



**HAL**  
open science

## CRF based context modeling for person identification in broadcast videos

Paul Gay, Sylvain Meignier, Jean-Marc Odobez, Paul Deléglise

► **To cite this version:**

Paul Gay, Sylvain Meignier, Jean-Marc Odobez, Paul Deléglise. CRF based context modeling for person identification in broadcast videos. *Frontiers in information and communication technologies*, 2016, 3, pp.9. 10.3389/fict.2016.00009 . hal-01433154

**HAL Id: hal-01433154**

**<https://hal.science/hal-01433154>**

Submitted on 21 Mar 2017

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

---

# CRF based context modeling for person identification in broadcast videos

Paul Gay<sup>1,2</sup>, Sylvain Meignier<sup>1</sup>, Paul Deléglise<sup>1</sup> and Jean-Marc Odobez<sup>2</sup>

<sup>1</sup>*LIUM Laboratory, Le Mans, France*

<sup>2</sup>*Idiap Research Institute, Martigny, Switzerland*

Correspondence\*:

Jean-Marc Odobez

Idiap Research Institute, Centre du Parc, Rue Marconi 19, Case Postale 592,  
CH-1920, Martigny, Switzerland, odobez@idiap.ch

## 2 ABSTRACT

3 We are investigating the problem of speaker and face identification in broadcast videos.  
4 Identification is performed by associating automatically extracted names from overlaid texts  
5 with speaker and face clusters. We aimed at exploiting the structure of news videos to solve  
6 name/cluster association ambiguities and clustering errors. The proposed approach combines  
7 iteratively two Conditional Random Fields (CRF). The first CRF performs the person diarization  
8 (joint temporal segmentation, clustering and association of voices and faces) jointly over the  
9 speech segments and the face tracks. It benefits from contextual information being extracted from  
10 the image backgrounds and the overlaid texts. The second CRF associates names with person  
11 clusters thanks to co-occurrence statistics. Experiments conducted on a recent and substantial  
12 public dataset containing reports and debates demonstrate the interest and complementarity of  
13 the different modeling steps and information sources: the use of these elements enables us to  
14 obtain better performances in clustering and identification, especially in studio scenes.

15 **Keywords:** Face identification, speaker identification, broadcast videos, conditional random field, face clustering, speaker diarization

16

## 1 INTRODUCTION

17 For the last two decades, researchers have been trying to create indexing and fast search and browsing  
18 tools capable of handling the growing amount of available video collections. Among the associated  
19 possibilities, person identification is an important one. Indeed, video contents can often be browsed through  
20 the appearances of their different actors. Moreover, the availability of each person intervention allows  
21 easier access to video structure elements such as the scene segmentation. Both motivations are especially  
22 verified in the case of news collections. The focus of this paper is therefore to develop a program able to  
23 identify persons in broadcast videos. That is, the program must be able to provide all temporal segments  
24 corresponding to each face and speaker.

25 Person identification can be supervised. A face and/or a speaker model of the queried person is then  
26 learned over manually labeled training data. However, this raises the problem of annotation cost. An  
27 unsupervised and complementary approach consists in using the naming information already present in  
28 the documents. Such resources include overlaid texts, speech transcripts and metadata. Motivated by this



**Figure 1.** Example frames from the REPERE corpus showing the variety of the visual conditions (pose, camera viewpoint, illumination) and the name face association challenges such as: multiface images (image c) and name propagation (from a to b). Images a and c show an example of OPNs.

29 opportunity, unsupervised identification has been investigated for 15 years from the early work of Satoh  
 30 et al. (1999) to the development of more complex news browsing systems exploiting this paradigm (Jou  
 31 et al. (2013)), or thanks to sponsored competitions (Giraudel et al. (2012)). Whatever the source of naming  
 32 information, it must tackle two main obstacles: associate the names to co-occurring speech and face  
 33 segments, and propagate this naming information from the co-occurring segments to the other segments of  
 34 this person.

35 There are several challenges related to this task. First, the named entities need to be recognized and an  
 36 association step must decide if the name corresponds to people co-occurring in the document. Ambiguities  
 37 arise when multiple audiovisual (AV) segments co-occur with one name. This is illustrated in Fig 1c where  
 38 there is more than one face in the image. This situation is becoming more common with modern video  
 39 editing. Regarding the identity propagation, it can be done with speaker and face diarization techniques  
 40 (detecting and clustering person interventions). However, these two tasks have been active fields of  
 41 research for more than a decade and thus are difficult problems to solve. Indeed, a person may appear in  
 42 different contexts thus introducing huge intraperson variabilities. We can distinguish them in function of  
 43 the modalities and the different types of videos. For the speaker diarization, the main challenge in broadcast  
 44 news is background noise such as music, or a noisy environment during outside reports. If we consider  
 45 debates in studio where the speech is more spontaneous, the bottleneck becomes the overlapping speech  
 46 and short speech segments. Regarding face diarization, report videos usually exhibit the largest variations  
 47 as location and time may change between two scenes, and so will be the illumination conditions. For the  
 48 debate and studio scenes, variations come essentially from changes in the facial poses.

49 In this paper, we assume that closed-captions are not available as this is the case in European medias.  
 50 Instead, we focus on Overlaid Person Names (OPNs) which are used to introduce the speakers as illustrated  
 51 in Fig. 1a. Such names are appealing since their extraction is much more reliable than pronounced names  
 52 obtained through Automatic Speech Recognition (ASR). Moreover, their association with face or speech  
 53 segments is in general easier than analyzing whether pronounced names in ASR transcripts refer to people  
 54 appearing in the video. The identification systems submitted at the recent REPERE campaign (Bredin et al.  
 55 (2013); Bechet et al. (2014); Poignant et al. (2014)) mainly rely on such names.

56 Our approach offers several advantages. Faces are identified by alternating between a clustering step of  
 57 faces and audio speech segments and a naming step of the resulting AV clusters. Each step is performed by a  
 58 dedicated CRF. The use of CRF enables us to include heterogeneous context cues in our modeling. The use  
 59 of such cues is challenging because they must use as little specific prior information as possible in order to

60 achieve generalization over the different types of videos. In this paper, we include different generic context  
61 cues. First, we have AV association scores which enable to associate overlapping speaker and face segments  
62 when they correspond to the same person. Then, we use uniqueness constraints between simultaneously  
63 appearing pairs of faces. Furthermore, one of the main contribution is a background recurrence descriptor  
64 which attributes a soft role to each segment. It enables to distinguish the persons which are announced  
65 by the OPNs such as guests or journalists from the anonymous persons appearing around them. Last but  
66 not least, the names contained in the OPNs are included to guide the clustering by using the probabilities  
67 obtained with the naming CRF. These different cues enable to improve the clustering by reducing errors due  
68 to monomodal intracluster variations such as facial pose or audio background noise. Eventually, the CRF  
69 formulation avoids hard local decisions by providing a joint probability distribution over all the segments.

70 The first CRF performs jointly the clustering of face tracks and speaker segments thanks to AV association  
71 as introduced in Gay et al. (2014c). In practice, AV association is initialized in a pre-processing step based  
72 on temporal co-occurrence and then refined inside the CRF thanks to talking head detection scores and  
73 the previously described contextual cues. The second CRF assigns a name to each cluster by using co-  
74 occurrence statistics and a uniqueness constraint preventing any two faces on the same image to receive  
75 the same name. In Gay et al. (2014b), this approach was designed for face identification. In the present  
76 case, we extend this approach for the AV case and provide results for the final evaluation of the REPERE  
77 campaign. Identification performances are discussed by investigating the algorithm behavior in different  
78 types of shows (reports, news, debates, celebrity magazines) and the relations with the clustering quality.

79 The rest of the article is organized as follows: Section 2 reviews related work on unsupervised  
80 identification. Then, Section 3 presents the proposed CRF-based system. Experiments and results are  
81 presented in section 4. Finally, Section 5 sums up our main findings and concludes the paper.

## 2 RELATED WORK

82 As stated in the introduction, unsupervised people identification must address the problems of local  
83 person/name association and propagation to the video parts where the names are absent. The association  
84 is conducted via the use of co-occurrence statistics between the names present in the document and the  
85 detected persons. The propagation can be seen as a clustering problem. Clustering methods can regularly  
86 benefit from new improvements in speaker and face representations. At the time of writing, the ivector  
87 approach is one of the most successful (Rouvier et al. (2013)) for the speaker diarization task. Regarding  
88 face representation, recent advances include encodings (Simonyan et al. (2013)), metric learning (Bhattarai  
89 et al. (2014)) and feature learning by deep Convolutional Neural Networks (Schroff et al. (2015)). However,  
90 most of the systems require explicit face alignment to obtain frontal views which is not always feasible.  
91 The work published in Zhang et al. (2015) suggests that using only face representation is a great limitation  
92 when dealing with unconstrained views of persons. For this reason, we believe that investigation into  
93 context-assisted clustering is justified, especially for broadcast news videos which exhibit a strong structure.

