# Inner Lips Parameter Estimation based on Adaptive Ellipse Model

*Li Liu[1,2], Gang Feng[1,2], Denis Beautemps[1,2]*

[1]Univ. Grenoble Alpes, Gipsa-lab, F-38000 Grenoble, France
[2]CNRS, Gipsa-lab, F-38000 Grenoble, France

`li.liu@gipsa-lab.grenoble-inp.fr,`
`feng.gang@gipsa-lab.grenoble-inp.fr`

## Abstract

In this paper, a novel automatic method using an adaptive ellipse model to estimate inner lips parameters (inner lips width *A* and height *B*) of speakers without any artifices is presented. Color based image processing is first applied to segment preliminary inner lips. A single discontinuity elimination combining horizontal and vertical filling are used to obtain a binary inner lips image as complete as possible. After the previous pre-processing steps, an optimal adaptive ellipse is determined to match the inner lips. The proposed method is evaluated on 4693 images of three French speakers including one Cued Speech (CS) speaker. It obtains RMSE of 3.37 mm for *A* parameter and of 0.84 mm for *B* parameter which outperform the baseline of inner lips parameter estimation in the state of the art. Moreover, CS recognition based on 34 French phonemes shows that using the estimated two parameters achieves an accuracy which is comparable to that using raw lips ROI.

**Index Terms**: Adaptive Ellipse Model, image processing, inner lips contour parameter, Cued Speech.

## 1. Introduction

Lips detection is an active research topic which is applied in many multimedia applications such as speech reading, speech production, audiovisual communication and robot speech application. Cued Speech (CS) [1-3] is a complement of lip reading to enhance speech perception from visual input including lips and hand. It is used by deaf people or hearing people communicating with deaf orally educated. This paper especially estimates inner lips width parameter A and height parameter B [4-6] in the CS case, in which lips may be occluded by hand. Moreover, the method is also used in non-CS case.

Lips (especially inner lips) tracking and parameter estimation remains difficult tasks. The variability of lips shape, extreme lips movements with the interference of tongue/teeth, inconstant lighting condition and hand occlusion (in CS case) are factors which can decrease the precision and robustness of the model. Several approaches to extracting lips in speech processing have been investigated in the literature. One of the most widely used techniques is model-based lips detection. Active Shape Model (ASM) [7] and Active Appearance Model (AAM) [8] were proposed to segment lips contour. Shape and appearance of lips are learned from training data with manually annotated faces. Lips configurations are described by a set of model parameters. Bandyopadhyay [9] investigated a lips feature extraction technique combined with ASM and used the contrast between lips and face to segment lips contour. Large training set and good initial condition are necessary for model-based technique. In [10], a method based on a statistical model of shape with local appearance Gaussian descriptors is proposed. Some authors used parametric models in which the lips shape is described by curves controlled by a limited set of parameters. Stillitano *et al.* [11] used both active contours and parametric models for lips contour extraction. This method needs prior knowledge of the lips shape. Henecke et al. [12] used a deformable model, but the lack of flexibility of the template can be a problem. Another technique is based on segmentation in color space [13, 14]. Color-based clustering assumes that there are only two classes, i.e., skin and lips, and this technique may not be efficient if facial hair or teeth exist. Also, it needs controlled lighting condition and good contrast between the color of lips and skin. Currently, Convolutional Neural Network is very popular in the feature extraction field. Hlavac presented a CNN method for lips landmarks detection in [15], which achieves a sub-pixel accuracy in landmarks detection, but some errors remain since no robust features around chin can be locked.

In the previous CS study of lips parameter estimation, Heracleous et al. [1] extracted the lips A and B parameters by painting blue color on subject's lips and tracking the blue lips. A dynamic correlation template method based on Constrained Local Neural fields (CLNF) [16] in our previous work (without any artifice color marks) [17,18] is investigated to estimate B parameter. CLNF is proposed by Baltusaitis *et al.* in 2013, which is robust for facial landmark detection in the wild. It is a novel instance of Constrained Local Model (CLM) [19] that deals with the issues of feature detection in complex scene. In [34], a 3D lips model from thirty control points for any lips shape was presented. A and B parameters precision are reported in their study, mean error 3.5 mm (standard deviation 4.5 mm) was obtained for A parameter and mean error 1.0 mm (standard deviation 1.0 mm) was obtained for B parameter.

