THE USE OF LIP INFORMATION FOR ROBUST SPEAKER IDENTIFICATION AND SPEECH RECOGNITION

H. E. Çetingül, E. Erzin, Y. Yemez and A. M. Tekalp

Multimedia, Vision and Graphics Laboratory College of Engineering, Koç University, Sarıyer, Istanbul, 34450, Turkey ecetingul,eerzin,yyemez,mtekalp@ku.edu.tr

ABSTRACT

This study investigates the benefits of multimodal fusion of audio, lip motion and lip texture modalities for speaker and speech recognition. The audio modality is represented by the well-known mel-frequency cepstral coefficients (MFCC) along with the first and second derivatives, whereas lip texture modality is represented by the 2D-DCT coefficients of the luminance component within a bounding box about the lip region. A new lip motion modality representation based on *discriminative analysis* of the dense motion vectors within the same bounding box is employed for speaker/speech recognition. The fusion of audio, lip texture and lip motion modalities is performed by the so-called *Reliability Weighted Summation* (RWS) decision rule. Experimental results show that inclusion of lip motion and lip texture modalities provides further performance gains in both speaker identification and speech recognition scenarios.

1. INTRODUCTION

Audio is probably the most natural modality to recognize speech content and a valuable source to identify a speaker [1]. Video also contains important biometric information, which includes face/lip texture and lip motion information that is correlated with the audio. Audio-only speaker/speech recognition systems are far from being perfect especially under noisy conditions. Furthermore, it is a known fact that the content of speech can be revealed partially through lip-reading. Performance problems are also observed in video-only speaker/speech recognition systems, where poor picture quality, changes in pose and lighting conditions, and varying facial expressions may have detrimental effects [2, 3]. Hence, robust solutions for both speaker and speech recognition should employ multiple modalities, such as audio, lip texture, and lip motion in a unified scheme.

The design of a multimodal recognition system requires addressing three basic issues: i) Which modalities to fuse, ii) How to represent each modality with a discriminative and low-dimensional set of features, and iii) How to fuse existing modalities. Speech content and voice can be interpreted as two different, though correlated, information existing in audio signals. Likewise, video signal can be split into different modalities, such as face/lip texture and lip motion. The second issue, representative feature selection, also includes modeling of classifiers through which each class is represented with a statistical model or a representative feature set. For the final issue, fusion problem, different strategies are possible: In the so-called "early integration", modalities are fused at data or feature level, whereas in "late integration" decisions or scores resulting from each unimodal recognition are combined to give the final conclusion. A comprehensive survey and discussion on classifier combination techniques can be found in [4, 5].

State-of-art speech recognition systems have been jointly using lip information with audio [6, 7, 8]. For speech recognition, it is usually sufficient to extract the principal components of the lip information and to match the mouth openings-closings with the phonemes of speech. Speaker identification using audio and lip information, on the other hand, has been addressed in only few works such as [9, 4, 10]. The main challenge is that the principal components of the lip information are not usually sufficient to discriminate between speakers. Non-principal components are also valuable especially when the objective is to model the biometrics. In the speaker/speech recognition literature, audio is generally modeled by mel frequency cepstral coefficients (MFCC). However for lip information, there are several approaches reported in the literature such as texture-based, motion-based, geometrybased and model-based. In texture-based approaches, pure or DCTdomain lip image intensity are used as features [7, 4]. Motionbased approaches compute motion vectors to represent the lip movement during speaking [9, 11]. Geometry-based and model-based approaches, in fact, utilize similar processing methods such as active shape models [12], active contours [13] or parametric models [14] to segment the lip region. They differ in feature selection such that model-based approaches assign the fitted model parameters as features, while shape features such as lengths of horizontal and vertical lip openings, area, perimeter, pose angle, etc, are selected for lip representation in geometric-based approaches. The speaker recognition schemes proposed in [9, 10] are basically opinion fusion techniques that combine multiple expert decisions through adaptive or non-adaptive weighted summation of scores, whereas in [15], fusion is carried out at feature-level by concatenating individual feature vectors so as to exploit the temporal correlations that may exist between audio and video signals.

