

# Multimodal Evolution Inference

## Temporal probabilistic models in multimodal perception

Xavier Alameda-Pineda

January 20, 2014

### Abstract

We derive an EM for a big graphical model.

## 1 Notation

### 1.1 Generics

**Scalar, vector and matrix** Italic means scalar, italic bold means vector and bold means matrix. For example, in the same order:

*A*   ***A***   **A**

**Collection of indexed variables** Whenever a variable has an index, such as  $A_n$  the following notations hold:

$A_{1:N}$  is the set of  $A_n$  for all values of  $n$  from 1 to  $N$ .

$A_{1:m:N}$  is the set of  $A_n$  for all values of  $n$  from 1 to  $N$  except from  $m$ .

**Vector and matrix components** Whenever  $\mathbf{A}$  is a vector  $\mathbf{A}[n]$  denotes the  $n$ -th component of the vector. The same intuition holds for  $\mathbf{A}[i, j]$  when  $\mathbf{A}$  is a matrix.

**Important sets** We need to consider different sets of matrices and vectors to simplify the notations:

- $\mathcal{CM}^D \subset \mathbb{R}^{D \times D}$  is the set of covariance matrices of dimension  $D$ . That is  $\mathbf{x} \in \mathcal{CM}^D$  if and only if  $\mathbf{x}$  is symmetric and positive definite.
- $\mathcal{SV}^D \subset \mathbb{R}^D$  is the set of stochastic vectors of dimension  $D$ . That is  $\mathbf{x} \in \mathcal{SV}^D$  if and only if  $\mathbf{x}[d] \geq 0 \forall d = 1, \dots, D$  and  $\sum_{d=1}^D \mathbf{x}[d] = 1$ .
- $\mathcal{SM}^{D \times K} \subset \mathbb{R}^{D \times K}$  is the set of row-stochastic matrices of dimension  $D \times K$ . That is  $\mathbf{x} \in \mathcal{SM}^{D \times K}$  if and only if  $\mathbf{x}[d, k] \geq 0 \forall d = 1, \dots, D, k = 1, \dots, K$  and  $\sum_{k=1}^K \mathbf{x}[d, k] = 1, \forall d = 1, \dots, D$ .

Finally, we will use  $\mathcal{D}(\cdot)$  to denote the dimensionality of the model.

### 1.2 In the model

$N$  Number of sources.

#### 1.2.1 Source Position

- $D$  dimension of the source space,  $\mathbb{R}^D$ .
- $\mathbf{S}_{t,n}$  position of the  $n^{\text{th}}$  source at time  $t$ .  $\mathbf{S}_t$  concatenation of the  $N$  position vectors.
- $\boldsymbol{\nu}_n, \boldsymbol{\Omega}_n$ , parameters of the Gaussian for  $\mathbf{S}_{1,n}$ ,  $\boldsymbol{\nu}_n \in \mathbb{R}^D$  and  $\boldsymbol{\Omega}_n \in \mathcal{CM}^D$ .
- $\boldsymbol{\Lambda}_{n,0}, \boldsymbol{\Lambda}_{n,1}$ , parameters of the diffusion dynamics (not moving and moving).  $\boldsymbol{\Lambda}_{n,0}, \boldsymbol{\Lambda}_{n,1} \in \mathcal{CM}^D$ .
- $\theta^S = \{\boldsymbol{\nu}_n, \boldsymbol{\Omega}_n, \boldsymbol{\Lambda}_{n,0}, \boldsymbol{\Lambda}_{n,1}\}_{n=1}^N$ .

#### 1.2.2 Auditory status

- $\mathbf{A}_t \in \{0, 1\}^N$  (or  $\mathbf{A}_t \in \{1, \dots, 2^N\}$ ), auditory status of the  $N$  sources at time  $t$ .
- $\mathbf{T}^A, \boldsymbol{\lambda}^A$ , transition matrix and prior vector for  $\mathbf{A}_t$ .
- $\mathbf{T}^A \in \mathcal{SM}^{2^N \times 2^N}$  and  $\boldsymbol{\lambda}^A \in \mathcal{SV}^{2^N}$ .
- $\theta^A = \{\mathbf{T}^A, \boldsymbol{\lambda}^A\}$ .

#### 1.2.3 Visibility and Motion status

The follow the same logic, changing  $F$  by  $M$ .

- $F_{t,n}$  is the visibility of the  $n^{\text{th}}$  source at time  $t$ .
- $\mathbf{T}_n^F, \boldsymbol{\lambda}_n^F$ , transition matrix and prior vector for  $F_{t,n}$ .
- $\mathbf{T}_n^F \in \mathcal{SM}^{2 \times 2}$  and  $\boldsymbol{\lambda}_n^F \in \mathcal{SV}^2$ .
- $\theta^F = \{\mathbf{T}_n^F, \boldsymbol{\lambda}_n^F\}_{n=1}^N$ .

#### 1.2.4 Mel Frequency Cepstral Coefficients [MFCC]

$D_V$  dimension of the MFCC,  $\mathbf{V} \in \mathbb{R}^{D_V}$ .

- $G_V$  number of components of the GMM.
- $K_t^V$  number of MFCC observations at time  $t$ .
- $Z_{t,k}^V \in \{1, \dots, G_V\}$  observation-to-component assignment variable of the  $k^{\text{th}}$  observation at time  $t$ .
- $\mathbf{V}_{t,k}$   $k^{\text{th}}$  MFCC observation at time  $t$ .
- There is one GMM per auditory status  $1, \dots, 2^N$ .