94 To identify the faces, most approaches try to solve the association and the propagation problems jointly.  
95 On one hand, co-occurrence statistics at cluster level are more discriminant and accurate than just describing  
96 a face locally with namedness features (like face position or talking activity) to assess whether the detected  
97 name should be associated. On the other hand, name/face co-occurrences are used as a contextual cue to  
98 improve the face clustering process. These principles have been used intensively since the seminal works  
99 of Berg et al. (2004) and Everingham et al. (2006) which applied to two representative use-cases: captioned  
100 images, as exemplified by the *Yahoo News!* dataset, and soap series with the buffy dataset. The first case

101 study consists of news articles with images illustrating the subject. The initial approach described in Berg  
102 et al. (2004) is an EM clustering where the update of the model parameters takes into account the name/face  
103 co-occurrences. In this context, the work of Ozkan and Duygulu (2010) exploits the fact that a textual query  
104 enables to retrieve faces where the queried person holds the majority. The problem of finding those faces is  
105 posed as finding the densest component in a graph. This idea was later extended in Guillaumin et al. (2010)  
106 where the distance within clusters is minimized with respect to a cannot-link constraint which implies  
107 that two faces must belong to different clusters if their captions contain different names. However, those  
108 co-occurrence statistics can fail when group of people co-occur in a similar fashion, a situation commonly  
109 encountered in TV programs. In soap series, the names of the speakers can be obtained with the transcripts  
110 and the subtitles. Works in Cour et al. (2011); Wohlhart et al. (2011); Bauml et al. (2013) use those names  
111 as weak labels to improve supervised classifiers. They choose a learning setting which takes into account  
112 the label ambiguities, for example: multiple instance learning (Wohlhart et al. (2011)) and semi-supervised  
113 strategies (Bauml et al. (2013)). Talking head detection (Everingham et al. (2006); Cour et al. (2011)) and  
114 dialogue cues (Cour et al. (2010)) are also used to solve the ambiguities in the face/name association. Note  
115 that in the previous two case studies, the naming co-occurrence statistics are quite different to those in  
116 broadcast videos where the OPNs are more sporadic. Indeed, the OPN of a given person only appears  
117 one or a few times (usually for the first time utterance). This scarcity increases the dependence of the  
118 identification performance on the clustering quality.

119 Originally, unsupervised speaker identification in broadcast news was conducted by first performing a  
120 speaker diarization (i.e. clustering) step of the audio track and then assigning the names extracted from  
121 the transcription to the speaker clusters by using semantic classification trees (Jousse et al. (2009)) or  
122 Maximum-entropy classifiers (Ma et al. (2007)). More recently, the idea of constrained speaker clustering  
123 has been exploited in Bredin and Poignant (2013); Poignant et al. (2014). The system described in Bredin  
124 and Poignant (2013) defines a graph where the nodes are speaker segments and OPNs. OPNs are used to  
125 express must-link and cannot-link constraints between the utterances. The clustering and the naming of  
126 those segments is done using an Integer Linear Programming formulation. As first investigated by Li et al.  
127 (2001), the case study of videos allows to exploit the complementarity of audio and video modalities. AV  
128 cues such as talking head detection scores can be used to match faces and speakers and to improve the  
129 monomodal speaker and face diarizations. The scores of such cues are computed by estimating motion in  
130 the region of the lips. In addition, features such as the face size, the face position or the number of faces in  
131 the image are extracted and given to a supervised classifier (El Khoury et al. (2012); Vallet et al. (2013)) to  
132 further refine the talking assessment. In order to bring corrections to the initial monomodal diarizations,  
133 the talking head detection scores should be reliable where monomodal errors are present. Moreover, the  
134 audio and video will also be more complementary if they make errors at different moments. In other words,  
135 the improvements of the AV diarization depend on the performances of the initial monomodal ones. The  
136 work of Noulas et al. (2012) integrates faces and speech segments in a factorial hidden markov model. The  
137 assignment of a segment to a cluster label is based on biometric model and on AV links with co-occurrent  
138 segments from the other modality. The use of a graphical model enables to express dependences between  
139 variables with a global probabilistic formulation which can then be optimized jointly. In order to jointly  
140 identify faces and speakers, the authors of Poignant et al. (2015) proposed a constrained multimodal  
141 clustering. They use the simple idea that two segments which co-occur with different names implies that  
142 they should be assigned to different clusters. The authors also showed that their multimodal clustering of  
143 faces and speakers can make use of talking head detection scores to correct errors present in the monomodal  
144 systems.

145 The work of Bechet et al. (2014), an interesting yet not detailed contribution to the field, reports the  
146 intensive use of multimodal scene understanding cues. First, speaker diarization is performed and speakers  
147 are identified using OPNs or pre-trained models. Then, identities are propagated from the speakers to the  
148 faces. Scene segmentation, role detection, and pre-trained visual models for each TV set (and sometimes  
149 for each camera) are used to indicate how many faces are present on screen and what their roles are. Such  
150 a fine-grain modeling enables them to report the best identification on the REPERE campaign. Indeed,  
151 it permits to tell which persons are present without detecting the faces by detecting the specific shot (up  
152 to which studio camera is used). Thus, profile views and persons seen from the back can be identified.  
153 However, to learn those models, manual annotations have been made for each show. This poses the problem  
154 of human labor cost and lack of generalization. More generally, several researchers focus on exploiting the  
155 context surrounding the faces. The work in Zhang et al. (2013) uses clothes, image background, cluster  
156 co-occurrences and attribute classifiers and Tapaswi et al. (2014) build must-link and cannot-link constraints  
157 deduced from shot threads (sequence of shots obtained from the same camera angle).  
158

159 **Contributions:** in this paper, we leverage on different contextual cues present in the state of the art,  
160 introduce new ones, and include them in our CRF model. First, instead of conducting speaker and face  
161 clustering separately [15, 5], we perform a joint clustering of face tracks and speaker segments which also  
162 benefits from the OPNs information. To be more precise, we compute Local Face visual Backgrounds  
163 (LFBs) around each face track and cluster them. This provides us with a signature for each face track  
164 characterizing the level of recurrence of its LFB in the data. Intuitively, recurrent a LFB correspond to  
165 people who are important and can be seen as a soft role assignment distinguishing faces to be named from  
166 faces of figurative people. Concretely, it enables to encourage faces tracks with recurrent LFBs to join  
167 named clusters, i.e. overlapping an OPN. Secondly, a naming CRF performs the joint identification of all  
168 person clusters, thus allowing to account for uniqueness constraints and co-occurrence statistics between  
169 clusters and OPNs. Unlike previous works which rely on extensive annotations (Bechet et al. (2014)), those  
170 elements of context have better generalization capabilities, since we can learn one single model over a  
171 large and diversified corpus, and require less annotations if we want to learn a new type of show. Thanks to  
172 the flexibility of the CRF formulation, new contextual cues could be added in the future to further improve  
173 the performances.

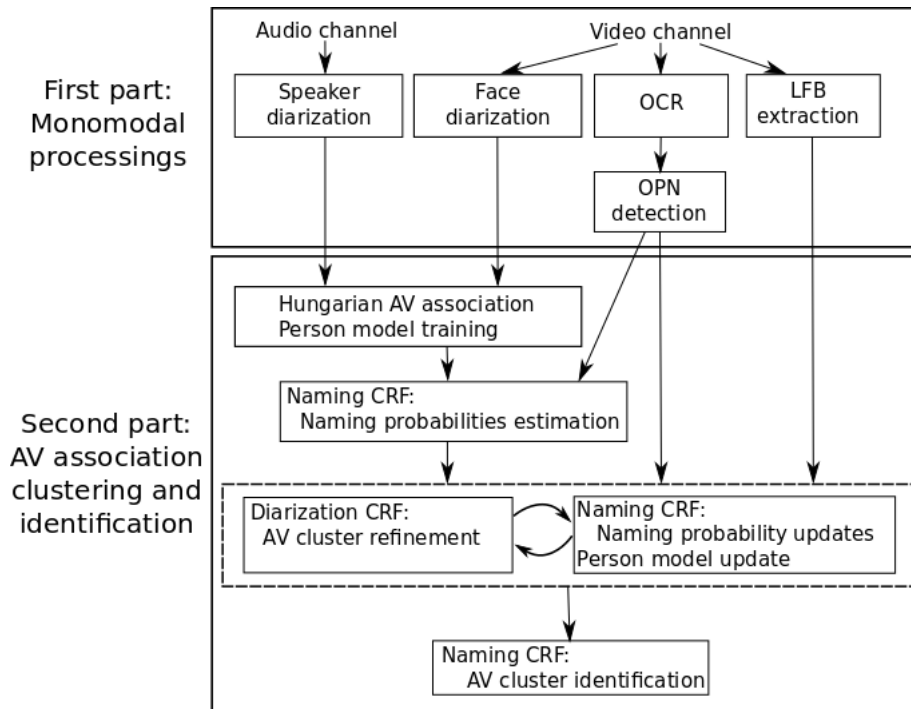
### 3 METHOD

174 The method will be first described globally in section 3.1. In the section 3.2, we introduce the notations.  
175 We then describe how we extract LFB and AV association features in sections 3.3 and 3.4. In section 3.5,  
176 the diarization CRF which clusters face and speech segments is presented, and in section 3.6 the naming  
177 CRF which is in charge of identifying the clusters. To conclude this part, we describe how the full system  
178 is used and optimized in section 3.7.

#### 179 3.1 Method overview

180 The general approach is summarized in Fig 2. First, the different modalities are processed separately:  
181 monomodal speaker (Rouvier et al. (2013)) and face (Khoury et al. (2013)) diarizations are performed,  
182 LFBs are extracted around each face and clustered, Optical Character Recognition (OCR) is performed to  
183 extract the overlaid texts (Chen and Odobez (2005)) and named entities are detected (Gay et al. (2014a)).