In this paper, we explore a new method based on the adaptive ellipse model to estimate A and B parameter of inner lips without generating a whole inner lips contour. In Section 2, experiment database is introduced. Section 3 describes each step of the proposed method. Firstly, an image processing is realized to segment inner lips as much as possible. To make the extracted inner lips more complete and connected, a single discontinuity smoothing and a horizontal vertical filling are applied. Then, an adaptive ellipse is used to match the inner lips and gives the best *A* and *B* parameters. An outline of this process is shown in Figure 1. Evaluation and results on the accuracy and robustness of the method are presented in Section 4. The proposed method is evaluated on 4693 images from three subjects. RMSE of 3.37 mm is reached for *A* parameter and RMSE of 0.84 mm is obtained for *B* parameter. Compared with CLNF which is an efficient facial (including inner lips) landmark detection model developed in computer vision, our

1

method significantly outperforms CLNF concerning inner lips parameter estimation. As one application, Principle Components Analysis (PCA) is used to extract good features on raw lips ROI and our parametric ellipse lips ROI. 34 French phonemes recognition is performed using both features. The recognition accuracy also confirms the high precision of the estimated $A$ and $B$ parameters.
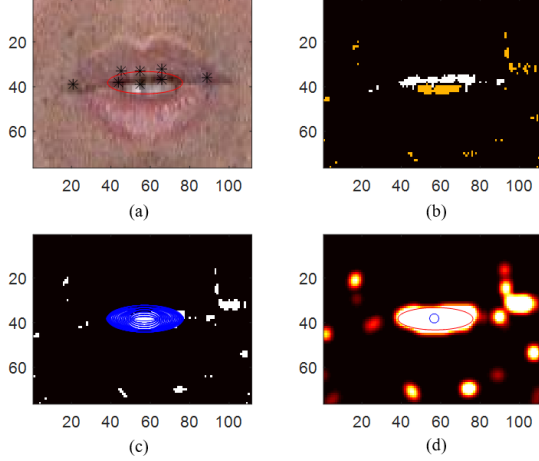


Figure 1. *Overview of adaptive ellipse model for inner lips parameters estimation. (a). raw lips image in ROI, with the optimal inner ellipse shown in red. Black stars are the inner lips landmarks given by CLNF (for comparison). (b). extraction of dark area (white region) and teeth (yellow region) using image processing. (c). Adaptive searching for optimal position and size of the ellipse. (d). The final optimal ellipse determined after smoothing and scaling post-processing.*

## 2. Database

The database contains videos of 50 French words made of numbers and daily words. Each corpus is uttered 10 times by 3 French subjects: one female CS speaker and two male speakers. The recording (RGB images 576*720*3, 50fps) is made in a sound-proof booth in Gipsa-lab, France. Words are annotated with Praat based on speech sound signal. We use the first repetition of the three speakers corresponding to all types of lips shape to evaluate $A$ and $B$ parameters. The image of our database for three subjects are 1377, 1744 and 1572, respectively (total 4693 images). To evaluate the performance of the proposed model, the ground truth inner lips contour is extracted manually by an expert placing several landmarks on lips. In the application of 34 CS phonemes [6] recognition, the temporal boundaries of each phoneme were extracted from the audio signal using a conventional Automatic Speech Recognition system and a forced-alignment procedure.

## 3. Methodology

The main objective of the proposed method is to realize a robust and efficient estimation of $A$ and $B$ parameters, without extracting the real inner lips contour. In fact, we try to adapt an ellipse to the inner lips region so that the parameters of this ellipse are a good descriptor of $A$ and $B$ parameters. We propose

an adaptive algorithm to optimize the fitting between the ellipse and the inner lips region.

