

# Effective Lip Localization and Tracking for Achieving Multimodal Speech Recognition

Wei Chuan Ooi, Changwon Jeon, Kihyeon Kim, David K. Han and Hanseok Ko

**Abstract**—Effective fusion of acoustic and visual modalities in speech recognition has been an important issue in Human Computer Interfaces, warranting further improvements in intelligibility and robustness. Speaker lip motion stands out as the most linguistically relevant visual feature for speech recognition. In this paper, we present a new hybrid approach to improve lip localization and tracking, aimed at improving speech recognition in noisy environments. This hybrid approach begins with a new color space transformation for enhancing lip segmentation. In the color space transformation, a PCA method is employed to derive a new one dimensional color space which maximizes discrimination between lip and non-lip colors. Intensity information is also incorporated in the process to improve contrast of upper and corner lip segments. In the subsequent step, a constrained deformable lip model with high flexibility is constructed to accurately capture and track lip shapes. The model requires only six degrees of freedom, yet provides a precise description of lip shapes using a simple least square fitting method. Experimental results indicate that the proposed hybrid approach delivers reliable and accurate localization and tracking of lip motions under various measurement conditions.

## I. INTRODUCTION

A MULTIMODAL speech recognition system is typically based on a fusion of acoustic and visual modalities to improve its reliability and accuracy in noisy environments. Previously, it has been shown that speaker lip movement is a significant visual component that yields linguistically relevant information of spoken utterances. However, there have been very few lip feature extraction methods that work robustly under various conditions. The difficulty is caused by variation of speakers, visual capture devices, lighting conditions, and low discriminability in lip and skin color.

Historically, there have been two main approaches [1] in extracting lip features from image sequences. The first method is called the Image-based approach. In this approach, image pixels (e.g. intensity values) around the lip region are used as features for recognition. For instance, these approaches are based on a DCT or a PCA method. The projected low dimensional features are used for speech recognition. The extracted features not only consist of lip features but also of other facial features such as tongue and jaw movement depending on ROI size. The drawback is that it is sensitive to rotation, translation scaling, and illumination

variation.

The second type is known as the model-based method. A lip model is described by a set of parameters (e.g. height and width of lips). These parameters are calculated from a cost function minimization process of fitting the model onto a captured image of the lip. The active contour model, the deformable geometry model, and the active shape model are examples of such methods widely used in lip tracking and feature extraction. The advantage of this approach is that lip shapes can be easily described by low order dimensions and it is invariant under rotation, translation, or scaling. However, this method requires an accurate model initialization to ensure that the model updating process converges.

In this paper, we propose a model-based method designed primarily to improve accuracy and to reduce the processing time. The model-based method requires good initialization to reduce the processing time. By integrating color and intensity information, our algorithm maximizes contrast between lip and non-lip regions, thus resulting accurate segmentation of lip. The segmented lip image provides initial position for our point based deformable lip model which has built in flexibility for precise description of symmetric and asymmetric lip shapes.

We describe in more detail of the PCA based color transformation method in section II. In section III, we present a new deformable model for lip contour tracking. We also describe cost function formulation and model parameter optimization in the same section. Experimental results and comparison with other color transformation are presented in Section IV. The conclusion is presented in Section V.

## II. COLOR TRANSFORMATION

### A. Lip Color and Intensity Mapping

Many methods have been proposed for segmenting the lip region that based on image intensity or color. We propose a new color mapping of the lips by integrating color and intensity information. Among color based lip segmentation methods are Red Exclusion [2], Mouth-Map [3], R-G ratio [4], and Pseudo hue [5]. Theoretically, pseudo hue method gives better color contrast, but we found that it is only useful in performing coarse segmentation which is not adequate for our purpose. Thus, we perform a linear transformation of RGB components in order to gain maximum discrimination of lip and non-lip colors. We employ a PCA to estimate the optimum coefficients of transformation. From a set of training images,  $N$  pixels of lip and non-lip are sampled and its distribution shown in Fig. 1(a). Each pixel is regarded as 3 dimensional vector  $x_i = (R_i, G_i, B_i)$ . The covariance matrix

Manuscript received May 1, 2008

The authors are with School of Electrical Engineering, Korea University, Seoul, Korea (e-mail: wcooi@ispl.korea.ac.kr; cwjeon@ispl.korea.ac.kr; khkim@ispl.korea.ac.kr), and the United States Naval Academy (e-mail: han@usna.edu). Hanseok Ko (corresponding author phone:82-2-3290-3239; Fax: 82-3290-2450; e-mail: hsko@korea.ac.kr)

is obtained from the three dimensional vector and the associated eigenvectors and eigenvalues are determined from the covariance matrix.