In this study, we use the lip motion features that are extracted by a novel discrimination analysis method [11]. Then we integrate lip texture, lip motion and audio features by a reliability-based decision fusion system reported in [4]. The main contribution of this paper is to investigate the fusion of audio modality with lip motion and lip texture representations for two distinct problems, speaker and speech recognition. The audio and lip features are presented in detail in Section 2. In Section 3, we describe the probabilistic framework that we use for the speaker/speech recognition problem, and present the reliability weighted summation rule for de-

This work has been supported by TUBITAK under the project EEEAG-101E026 and by the European FP6 Network of Excellence SIMI-LAR (http://www.similar.cc).

cision fusion of the multimodal system. Experimental results are presented and discussed in Section 4, and finally concluding remarks are given in Section 5.

2. MODALITIES AND FEATURES

In this paper, audio, lip texture and lip motion are considered as different modalities. The mel-frequency cepstral coefficients are used as features for the audio modality. The audio feature vector f_A is formed as a collection of MFCC vector along with the first and second derivatives. The features for the lip texture modality are 2-D DCT coefficients of the luminance component, and features for the lip motion modality are based on dense motion vectors within a rectangular box about the lip region.

A preprocessing step is required to locate the lip region and eliminate the global motion of the head between the frames so that the extracted motion features within the lip region provides us with the pure movement of the speaking act. To this effect, each face frame is aligned with the first frame of the sequence using a 2D parametric motion estimator. For every two consecutive face images, global head motion parameters are calculated using hierarchical Gaussian image pyramids and 12-parameter quadratic motion model [16]. The face images are successively warped according to these calculated parameters [11]. In the resulting aligned image sequence, the location of the lip region remains almost unchanged except for local movements. Thus, by only hand-labeling the mid-point of the lip region on the first frame, we automatically extract a region of interest around this point so as to obtain a sequence of lip frames of size 128×80 .

2.1. Features for Lip Texture Modality

It has been a common practice to use intensity-based features for the representation of lip texture [7, 4]. There are certain advantages and draw-backs of the intensity-based lip features, such as representing texture information as well as shape but being sensitive to illumination changes. The intensity-based lip features f_{L_t} are extracted by the Bayesian discrimination [11] from the zig-zag scan of 2D-DCT coefficients.

2.2. Features for Lip Motion Modality

Although lip movement is considered as the primary source for visual speech applications, it is rarely represented by its pure motion features. There are few studies incorporating the pure lip motion as the visual feature [9]. In [9], the lip motion is represented by the full set of 2D-DCT coefficients of the vectors. In this study the best lip motion representation that is found in [11, 17] is employed. A brief summary of this representation is presented in the following.

After performing global head motion compensation and lip region extraction, the use of a dense uniform grid of size 64×40 on the intensity lip image is considered. This grid definition allows us to analyze the whole motion information contained within the rectangular mouth region and it has proven its identification performance [17]. We use hierarchical block matching to estimate the lip motion in quarter-pixel accuracy by interpolating the original lip image with appropriate 6-tap Wiener and bilinear filters as used in H.264/MPEG-4 AVC [18]. The motion estimation procedure yields two 64×40 2D matrices V_x and V_y , each of which stores the motion vector components at grid points of the mouth region. The x and y components of the motion vector computed at the grid point (i, j) is given by the (i, j)-th entries of V_x and V_y , respectively. The motion matrices, V_x and V_y , are separately transformed via 2D-DCT. The first 50 DCT coefficients of the zig-zag scan both on x and y directions are combined to form a feature vector f_m of dimension 100.

In [11], we proposed a two-stage discriminative feature selection approach to determine best lip motion features. It takes into account the temporal discrimination information as well as the intra-class and inter-class distribution of individual single-frame lip feature vectors. At the first stage, we achieve discrimination in the Bayesian sense using a probabilistic measure that maximizes the ratio of intra-class and inter-class probabilities. The most discriminative features among the whole set of features f_m are selected. At the second stage, the successively concatenated lip feature vectors are created as a new sequence of higher dimensional feature vectors, each centered at the current frame instant. Then, they are projected to a lower dimensional feature space using linear discriminant analysis (LDA). The resulting lower dimensional feature vector representing the dense grid motion will be denoted by f_{L_m} .

