

On dynamic stream weighting for Audio-Visual Speech Recognition

Virginia Estellers*, *Student Member, IEEE*, Mihai Gurban and Jean-Philippe Thiran, *Senior Member, IEEE*,

Abstract—The integration of audio and visual information improves speech recognition performance, specially in the presence of noise. In those circumstances it is necessary to introduce audio and visual weights to control the contribution of each modality in the recognition task. We present a method to weight each modality according to its reliability for speech recognition, allowing those weights to change with time and adapt to different noise and working conditions. In our dynamic weighting scheme the weights are derived from several measures of the stream reliability, some specific to speech processing and others inherent to any classification task. In this paper we propose a new confidence measure, compare it to existing ones and point out the importance of the correct detection of silence utterances in Audio-Visual Speech Recognition. Experimental results are used to analyse the performance of several confidence measures and weighting techniques on different system’s architectures. We also prove the special role of the silence class for the definition of stream weights, showing that the inclusion of a voice activity detector in the weighting schema improves performance over different system architectures and confidence measures, leading to an increase of speech recognition more relevant than any difference between confidence measures.

Index Terms—Audio-Visual Speech Recognition, audio-visual fusion, weighted classifier, multi-stream HMM

I. INTRODUCTION

THE performance of Automatic Speech Recognition (ASR) systems degrades heavily in the presence of noise, compromising their use in real world scenarios. In those circumstances, ASR systems can benefit from the use of other sources of information complementary to the audio signal and yet related to speech. Visual speech constitutes such a source of information. Mimicking human lipreading, visual ASR systems are designed to recognize speech from images and videos of the speaker’s mouth. This fact gives rise to Audio-Visual Automatic Speech Recognition (AV-ASR), combining the audio and visual modalities of speech to improve the performance of audio-only ASR, specially in presence of noise [1], [2]. In those situations we can not trust the corrupted audio signal and must rely on the visual modality of speech to guide recognition, that is, give more importance to the visual than the audio cues when taking the decisions about the speech classes. Consequently, the problem of weighting each of the modalities for speech classification naturally arises.

The weight assigned to each modality should be related to its reliability to classify speech. In a quiet environment

with ideal audio and visual signals, higher weight should be given to the audio stream, reflecting the fact that the audio modality is more reliable than the video when it comes to recognize speech. When one of the modalities is degraded (due to background noise in the audio channel or an occlusion of the speaker’s mouth) the importance assigned to it should decrease and reflect the confidence we have on that modality under such circumstances.

In general terms the problem can be formulated as the combination of different streams of information in a classification task and is therefore not limited to AV-ASR. Indeed, it has been introduced in biometric person identification [3] to include feature streams from different modalities and in multi-band speech recognition [4] to consider different processing techniques applied to the same audio signal.

In our work we focus on the integration of the audio and visual information for the recognition of speech. We propose a dynamic scheme where weights are derived from instantaneous measures of the stream reliability, some specific to speech processing and others inherent to any classification task. The use of fixed weighting schemas has already been addressed in AV-ASR literature [2], [5]–[19], but only a few works [20]–[24] focus on dynamic weights adapting the system to changing environmental conditions. Moreover, some the results reported in literature for dynamic weights seem contradictory [20]–[22] and conclusions can not be derived because different confidence measures have been tested with different AV-ASR architectures, recognition criteria and databases. In that sense, the first contribution of the paper is a fair comparison of existing stream reliability measures in the estimation of the optimal weights. To that purpose we adopt the same form for the measure-to-weight mapping and optimization criteria and test the different confidence measures in both standard hidden Markov models and artificial neural network systems. The main contributions of the paper are, however, a new confidence measure inspired by the Viterbi algorithm and the introduction of a voice activity detector (VAD) in the weighting schema, taking into account the special role of the silence class in the definition of stream weights and AV-ASR systems. In fact, our experiments show that the improvement associated to the introduction of a VAD in the definition of stream weights is more relevant than any difference of performance between the proposed stream confidence measures.

The rest of the paper is organized as follows. In Section II we explain how the audio and visual integration takes place in ASR systems, review different stream weighting techniques proposed in literature and justify the necessity of dynamic weights adapting the system to changing envi-

The authors are with the Signal Processing Laboratory LTS5, Ecole Polytechnique Fédérale de Lausanne (EPFL), Ecublens 1015, Switzerland. (email: virginia.estellers@epfl,mihai.gurban@epfl.ch,jp.thiran@epfl.ch)

This work is supported by the Swiss SNF grant number 200021-130152
Manuscript received...

ronmental conditions. In the context of dynamic weights, in Section III we present existing techniques to estimate their correct value as a function of the stream reliability, propose a new confidence measure and the use of different weighting strategies for the speech and silence intervals. In Section V we report experiments with a reference database, comparing the performance and limitations of the different weighting techniques and in Section ?? we analyse the experimental results and discuss in more detail the role of the silence detection. Finally, conclusions are drawn in Section VI.

II. MULTI-MODAL FUSION FOR AV-ASR

In this Section we present the state-of-the-art for audio-visual fusion and stream weighting in ASR. We do not attempt to review the literature of ASR classification and refer the reader to [25]–[29] for a more complete overview on speech classification models. We simply justify the weighting model adopted in speech recognition and explain how it is included in the audio-visual models. The last part of the Section focuses on the weights associated to each stream as parameters of the model and reviews the existing techniques for their estimation.

A. Weighted multi-stream classifiers

In statistical classification it is common to assume that features of different streams are independent of each other. In that case, the statistical models factorize the joint probability distribution into single-stream distributions and reduce the complexity of the system. In the case of the audio-visual speech, perceptive studies showed that humans treat the streams as class conditionally independent [30], [31], that is, audio and visual features o_A , o_V are independent given that the speech class $q = q_i$. Under that hypothesis, speech recognition can be improved by introducing stream weights λ_A , λ_V and probability combination rules [32]. The model commonly used is a weighted geometrical combination of the audio and visual likelihoods

$$p(o_A o_V | q = q_i) = p(o_A | q = q_i)^{\lambda_A} p(o_V | q = q_i)^{\lambda_V} \quad (1)$$

controlling the importance of each modality in the classification task with its associated weight [33] and including the hypothesis that audio and visual modalities are class conditionally independent (equal unit weights).

Introducing the previous weighting schema into the statistical models used in ASR leads to the definition of multi-stream hidden Markov models (HMM) [25]. In single-stream HMMs, a discrete state variable $q(t)$ evolves through time as a first-order Markov process and controls the observed features $o(t)$ by defining a statistical model for the emission likelihoods $p(o(t) | q(t) = q_i)$ ¹. An HMM therefore factorizes the problem into the estimation of transition probabilities between states, which encode the temporal evolution of speech, and emission likelihoods associated to each state.

¹We use P to indicate the estimated probability value from a discrete distribution and p for the value taken by a probability density distribution of a continuous variable

In multi-stream HMMs, only the emission likelihoods are affected by the inclusion of different streams. The likelihoods are now computed independently for each stream and combined at certain level, which depends on the integration technique. In early integration the streams are assumed to be state synchronous and the likelihoods are combined at state level as indicated by (1). Late integration, on its turn, combines the likelihoods at utterance level, while in intermediate integration the combination takes place at intermediate points of the utterance. The weighting schema, nonetheless, remains the same and early or intermediate integration are generally adopted as leading to better results and finer control of the stream integration [21]. A common restriction is that the weights λ_A , λ_V sum up to one, so that the ratio between the lg observation and transition probabilities is kept the same as in single-stream HMMs.

