

# Active Speaker Detection as a Multi-Objective Optimization with Uncertainty-based Multimodal Fusion

Baptiste Pouthier<sup>1,2</sup>, Laurent Pilati<sup>1</sup>, Leela K. Gudupudi<sup>1</sup>, Charles Bouveyron<sup>2</sup>, Frederic Precioso<sup>2</sup>

<sup>1</sup>NXP Semiconductors, France

<sup>2</sup>Université Côte d’Azur, Inria, CNRS, LJAD, I3S, Maasai, France

{baptiste.pouthier, laurent.pilati, leela.k.gudupudi}@nxp.com  
{charles.bouveyron, frederic.precioso}@univ-cotedazur.fr

## Abstract

It is now well established from a variety of studies that there is a significant benefit from combining video and audio data in detecting active speakers. However, either of the modalities can potentially mislead audiovisual fusion by inducing unreliable or deceptive information. This paper outlines active speaker detection as a multi-objective learning problem to leverage best of each modalities using a novel self-attention, uncertainty-based multimodal fusion scheme. Results obtained show that the proposed multi-objective learning architecture outperforms traditional approaches in improving both mAP and AUC scores. We further demonstrate that our fusion strategy surpasses, in active speaker detection, other modality fusion methods reported in various disciplines. We finally show that the proposed method significantly improves the state-of-the-art on the AVA-ActiveSpeaker dataset.

**Index Terms:** active speaker detection, audiovisual, multimodal fusion, multi-objective

## 1. Introduction

Active Speaker Detection (ASD) contemplates on identifying active speakers in a video by analyzing both visual and audio features. Hence, ASD is inherently multimodal in nature, where video and audio data are essential attributes. In recent years there has been considerable interest in the ASD methods based upon audio-visual cues [1, 2, 3, 4, 5]. Despite this interest, the lack of large-scale in-the-wild datasets impeded scientific progress in the field. Unfortunately, most of the prior works are challenged by skewed results owing to the poor quality of the considered datasets.

The recent AVA-ActiveSpeaker dataset [6] can potentially overcome these limitations and reshape the ASD field of study. Together with the large in-the-wild dataset, the authors presented a baseline model based on a two-stream network that merges video and audio modalities in an end-to-end fashion. Within the annual AVA-ActiveSpeaker challenge [7], Chung [8] and Zhang et al. [9] improved this baseline using hybrid 3D-2D CNNs pre-trained on large-scale multimodal datasets [10]. Unfortunately, this approach encounters two major challenges in practice, as illustrated in Fig.1: (1) multi-speaker scenario where multiple persons in a video frame are speaking, and (2) low-resolution and/or indiscernible faces in video frames. In [11], Alcàzar et al. addressed the multi-speakers scenario by learning long-term relationships between speakers, and Huang et al. [12] handled the uncertainty in video modality by adding optical flow to raw pixel representation to strengthen face characterisation. Nevertheless, these studies focused on ad-hoc objectives with limited scope and little attention has been paid to the aggregated approaches which exploit correlation between both the challenges. In this paper, we propose to learn this correlation using a cross-modal fusion that involves simultaneous learning of the uncertainty in both the

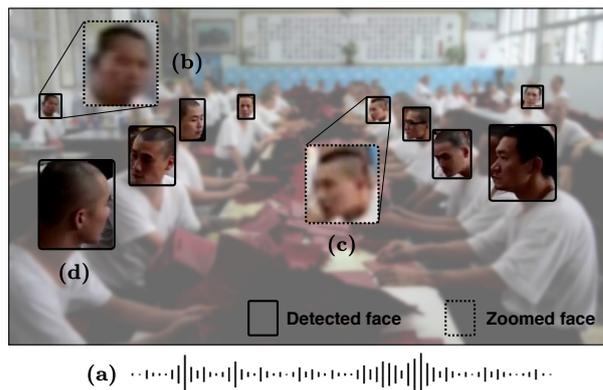


Figure 1: This is an illustration of how ambiguous both video and audio modalities are within the ASD problem. (a) represents the scene’s audio track that contains speech. But the uncertainty here is who is speaking? It is difficult to say it from video also, because some face resolutions are insufficient (b,c) or characters lips may be partially (c) or entirely (d) concealed.

modalities using a multi-objective learning scheme as described in Section 2.1.