### 3.1. Image process for segmenting inner lips area

Firstly, a lips ROI is extracted based on CLNF which uses the neural network layer and convolution kernel to track facial landmarks (including eight inner lips landmarks) even with hand occlusion. In order to estimate inner lips parameters, the most important thing is to determine the inner region of lips. More precisely, we propose to first extract the inner region of the lips instead of directly finding the lips. As we can see, teeth and the darker area inside inner lips have different color properties with lips (see Figure 2), a color-based method is applied to segment inner lips and non-inner lips [11, 20 and 21] as much as possible.

In YCbCr space, the dark area has a lower luminance, a threshold of Y value can be used to distinguish it. In our experiment, a threshold of 70 gives satisfactory performance. In RGB color space, teeth have a much whiter color than non-teeth regions. In comparison, lips have a red-dominant color component. We have found a ratio R/G < 1.25 (coupled with G+B>160) permits to extract teeth efficiently. It should be noted that the thresholds may be varied a little bit for different subjects. The final performance of inner lips segmentation is illustrated in Figure 2. In fact, the color of the tongue has very similar color as lips and skin, it causes difficulties to segment tongue and extract a complete inner lips region only by the color-based approach.

After this preliminary image processing procedure, the pixel value of the detected inner lips region is set to 255 (white pixel) and the other regions set to 0 (black pixel) inside lips ROI.
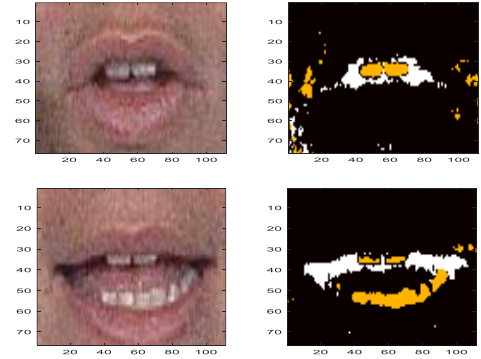


Figure 2. *Teeth and dark area extraction. Left: lips ROI. Right: the white part is the detected dark area, and the yellow area is the detected teeth. Note that tongue are not detected.*

### 3.2. Single discontinuity smoothing and horizontal vertical filling

The previous image processing allows extract the teeth and some dark area insides inner lips. However, it is still not enough to form a whole inner contour since tongue and fuzzy invisible teeth are not detected. Matlab function *imfill* can be used to reduce the discontinuity only when it is fully surrounded by while pixels. It is not suitable when the extracted region is not connected. In order to solve this problem, we propose two methods: a single discontinuity smoothing and a horizontal vertical block filling.

The single discontinuity smoothing aims at eliminating single pixels which are the pixels insides inner lips but not be detected. In the following ten cases (Figure 3), the central pixel (the single rupture) is set to be "white". This procedure is implemented from top to bottom and then from left to right.

In some case, one or several blocks of the black pixel can still remain after the single discontinuity filling (the second case of Figure 2 gives an example of this phoneme). In order to solve this problem, a horizontal vertical block filling is proposed. Take the horizontal block filling as an example (Figure 4). We first examine from top to bottom of each rows in the lips rectangle ROI. In each row, the line is divided by several black intervals separated by white intervals. If a black interval length is less than the sum of two adjacent white interval ($b_i \ll l_i + l_{i+1}$ in Figure 4), the black interval will be filled as a white interval. The vertical filling uses the same principle with the horizontal filling. Using these two methods, most of the block-ruptures can be filled efficiently.
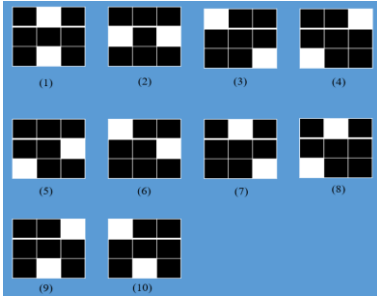


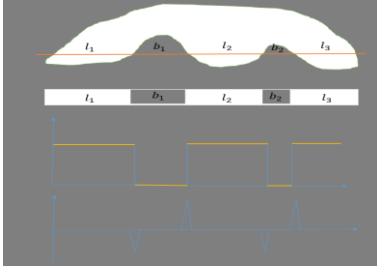Figure 3. *Ten cases of single discontinuity filling.*



Figure 4. *Horizontal filling. The vertical filling has the same principle with the horizontal case.*

### 3.3. Adaptive Ellipse model

After the previous pre-processing steps, a relatively well filled inner lips region is obtained. It is a binary image where the detected inner lips region is white and other area is black. However, uncontrolled lighting condition when the hand gets close to face, highly deformable of lips especially when the subject is speaking and the appearance of tongue (which has almost the same color as skin) inside the mouth cause great difficulties to extract a complete inner lips region only by the color-based approach with filling.