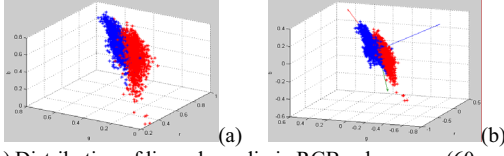


Fig. 1(a) Distribution of lip and non-lip in RGB color space (60 people under different lighting conditions), (b) Eigenvectors of distribution in Fig. 1(a)

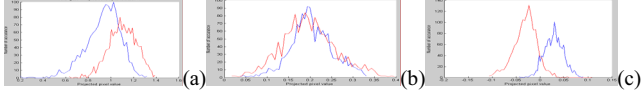


Fig. 2(a) Histogram of projected pixels onto first principle component, (b) Histogram of projected pixels onto second principle component, (c) Histogram of projected pixels onto third principle component

$V_3 = (v_1, v_2, v_3)$  is an eigenvector corresponding to the third smallest eigenvalue where lip and non-lip pixels are the least overlapping as shown in Fig. 2(c).

Experimentally,  $v_1 = 0.2$ ,  $v_2 = -0.6$ ,  $v_3 = 0.3$  are obtained. Thus a new color space,  $C$  is defined as

$$C = 0.2 \times R - 0.6 \times G + 0.3 \times B \quad (4)$$

The new color space  $C$  is normalized as

$$C_{norm} = \frac{C - C_{min}}{C_{max} - C_{min}} \quad (5)$$

Note that after normalization, the lip region shows higher value than the non-lip region. By squaring the  $C_{norm}$ , we can further increase the dissimilarity between these two clusters as shown in Fig. 3. A similar conversion of RGB values using the Linear Discriminant Analysis (LDA) was employed by Chan [6] to direct the evolution of snake. PCA based method is simpler compared to LDA especially in dealing with 3-D RGB components.

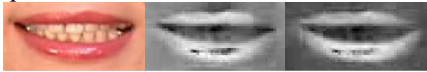


Fig. 3(a) Original image (b) Transformed image,  $C$  (c)  $C$  squared image After the color transformation, the  $C$  squared image may still show low contrast in the upper lip region. This problem can be resolved by using the intensity information  $I$ . The upper lip region typically consists of lower intensity values. So by combining the  $C$  squared image (which is well separable in the lower lip) and intensity image (which has a stronger boundary in the upper lip), we can obtain an enhanced version of the lip color map  $C_{map}$  as follows.

$$C_{map} = \alpha C_{sqr} + \gamma \frac{1}{I} \quad \text{where } \alpha + \gamma = 1 \quad (6)$$

Empirically,  $\alpha = 0.75$ ,  $\gamma = 0.25$  are derived. Higher weight is given to the  $C$  squared image since it captures most of the lip shape except the upper part and corners of lips.



Fig. 4 (a)  $C_{map}$  image (b)  $C_{map}$  negative image (c) Gray-scale image

### B. Threshold Selection

In this paper, the global threshold is selected based on Otsu [7] method. The optimal threshold  $T_{opt}$  is chosen so that

between classes variance  $\sigma_B^2$  is maximized.

$$T_{opt} = \underset{0 < T < 1}{\text{Arg max}} \sigma_B^2(T) \quad (7)$$



Fig. 5 Segmented images

## III. LIP CONTOUR EXTRACTION

### A. Lip Model

Most of the deformable geometric models are established using quadratic fittings (e.g. parabolic) with a prior assumption of lip shape being always symmetric about the center axis. Our lip model is an enhanced version of the proposed method in [8]. In [8] the writers integrated flexibility and constrained deformable template with point distribution model in order to reduce computations. However the geometric model in the paper is described by 15 parameters resulting in significant computation for the parameter updating process. Our proposed lip model is established by 6 parameters and is composed of 3 curves defined as follows:

i) Lower lip,  $\{0 < x < 1\}$

$$y_{low} = \alpha_{low} \cdot x \cdot (\log_2 x) + \beta_{low} \cdot (1-x) \cdot (\log_2 (1-x)) + \gamma_{low} \left( \frac{(x-0.5)^4}{0.5^4} - \frac{(x-0.5)^2}{0.5^2} \right) \quad (8)$$

ii) Upper right lip,  $\{0.5 \leq x < 1\}$

$$y_{up\_r} = -3.148 \cdot \alpha_{up\_r} \cdot (x-0.4)^2 \cdot (\log_2 x)^{\frac{1}{2}} + \gamma_{up\_r} \left( \frac{(x-0.5)^4}{0.5^4} - \frac{(x-0.5)^2}{0.5^2} \right) \quad (9)$$

iii) Upper left lip,  $\{0 < x \leq 0.5\}$

$$y_{up\_l} = -3.148 \cdot \alpha_{up\_l} \cdot (0.6-x)^2 \cdot (\log_2 (1-x))^{\frac{1}{2}} + \gamma_{up\_l} \left( \frac{(x-0.5)^4}{0.5^4} - \frac{(x-0.5)^2}{0.5^2} \right) \quad (10)$$

