

Use of Bimodal Coherence to Resolve the Permutation Problem in Convolutional BSS

Qingju Liu^{1,*}, Wenwu Wang^{2,*}, Philip Jackson^{3,*}

Centre for Vision, Speech and Signal Processing, Faculty of Engineering and Physical Sciences, University of Surrey, Guildford, GU2 7XH, United Kingdom

Abstract

Recent studies show that facial information contained in visual speech can be helpful for the performance enhancement of audio-only blind source separation (BSS) algorithms. Such information is exploited through the statistical characterisation of the coherence between the audio and visual speech using, e.g., a Gaussian mixture model (GMM). In this paper, we present three contributions. With the synchronized features, we propose an adapted expectation maximization (AEM) algorithm to model the audio-visual coherence in the off-line training process. To improve the accuracy of this coherence model, we use a frame selection scheme to discard nonstationary features. Then with the coherence maximization technique, we develop a new sorting method to solve the permutation problem in the frequency domain. We test our algorithm on a multimodal speech database composed of different combinations of vowels and consonants. The experimental results show that our

*Corresponding authors. E-mails: {Q.Liu, W.Wang, P.Jackson}@surrey.ac.uk.

Fax: +44 (0) 1483 686031

¹Tel: +44 (0) 1483 683413

²Tel: +44 (0) 1483 686039

³Tel: +44 (0) 1483 686044

proposed algorithm outperforms traditional audio-only BSS, which confirms the benefit of using visual speech to assist in separation of the audio.

Keywords: convolutive blind source separation (BSS), audio-visual coherence, Gaussian mixture model (GMM), feature selection and fusion, adapted expectation maximization (AEM), indeterminacy

1. Introduction

Human speech perception is essentially bimodal as speech is perceived by the interactions of auditory and visual sensory processing [1, 2]. Looking at the speaker’s lips improves the intelligibility of human speech embedded in cocktail party noise due to the contribution of the complementary visual information [2]. There is a complex non-linear relationship between the auditory and visual streams, usually referred to as the audio-visual coherence or correlation [3]. In feature space, the coherence can be coded by audio-visual atoms or dictionaries [4, 5] with matching pursuit [6] techniques, or characterised statistically with models such as Gaussian mixture models (GMM) [7]. Exploiting these cross-modal interactions, the visual stream has proven a success in improving the robustness to noise in many fields of applications, including automatic speech recognition [8], speaker localization [4, 9], speech enhancement or audio filtering [10, 11], and blind source separation [3, 5, 12–16].

In traditional blind source separation (BSS) for auditory mixtures, typically only audio signals are considered. Under the framework of independent component analysis (ICA) [17], the BSS problems have been extensively studied and many classical algorithms have been proposed for the instantaneous

20 mixing model such as the “J-H” algorithm [18], JADE [19], Infomax [20],
21 SOBI [21] and FastICA [22] algorithms. For the more complex convolutive
22 mixing model, one can apply either the time domain deconvolution algo-
23 rithms [23–25], or the frequency domain separation algorithms [12–15, 26–31],
24 which often suffer from the permutation and scaling ambiguity problems.

25 Considering the bimodal nature of human speech, we could potentially
26 improve the separation of the source signals from their audio mixtures utiliz-
27 ing the audio-visual coherence obtained by the integration of visual speech.
28 This is known as audio-visual or bimodal BSS [3, 5, 12, 13, 15, 16], a recent
29 development in multi-modal signal processing. Soderoy et al. [3] addressed
30 the separation problem for an instantaneous mixture of decorrelated sources,
31 with no further assumptions on independence or non-Gaussianity. Wang et
32 al. [13] implemented a similar idea by applying the Bayesian framework to the
33 fused feature observations for both instantaneous and convolutive mixtures.
34 Rivet et al. [12] proposed a new statistical tool utilizing the log-Rayleigh
35 distribution for modelling the audio-visual coherence, and then used the co-
36 herence to address the permutation and scaling ambiguities in the spectral
37 domain. Casanovas et al. [5] detected temporal audio-visual structures rep-
38 resented by atoms taken from redundant dictionaries, and extracted sources
39 from a soundtrack. Naqvi et al. [16] utilized beamforming in the frequency
40 domain for moving sources in the teleconference-like scenario, incorporating
41 the geometrical model derived on the basis of the beamforming theory.

42 Despite being promising, these approaches are also limited in some sit-
43 uations. For example, the algorithm proposed in [3] was designed only for
44 instantaneous mixtures. The method in [13] considered a convolutive model

45 with a relatively small number of taps for the mixing filters. The approach
46 in [12] modelled the audio-visual coherence in a high dimensional feature
47 space, which often results in an over-fitting problem and therefore is sensi-
48 tive to outliers. Cross-modal correlation was not exploited in the separation
49 stage in [5], where visual information was used only for voice activity detec-
50 tion. In [16], the video provided the position information about the distance
51 and azimuth angles between the moving speakers and the microphone array,
52 however, source separation was still performed in the audio domain.

53 In this paper, we attempt to address some of these limitations. Moti-
54 vated by the work in [12, 13], we follow a similar two-stage framework which
55 includes off-line training and online separation. In particular, we consider a
56 convolutive mixing model and address the permutation problem associated
57 with the frequency domain BSS (FD-BSS). In the off-line training stage, we
58 build a model to statistically characterise the audio-visual coherence in the
59 feature space. This coherence is built on the audio-visual features extracted
60 from the target speech. Mel-frequency cepstral coefficients (MFCCs) are used
61 as the audio features, and the lip width and height as visual features, which
62 are synchronised with the audio features on a frame-by-frame basis before
63 statistical training. In the separation stage, coherence maximization is ap-
64 plied for the alignment of the ICA-separated spectral components. Different
65 from [12, 13], however, we have proposed three new techniques to improve
66 the training and separation processes. First, a frame selection scheme is pro-
67 posed to remove the non-stationary features which consequently improves
68 the robustness and accuracy of the estimation of the audio-visual coherence.
69 Second, the classical expectation maximisation (EM) algorithm is modified

70 to take into account the different influences of the audio features, resulting in
71 an adapted EM (AEM) algorithm, which further improves the estimation of
72 the joint audio-visual probability. Third, a novel sorting scheme is proposed
73 to address the permutation problem. A preliminary version of this work was
74 presented in [15]. Different from [15], in this paper, we have developed a ro-
75 bust feature selection scheme for audio-visual modelling as mentioned above.
76 In addition, we have further improved the audio feature representation as
77 described in Section 3.1. Moreover, here we have performed systematic eval-
78 uations on real recordings, and compared the performance of the proposed
79 method with the state-of-the-art methods.

80 The remainder of the paper is organised as follows. An overview of tradi-
81 tional frequency domain convolutive BSS and the framework of the proposed
82 audio-visual BSS system are presented in Section 2. Then Section 3 in-
83 troduces the feature extraction and fusion method for the modelling of the
84 cross-model correlation, including a new frame selection approach and an
85 adapted expectation maximisation algorithm to improve the accuracy of this
86 model. The proposed de-permutation algorithm exploiting the audio-visual
87 coherence is presented in Section 4. The simulation results are analysed and
88 discussed in Section 5, followed by the conclusions.

89 **2. BSS for Convolutional Mixtures**

90 *2.1. Convolutional Model*

BSS aims to recover sources from their mixtures without any or with little prior knowledge about the sources or the mixing process. Consider a cocktail party scenario, the observation at each sensor is the sum of K filtered source

signals, which can be approximated by the convolutive model:

$$\begin{aligned}
 x_p(n) &= \sum_{k=1}^K \sum_{m=0}^{+\infty} h_{pk}(m) s_k(n-m) + \xi_p(n), \\
 \mathbf{x}(n) &= \mathbf{H} * \mathbf{s}(n) + \boldsymbol{\xi}(n),
 \end{aligned} \tag{1}$$

91 where h_{pk} represents the room impulse response filter from source k to sensor
 92 p . We denote $\mathbf{x}(n) = [x_1(n), \dots, x_P(n)]^T$ as the observation vector at the
 93 discrete time index n ; $\mathbf{s}(n) = [s_1(n), \dots, s_K(n)]^T$ the source vector and $\boldsymbol{\xi}(n) =$
 94 $[\xi_1(n), \dots, \xi_P(n)]^T$ the additive noise vector, where T is vector transpose. \mathbf{H} is
 95 the mixing matrix whose elements are filters h_{pk} and $*$ denotes convolution.

Convolutive BSS aims to find a set of separation filters $\{w_{kp}\}$ that satisfy:

$$\begin{aligned}
 \hat{s}_k(n) = y_k(n) &= \sum_{p=1}^P \sum_{m=0}^{+\infty} w_{kp}(m) x_p(n-m), \\
 \hat{\mathbf{s}}(n) = \mathbf{y}(n) &= \mathbf{W} * \mathbf{x}(n),
 \end{aligned} \tag{2}$$

96 where \mathbf{W} is the separation matrix whose entry w_{kp} is the impulse response
 97 filter from observation p to the estimate of source k ($\mathbf{y}(n)$ or $\hat{\mathbf{s}}(n)$ repre-
 98 sents the estimated version of $\mathbf{s}(n)$). We consider a time-invariant system
 99 where both the mixing filter \mathbf{H} and separation filter \mathbf{W} are assumed to be
 100 time-invariant. In practice, a finite impulse response (FIR) filter is used to
 101 implement h_{pk} and w_{kp} .