Our idea is to fit an adaptive ellipse which can match the detected white area as good as possible. The ellipse starts moving and growing up from the center of the image so that it can overcome the influence of the above mentioned noise. In this way, the above mentioned noise can be efficiently

eliminated. The adaptive ellipse model contains the following steps.

#### 3.3.1.  Initial ellipse center and radius determination

Firstly, a small rectangle (30*10) is established using the center of lips ROI to delimit the noise. This size is then adjusted until the white area in this rectangle takes a large weight. More precisely, the small rectangle grows towards left and right by 3 pixels, and then towards up and down by 3 pixels. If the increased white pixels are more than 20% of the increased rectangle area, it continues expanding without exceeding the boundary of lips ROI. Otherwise, it stops. This rectangle thus contains the main part of inner lips region.

We then determine the center of inner lips region contained in the rectangle, this center noted as $(x_0, y_0)$, which will be the initial center of the adaptive ellipse, is determined using (1), where $P_x(i)$ is the sum of all the horizontal luminance and $P_y(j)$ sum of all the vertical luminance.

$$x_0 = \frac{\sum_i i P_x(i)}{\sum_i P_x(i)}, \quad y_0 = \frac{\sum_j j P_y(j)}{\sum_j P_y(j)} \tag{1}$$

To determine the parameters of initial ellipse (the semi-major axis $a$ and the semi-minor axis $b$), we first calculate the inertial moments about the two axes of the inner lips region in the rectangle by (2),

$$\sigma_x^2 = \frac{\sum_i (i - x_0)^2 P_x(i)}{\sum_i P_x(i)}, \quad \sigma_y^2 = \frac{\sum_j (j - y_0)^2 P_y(j)}{\sum_j P_y(j)} \tag{2}$$

$a$ and $b$ are then determined by:

$$a = \sigma_x / 2, \quad b = \sigma_y / 2 \tag{3}$$

#### 3.3.2.  Optimal ellipse searching

By using the center determined by (1) and the initial parameters $(a, b)$ given by (3), the initial ellipse is proportional to the inner lips region but still sufficiently small. The small ellipse starts to move and grow up successively in the four directions (up, right, down and left). With the ellipse growing up in each direction, the update of its radius and the center position (as illustrated in Figure 5), it will convergences to its optimal position and size which match the inner lips best. This adaptation can be summarized as follows:

$$a_{n+1} = a_n + \Delta a, \quad b_{n+1} = b_n + \Delta b \tag{4}$$

$$x_0^{n+1} = x_0^n + \Delta x_0, \quad y_0^{n+1} = y_0^n + \Delta y_0 \tag{5}$$

where $\Delta a, \Delta b, \Delta x_0$ and $\Delta y_0$ are strides of $a, b, x_0$ and $y_0$. Note that in each iteration, only one direction, i.e. $a$ and $x_0$ or $b$ and $y_0$ are updated at the same time. The sign of $\Delta x_0$ and $\Delta y_0$ are changed according to the moving direction (for example, $\Delta x_0$ is positive if it moves to right, negative to left, etc.). In our experiment, $\Delta x_0$ are fixed to 0.5 and $\Delta y_0$ to 0.2.

$\Delta a$ and $\Delta b$ describe the update rate of the ellipse size in each iteration. These parameters play an important role in the optimal ellipse searched process. In our experiment, the value of $\Delta a$ is fixed to 0.5 and let $b$ equals to $k \times \Delta a$, where $k$ equals

to the ratio $b_n/a_n$ in each iteration.