184 In the second part, we perform the AV clustering and the naming of the persons. Initially, we use the  
185 Hungarian algorithm to associate face and speaker clusters based on their temporal overlap. Naming



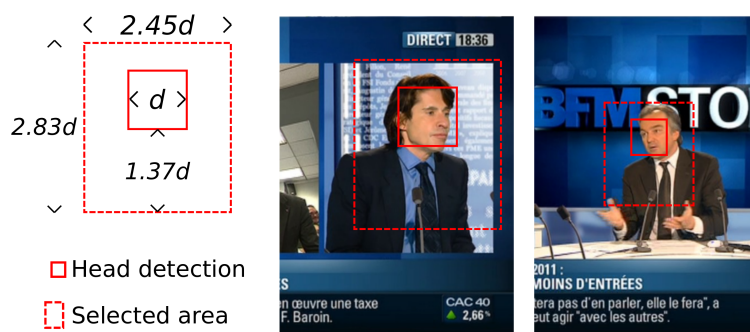
**Figure 2.** Overview of the system. First, face tracks and speech utterances are detected and clustered separately. They are then used with the OPNs in an iterative combination of the two CRFs. A first one to refine the diarization, and the other one to identify the clusters. The latter is eventually used for the final cluster identification.

186 probabilities are then computed onto those AV clusters with the naming CRF. Lastly, the system iterates  
 187 over a clustering step and a naming step. In the clustering step, the diarization CRF infers a cluster label  
 188 for each face track and utterance given the naming probabilities, an acoustic and visual person model for  
 189 each cluster label, and various context clues including the LFBs. In the naming step, person models and  
 190 naming probabilities are updated as a result of the new diarization. The motivation factor being that the  
 191 diarization CRF is able to use contextual clues to correct potential clustering errors made by the monomodal  
 192 diarizations and thus improve the final identification. Lastly, a name is associated to each cluster with the  
 193 naming CRF.

### 194 3.2 Notations

195 The pre-processing includes obtaining initial monomodal face and speaker clusters, a set of OPNs and  
 196 extracting the features from those elements. First, faces are detected (Viola and Jones (2004)) and tracked  
 197 within each shot, resulting in a set of face tracks denoted as  $V = \{V_i, i = 1 \dots N^V\}$ . Each face track  $V_i$   
 198 is characterized by a set of visual features  $x_i^{\text{surf}}$  (set of Speeded-Up-Robust Features (SURF) extracted in up  
 199 to 9 images of the face track (El Khoury et al. (2010))) and a set of boolean features  $\{x_i^{\text{lfbv}}(k), k \in K\}$   
 200 indicating whether  $V_i$  corresponds to a recurrent LFB as explained in the next section 3.3.

201 Second, OCR (Chen and Odobez (2005)) and Named Entity Detection techniques based on string  
 202 matching against external resources (predefined lists, freebase database, Google hits...) are applied as  
 203 described in Gay et al. (2014a) to extract the set  $O = \{O_i, i = 1 \dots N^O\}$  of OPNs. Each OPN  $O_i$   
 204 is characterized by its duration  $d_i^{\text{opn}}$  and its name  $x_i^{\text{opn}} \in M$  where  $M = \{n_j, j = 1 \dots N^M\}$  denotes the set  
 205 of unique names extracted from the video.



**Figure 3.** Left: predefined spatial template for the background face selection area given a head detection. Center and right: two examples with head detections and the selected background areas.



**Figure 4.** Top: original images. Bottom: clusters of recurrent local face backgrounds (LFB) automatically cropped from each image and clustered given the localised faces. Each row corresponds to a local background cluster. Clusters 1 and 2 are recurrent backgrounds and correspond to speakers. Clusters 3 and 4 are not recurrent. They contain non-speaking faces which appear occasionally in the video.

206 Finally, the audio stream is segmented into a set  $A = \{A_i, i = 1 \dots N^A\}$  of continuous speech segments  
 207 called utterances, each described by a set of acoustic features  $x_i^a$ . Features are 12 MFCCs with first order  
 208 derivatives. Each frame is normalized with a short-term windowed mean and variance. Feature warping is  
 209 also applied. In addition, a set of boolean features  $\{x_i^{lfb}(k), k \in K\}$  is extracted indicating whether  $A_i$  is  
 210 co-occurring with a recurrent LFB, as described in the next section. Finally, talking head detection features  
 211  $x_{ij}^{av}$  are extracted between each couple  $(A_i, V_j)$  with a non-zero overlap as described in section 3.4.

212 **3.3 Local Face Background recurrence**

213 We want to capture whether a face appears with a recurrent visual background. This feature will be  
 214 included in the diarization CRF. To this end, we focus on an area around each face track  $V_i$  to capture the  
 215 background context of this face. We do not consider full images as the same image might include different  
 216 face visual contexts (see the first, fourth and fifth images from the left in the top row of Fig 4). Instead,  
 217 we select a rectangle area around each face as local face background (LFB) representative by following a  
 218 predefined spatial template between the face and this rectangle as can be seen in Fig 3. In practice, the  
 219 fixed proportions were chosen manually so as to avoid a potential overlap with other parts of the images in  
 220 typical edited videos like in the 4th and 5th images from the left on the top of Fig 4. We then characterize



221 each obtained rectangle area with SURF features and, in order to cluster them, we use a hierarchical  
 222 clustering approach (El Khoury et al. (2010)). Then, we set  $x_i^{\text{fbv}}(k)$  to true if face track  $V_i$  belongs to a  
 223 local background cluster whose number of elements is higher than  $k$ . In practice, multiple values of  $k$  can  
 224 be used to characterize different levels of recurrence and reduce the importance of the stopping criterion of  
 225 the hierarchical clustering. Fig 4 shows examples of obtained recurrent and non-recurrent patterns.

### 226 3.4 Talking head detection features

227 In order to integrate AV association information in the CRF, we detect talking heads. To characterize  
 228 talking heads, we use the following measures. These features are extracted for each overlapping  
 229 utterance/face track couple and include:

- 230 • **lip activity:** the lip activity of a given face at frame  $k$  is computed as described in El Khoury et al.  
 231 (2012) and consists in the mean intensity difference between frame  $k$  and  $k + 1$  after local image  
 232 registration in predefined regions corresponding to the lips. In addition, we focus on the relative lip  
 233 activity by dividing by the sum of all the lip activities measured from all people in the image.
- 234 • **Head size:** the interest of this feature relies on the hypothesis that the face of the speaker is usually  
 235 larger than the faces of other people in the image. Put simply, we take the diagonal size of the detection  
 236 bounding boxes. We also use the relative head size.

237 The previous features are computed from each frame of the face track. Eventually, the final feature  $x_{ij}^{\text{av}}$  is an  
 238 average over all values from the frames included in the overlap between the utterance  $A_i$  and the face track  
 239  $V_j$ . This corresponds to the method used in Gay et al. (2014c). To assess whether a couple of utterance/face  
 240 track corresponds to a talking head given the features, we use an SVM with gaussian kernel denoted as  $h$ .

### 241 3.5 Audio-visual (AV) person diarization CRF

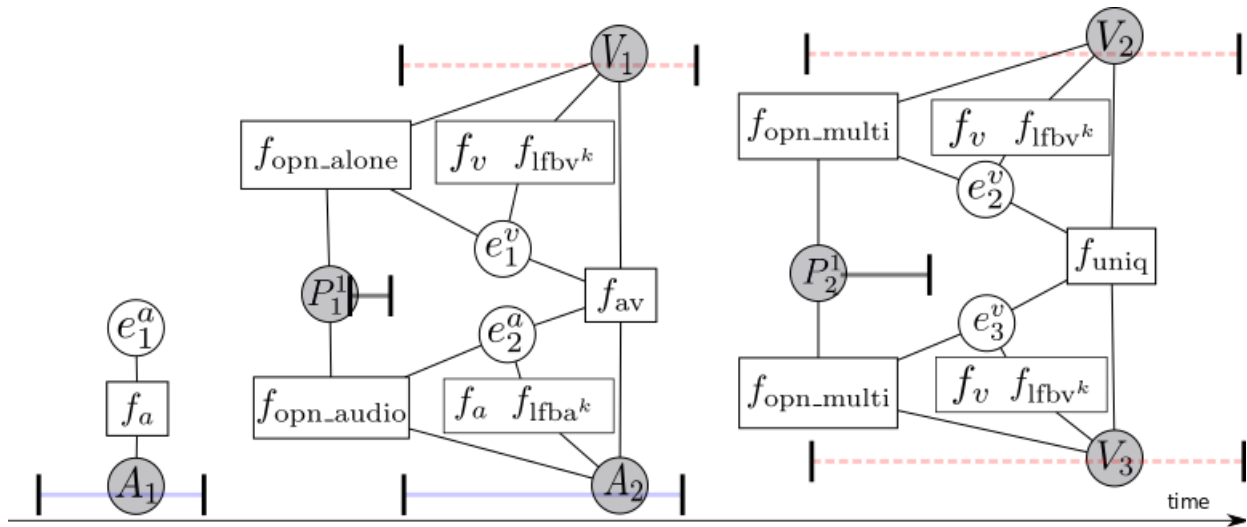
242 The clustering of face tracks and utterances defines itself by estimating the label field  $E^d = \{e_i^a, i =$   
 243  $1 \dots N^A, e_j^v, j = 1 \dots N^V\}$  as such, the same person index is used for  $e_i^a$  and  $e_j^v$  when the utterance  $A_i$   
 244 and the face track  $V_j$  correspond to the same person. The labels  $e_i^a$  and  $e_j^v$  take value in the set of possible  
 245 person indices denoted as  $P$ . To achieve this, let  $G$  be an undirected graph over the set of random variables  
 246  $A, V, O$ , and  $E^d$ . We then seek to maximize the CRF posterior probability formulated as:

$$P(E^d|A, V, O) = \frac{1}{Z(A, V, O)} \times \exp\left\{ \sum_{i \in \mathcal{F}} \sum_{c \in G_i} \lambda_i f_i(A_c, V_c, O_c, E_c^d) \right\} \quad (1)$$

247 where each triplet  $(f_i, G_i, \lambda_i)$  is composed of a feature function  $f_i$ , a weight  $\lambda_i$  learned at training time  
 248 and the set  $G_i$  of cliques where this function is defined.  $(A_c, V_c, O_c, E_c)$  denotes the set of nodes contained  
 249 in the clique  $c$ .  $\mathcal{F}$  is a set of abstract functions indices. We use 6 types of feature functions which will be  
 250 described in the next sections. A graphical representation of this model is illustrated on Fig 5.