102 2.2. Frequency domain BSS

Convolutive BSS can be directly performed in the time domain [23–25] by deconvolution, but the computational complexity is high especially when the mixing filters have long taps. Based on the short-time stationarity of the speech signals and the linear time-invariance of the mixing process, an

alternative is to perform convolutive BSS in the frequency domain by applying the short-time Fourier transform (STFT) to the observations. In each frequency bin f , we get an instantaneous mixing model ignoring the noise:

$$\mathbf{X}(f, t) = \mathbf{H}(f)\mathbf{S}(f, t), \quad (3)$$

103 where $\mathbf{X}(f, t) = [X_1(f, t), \dots, X_P(f, t)]^T$ is the observation vector in frequency
 104 bin f and time frame t , and $\mathbf{H}(f)$ is the Fourier transform of \mathbf{H} .

Then in each frequency bin f , separate ICA [17] algorithms for instantaneous models are applied to obtain the independent outputs $\mathbf{Y}(f, t) = [Y_1(f, t), \dots, Y_K(f, t)]^T$, assumed to be the source estimates:

$$\mathbf{Y}(f, t) = \mathbf{W}(f)\mathbf{X}(f, t) = \hat{\mathbf{S}}(f, t). \quad (4)$$

105 The technique of the frequency domain BSS is depicted in the upper dashed
 106 box of Figure 1. In this paper, a determined system is considered, i.e.,
 107 $K = P = 2$.

108 2.3. Scaling and Permutation Indeterminacy

However, the ICA algorithms can estimate the sources only up to a permutation matrix $\mathbf{P}(f)$ and a diagonal matrix $\mathbf{D}(f)$:

$$\hat{\mathbf{S}}(f, t) = \mathbf{Y}(f, t) = \mathbf{P}(f)\mathbf{D}(f)\mathbf{S}(f, t). \quad (5)$$

109 These are the so-called permutation ($\mathbf{P}(f)$) and scaling ($\mathbf{D}(f)$) ambiguities,
 110 which present severe problems when reconstructing the separated sources in
 111 the time domain. The scaling ambiguity can be greatly mitigated by the
 112 normalization of the separation matrices based on the minimal distortion
 113 principle (MDP) [29] **which is also used here.** In this paper, we only

114 consider the permutation problem, where the order of the recovered source
 115 components at each frequency bin may not be consistent with each other
 116 (e.g. $\begin{bmatrix} Y_1(f_1,t) \\ Y_2(f_1,t) \end{bmatrix} = \begin{bmatrix} S_1(f_1,t) \\ S_2(f_1,t) \end{bmatrix}, \begin{bmatrix} Y_1(f_2,t) \\ Y_2(f_2,t) \end{bmatrix} = \begin{bmatrix} S_2(f_2,t) \\ S_1(f_2,t) \end{bmatrix}, \begin{bmatrix} Y_1(f_3,t) \\ Y_2(f_3,t) \end{bmatrix} = \begin{bmatrix} S_1(f_3,t) \\ S_2(f_3,t) \end{bmatrix}, \dots$
 117).

118 To address the permutation problem, many algorithms [26–28, 30–33]
 119 have been proposed. For example, the approach in [26, 30] utilizes the conti-
 120 nuity of the spectral components in adjacent frequency channels while [27, 28]
 121 use direction of arrival estimation, [31] combines both previous techniques,
 122 and [32] utilizes statistical signal models. However, the performance of these
 123 algorithms can be degraded by acoustical noise. Information from the video
 124 has been shown to be useful for improving the performance of automatic
 125 speech recognition systems [8]. The potential of using visual information for
 126 audio source separation problems, such as the permutation problem, has not
 127 been fully investigated, which motivates this study, as now discussed.

128 *2.4. Overview of the Proposed System*

129 To address the permutation problem in FD-BSS, we use an audio-visual
 130 BSS system, shown in Figure 1. As mentioned in Section 1, the system con-
 131 tains two stages: training stage and separation stage. The training stage is
 132 shown in the lower dashed box of Figure 1, which includes feature extrac-
 133 tion and feature fusion. First, we extract the audio features $\mathbf{a}(t)$ from the
 134 training audio data $s(t)$, and some geometric-type features $\mathbf{v}(t)$ from the
 135 training video stream $v(t)$. Second, we use the GMM to statistically charac-
 136 terise the audio-visual coherence $p(\mathbf{a}(t), \mathbf{v}(t))$, and then an AEM algorithm
 137 is applied to estimate the parameters of the GMM model. The separation
 138 stage is shown in the upper dashed box, which is performed in the audio

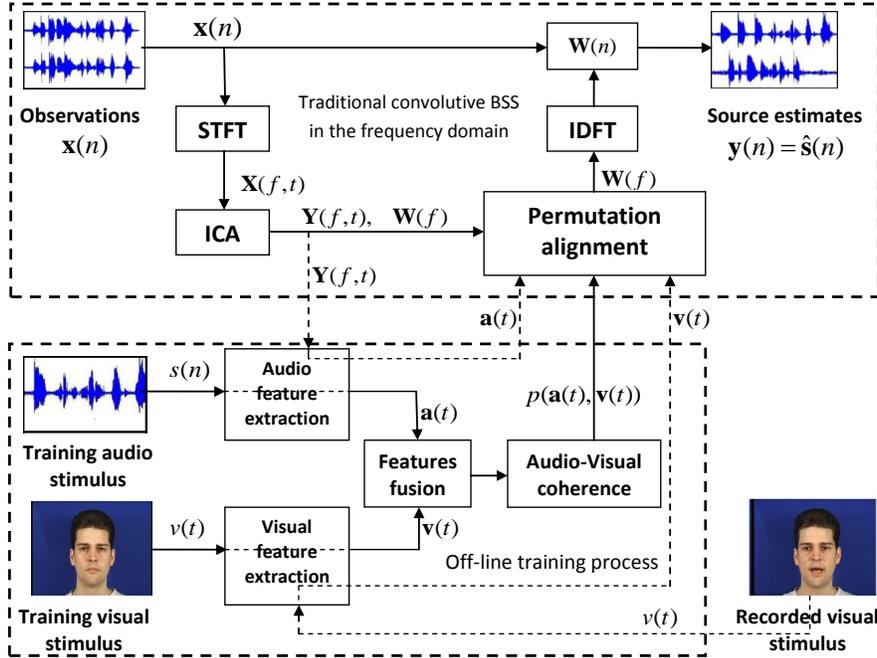


Figure 1: Flow of the proposed audio-visual BSS system. The upper dashed box illustrates the traditional audio-only BSS in the frequency domain. The lower dashed box shows the training process, which aims to estimate the parameters for the audio-visual coherence after audio-visual fusion. Before the reconstruction of the recovered signals in the time domain, we resolve the permutation indeterminacy by audio-visual coherence maximization using the feature vectors $\mathbf{a}(t)$ and $\mathbf{v}(t)$ in dashed lines (with arrows).

139 domain. To address the permutation problem in the separation stage, we use
 140 the information (i.e. the audio-visual joint probability) obtained from the
 141 training stage. Specifically, we align the permutations across the frequency
 142 bands based on a sorting scheme which iteratively re-estimates the audio-
 143 visual coherence probability from the separated source components at each
 144 frequency bin based on the trained audio-visual model.

145 3. Feature Extraction and Fusion

146 Our algorithm is based on the fact that there is a relationship between
147 the video signal and the corresponding contemporary audio signal, which is
148 the so-called audio-visual coherence. We model the coherence in the feature
149 level, for which the features are extracted from the audio and video data re-
150 spectively. Since the two types of signals are recorded with different sampling
151 rates and dimensions, we need to extract the features from them separately.

152 3.1. Extraction of Audio and Visual Features

153 We take the Mel-frequency cepstral coefficients (MFCCs) as audio fea-
154 tures as in [13] with some modifications. The MFCCs exploit the non-linear
155 resolution of the human auditory system across an audio spectrum, which
156 are the Discrete Cosine Transform (DCT) results of the logarithm of the
157 short time power spectrum on a Mel-frequency scale. We then apply a lifter
158 for the reweighing of cepstral coefficients and obtain a $(L + 1)$ -dimensional
159 MFCC vector $[c_0(t), c_1(t), \dots, c_L(t)]^T$, where t is the time frame index, as in
160 equation 3. Compared to our early work in [15], where the first compo-
161 nent is the logarithmic power, we remove the first coefficient $c_0(t)$ to avoid
162 the influence of the magnitude, and obtain the L -dimensional audio feature
163 $\mathbf{a}(t) = [c_1(t), \dots, c_L(t)]^T$. For simplicity, we denote the training audio feature
164 vector as $\mathbf{a}(t) = [a_1(t), \dots, a_L(t)]^T$.