In speech, the transition probabilities are chosen to force the HMM to evolve from left to right while either generative or discriminative strategies are used to estimate the probability distributions of the observed features. In generative systems a separate probabilistic model, usually a Gaussian Mixture Model (GMM) [27], is assumed for $p(o(t) | q(t) = q_i)$ and the corresponding parameters of the model are separately estimated for each class q_i . On the other hand, discriminative models use a single artificial neural network (ANN) or support vector machine (SVM) to assign a class probability distribution to the observed data $P(q(t) = q_i | o(t))$ and are thus designed to discriminate between classes, not generate class models. In that sense, training ANNs or SVMs to classify speech from different classes is more complex than estimating independently the GMMs for each class, but lead to models computationally simpler at testing stage than a large collection of GMMs. We will see that the use of GMM or ANN also affects the definition of reliability measures for the streams weights based on the performance of the classifier.

B. Weight estimation criteria

If we assume that the weights are fixed parameters of our models, we can estimate their optimal value with training or held-out data. In this case, the trained weights will only be relevant for the particular environmental conditions in which that data was acquired.

Ideally, we want to choose the weights that minimize the final Word Error Rate (WER) of our classifier, which is the natural measure of performance in ASR. However, the WER is not a smooth function of the training data, as its computation involves finding the most likely path between all possible state sequences and penalizing different types of errors (insertions, deletions and substitutions). Therefore, using the minimum WER as optimization criteria leads to simple grid-search methods choosing the weights with minimum WER in a training dataset, as reported in [7], [8].

The WER is a global measure of the performance of the system. It gives a score for each utterance, but it does not reflect the temporal evolution of the error within the speech sequence or how the weights affect the likelihood of the speech models used for classification. To overcome that

issue, some authors have proposed different smooth measures of the system's performance [5]–[8], [24], allowing the use of standard iterative techniques and optimization criteria on the training dataset. Those techniques usually minimize the frame error rate and maximize the discrimination between the different hypothesis of the classifiers. For instance, in [5], [7], HMM models are used to find the n most likely state alignments for the training data and their associated audio and visual likelihoods. The weights are then chosen to maximize the discrimination between those state alignments and the correct one in terms of their joint audio-visual likelihoods. It is not clear, however, how those measures of the system's performance relate to the final WER. Indeed, in [7], the authors point out that the minimum WER of their training dataset and the optimum of their proposed smooth function are not obtained for the same value of weights.

Other methods do not involve a training procedure. The weights are not chosen to optimize the WER or any function of the system's performance on some training data, but are set at testing to adapt the system to the working conditions based on the data itself. In [9], the authors use previous theoretical results [10] to estimate the optimal stream weights as inversely proportional to the single stream misclassification error. To that purpose, they build class specific models and anti-models and use them in a small amount of unlabelled data to compute inter and intra-class distances for each stream, from which they estimate their classification error and the corresponding optimal weights. Another criterion is proposed in [11], where the weights are chosen to maximize the dispersion of the test emission likelihoods and lead to a more discriminative classification, even though they might cause a wrong recognition. An extension of this algorithm is based on output likelihood normalization [12], where class-dependent weights are computed as the ratio between the average class-likelihoods over a time period. Note that here the weights become dynamic, as they are defined to normalize the class likelihoods at each time instant.

Dynamic stream weights, however, are usually introduced to adapt the system to changing environmental conditions due, for instance, to the temporary presence of noise in one stream. In this case we can not estimate the weights as fixed parameters of the system, but we have to make them evolve as a function of the estimated noise on each channel and the reliability of that stream for classification. For each noise level or estimated stream reliability, the weights can be chosen to optimize different measures of the performance of the system: recognition of isolated words [13]–[15], WER of continuous speech [2], [22], [23] or frame classification error [20], [21]. The weights can then be adjusted based on the estimated signal-to-noise ratio (SNR), as in [2], [13]–[17] or the voicing index in the case of speech recognition [18], [22], [34]. Other weighting methods applied to recognition tasks are modality-independent, as they are determined by the classifier's confidence [19]–[21] and can be used indifferently with the audio, video or any other stream. Only few of the previous works [20]–[24] applied dynamic weights to AV-ASR, that is, allowing the weights to change at frame level and adapt dynamically to the different noise conditions within

a sequence. It is also interesting to note that a few works [23], [24], [35], [36] combined audio and visual estimates of the stream reliability to define audio and visual weights.

In our work, we focus on the use of dynamic weights in different AV-ASR systems (HMM-ANN and HMM-GMM) when the audio stream is subject to noise. We present existing confidence measures, propose a new one computationally simpler and study how they map to the weights leading to minimum WER in a noisy training dataset. The minimum WER criterion is chosen as it is the final measure used to evaluate the system's performance and, as we have already said, it is not clear how other criteria relate to it. The questions that naturally arise are, for instance, deciding how to weight the audio and video streams during the silence periods inherent to speech, how quickly to adapt those weights in relation to the variations of the noise present on the stream and how the final WER is affected by the use of dynamic weights in a controlled noisy environment. The current paper shows the advantages, limitations and restrictions necessary to apply dynamic weights with changing levels and types of noise and points out the importance of silence recognition in the weighting scheme.

III. MEASURING STREAM CONFIDENCE

Two main strategies exist to estimate the reliability of the audio stream during fusion, either estimating a measure of the noise present on the channel directly from the audio signal or analysing the estimated posterior probability distributions of the classifiers. Figure ??? shows a block diagram of the proposed schemas. In this section, we present and compare some of the existing techniques, introduce a new one and propose a combination of them with a VAD to improve recognition results.

Based on the speech signal itself we estimate the SNR present on the audio channel by means of a VAD and simple power estimates. To measure the classifiers confidence we use the dispersion or the entropy of the class posterior probabilities and HMM emission likelihoods. We show how each measure is implemented and suits HMM-ANN or HMM-GMM systems and propose a new measure common and suitable for both architectures. Finally we study how the performance of the different measures can be improved by taking into account the classifiers behaviour for the different speech classes and the role of the silence detection in AV-ASR.

A. SNR of the audio signal

Measuring the SNR in ASR requires estimating the power of speech (signal of interest) and the power of the noise present on the audio signal. Due to the bursting nature of speech, we can obtain estimates of the noise power on the silence intervals inherent to any speaking utterance and derive from it an estimate of the power of the speech signal. The estimated SNR, denoted as \mathcal{S} , is then computed as the ratio of the speech and noise power estimates. To that purpose, we must first detect the silence and speech intervals with a VAD. At each time instant, we compute the power of the audio signal and assign it to speech or non-speech. If the sample is associated to

non-speech, we use it to update an estimate of the noise power. Otherwise the sample is detected as speech and, assuming noise and speech to be independent, the power of speech is estimated subtracting the previous estimated power of noise from the power of the audio signal. Note that estimating the power of noise during silence intervals defines an artificial low SNR for the silence intervals, when actually no speech is present and the SNR is ill-defined. This has a non-negligible effect on the weighting strategy, as the experiments will show.

As the aim of our work is not the study of VAD, we choose the technique best suited to our system and obtaining state-of-the-art results. We justify our choice as follows. VAD have a training stage where speech and non-speech models are built, usually assuming different Gaussian probability distributions of speech-related features for the speech and noise samples. On testing, a hypothesis test is used to estimate the likelihood of each sample belonging to speech or non-speech [37] and the results of the classification are afterwards smoothed in time. As the results of this instantaneous classification can be smoothed with an HMM, in fact we can directly use an audio-only HMM-GMM (which is part already of our AV-ASR system) to segment the signal into speech and non-speech intervals.

Note that using the estimated SNR as confidence measure is not particular to the audio channel, but its computation by means of a VAD is limited to speech signals. Moreover, it assumes a non-speech nature of noise and is therefore not designed to cope with babble noise. The use of audio-only HMMs, instead of other energy or GMM-based VAD techniques, enables the detection of speech/silence intervals subject to babble noise originated from a different vocabulary or grammar models than the trained audio HMM and provides a robust SNR estimate.