- $\boldsymbol{\pi}_j^V$  prior vector of the mixing coefficients,  $j = 1, \dots, 2^N$ .
- $\boldsymbol{\mu}_{j,g}^V, \boldsymbol{\Sigma}_{j,g}^V$  mean and covariance matrix of the  $g^{\text{th}}$  component of the  $j^{\text{th}}$  GMM,  $g = 1, \dots, G_V, j = 1, \dots, 2^N$ .
- $\theta^{ZV} = \{\boldsymbol{\pi}_1^V, \dots, \boldsymbol{\pi}_{2^N}^V\}, \boldsymbol{\pi}_j^V \in \mathcal{SV}^{G_V}$ .
- $\theta^V = \{\boldsymbol{\mu}_{j,g}^V, \boldsymbol{\Sigma}_{j,g}^V\}_{g=1, j=1}^{G_V, 2^N}, \boldsymbol{\mu}_{j,g}^V \in \mathbb{R}^{G_V}, \boldsymbol{\Sigma}_{j,g}^V \in \mathcal{CM}^{D_V}$ .

### 1.2.5 Spectrogram [SPEC]

- $\mathbf{P}^L, \mathbf{P}^R$  left and right spectrograms.  $\mathbf{P}^L, \mathbf{P}^R \in \mathbb{C}^{P_F \times P_I}$ .
- $G_P$  number of components generating the spectrograms.
- $\mathbf{C}_{t,g}$ ,  $g^{\text{th}}$  component at time  $t$ .
- $\boldsymbol{\omega}_g, \boldsymbol{\xi}_g$ , spectral pattern and temporal activation pattern of the  $g^{\text{th}}$  component.
- $\boldsymbol{\Xi}^L, \boldsymbol{\Xi}^R$  left and right mixing matrices.  $\boldsymbol{\Xi}^L, \boldsymbol{\Xi}^R \in \mathbb{C}^{P_F \times G_P}$ .
- $\boldsymbol{\Sigma}_f^{PL}, \boldsymbol{\Sigma}_f^{PR}$  noise variances of the  $f^{\text{th}}$  frequency bin (L&R).
- $\theta^C = \{\boldsymbol{\omega}_g, \boldsymbol{\xi}_g\}_{g=1}^{G_V}$ .
- $\theta^P = \{\boldsymbol{\Xi}^L, \boldsymbol{\Xi}^R, \{\boldsymbol{\Sigma}_f^{PL}, \boldsymbol{\Sigma}_f^{PR}\}_{f=1}^{P_F}\}$ .

### 1.2.6 Binaural SSL [BINAURAL] [MOTION]

Motion and binaural features follow the very same model, replacing  $B$  by  $Y$ , and conditioning with respect to  $\mathbf{M}_t$  instead of  $\mathbf{A}_t$ .

- $K_t^B$  number of binaural observations at time  $t$ .
- $\mathbf{B}_{t,k}$   $k^{\text{th}}$  binaural observation at time  $t$ .
- $\mathbf{B}_{t,k}$  follows a  $N$ -component GMM.
- $\mathbf{W}_t^B$  mixing priors, following a Dirichlet.
- $\boldsymbol{\alpha}_n^B$ , parameters of the  $n^{\text{th}}$  component of the Dirichlet.
- $Z_{t,k}^B$  observation-to-component (source).
- $\boldsymbol{\Sigma}_{t,n}^B$ , covariance matrix of the  $n^{\text{th}}$ ,  $\boldsymbol{\Sigma}_{t,n}^B \in \mathcal{CM}^D$ .
- $\theta^{WB} = \{\boldsymbol{\alpha}_1^B, \dots, \boldsymbol{\alpha}_N^B\}$ .
- $\theta^B = \{\boldsymbol{\Sigma}_{t,n}^B\}_{t=1, n=1}^{T, N}$ .

### 1.2.7 [FACIAL]

- $K_t^F$  number of detected faces at time  $t$ .
- $e$  face detector false positive rate.
- $\mathbf{Q}_{t,k} \in \mathbb{R}^D$   $k^{\text{th}}$  detection position at time  $t$ .
- $\boldsymbol{\Sigma}^Q$  covariance matrix (face detector localization error).
- $Z_{t,k}^R$  pose detection (among  $G_R$  poses).
- $G_F$  number of facial landmarks.
- $\mathbf{R}_{t,k}^G \in \mathbb{R}^{G_F D}$  positions of the landmarks.
- $\boldsymbol{\Sigma}_{n,m}^G$  covariance of the Gaussian describing  $\mathbf{R}^G$  for the  $m^{\text{th}}$  pose of the  $n^{\text{th}}$  source,  $\boldsymbol{\Sigma}_{n,m}^G \in \mathcal{CM}^{G_F D} \forall n, m$
- $D_A$  dimension of the landmark descriptors.
- $\mathbf{R}_{t,k,g}^A \in \mathbb{R}^{D_A}$  descriptor of the  $g^{\text{th}}$  landmark.
- $\boldsymbol{\mu}_{n,m,g}^A, \boldsymbol{\Sigma}_{n,m,g}^A$  mean and covariance of the Gaussian describing the  $g^{\text{th}}$  landmark of the  $m^{\text{th}}$  pose of the  $n^{\text{th}}$  source.
- $G_L$  number of lip landmarks.
- $\mathbf{L}_{t,k,g} \in \mathbb{R}^{D_L}$  lip movement descriptor of the  $g^{\text{th}}$  landmark.
- $\boldsymbol{\mu}_{n,m,g}^L, \boldsymbol{\Sigma}_{n,m,g}^L$  mean and covariance of the Gaussian describing the  $g^{\text{th}}$  lip landmark of the  $m^{\text{th}}$  pose of the  $n^{\text{th}}$  source.