Traditionally, in multi-task learning, uncertainty is handled by learning adaptive weighting between each task’s loss functions [13, 14]. In the field of Automatic Video Description, researchers have investigated attention-based [15] fusion mechanisms that capture the importance of each of the modalities. Hore et al. [16] proposed a model that selectively uses features from different modalities. This approach was later improved using hierarchical attention fusion in [17]. Despite the growing interest of these fusion mechanisms in other disciplines, to the best of our knowledge, no studies have been conducted on ASD problem. The present paper aims to investigate different fusion schemes to solve the ASD problem. We also propose a novel audiovisual fusion mechanism as a first attempt to enrich the self-attention model with uncertainty information to disentangle the practical challenges.

The proposed method consists of learning a self-attention [15] and uncertainty-based fusion mechanism that weights video and audio embeddings in an end-to-end fashion. We use the typical two-stream DNN architecture from the literature [6, 11], with a stream per each modality, to encode both embeddings. By characterizing uncertainty of each modality, we aim at disentangling the problem presented in Fig.1. Indeed, the fusion scheme we propose determines the viability of each modality at any given time to learn a comprehensive understanding of every situation towards ASD disambiguation. Our solution significantly outperforms state-of-the-art in the AVA-ActiveSpeaker dataset by 4.8% and 1.7% on validation and test sets, respectively.



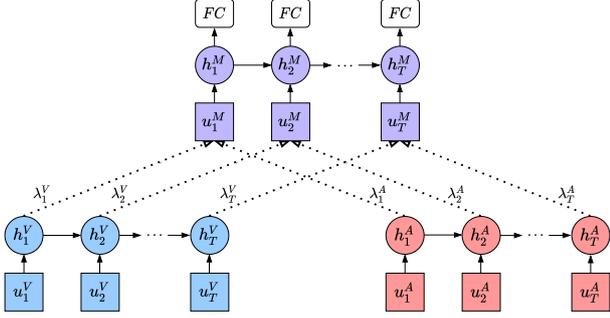


Figure 4: Multimodal fusion scheme where video and audio embedded representations are merged. FC represents a 2-dim Fully Connected layer with a softmax function.

probabilities of target over the entire training data. Five bins are used to pre-process the number of detected faces (four bins for  $\eta \leq 4$  plus a bin handling  $\eta \geq 5$ ). Face-size values are discretized within eight bins with first and last bins being left-open and right-open intervals, respectively.

**Auxiliary Uncertainty:** We leverage auxiliary classifiers predictions to estimate the uncertainty in video and audio embeddings. Intuitively, the data leading to a highly polarized auxiliary decision can be considered as reliable. The output of a softmax layer cannot be used reliably as true probabilities [19, 20, 21]. Therefore, we use a simple yet effective workaround called *temperature scaling* [20] to generate relevant auxiliary output scores using a calibrated softmax. We denote  $\delta_t^V$  and  $\delta_t^A$  the uncertainty values linked to the video and audio auxiliary predictions, respectively. Given  $\tilde{y}_t^{\{V,A\}}$  the temperature-scaled probabilities of the video or audio auxiliary output, we define  $\delta_t^{\{V,A\}} = 2(\max \tilde{y}_t^{\{V,A\}} - 0.5)$ .

**Self-Attention:** We aim at using dot-product global self-attention [15] mechanism to evaluate how important video and audio embedded representations are within their local temporal context. Given  $H \in \mathbb{R}^{T \times 128}$  as defined in Eq.3, we compute the self-attention scores  $S^{\{V,A\}} \in \mathbb{R}^{T \times T}$  using Eq.4.

$$H^{\{V,A\}} = [h_1^{\{V,A\}}, \dots, h_T^{\{V,A\}}] \quad (3)$$

$$S^{\{V,A\}} = \text{softmax}(H^{\{V,A\}} H^{\{V,A\}T} + B) \quad (4)$$

$B \in \mathbb{R}^{T \times T}$  is the bias matrix that limits the temporal context: the future timestamps are masked to preserve self-attention causality, and the distant past timestamps are masked too, to keep only the recent past events. The motivation is to transform the score matrix to a one-dimensional array by assigning a scalar value to each intermediate feature. We extract the diagonal values of  $S^{\{V,A\}}$  such as  $a_t^{\{V,A\}} = \text{diag}(S^{\{V,A\}})_t$  is the attention scalar that characterizes  $h_t^{\{V,A\}}$  within its local temporal context  $[t-3, \dots, t]$ . Here, the normalization effect of the softmax is critical, as each diagonal element will be scaled according to its relative importance within its context.