B. Confidence measures of the classifier

Generally it is advantageous to use stream confidence measures based on the classifier itself, as they convey information about audio and visual reliability for the classification task, can be applied to any data stream and are not specific to audio or speech signals. In ASR systems, the distribution of the posterior class probabilities or data likelihoods are the most common confidence measures based on classifiers. These measures assume that if the classifier assigns a very high probability to certain class while the rest present low probabilities, then the sample being tested fits correctly one of the trained models and the classification can be considered reliable. Conversely, when all classes have similar probabilities or emission likelihoods, the sample does not seem to distinctively fit any particular class and we assume it is corrupted by noise or due to an unreliable stream. Note that it is a reasonable assumption for speech classifiers trained with generative criteria (GMM usually), but its validity with ANN systems trained to discriminate between classes is less clear and, in fact, has proved false. Nevertheless, two measures have been proposed in ASR literature to capture this information: the entropy and dispersion of the posterior class probabilities or data likelihood of single-stream HMM classifiers.

In HMM systems, the dispersion of the emission log-likelihoods was first proposed in [13] to measure the difference on the probability scores of the N most likely states. Formally, if $\{q^1(t) \dots q^N(t)\}$ are the sorted N most likely states for the audio stream at time t , then the log-likelihood dispersion is:

$$\mathcal{D}(t) = \frac{2}{N(N-1)} \sum_{m=1}^N \sum_{n=m+1}^N \log \frac{p(o(t)|q(t)=q^m(t))}{p(o(t)|q(t)=q^n(t))}.$$

The dispersion measures of the confusability of the different classes in terms of emission likelihoods. High values of dispersion are associated to a small level of confusability and a reliable stream, whereas a low dispersion is encountered when all the likelihoods take similar values and the classes are highly confusable. An equivalent dispersion measure has been defined with class posterior distributions instead of emission likelihoods [22].

The entropy of the state posteriors has also been used both in HMM-GMM [20], [21] and in HMM-ANN systems for multi-band [38] and Audio-Visual ASR [22]. It is defined as

$$\mathcal{H}(t) = - \sum_{i=1}^C P(q(t)=q_i|o(t)) \log P(q(t)=q_i|o(t)),$$

where $\{q_1 \dots q_C\}$ are now all possible HMM states. It is important to note that in HMM-GMM systems the estimation of state posteriors requires the use of the Bayes rule and a priori class estimates from the training dataset. Contrarily to the dispersion, the entropy reaches its maximum value for equiprobable classes and has low values when the sample seems to specially fit one of the classes.

Even though both entropy and dispersion are based on measuring the peakiness of the probability distributions, it is not clear which one is better suited for the task of stream reliability estimation. In [20], [21] HMM-GMM systems obtained better performance for the dispersion than the entropy measures, while the contrary was observed with HMM-ANNs [22]. Those apparently contradictory results are due to the effect of estimating class prior probabilities when computing the entropy in HMM-GMM systems and to the different training strategies (generative training of GMMs compared to a discriminative training of ANNs). In terms of implementation, the entropy can be directly computed in the ANN models, which have class posterior probabilities as output. In the HMM-GMMs, however, the Bayes rule must be first applied for the computation of the entropies, while the dispersion can be directly computed from the emission log-likelihoods. In that sense, entropy seems more adequate for the ANN than the GMM architecture, while dispersion of posteriors or log-likelihoods suits both models equally. Nevertheless, computing the dispersion requires sorting the instantaneous likelihoods or probabilities and is computationally more expensive than the entropy. In our work we use both measures, dispersion and entropy, to measure the reliability of ANN and GMM-based HMMs and propose a new measure suited to both architectures and computationally simpler.

C. Proposed confidence measure of the HMM classifier

We observe that both GMM and ANN multi-stream systems share the same HMM structure to control the time evolution of the speech, but that only the GMM and ANN outputs were used in the definition of entropy and dispersion. We propose a new measure of the classifiers confidence not based on the values taken by the GMM or ANN's emission likelihoods, but on the transition probabilities of their common HMM structure. The proposed measure is inspired by the Viterbi decoder, where the transition probabilities between neighbouring sequence states are combined with their emission likelihoods to find the most likely sequence of states, naturally including the left-to-right property of speech HMMs and vocabulary restrictions. Our measure takes into account both the data likelihood and the time evolution constraints inherent to speech and exploits single-stream classifiers in terms of GMM/ANN models and HMM transition probabilities. These two terms can also be understood as a measure of data fidelity (emission likelihoods or class posterior probabilities associated to each sample) and a regularity constraint (transition probabilities associated to the most likely state of consecutive samples for each stream).

During recognition and for each stream we keep track of the most likely state in the single-stream ANN or GMM at each time instant $q^{ML}(t)$ (different from the most-likely state in the multi-stream HMM, whose computation requires the definition of weights) and update a counter for the stream $\mathcal{C}(t)$ with the transition probability between the previous and the instantaneous most likely state

$$\mathcal{C}(t) = \mathcal{C}(t-1) + p(q(t) = q^{ML}(t) | q(t-1) = q^{ML}(t-1)).$$

In practice, we do not keep track of the whole history of $q^{ML}(t)$, but define a limited memory to adapt the system to changing conditions. The counter is then implemented as a moving average of the transition probabilities between the instantaneous most-likely state for each stream. Note that a similar procedure of tracking the most-likely state and updating the log-likelihood of the path with the associated transition probability is done on the Viterbi decoder when recognition is performed with the single-stream HMM. In our counter, we do not keep track of the GMM/ANN emission likelihoods and simply use them to select $q^{ML}(t)$ and update the counter with the corresponding transition probability. The proposed measure is then easier to implement. Compared to entropy or dispersion, it does not require sorting or additional functions of the emission likelihoods, a max search and a single addition are enough to update of the transition counter.

The reliability associated to the stream increases with the value of the counter. If there is noise in the audio stream, its most likely state at each time instant will jump between states not matching the time evolution of the trained models and the associated transition probabilities will remain close to zero. We have experimentally observed that in presence of noise, $q^{ML}(t)$ mostly jumps between states corresponding to impossible transitions for the system (from the first state of phoneme A to the second state of phoneme B, transitions not allowed by left-to-right HMMs and vocabulary restrictions) and thus the

counter is updated with a transition probability equal to zero. This fact also justifies our choice of directly adding transition probabilities instead of their logarithms (compared to Viterbi), which does not provide a physical meaning to our measure, but results in a more stable confidence measure and avoids the instabilities and overflows of a counter regularly updated with the logarithm of transition probabilities close to zero. In terms of a regularity constraint, it corresponds to choosing a penalty function which allows punctual misfitting of the data to the models (transition probabilities 0) instead of introducing a large penalty for them.

Compared to the entropy or dispersion, the proposed counter takes also into account the temporal evolution of speech and vocabulary restrictions (sequences not allowed in the vocabulary have transition probabilities equal to zero). Entropy or dispersion only consider the emission likelihoods of the GMM or ANN, that is, how a sample instantaneously fit the observation models but not how the sequence of features fit the time evolution of speech. In our proposed method, the observation models are used to choose the most likely state at each time instant (data fidelity), while the time evolution of speech is taken into account by the transition probabilities (regularity term). Moreover, the transition counter suits HMM-GMM and HMM-ANN architectures and is computationally cheaper to compute than the entropy or dispersion.

IV. FROM STREAM RELIABILITIES TO WEIGHTS

Once the different stream reliability measures have been proposed, we have to map them to the stream weights leading to minimum WER of our system. To that purpose, we train the weighting system with an evaluation dataset subject to different known noise conditions. In the evaluation set we have artificially added different kinds and levels of noise to the audio stream. For each SNR level $n \in \mathcal{N}$ and type of noise, we do a grid search from the possible weights satisfying $\lambda_A + \lambda_V = 1$ and assign to that SNR level the weights λ_A^n, λ_V^n obtaining minimum WER on the evaluation data. For the same dataset, we compute the mean value of the different stream reliability measures at each SNR level $\bar{\mathcal{S}}^n, \bar{\mathcal{D}}^n, \bar{\mathcal{H}}^n, \bar{\mathcal{C}}^n$ and define the mappings f_S, f_D, f_H, f_C as continuous functions minimizing the mean square error (MSE) over all noise levels. For instance, for the transition counter we write

$$f_C = \arg \min_{f_C} \sum_{n \in \mathcal{N}} |\lambda_A^n - f_C(\bar{\mathcal{C}}^n)|^2. \quad (2)$$

In order to fairly compare all the stream reliability measures, we define mapping functions of the same complexity for the different reliability measures. From the values taken by those measures and the optimal weights in the evaluation data, we chose an exponential form for the mapping function $f(x) = A e^{Bx} + C e^{Dx}$, whose parameters A, B, C, D are estimated iteratively with a region-trust method [39]. We choose that form for the mapping based on some preliminary results with the evaluation data, as the performance was similar to the obtained with sigmoid [22] or piecewise linear functions [20] and the values of the resulting MSE (2) for the different measures where more similar between each other.