165 Unlike the appearance-based visual features used in [8, 13], which are
166 sensitive to the lighting conditions, we use the same front geometric visual
167 features as in [3, 12]: the lip width (LW) and height (LH) from the in-
168 ternal labial contour. The geometric features $\mathbf{v}(t) = [\text{LW}(t), \text{LH}(t)]^T$ are

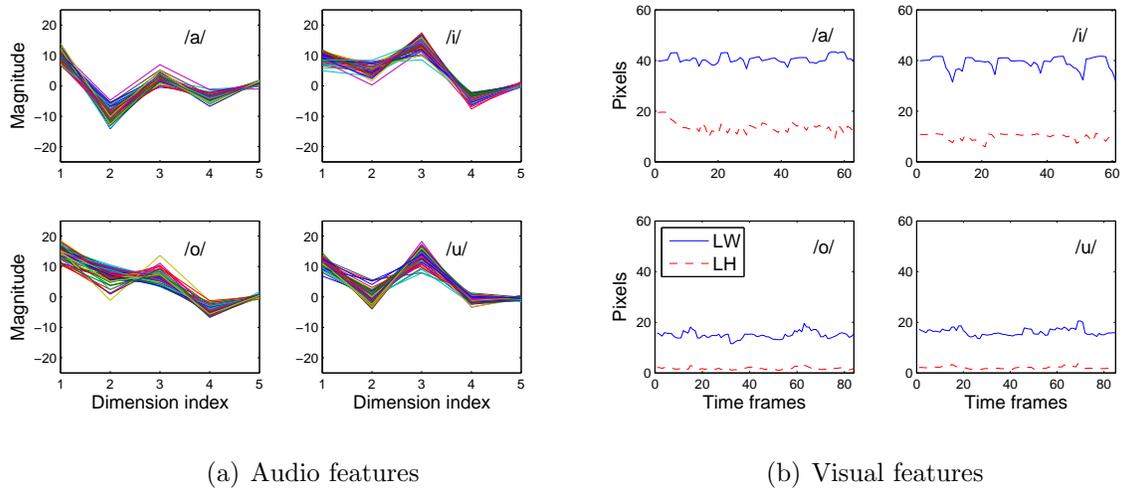


Figure 2: The audio-visual features of four vowels /a/, /i/, /o/, /u/. The vowels are obtained by selecting and concatenating the first several frames in each audio sequence in the data corpus. (a) Each curve corresponds to one time frame, associated to the time frame in the right half plot. We choose $L = 5$, therefore the audio feature is 5-dimensional. (b) The visual features last about 60–80 frames at 50 Hz.

169 low-dimensional and robust to luminance, which mitigates the over-fitting
 170 problem encountered in complex systems with a large number of parameters
 171 and greatly reduces the computational complexity. Figure 2 shows the typi-
 172 cal features extracted from four vowels: /a/, /i/, /o/ and /u/, which were
 173 obtained from the multimodal database, as depicted in Section 5, which is
 174 also used in Figures 3 and 4. After the feature extraction process, we concate-
 175 nate synchronized audio and visual features to build the $(L + 2)$ -dimensional
 176 audio-visual space $\mathbf{u}(t) = [\mathbf{v}(t); \mathbf{a}(t)]$, which will be used for training.

177 *3.2. Robust Feature Frame Selection*

178 If someone utters an isolated speech sound such as /a/, the visual features
 179 will likely be stationary with minimal fluctuation. However in the transition
 180 periods from one phone to another in fluent speech, the visual parameters
 181 fluctuate drastically with a large variance, which can produce ambiguous
 182 visual features typical of another speech sound or phone. For instance, in
 183 the transition process from /a/ to /b/, several frames of the mouth shape
 184 may look like the utterance of /o/. Also, these transitions are not stationary
 185 in the audio signal. Therefore, to improve the estimate of audio-visual co-
 186 herence, we propose a feature frame selection scheme based on the dynamic
 187 characteristics [34] of the visual features.

At each time frame centred by the visual feature $\mathbf{v}(t) = [\text{LW}(t), \text{LH}(t)]^T$, we extract a short time period with $2Q + 1$ frames, then calculate

$$\gamma_{\text{LW}}(t) = \sigma(\text{LW}(t)) + \alpha_{\text{LW}} \|\text{LW}(t + Q) - \text{LW}(t - Q)\|, \quad (6)$$

where $\sigma(\cdot)$ is the standard deviation and α_{LW} is a weighting coefficient, chosen between 0 and 1. Then we define a Boolean variable to determine the stationarity of this frame

$$\mathcal{F}_{\text{LW}}(t) \stackrel{\text{def}}{=} \begin{cases} 1, & \gamma_{\text{LW}}(t) < \delta_{\text{LW}} \overline{\text{LW}(t)} \\ 0, & \text{otherwise} \end{cases} \quad (7)$$

where δ_{LW} is a comparison coefficient, typically chosen as 0.5, and $\overline{\text{LW}(t)}$ is the mean over the $2Q + 1$ frames, defined as $\overline{\text{LW}(t)} = \frac{1}{2Q+1} \sum_{q=-Q}^Q \text{LW}(t+q)$. Then, we smooth the binary variable between adjacent frames

$$\mathcal{F}_{\text{LW}}^{\text{S}}(t) = \mathcal{F}_{\text{LW}}(t-1) \vee \mathcal{F}_{\text{LW}}(t) \vee \mathcal{F}_{\text{LW}}(t+1), \quad (8)$$

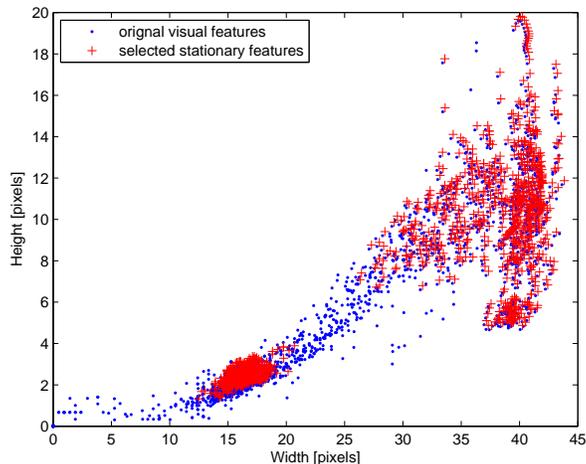


Figure 3: Visual frame selection scheme. Each dot (or cross) represents a two dimensional visual feature vector at one time frame. After frame selection, the transitions from one syllable to another have been removed. In other words, the features are better clustered (denoted by crosses), as compared to the original scatter plot of all the features (denoted by dots).

where \vee denotes disjunction, i.e., a logical OR operator. In the same way, we can determine $\mathcal{F}_{\text{LH}}^{\text{S}}(t)$, and the final decision is

$$\mathcal{F}(t) = \mathcal{F}_{\text{LW}}^{\text{S}}(t) \wedge \mathcal{F}_{\text{LH}}^{\text{S}}(t), \quad (9)$$

188 where \wedge denotes conjunction, i.e., a logical AND operator.

189 If $\mathcal{F}(t) = 1$, the frame will be chosen, otherwise it will be discarded. The
 190 audio-visual features associated with the selected frames are used in both
 191 the training and separation stages. Figure 3 shows an example of the visual
 192 features before (dot) and after (cross) selection. The feature selection has
 193 essentially removed the redundant unstable features and hence improves the
 194 spatial distribution of the clusters of the lip features.

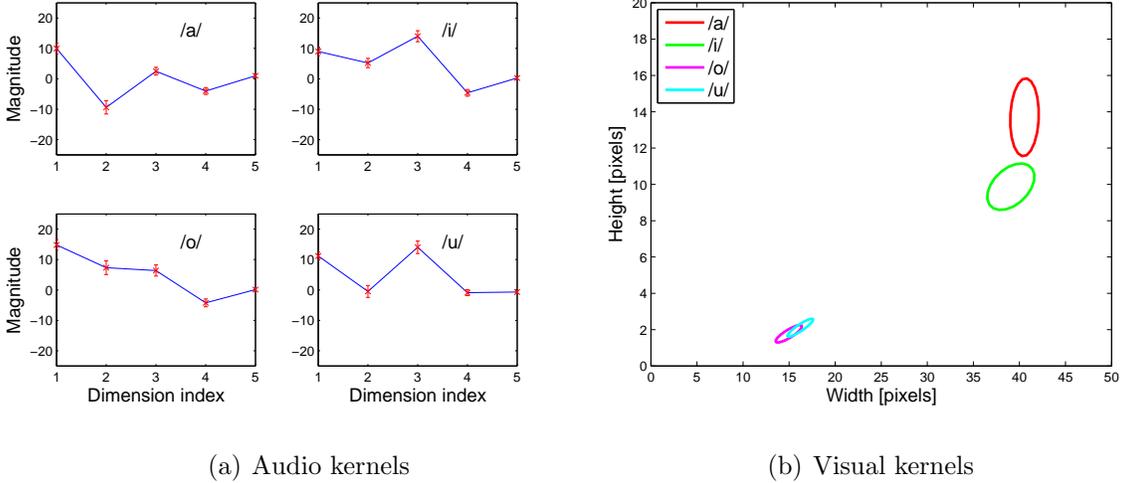


Figure 4: The audio-visual kernels for four vowels: /a/, /i/, /o/ and /u/. (a) Distributions of audio features described by mean (solid line) and the variance of one standard deviation (bar). (b) The visual kernels show the spacial distributions of the four vowels.