## 2 The Model

We present in this section the graphical model, that is the probability description of the relationship between the observed variables and the hidden variables and how this relationship is parametrized. The full model is shown in Figure 1. In the following, we describe the variables and its relations little by little.

We assume the existence of  $N$  sources (potential speakers).

### 2.1 Dynamics

#### 2.1.1 Positions

The  $N$  sources are placed in  $\mathbb{R}^D$ .  $\mathbf{S}_{t,n} \in \mathbb{R}^D$  denotes the position of the  $n^{\text{th}}$  source at time  $t$ , and  $\mathbf{S}_t$  is the concatenation of the  $N$  vectors, so  $\mathbf{S}_t \in \mathbb{R}^{ND}$ . The dynamics of the positions are modelled separately:

$$\mathbf{S}_{t,n} | \mathbf{S}_{t-1,n}, M_{t-1,n} \sim \mathcal{N}(\mathbf{S}_{t,n}; \mathbf{S}_{t-1,n}, \boldsymbol{\Lambda}_{n, M_{t-1,n}}) \quad (2.1)$$

$$\mathbf{S}_{1,n} \sim \mathcal{N}(\mathbf{S}_{1,n}; \boldsymbol{\nu}_n, \boldsymbol{\Omega}_n). \quad (2.2)$$

This means that for each source there are two possible dynamics, either moving  $M_{t-1,n} = 1$  or not moving  $M_{t-1,n} = 0$  (see Section 2.1.4). Therefore, we consider two transition matrices per source: dynamic  $\boldsymbol{\Lambda}_{n,1}$  and static  $\boldsymbol{\Lambda}_{n,0}$ . This holds for all time steps, except for the first one, in which is modelled as a multivariate Gaussian with mean  $\boldsymbol{\nu}_n$  and covariance  $\boldsymbol{\Omega}_n$ .

The parameters are:  $\theta^S = \{\boldsymbol{\nu}_n, \boldsymbol{\Omega}_n, \boldsymbol{\Lambda}_{n,1}, \boldsymbol{\Lambda}_{n,0}\}_{n=1}^N$ , with  $\boldsymbol{\nu}_n \in \mathbb{R}^D, \boldsymbol{\Lambda}_{n,0:1}, \boldsymbol{\Omega}_n \in \mathcal{CM}^N$  for all  $n$ . The dimensionality of the model for one source is:  $D + 3D(D+1)/2$  (mean vector and three covariance matrices), so  $\mathcal{D}(\mathbf{S}_{1:T,1:N}) = ND(3D+5)/2$ .

#### 2.1.2 Auditory status

The auditory states are represented in  $\mathbf{A}_t$ , which is a  $N$ -dimensional binary variable, i.e.,  $\mathbf{A}_t \in \{0, 1\}^N, \forall t$ .  $A_{t,n} = 1$  if the  $n^{\text{th}}$  source is speaking at time  $t$  and  $A_{t,n} = 0$  otherwise. Often, we abuse the notation by writing  $\mathbf{A}_t = i$ , being  $i = 1, \dots, 2^N$  instead of the binary representation of  $i - 1$ . The auditory states are modelled as an HMM:

$$p(A_t = i | A_{t-1} = j, \mathbf{T}^A) = \mathbf{T}^A[i, j] \quad (2.3)$$

$$p(A_1 = i | \boldsymbol{\lambda}^A) = \boldsymbol{\lambda}^A[i], \quad i, j = 1, \dots, 2^N, \quad (2.4)$$

With  $\theta^A = \{\mathbf{T}^A, \boldsymbol{\lambda}^A\}$ , being  $\mathbf{T}^A \in \mathcal{SM}^{2^N \times 2^N}$ , the transition matrix and  $\boldsymbol{\lambda}^A \in \mathcal{SV}^{2^N}$  the initial probabilities. The dimensionality of the model is  $\mathcal{D}(\mathbf{A}_{1:T}) = 2^N - 1 + 2^N(2^N - 1)$ , from the stochastic vector and matrix, that is  $\mathcal{D}(\mathbf{A}_{1:T}) = 4^N - 1$ .