### B. Parameters Description

There are 6 parameters to fully model the lower and upper lips. Parameters  $\alpha_{low}$  and  $\beta_{low}$  control vertical height and skewness of lower lip as shown in Fig. 6.

- 1) if  $\alpha_{low} = \beta_{low}$ , lip (lower) is symmetric with respect to center point, then center height  $= \alpha_{low} = \beta_{low}$
- 2) if  $\alpha_{low} > \beta_{low}$ , lip shape slides to the left
- 3) if  $\alpha_{low} < \beta_{low}$ , lip shape slides to the right

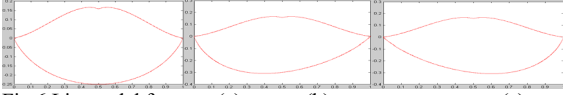


Fig.6 Lip model for case (a) (b) (c)

The parameter  $\gamma_{low}$  controls curvature of lower lip shape with values between -0.15 to 0.15 as shown in Fig. 7.

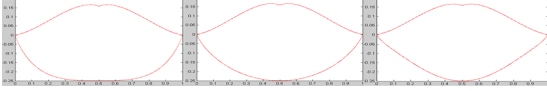


Fig. 7 (a)  $\gamma_{low} = -1.5$  (b)  $\gamma_{low} = 0$  (c)  $\gamma_{low} = 1.5$

Compared to the lower lip, upper lip shape remains relatively symmetric. This is due to the fact that the lower lip motion is a result of the mandible movement. The articulation available of the jaw joint allows the mandible movement of left or right of the centerline of the face resulting in asymmetric movement of the lower lip. Thus, we assume the upper lip always remains symmetric.

- 1)  $\alpha_{up\_l} = \alpha_{up\_r}$  = center height of upper lip
- 2)  $\gamma_{up\_r}$  and  $\gamma_{up\_l}$  control curvature of upper lip with value between 0.2 to -0.3

Lip shapes according to the variations of the upper lip parameters are shown in Fig. 8.

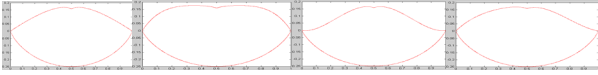


Fig. 8 (a)  $\gamma_{up\_r} = \gamma_{up\_l} = 0$  (b) -0.3, -0.3 (c) 0.2, 0.2 (d) -0.2, -0.1

### C. Model Initialization and Normalization

Our model initialization is based on a segmented lip shape image. In order to reduce the computation time, we limited our model with just 16 points of  $p = \{p_1, p_2, \dots, p_{16}\}$  where  $p_i = (x_i, y_i)$ . The lip points are labeled in anti-clockwise direction starting from the left corner. These points are divided into three groups where  $p_1, \dots, p_9$  describe lower lip,  $p_9, \dots, p_{13}$  describe upper right points, and  $p_{13}, \dots, p_{16}$  describe upper left. The contour points normalization process is applied to reduce processing time and to simplify the curve fitting process. The left corner point is fixed as the origin. Rotation and scaling transformations are employed to normalized all points so that  $p_1$  is at (0, 0) and  $p_9$  is at (1, 0). Reverse normalization is applied after curve fitting for obtaining the original coordinates.

### D. Model Optimization

The optimization procedure is an iterative process and the lip points are adjusted in order to reduce the cost function in each iteration process. Our cost function  $F$  is defined in (11)

$$F = \arg \min_p \sum_{i=1}^{16} \alpha E_{int}(p_i) + \beta E_{ext}(p_i) + \gamma E_{bal}(p_i) \quad (11)$$

$E_{int}(p_i)$  is an energy function dependent on the shape of the contour points and it is the continuity energy that enforces the

shape of the contour.  $E_{ext}(p_i)$  is an energy function based on image properties (we use gradient in this paper).  $E_{bal}(p_i)$  is a balloon force that causes the contour to expand (or shrink). In most cases, our binary image provides a good model initialization. Hence, the model usually takes only 5 to 8 iterations to converge to lip shape in a given image.