195 *3.3. Feature-Level Fusion*

The audio-visual coherence can be statistically characterized by a GMM with I kernels:

$$p_{AV}(\mathbf{u}(t)) = \sum_{i=1}^I \gamma_i \mathcal{N}(\mathbf{u}(t) \mid \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i), \quad (10)$$

where γ_i is the weighting parameter, $\boldsymbol{\mu}_i$ is the mean vector and $\boldsymbol{\Sigma}_i$ is the full covariance matrix of the i -th kernel. Every kernel of this mixture represents one cluster of the audio-visual data modelled by a multivariate normal distribution:

$$\mathcal{N}(\mathbf{u}(t) \mid \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) = \frac{\exp\{-\frac{1}{2}(\mathbf{u}(t) - \boldsymbol{\mu}_i)^T \boldsymbol{\Sigma}_i^{-1} (\mathbf{u}(t) - \boldsymbol{\mu}_i)\}}{\sqrt{(2\pi)^{L+2} |\boldsymbol{\Sigma}_i|}}. \quad (11)$$

196 We denote $\lambda_i = \{\gamma_i, \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i\}$ as the parameter set, which can be estimated
 197 by the expectation maximization (EM) algorithm. In the traditional EM

198 training process, all the components of the training data are treated equally
 199 whatever their magnitudes. Nevertheless, some components of the audio
 200 vector with large magnitudes are actually more informative about the audio-
 201 visual coherence than the remaining components (consider, for instance, the
 202 case of lossy compression of audio and images using DCT where small com-
 203 ponents can be discarded). For example, the first component of the audio
 204 vector ($a_1(t)$) should play a more dominant role in determining the proba-
 205 bility $p_{AV}(\mathbf{u}(t))$ than the last one. Also, the components of the audio vector
 206 having very small magnitudes are likely to be affected by noise. Therefore,
 207 considering these factors, we propose an adapted expectation maximization
 208 (AEM) algorithm.

209 I. Initialize the parameter set $\{\lambda_i\}$ with the K-means algorithm.

210 II. Run the following iterative process:

i. Compute the influence parameters $\beta_i(\cdot)$ of $\mathbf{u}(t)$ for $i = 1, \dots, I$.

$$\beta_i(\mathbf{u}(t)) = 1 - \frac{\|\mathbf{u}(t) - \boldsymbol{\mu}_i\|}{\sum_{j=1}^I \|\mathbf{u}(t) - \boldsymbol{\mu}_j\|}, \quad (12)$$

ii. Calculate the probability of each cluster given $\mathbf{u}(t)$.

$$p_i(\mathbf{u}(t)) = \frac{\gamma_i p_G(\mathbf{u}(t) | \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) \beta_i(\mathbf{u}(t))}{\sum_{j=1}^I \gamma_j p_G(\mathbf{u}(t) | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j) \beta_j(\mathbf{u}(t))}. \quad (13)$$

iii. Update the parameter set $\{\lambda_i\}$:

$$\begin{aligned} \boldsymbol{\mu}_i &= \frac{\sum_t \mathbf{u}(t) p_i(\mathbf{u}(t))}{\sum_t p_i(\mathbf{u}(t))}, \quad \gamma_i = \frac{\sum_t p_i(\mathbf{u}(t))}{\sum_t 1}, \\ \boldsymbol{\Sigma}_i &= \frac{\sum_t (\mathbf{u}(t) - \boldsymbol{\mu}_i)(\mathbf{u}(t) - \boldsymbol{\mu}_i)^T p_i(\mathbf{u}(t)) + c \{\boldsymbol{\Sigma}_p\}_i}{\sum_t p_i(\mathbf{u}(t)) + c}, \end{aligned} \quad (14)$$

211 where $\|\cdot\|$ denotes the squared Euclidean distance, and c is a constant
 212 chosen to be proportional to the number of samples and $\{\boldsymbol{\Sigma}_p\}_i$ is the penalty

213 term. Different from the traditional EM algorithm, we have added an in-
 214 fluence parameter $\beta_i(\mathbf{u}(t))$ in the expectation step of the AEM algorithm,
 215 where $\beta_i(\mathbf{u}(t))$ takes into account the various distance between the sample
 216 (i.e. the vector $\mathbf{u}(t)$) to different kernel centres, similar to the idea used in
 217 the classical K-means algorithm. Therefore, the different distance between
 218 the vector $\mathbf{u}(t)$ and a candidate cluster $\boldsymbol{\mu}_i$ will have a different impact on
 219 the probability $p_{AV}(\mathbf{u}(t))$. In addition, unlike our preliminary work in [15],
 220 we have added the penalty term $\boldsymbol{\Sigma}_p$ to avoid the covariance matrix becom-
 221 ing singular, which can occur when a kernel converges on one or two sample
 222 points (typically outliers). This has a similar effect to a variance floor. With-
 223 out such a safeguard, the probability would approach infinity in such cases,
 224 leading to numerical stability problems in practical implementation. The pa-
 225 rameter $\boldsymbol{\Sigma}_p$ can be chosen as a diagonal matrix and the subscript p denotes
 226 penalization.

227 Figure 4 shows the audio-visual kernels for the vowels used in Figure
 228 2. It can be observed that the visual kernels of /o/ and /u/ overlap with
 229 each other, but the related audio kernels differ greatly. On the other hand,
 230 the audio kernels of /i/ and /u/ look similar, but their visual kernels do
 231 not have overlap at all. This implies that the audio and visual features are
 232 complementary to each other.

233 4. Resolution of Permutation Problem

As $y_k(n)$ is the estimate of $s_k(n)$, $y_k(n)$ will have a maximum coherence
 with the corresponding video signal $v_k(t)$. Therefore we can maximize the
 following criterion in the frequency domain to address the permutation prob-

lem:

$$\hat{\mathbf{P}}(f) = \arg \max_{\mathbf{P}(f)} \sum_t \sum_{k=1}^K p_{AV}(\mathbf{u}_k(t)), \quad (15)$$

where $\mathbf{u}_k(t) = [\mathbf{a}_k(t); \mathbf{v}_k(t)]$ is the audio-visual feature extracted from the profile $\hat{S}_k(\cdot, t) = Y_k(\cdot, t)$ and the recorded video associated with the k -th speaker. If we are only interested in an estimate of $s_1(n)$ from the observations, we can get the first separation vector $\mathbf{p}(f)$ by maximizing:

$$\hat{\mathbf{p}}(f) = \arg \max_{\mathbf{p}(f)} \sum_t p_{AV}(\mathbf{u}_1(t)). \quad (16)$$

234 Since the permutation problem is the main factor in the degradation of
 235 the recovered sources, we focus on permutation indeterminacy cancellation
 236 for a two-source two-mixture case (i.e., $K = P = 2$). Assuming the spectral
 237 analysis window employs an fast Fourier transform (FFT) size of N samples
 238 from the audio mixture signal, based on the symmetry property, we will only
 239 need to consider the positive $M = N/2$ bins. We denote $\mathbf{v}_1(t)$ as the visual
 240 feature that we have extracted from the recorded video signal associated
 241 with the target speaker. We also generate a complementary variable $Y_1^\dagger(f, t)$
 242 spanning the same frequency and time-frame space as $Y_1(f, t)$. The proposed
 243 sorting scheme for the alignment of the permutation is summarized in the
 244 following table, and also shown in Figure 5.

245 **Input:** spectral components $\mathbf{Y}(f, t)$, separation matrix $\mathbf{W}(f)$, GMM pa-
 246 rameter set $\{\lambda_i\}$, and visual feature $\mathbf{v}_1(t)$.

247 **Output:** aligned $\mathbf{W}(f)$ and $\mathbf{Y}(f, t)$.

248 Initialize the degree of frequency division $I_{\max} = 5$.

249 **For each** $i = 1, 2, \dots, I_{\max}$, **do**

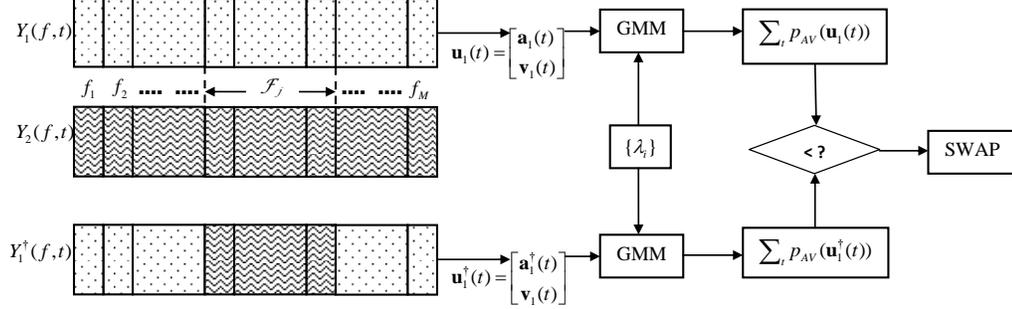


Figure 5: Diagram of the sorting scheme. The j -th loop in the i -th iteration. Spectral components of $Y_1^\dagger(f, t)$ come from $Y_1(f, t)$ (denoted by the dots) except those in the band \mathcal{F}_j which are from $Y_2(f, t)$ (denoted by the curves). We compare the coherence of $Y_1(f, t)$ and $Y_1^\dagger(f, t)$ to determine whether to swap components of $Y_1(f, t)$ and $Y_2(f, t)$ in \mathcal{F}_j .