**Attributes Combination:** Let  $v$  be the combination of all the uncertainty indicators such that  $v_t = \{\eta_t, \mu_t, \delta_t^V, \delta_t^A, a_t^V, a_t^A\}$ .  $\lambda_t^V$  and  $\lambda_t^A$  in Eq.2 are dynamically computed using Eq.5:

$$\lambda_t^{\{V,A\}} = W^{\{V,A\}} v_t + b^{\{V,A\}} \quad (5)$$

where  $W^{\{V,A\}}$  and  $b^{\{V,A\}}$  are trainable weights and biases, respectively, that are updated during the end-to-end training of the whole architecture.

Table 1: Performance comparison of the proposed method with the state of the art. Results are reported on both validation and ActivityNet Challenge hidden test sets.

Method	mAP	AUC
<i>Validation subset</i>		
<b>Our method</b>	<b>91.9</b>	<b>96.3</b>
Active Speakers Context [11]	87.1	-
Huang et al. [12]	-	93.2
Google Baseline [6]	86.3	92.0
Naver Corp. (Temporal Convolutions) [8]	85.5	-
Zhang et al. [9]	84.0	-
<i>ActivityNet Challenge Leaderboard</i>		
<b>Our method</b>	<b>89.5</b>	-
Naver Corp. [8]	87.8	-
Active Speakers Context [11]	86.7	-
Zhang et al. [9]	83.5	-
Google Baseline [6]	82.1	-

### 3. Evaluation and Analysis

#### 3.1. Experimental Setup

**AVA-ActiveSpeaker Dataset** [6]: It is the most comprehensive, largest, and challenging publicly available dataset for audiovisual ASD problem. It consists of 262 movies divided into training (120), validation (33) and test (109) sets. In total, 5,498K faces are labeled with normalized bounding boxes. For training and validation sets, ground-truths on whether someone is speaking are also provided. The ground-truths for the test set are withheld for the annual ActivityNet Challenge [7].

**Training Strategy:** The ADAGRAD optimizer [22] is used with a learning rate of 0.015 to train the network in end-to-end fashion for 20 epochs with mini-batches of 16 sequences and without any pre-training. Roth et al. [6] demonstrated that stacking few consecutive frames as input of the first 2D convolutional layer is beneficial to learn short temporal motion. Therefore, the input to the visual network (Fig.3b) is a stack of 3 consecutive grayscale face thumbnails. The faces are extracted using the provided bounding boxes and resized to 224x224. We feed the audio embedding network (Fig.3c) with 13 MFCC features extracted from the preceding 0.5s of audio with a 25ms window and a 10ms step. Our BiGRUs are trained with 1.12s long segments (28 frames) to capture most of the different speech patterns within the AVA-Active-Speaker dataset, the average continuous speech duration being 1.11s [6].

**Metrics:** For ease of comparison with previous studies, we evaluate the proposed method using the official ActivityNet Challenge evaluation script [6] that computes the mean Average Precision (mAP) score. When available, Area Under Receiver Operating Characteristic Curve (AUC) score is also provided in Table 1. As the AVA-ActiveSpeaker test set ground-truths are kept private for the official challenge, most of our performance analysis is conducted on the validation set.

#### 3.2. Comparison with State Of The Art

Experimental results in Table 1 show that the proposed method outperforms all existing approaches on both test and validation sets, in terms of mAP and AUC scores. It is worth highlighting that the proposed method outperforms one of the best approaches in [11] by a significant margin of 4.8% and 2.8% in mAP score on the validation and test sets respectively. Furthermore, our approach surpasses Naver’s method [8] by 1.7% and ranks first in the AVA-ActiveSpeaker Leaderboard.