The average measures behaved as expected, with the dispersion and transition counter decreasing as the SNR level decreases and the entropy increasing for noisy data. It is to note, however, that the computation of the entropy in the GMM case is considerably sensitive to the estimation of the state prior probabilities and that the correct performance of that method requires a fine estimation of these probabilities. The best results are obtained using the time durations of phonemes in the training data to compute the state priors, while assuming equal class probabilities for all the states leads to a considerably poorer performance. Actually, as the distribution of phonemes in the training, evaluation and testing data is the same, the estimated priors match the testing ones. However it is not generally the case in real scenarios and it is more advisable to use the transition counter than the entropy in HMM-GMM systems, which performs similarly, is simpler to compute and does not require the estimation of prior probabilities.

It is important to note that we try to learn a mapping to be applied dynamically and yet we estimate it by experiments with fixed weights. For the evaluation dataset that is justified because the SNR is carefully kept fixed through all the sentences, which is achieved by artificially adding noise to the clean audio sequences. The value of the stream reliability indicator, however, varies within the sequence and we need to average it to define the mapping. In a real system, actually, the reliability measures and stream weights change instantaneously and the mapping learned from fixed weights might be incorrect. Smoothing the stream reliabilities through the testing sequences can give a similar behaviour as the one seen on training, but it does not ensure that the weights are instantaneously the best ones. In fact, it is necessary to study how each stream reliability measure evolves through an evaluation sequence with a fixed SNR level. If the confidence measure takes a relatively constant value throughout the sequence, then the mapping defined with the evaluation dataset between the mean value of this measure and the fixed optimal weights can be used. Otherwise, if the variations of the reliability measure on a fixed SNR sequence are comparable with the variations between different SNR levels, then the mapping can not be directly used without a large smoothing of the confidence measure on the testing sequences, which hinders a quick adaptation to changing noise conditions.

Analysing the evolution of the different reliability measures on the evaluation dataset, we state that the estimated SNR $\mathcal{S}(t)$ requires a considerable smoothing, mainly due to the estimation of SNR during the silence periods between words. In those periods, the SNR ratio is small as there is no speech signal present. For the same sequence, the variations of $\mathcal{S}(t)$ between the speech and non-speech intervals are higher than between the different SNR levels. Nevertheless, as $\overline{\mathcal{S}}^n$ evolves coherently with the SNR levels $n \in \mathcal{N}$, that issue can be solved assuring that the smoothing applied on testing always includes speech and silence intervals. The measures based on the classifiers confidence $\mathcal{D}, \mathcal{H}, \mathcal{C}$ show also different mean values for the silence and speech utterances, as shown in Fig.IV. The variations here are not caused by the SNR estimation procedure and indicate that the classifiers show

different behaviour for the speech and silence intervals.

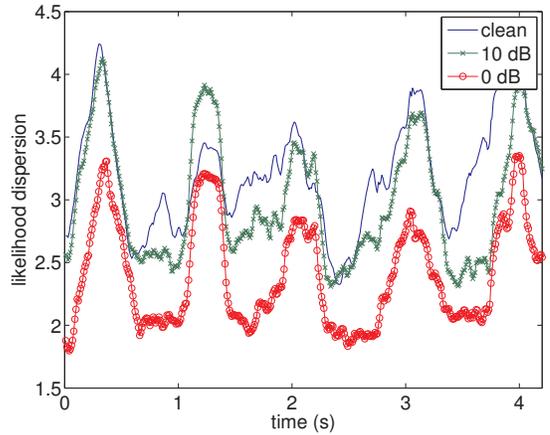


Fig. 1. Time evolution of the GMM likelihood dispersion for the same evaluation sequence and different levels of white audio noise artificially added. The likelihood increases at 1.2, 2, 2.9 and 3.8 seconds due to word utterances compared to its value during inter-word silences.

This analysis suggests that the silence intervals inherent to speech might play an important role in the definition of proper stream weights. Indeed, in [22] the authors point out that for very noisy environments, training the weights with the minimum WER criterion leads to choosing the modality better suited for the detection of the silences existent between the words in the utterances. To avoid that kind of behaviour, we developed a second strategy shown in figure ???. Using the same VAD used for the SNR estimation, we first assign each sample to speech or silence and then set the weights accordingly with different speech and silence mapping functions. To that purpose, we define two mappings from each reliability measure to the optimal stream weights, one for the recognition of silences and another for the recognition of speech. To train this combined weighting system, we split the evaluation dataset into speech and silence examples based on the available labels, we concatenate them into continuous speech and silence utterances and learn the corresponding mappings as previously explained. Note that we learn the mappings from continuous speech recognition experiments with sequences containing only continuous speech or silence examples. Learning those mappings from isolated word recognition tests would define a weight threshold leading to correct or wrong recognition of each word instead of maximizing the WER, which is the performance criterion used in continuous speech recognition. Now the stream reliability measures, specially the SNR estimator, are not influenced by the presence of silence intervals and are more stable within the evaluation sequences. Moreover, considering only speech or silence utterances to define the optimal stream weights, we obtain a mapping better suited to classify speech while the decision about silence or speech intervals is taken with a VAD designed to that purpose. We refer to those mapping strategies as f_S^{VAD} , f_D^{VAD} , f_H^{VAD} and f_C^{VAD} .

At testing stage, both single and double mapping weighting schemas first smooth the different reliability measures, then use the estimated mappings to compute the corresponding

weights and use them in the multi-stream HMM classifier to recognize speech. Compared to the training stage, the confidence measures are not averaged per sequence but only smoothed in time. A moving average is used for smoothing the obtained \mathcal{H} , \mathcal{D} and \mathcal{S} , while it is intrinsic to the definition of \mathcal{C} as an accumulated transition counter. The window size adopted for the smoothing has been chosen based on preliminary experiments with the evaluation dataset, where a 10 ms window obtained good results across the different reliability measures. The other main difference between testing and training affects only the double-mapping schemes, where a VAD is used to classify samples as speech or silence. During training this decision is taken based on the labels of the evaluation data, while at testing a VAD is used for that purpose, which introduces a possible source of error.

V. EXPERIMENTS

We perform continuous speech recognition experiments on the CUAVE database [40]. We use the static portion of the individuals Section of the database, consisting of 36 static subjects repeating the digits five times in front of a camera. We do speaker independent experiments with 6-fold cross validation, using 30 speakers for training, 3 for evaluation and 3 for testing. The results are given in terms of speaker independent WER and the statistical significance of the results is evaluated in a paired manner comparing the different confidence measures.

A. Feature streams

Normalized mel-frequency cepstral coefficients (MFCC) are extracted as audio features, with their first and second temporal derivatives. Thirteen MFCC features are computed with a 30 ms window, with 10 ms overlap, leading to an audio rate of 100 feature vectors per second. We train any HMM parameters on clean audio data and artificially add white and babble noise from the NOISEX database [41] on testing. Different levels of noise are added in order to show how our dynamic weighting algorithm performs across a large range of SNRs, from clean to -10 dB. Adding noise to the recorded signal instead of adding it during the recordings does not take into account the changes in articulation speakers produced when background noise is present [42] and therefore generates somehow non-realistic scenarios. On the other hand it enables to test exactly the same utterances in different noise conditions, facilitates the recordings of the data and a complete control of the noise conditions in the evaluation set.