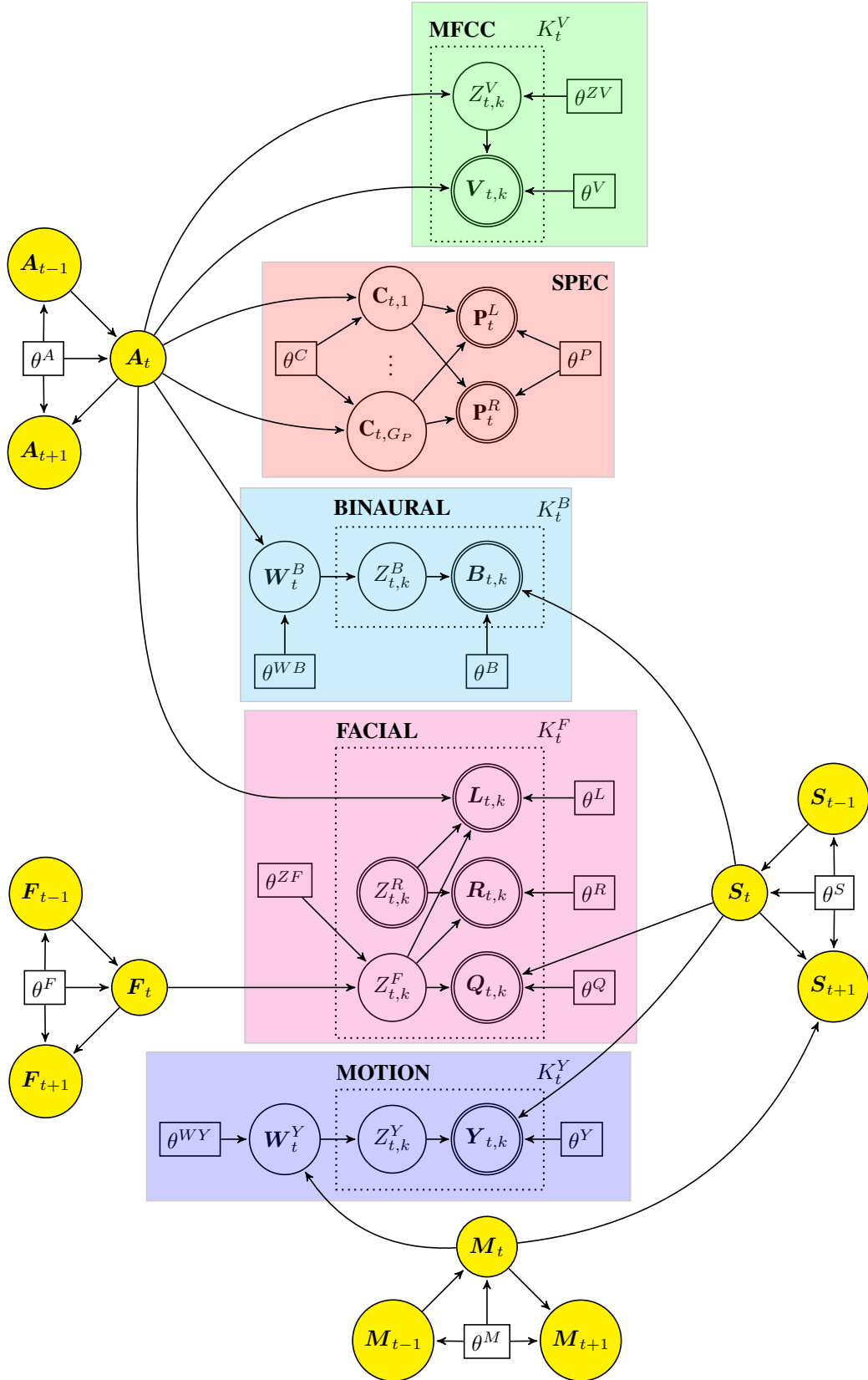


Figure 1: Full graphical model at time  $t$ . Rectangular nodes mean parameters, single circular nodes mean hidden variables and double circular nodes mean observations. Variables are grouped into six colors. The yellow ones,  $A_t$ ,  $F_t$ ,  $M_t$  and  $S_t$ , correspond to dynamic variables, i.e., those which their values across time are linked. The green variables,  $Z_{t,k}^V$  and  $V_{t,k}$ , are related to the modelling of the Mel Frequency Cepstral Coefficients (MFCC). Red variables,  $C_{t,1}, \dots, C_{t,J}$ ,  $P_t^L$  and  $P_t^R$ , model the spectrogram (SPEC). Cyan variables,  $W_t^B$ ,  $Z_{t,k}^B$  and  $B_{t,k}$ , model the extraction of sound source locates from binaural cues. (BINAURAL). Magenta variables,  $Z_{t,k}^F$ ,  $L_{t,k}$ ,  $R_{t,k}$  and  $Q_{t,k}$  are related to the facial detections (FACIAL). Finally, blue variables,  $W_t^Y$ ,  $Z_{t,k}^Y$  and  $Y_{t,k}$ , model the motion features (MOTION).

### 2.1.3 Visibility

The visibility variable indicates whether the face of the  $n^{\text{th}}$  speaker is visible at time  $t$  ( $F_{t,n} = 1$ ) or not ( $F_{t,n} = 0$ ). They are modelled as independent HMM:

$$p(F_{t,n} = i | F_{t-1,n} = j, \mathbf{T}_n^F) = \mathbf{T}_n^F[i, j] \quad (2.5)$$

$$p(F_{1,n} = i | \lambda_n^F) = \lambda_n^F[i]. \quad (2.6)$$

We notice that  $\theta^F = \{\mathbf{T}_1^F, \dots, \mathbf{T}_N^F, \lambda_1^F, \dots, \lambda_N^F\}$  with  $\mathbf{T}_n^F \in \mathcal{SM}^{2 \times 2}$  and  $\lambda_n^A \in \mathcal{SV}^2, \forall n = 1, \dots, N$ . Therefore, the dimensionality of the model is  $\mathfrak{D}(\mathbf{F}_{1:T}) = N(2(2-1)) + N(2-1) = 3N$ .

### 2.1.4 Motion

The motion variable indicates whether the  $n^{\text{th}}$  speaker is moving or not at time  $t$ ,  $M_{t,n} \in \{0, 1\}$ .

$$p(M_{t,n} = i | M_{t-1,n} = j, \mathbf{T}_n^M) = \mathbf{T}_n^M[i, j] \quad (2.7)$$

$$p(M_{1,n} = i | \lambda_n^M) = \lambda_n^M[i]. \quad (2.8)$$

Therefore  $\theta^M = \{\mathbf{T}_1^M, \dots, \mathbf{T}_N^M, \lambda_1^M, \dots, \lambda_N^M\}$ , the same as in the previous case. Obviously, the dimensionality of the model is also  $\mathfrak{D}(\mathbf{M}_{1:T}) = 3N$ .