We denote $S_w$ the area of the white region in the current ellipse and $S_e$ the current ellipse area (Figure 6). For each direction, the ellipse expansion will be stopped moving and expanding if $S_e$ and $S_w$ satisfy the following condition (6), or the searching region exceeds the lips ROI. When the expansion of one direction is stopped, the expansion of other directions can still continue. The optimal ellipse is obtained until it convergences in all four directions (Figure 7).

$$S_w < S_e \times 0.7 \tag{6}$$

The size of the final ellipse is used to estimate the expected inner lips parameters using the following formulas:

$$A = \gamma \times 2a, \quad B = \gamma \times 2b \tag{7}$$

Given the stop criteria (6), $\gamma$ is logically equal to $\sqrt{0.7}$. However, $\gamma=0.87$ is chosen experimentally making the estimation results more accurate with respect to ground truth values.
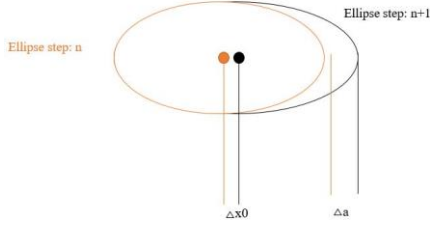


Figure 5. *Right expansion with the movement of the ellipse. Note that the major axis of the ellipse and also its center position are updated at the same time.*
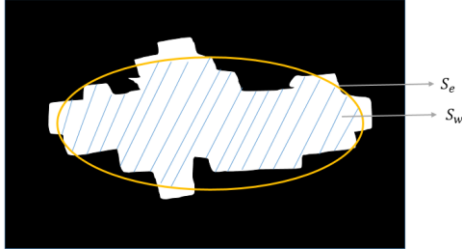


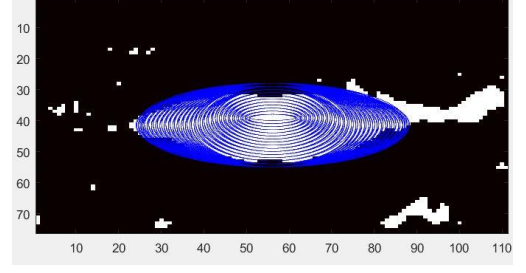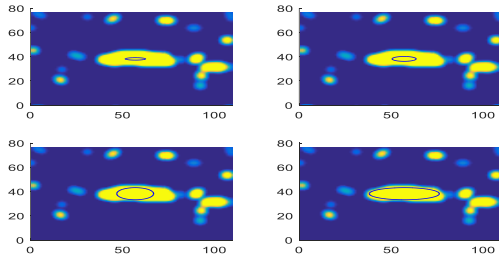Figure 6. *Ellipse expansion stop criterion*





Figure 7. *The iteration process of ellipse finding the optimal inner lips parameters (top), and a complete expansion procedure of ellipse model (bottom) in the direction of right, bottom, left and up.*

## 4.  Evaluation and Results

The adaptive ellipse model efficiently estimates the *A* and *B* parameters of inner lips for all kind of lips shape. The performance of three subjects can be visually shown in Figure 8.
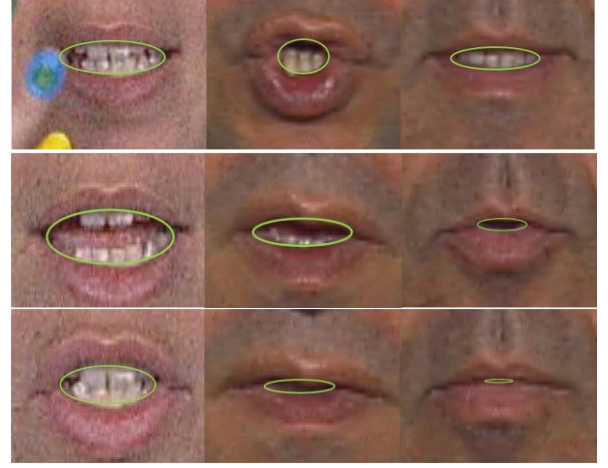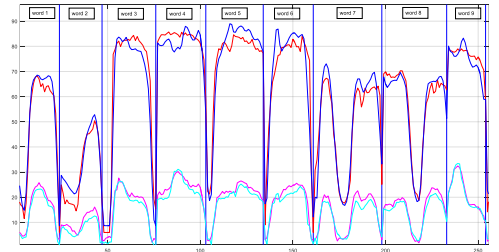