The visual features are selected Discrete Cosine Transform (DCT) coefficients from a region of interest defined around the mouth, which consists of a 128×128 image of the speaker's mouth, normalized for size, centred and rotated. The DCT coefficients are the 13 most important ones taken in a zig-zag order, as in the MPEG/JPEG standard, together with their first and second temporal derivatives with their means removed. More details about that visual feature extraction system can be found in [43]. The temporal resolution of the visual features is then increased through interpolation to reach the audio rate, since synchronous audio and visual feature streams are

required by the classifiers. No noise is added to the visual features as AV-ASR with non-ideal visual conditions requires the development of new visual feature extraction methods before audio-visual integration can be successfully studied [44], [45].

B. Speech Classifiers and VAD

We use the HTK library [46] for the HMM-GMM implementation of 3 state left-to-right phoneme models. Each state has 3 Gaussians for the audio and one for the visual stream, all with diagonal covariance matrices. We start by training separately audio and visual HMM-GMM models, we then build the multi-stream models and jointly re-estimate their parameters setting the audio and visual weights to one during training².

In the ANN case, the emission likelihoods are replaced by posterior phoneme probabilities estimated with a Multi-Layer-Perceptron. The audio and visual neural networks are implemented as feed-forward ANN with two neural layers and 10000 neurons. One feature vector is feed to the ANN each time with sigmoid functions used in the input layer. The ANN has an output node for each class, with softmax functions used to provide an estimate of the class posterior probabilities associated to the input sample. The values of the transition probabilities from the HMM-GMM case are kept for the HMM-ANN system, as they correspond to a time model of the duration of phonemes learned from the same training data.

Recognition is based on phonemes, which are concatenated to form words and sentences by means of a dictionary and grammar. In our case, as the testing correspond to sentences containing sequences of numbers, no grammar is used and the dictionary includes only the phonetic transcription the English digits.

Designing the VAD we must compromise between having voice detected as noise or noise detected as voice (between false positive and false negative). In our case, the VAD must be able to detect speech under several types and levels of background noise and we design it to be fail-safe, that is, to detect speech when the decision is in doubt and lower the chance of losing speech segments. As already explained, we use the audio-only HMM systems to classify features as speech or silence. Single-stream HMM are also used to estimate the entropy, dispersion and transition counter reliability measures of the audio stream. The obtained confidence measures are used to compute the weights of the multi-stream HMM systems, taking also into account the decision of the VAD in the case of double-map weighting schemes.

C. Evaluation of the results

In our experiments we compare different weighting strategies learned and tested on the same data and the results, therefore, reflect differences between the weighting strategies rather than differences in the test datasets. In that case, the

²Experiments have shown that final performance of the system is dominated by the value of the weights at testing and not during training [47]

statistical significance of the results can not be evaluated by means of confidence intervals associated to the performance of each method independently, but requires the comparison of the different methods in a one-to-one basis for the same sentences, speakers and train/test datasets.

In speech recognition a small modification to a system will alter the recognition results in a few sentences or speakers only. Intuitively we would acknowledge a probability of reducing the errors of 10% if the number of errors drops on 10% of the sentences while the others remain unchanged. On the other hand, an overall improvement of the word error rate should be considered random if 50% of the sentences improved while 50% degraded. More formally, we estimate the “probability of error reduction” p_e between two systems A and B measuring the number of independent testing samples that favour system A over B while leaving the rest of the samples unchanged. The actual computation of p_e involves estimating a probability distribution associated to the paired comparison of systems. To that purpose, we bootstrap the WER obtained by the different weighting methods for independent samples (sentences) and perform a paired hypothesis test to obtain p_e . Bootstrapping allows us to estimate the unknown distributions associated to the WER (whose computation involves the ratio between several types of errors with unknown distributions) and also to obtain an estimate p_e , which does not depend on the number of sentences used in the comparison. We use p_e as an evaluation tool and refer the reader to [48] for more details. On the following only the values of p_e relevant to assess if one method significantly outperforms another are given.

D. Experimental Results

The aim of this work is not to compare HMM-ANN and HMM-GMM architectures, so the results of the stream reliability measures will be compared for each of the systems separately. We include results for 3 extra baseline systems that we use to analyse the improvement obtained with a weighting strategy.

As mentioned previously, the use of a VAD for the computation of some reliability measures and weights, also includes a font of error due to failing voice activity recognition. The performance of the VAD, shown in table I, should also be taken into account in the analysis of the results. We observe that the VAD works reasonably well for white noise and up to 5-0 dB for babble noise, when performance drops under 70%.

TABLE I

PERFORMANCE OF THE VAD FOR DIFFERENT TYPES AND LEVELS OF ACOUSTIC BACKGROUND NOISE. RESULTS ARE GIVEN IN TERMS OF TRUE POSITIVE RATE OR SENSITIVITY (TPR) AND FALSE POSITIVE RATE (FPR) AND ACCURACY (ACC).

white	clean	25dB	20 dB	15 dB	10 dB	5 dB	0 dB	-5 dB	-10 dB
TPR	97.3	94.4	91.7	88.1	82.9	76.7	70.2	58.0	36.8
FPR	4.6	2.9	2.4	1.7	1.3	0.9	0.6	0.6	1.1
acc	96.3	95.8	94.8	93.4	91.2	88.4	85.5	79.7	69.7
babble	clean	25dB	20 dB	15 dB	10 dB	5 dB	0 dB	-5 dB	-10 dB
TPR	97.1	94.8	92.7	89.5	82.5	73.9	61.2	45.9	36.1
FPR	4.6	3.6	3.5	5.0	8.4	15.3	23.5	28.8	29.9
acc	96.2	95.6	94.6	92.3	87.1	79.4	69.1	59.3	54.53

The first baseline system is an audio-only ASR system, showing the gain obtained by inclusion of the visual modality, and the other two are AV-ASR systems with fixed unit and “cheat” weights, that is weights assuming class conditional independence of audio and visual streams and the weights obtaining minimum WER for each SNR with the test dataset (not considering the difference between silence/speech in the weight definition). Fixed units weights corresponds to a system where no stream weighting is used and audio and visual features are just considered conditionally independent. Comparing to such a system shows the improvement obtained by a weighting strategy under different noise circumstances. In Fig.2a and 2b we see that a weighting strategy is significantly useful below 25 and 20 dB of audio SNR for the HMM-GMM and ANN case. The probability of error reduction p_e ranges from 0.7 to 1.0 for the different SNR levels and it defines the range of noise levels, where the comparison of the different weighting strategies is relevant. In that case, it is also important to note the performance of 60.4% obtained with a HMM-GMM visual-only system and 56.2% for the HMM-ANN, which specially justifies the inclusion of the visual modality under 5-0 dB of SNR for the different systems and the study of AV-ASR in these circumstances.

Comparing to the fixed “cheat” system shows how far we are from the best behaviour under the assumption that the weights only depend on the SNR of the stream. In that sense, when a dynamic system outperforms the fixed one, it is due to the fact that the silence/speech class should also be taken into account for the weight definition, which is not the case with the “cheat” weights. Such is the case for the HMM-ANN system under 10 dB of SNR, see Fig.4b and Fig.4d, with a probability of error reduction p_e over 0.9 for all the noise kinds and levels. The same behaviour is observed under -5 dB and 5 dB SNR for the HMM-GMM case subject to white and babble noise, see Fig.3b and Fig.3d. In those cases the optimal fixed weights choose the modality better suited for the silence detection, while the confidence measures including the VAD define different mappings for the silence and speech intervals. In that case, results show that the performance of the AV-ASR is limited by the performance of the VAD.

