

# AUDIO-VISUAL KEYWORD SPOTTING BASED ON MULTIDIMENSIONAL CONVOLUTIONAL NEURAL NETWORK

Runwei Ding, Cheng Pang and Hong Liu<sup>†</sup>

Key Laboratory of Machine Perception, Shenzhen Graduate School, Peking University, China  
dingrunwei@pkusz.edu.cn; chengpang@sz.pku.edu.cn; hongliu@pku.edu.cn

## ABSTRACT

The fusion of audio and visual information is one of the most promising solutions for reliable keyword spotting (KWS), particularly when audio is corrupted by noise. KWS aims to detect a specific word in an audio stream, which still remains a challenging problem under noisy environments. In this paper, an audio-visual neural network based on multidimensional convolutional neural network (MCNN) is proposed to perform audio-visual KWS. Firstly, the log mel-spectrogram and lip area sequence are extracted, respectively, from the audio and visual streams, and are taken as the input of the audio-visual neural network. Then, an audio-visual neural network based on MCNN consisting of 2D CNN and 3D CNN is used to model the time-frequency feature of the log mel-spectrogram and the spatiotemporal feature of the lip area sequence, respectively. Finally, the outputs of the audio and visual networks are combined for KWS through decision fusion. Experimental results on the PKU-AV database under complex acoustic conditions demonstrate that the proposed method achieves preferable performance compared to other state-of-the-art methods.

**Index Terms**— Audio-visual, keyword spotting, multidimensional neural network, decision fusion.

## 1. INTRODUCTION

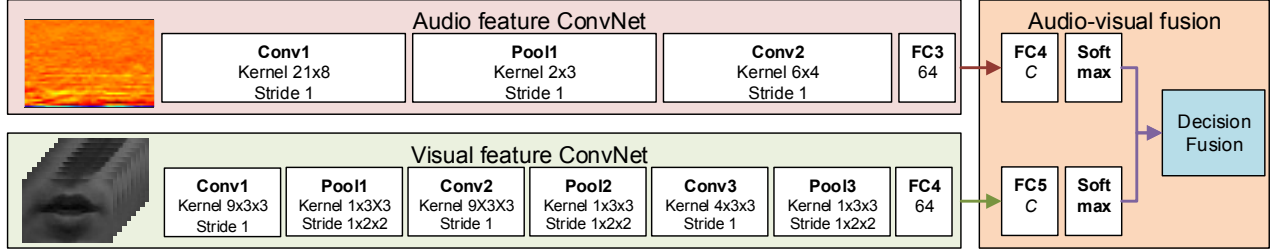
Keyword spotting (KWS) aims to detect a specific word in a speech signal, which has wide applications such as wake-word recognition, speech data mining, human-robot interaction, etc [1–4]. Although continuous speech recognition (CSR) [5] has been widely researched in the past decades, in some scenarios, its complete transcription is unnecessary because the key information usually lies in only part of the input utterance. In general, a KWS system has a much lower complexity than a CSR system, and is thus likely to have a better performance.

<sup>†</sup>: Corresponding author

This work is supported by National Natural Science Foundation of China (NSFC, No. U1613209), Scientific Research Project of Shenzhen City (No. JCYJ20170306164738129, CKCY2017050810242781).

In the past decades, significant research efforts have been made to perform KWS [6–10], nevertheless, the performance degrades substantially in the presence of acoustic interference, e.g. noise and reverberation. Motivated by the fact that human speech perception involves not only hearing but also about seeing, the visual information of lip/face is used and verified to be helpful for many speech-related applications, such as CSR, speaker or emotion recognition, etc [11–14]. Audio-visual keyword spotting (AV-KWS) using both acoustic and visual information can complementarily achieve more robust KWS. Since visual speech components are not affected by acoustic noise, they can compensate for the performance degradation in audio-only KWS under acoustically noisy conditions. In the past few years, few related works have been proposed to achieve AV-KWS. In [15], English word spotting has been achieved based on hidden Markov model (HMM) with the feature fusion of audio and face/mouth information. A novel lip descriptor that includes spatiotemporal information, was proposed to enhance the HMM-based AV-KWS for Mandarin under diverse noise conditions [16]. Then, a parallel two-step decision-fusion strategy was designed to combine the lip information and audio information for a more robust AV-KWS.