Figure 8. *Results of the inner lips ellipse in different cases (three speakers), the green ellipse is the optimal ellipse which can determine the A and B parameters. Left is CS subject, and the middle and right show the other two male subjects. In CS case (the first image), marked hand appears.*
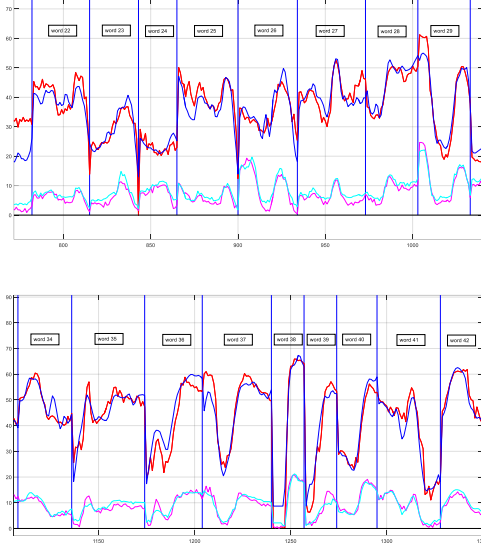
Figure 9. *The x-axis is the number of images, and the y-axis is the value in pixels unit. Red curve: the ground truth A parameter determined by expert. Blue curve: A values estimated by the proposed method. Cyan line: B values by proposed method. Magenta curve: Expert determined B parameter. The female speaker (top), other two male speakers (middle and bottom). Blue vertical lines show the boundary of each word. To observe clearly, we randomly choose several words interval of three speakers.*
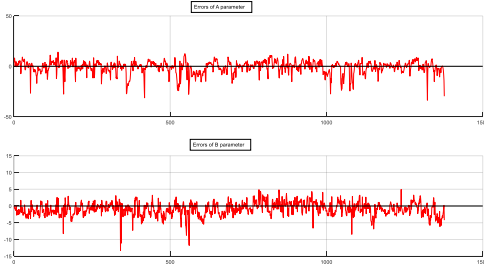


Figure 10. *Estimation errors of A (top) and B (bottom) parameters for CS speaker (total 1377 images in the first repetition).*

The proposed method is evaluated by comparing the *A* and *B* parameters with the ground truth value. In Figure 9, we can see that the estimated *A* and *B* parameters curves are quite close to the groud truth value curves for three speakers. In Figure 10, estimation errors of *A* and *B* parameters for the CS speaker are shown. We can see that except few jumped singularities, a global uniform distribution is presented without any evident dependence on the diversification of lips shape.

Mean values and RMSE (std) of *A* and *B* parameters are calculated for these errors. Results are shown in Table 1 and Table 2. As a comparison, values obtained by CLNF estimation are also given. We can see that the total mean errors of the *A* parameter are just half a pixel. The total RMSE is about 6 pixels (3.37 mm) for *A* parameter. It shows a better performance than [5] with a mean value of 3.5 mm and a RMSE of 4.5 mm.

Moreover, the proposed method significantly outperform the estimation results of the CLNF.

For *B* parameter, the total mean value is close to zero. Even though each speaker is considered separately, a mean value of 1 pixel (0.5 mm) remains a satisfactory result. These results are better than that in [5], for a mean error of 1.0 mm and a standard deviation of 1.0 mm. Meanwhile, the proposed method gives much better accuracy compared with CLNF.

In summary, our results show a far superior performance of A and B parameter estimation to that in the state of the art, and also are comparable to our previous work [18] which only addressed *B* parameter estimation.