251 The association function  $f_{\text{av}}$  favors the association of talking heads to utterances. The function is defined  
 252 on all overlapping utterance/face track couples  $\{(i, j)/t(A_i, V_j) \neq 0\}$  where  $t(A_i, V_j)$  is the overlapping  
 253 time duration between segments  $A_i$  and  $V_j$ :

$$f_{\text{av}}(A_i, V_j, e_i^a, e_j^v) = \begin{cases} t(A_i, V_j)h(x_{ij}^{\text{av}}) & \text{if } e_i^a = e_j^v \\ -t(A_i, V_j)h(x_{ij}^{\text{av}}) & \text{otherwise} \end{cases} \quad (2)$$



**Figure 5.** Factor graph illustrating the diarization CRF using talking head information ( $f_{av}$ ) and the context from the OPNs ( $f_{opn}$ ,  $f_{lfb}$ ). The blank circle nodes correspond to hidden variables, the shaded circle nodes correspond to observations while the squares represent the feature functions. The x-axis represents time and the drawing shows also segments corresponding to the time intervals during which a specific observation (track, utterance, OPN) occurs. Red dot segments illustrate the face track temporal segments while blue plain segments the utterances.

255 where  $h(x_{ij}^{av})$  represents the binary output of the SVM classifier introduced in section 3.4. It corresponds  
 256 to 1 when the face and the speaker correspond to the same person and -1 otherwise. We chose a SVM  
 257 classifier since it shows good results in El Khoury et al. (2012); Vallet et al. (2013). Other techniques could  
 258 be employed but we leave this problem for future research.

259 The visual feature function  $f_v(V_i, e_i^v)$ , defined for all face tracks  $V_i \in V$ , indicates how likely the visual  
 260 features  $x_i^{surf}$  of  $V_i$  should be labeled with the person index  $e_i^v$ . This is a face modeling task in which for  
 261 each label  $e_i$ , we need to define a visual model that is learned from the data currently associated to the label.  
 262 Practically,  $f_v$  computes as score between  $V_i$  and a label  $e_i^v$  the 10th percentile SURF vector distances  
 263 between  $x_i^{surf}$  and all the SURF features of the current face tracks associated with this label. The distance  
 264 between two face tracks is computed following (El Khoury et al. (2010)). Although the use of SURF  
 265 features could be discussed regarding other more modern representations, we observe that their matching  
 266 power is useful for similar faces of the same person viewed from a similar view point. The previous  
 267 work in Gay et al. (2014b) uses an average of the distances. By using the percentile, we found a slight  
 268 improvement for the diarization task (0.2 points on the development REPERE corpus). We believe that the  
 269 use of a percentile instead of averaging enables to merge 2 clusters of the same identity but containing  
 270 samples whose poses are dominantly from different poses.

271 The acoustic function  $f_a(A_i, e_i^a)$ , defined over all utterances  $A_i \in A$ , is the audio equivalent of  $f_v$ . We  
 272 chose a 512 GMM-UBM with diagonal covariance following Ben et al. (2004). We did not use ivectors  
 273 since we might need to learn a model on small clusters containing only a few seconds of speech.  $f_a(A_i, e_i^a)$   
 274 computes the likelihood score of the features  $x_i^a$  given the GMM model learned over the data currently  
 275 associated to the cluster label  $e_i^a$ .

276 The LFB feature function is driven by the assumption that faces inside a recurrent LFB are likely to  
 277 correspond to a speaker announced by an OPN. To favor face tracks identified as recurrent LFB to join a

278 person cluster which could be named, we define the following feature function. For each face track  $V_i$ ,

$$f_{\text{lfbv}^k}(V_i, e_i^v) = \begin{cases} 1 & \text{if } x_i^{\text{lfbv}}(k) \text{ and } e_i^v \in \mathcal{E}^{\text{opn}} \\ 0 & \text{otherwise} \end{cases} \quad (3)$$

279 where  $\mathcal{E}^{\text{opn}}$  is the set of person clusters indices co-occurring with an OPN, i.e. the set of clusters which are  
280 currently associated with a name.

281 This principle is extended to each utterance  $A_i$  with the function  $f_{\text{lfbv}^k}$  which employs the feature  $x_i^{\text{lfbv}}(k)$ .  
282 To this end, we assume that the utterances co-occurring with a recurrent LFB should be assigned a cluster  
283 label from the set  $\mathcal{E}^{\text{opn}}$ . Thus, as discussed in section 3.3,  $x_i^{\text{lfbv}}(k)$  is set to true if utterance  $A_i$  is overlapping  
284 with a face track  $V_j$  such that  $x_j^{\text{lfbv}}(k)$  is true. We then introduce the same function as in the video case:

$$f_{\text{lfbv}^k}(A_i, e_i^a) = \begin{cases} 1 & \text{if } x_i^{\text{lfbv}}(k) \text{ and } e_i^a \in \mathcal{E}^{\text{opn}} \\ 0 & \text{otherwise} \end{cases} \quad (4)$$

285

286 Interestingly, these functions act as a namedness feature (Pham et al. (2008)) in the sense that they favor  
287 the naming of the corresponding face tracks and utterances. They also softly constrain the number of  
288 clusters. In other words, the clusters whose labels belong to  $\mathcal{E}^{\text{opn}}$  will attract the segments identified as  
289 recurrent LFB. Note that if the constraint was strictly enforced, each concerned audio or visual segment  
290 would only be assigned to a member of  $\mathcal{E}^{\text{opn}}$ .

291 The OPN feature functions bring a special treatment to the segments co-occurring with OPNs. The idea  
292 is to favor segments (face tracks or utterances) co-occurring with an OPN  $O_j$  to be assigned to a person  
293 cluster likely to be labeled with the name  $x_j^{\text{opn}}$ . Thus, we define:

$$f_{\text{opn\_alone}}(V_i, O_j, e_i^v) = \begin{cases} p(e_i^c = x_j^{\text{opn}} | C, P) & \text{if } V_i \text{ is alone in the image and co-occurs with OPN } O_j \\ 0 & \text{otherwise} \end{cases} \quad (5)$$

294 where  $p(e_i^c = x_j^{\text{opn}} | C, P)$  is the probability that the name contained in the OPN  $O_j$  corresponds to the  
295 cluster label  $e_i^c$  given the clustering  $C$  and the set of OPNs  $P$ . Here, we denote as  $e_i^c$  the naming label of  
296 cluster label  $e_i^v$ . This probability is computed with the naming CRF as defined in section 3.6.

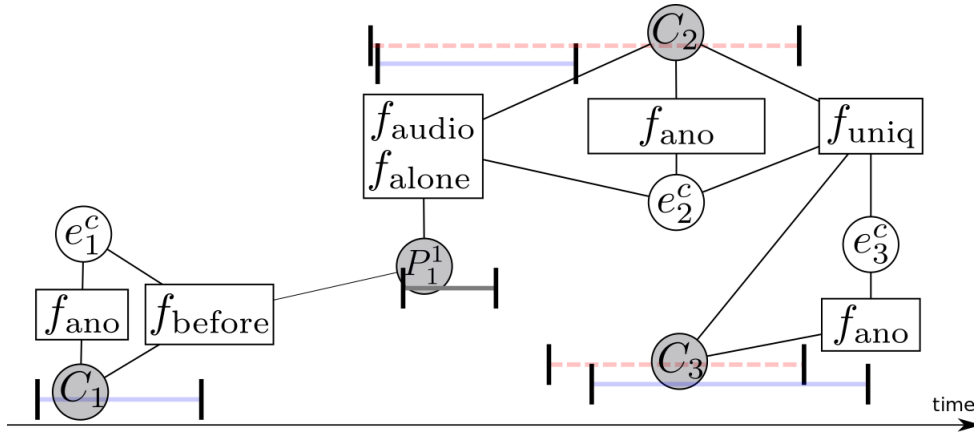
297 Similarly, we use  $f_{\text{opn\_multi}}$  if  $V_i$  co-occurs with other faces:

$$f_{\text{opn\_multi}}(V_i, O_j, e_i^v) = \begin{cases} p(e_i^c = x_j^{\text{opn}} | C, P) & \text{if } V_i \text{ co-occurs with OPN } O_j \\ 0 & \text{otherwise} \end{cases} \quad (6)$$

298 We also define  $f_{\text{opn\_audio}}$  for each co-occurring couple  $(A_i, O_j)$ :

$$f_{\text{opn\_audio}}(A_i, O_j, e_i^a) = \begin{cases} p(e_i^c = x_j^{\text{opn}} | C, P) & \text{if } A_i \text{ co-occurs with OPN } O_j \\ 0 & \text{otherwise} \end{cases} \quad (7)$$

299 Differentiating these 3 cases enables to learn specific  $\lambda$  weights so that the model behavior is adapted to  
300 each situation.