## 2.2 Modalities

### 2.2.1 Mel Frequency Cepstral Coefficients

The Mel Frequency Cepstral Coefficients (MFCC) are denoted by  $\mathbf{V} \in \mathbb{R}^{D_V}$ . They correspond to the green box in Figure 1, also shown in Figure 2.  $\mathbf{V}$  are modelled using Gaussian Mixture Models (GMM). The set of parameters of the GMM, that is mixing proportions, mean vectors and covariance matrices, depend on the auditory status  $\mathbf{A}_t$ . We assume that  $K_t^V$  MFCC vectors are extracted at time  $t$ , and that they are independent and identically distributed realizations. GMMs have  $G_V$  components. For the auditory state  $\mathbf{A}_t = j$ , the GMM is parametrized by the mixing parameters  $\pi_j^V \in \mathcal{SV}^{G_V}$ , the mean vectors  $\mu_{j,1}^V, \dots, \mu_{j,G_V}^V \in \mathbb{R}^{D_V}$  and the covariance matrices  $\Sigma_{j,1}^V, \dots, \Sigma_{j,G_V}^V \in \mathcal{CM}^{D_V}, \forall j = 1, \dots, 2^N$ . The hidden variable  $Z_{t,k}^V$  represents the assignment of the  $k^{\text{th}}$  MFCC observation at time  $t$ ,  $\mathbf{V}_{t,k}$ , to one of the mixture's components. Therefore,  $Z_{t,k}^V \in \{1, \dots, G_V\}$ .

$$Z_{t,k}^V | \mathbf{A}_t = j \sim \text{Mult}(\pi_j^V) \quad (2.9)$$

$$\mathbf{V}_{t,k} | \mathbf{A}_t = j, Z_{t,k}^V = g \sim \mathcal{N}(\mu_{j,g}^V, \Sigma_{j,g}^V) \quad (2.10)$$

We now denote  $\theta^{ZV} = \pi_{1:2^N}^V$  and  $\theta^V = \{\mu_{j,g}^V, \Sigma_{j,g}^V\}_{j=1, g=1}^{2^N, G_V}$ . The dimensionality of the MFCC model per each value of the auditory state variable  $\mathbf{A}_t$  is  $G_V(D_V + D_V(D_V + 1)/2) = G_V D_V(D_V + 3)/2$ , so the dimensionality is:  $\mathfrak{D}(\mathbf{V}_{1:T, 1:K^V}) = 2^N G_V D_V(D_V + 3)/2$ .

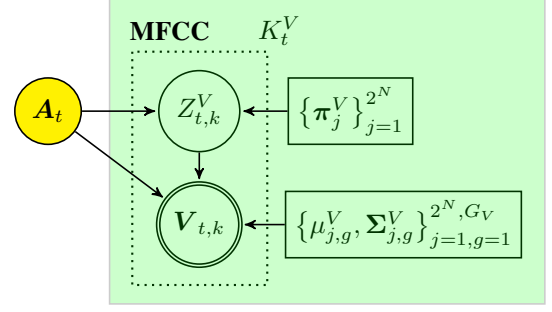


Figure 2: Graphical model of the MFCC. At each time interval we assume  $K_t^V$  i.i.d. realizations of the same GMM. The parameters of this GMM depend on the auditory state  $\mathbf{A}_t$ . The assignment variable  $Z_{t,k}^V$  describes which of the  $G_V$  Gaussian components has generated the  $k^{\text{th}}$  observation  $\mathbf{V}_{t,k}$ . The parameters corresponding to this part of the model are the mixing coefficients  $\pi_j^V$ , the mean vectors  $\mu_{j,g}^V$  and the covariance matrices  $\Sigma_{j,g}^V$ .

### 2.2.2 Spectrogram

The spectrogram is modelled using mixtures of proper<sup>1</sup> complex Gaussians (see Figure 3). We inspired from the works of [3, 5]. The mixture has  $G_P$  components, assigned to the different speakers. Let  $\{\mathcal{G}_n\}_{n=1}^N$  be a non-trivial partition of  $\{1, \dots, G_P\}$ , then component  $g$  is assigned to source  $n$  if and only if  $g \in \mathcal{G}_n$ . Each Time-Frequency point of each component's spectrogram follows a proper complex Gaussian:

$$\mathbf{C}_{t,g}[f, i] | \mathbf{A}_{t,n_g} \sim \mathcal{A}_{t,n_g} \mathcal{N}_c(0, \omega_g[f] \xi_g[i]) \quad (2.11)$$

with  $f = 1, \dots, P_F, i = 1, \dots, P_I$ .  $n_g$  is the source to which the  $g^{\text{th}}$  component is assigned (i.e.  $n_g = n \Leftrightarrow g \in \mathcal{G}_n$ ). The two vectors  $\omega_g \in (\mathbb{R}^+)^{P_F}$  and  $\xi_g \in (\mathbb{R}^+)^{P_I}$  model the power of the spectrogram of  $\mathbf{C}_{t,g}$ . The final left spectrogram is:

$$\mathbf{P}_t^L[f, i] | \mathbf{C}_{t,1:G_P}[f, i] \sim \mathcal{N}_c\left(\sum_{g=1}^{G_P} \Xi^L[f, g] \mathbf{C}_{t,g}[f, i], \Sigma_f^{PL}\right), \quad (2.12)$$