In this paper, a novel audio-visual KWS method based on multidimensional convolutional neural network (MCNN) is proposed. The log mel-spectrogram and lip area sequence respectively extracted from the audio and visual streams are taken as the audio and visual features, respectively. To make full use of each dimension information in the audio and visual features, an audio-visual neural network based on MCNN that consists of 2D CNN and 3D CNN, is proposed to model the audio and visual features. The 2D CNN is used to simultaneously learn the time-frequency features of the log mel-spectrogram through 2D convolutional operation. Similarly, the 3D CNN is utilised to learn the temporal and spatial features of the lip area sequence through 3D convolutional operation. Finally, the outputs of the audio and visual network are combined through decision fusion to estimate the posterior probability of each keyword. Experimental results on the PKU-AV database demonstrate that the proposed method can obtain more robust performance compared to other popular methods.



**Fig. 1.** The architecture of the MCNN-based audio-visual neural network. The upper part is the 2D CNN architecture for the audio network. The lower part is the 3D CNN architecture for the visual network.

## 2. AUDIO AND VISUAL FEATURE EXTRACTION

To more conveniently analyze the audio signal, it is enframed by a window of 30 ms with a frame shift of 10 ms. At each audio frame, 40 dimensional log mel-spectrogram features are extracted. In order to exploit the temporal relation of speech signal, and be consistent to the visual frame rate, the context frames with a duration of 1 s are used for the each individual frame. Finally, for the current frame, the log mel-spectrogram features of all the context frames are stacked to generate the audio feature  $A \in \mathcal{R}^{101 \times 40}$ , where 101 and 40 denote the feature dimensions in time and frequency, respectively.

Lipreading which recognizes utterances by analyzing the visual recordings of a speaker’s mouth without audio information, is the core of visual speech recognition. Here, the lip areas in visual stream are extracted as the visual feature. The human face in each video frame is firstly extracted by using the histogram of oriented gradients (HOG) feature combined with a linear classifier, which is then transformed to a gray-scale image. Then, the lip area for each video frame is extracted by estimating the face pose with 68 landmarks [17]. The extracted lip areas are resized to a fixed size of  $60 \times 100$ , and sequentially concatenated to constitute the lip area sequence. The length of the lip area sequence is set to 20, i.e. to 1 s, as in the case for audio sequence. Let  $V$  denote the extracted lip area in one frame, the lip area sequence is written as  $V = [V_1, V_2, \dots, V_{20}]^T$ ,  $V \in \mathcal{R}^{20 \times 60 \times 100}$ , where  $\top$  denotes transpose operation.

## 3. AUDIO-VISUAL NEURAL NETWORK

The architecture of the audio-visual neural network which is composed of multidimensional convolutional neural networks (MCNNs), is shown in Fig. 1. The MCNN concludes coupled 2D and 3D CNNs. For both the audio and visual networks, each convolutional layer is followed by a Rectified Linear Unit (ReLU) activation except for the last layer.

### 3.1. Audio Network

Due to the real-time requirement for KWS, the number of layers in the audio network should be as small as possible to

reduce its computational complexity while maintaining satisfactory KWS performance. To this end, we design an audio network that includes two 2D convolutional layers, one 2D max-pooling layer and one fully connected (FC) layer. As shown in the upper part of Fig. 1, the audio feature  $A$  is first put into a 2D convolutional layer with kernel size of  $21 \times 8$ . Then, a 2D max-pooling layer with kernel size of  $2 \times 3$  is used to reduce the variability in the time-frequency domain that caused by the speaking style, channel distortions, etc. The pooling operation performs sub-sampling to reduce the dimension of time-frequency audio feature. After the pooling operation, a 2D convolutional layer with kernel size of  $6 \times 4$  is used to weight the audio features. Finally, a fully connected layer is used to compress the output of the previous layer into 64 output units. In this network, the strides for the 2D convolutional and pooling layers are 1. Zero-padding is not adopted, because it would introduce extra virtual zero-energy coefficients that are meaningless in the sense of local feature extraction. Non-square kernels are used in the CNN layers to learn more time-domain information with the number-limited layers.