**Figure 6.** Factor graph illustrating the naming CRF using the co-occurrence functions with the OPNs and the uniqueness constraint. The conventions are the same as in Fig 5.

301 *The uniqueness feature function* ensures two faces that co-occur in the same shot to have different  
 302 labels (Berg et al. (2004); Pham et al. (2013)). For such a pair  $V_i, V_j$ :

$$f_{\text{uniq}}(V_i, V_j, e_i^v, e_j^v) = \begin{cases} -\text{Inf} & \text{if } e_i^v = e_j^v \\ 0 & \text{otherwise} \end{cases} \quad (8)$$

303 It is crucial to use this function because due to the OPN feature functions, multiple faces co-occurring  
 304 with the same OPN will tend to be assigned to the same person cluster.

305 **3.6 Cluster identification**

306 The previous diarization CRF provides us a set of AV person clusters  $C = \{C_i, i = 1 \dots N^C\}$ . Thus,  
 307 in the naming step, the goal incorporates estimating the label field  $E^N = \{e_i^c, i = 1 \dots N^C\}$  such that  
 308 the label  $e_i^c$  corresponds to the name of the cluster  $C_i$ . The label  $e_i^c$  takes value in the set of names  $M$   
 309 augmented by an anonymous label which should be assigned to anonymous persons. For this naming CRF,  
 310 the posterior probability uses 6 feature functions:

$$P(E^N | C, O) = \frac{1}{Z(C, O)} \times \exp \left\{ \sum_{i=1}^6 \sum_{c \in G_i} \lambda_i f_i(E_c^N, C_c, O_c) \right\} \quad (9)$$

311 Fig 6 represents an illustration of this. This naming model exploits four different co-occurrence statistics  
 312 between clusters and OPNs. The first function  $f_{\text{alone}}$  is defined over each triplet  $(e_i^c, C_i, O_j)$ , where the  
 313 OPN  $O_j$  must co-occur with a face track which belongs to  $C_i$  and which is alone in the image. Let us  
 314 denote as  $\delta(C_i, O_j)$  the co-occurring time between the face tracks which occurs alone in the cluster  $C_i$  and  
 315  $O_j$ . Then, we have:

$$f_{\text{alone}}(e_i^c, C_i, O_j) = \begin{cases} \frac{\delta(C_i, O_j)}{d_j^{\text{opn}}} & \text{if } x_j^{\text{opn}} = e_i^c \\ 0 & \text{otherwise} \end{cases} \quad (10)$$

316 As for the OPN diarization model components, we define similarly two other functions  $f_{\text{multi}}$  and  $f_{\text{audio}}$   
 317 which measure the overlapping time between  $O_j$  and the face tracks of  $C_i$  which occur with other faces

318 on one hand, and with the audio segments of  $C_i$  on the other hand. Moreover, we exploit the assumption  
 319 that a person does not usually appear or speak before the first apparition of his name in an OPN to define  
 320  $f_{\text{before}}(e_i^c, C_i, O_j)$ , which returns the number of audio segments from cluster  $C_i$  that occur before the first  
 321 apparition of the name  $x_j^{\text{opn}}$  associated to the OPN  $O_j$ .

$$322 \quad f_{\text{before}}(e_i^c, C_i, O_j) = \begin{cases} \#\{A_i \in C_i, \text{end}(A_i) < \text{start}(O_j)\} & \text{if } x_j^{\text{opn}} = e_i^c \\ 0 & \text{otherwise} \end{cases} \quad (11)$$

323 We also introduce prior knowledge over the anonymous label by defining a fifth feature function  
 324  $f_{\text{ano}}(e_i^c, C_i)$  which returns 1 if  $e_i^c$  is the anonymous label. When applied, it allows the model to penalize the  
 325 fact of not identifying a person and improves the recall.

326 Lastly, we define a uniqueness function  $f_{\text{uniq}}(e_i^c, C_i, e_j^c, C_j)$  over visually overlapping clusters just as in  
 327 the diarization step. For each cluster pair  $(C_i, C_j)$  with overlapping face tracks:

$$f_{\text{uniq}}(e_i^c, C_i, e_j^c, C_j) = \begin{cases} -\infty & \text{if } e_i^c = e_j^c \\ 0 & \text{otherwise} \end{cases} \quad (12)$$

### 328 3.7 Optimization

329 The joint use of the two CRFs is conducted by applying the following steps: i) the diarization labels are  
 330 firstly initialized by separately performing audio and video clustering and then associating the clusters to  
 331 obtain the potential AV person labels  $P$  (audio and face cluster couples). The association is conducted  
 332 using the Hungarian algorithm (Kuhn (1955)) where the cost for a cluster couple is defined as the sum of  
 333 the scores from the function  $f_{\text{av}}$  over all its utterance/face track pairs. ii) For each resulting person label  
 334  $p_i$ , biometric models are learned from their associated data and naming probabilities for each label are  
 335 estimated by using the naming CRF. iii) Given these models, we run the loopy belief propagation inference  
 336 to get the most probable diarization labels  $E^d$  by solving  $E^d = \arg \max_{E^d} P(E^d | A, V, O)$ .

337 Eventually, Steps ii) and iii) are iterated in a Expectation-Maximization style by alternating model updates  
 338 and inference. Ideally, one would iterate until convergence, i.e. when the label for each segment becomes  
 339 stable. In practice, as there is no guarantee that the algorithm converges, a fixed number of iterations is  
 340 tuned over the development set since we observe only small modifications after a few iterations.

341 The computational bottleneck with the Loopy Belief Propagation algorithm occurs in the presence of big  
 342 cliques. This is the case, in our graphs, when the uniqueness constraint is applied to images where there are  
 343 more than 20 faces. In such cases, uniqueness constraints can be dropped from the graph during inference  
 344 and enforced in a post-processing step.

## 4 RESULTS AND DISCUSSION

345 This section will firstly present our experimental set-up: the corpus (section 4.1), implementation details  
 346 (section 4.2) and the metrics used for the evaluation (section 4.3). Then in section 4.4, we present our  
 347 results showing identification and clustering performances in function of the different parts of the model.

### 348 4.1 Corpus description

349 We used the REPERE corpus (Giraudel et al. (2012)) for our experiments. It involves broadcast data  
 350 videos containing 4 main types of shows: i) debates in indoor studio (Fig 1a,b); ii) modern format

351 information shows which contain reports and interviews with dynamic picture compositions (Fig 1c,d); iii)  
 352 extracts from parliamentary sessions "Questions to the government" (Fig 1e) iv); celebrity news (Fig 1f).

353 We evaluate our approach on the final test set which contains 37 hours during which 10 are annotated. A  
 354 development set is used to optimize the number of LFB functions and the number of iterations between  
 355 the two CRFs. It consists of 28 hours among which 6 are annotated. The SVM  $h$  used in the  $f_{av}$  feature  
 356 function and the CRF parameters are learned on the test set of the first REPERE evaluation composed of 3  
 357 hours of annotated data.

## 358 4.2 Parameter settings and algorithm details

359 We set the K value to  $\{3, 4, 5\}$  for the LFB feature functions. We set the number of iterations between the  
 360 two CRFs to 3 as we noticed that no major changes usually occur after that point. It is important to note  
 361 that these CRF parameters are learned on automatic detections and automatic clusters and not on cleanly  
 362 segmented ones. Therefore, it enables us to take into account the noise present at test time. We use the  
 363 GRMM toolbox (McCallum" (2002)) for the CRF implementation.

364 The initial speaker diarization system is the LiumSpkDiarization toolbox<sup>1</sup> which combines ivector  
 365 representation and ILP clustering (Rouvier et al. (2013)). It has achieved state-of-the-art results in several  
 366 speaker diarization benchmarks (Rouvier and Meignier (2012)). The initial face diarization uses the system  
 367 described in (Khoury et al. (2013)) which combines SURF based distances and DCT features whose  
 368 distribution is modeled with GMMs. This system has been evaluated on the public Buffy dataset (Cinbis  
 369 et al. (2011)) and compares favorably to other metric learning methods. The use of state-of-the-art systems  
 370 enables us to verify that our CRF is able to correct errors which are proven difficult to solve in the  
 371 monomodal case.

## 372 4.3 Performance measures

373 The overall identification performance is measured with the Estimated Global Error Rate (EGER) which  
 374 is the REPERE evaluation metric. It is defined as follows:

$$\text{EGER} = \frac{\# \text{ conf} + \# \text{ miss} + \# \text{ false}}{\# \text{ total}} \quad (13)$$

375 where  $\# \text{ conf}$  is the number of wrongly identified persons,  $\# \text{ miss}$ , the number of missed persons,  $\# \text{ false}$ , the  
 376 number of false alarms and  $\# \text{ total}$ , the total number of persons to be detected. It should be noted that the  
 377 metric ignores the spatial position of the faces and simply uses a person list for each annotated image. The  
 378 behavior of this metric is illustrated on figure 7. Wrong predictions are counted as false alarms only if the  
 379 number of predictions exceeds the number of persons in the annotation. Otherwise, they are counted as  
 380 confusions. Similarly, missing persons are reported only if the number of predictions is smaller than the  
 381 number of persons.