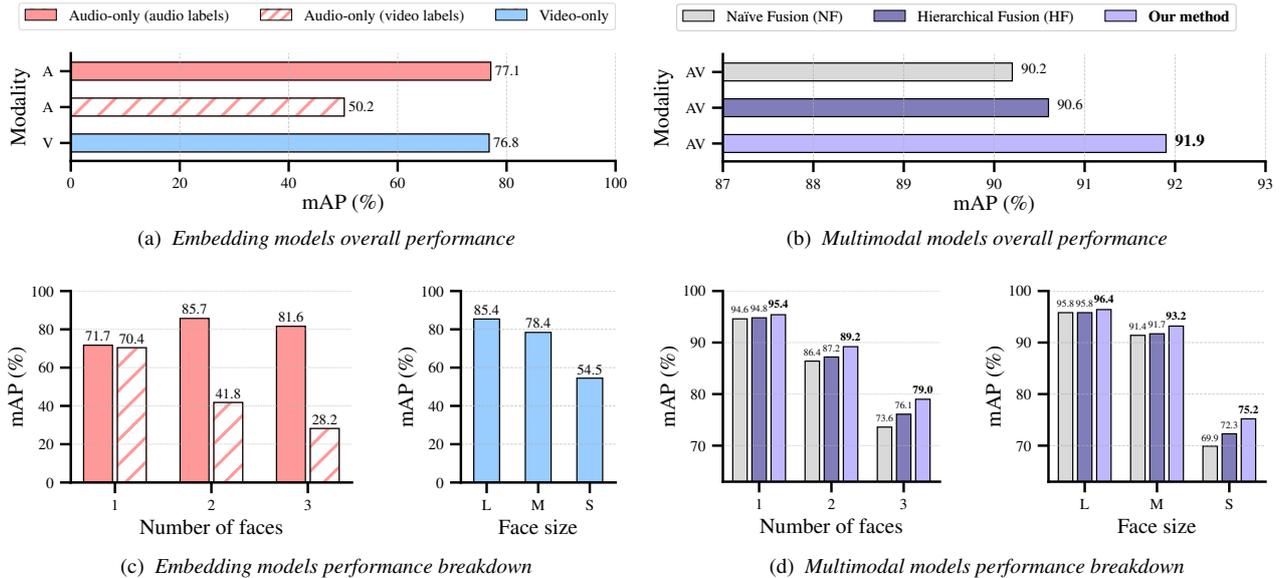


Figure 5: Performance comparison of the embedding (a,c) and the multimodal (b,d) models on the validation set. (c) and (d) present a performance analysis according to the number of detected faces (left) and the face size (right). For the number of faces, we split the validation set into three subsets by gathering one, two, and three faces frames. Altogether, these three cases cover more than 90% of the AVA-ActiveSpeaker dataset. For the face sizes, we sort the validation dataset by ascending face-size order and split it in three equal parts denoted Small (S), Medium (M) and Large (L) respectively.

Table 2 compares the number of parameters of the top-ranked models on the ActivityNet Challenge Leaderboard. Our results are very favorable since our approach has significantly fewer parameters compared to the state of the art and does not necessitate any pre-training and/or ensemble-models.

Table 2: Comparison of the gross number of parameters.

Method	#params	pre-training
<b>Our method</b>	<b>2M</b>	<b>×</b>
Naver Corp. [8]	13M	✓
Active Speakers Context [11]	22M	✓
Zhang et al. [9]	22M	✓

### 3.3. Performance Breakdown

Multi-objective learning, as formulated in Section 2.1, aims at (1) learning accurate representations of each modality, and (2) allowing unbiased estimation of the uncertainty of video and audio intermediate features. We therefore evaluate the discriminative power of video and audio embedded representations. Figures 5a and 5c detail the performance of the embedding networks detailed in Fig.3 when used alone. The presented results are with 128-dim BiGRU added on top of each embedding network. We compare the performance of the audio embedding network optimized towards either video or audio ground-truth labels as discussed in Section 2.1. As expected, the performance of the audio embedding network trained with video ground-truths strongly degrades on the number of speakers. It suffers in the ambiguous scenario presented in Fig.1 where multiple persons share the same audio track. On the contrary, the performance of the audio network trained with audio labels is almost constant while increasing the number of speakers. We also observed a major mAP score increase of 26.9% compared with using video labels. Thus, the multi-objective approach using an independent audio-based labels allows the reliability on the audio network in difficult/ambiguous scenarios.

Figures 5b and 5d compare the performance of our fusion method with the Naïve and Hierarchical [17] Fusion schemes in different scenarios. Our multi-objective model is first evaluated using a naïve concatenation fusion (NF) to combine video and audio modalities. Here it is crucial to note that the method improves Active Speakers Context [11] mAP by 3.1% on validation subset (Table 1). This result highlights the effectiveness of our end-to-end multi-objective learning. Hierarchical Fusion (HF) refers the fusion scheme presented in [17]. Note that our adaptation implies a slight modification of the initial method to match our "many-to-many" architecture. As shown in Fig.5b, the proposed fusion scheme clearly has an advantage over NF and HF. Fig.5d presents additional comparative analysis results by varying the number of faces detected (left) and the size of the face thumbnails (right). The proposed method clearly outperforms both NF and HF schemes, especially in challenging scenarios.

## 4. Conclusion

In this paper, we proposed a self-attention and uncertainty-based fusion mechanism that learns a comprehensive understanding of every situation towards ASD disambiguation. Our approach catered multi-objective optimization to encourage the learning of unbiased multimodal features. Experimental results on the challenging AVA-ActiveSpeaker dataset demonstrate that the proposed method achieves superior performance than existing methods. Besides, the proposed method outperformed the state-of-the-art on both validation and test datasets and ranked first in the ActivityNet Challenge, despite having fewer parameters and without any pre-training.