- 250 Equally divide M bins into 2^{i-1} parts.
- 251 **For each** $j = 1, \dots, 2^{i-1}$, **do**
- 252 1) Let the j -th frequency band \mathcal{F}_j span $\{f_{\frac{M}{2^{i-1}}(j-1)+1}, \dots, f_{\frac{M}{2^{i-1}}j}\}$. For
- 253 $f \in \mathcal{F}_j$, let $Y_1^\dagger(f, \cdot) = Y_2(f, \cdot)$; otherwise, $Y_1^\dagger(f, \cdot) = Y_1(f, \cdot)$.
- 254 2) Extract the audio feature $\mathbf{a}_1(t)$ and $\mathbf{a}_1^\dagger(t)$ from $Y_1(\cdot, t)$ and $Y_1^\dagger(\cdot, t)$,
- 255 respectively. Let $\mathbf{u}_1(t) = [\mathbf{a}_1(t); \mathbf{v}_1(t)]$, $\mathbf{u}_1^\dagger(t) = [\mathbf{a}_1^\dagger(t); \mathbf{v}_1(t)]$.
- 256 3) Calculate the audio-visual probability $p_{AV}(\mathbf{u}_1(t))$ and $p_{AV}(\mathbf{u}_1^\dagger(t))$ re-
- 257 spectively, based on the GMM model in equation (10) and the pa-
- 258 rameter set $\{\lambda_i\}$ that has been estimated by the AEM algorithm.
- 259 4) If $\sum_t p_{AV}(\mathbf{u}_1(t)) > \sum_t p_{AV}(\mathbf{u}_1^\dagger(t))$, do nothing; otherwise, swap the
- 260 rows of $\mathbf{W}(f)$ and $\mathbf{Y}(f, \cdot)$ for $f \in \mathcal{F}_j$.
- 261 **End** j
- 262 **End** i
-

263 This scheme can reach a high resolution, which is determined by the num-
264 ber of partitions $2^{I_{\max}-1}$ at the final division, and the larger the number I_{\max} ,
265 the higher the resolution. However, most permutations occur contiguously
266 in practical situations, therefore even if we stop running the algorithm at a
267 very ‘coarse’ resolution (corresponding to a small I_{\max}), the permutation
268 ambiguity can still be substantially reduced (e.g. stop the iteration when the
269 algorithm has divided the positive frequency bins into 16 parts, i.e. $I_{\max} = 5$
270 and a frequency band of 500Hz).

271 5. Experimental Results

272 5.1. Data, Parameter Setup and Performance Metrics

273 Very similar to the database used in [3, 12], the corpus⁴ used in our
274 research contains sequences of “V1-C-V2”, where “V1” and “V2” are vowels
275 from /a/, /i/, /o/, /u/, and “C” stands for the consonant from /p/, /t/,
276 /k/, /b/, /d/, /g/ or no plosive (in the case of no plosive, the sequences are
277 “V1-V2”). There are 112 combinations recorded twice, one for training and
278 another for testing. The audio sequences are sampled at 16 kHz in mono, 16
279 bit PCM wave files, while the video sampling rate is 50 Hz and the associated
280 visual features are extracted by a chrome based system with 2496 frames for
281 training and another 2547 frames for testing.

282 In our experiment, the audio data were obtained by concatenating in-
283 dependent sequences, with each sequence lasting an integer multiple of 20
284 ms. We concatenated the 112 isolated sequences to obtain approximately

⁴Thanks to Bertrand Rivet in GIPSA-Lab for providing us with this multimodal database.

285 50 seconds audio for training. In the same way, we chose the beginning 400
 286 frames (approximately 8 s) of the testing data to demonstrate our algorithm.
 287 128 ms sliding Hamming window (2048 taps) with 108 ms overlap (20 ms
 288 step size, to be synchronised with the visual features) was applied in STFT.
 289 5-dimensional ($L = 5$) MFCCs as audio features were computed from 24
 290 Mel-scaled filter banks, thus the audio-visual feature was 7-dimensional. We
 291 set the FFT size as 2048 (i.e., $M=1024$). In the frame selection process,
 292 we reserved 62% of features by assigning $\delta_{LW} = 0.35$, $\delta_{LH} = 0.5$, $Q = 2$,
 293 and $\alpha_{LW} = \alpha_{LH} = 1$. For simplicity, we only used 20 ($I = 20$) kernels to
 294 approximate the audio-visual coherence.

295 The algorithm was tested on convolutive mixtures synthesized on com-
 296 puter. We used the real room recordings from the binaural room impulse
 297 response database (i.e. AIR database) [35] for the mixing filters. The mea-
 298 surements were recorded with or without dummy head in a low-reverberant
 299 studio booth, an office room, a meeting room and a lecture room respectively,
 300 of which we chose the meeting room scenario⁵. For a 2×2 mixing system, we
 301 chose two positions for two source signals, and used the corresponding room
 302 impulse responses as mixing filters. As a result, $C_5^2 = 10$ combinations of
 303 mixing filters were used in the following performance evaluations. The two
 304 source signals, one was the previously mentioned 8-second truncated segment

⁵In the meeting room scenario, a dummy head was placed on a fixed position and the room impulse responses were captured for five different positions opposite of the head. Each room impulse has 10923 taps (approximately 700 ms) where sampling frequency is 16 kHz, and the reverberation time (T_{60}) is about 300 ms.

305 from the test audio, and another source signal ⁶ was continuous speech from
 306 the XM2VTS database [36]. Gaussian white noise was added to both mix-
 307 tures at different signal to noise ratios (SNRs). Figure 6(a) shows the source
 308 signals used in our experiments (note that only 2 s are shown here).

In the frequency domain we use the global filters and signal to interference and noise ratio (SINR) as criteria to evaluate the performance of our bimodal BSS algorithm at different signal to noise ratios (SNRs). Suppose $s_1(n)$ is the target source, then in the 2×2 case:

$$\mathbf{G}(f) = \begin{bmatrix} \mathbf{G}_{11}(f) & \mathbf{G}_{12}(f) \\ \mathbf{G}_{21}(f) & \mathbf{G}_{22}(f) \end{bmatrix} = \mathbf{W}(f)\mathbf{H}(f), \quad (17)$$

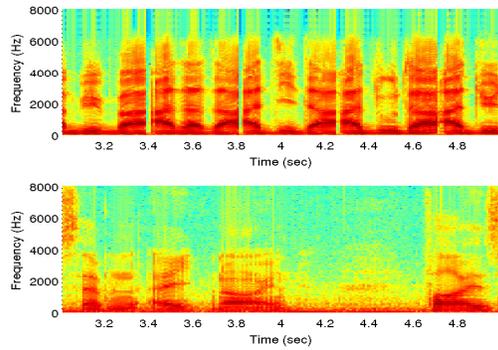
$$\text{SINR} = 10 \log \frac{P_{s_1}}{P_{\hat{s}_1 - s_1}} = 10 \log \frac{\sum_n \left\| \sum_{p=1}^P w_{1p} * h_{p1} * s_1(n) \right\|^2}{\sum_n \left\| \hat{s}_1(n) - \sum_{p=1}^P w_{1p} * h_{p1} * s_1(n) \right\|^2}. \quad (18)$$

309 5.2. Experimental Results

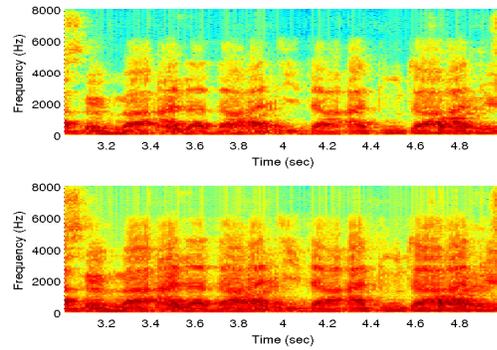
310 5.2.1. An example

311 In the first experiment, we show the effectiveness of our algorithm by
 312 comparing it with the FD-BSS method without any permutation alignment
 313 (denoted as “No Alignment”). For the FD-BSS, **the joint approximate**
 314 **diagolization method proposed in [30] is used for source separation**
 315 **at each frequency bin.** Note that the same ICA algorithm was used here
 316 for the contrast methods [30] and [33] in our comparisons in Section 5.2.2.
 317 **In addition, the scaling problem is addressed based on the minimal**
 318 **distortion principle [29] for all these methods before performing the**
 319 **permutation alignment.** We used the same parameter set-up as described

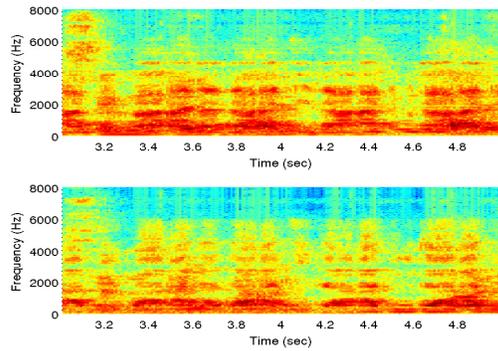
⁶The source signals can also be chosen from a same dataset as done in our work [15].