3. MULTIMODAL FUSION

When more than one information source is available, the fusion of information from different sources can reduce overall uncertainty and increase the robustness of a classification system. Various alternative approaches have been proposed in the literature to the product rule such as max rule, min rule and reliability-based weighted summation. In fact, the most generic way of computing joint ratios (or scores) can be expressed as a weighted summation:

$$\rho(\lambda_r) = \sum_{n=1}^{N} \omega_n \rho_n(\lambda_r) \quad \text{for } r = 1, 2, ..., R,$$
(1)

where $\rho_n(\lambda_r)$ is the log-likelihood of the class-conditional probability, $\log P(f_n|\lambda_r)$, for the *n*-th modality f_n with class λ_r , and ω_n denotes the weighting coefficient for modality *n*, such that $\sum_n \omega_n = 1$. Then, the fusion problem becomes finding the optimal weight coefficients. Note that when $\omega_n = \frac{1}{N} \quad \forall n$, (1) is equivalent to the product rule. Since the ω_n values can be regarded as the reliability values of the classifiers, we referred to this combination method as RWS (Reliability Weighted Summation) rule in [4]. The statistics and the numerical range of these likelihood scores mostly vary from one classifier to another, and thus using sigmoid and variance normalization as described in [4], they can be normalized into (0, 1) interval before the fusion process.

The RWS rule is employed for the fusion of audio, lip texture and lip motion modalities for speaker and speech recognition problems, using the reliability value estimation, which is described in Section 3.3.

3.1. Speaker Recognition

Recognition task can be formulated as either verification or identification problem. The latter can further be classified as open-set or closed-set identification. In the closed-set identification problem, a reject scenario is not defined and an unknown observation is classified as belonging to one of the R registered pattern classes. In the open-set problem, the objective is, given the observation from an unknown pattern, to find whether it belongs to a pattern class registered in the database or not; the system identifies the pattern if there is a match and rejects otherwise. Hence, the problem can be thought of as an R + 1 class identification problem, including also a reject class. Open-set identification has a variety of applications such as the authorized access control for computer and communication systems, where a registered user can log onto the system with her/his personalized profile and access rights. In this paper, we formulate the speaker recognition problem in an openset identification framework, which is a more challenging and realistic way of addressing the problem as compared to closed-set speaker identification and verification. Note that verification is a special case of the general open-set identification problem.

In the open-set identification problem, an imposter class λ_{R+1} is introduced as the R+1'th class. Since it is difficult to accurately model the imposter class, λ_{R+1} , we employ the following solution which includes a reject strategy through the definition of the like-lihood ratio:

$$\bar{\rho}(\lambda_r) = \log \frac{P(\boldsymbol{f}|\lambda_r)}{P(\boldsymbol{f}|\lambda_{R+1})} = \log P(\boldsymbol{f}|\lambda_r) - \log P(\boldsymbol{f}|\lambda_{R+1}).$$
(2)

Then, the decision strategy of the open-set identification can be implemented in two steps. First, determine

$$\lambda_* = \underset{\lambda_1, \dots, \lambda_R}{\arg \max} \bar{\rho}(\lambda_r), \tag{3}$$

and then

if
$$\bar{\rho}(\lambda_*) \ge \tau$$
 accept
otherwise reject (4)

where τ is the optimal threshold which is usually determined experimentally to achieve the desired false accept or false reject rate.

Computation of class-conditional probabilities needs a prior modeling step, through which a probability density function of feature vectors is estimated for each class r = 1, 2, ..., R by using available training data. A common and effective approach to model the impostor class is to use a universal background model, which is estimated by using all available training data regardless of which class they belong to.

3.2. Speech Recognition

Speech recognition task can be formulated to identify a specific utterance, such as in the isolated word recognition task. Therefore the closed-set identification framework can be used to address the speech recognition problem with an isolated word dictionary.