### 3.2. Visual Network

Following the similar principle for the audio network, the visual network is composed of three 3D convolutional layers, three 3D max-pooling layers and one fully connected layer. As shown in the lower part of Fig. 1, first, the visual feature  $V$  is put into a 3D convolutional layer with kernel size of  $9 \times 3 \times 3$ , and a 3D max-pooling layer with kernel size of  $1 \times 3 \times 3$  is subsequently used to achieve spatial feature pooling. Next, the same 3D convolution and max-pooling operations are repeated one additional time. Then, a 3D convolutional layer with kernel size of  $4 \times 3 \times 3$  and a 3D max-pooling layer with kernel size of  $1 \times 3 \times 3$  are applied. Finally, a fully connected layer is used to compress the output of the previous layer into 64 output units. In this network, the 3D convolutional operations are performed to find the correlation of spatiotemporal lip features. The strides in the 3D convolutional layers are 1. To improve the robustness to the moving lip effect, the pooling stride in the 3D max-pooling layers is set to 2, to maintain lip movement features in the neighborhood of the pooling kernel.

## 4. AUDIO-VISUAL FUSION

Audio and visual information can be complementary for KWS, and they are integrated at the decision level in this work. As the audio-visual fusion shown in Fig. 1, it includes two layers, namely, a fully connected layer and a softmax layer. The fully connected layer is used to separately compress the 64 output units of the audio and visual networks, whose output size is the number of keywords  $C$ . The compressed outputs of this layer are respectively taken as the scores of audio feature  $A$  and visual feature  $V$  over  $C$  candidate keywords, which are labeled as  $coa$  and  $cov$ , respectively. Then, a softmax layer is adopted, whose outputs are taken as the predicted probability from the audio and visual networks with the given audio features  $A$  and visual features  $V$ .

The audio-visual fusion is achieved by weighted addition of the predicted probabilities from the audio and visual networks. Let  $l$  denote the index (i.e., label) of keyword  $x$ , the final predicted probability is obtained by summing the predicted probabilities from the audio and visual networks:

$$P(x_l|A, V, W) = \alpha P(x_l|A, W_a) + (1 - \alpha) P(x_l|V, W_v) \\ = \alpha \frac{e^{coa_l}}{\sum_{i=1}^C e^{coa_i}} + (1 - \alpha) \frac{e^{cov_l}}{\sum_{i=1}^C e^{cov_i}}, \quad (1)$$

where  $P(x_l|A, W_a)$  and  $P(x_l|V, W_v)$ , respectively, denote the predicted probabilities from audio and visual networks,  $W_a$ ,  $W_v$  and  $W$ , respectively, denote the learnable weight matrixes for audio-only, visual-only and audio-visual networks,  $\alpha$  is the weighting factor for the audio information. The keyword label estimated by audio-visual fusion is obtained through:

$$\hat{l} = \arg \max_l (P(x_l|A, V, W)). \quad (2)$$

Similarly, the class label estimated by the audio-only network (or visual-only network) can be obtained via the ‘‘argmax’’ operation for  $P(x_l|A, W_a)$  (or  $P(x_l|V, W_v)$ ).

A combined loss function is designed to train the audio-visual neural network. The cross-entropy loss is used:

$$L(y, l) = -\log \frac{e^{yl}}{\sum_{i=1}^C e^{y_i}}, \quad (3)$$

where  $y$  is the prediction. The loss function used for the training of audio-visual network is formulated as:

$$L^{av} = \beta L^a + (1 - \beta) L^v \\ = -\beta \log \frac{e^{coa_l}}{\sum_{i=1}^C e^{coa_i}} - (1 - \beta) \log \frac{e^{cov_l}}{\sum_{i=1}^C e^{cov_i}}, \quad (4)$$

where  $L^{av}$  denotes the loss of audio-visual fusion,  $L^a$  and  $L^v$  denote the losses of audio and visual information, respectively,  $\beta$  is a hyperparameter to control the importance of audio and visual information. The audio information is generally more effective than visual information for speech recognition, therefore  $\alpha$  and  $\beta$  are set to 0.7 and 0.7, respectively. The posterior handling module in [9] is adopted to combine the frame-level posterior scores into a single score for each keyword, which is used for the final word detection.

