successively more complex classifiers in a cascade structure which increases the speed of detector by focusing attention on promising regions of the image. The cascaded classifier is composed of stages. Each one contains a strong classifier. The job of each stage is to determine whether a given sub-window is definitely negative (non-face) or maybe positive. The detection process is shown in figure 2.14. This figure is taken from the article of Viola–Jones.

The strong classifiers are arranged in order of complexity. If a sub-window passes the first classifier, it will be evaluated at the second classifier. A positive result from the second classifier triggers a third classifier, etc. The

⁷File facedetect.cpp located in the folder OpenCV-2.3.1/samples/c

⁸<http://en.wikipedia.org/wiki/AdaBoost>

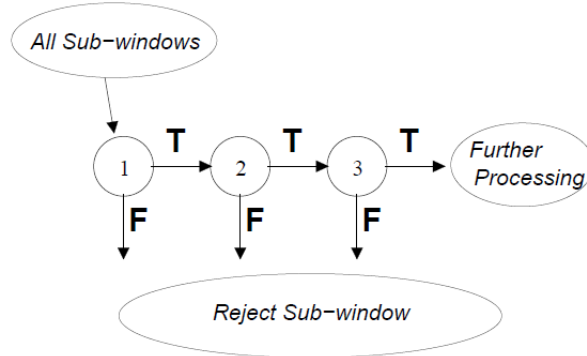


Figure 2.14: Cascade Architecture

initial classifier discards a large number of negative sub-windows with very little processing. Successive classifier discards additional negatives but require more processing. Through this model, the number of sub-windows have been decreased considerably.

The cascade structure reflects the fact that within an image an excessive amount of sub-windows are negative. In stead of finding faces, the algorithm should discard as many negatives as possible at the earliest stage. While a positive instance will trigger the evaluation of each classifier, this event happens rarely.

2.4.2 Gaussian Map Construction

This task is done after we had the result of automatic face detection. Therefore, we had the face coordinate in video frames via automatic detection. For this work, we also use the classes described in table 2.1 and perform the similar processing steps in figure 2.9. By comparing the maps produced in this task and those produced in section 2.2.3, we can know how correct the detection is.

2.4.3 Statistical Results with Automatic Detection

In section 2.3.2, we have shown that faces attract more viewer's attention and thus face is a good idea for semantic maps. The result of this work interests us because we want to know if automatic face detection is good enough, compared to manual ground-truth face detection. In this section, we perform the same experiment that we have done in section 2.3. We also use the similar input data as in previous experiment. The implementation of Viola-Jones

algorithm in OpenCV have been modified in order to produce the BBoxes file. After the face detection is done, we use the produced BBoxes file as input data of computation program. The content of this file is described in section 2.1.2.

One problem with face detection is that a detector trained on frontal faces is unable to detect profile faces. Therefore, instead of using a frontal-face detector or a profile-face detector, we use both of them to detect as many faces as possible. In this case, we will have two BBoxes files generated by both detectors. Thus, we will have to merge these files into one file.

Observing the graph in figure 2.15 and the one in figure 2.11 page 31, we see that the variation of mean values is quite considerable. For all methods, the mean values in the former are greater than in the latter. The percentage of variation is shown in table 2.2.