Comparing the two mapping strategies (f_S^{VAD} against f_S , $f_{\mathcal{H}}^{\text{VAD}}$ against $f_{\mathcal{H}}$, etc.), see Fig.3 and Fig.4, we observe that a considerable improvement is obtained when different mappings are used for the classification of speech and silences. The improvement is more remarkable in the SNR estimator (with p_e over 0.9 for the different SNR levels), whose estimation is based on the correct detection of speech and silence intervals, while the entropy, dispersion and transition counter benefit less from the inclusion of a VAD in the weighting system (p_e between 0.7 and 0.8) as they already convey information about the confidence that should be given to the classifier during silence intervals. The gain is also more clear for low SNR levels, when the use of only one mapping mainly shifts the weights to use the modality better suited to the silence detection. In this case, the use of two mappings relies on the VAD to detect silence and speech intervals and then uses the corresponding speech or silence mapping for each confidence measure. As a result, for low SNR the dynamic weights even

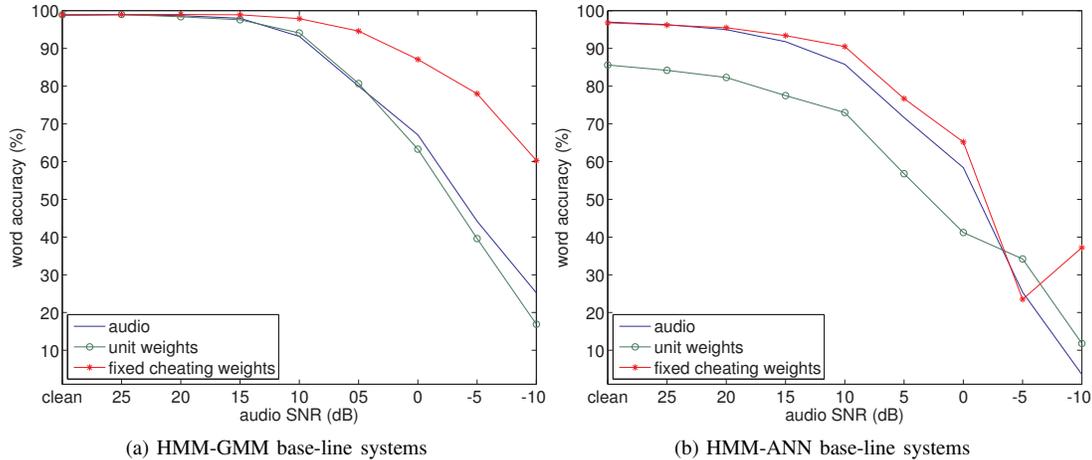


Fig. 2. Baseline 1: performance of single and multi-stream HMM systems for different SNR levels of audio white noise. The weighting strategy is significantly useful below 25 and 20 dB of audio SNR for the HMM-GMM and ANN case.

outperform the fixed cheating weights, which do not consider the fact that the proper detection of silences might require a different weight than the speech intervals.

Using only one mapping in the HMM-ANN system, different measures seem to obtain better results for different working conditions. The SNR estimator performs well for high and medium SNR values, while the measures based on the classifiers confidence gain in very noisy environments. In the GMM system, on its turn, the entropy and transition counter do slightly better than the other confidence measures. When the VAD is included in the weighting strategy, however, the different confidence measures perform similarly and the differences in performance are not statistically significant (in terms of p_e) for any of the measures.

To summarize, we see that the improvement obtained by the inclusion of the VAD in the weighting strategy is more relevant than any differences in performance between the confidence measures. Without the use of a VAD, different confidence measures obtain better performance for different systems and levels of noise, while the introduction of a VAD into the weighting system improves the performance of all the confidence measures and leads to statistically equivalent results for the different measures. In that case, the proposed transition counter f_C^{VAD} performs equivalently to other classifier's derived measures f_H^{VAD} or f_D^{VAD} and is computationally simpler. Similarly, the estimated SNR measure f_S^{VAD} provides good results and is easier to compute than the entropy or dispersion of emission likelihoods.

VI. SUMMARY AND CONCLUSIONS

We presented our work on stream weighting for AV-ASR systems, where weights are introduced to control the contribution of each stream to the the recognition task. We focus on the use of dynamic weights in changing environmental conditions, setting the weights proportional to different measures of the confidence associated to the stream. The main contributions of the paper are the following: the experimental investigation of dynamic weighting schemes in different noisy environments and system architectures, the effectiveness of introducing a

VAD in the weighting schema and the proposal of a new confidence measure.

Based on the signal itself we estimate the SNR present on the audio channel, while we measure the classifier's confidence associated to the stream in terms of the dispersion and the entropy of the class probability distributions. We show how each measure is implemented and suits HMM-ANNs or HMM-GMMs systems and propose a new measure based on the transition probabilities common to both HMM architectures. Evaluating the different stream confidence measures and taking into account the classifiers behaviour for the different speech classes, we improve recognition results by the introduction of different mappings for the speech and silence classes.

Experimental results show that dynamic weights perform well in a variety of conditions. For high and medium SNR ratios, a weighting algorithm based on the classifier's reliability estimators performs well because audio and visual streams incur in uncorrelated errors that can be avoided by the audio-visual system. For very noisy environments, however, the confusion with the silence class is the main cause of failure of the systems and the weighting should first avoid the confusions with the silence class and then focus on recognition of speech. In fact, statistical analysis of the results show that the increase in performance associated to differentiating between silence and speech on the definition of the stream weights is more relevant than any difference in performance between different reliability measures.

REFERENCES

- [1] G. Potamianos, C. Neti, J. Luetin, and I. Matthews, "Audio-visual automatic speech recognition: an overview," in *Issues in audio-visual speech processing*, G. Bailly, E. Vatikiotis-Bateson, and P. Perrier, Eds. MIT Press, 2004.
- [2] S. Dupont and J. Luetin, "Audio-visual speech modeling for continuous speech recognition," in *IEEE Transactions on Multimedia*, vol. 2, 2000, pp. 141–151.
- [3] N. Fox, R. Gross, J. Cohn, and R. Reilly, "Robust biometric person identification using automatic classifier fusion of speech, mouth, and face experts," *IEEE Transactions on Multimedia*, vol. 9, no. 4, pp. 701–714, 2007.