## 5. Acknowledgements

This work has been supported by the French government, through the 3IA Côte d'Azur Investment in the Future project managed by the National Research Agency (ANR) with the reference numbers ANR-19-P3IA-0002.

## 6. References

- [1] R. Cutler and L. S. Davis, "Look who's talking: speaker detection using video and audio correlation," in *IEEE International Conference on Multimedia and Expo. ICME2000. Proceedings. Latest Advances in the Fast Changing World of Multimedia (Cat. No.00TH8532)*, 2000.
- [2] P. Chakravarty, S. Mirzaei, T. Tuytelaars, and H. Van hamme, "Who's speaking? audio-supervised classification of active speakers in video," in *Proceedings of the 2015 ACM on International Conference on Multimodal Interaction*, 2015, p. 87–90.
- [3] P. Chakravarty and T. Tuytelaars, "Cross-modal supervision for learning active speaker detection in video," *ArXiv*, 2016.
- [4] P. Chakravarty, J. Zegers, T. Tuytelaars, and H. V. hamme, "Active speaker detection with audio-visual co-training," *Proceedings of the 18th ACM International Conference on Multimodal Interaction*, 2016.
- [5] J. S. Chung and A. Zisserman, "Out of time: automated lip sync in the wild," in *Workshop on Multi-view Lip-reading, ACCV*, 2016.
- [6] J. Roth, S. Chaudhuri, O. Klejch, R. Marvin, A. C. Gallagher, L. Kaver, S. Ramaswamy, A. Stopczynski, C. Schmid, Z. Xi, and C. Pantofaru, "Ava-activespeaker: An audio-visual dataset for active speaker detection," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2020.
- [7] B. G. Fabian Caba Heilbron, Victor Escorcia and J. C. Niebles, "Activitynet: A large-scale video benchmark for human activity understanding," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015.
- [8] J. S. Chung, "Naver at activitynet challenge 2019 – task b active speaker detection (ava)," 2019.
- [9] Y.-H. Zhang, J.-Y. Xiao, S. Yang, and S. Shan, "Multi-task learning for audio-visual active speaker detection," 2019.
- [10] S.-W. Chung, J. S. Chung, and H.-G. Kang, "Perfect match: Improved cross-modal embeddings for audio-visual synchronisation," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2019.
- [11] J. L. Alcazar, F. Caba, L. Mai, F. Perazzi, J.-Y. Lee, P. Arbelaez, and B. Ghanem, "Active speakers in context," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.
- [12] C. Huang and K. Koishida, "Improved active speaker detection based on optical flow," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, 2020.
- [13] S. L. Happy, A. Dantcheva, A. Das, F. Bremond, R. Zeghari, and P. Robert, "Apathy classification by exploiting task relatedness," in *FG2020 - 15th IEEE International Conference on Automatic Face and Gesture Recognition*, 2020.
- [14] A. Kendall, Y. Gal, and R. Cipolla, "Multi-task learning using uncertainty to weigh losses for scene geometry and semantics," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [15] T. Luong, H. Pham, and C. D. Manning, "Effective approaches to attention-based neural machine translation," in *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, 2015.
- [16] C. Hori, T. Hori, T. Lee, Z. Zhang, B. Harsham, J. R. Hershey, T. K. Marks, and K. Sumi, "Attention-based multimodal fusion for video description," in *IEEE International Conference on Computer Vision (ICCV)*, 2017.
- [17] M. Sanabria, F. Precioso, and T. Menguy, "Hierarchical Multimodal Attention for Deep Video Summarization," in *25th International Conference on Pattern Recognition*, 2021.
- [18] D. Micci-Barreca, "A preprocessing scheme for high-cardinality categorical attributes in classification and prediction problems," *SIGKDD Explor. Newsl.*, vol. 3, p. 27–32, 2001.
- [19] Y. Gal, "Uncertainty in deep learning," Ph.D. dissertation, University of Cambridge, 2016.
- [20] C. Guo, G. Pleiss, Y. Sun, and K. Q. Weinberger, "On calibration of modern neural networks," ser. *Proceedings of Machine Learning Research*, vol. 70, 2017, pp. 1321–1330.
- [21] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus, "Intriguing properties of neural networks," in *International Conference on Learning Representations*, 2014.
- [22] J. Duchi, E. Hazan, and Y. Singer, "Adaptive subgradient methods for online learning and stochastic optimization," *Journal of Machine Learning Research*, vol. 12, pp. 2121–2159, 2011.