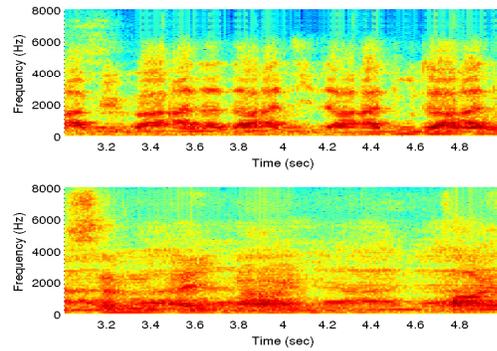
(a) Source signals



(b) Mixtures



(c) Estimates (no alignment)

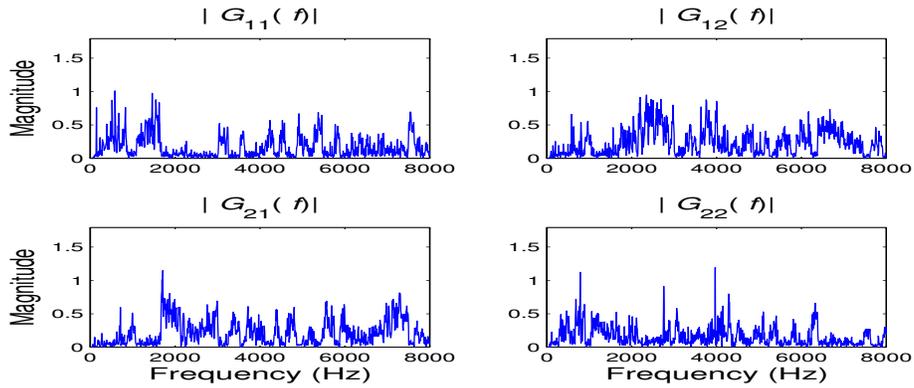


(d) Estimates (audio-visual)

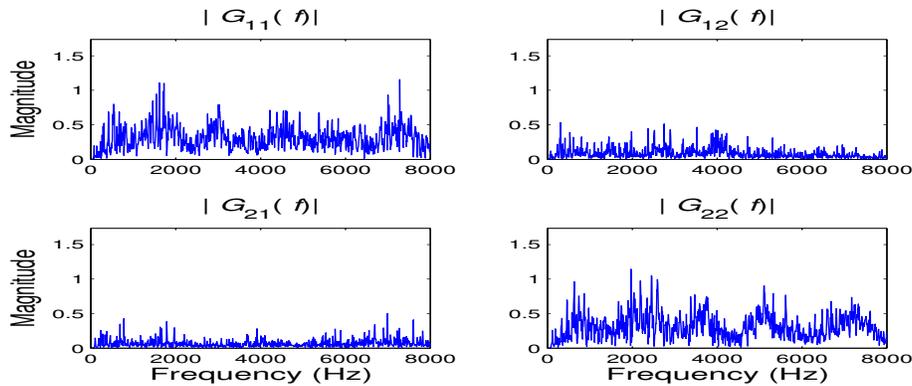
Figure 6: Spectrograms of the sources (a), the mixtures (b) and the estimated sources without permutation alignment (c) where $\text{SINR}=0.48$ dB and the estimated sources by the proposed algorithm (d) where $\text{SINR}=7.91$ dB. The upper sub-plot in (a) is the target signal.

320 in Section 5.1. The room impulse responses obtained from the AIR database
 321 are used as mixing filters, where the sources are placed respectively on the
 322 first and fourth position of the meeting room. Each filter has 10923 taps,
 323 and the reverberation time (T_{60}) of the meeting room is about 300 ms. We
 324 set $N = 2048$ and $I_{max} = 5$. No noise is added to the mixtures. The
 325 generated mixtures are shown in Figure 6(b). Figure 6(c) and 6(d) show
 326 the separated sources from the mixtures without and with our proposed
 327 algorithm respectively, where it can be clearly observed that the audio-visual
 328 approach has improved the quality of the separated speech.

329 Figure 7 shows typical global filters obtained from this experiment. Ide-
 330 ally the global filter should be an identity matrix for each frequency bin.
 331 From this figure, it can also be seen that the permutation problem is well ad-
 332 dressed by the audio-visual approach, as the magnitudes of G_{12} and G_{21} have
 333 been reduced considerably. Alternatively, based on the diagonal property of
 334 the global filter matrix, we can also use the criterion $|G_{11}| \cdot |G_{22}| - |G_{12}| \cdot |G_{21}|$
 335 to evaluate the consistency of the permutations across the frequency bins,
 336 as shown in Figure 8. **By comparing Figure 8(d) with Figure 8(a),**
 337 **it can be observed that our algorithm successfully corrected the**
 338 **permutation ambiguities at most frequency bins, as the values of**
 339 **$|G_{11}| \cdot |G_{22}| - |G_{12}| \cdot |G_{21}|$ are consistently positive across the frequency**
 340 **bins. In addition, our algorithm performs better than the two base-**
 341 **line methods shown in Figures 8(b) and 8(c). To observe the effect**
 342 **of noise on the permutation errors, we have also plotted the quan-**
 343 **tity $|G_{11}| \cdot |G_{22}| - |G_{12}| \cdot |G_{21}|$ for the case where 10 dB white Gaussian**
 344 **noise was added to the mixtures, as shown in Figure 9, both the**



(a) No alignment



(b) Audio-visual

Figure 7: Global filters in the frequency domain obtained by the FD-BSS without permutation alignment (a) and the audio-visual alignment (b). The same mixtures and parameters as those in Figure 6 were used. The permutation ambiguities occur in many frequency bands in (a). Therefore, the spectrum of the recovered signal $y_1(t)$ contains distortions from the other source signal $s_2(t)$, which can also be observed in Figure 6(c). With audio-visual alignment, the permutation ambiguities have been significantly reduced, as shown in (b) as well as Figure 6(d).

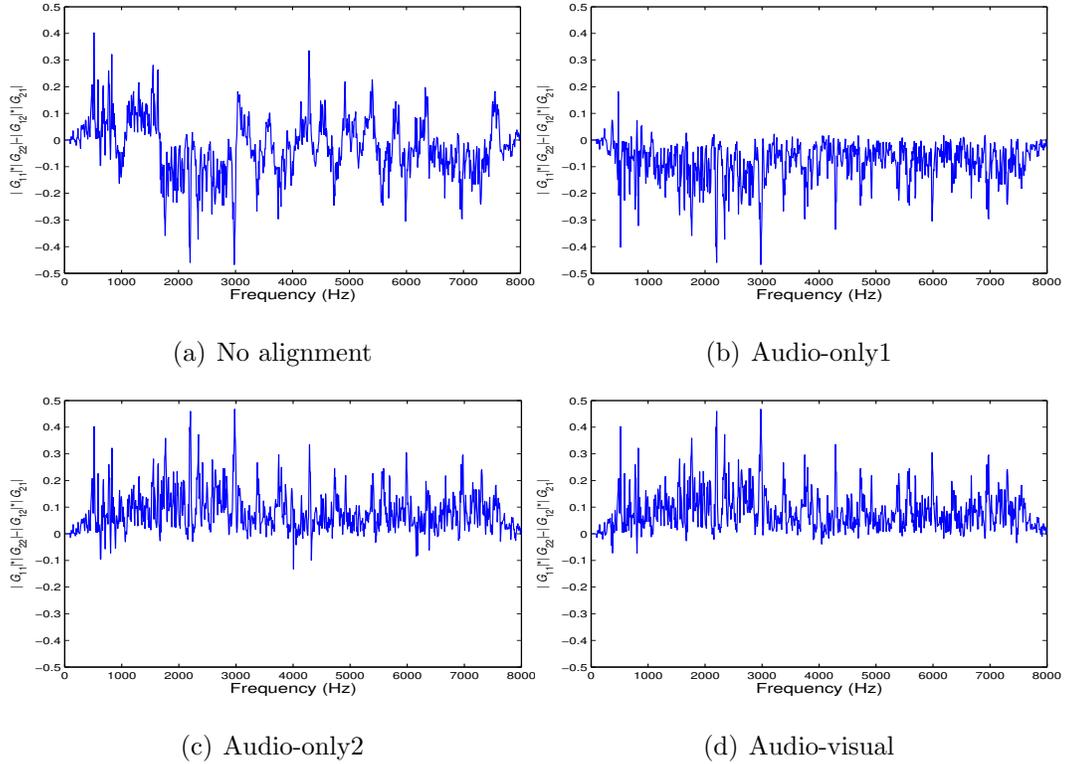


Figure 8: **The plot of $|G_{11}| \cdot |G_{22}| - |G_{12}| \cdot |G_{21}|$ for the noise-free case, when no permutation alignment is conducted (a) and the permutations are aligned by Audio-only1 (b) in [30], Audio-only2 (c) in [33] and the audio-visual approach (d). $|G_{11}| \cdot |G_{22}| - |G_{12}| \cdot |G_{21}|$ should be consistently larger or smaller than 0, if the permutations across the frequency bins are correctly aligned.**