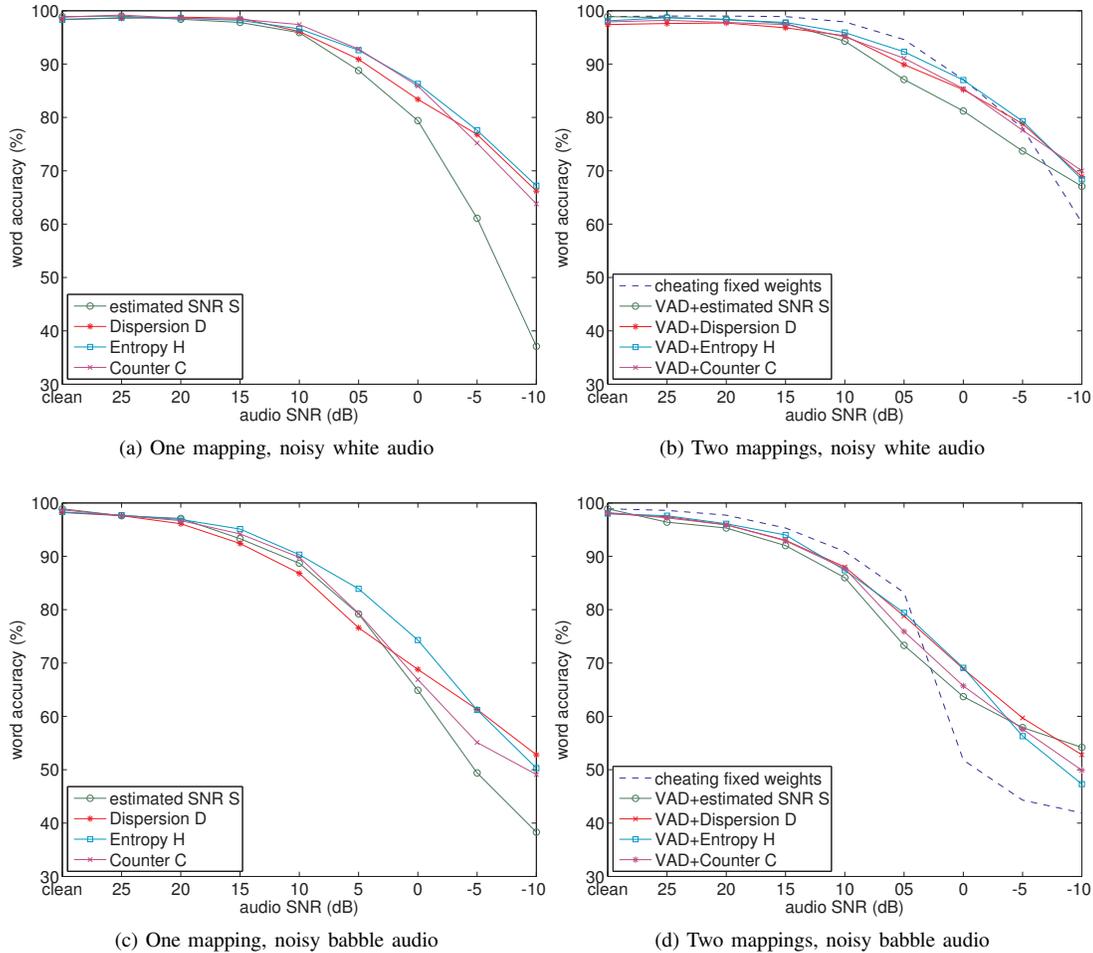


Fig. 3. Performance of cheat fixed weights HMM-GMM system and dynamic weights obtained with one and two mappings. Comparing to the fixed “cheat” system shows how far we are from the best behaviour under the assumption that the weights depend only on the SNR. When a dynamic system outperforms the fixed one, it is due to the fact that the silence/speech class should also be taken into account for the weight definition. Comparing Fig.3a with 3b and 3c with 3d we observe the improvement obtained by the inclusion of the VAD in the GMM system for white and babble noise.

- [4] A. Morris, A. Hagen, H. Glotin, and H. Bourlard, “Multi-stream adaptive evidence combination for noise robust ASR,” *Speech Communication*, vol. 34, no. 1-2, pp. 25–40, 2001.
- [5] G. Potamianos and H. Graf, “Discriminative training of HMM stream exponents for audio-visual speech recognition,” in *International Conference on Acoustics, Speech and Signal Processing*, 1998, pp. 3733–3736.
- [6] S. Nakamura, H. Ito, and K. Shikano, “Stream weight optimization of speech and lip image sequence for audio-visual speech recognition,” *International Conference on Spoken Language Processing*, vol. III, pp. 20–23, 2000.
- [7] G. Gravier, S. Axelrod, G. Potamianos, and C. Neti, “Maximum entropy and MCE based HMM stream weight estimation for audio-visual ASR,” *International Conference on Acoustics, Speech and Signal Processing*, 2002.
- [8] C. Miyajima, K. Tokuda, and T. Kitamura, “Audio-visual speech recognition using MCE-based HMMs and model-dependent stream weights,” *International Conference on Spoken Language Processing*, vol. II, pp. 1023–1026, 2000.
- [9] E. Sánchez-Soto, A. Potamianos, and K. Daoudi, “Unsupervised stream-weights computation in classification and recognition tasks,” *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 17, no. 3, pp. 436–445, 2009.
- [10] A. Potamianos, E. Sanchez-Soto, and K. Daoudi, “Stream weight computation for multi-stream classifiers,” in *International Conference on Acoustics, Speech and Signal Processing*, vol. 1, 2006.
- [11] S. Tamura, K. Iwano, and S. Furui, “A stream-weight optimization method for audio-visual speech recognition using multi-stream HMMs,” *International Conference on Acoustics, Speech and Signal Processing*, vol. 1, pp. 857–860, 2004.
- [12] —, “A stream-weight optimization method for multi-stream HMMs based on likelihood value normalization,” in *International Conference on Acoustics, Speech and Signal Processing*, 2005, pp. 468–472.
- [13] A. Adjoudani and C. Benoît, “On the integration of auditory and visual parameters in an HMM-based ASR,” in *Speechreading by humans and machines*, D. Stork and M. Hennecke, Eds. Springer, 1996, pp. 461–471.
- [14] S. Cox, I. Matthews, and A. Bangham, “Combining noise compensation with visual information in speech recognition,” *European Tutorial Workshop on Audio-Visual Speech Processing*, 1997.
- [15] S. Gurbuz, Z. Tufekci, E. Patterson, and J. Gowdy, “Application of affine-invariant Fourier descriptors to lipreading for audio-visual speech recognition,” *International Conference on Acoustics, Speech and Signal Processing*, pp. 177–180, 2001.
- [16] P. Teissier, J. Robert-Ribes, and J. Schwartz, “Comparing models for audiovisual fusion in a noisy-vowel recognition task,” *IEEE Transactions on Speech and Audio Processing*, vol. 7, pp. 629–642, 1999.
- [17] U. Meier, W. Hurst, and P. Duchnowski, “Adaptive bimodal sensor fusion for automatic speechreading,” *International Conference on Acoustics, Speech and Signal Processing*, pp. 833–836, 1996.
- [18] H. Glotin, D. Vergyri, C. Neti, G. Potamianos, and J. Luetin, “Weighting schemes for audio-visual fusion in speech recognition,” *International Conference on Acoustics, Speech and Signal Processing*, pp. 173–176, 2001.
- [19] R. Seymour, D. Stewart, and J. Ming, “Audio-visual integration for robust speech recognition using maximum weighted stream posteriors,” *Proceedings of Interspeech*, 2007.