where  $\Xi^L \in \mathbb{C}^{P_F \times G_P}$  is the mixing matrix and  $\Sigma_f^{PL} \in \mathcal{CM}^1 = \mathbb{R}^+$  is the noise variance. The same formula holds for the right spectrogram with  $\Xi^R$  and  $\Sigma_f^{PR}$ . Notice that  $\Xi^L$  and  $\Xi^R$  are structured matrices. Indeed, the rows corresponding to the same source must be identical. Therefore, the parameters are:  $\theta^C = \{\omega_g, \xi_g\}_{g=1}^{G_P}$  and  $\theta^P = \{\Xi^L, \Xi^R, \{\Sigma_f^{PL}, \Sigma_f^{PR}\}_{f=1}^F\}$ ; and the dimensionality of the model is  $\mathfrak{D}(\mathbf{P}_{1:T}^L, \mathbf{P}_{1:T}^R) = G_P(P_F + P_I) + 2(2G_P P_F + P_F)$ , from the  $\omega$ 's and the  $\xi$ 's, the two mixing matrices  $\Xi^L, \Xi^R$  and the noise  $\Sigma_f^{PL}, \Sigma_f^{PR}$ .

### 2.2.3 Binaural SSL cues

We will use binaural cues to help in the source localisation task. Binaural cues, denoted by  $\mathbf{B}$  follow a GMM with a Dirichlet prior

<sup>1</sup>Proper complex Gaussians, also called circular symmetric complex Gaussians, have zero mean and zero relation matrix [4].

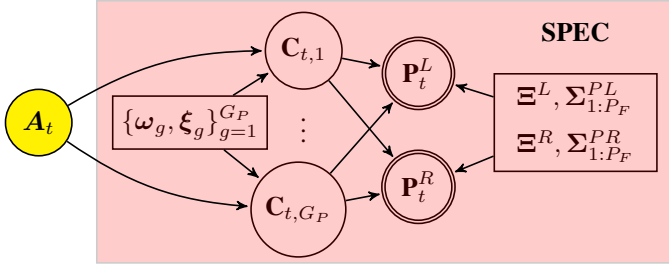


Figure 3: Graphical model of the spectrogram. The  $G_P$  components  $C_{t,g}$  are assigned to the different sources and the power spectrum at each time frequency point is modelled using  $\omega_g$  and  $\xi_g$ . This hidden components are mixed with the matrices  $\Xi^L$  and  $\Xi^R$ . Finally stationary noise is added in form complex Gaussian noise (with covariance matrices  $\Sigma_f^{PL}$  and  $\Sigma_f^{PR}$  respectively).

on the mixing coefficients (see Figure 4). We assume  $K_t^B$  observations at time  $t$ , all i.i.d., following the same GMM. The variable  $\mathbf{W}_t^B \in \mathcal{S}\mathcal{V}^N$  follows a Dirichlet distribution whose parameters depend on the auditory state:

$$\mathbf{W}_t^B | \mathbf{A}_t \sim \mathcal{D}(\alpha^B(\mathbf{A}_t)) \quad (2.13)$$

with

$$\alpha^B(\mathbf{A}_t)[n] = \begin{cases} \alpha_n^B[0] & \text{if } A_{t,n} = 0 \\ \alpha_n^B[1] & \text{if } A_{t,n} = 1 \end{cases} \quad n = 1, \dots, N. \quad (2.14)$$

The vector  $\mathbf{W}_t^B$  is used as mixing coefficients of the GMM [1] (Chapter 10), so the parameters of the multinomial distribution generating the assignment variable  $Z_{t,k}^B \in \{1, \dots, N\}$ :

$$Z_{t,k}^B | \mathbf{W}_t^B \sim \text{Mult}(\mathbf{W}_t^B). \quad (2.15)$$

Finally, the observations  $\mathbf{B}_{t,k} \in \mathbb{R}^D$  follow a Gaussian distribution:

$$\mathbf{B}_{t,k} | Z_{t,k}^B = n, \mathbf{S}_t \sim \mathcal{N}(\mathbf{S}_{t,n}, \Sigma_{t,n}^B) \quad (2.16)$$

Therefore,  $\theta^{WB} = \alpha_{1:N}^B$  and  $\theta^B = \Sigma_{1:T,1:N}^B$ ; and the dimensionality of the model is  $\mathfrak{D}(\mathbf{B}_{1:T,1:K^B}) = 2N + TND(D+1)/2$ , from the parameters of the Dirichlet, plus the covariance matrices ( $D(D+1)/2$  per person per time interval). If there are not enough samples to correctly estimate the covariance matrices, we will use one covariance per source (so the same for all time intervals).

## 2.2.4 Facial cues

Facial cues consist on four observations: face position  $\mathbf{Q}$ , face descriptor  $\mathbf{R}$ , pose detection  $Z^R$  and lip movement detection  $\mathbf{L}$ . We assume the existence of  $K_t^F$  face detections at interval  $t$ . The assignment variable,  $Z_{t,k}^F$  indicates to which source the  $k^{\text{th}}$  detected face is assigned, and follows a multinomial with  $N+1$  elements (since we account for false detections):

$$Z_{t,k}^F | \mathbf{F}_t \sim \text{Mult} \left( \frac{(1-e)F_{t,1}}{N_t^F}, \dots, \frac{(1-e)F_{t,N}}{N_t^F}, e \right), \quad (2.17)$$

being  $N_t^F = \sum_{n=1}^N F_{t,n}$  the number of visible people and  $e$  the false error rate of the face detector.