The identification problem is formalized within the maximum likelihood framework. We can employ the maximum likelihood solution, which maximizes the class-conditional probability, $P(f|\lambda_r)$, for r = 1, ..., R. Hence a decision in the closed-set identification is taken as,

$$\lambda_* = \underset{\lambda_1, \dots, \lambda_R}{\arg \max} \log P(\boldsymbol{f} | \lambda_r) = \underset{\lambda_1, \dots, \lambda_R}{\arg \max} \rho(\lambda_r).$$
(5)

3.3. The Reliability Estimation for the RWS

Among various reliability estimation techniques existing in the literature, we favor the one proposed in [4], since it is better suited to the open-set speaker identification problem by assessing both accept and reject decisions of a classifier, and it can easily be defined for the closed-set identification problem. The RWS rule combines likelihood ratio values of the N modalities weighted by their reliability values ω_n as in (1). The reliability value ω_n is estimated based on the difference of likelihood ratios of the best two candidate classes λ_* and λ_{**} , that is, $\Delta_n = \rho_n(\lambda_*) - \rho_n(\lambda_{**})$, for modality n. In the absence of reject class, that is for closed-set identification, the likelihood difference of the best two candidates, Δ_n , can be used as the reliability value. However, in the presence of a reject class, one would expect a high likelihood ratio $\rho_n(\lambda_*)$ and a high Δ_n value for true accept decisions, and a low likelihood ratio $\rho_n(\lambda_*)$ and a low Δ_n value for true reject decisions. Hence, a normalized reliability measure ω_n can be estimated by,

 $\omega_n = \frac{1}{\sum_i \gamma_i} \gamma_n,$

where

$$\gamma_n = \begin{cases} \Delta_n & \text{closed} - \text{set} \\ (e^{\varrho_n} - 1) + (e^{\kappa - \varrho_n} - 1) & \text{open} - \text{set} \end{cases}$$
(7)

and

$$\varrho_n = \rho_n(\lambda_*) + \Delta_n. \tag{8}$$

(6)

The first and second terms for open-set identification in γ_n are associated with the true accept and true reject, respectively. The symbol κ stands for an experimentally determined factor to reach the best compromise between accept and reject scenarios. The κ value is set to 0.65 as it is found to be optimal for open-set speaker identification task in [4].

4. EXPERIMENTAL RESULTS

Hidden Markov Models (HMM) are known to be as effective structures to model the temporal behavior of the speech signal, and thus they are widely used both in audio-based speaker identification and speech recognition applications [1]. The speaker identification problem can further be classified as text-dependent and textindependent depending on the audio content. In the text-independent problem, identification is performed over a content free utterance of the speakers, whereas in the text-dependent case, each speaker is expected to utter a personalized secret phrase for the identification task. State-of-the-art systems use HMMs for textdependent and Gaussian Mixture Models (GMM) for text-independent speaker identification [19]. HMM-based techniques are preferred in text-dependent scenarios since HMM structures can successfully exploit the temporal correlations of a speech signal. Since lip motion is strongly coupled with audio utterance, HMMs can also be employed for temporal characterization of lip features. Hence, class-conditional probability modeling is performed using HMM architectures in our experiments.

In this work, we consider a text-dependent scenario for the speaker recognition problem and address it in the open-set identification framework, whereas for the speech recognition problem the closed-set identification framework is employed. The database consists of audio and video signals belonging to individuals of a certain population. Thus in our system the temporal characterization of the lip-motion modality is performed using HMMs. We use word-level continuous-density HMM structures for both the speaker identification and the speech recognition tasks.

The performance of the speaker verification systems are often measured using the equal error rate (EER) figure. The EER is calculated as the operating point where false accept rate (FAR) equals false reject rate (FRR). The performance of speech recognition system on the other hand is presented with the recognition rate, that is the ratio of the true matches to the total number of trials.

The speaker and speech recognition experiments have been conducted using the MVGL-AVD audio-visual database [20]. The database includes 50 subjects and considers two distinct text-dependent speaker identification scenarios, which are the name (\mathcal{D}_n) and the digit (\mathcal{D}_d) scenarios. In the name scenario, each subject utters ten repetitions of her/his name as the secret phrase. A set of impostor data is also collected with each subject in the population uttering five different names from the population. In the digit scenario, each subject utters ten repetitions of a fixed digit password 348 572. Although we have a limited variation in the name scenario, each name is considered as an isolated word, and a subset of the name scenario, $\mathcal{D}_s \subset \mathcal{D}_n$, which includes each name utterance with more than 12 repetitions, is considered as the testbed of the speech recognition experiments.

The audio recordings are perturbed with varying levels of additive noise during the testing sessions to simulate adverse environmental conditions. The additive acoustic noise is picked to be vehicle noise. Abbreviations and descriptions for the modalities and fusion techniques are given in Table 1.