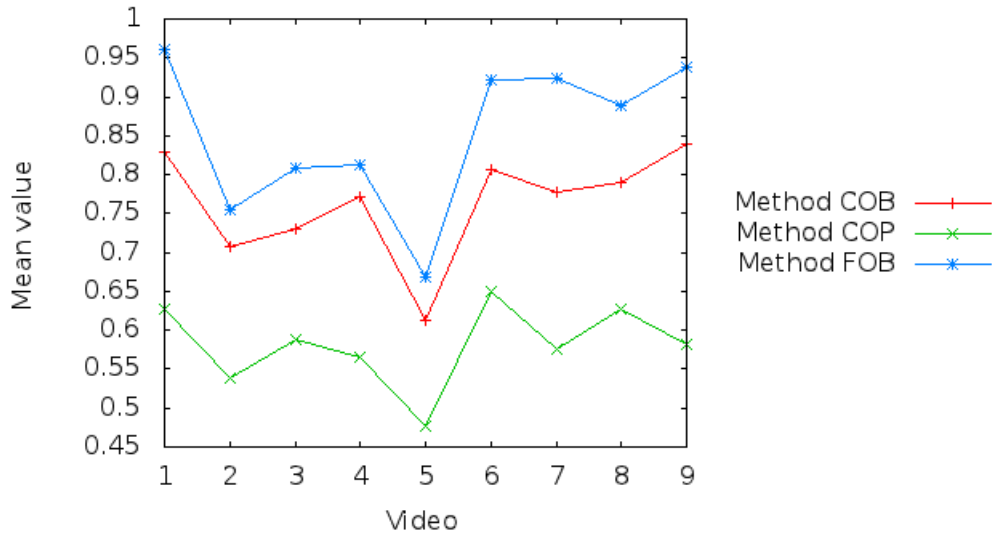


Figure 2.15: Statistics of faces (automatic detection) in saliency maps

Video	Method COB (%)	Method COP (%)	Method FOB (%)
1	15.5218	34.5074	2.269
2	25.1572	41.8283	20.3616
3	27.031	40.1353	19.1257
4	22.9167	43.0406	18.75
5	34.3675	47.6441	28.6868

6	17.8407	33.4175	6.1338
7	20.6969	39.3025	6.0099
8	21.0391	37.1462	11.1111
9	14.4102	39.3511	4.6759

Table 2.2: Percentage of variation in mean values

With the comparison in above table, we can see that the reliability of automatic face detection is not high. By looking at “Method COP” curve in figure 2.15, we see that the mean values approximately range from 0.45 to 0.65. In figure 2.11, these values approximately range from 0.95 to 1. These ranges let us know that sometimes the detection is false. Using the results of unreliable detection, we will compute the lower mean values. For the two remaining statistics, the mean values are slightly different from those in previous experiment (manual detection).