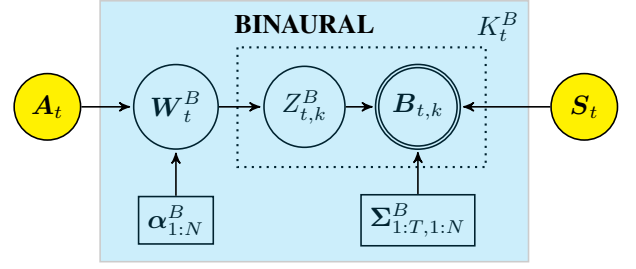


Figure 4: Graphical model for the binaural cues.  $\mathbf{W}_t^B$  follows a Dirichlet distribution whose parameters,  $\alpha_{1:N}^B$ , depend on the auditory state  $\mathbf{A}_t$ . One realization of the variable corresponds to the mixing coefficients for the GMM at time  $t$ . The means of the GMM correspond to the source positions  $\mathbf{S}_{t,1:N}$  and the covariances,  $\Sigma_{1:T,1:N}^B$  need to be estimated.

The face is detected using the method proposed in [7]. This method provides the position of the face  $\mathbf{Q}_{t,k}$ , an estimation of the face pose  $Z_{t,k}^R \in \{1, \dots, G_R\}$  and a set of  $G_F$  landmarks (the number may depend on the pose). The position of the face follows a Gaussian distribution:

$$\mathbf{Q}_{t,k} | Z_{t,k}^R = n, \mathbf{S}_t \sim \mathcal{N}(\mathbf{S}_{t,n}, \Sigma^Q) \quad (2.18)$$

where  $\Sigma^Q$  describes the distribution of the localisation error of the face detector. Therefore  $\theta^Q = \Sigma^Q$  and the dimensionality of the model is  $\mathfrak{D}_{1:\mathfrak{T},1:\mathfrak{R}^{\mathfrak{S}}}(\mathbf{Q}) = D(D+1)/2$ .

The extracted landmarks are used to build a face descriptor which consists on the concatenation of the landmarks' positions together with a  $D_A$ -dimensional local descriptor per landmark. Therefore,  $\mathbf{R}_{t,k} \in \mathbb{R}^{G_F(D+D_A)}$ . This observation is modelled as a Gaussian, in which the geometric part (the landmark's positions) is independent of the appearance part and the appearance observations are independent between them. We denote the geometric part of  $\mathbf{R}$  as  $\mathbf{R}^G \in \mathbb{R}^{G_FD}$  and the appearance parts as  $\mathbf{R}_g^A \in \mathbb{R}^{D_A}$ , for  $g = 1, \dots, G_F$ . Therefore  $\mathbf{R} = \left( (\mathbf{R}^G)^\top, (\mathbf{R}_1^A)^\top, \dots, (\mathbf{R}_{G_F}^A)^\top \right)^\top$ , modelled by:

$$\mathbf{R}_{t,k}^G | Z_{t,k}^F = n, Z_{t,k}^R = m \sim \mathcal{N}(\boldsymbol{\mu}_{n,m}^G, \Sigma_{n,m}^G) \quad (2.19)$$

$$\mathbf{R}_{t,k,g}^A | Z_{t,k}^F = n, Z_{t,k}^R = m \sim \mathcal{N}(\boldsymbol{\mu}_{n,m,g}^A, \Sigma_{n,m,g}^A), \quad (2.20)$$

for  $g = 1, \dots, G_F$ . We inspired from [2], but for the time being we remove the priors on the matrices and means. Therefore  $\theta^R = \left\{ \boldsymbol{\mu}_{n,m}^G, \Sigma_{n,m}^G, \{ \boldsymbol{\mu}_{n,m,g}^A, \Sigma_{n,m,g}^A \}_{g=1}^{G_F} \right\}_{n=1,m=1}^{N,G_R}$ . In all,  $\mathfrak{D}(\mathbf{R}_{1:T,1:K^F}) = NG_R(G_FD(G_FD+3) + G_FD_A(D_A+3))/2$ .

The lip movement descriptor will consist on the concatenation of spatio-temporal features (e.g. Laptev or [6]). Therefore, we consider  $G_L$  (one per lip landmark) independent  $D_L$ -dimensional features all of them following a different Gaussian:

$$\mathbf{L}_{t,k,g} | Z_{t,k}^F = n, Z_{t,k}^R = m, A_{t,n} = a \sim \mathcal{N}(\boldsymbol{\mu}_{n,m,g,a}^L, \Sigma_{n,m,g,a}^L), \quad (2.21)$$

with  $g = 1, \dots, G_L$ . Therefore  $\theta^L = \left\{ \boldsymbol{\mu}_{n,m,g,a}^L, \Sigma_{n,m,g,a}^L \right\}_{n=1,m=1,g=1,a=0}^{N,G_F,G_L,1}$  and the dimensionality of  $\mathbf{L}$ :  $\mathfrak{D}(\mathbf{L}_{1:T,1:K^F}) = 2NG_RG_LD_L(D_L+3)/2$ .



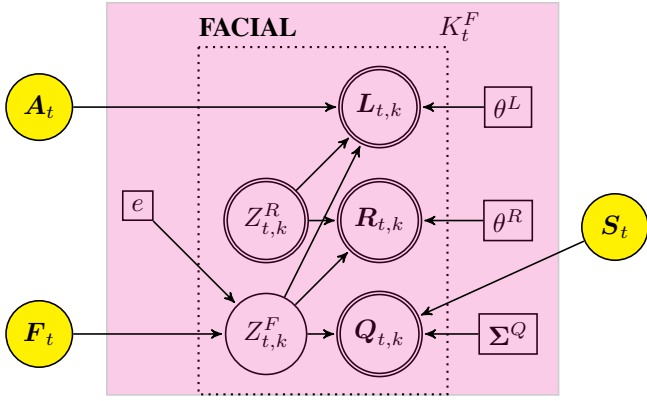


Figure 5: Graphical model of the facial cues.  $K_t^F$  faces are detected at time  $t$ .  $Z_{t,k}^F$  assigns the faces to the sources, and  $Z_{t,k}^R$  is the estimated pose.  $Q_{t,k}$  is the position of the face,  $R_{t,k}$  describes the geometry and the appearance of the face landmarks and  $L_{t,k}$  describes the local dynamics of the lip landmarks.