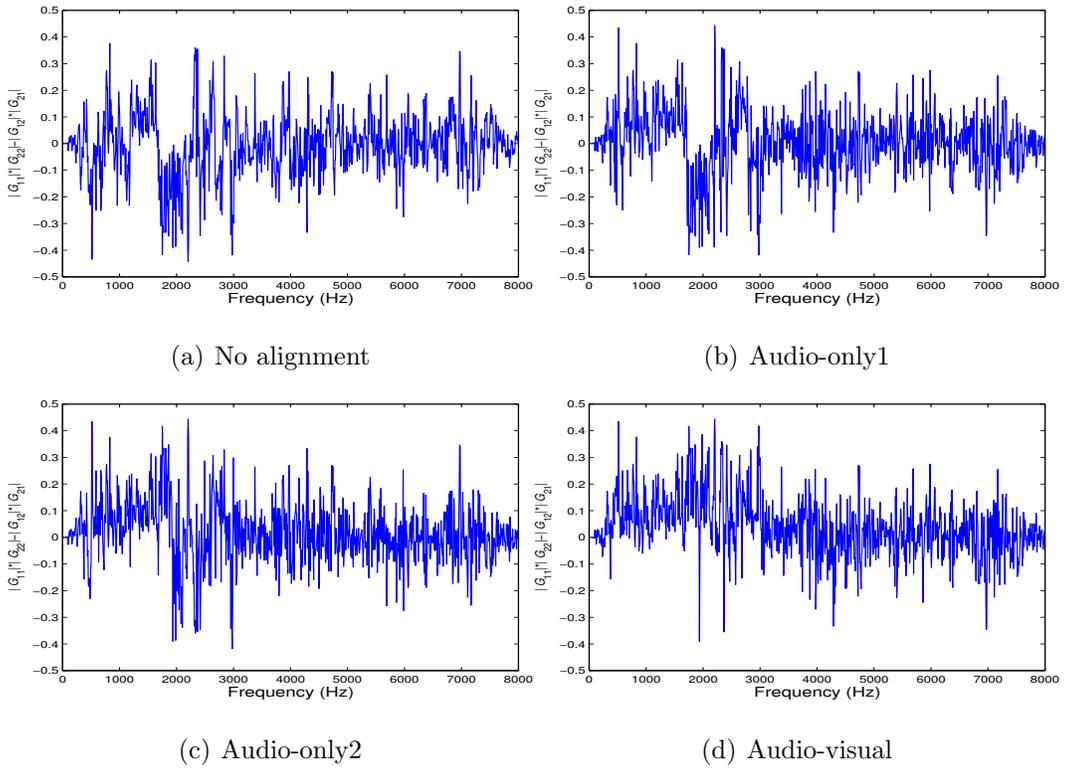


Figure 9: The plot of $|G_{11}| \cdot |G_{22}| - |G_{12}| \cdot |G_{21}|$ for the case where 10 dB Gaussian white noise was added to the mixtures. The meanings of the sub-plots correspond similarly to those in the noise-free case shown in Figure 8.

345 audio-only methods failed to group the spectral components accu-
346 rately in the frequency band between approximately 1700 Hz and
347 3000 Hz, which lies in the range of frequencies that are essential to
348 human intelligibility of speech. The audio-visual method provides
349 more accurate alignment in this case. The results and comparisons
350 in terms of SINR measurements are provided in the next section.

351 The magnitude spectra in Figures 7, 8, and 9 appear somewhat
352 peaky, which implies that a certain amount of scale ambiguity re-
353 mains despite having applied the MDP-based post-processing to
354 the unmixing filters. However, our informal listening tests show
355 that the effect of the peakedness of the frequency spectrum on the
356 perceptual quality of the recovered speech signals is negligible. In
357 other words, we didn't observe strong distortions of the signal as
358 a result of the residual spectral peaks.

359 5.2.2. Comparisons

360 First, we compare the effect of different frequency partitions (i.e. the
361 choice of I_{max}) on the performance of the proposed method for permutation
362 alignment. To this end, we used the same experimental set-up as described
363 in Sections 5.1 and 5.2.1, except the following two changes. First, we change
364 the values of I_{max} by selecting it from [3 4 5 6]. Second, we perform 10 exper-
365 iments in each of which the mixing filters were selected from the 10 sets of the
366 room impulse responses described in Section 5.1. The average results over
367 the 10 different sets of mixtures are shown in Table 1. From this table, it can
368 be observed that the different number of frequency partitions does influence
369 the separation performance. In general, $I_{max} = 4$ or 5, which corresponds

Table 1: The effect of the choice of I_{max} on the separation performance measured by SINR in dB

I_{max}	3	4	5	6
SINR(dB)	6.39	8.98	8.21	6.40

370 **to 8 or 16 partitions respectively**, gives reasonably good performance.
 371 More partitions however do not increase the system performance. We choose
 372 $I_{max} = 5$ for subsequent performance comparisons.

373 We then compare our algorithm with other de-permutation methods us-
 374 ing only audio signals. We use the method in [30] (denoted as “Audio-
 375 only1”) and the approach in [33] (denoted as “Audio-only2”) as contrast
 376 algorithms for permutation alignment. Audio-only1 integrates information
 377 across different frequencies with the assumption that signal profiles in differ-
 378 ent bins undergo interrelated changes, even for distant frequency channels.
 379 Audio-only2 exploits the cross-frequency correlation between neighbouring
 380 frequency bands based on a hierarchical structure.

381 We evaluate the performance of our algorithm and the above baseline al-
 382 gorithms with respect to different signal to noise ratios (SNRs) and different
 383 FFT sizes. First, we fix the parameters as used in above sections, and only
 384 change the values of FFT size among [512 1024 2048 4096]. For each FFT
 385 size N , ten independent tests were run on the ten different sets of mixtures
 386 for each of the algorithms under comparison. No noise was added to these
 387 mixtures. The average SINR (dB) results are shown in Table 2. From this
 388 table, it can be observed that the choice of N has impacts on the separation
 389 performance. For example, a smaller N , such as $N = 512$, gives relatively
 390 lower performance for almost all the tested algorithms. Our proposed algo-

391 rithm offers competitive performance in all the test cases. $N = 4096$ appears
 392 to offer better performance for the majority of the algorithms. However, a
 393 higher computational cost is usually involved for a larger N when performing
 permutation alignment.

Table 2: SINR measurements for different FFT size (i.e. N).

N	512	1024	2048	4096
No alignment	2.19	1.85	3.53	3.87
Audio-only1	4.23	6.36	6.12	7.55
Audio-only2	2.38	2.76	8.69	8.88
Audio-visual	4.15	4.56	8.21	8.98

394
 395 Then, we changed the SNR level from 5 dB to 30 dB while maintaining
 396 $N = 2048$, $I_{max} = 5$, and other parameters as set in the above two sections.
 397 For each SNR, again ten independent sets of tests were run for each of the
 398 algorithms under comparison. The average results measured by SINR are
 399 shown in Figure 10. It is shown from this figure that the proposed algorithm
 400 tends to offer better performance for a higher noise level. For example, when
 401 $\text{SNR} = 10$ dB, our algorithm provides 0.73 dB improvement over Audio-
 402 only1 and 1.75 dB over Audio-only2. For the noise-free case (not plotted on
 403 this figure), our algorithm gains 2.08 dB improvement over Audio-only1, but
 404 is 0.48 dB lower than Audio-only2. This suggests that the advantage of the
 405 audio-visual method is more prominent for noisy mixtures than the noise-free
 406 ones.

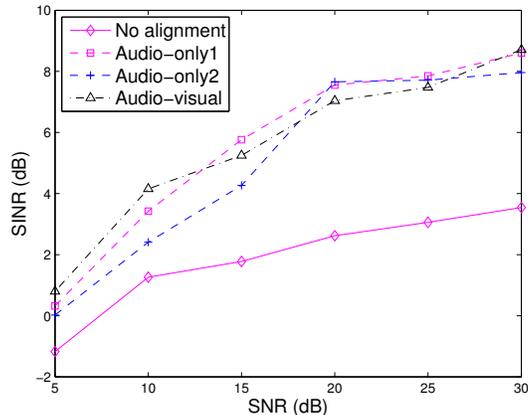


Figure 10: Average SINR computed over the independent tests for the 4 algorithms, where $N = 2048$ and $I_{max} = 5$.