382 We also use the clustering error rate (CER) to study the correlation between clustering and identification  
 383 performances as our work is motivated by an interdependence between those two tasks. Initially, the CER  
 384 has been introduced for the speaker clustering task (NIST (2003)) and is defined as:

$$\text{CER} = \frac{\sum_{\text{seg} \in \text{Segs}} \text{dur}(\text{seg})(\min(N_{\text{Ref}}(\text{seg}), N_{\text{Sys}}(\text{seg})) - N_{\text{Correct}}(\text{seg}))}{\sum_{\text{seg} \in \text{Segs}} \text{dur}(\text{seg})N_{\text{ref}}(\text{seg})} \quad (14)$$

<sup>1</sup> <http://www-lium.univ-lemans.fr/diarization>

Show: BFMStory\_12 frame: 4312  
 Head Ref: Barack\_OBAMA Augusta\_ADA\_KING  
 Head Hyp: Augusta\_ADA\_KING David\_HAMILTON Alan\_TURING

**Figure 7.** Extract of an evaluation file for face identification. The second row is the reference name list and the third row is the predicted list. Augusta\_ADA\_KING will be counted as correct. One of the two remaining names will be counted as confusion with Barack\_OBAMA, and the third one will be a false alarm. Since there are 2 persons in the reference and the system made 2 errors, the corresponding EGER of this example is 1.

Speaker diarization results			
	Initial monomodal	CRF Dia	CRF Dia without OPNs
News	6.9%	7.0%	<b>6.8%</b>
Debates	6.6%	<b>4.0%</b>	6.5%
Parliament	6.9%	<b>5.0%</b>	9.5%
Celebrity	<b>14.6%</b>	15.1%	<b>14.6%</b>
All	7.4	<b>6.8%</b>	7.4%

Face diarization results			
	Initial monomodal	CRF Dia	CRF Dia without OPNs
News	<b>4.8%</b>	5.4%	5.9%
Debates	4.6%	<b>1.9%</b>	4.4%
Parliament	11.2%	<b>10.4%</b>	13.7%
Celebrity	<b>3.5%</b>	7.9%	6.4%
All	5.2%	<b>5.0%</b>	6.1%

**Table 1.** Speaker and face diarization performances in terms of CER. The first column presents the initial monomodal systems Khoury et al. (2013); Rouvier et al. (2013). The second one is the diarization CRF presented in this paper. The third one is the same as the second one, however, we remove the OPN related functions  $f_{\text{fb}}$  and  $f_{\text{opn}}$ .

385 where the audio file is divided in continuous segments at each speaker change and:

- 386 •  $\text{dur}(\text{seg})$  is the duration of the segment  $\text{seg}$ .
- 387 •  $N_{\text{Ref}}(\text{seg})$  is the number of active speakers during segment  $\text{seg}$ .
- 388 •  $N_{\text{Sys}}(\text{seg})$  is the number of speakers detected by the system.
- 389 •  $N_{\text{Correct}}(\text{seg})$  is the number of speakers correctly detected by the system. A match needs to be made
- 390 between the clusters and the speaker references in order to compute this term.

391 We applied this measure to the face clustering task. With the audio CER, a detected speech segment is

392 matched to a reference during their temporal overlap. The only modification to tackle visual modality is

393 that a face detection must have a temporal AND spatial overlap to be matched with a reference. In addition,

394 note that we do not consider false alarms and missed detections that are usually considered in NIST to

395 compare the effects of the different systems since the only error that changes with methods given the setup

396 (fixed face tracks and utterances) is due to the final clustering of the face and speech segments. Thus, miss

397 detections and false alarms are identical.

#### 398 4.4 Identification and clustering results with the CRF combination

399

400 **Diarization results:** we first describe the diarization results presented in Table 1. We can see that the full

401 CRF model has a slightly lower error rate over the whole corpus than the initial monomodal systems (6.8%

402 vs 7.4% for the speakers and 5.0% vs 5.2% for the faces). On the other hand, the performances depend

function	$f_a$	$f_v$	$f_{av}$	$f_{\text{fbv}^k}$	$f_{\text{fba}^k}$
$\lambda$	$\lambda_a$	$\lambda_v$	$\lambda_{av}$	$\lambda_{\text{fbv}^k}$	$\lambda_{\text{fba}^k}$
$\lambda$ value	0.4	1.8	0.2	1.9	1.7

**Table 2.** The  $\lambda$  parameter values for some of the feature functions used by the diarization CRF. For  $\lambda_{\text{fbv}^k}$  and  $\lambda_{\text{fba}^k}$ , the value of  $k$  is 5, which is the highest parameter value. It corresponds to the most common case as 90% of the segments are inside background clusters which contain more than 5 elements.

403 strongly on the type of shows. For instance, an important part of the global improvement comes from the  
 404 debate videos (4.0% vs 6.6% for the speakers and 1.9% vs 4.9% for the faces). In debates, most of the  
 405 scenes are in the same studio thereby reducing the visual variability of the background image and most  
 406 of the persons present are speakers announced by an OPN. Thus, most faces and utterances are featured  
 407 as recurrent (i.e.  $x_i^{\text{fbv}}$  is set to true) and the  $f_{\text{fb}}$  functions have a positive impact on the diarization. They  
 408 enable to solve clustering confusion errors by constraining the number of clusters toward the number of  
 409 detected OPNs. Indeed, if we remove the OPN related functions  $f_{\text{fb}}$  and  $f_{\text{opn}}$  (cf third column of the  
 410 Table 1), most of the improvements are lost. It appears that the use of multimodality does not help to correct  
 411 clustering errors. This is somewhat surprising as past works (Gay et al. (2014c)) reports improvements  
 412 in the audio modality with this very system on the same type of data. The difference with this previous  
 413 work is that our initial monomodal speaker diarization system has become much more efficient, essentially  
 414 thanks to a careful selection of the data used to train the generic speaker model UBM. This way, there are  
 415 much fewer errors to correct.

416 In the case of celebrity magazines, the diarization CRF increases the error rate (15.1% vs 14.6% for  
 417 the speakers and 7.9% vs 3.5% for the faces). Those videos contain very few OPNs and essentially short  
 418 outdoor scenes. Thus, the  $f_{\text{fb}}$  functions cannot help the CRF to take appropriate decisions. Moreover,  
 419 previous experiments reported in Gay et al. (2014c) showed that the use of the biometric person models  
 420 inside the CRF framework appears to be less efficient than when it is used in the hierarchical monomodal  
 421 systems.

422 The importance of the OPN related functions is also visible if we consider the  $\lambda$  parameters learned by  
 423 the CRF in Table 2. During training, the weight  $\lambda_{av}$  are indeed set to a relatively low value as compared to  
 424 the other terms (although those values are ponderated by the amplitude of the feature functions). We have  
 425 found that for a majority of segments, the  $f_{\text{fb}}$  function is dominant. This is further illustrated in table 4.  
 426

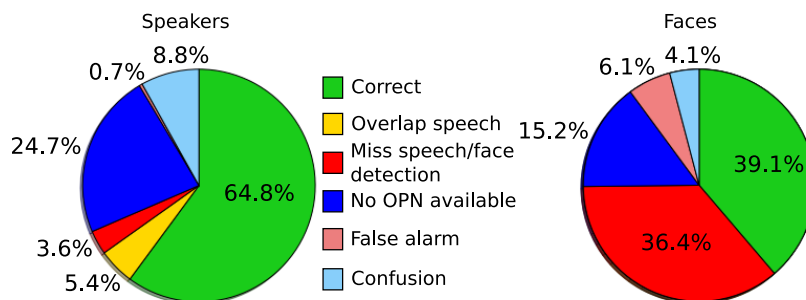
427 **Identification results:** We now turn to the identification results reported in Table 3. We compare 3 systems:  
 428 we denote by  $N$  the naming CRF applied on top of the initial monomodal diarizations described in Khoury  
 429 et al. (2013); Rouvier et al. (2013),  $N + D$  is the joint use of the naming and the diarization CRF, and  
 430 the last one is an oracle. Note that the oracle still produces errors, since, as we deal with automatic face  
 431 detection and tracking, there are errors that a perfect clustering and naming cannot correct: false alarms,  
 432 missed faces and face tracks for which the identity is not introduced by an OPN (see more about this in  
 433 Fig 8). Adding the diarization CRF permits to globally reduce the error rates in both modalities (31.4% vs  
 434 33.4% for the speakers and 52.2% vs 54.5% for the faces), especially for debate and parliament videos.  
 435 This is not surprising as we previously showed that the diarization CRF have less confusion errors for  
 436 studio scenes than the initial monomodal systems.

437 Regarding news videos, although we saw that clustering confusion errors were not reduced globally,  
 438 the use of the diarization CRF also improves the identification. This is probably due to the correction of  
 439 confusion errors in studio scenes which have a greater impact on the identification than errors concerning  
 440 anonymous persons in reports.



	Audio			Visual		
	$N$	$N + D$	Oracle	$N$	$N + D$	Oracle
News	31.6%	30.8%	25.7%	58.2%	56.4%	37.7%
Debates	18.0%	14.0%	11.3%	42.0%	38.0%	35.6%
Parliament	11.3%	8.7%	5.2%	62.2%	59.6%	47.4%
Celebrity	85.6%	85.8%	82.1%	83.9%	86.6%	75.3%
All	33.4%	31.4%	27.2%	54.5%	52.2%	40.2%

**Table 3.** Identification performances measured in EGER. The system  $N$  is the naming CRF on top of the monomodal diarizations and the system  $N + D$  is the naming and diarization CRF combination.