Table 1: *RMSE values of A parameter for adaptive ellipse model and CLNF, expressed in pixels and in mm.*

| | Accuracy | Speaker 1 | Speaker 2 | Speaker 3 | Total error |
|---|---|---|---|---|---|
| **Proposed method** | mean | -1.06 (-0.60mm) | -0.38 (-0.21mm) | -0.31 (-0.17mm) | -0.55 (-0.31mm) |
| | std | 6.26 (3.54mm) | 6.18 (3.50mm) | 5.41 (3.06mm) | 5.95 (3.37mm) |
| **CLNF** | mean | 14.70 (8.32mm) | 32.99 (18.67mm) | 19.15 (10.84mm) | 23.01 (13.02mm) |
| | std | 19.38 (10.97mm) | 12.15 (6.88mm) | 13.79 (7.81mm) | 15.11 (8.55mm) |

Table 2: *RMSE values of B parameter for adaptive ellipse model and CLNF, expressed in pixels and in mm.*

| | Accuracy | Speaker 1 | Speaker 2 | Speaker 3 | Total error |
|---|---|---|---|---|---|
| **Proposed method** | mean | -0.92 (-0.48mm) | 0.91 (0.47mm) | -0.08 (-0.04mm) | -0.04 (-0.02mm) |
| | std | 1.97 (1.03mm) | 1.49 (0.78mm) | 1.36 (0.71mm) | 1.61 (0.84mm) |
| **CLNF** | mean | -2.04 (1.07mm) | -5.32 (-2.78mm) | -5.50 (-2.87mm) | -4.42 (-2.31mm) |
| | std | 3.83 (2.00mm) | 4.59 (2.40mm) | 4.10 (2.14mm) | 3.80 (1.99mm) |

To further evaluate the performance of the proposed method, we applied the estimated *A* and *B* parameters of inner lips to CS recognition on 34 French phonemes. 500 words contains about 1500 phonemesre used. HMM-GMM decoder (HTK 3.4) is used for recognition with a PCA-based feature extraction. An ellipse of white color with size (*A* and *B*) is superimposed on lips ROI (see Figure 11). PCA was used to extract good features on raw lips ROI and on the white filled parametric ellipse lips. Results (Table 3) show that the recognition score using white filled parametric ellipse lips is slightly higher than that using row images. This confirms the high precision of the estimated A and B parameters.
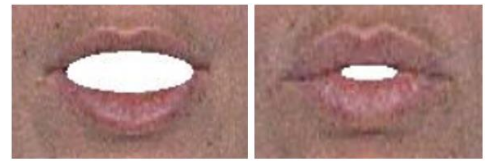


Figure 11. *The white filled parametric ellipse in lips ROI*

Table 3: *34 French phonemes recognition using PCA features on parametric ellipse lips ROI and raw lips ROI, 30 PCA components are used.*

| Features | Accuracy |
|---|---|
| PCA on parametric ellipse lips ROI | 62.0% |
| PCA on raw lips ROI | 59.8% |

# 5. Conclusion

In this paper, an efficient inner lips parameters estimation method based on the adaptive ellipse model is presented. We deal with the parametric extraction of inner lips from video without using any artifices. This method first extract the inner region of lips with an image processing combining single discontinuity elimination, and horizontal vertical filling. Then, an adaptive ellipse is used to match the inner lips region and give the best estimation of *A* and *B* parameters. Numerical precision is evaluated on 4693 images of three French speakers including CS speaker. The proposed method permits to obtain RMSE of 3.37 mm for *A* parameter and of 0.84 mm for *B* parameter, which outperform the baseline of inner lips parameters estimation in the state of the art. 34 French phonemes based CS recognition also confirms the superior performance. As a future work, CNN could be trained on our database to capture the non-linear relationships of inner lips to decrease the dependence of thresholds for different subjects.

# 6. Acknowledgements

# 7. References

[1] Panikos Heracleous, Denis Beautemps and Noureddine Aboutabit, ''Cued Speech Automatic Recognition in Normal Hearing and Deaf Subjects.'' *Speech Communication*, vol.52, Issue 6, pp. 504-512, 2010.

[2] R. O. Cornett, "Cued Speech", *American Annals of the Deaf*, vol. 112, pp. 3–13, 1967.