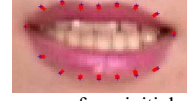


Fig. 9 Lip contour points converge from initial position to optimum position

### E. Model Fitting to Contour Points

We used least square approach to fit the model onto optimum contour points. Least square fitting is constructed from the given optimum lip contour points. For example, in lower lip model that employs 3 parameters,  $\theta = \{\alpha_{low}, \beta_{low}, \gamma_{low}\}$ .

$$H = \begin{bmatrix} x_2 \cdot (\log_2 x_2) & (1-x_2) \cdot (\log_2 (1-x_2)) & \left( \frac{(x_2-0.5)^4}{0.5^4} - \frac{(x_2-0.5)^2}{0.5^2} \right) \\ \vdots & \vdots & \vdots \\ x_8 \cdot (\log_2 x_8) & (1-x_8) \cdot (\log_2 (1-x_8)) & \left( \frac{(x_8-0.5)^4}{0.5^4} - \frac{(x_8-0.5)^2}{0.5^2} \right) \end{bmatrix}$$

$$\theta = \begin{bmatrix} \alpha_{low} \\ \beta_{low} \\ \gamma_{low} \end{bmatrix}, Y = \begin{bmatrix} y_{low2} \\ \vdots \\ y_{low8} \end{bmatrix}$$

$$\therefore H\theta = Y$$

$$\Rightarrow \theta = (H^T H)^{-1} H^T Y \quad (12)$$

Note that  $p_1$  and  $p_9$  are not included since these two points are fixed in the normalization process. The same curve fitting process is applied to the remainder of the lip model.

## IV. EXPERIMENTAL RESULTS

### A. Lip Contour Extraction Results

In order to test the performance of our proposed hybrid procedure, we use 2000 lip images with different sizes over 50 people (not including images that were used in color space training).

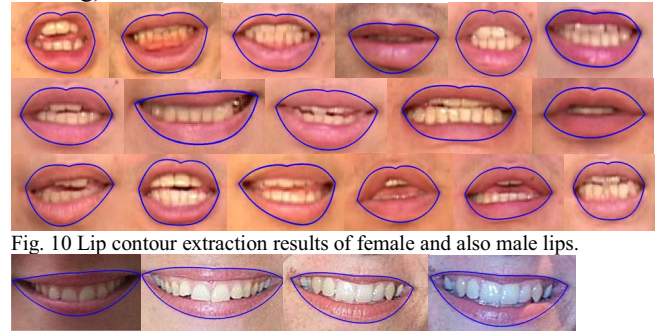


Fig. 10 Lip contour extraction results of female and also male lips.



Fig. 11 Lip contour extraction results under different lightning conditions Overall, 97% of the lip contours are accurately extracted. Fig. 10 and 11 show examples of such images. We also show that our proposed lip color and intensity mapping have successfully improved the lip contour extraction performance under different lightning conditions. With our algorithm implemented in Matlab, the average computation time for 85 x 100 size images was approximately 0.9 sec.

### B. Comparative Studies of Color Transformation

We apply a quantitative technique to evaluate the performance of our color space transformation algorithm. Since no ground truth is available, we manually draw the boundaries of 25 lip images. The first measurement method is the degree of overlap (DOL) between the lip and the non-lip histograms. DOL [9] is used to measure discriminability of the transformed color spaces for differentiating the lip and the non-lip colors. A lower percentage of DOL means a higher contrast between the lip and the non-lip regions.

$$DOL = \sum_{i=0}^l \min(P_{lip}(i), P_{nonlip}(i)) \quad (13)$$

where  $P_{lip}(i) = Num_{lip}(i) / Total_{lipPixel}$

$$P_{nonlip}(i) = Num_{nonlip}(i) / Total_{nonlipPixel}, \{0 \leq i \leq l\}$$

The second method is Classification Error (CE) which is the average of the False Positive (FP) rate and the False Negative (FN) rate. FP is error rate of classifying a non-lip as a lip pixel. FN is the error rate in classifying the lip as non-lip pixel.

$$CE = (FN + FP) / 2$$

$$\text{where } FN = False_{nonlip} / (False_{nonlip} + True_{lip}) \quad (14)$$

$$FP = False_{lip} / (False_{lip} + True_{nonlip})$$

TABLE I  
COMPARISON OF DOL AND CE BASED ON FIVE COLORS TRANSFORM METHODS FOR BELLOW 3 IMAGES

	Image 1 (Fig. 12)		Image2 (Fig. 13)		Image 3 (Fig. 14)	
	DOL	CE	DOL	CE	DOL	CE
MM	0.274	0.195	0.202	0.237	0.232	0.229
<b>Our</b>	0.182	0.150	0.201	0.143	0.173	0.144
RE	0.372	0.466	0.423	0.319	0.291	0.5
PH	0.223	0.243	0.294	0.188	0.210	0.162
RG	0.977	0.496	0.318	0.500	0.993	0.499