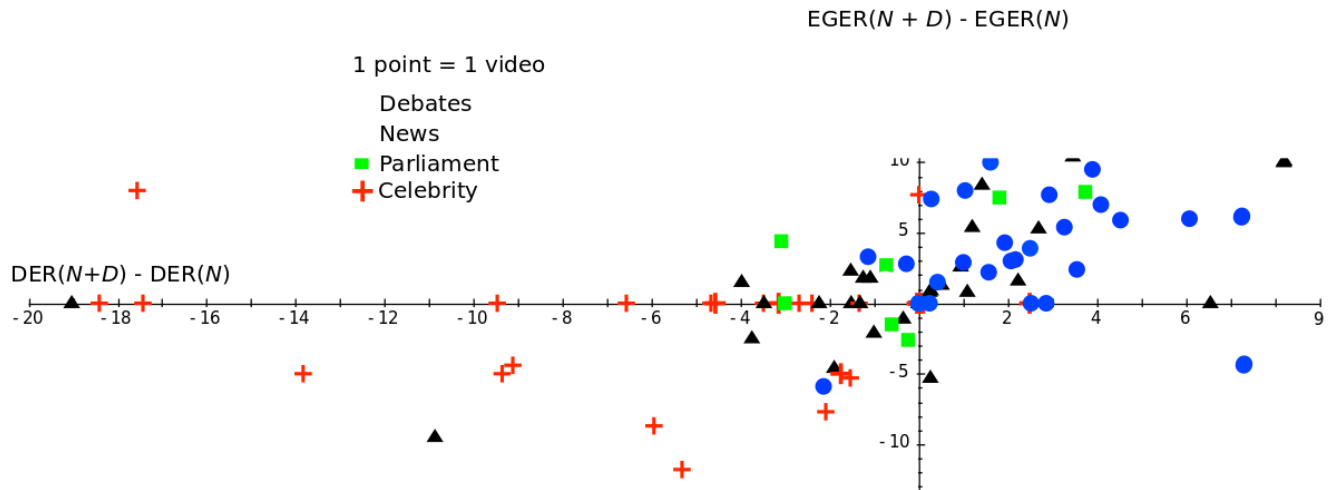
**Figure 8.** Different errors for the speaker (left) and face (right) identification tasks. Percentages are expressed relatively to the number of annotations.

441 The structure of celebrity magazines differs from the other shows as it contains very few OPNs and  
 442 recurrent LFB. In those cases, the diarization CRF degrades both diarization and identification performances.  
 443 We design an oracle on the diarization and the identification to measure the potential improvements. It uses  
 444 automatic face/speech segment detections and automatic OPN extraction. Then, the association between  
 445 these segments and the OPNs is done with the manual reference. Thus, the errors made by the oracle  
 446 correspond to missing OPNs or missing segment detections. In the case of celebrity shows, with an error  
 447 rate of 75.3%, the OPN-based approach is clearly not suitable.

448  
 449 **Error analysis:** the proportion of the different error types can be visualized globally on the pie charts in  
 450 Fig 8. Regarding the speaker identification task, the lack of OPNs explains most of the errors as 24.7% of  
 451 the annotated persons are not announced, most of them being journalists. As for the faces, the detection  
 452 step is more crucial as 36.4% of the persons faces are not detected. This corresponds usually to profile  
 453 faces or persons seen from the back. Most of the false alarms are anonymous persons incorrectly identified.

454 We also illustrate the correlation between diarization and identification performances in figure 9. We  
 455 plot the performance differences for each video between the full system ( $N + D$ ) and the CRF naming  
 456 alone ( $N$ ). We observe that they are unique to their type of show. The debate videos appear in the top-right  
 457 part of the plane, which means that the diarization CRF improves the diarization and the identification.  
 458 Concerning news and parliament videos, the correlation between CER and DER is not as strong. The  
 459 presence of anonymous persons and off voices imply that a change in the diarization does not necessarily  
 460 correspond to a change in identification performances.

461 Finally, the table 4 shows the performance of the model when adding the different components of the  
 462 diarization CRF one by one. If we focus on the first and second lines, we see that the CRF with only 3  
 463 feature functions degrades the performances compared to the monomodal diarizations. We find that, used  
 464 alone, the monomodal representations present in the CRF (see the  $f_a$  and  $f_v$  functions) do not compare  
 465 favourably with the monomodal diarization frameworks. This could be improved in a future work by using  
 466 better person representations. However, each other component enables to reduce the error rate and the



**Figure 9.** The Y-axis is the EGER difference between the CRF combination and the naming CRF alone measured for the faces. The X-axis is the DER difference between the diarization CRF and the initial monomodal face diarization Khoury et al. (2013).

	Audio	Visual
$N$	33.4%	54.5%
$N + D (f_a + f_v + f_{av})$	34.1%	56.2%
$N + D (f_a + f_v + f_{av} + f_{opn})$	33.9%	56.4%
$N + D (f_a + f_v + f_{av} + f_{opn} + f_{uniq})$	33.9%	55.6%
$N + D (f_a + f_v + f_{av} + f_{opn} + f_{uniq} + f_{lfb})$	31.4%	52.7%

**Table 4.** Contribution of the different diarization model components on the naming task (results in EGER). As in table 3, the system  $N$  is the naming CRF with the monomodal diarizations. The other lines correspond to the combination of the naming and the diarization CRF, using as feature functions in the diarization CRF those given in parenthesis.

467 full model provides the best performances. It should also be noticed that, although it might generates big  
 468 cliques in some cases, the uniqueness function is essential to benefit from the  $f_{opn}$  feature functions. If not  
 469 applied, an OPN will be propagated to all the faces overlapping with him.  
 470

471 **Comparison with state of the art and discussion:** on the same dataset, the system described in Bechet  
 472 et al. (2014) obtains an EGER of 30.9% for the speakers and 39.4% for the faces. Thus, it proves to  
 473 have a better performance especially regarding the faces. This is possible with the help of pre-trained  
 474 models for each show which enable to indicate how many faces should be present on screen and what their  
 475 roles are. For instance, when it detects the configuration shown in Fig 1c, it deduces that the announced  
 476 guest is present on the right even if no faces have been detected. In fact, this approach does not even  
 477 use a face diarization module. However, it requires a large amount of learning and *a priori* information.  
 478 By comparison, our method is much simpler to implement, especially since it has better generalization  
 479 capabilities, we learn one single model over a large and diverse corpus, and what is more, it requires less  
 480 annotations if we need to process a new type of show.

481 The constrained hierarchical clustering detailed in Poignant et al. (2015) obtains an EGER of 35.9% for  
 482 the speakers and 44.3% for the faces. Compared with our system, it has better performances on the faces,  
 483 but worst for the speakers. As we do, they only rely on OPNs without other specific supervised information  
 484 on the show. According to their paper, it seems that their constrained multimodal clustering, which avoids  
 485 clustering together faces which co-occur with different OPN names, is one of the contributions which  
 486 improves results and that we do not use, and could explain the difference. Nevertheless, the influence of

487 each pre-processing (speaker and face detections, monomodal clusterings and OPN detection) makes it  
488 hard to analyse the performance difference.

## 5 CONCLUSION

489 In this paper, we presented our contribution for AV person diarization and identification from OPNs. Our  
490 system uses an iterative combination of 2 CRFs. One performing the AV diarization at a person level, and a  
491 second one associating the names and the clusters. Several context modeling cues are used to solve the  
492 person/name association problem and the diarization issues. While it is clear that more supervised learning  
493 and *a priori* information on the context can improve the performances, our approach provides an interesting  
494 trade-off between performance on one hand and generalization/low annotation cost on the other hand. The  
495 principal contextual cue consists in the face image background. It allows us to distinguish the faces and the  
496 speakers which are announced by OPNs and guide the clustering accordingly.

497 In this work, we did not address the issue of non-frontal face detection. As a short term perspective,  
498 it would be interesting to increase the recall of the face detector, for instance, by adding a profile view  
499 detector. This would render the face clustering task more challenging and the potential benefit from context  
500 modeling would be greater. Secondly, our context modeling assumes that speakers are announced by an  
501 OPN the first time they talk. For the REPERE dataset, this is the case. However, this assumption could  
502 be sensible to broadcaster's editing policies. Actually, the optimal choice of the context for unsupervised  
503 person identification is a difficult problem if we want to avoid the need for specific annotations for each  
504 show. One solution to consider is to learn the setting of each show or a part of the setting from a corpus in  
505 an unsupervised way.

## REFERENCES

- 506 Bauml, M., Tapaswi, M., and Stiefelhagen, R. (2013). Semi-supervised learning with constraints for person  
507 identification in multimedia data. In *Proceedings of the IEEE Conference on Computer Vision and*  
508 *Pattern Recognition*, pages 3602–3609, Portland, USA.
- 509 Bechet, F., Bendris, M., Charlet, D., Damnati, G., Favre, B., Rouvier, M., Auguste, R., Bigot, B., Dufour,  
510 R., Fredouille, C., Linares, G., Martinet, J., Senay, G., and Tirilly, P. (2014). Multimodal understanding  
511 for person recognition in video broadcasts. In *Proceedings of InterSpeech*, pages 146–151, Singapore.
- 512 Ben, M., Betsler, M., Bimbot, F., and Gravier, G. (2004). Speaker diarization using bottom-up clustering  
513 based on a parameter-derived distance between adapted GMMs. In *Proceedings of International*  
514 *Conference on Spoken Language Processing*, pages 523–538, Jeju, Island.
- 515 Berg, T. L., Berg, A. C., Edwards, J., Maire, M., White, R., Teh, Y. W., Learned-Miller, E. G., and Forsyth,  
516 D. A. (2004). Names and faces in the news. In *Proceedings of the IEEE Conference on Computer Vision*  
517 *and Pattern Recognition*, volume 2, pages 836–848, Washington, USA.
- 518 Bhattarai, B., Sharma, G., Jurie, F., and Pérez, P. (2014). Some faces are more equal than others:  
519 Hierarchical organization for accurate and efficient large-scale identity-based face retrieval. In  
520 *Proceedings of the European Conference on Computer Vision*, pages 160–172, Zurich, Switzerland.
- 521 Bredin, H. and Poignant, J. (2013). Integer linear programming for speaker diarization and cross-modal  
522 identification in TV broadcast. In *Proceedings of InterSpeech*, pages 49–54, Lyon, France.
- 523 Bredin, H., Poignant, J., Tapaswi, M., Fortier, G., Le, V., Napoleon, T., Gao, H., Barras, C., Rosset, S.,  
524 Besacier, L., et al. (2013). QCOMPERE@REPERE 2013. In *Workshop on Speech, Language and Audio*  
525 *in Multimedia*, pages 49–54, Marseille, France.