 Table 1. Abbreviations and descriptions for modalities and fusion techniques

A	Audio modality
L_t	Lip texture modality
L_m	Lip motion modality
+	Product rule
\oplus	RWS rule

4.1. Speaker Recognition: Name Scenario

The \mathcal{D}_n database is partitioned into two sets namely $\{\mathcal{D}_{n_A} \text{ and } \mathcal{D}_{\bar{n}_A}\}$, where \mathcal{D}_{n_A} and $\mathcal{D}_{\bar{n}_A}$ are mutually exclusive sets each having five repetitions from each subject in the database. The subsets \mathcal{D}_{n_A} and $\mathcal{D}_{\bar{n}_A}$ are used for training and testing, respectively. Since there are 50 subjects and five repetitions for each true and imposter client tests, the resulting total number of trials for the true accepts and true rejects become respectively $N_a = 250$ and $N_r = 250$.

Table 2 presents the EER performance of the unimodal and multimodal open-set speaker identification systems for audio, lip texture and lip motion modalities. The EER performances of the lip texture and lip motion modalities are 5.6% and 5.2%, which are close to each other and better than the audio modality at 10 dB SNR and below. When the product rule and the RWS rule are applied to fuse pair of modalities or all the three modalities, the EER performance increases significantly. The RWS rule is observed to perform better than product rule, especially under noisy conditions. The best EER performance is achieved with the fusion of all three modalities at 10 dB SNR and below. Above 10 dB SNR, the best performance is achieved with the fusion of lip texture and audio modalities.

4.2. Speaker Recognition: Digit Scenario

The \mathcal{D}_d database is partitioned into two sets $\{\mathcal{D}_{d_A} \text{ and } \mathcal{D}_{\bar{d}_A}\}$, where \mathcal{D}_{d_A} and $\mathcal{D}_{\bar{d}_A}$ are mutually exclusive sets each having five repetitions of the same 6-digit number from each subject in the database. The subsets \mathcal{D}_{d_A} and $\mathcal{D}_{\bar{d}_A}$ are used for training and testing, respectively. Note that, in the digit scenario no imposter recordings are performed since every subject utters the same 6digit number. Hence, the imposter clients are generated by the *leave-one-out* scheme, where each subject, let us denote her/him by S, becomes the imposter of the remaining R-1 subjects in the population. Since, the class S is out of the population during the imposter tests, every test utterance that belongs to S becomes an imposter test. Having R = 50 subjects and five testing repetitions the resulting total number of trials for the true accepts and true rejects (imposters) become respectively $N_a = 250$ and $N_r = 250$.

Table 3 presents the EER performance of the unimodal and multimodal open-set speaker identification systems for audio, lip texture and lip motion modalities. The EER performances of the lip texture and lip motion modalities are 1.7% and 5.2%. Since, in the digit scenario every subject utters the same six digit password, the audio modality suffers and the lip texture modality benefits with respect to the name scenario. When the product rule and the RWS rule are applied to fuse pair of modalities or all the three modalities, the EER performance increases significantly. The RWS rule is observed to perform better than product rule. The best EER performance is achieved with the fusion of all three modalities at all SNR levels.

4.3. Speech Recognition

In this scenario, the database \mathcal{D}_s includes 35 different phrases (isolated words) where each phrase is actually names of the subjects in the database and repeated at least twelve times. The \mathcal{D}_s database is partitioned into two sets \mathcal{D}_{s_A} and $\mathcal{D}_{\bar{s}_A}$, where they are mutually exclusive sets each having equal number of utterance repetitions. The subsets \mathcal{D}_{s_A} and $\mathcal{D}_{\bar{s}_A}$ are used for training and testing, respectively.

Table 4 presents the recognition performance of the unimodal and multimodal speech recognition systems for audio, lip texture and lip motion modalities. The recognition performances of the lip texture and lip motion modalities are 62.86% and 72.86%. Since, the lip texture modality suffers to capture lip reading related information, the recognition rate of this modality is relatively poorer than the lip motion and audio modalities. When the product rule and the RWS rule are applied to fuse pair of modalities or all the three modalities, the recognition performance increases if the lip texture modality is not in the fusion. The best recognition performance is achieved with the RWS fusion of audio and lip motion modalities at all SNR levels.