TABLE II  
COMPARISON FOR AVERAGE DOL AND CE FOR BELLOW 3 IMAGES AND ADDITIONAL 22 TESTING IMAGES

	MM	<b>Our</b>	RE	PH	RG
<b>Average DOL (%)</b>	23.8	16.4	36.7	22.9	55.7
<b>Average CE (%)</b>	21.4	12.5	35.8	17.2	36.4

From the results, we can see that our proposed lip color transformation method gives the lowest DOL and CE.

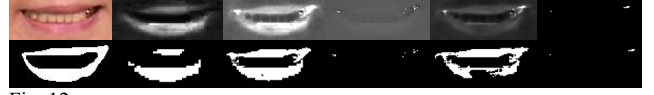


Fig. 12

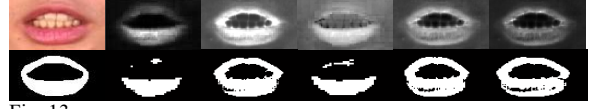


Fig. 13

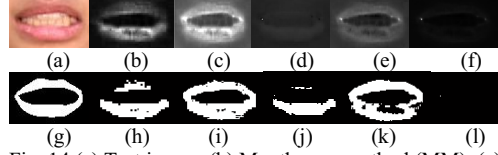


Fig. 14 (a) Test image (b) Mouth-map method (MM) (c) Our method (d) Red Exclusion method (RE) (e) Pseudo Hue method (PH) (f) R-G ratio method (RG) (g) Ground truth image (i) is segmented image of (c) based on Otsu thresholding (h),(j),(k),(l) are segmented images of (b),(d),(e),(f) with the threshold values proposed by corresponding previous methods.

### V. CONCLUSION

In this paper, we describe a new hybrid approach to improve lip localization and tracking. The first part of our proposed algorithm is lip mapping based on color and intensity information. From experimental results, our proposed mapping method successfully enhances the contrast between lip and non-lip regions. Results from the contrast enhancement process allowed more accurate lip region segmentation. In the second part, a new flexible while constrained deformable geometric model is established to accurately locate and track lip shape. Overall, our implemented hybrid approach has shown high reliability and is able to perform robustly under various conditions.

### ACKNOWLEDGMENT

This research was supported by MKE (Ministry of Knowledge Economy), Korea, under the ITFSIP (IT Foreign Specialist Inviting Program) supervised by the IITA.

### REFERENCES

- [1] G.Potamianos, C.Neti, G.Gravier, A.Garg, and A.W.Senior, "Recent advances in the automatic recognition of audio-visual speech," *Invited, IEEE Proc.*, vol. 91, pp. 1306-1326, 2003.
- [2] Lewis.T.W, Powers.D.M., "Lip Feature Extraction Using Red Exclusion," *Proc. Selected papers from Pan-Sydney Workshop on Visual Information Processing*, pp.61-67, 2000.
- [3] R.L.Hsu, M.Abdel, A.K.Jain, "Face Detection in Color Images," *IEEE Trans. on Pattern Analysis & Machine Intelligence*, 2002.
- [4] S.Igawa, A.Ogihara, A.Shintani, and S.Takamatsu, "Speech recognition based on fusion of visual and auditory information using full-frame color image," *ZEZCE Trans. Fundamentals*, 1996.
- [5] A.Hulbert and T.Poggio, "Synthesizing a Color Algorithm From Examples", *Science*, vol. 239, pp. 482-485, 1998.
- [6] Chan, M.T., "Automatic lip model extraction for constrained contour-based tracking," *ICIP*, pp. 848-851 1999.
- [7] N.Otsu, "A threshold selection method from gray-level histograms," *IEEE Trans. on Systems Man Cybernet*, pp. 62-66, 1979.
- [8] Wang, S.L, Lau,W.H, Leung,S.H, "A new real-time lip contour extraction algorithm," *ICASSP*, pp. 217-220, 2003.
- [9] Terrillon,T.C, Shirazi,M.N, Fukamachi,H, "Comparative performance of different chrominance skin chrominance models and chrominance spaces for the automatic detection of human faces in color images," *Proc. IEEE Int. Conf. Automatic Face and Gesture Recognition*, pp. 54-61, 2000.