### 2.2.5 Movement cues

Movement cues have a model very similar to the binaural cues. We will use movement cues to help in the source localisation task. Motion cues,  $\mathbf{Y}$  will be assumed to follow a GMM with a Dirichlet prior on the mixing coefficients (see Figure 6). We assume  $K_t^Y$  observations at time  $t$ , all i.i.d., following the same GMM. The variable  $\mathbf{W}_t^Y \in \mathcal{S}\mathcal{V}^N$  follows a Dirichlet distribution whose parameters depend on the auditory state:

$$\mathbf{W}_t^Y | \mathbf{M}_t \sim \mathcal{D}(\boldsymbol{\alpha}^Y(\mathbf{M}_t)) \quad (2.22)$$

with

$$\boldsymbol{\alpha}^Y(\mathbf{M}_t)[n] = \begin{cases} \alpha_n^Y[0] & \text{if } M_{t,n} = 0 \\ \alpha_n^Y[1] & \text{if } M_{t,n} = 1 \end{cases} \quad n = 1, \dots, N. \quad (2.23)$$

The vector  $\mathbf{W}_t^Y$  is used as mixing coefficients of the GMM, so the parameters of the multinomial distribution generating the assignment variable  $Z_{t,k}^Y \in \{1, \dots, N\}$ :

$$Z_{t,k}^Y | \mathbf{W}_t^Y \sim \text{Mult}(\mathbf{W}_t^Y). \quad (2.24)$$

Finally, the observations  $\mathbf{Y}_{t,k} \in \mathbb{R}^D$  follow a Gaussian distribution:

$$\mathbf{Y}_{t,k} | Z_{t,k}^Y = n, \mathbf{S}_t \sim \mathcal{N}(\mathbf{S}_{t,n}, \boldsymbol{\Sigma}_{t,n}^Y) \quad (2.25)$$

Therefore,  $\theta^{WY} = \boldsymbol{\alpha}_{1:N}^Y$  and  $\theta^Y = \boldsymbol{\Sigma}_{1:T,1:N}^Y$ ; and the dimensionality of the model is  $\mathfrak{D}(\mathbf{Y}_{1:T,1:K^Y}) = 2N + TND(D+1)/2$ , from the parameters of the Dirichlet, plus the covariance matrices ( $D(D+1)/2$  per person per time interval). If there are not enough samples to correctly estimate the covariance matrices, we will use one covariance per source (so the same for all time intervals).

## References

[1] Christopher M Bishop and Nasser M Nasrabadi. *Pattern recognition and machine learning*, volume 1. springer New York, 2006.

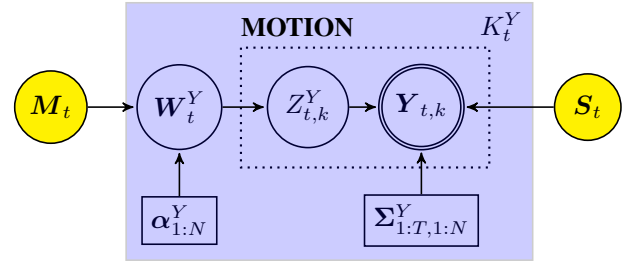


Figure 6: Graphical model for the binaural cues.  $\mathbf{W}_t^Y$  follows a Dirichlet distribution whose parameters,  $\boldsymbol{\alpha}_{1:N}^Y$ , depend on the auditory state  $\mathbf{M}_t$ . One realization of the variable corresponds to the mixing coefficients for the GMM at time  $t$ . The means of the GMM correspond to the source positions  $\mathbf{S}_t$  and the covariances  $\boldsymbol{\Sigma}_{1:T,1:N}^Y$  need to be estimated.

- [2] Li Fei-Fei, Rob Fergus, and Pietro Perona. Learning generative visual models from few training examples: An incremental bayesian approach tested on 101 object categories. *Computer Vision and Image Understanding*, 106(1):59–70, 2007.
- [3] Cédric Févotte, Nancy Bertin, and Jean-Louis Durrieu. Non-negative matrix factorization with the itakura-saito divergence: With application to music analysis. *Neural computation*, 21(3):793–830, 2009.
- [4] NR Goodman. Statistical analysis based on a certain multivariate complex gaussian distribution (an introduction). *The Annals of Mathematical Statistics*, 34(1):152–177, 1963.
- [5] Alexey Ozerov and Cédric Févotte. Multichannel nonnegative matrix factorization in convolutive mixtures for audio source separation. *Audio, Speech, and Language Processing, IEEE Transactions on*, 18(3):550–563, 2010.
- [6] Guoying Zhao, Mark Barnard, and Matti Pietikainen. Lipreading with local spatiotemporal descriptors. *Multimedia, IEEE Transactions on*, 11(7):1254–1265, 2009.
- [7] Xiangxin Zhu and Deva Ramanan. Face detection, pose estimation, and landmark localization in the wild. In *Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on*, pages 2879–2886. IEEE, 2012.