5. CONCLUSIONS

A multimodal speaker/speech recognition system that integrates audio, lip texture and lip motion modalities is investigated. The lip motion modality is represented by dense-motion based features within a rectangular grid. We emphasize that lip motion modality carries additional useful information over that is present in the lip texture modality for both speaker and speech recognition applications. Hence, fusion of lip motion with audio and lip texture modalities is observed to provide additional performance gains.

Table 2. Speaker identification results for name scenario: Equal error rates at varying vehicle noise levels for different modalities.

EER (%)									
Source	Noise Level (dB SNR)								
Modality	clean	25	20	15	10	5	0		
A	1.0	1.2	1.6	4.8	13.2	22.4	30.8		
L_t	5.6								
L_m	5.2								
$L_m + A$	2.6	2.2	2.8	3.4	6.0	12.0	15.1		
$L_m \oplus A$	0.8	0.8	1.2	2.0	4.8	9.6	13.6		
$L_t + A$	0.4	0.8	0.8	1.2	3.6	7.0	12.0		
$L_t \oplus A$	1.0	0.8	1.2	2.0	2.8	5.0	6.8		
$L_m + L_t + A$	1.6	1.2	1.2	2.0	1.6	2.8	4.8		
$L_m \oplus L_t \oplus A$	1.2	0.8	1.0	1.6	1.2	2.4	4.0		

 Table 3. Speaker identification results for digit scenario: Equal error rates at varying vehicle noise levels for different modalities.

 FER (%)

$\operatorname{ELK}(\mathcal{N})$								
Source	Noise Level (dB SNR)							
Modality	clean	25	20	15	10	5	0	
A	2.4	2.6	2.8	5.6	11	24.2	37.2	
L_t	1.7							
L_m	5.2							
$L_m + A$	2.4	2.4	2.8	2.8	5.8	8.6	17.1	
$L_m \oplus A$	2.4	2.2	2.0	2.0	3.8	7.86	18.8	
$L_t + A$	0.4	0.4	0.4	0.4	1.4	5.2	13.8	
$L_t \oplus A$	0.4	0.4	0.4	0.4	0.8	2.8	10.8	
$L_m + L_t + A$	0.8	0.8	0.4	0.4	0.4	1.2	1.6	
$L_m \oplus L_t \oplus A$	0.4	0.4	0.0	0.0	0.0	0.6	1.6	

Furthermore, the lip motion is found to be more valuable than the lip texture modality for speech recognition. The fusion of audio, lip texture and lip motion modalities is performed by the so-called *Reliability Weighted Summation* (RWS) decision rule, which is observed to perform better than product rule.

6. REFERENCES

- J.P. Campbell, "Speaker recognition: A tutorial," *Proceedings of the IEEE*, vol. 85, no. 9, pp. 1437–1462, September 1997.
- [2] Y. Yan J. Zhang and M. Lades, "Face recognition: Eigenface, elastic matching, and neural nets," *Proceedings of the IEEE*, vol. 85, no. 9, pp. 1423–1435, September 1997.
- [3] M. Turk and A. Pentland, "Eigenfaces for recognition," *Journal of Cognitive Neuroscience*, vol. 3, no. 1, pp. 586–591, September 1991.
- [4] E. Erzin, Y. Yemez, and A.M. Tekalp, "Multimodal speaker identification using an adaptive classifier cascade based on modality reliability," *accepted for publication on IEEE Transactions on Multimedia*, 2004.
- [5] J. Kittler, M. Hatef, R. Duin, and J. Matas, "On Combining Classifiers," *IEEE Trans. on Pattern Analysis and Machine Intelligence*, vol. 20, no. 3, pp. 226–239, March 1998.