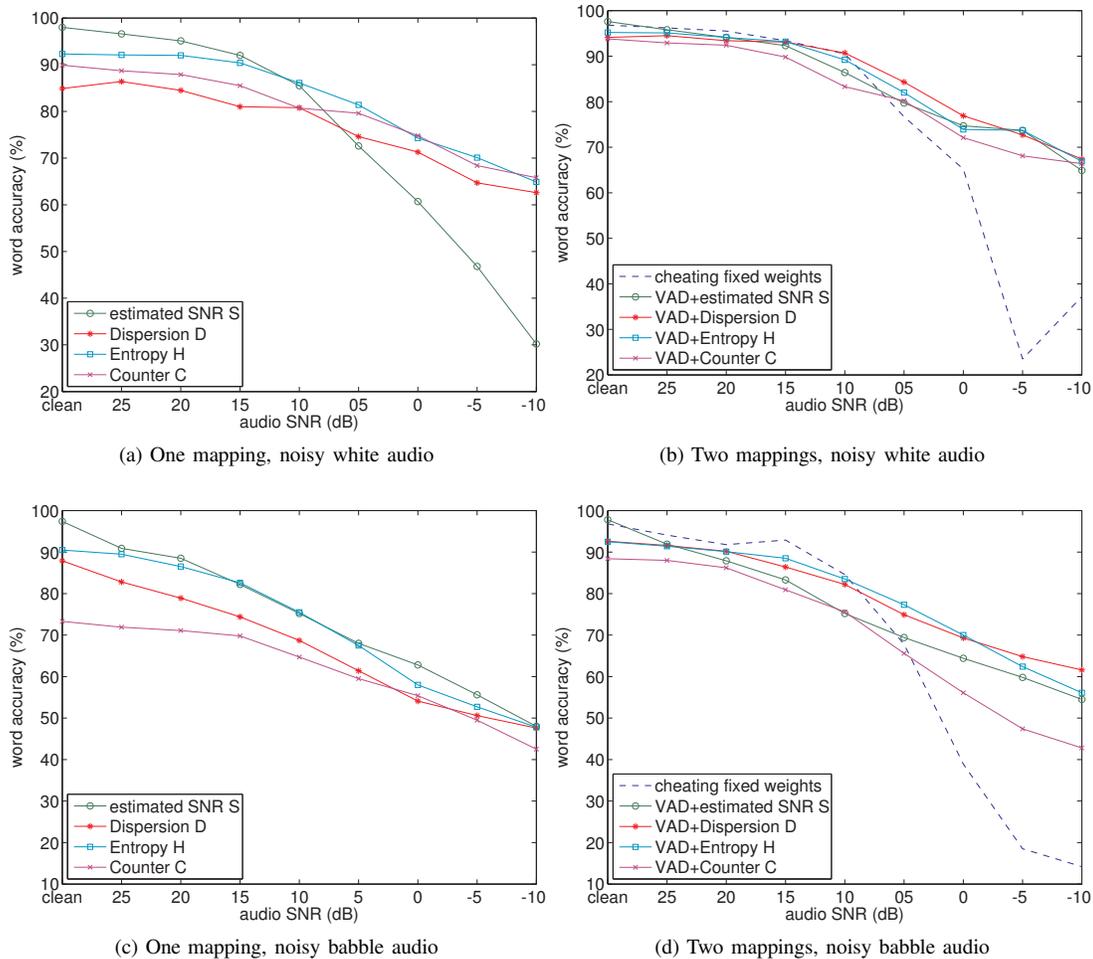


Fig. 4. Performance of cheat fixed weights HMM-ANN system and dynamic weights obtained with one and two mappings. Comparing to the fixed “cheat” system shows how far we are from the best behaviour under the assumption that the weights depend only on the SNR. When a dynamic system outperforms the fixed one, it is due to the fact that the silence/speech class should also be taken into account for the weight definition. Comparison of Fig.4a with 4b and 4c with 4d show the improvement obtained by the inclusion of the VAD in the ANN system.

- [20] G. Potamianos and C. Neti, “Stream confidence estimation for audio-visual speech recognition,” *International Conference on Spoken Language Processing*, 2000.
- [21] G. Potamianos, C. Neti, G. Gravier, A. Garg, and A. Senior, “Recent advances in the automatic recognition of audio-visual speech,” *Proceedings of the IEEE*, vol. 91(9), 2003.
- [22] M. Heckmann, F. Berthommier, and K. Kroschel, “Noise adaptive stream weighting in audio-visual speech recognition,” *EURASIP Journal on Applied Signal Processing*, vol. 2002, pp. 1260–1273, 2002.
- [23] E. Marcheret, V. Libal, and G. Potamianos, “Dynamic stream weight modeling for audio-visual speech recognition,” in *International Conference on Acoustics, Speech and Signal Processing*, vol. 4, 2007.
- [24] A. Garg, G. Potamianos, C. Neti, and T. Huang, “Frame-dependent multi-stream reliability indicators for audio-visual speech recognition,” in *International Conference on Multimedia and Expo*, 2003, pp. 605–608.
- [25] L. Rabiner and B. Juang, “An introduction to Hidden Markov Models,” *IEEE ASSP Magazine*, vol. 3, no. 1, pp. 4–16, 1986.
- [26] L. Rabiner and B.-H. Juang, *Fundamentals of Speech Recognition*. Prentice-Hall Signal Processing Series, 1993.
- [27] C. Bishop, *Neural Networks for Pattern Recognition*. Oxford University Press, 1995.
- [28] N. Morgan and H. Bourlard, “Continuous speech recognition, an introduction to the hybrid HMM/connectionist approach,” *IEEE Signal Processing Magazine*, vol. 12, no. 3, pp. 25–42, 1995.
- [29] A. Ganapathiraju, J. Hamaker, and J. Picone, “Hybrid SVM/HMM architectures for speech recognition,” *International Conference on Spoken Language Processing*, vol. 4, pp. 504–507, 2000.
- [30] J. Movellan and G. Chadderdon, “Channel separability in the audio-visual integration of speech: A Bayesian approach,” *NATO ASI SERIES F COMPUTER AND SYSTEMS SCIENCES*, vol. 150, pp. 473–488, 1996.
- [31] D. Massaro and D. Stork, “Speech recognition and sensory integration,” *American Scientist*, vol. 86, no. 3, pp. 236–244, 1998.
- [32] J. Kittler, M. Hatef, R. Duin, and J. Matas, “On combining classifiers,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 20, no. 3, pp. 226–239, 1998.
- [33] K. Kirchhoff and J. Bilmes, “Dynamic classifier combination in hybrid speech recognition systems using utterance-level confidence values,” *International Conference on Acoustics, Speech, and Signal Processing*, pp. 693–696, 1999.
- [34] F. Berthommier and H. Glotin, “A new SNR-feature mapping for robust multistream speech recognition,” *International Congress on Phonetic Sciences*, pp. 711–715, 1999.
- [35] L. Terry, D. Shiell, and A. Katsaggelos, “Feature space video stream consistency estimation for dynamic stream weighting in audio-visual speech recognition,” in *International Conference on Image Processing*, 2008, pp. 1316–1319.
- [36] X. Shao and J. Barker, “Stream weight estimation for multistream audio-visual speech recognition in a multispeaker environment,” *Speech Communication*, vol. 50, no. 4, pp. 337–353, 2008.
- [37] J. Sohn, N. Kim, and W. Sung, “A statistical model-based voice activity detection,” *IEEE Signal Processing Letters*, vol. 6, no. 1, pp. 1–3, 1999.
- [38] H. Misra, H. Bourlard, and V. Tyagi, “New entropy based combination rules in HMM/ANN multi-stream ASR,” in *International Conference on Acoustics, Speech, and Signal Processing*, 2003.

- [39] J. Moré and D. Sorensen, "Computing a trust region step," *SIAM Journal on Scientific and Statistical Computing*, vol. 4, p. 553, 1983.
- [40] E. Patterson, S. Gurbuz, Z. Tufekci, and J. Gowdy, "Moving-talker, speaker-independent feature study and baseline results using the CUAVE multimodal speech corpus," *EURASIP Journal on Applied Signal Processing*, vol. 2002(11), pp. 1189–1201, 2002.
- [41] A. Varga, H. Steeneken, M. Tomlinson, and D. Jones, "The NOISEX-92 study on the effect of additive noise on automatic speech recognition," *DRA Speech Research Unit, Malvern, England, Tech. Rep.*, 1992.
- [42] E. Lombard, "Le signe de l'élevation de la voix," *Ann. Maladies Oreille, Larynx, Nez, Pharynx*, vol. 37, no. 101-119, p. 25, 1911.
- [43] M. Gurban and J.-P. Thiran, "Information theoretic feature extraction for audio-visual speech recognition," *IEEE Transactions on Signal Processing*, vol. 57, no. 12, pp. 4765 – 4776, 2009.
- [44] G. Potamianos and C. Neti, "Audio-visual speech recognition in challenging environments," in *Eighth European Conference on Speech Communication and Technology*, 2003.
- [45] K. Livescu, O. Cetin, M. Hasegawa-Johnson, S. King, C. Bartels, N. Borges, A. Kantor, P. Lal, L. Yung, A. Bezman *et al.*, "Articulatory feature-based methods for acoustic and audio-visual speech recognition," in *Final Workshop Report, Center for Language and Speech Processing, John Hopkins University*, vol. 4, 2006.
- [46] S. Young, D. Kershaw, J. Odell, D. Ollason, V. Valtchev, and P. Woodland, *The HTK Book*. Cambridge, Entropic Ltd., 1999.
- [47] D. Dean, P. Lucey, S. Sridharan, and T. Wark, "Weighting and normalisation of synchronous HMMs for audio-visual speech recognition," *Auditory-Visual Speech Processing, Hilvarenbeek, The Netherlands, September*, pp. 110–115, 2007.
- [48] M. Bisani and H. Ney, "Bootstrap estimates for confidence intervals in ASR performance evaluation," in *IEEE International Conference on Acoustics, Speech, and Signal Processing*, 2004.