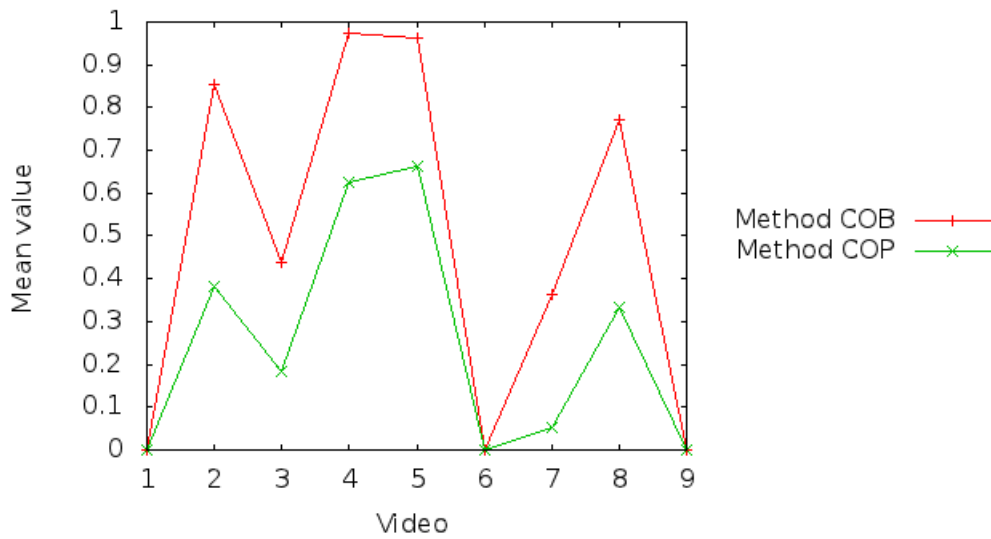


Figure 2.16: Statistics of faces (automatic detection) in error maps

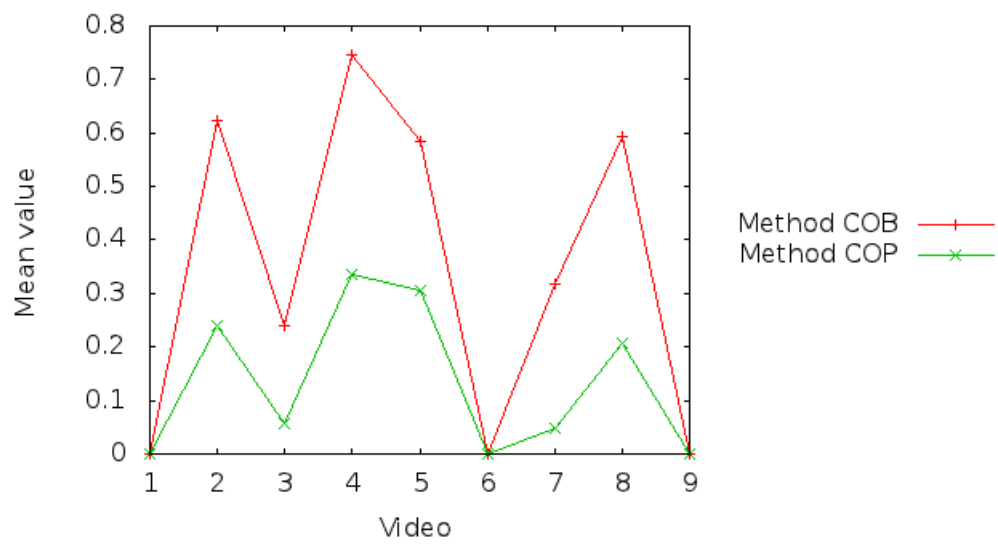


Figure 2.17: Statistics of faces (automatic detection) in combination of saliency maps and error maps

Chapter 3

Conclusion and Future Work

This chapter gives a summary of what I have done in this project and provides future work for improving the current objective quality assessment metric.

3.1 Conclusion

The purpose of my internship was to study if a semantic saliency, based on faces, could improve this metric. This goal has been achieved by:

- **Modifying BBBox software:** BBBox is used for annotating the human faces in each frame of a video. There are some problems with the actual version of BBBox. This task must be done so that this software is able to annotate High-definition videos.
- **Annotating videos:** I use the modified BBBox software to perform this task. The results are BBoxes files which are created by BBBox software.
- **Generating Gaussian map:** This map is generated by using the results of annotating videos and the formula of Gaussian function. It lets us know the positions on a video frame where human faces are located in. In Gaussian map, the white areas show the salient locations in the image.
- **Examining statistics of faces in saliency map:** Using the methods, we compute the values for all videos and plot all these numbers on one diagram for comparison. This statistics will let us know if human faces attract viewers' attention. Or how many percent of errors appearing on faces and where the errors happen on frames would actually attract viewers' attention?

- **Automatic Face Detection:** This task helps us detect human faces automatically. The produced results are the BBoxes files containing locations of faces.
- **Examining statistics of automatically-detected faces in saliency map:** We are interested in the output of this task. Based on this output, we can compare with the output in previous statistics. As a result, we can evaluate the effectiveness of using the automatic detection.

3.2 Future Work

In this part the word “we” denotes the members in AIV group (excluding me) who will perform the future work.

Although I have used both frontal-face detector and profile-face detector, the result of face detection is not as good as we desire. Sometimes the algorithm detects other objects (not faces), and sometimes it cannot determine the faces. In the future, we plan to improve the automatic detection by using face tracking algorithm. I have already made some programs to test the CAMSHIFT object tracking algorithm [2]. This algorithm is implemented in an example program¹ coming with OpenCV library. One problem with face tracking is that it still keeps tracking although tracked faces do not appear in current frame. Moreover, it also detects other objects if those ones have similar color as faces. One idea of improving tracking algorithm is that we make a stop condition. This condition will terminate the tracking if the current frame does not contain the detected faces in previous frame. We also have to deal with another problem. Suppose that in current frame, the detected face is not located near the old position in previous frame. Therefore, the tracking should be correct in this case.

Another part of future work is to evaluate if the automatically-produced saliency maps, including semantic information, allow better video quality assessment.

¹File camshiftdemo.cpp located in the folder OpenCV-2.3.1/samples/cpp

Bibliography

- [1] Jasmin Blanchette and Mark Summerfield. *C++ GUI Programming with Qt 4*. Prentice Hall PTR, Upper Saddle River, NJ, USA, 2006.
- [2] Gary R. Bradski. Computer Vision Face Tracking For Use in a Perceptual User Interface. 1998.
- [3] H. Boujut, J. Benois-Pineau, O. Hadar, T. Ahmed, and P. Bonnet. Weighted-MSE based on Saliency map for assessing video quality of H.264 video streams. *IS&T / SPIE Electronic Imaging, San Francisco: United States*, 2011.
- [4] L. Itti, C. Koch, and E. Niebur. A Model of Saliency-Based Visual Attention for Rapid Scene Analysis. *IEEE Trans. on Pattern Analysis and Machine Intelligence*, 20(11):1254–1259, 1998.
- [5] Paul Viola and Michael Jones. Robust Real-time Object Detection. In *International Journal of Computer Vision*, 2001.
- [6] X. Feng, T. Liu, D. Yang, and Y. Wang. Saliency Based Objective Quality Assessment of Decoded Video Affected by Packet Loss. In *ICIP*, pages 2560–2563, 2008.