- [6] X. Zhang, C.C. Broun, R.M. Mersereau, and M.A. Clements, "Automatic Speechreading with Applications to Human-Computer Interfaces," *EURASIP Journal on Applied Signal Processing*, pp. 1228–1247, 2002.
- [7] G. Potamianos, C. Neti, G. Gravier, A. Garg, and A.W. Senior, "Recent Advances in the Automatic Recognition of Audio-Visual Speech," *Proc. of the IEEE*, vol. 91, no. 9, September 2003.
- [8] J.F.G. Perez, A.F. Frangi, E.L. Solano, and K. Lukas, "Lip reading for robust speech recognition on embedded devices," *Proc. of the Int. Conf. on Acoustics, Speech and Signal Processing 2005 (ICASSP '05)*, vol. I, pp. 473–476, 2005.
- [9] R.W. Frischholz and U. Dieckmann, "BioID: A multimodal biometric identification system," *Journal of IEEE Computer*, vol. 33, no. 2, pp. 64–68, February 2000.
- [10] T. Wark and S. Sridharan, "Adaptive fusion of speech and lip information for robust speaker identification," *Digital Signal Processing*, vol. 11, no. 3, pp. 169–186, July 2001.
- [11] H.E. Cetingul, Y. Yemez, E. Erzin, and A.M. Tekalp, "Discriminative Lip-Motion Features for Biometric Speaker Identification," *Proc. of the Int. Conf. on Image Processing* 2004 (ICIP 2004), pp. 2023–2026, October 2004.
- [12] S.L. Wang, W.H. Lau, S.H. Leung, and H. Yan, "A real-time

Recognition Rate (%)									
Source	Noise Level (dB SNR)								
Modality	clean	25	20	15	10	5	0		
A	88.10	87.14	85.71	85.24	78.57	68.57	51.43		
L_t	62.86								
L_m	72.86								
$L_m + A$	84.76	83.81	82.38	81.43	75.71	73.33	67.14		
$L_m \oplus A$	90.00	88.57	86.67	85.71	82.85	75.24	73.33		
$L_t + A$	77.62	78.10	77.14	76.66	74.76	68.57	61.90		
$L_t \oplus A$	77.14	76.66	76.66	76.19	75.23	74.28	69.05		
$L_m + L_t + A$	79.52	80.00	79.05	77.62	76.66	76.19	73.81		
$L_m \oplus L_t \oplus A$	79.52	77.14	77.14	75.71	76.66	74.76	74.28		

Table 4. Speech Recognition results: Recognition rates at varying vehicle noise levels for different modalities.

automatic lipreading system," Proc. of the 2004 Int. Symp. on Circuits and Systems (ISCAS 2004), vol. 2, pp. 101–104, 2004.

- [13] T. Wakasugi, M. Nishiura, and K. Fukui, "Robust lip contour extraction using separability of multi-dimensional distributions," *Proc. of 6th IEEE Int. Conf. on Automatic Face and Gesture Recognition (FGR'04)*, pp. 415–420, May 2004.
- [14] N. Eveno, A. Caplier, and P.-Y. Coulon, "Accurate and Quasi-Automatic Lip Tracking," *IEEE Trans. on Circuits* and Systems For Video Technology, vol. 14, no. 5, pp. 706– 715, May 2004.
- [15] U.V. Chaudhari, G.N. Ramaswamy, G. Potamianos, and C. Neti, "Information fusion and decision cascading for audio-visual speaker recognition based on time-varying stream reliability prediction," *Proc. of the Int. Conf. on Multimedia & Expo 2003 (ICME2003)*, vol. 3, pp. 9–12, July 2003.
- [16] J.-M. Odobez and P. Bouthemy, "Robust Multiresolution Estimation of Parametric Motion Models," *Journal of Visual Communication and Image Representation*, vol. 6, no. 4, pp. 348–365, December 1995.
- [17] H.E. Cetingul, Y. Yemez, E. Erzin, and A.M. Tekalp, "Robust lip-motion features for speaker identification," *Proc. of the Int. Conf. on Acoustics, Speech and Signal Processing* 2005 (ICASSP '05), vol. I, pp. 509–512, March 2005.
- [18] A. Puri, X. Chen, and A. Luthra, "Video coding using the H.264/MPEG-4 AVC compression standard," *Signal Processing: Image Communication*, vol. 19, pp. 793–849, 2004.
- [19] D.A. Reynolds, "Speaker identification and verification using gaussian mixture speaker models," *Speech Communications*, vol. 17, pp. 91–108, 1995.
- [20] E. Erzin, Y. Yemez, and A.M. Tekalp, DSP in Mobile and Vehicular Systems, chapter Joint Audio-Video Processing for Robust Biometric Speaker Identification in Car, Kluwer Academic Publishers, October 2004.