407 6. Conclusions

408 We have presented a new audio-visual convolutive BSS system. In this
 409 system, we have used the MFCCs as audio features, which were combined
 410 with geometric visual features to form an audio-visual feature space. We
 411 have considered using some dynamic features in video for robust feature
 412 selection in audio-visual space. An adapted EM algorithm has been proposed
 413 by exploiting the different influences of the audio features for statistically
 414 modelling the audio-visual coherence. A new recursive sorting scheme based
 415 on the maximization of the audio-visual coherence has also been developed
 416 to solve the permutation problem.

417 Acknowledgement

418 This work was supported by the Engineering and Physical Sciences Re-
 419 search Council (EPSRC) (Grant number EP/H012842/1) and the MOD Uni-

420 versity Defence Research Centre on Signal Processing (UDRC). **The au-**
421 **thors would like to thank the anonymous reviewers and the guest**
422 **editors for their insightful comments that considerably improved**
423 **the quality of this paper.**

424 **References**

- 425 [1] D. A. Bulkin, J. M. Groh, Seeing sounds: visual and auditory interac-
426 tions in the brain, *IEEE Transactions on Neural Networks* 16 (4) (2006)
427 415–419.
- 428 [2] J.-L. Schwartz, F. Berthommier, C. Savariaux, Seeing to hear better:
429 evidence for early audio-visual interactions in speech identification, *Cog-*
430 *nition* 93 (2004) B69–B78.
- 431 [3] D. Sodoyer, J.-L. Schwartz, L. Girin, J. Klinkisch, C. Jutten, Separation
432 of audio-visual speech sources: a new approach exploiting the audio-
433 visual coherence of speech stimuli, *EURASIP Journal on Applied Signal*
434 *Processing* 2002 (11) (2002) 1165–1173.
- 435 [4] G. Monaci, P. Vandergheynst, F. T. Sommer, Learning bimodal struc-
436 ture in audio-visual data, *IEEE Transactions on Neural Networks* 20 (12)
437 (2009) 1898–1910.
- 438 [5] A. L. Casanovas, G. Monaci, P. Vandergheynst, R. Gribonval, Blind au-
439 diovisual source separation based on sparse redundant representations,
440 *IEEE Transactions on Multimedia* 12 (5) (2010) 358–371.

- 441 [6] S. G. Mallat, Z. Zhang, Matching pursuits with time-frequency dictio-
442 naries, *IEEE Transactions on Signal Processing* 41 (12) (1993) 3397–
443 3415.
- 444 [7] D. Ormoneit, V. Tresp, Averaging, maximum penalized likelihood and
445 bayesian estimation for improving gaussian mixture probability density
446 estimates, *IEEE Transactions on Neural Networks* 9 (4) (1998) 639–650.
- 447 [8] G. Potamianos, C. Neti, G. Gravier, A. Garg, A. W. Senior, Recent
448 advances in the automatic recognition of audio-visual speech, in: *Proc.*
449 *IEEE*, 2003, pp. 1306–1326.
- 450 [9] M. Gurban, Multimodal speaker localization in a probabilistic frame-
451 work, in: *Proc. of EUSIPCO*, 2006.
- 452 [10] L. Girin, J.-L. Schwartz, G. Feng, Audio-visual enhancement of speech
453 in noise, *The Journal of the Acoustical Society of America* 109 (6) (2001)
454 3007–3020.
- 455 [11] Q. Liu, W. Wang, P. Jackson, Bimodal coherence based scale ambiguity
456 cancellation for target speech extraction and enhancement, in: *Proc.*
457 *Interspeech*, 2010, pp. 438–441.
- 458 [12] B. Rivet, L. Girin, C. Jutten, Mixing audiovisual speech processing and
459 blind source separation for the extraction of speech signals from con-
460 volutive mixtures, *IEEE Transactions on Audio, Speech and Language*
461 *Processing* 15 (1) (2007) 96–108.
- 462 [13] W. Wang, D. Cosker, Y. Hicks, S. Sanei, J. Chambers, Video assisted
463 speech source separation, in: *Proc. ICASSP*, 2005, pp. 425–428.

- 464 [14] Q. Liu, W. Wang, P. Jackson, Audio-visual convolutive blind source
465 separation, in: Proc. Sensor Signal Processing for Defence (SSPD), 2010.
- 466 [15] Q. Liu, W. Wang, P. Jackson, Use of bimodal coherence to resolve spec-
467 tral indeterminacy in convolutive BSS, in: Proc. LVA/ICA, 2010, pp.
468 131–139.
- 469 [16] S. M. Naqvi, M. Yu, J. A. Chambers, A multimodal approach for blind
470 source separation of moving sources, IEEE Journal Selected Topics in
471 Signal Processing 4 (5) (2010) 895–910.
- 472 [17] P. Comon, Independent component analysis, a new concept?, Signal
473 Processing 36 (3) (1994) 287–314.
- 474 [18] C. Jutten, J. Herault, Blind separation of sources, part i: An adaptive
475 algorithm based on neuromimetic architecture, Signal Processing 24 (1)
476 (1991) 1–10.
- 477 [19] J.-F. Cardoso, A. Souloumiac, Blind beamforming for non-gaussian sig-
478 nals, Radar and Signal Processing, IEE Proceedings F 140 (6) (1993)
479 362–370.
- 480 [20] A. J. Bell, T. J. Sejnowski, An information-maximization approach to
481 blind separation and blind deconvolution, Neural Computation 7 (6)
482 (1995) 1129–1159.
- 483 [21] A. Belouchrani, K. Abed-Meraim, J.-F. Cardoso, E. Moulines, A blind
484 source separation technique using second-order statistics, IEEE Trans-
485 actions on Signal Processing 45 (2) (1997) 434–444.

- 486 [22] A. Hyvärinen, E. Oja, A fast fixed-point algorithm for independent com-
487 ponent analysis, *Neural Comput.* 9 (1997) 1483–1492.
- 488 [23] S. Amari, S. C. Douglas, A. Cichocki, H. H. Yang, Multichannel blind
489 deconvolution and equalization using the natural gradient, in: *Proc.*
490 *IEEE Int. Workshop on Wireless Commun.*, 1997, pp. 101–104.
- 491 [24] J. Thomas, Y. Deville, S. Hosseini, Time-domain fast fixed-point algo-
492 rithms for convolutive ICA, *IEEE Signal Process. Lett.* 13 (4) (2006)
493 228–231.
- 494 [25] T. Mei, J. Xi, F. Yin, A. Mertins, J. F. Chicharo, Blind source sepa-
495 ration based on time-domain optimization of a frequency-domain inde-
496 pendence criterion, *IEEE Transactions on Audio, Speech, and Language*
497 *Processing* 14 (6) (2006) 2075–2085.
- 498 [26] J. Anemüller, B. Kollmeier, Amplitude modulation decorrelation for
499 convolutive blind source separation, in: *Proc. ICA*, 2000, pp. 215–220.
- 500 [27] M. Z. Ikram, D. R. Morgan, A beamforming approach to permutation
501 alignment for multichannel frequency-domain blind speech separation,
502 in: *Proc. ICASSP*, 2002, pp. 881–884.
- 503 [28] F. Nesta, M. Omologo, P. Svaizer, A novel robust solution to the permu-
504 tation problem based on a joint multiple TDOA estimation, in: *Proc.*
505 *IWAENC*, 2008.
- 506 [29] K. Matsuoka, Minimal distortion principle for blind source separation,
507 in: *Proc. SICE*, Vol. 4, 2002, pp. 2138–2143.

- 508 [30] D.-T. Pham, C. Servière, H. Boumaraf, Blind separation of speech mix-
509 tures based on nonstationarity, in: Proc. ISSPA, Vol. 2, 2003, pp. 73–76.
- 510 [31] H. Sawada, R. Mukai, S. Araki, S. Makino, A robust and precise method
511 for solving the permutation problem of frequency-domain blind source
512 separation, *IEEE Transactions on Speech and Audio Processing* 12 (5)
513 (2004) 1063–6676.
- 514 [32] R. Mazur, A. Mertins, An approach for solving the permutation problem
515 of convolutive blind source separation based on statistical signal models,
516 *IEEE Transactions on Audio, Speech and Language Processing* 17 (1)
517 (2009) 117–126.
- 518 [33] K. Rahbar, J. P. Peilly, A frequency domain method for blind source
519 separation of convolutive audio mixtures, *IEEE Transactions on Speech
520 and Audio Processing* 13 (5) (2005) 832–844.
- 521 [34] D. Sodoyer, B. Rivet, L. Girin, C. Savariaux, J.-L. Schwartz, C. Jutten,
522 A study of lip movements during spontaneous dialog and its applica-
523 tion to voice activity detectio, *The Journal of the Acoustical Society of
524 America* 125 (2) (2009) 1184–1196.
- 525 [35] M. Jeub, M. Schafer, P. Vary, A binaural room impulse response
526 database for the evaluation of dereverberation algorithms, in: 16th In-
527 ternational Conference on Digital Signal Processing, 2009, data online:
528 <http://www.ind.rwth-aachen.de/AIR>.
- 529 [36] XM2VTS, Website, <http://www.ee.surrey.ac.uk/CVSSP/xm2vtsdb/>.