- 526 Chen, D. and Odobez, J.-M. (2005). Video text recognition using sequential monte carlo and error voting  
527 methods. *Pattern Recognition Letters*, 26(9):1386–1403.
- 528 Cinbis, R. G., Verbeek, J., and Schmid, C. (2011). Unsupervised metric learning for face identification in  
529 TV video. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1559–1566,  
530 Barcelona, Spain.
- 531 Cour, T., Sapp, B., Nagle, A., and Taskar, B. (2010). Talking pictures: Temporal grouping and dialog-  
532 supervised person recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern  
533 Recognition*, pages 1014–1021, San Francisco, USA. IEEE.
- 534 Cour, T., Sapp, B., and Taskar, B. (2011). Learning from partial labels. *The Journal of Machine Learning  
535 Research*, 12:1501–1536.
- 536 El Khoury, E., Senac, C., and Joly, P. (2010). Face-and-clothing based people clustering in video content.  
537 In *Proceedings of the International Conference on Multimedia Information Retrieval*, pages 295–304,  
538 Philadelphia, USA.
- 539 El Khoury, E., Sénac, C., and Joly, P. (2012). Audiovisual diarization of people in video content. *Multimedia  
540 Tools and Applications*, 68(3):747–775.
- 541 Everingham, M., Sivic, J., and Zisserman, A. (2006). Hello! my name is... buffy—automatic naming of  
542 characters in tv video. In *Proceedings of the British Machine Vision Conference*, volume 2, pages  
543 2365–2371, Edinburgh, Scotland.
- 544 Gay, P., Dupuy, G., Lailler, C., Odobez, J.-M., Meignier, S., and Deléglise, P. (2014a). Comparison of  
545 two methods for unsupervised person identification in TV shows. In *Proceedings of The Content Based  
546 Multimedia Indexing Workshop*, pages 1–6, Klagenfurt, Austria.
- 547 Gay, P., Elie, K., Sylvain, M., Jean-Marc, O., and Paul, D. (2014b). A conditional random field approach  
548 for face identification in broadcast news using overlaid text. In *Proceedings of the IEEE International  
549 Conference on Image Processing*, pages 318–322.
- 550 Gay, P., Khoury, E., Meignier, S., Odobez, J.-M., and Deleglise, P. (2014c). A conditional random field  
551 approach for audio-visual people diarization. In *Proceedings of the IEEE International Conference on  
552 Acoustics, Speech and Signal Processing*, pages 116–120, Florence, Italy.
- 553 Giraudel, A., Carré, M., Mapelli, V., Kahn, J., Galibert, O., and Quintard, L. (2012). The REPERE  
554 corpus: a multimodal corpus for person recognition. In *Proceedings of The International Conference on  
555 Language Resources and Evaluation*, pages 1102–1107, Istanbul, Turkey.
- 556 Guillaumin, M., Verbeek, J., and Schmid, C. (2010). Multiple instance metric learning from automatically  
557 labeled bags of faces. In *Proceedings of the European Conference on Computer Vision*, pages 634–647.  
558 Crete, Greece.
- 559 Jou, B., Li, H., Ellis, J. G., Morozoff-Abegauz, D., and Chang, S. (2013). Structured exploration of who,  
560 what, when, and where in heterogeneous multimedia news sources. In *Proceedings of ACM International  
561 conference on Multimedia*, pages 357–360, Barcelona, Spain.
- 562 Jousse, V., Petit-Renaud, S., Meignier, S., Esteve, Y., and Jacquin, C. (2009). Automatic named  
563 identification of speakers using diarization and ASR systems. In *Proceedings of the IEEE International  
564 Conference on Acoustics, Speech and Signal Processing*, pages 4557–4560, Taipei, Taiwan.
- 565 Khoury, E., Gay, P., and Odobez, J. (2013). Fusing matching and biometric similarity measures for face  
566 diarization in video. In *Proceedings of the IEEE International Conference on Multimedia Retrieval*,  
567 pages 97–104, Dallas, USA.
- 568 Kuhn, H. W. (1955). The hungarian method for the assignment problem. *Naval research logistics quarterly*,  
569 2:83–97.

- 570 Li, D., Wei, G., Sethi, I. K., and Dimitrova, N. (2001). Person identification in tv programs. *Journal of*  
571 *Electronic Imaging*, 10(4):930–938.
- 572 Ma, C., Nguyen, P., and Mahajan, M. (2007). Finding speaker identities with a conditional maximum  
573 entropy model. In *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal*  
574 *Processing*, volume 4, pages 253–261, Honolulu, USA.
- 575 McCallum”, A. K. (2002). ”MALLET: A machine learning for language toolkit”.  
576 ”<http://mallet.cs.umass.edu>”.
- 577 NIST (2003). The rich transcription spring 2003 (rt-03s) evaluation plan.
- 578 Noulas, A., Englebienne, G., and Krose, B. J. (2012). Multimodal speaker diarization. *IEEE Transactions*  
579 *on Pattern Analysis and Machine Intelligence*, 34(1):79–93.
- 580 Ozkan, D. and Duygulu, P. (2010). Interesting faces: A graph-based approach for finding people in news.  
581 *Pattern Recognition*, 43(5):1717–1735.
- 582 Pham, P., Moens, M.-F., and Tuytelaars, T. (2008). Linking names and faces: Seeing the problem in  
583 different ways. In *Proceedings of the European Conference on Computer Vision*, pages 68–81, Marseille,  
584 France.
- 585 Pham, P. T., Deschacht, K., Tuytelaars, T., and Moens, M.-F. (2013). Naming persons in video: Using  
586 the weak supervision of textual stories. *Journal of Visual Communication and Image Representation*,  
587 24(7):944–955.
- 588 Poignant, J., Besacier, L., and Quénot, G. (2014). Unsupervised speaker identification in TV broadcast  
589 based on written names. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*,  
590 23(1):57–68.
- 591 Poignant, J., Fortier, G., Besacier, L., and Quénot, G. (2015). Naming multi-modal clusters to identify  
592 persons in tv broadcast. *Multimedia Tools and Applications*, pages 1–25.
- 593 Rouvier, M., Dupuy, G., Gay, P., Khoury, E., Merlin, T., and Meignier, S. (2013). An open-source  
594 state-of-the-art toolbox for broadcast news diarization. In *Proceedings of InterSpeech*, pages 547–552,  
595 Lyon, France.
- 596 Rouvier, M. and Meignier, S. (2012). A global optimization framework for speaker diarization. In  
597 *Proceedings of the Odyssey workshop*, pages 546–552, Singapore.
- 598 Satoh, S., Nakamura, Y., and Kanade, T. (1999). Name-it: Naming and detecting faces in news videos.  
599 *IEEE Multimedia*, 6(1):22–35.
- 600 Schroff, F., Kalenichenko, D., and Philbin, J. (2015). Facenet: A unified embedding for face recognition  
601 and clustering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*,  
602 pages 815–823, Boston, USA.
- 603 Simonyan, K., Parkhi, O. M., Vedaldi, A., and Zisserman, A. (2013). Fisher vector faces in the wild. In  
604 *Proceedings of the British Machine Vision Conference*, pages 867–879, Bristol, United Kingdom.
- 605 Tapaswi, M., Parkhi, O. M., Rahtu, E., Sommerlade, E., Stiefelhagen, R., and Zisserman, A. (2014).  
606 Total cluster: A person agnostic clustering method for broadcast videos. In *Proceedings of the Indian*  
607 *Conference on Computer Vision Graphics and Image Processing*, pages 7–15, Bengaluru, India.
- 608 Vallet, F., Essid, S., and Carrive, J. (2013). A multimodal approach to speaker diarization on TV talk-shows.  
609 *IEEE Transactions on Multimedia*, 15(3):509–520.
- 610 Viola, P. and Jones, M. J. (2004). Robust real-time face detection. *International Journal of Computer*  
611 *Vision*, 57(2):137–154.
- 612 Wohlhart, P., Köstinger, M., Roth, P. M., and Bischof, H. (2011). Multiple instance boosting for face  
613 recognition in videos. *Pattern Recognition*, 6835:132–141.

- 614 Zhang, L., Kalashnikov, D. V., and Mehrotra, S. (2013). A unified framework for context assisted face  
615 clustering. In *Proceedings of the IEEE International Conference on Multimedia Retrieval*, pages 9–16,  
616 Dallas, USA.
- 617 Zhang, N., Paluri, M., Tagiman, Y., Fergus, R., and Bourdev, L. (2015). Beyond frontal faces: Improving  
618 person recognition using multiple cues. *Proceedings of the IEEE Conference on Computer Vision and*  
619 *Pattern Recognition*, pages 4804–4813.