[3] Attina V., Beautemps D., Cathiard M.-A. & Odisio M. (2004). ''A pilot study of temporal organization in Cued Speech production of French syllables: Rules for a Cued Speech synthesizer". *Speech Communication*, vol. 44, pp. 197-214, 2004.

[4] Lallouache T., ''Un poste Visage-Parole. Acquisition et traitement des contours labiaux,'' *Actes des Journées d'Etudes de la Parole*, Montréal, 1990.

[5] Lionel Reveret, Christian Benoit, "A new 3D lip model for analysis and synthesis of lips motion in speech production", *ESCA workshop on Audio-visual speech processing, AVSP*, Australia, December, 1998.

[6] Noureddine Aboutabit, " Reconnaissance de la Langue Francaise Parlée Complétée (LPC) : Décodage phonétique des gestes main-lèvres," PhD dissertation, INPG, Gipsa-lab, Université Grenoble Alpes in Grenoble, France, 2007.

[7] Tim Cootes, ''An Introduction to Active Shape Models,'' *Model-Based Methods in Analysis of Biomedical Images in ''Image Processing and Analysis''*, Oxford University Press, pp. 223-248, 2000.

[8] T.F. Cootes, G.J. Edwards and G.J. Taylor, ''Active Appearance Model,'' *Proc. European on Computer Vision*, vol.2, pp. 484-498.1998.

[9] Samir K. Bandyopadhyay, ''Lip Contour Detection Techniques Based on Front View of Face,'' *Journal of Global Research in Computer Science*, vol. 2, No. 5, 2011.

[10] Pierre Gacon, Pierre-Yves Coulon and Gérard Bailly, "Non-linear active model for mouth inner and outer contours detection", *13th European Signal Processing Conference,* 2005.

[11] Stillitano S., Girondel V., Caplier C., ''Lip contour segmentation and tracking compliant with lip-reading application constraints,'' *Machine Vision and Applications*, vol. 24, Issue 1, pp. 1-18, 2013.

[12] M. Hennecke, V. Prasad, and D. Stork. "Using deformable templates to infer visual speech dynamics", *28th Annual Asimolar Conference on Signals, Systems, and computer, volume 2, IEEE Com*puter, Pacific Grove, pp. 576-582,1994.

[13] Jian-Ming Zhang, Liang-Min Wang, De-Jiao Niu, and Yong-Zhao Zhan, ''Research and implementation of a real time approach to lip detection in video sequence,'' *IEEE. Int. Conf. on Machine Learning and Cybernetics*, 2003.

[14] Evangelos Skodras and Nikolaos Fakotakis, ''An unconstrained method for lip detection in color images,'' *ICASSP*, 2011.

[15] Miroslav Hlavac, "Lips landmark detection using CNN", *Studentská vědecká conference*, 2016.

[16] Baltrusaitis T., Morency L.-P., and Robinson P. ''Constrained local neural fields for robust facial landmark detection in the wild,'' *IEEE, Computer Vision Workshops (ICCV-W)*, Sydney, Australia, 2013.

[17] Li Liu, Gang Feng, Beautemps D., ''Extraction automatique de contour de lèvre à partir du modèle CLNF'', *Actes des Journées d'Etudes de la Parole*, 2016.

[18] Li Liu, Gang Feng, Beautemps D., ''Automatic tracking of inner lips based on CLNF'', *ICASSP*, 2017.

[19] Cristinacce D. and Cootes T., ''Feature detection and tracking with Constrained Local Models,'' *Actes de British Machine Vision Conference*, vol. 3, pp. 929-938, 2006.

[20] Mehrdad Shemshaki and Roya Amjadifard, "Lip Segmentation Using Geometrical Model of Color Distribution", *Machine Vision and Image Processing (MVIP),* 2011.

[21] C. C. Chiang,W. K. Tai,M. T. Yang,Y. T. Huang,C. J. Huang, "A novel method for detecting lips, eyes and faces in real time" , *Real-Time Imaging*, vol. 9, pp. 277 – 287, 2003.