## 5. EXPERIMENTS AND ANALYSES

### 5.1. Database

The datasets used in our experiments is an audio-visual database collected by ourselves, called the PKU-AV database. The PKU-AV database was collected in an acoustically quiet environment with controlled normal light, which recorded by 20 subjects (12 males and 8 females). Three hundred Chinese Mandarin utterances are spoken by each person, which is recorded by a video camera at 20 frames per second with a resolution of  $640 \times 480$  under the restriction that the mouth region is not occluded. The corresponding speech audio is synchronously recorded at a sampling frequency of 16 kHz with 16 bits per sample. We define 30 keywords/phrases that are frequently used in daily life. In each subject, there are 5 sample sentences for each keyword. So there are 100 sample sentences for each keyword and 3000 negative sample sentences which do not conclude the keywords.

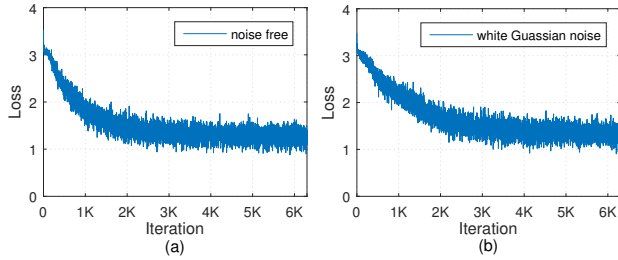
### 5.2. Experimental Setting

In this work, all the network training and evaluations are achieved through PyTorch<sup>1</sup> using one NVIDIA GeForce GTX 1080 GPU. The batch size is set to 64 for all the experiments. Each dataset is split with the ratio of 8:1:1 for training, validation and testing, respectively. The stochastic gradient descent with a momentum of 0.9 is used to train all the networks. The learning rate is set to 0.001, and Dropout [18] with a probability of 0.5 is used to alleviate overfitting.

Noisy data are generated by summing the audio in the original database and pure noise. Two types of noise, i.e., white Gaussian noise and speech babble noise, from the Noisex92 database [19] are adopted. Their sampling frequency is 16 kHz. The two noise signals are added to the original audio signals with different signal-to-noise ratios (SNRs). The noises are added to the original audio stream with the SNRs of 20 dB, 10 dB and 0 dB for the training and validation data, and with the SNRs of 20 dB, 15 dB, 10 dB, 5 dB, 0 dB for the testing data. For visual preprocessing, all the faces in the two databases are detected in advance. Figure of merit (FOM) is used to measure the experimental results. FOM is defined as the average percentage of correctly detected keywords as the threshold is varied from one to 10 false alarms per keyword per hour [20].

In order to allow speaker-independent recognition, the original PKU-AV database is also divided according to subject: training sets: data from 16 subjects are randomly selected from the PKU-AV database and used for training models; validation sets: data from 2 subjects are randomly selected from the remaining data in the PKU-AV database to validate the training; testing sets: data from the last 2 subjects are used to test the performance of the trained models. The average FOM performances over 10 trials are collected as the experimental results.

<sup>1</sup><http://pytorch.org>



**Fig. 2.** The loss curves of the audio-visual network for the training data (a) no noise (b) white Gaussian noise.

### 5.3. Audio-Visual KWS

The performances of audio-only system, visual-only system and their combination are evaluated for keyword spotting. Besides, the comparison with the method in [16] is performed to evaluate the effectiveness of our proposed method. Fig. 2 shows the training loss curves of the audio-visual neural network for the training data with and without noise. It can be seen that the training loss curve for the training data without noise decreases and converges faster than those with noise.

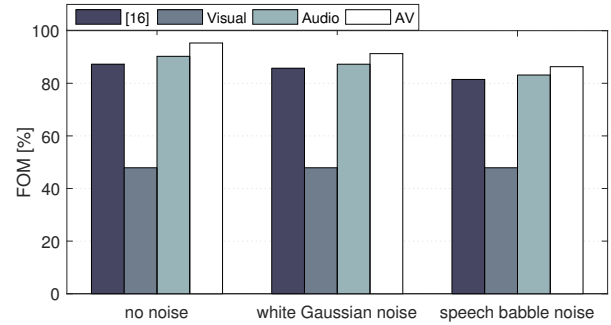
**Table 1.** FOM performances of the audio -only, visual-only, our AV-KWS system (AV) and the AV-KWS method in [16] using noise-free training data under the test condition with white noise.

SNR(dB)	20	15	10	5	0
[16]	81.93%	75.68%	65.57%	49.23%	45.37%
Audio	78.95%	68.73%	44.26%	31.56%	22.17%
Visual	47.89%	47.89%	47.89%	47.89%	47.89%
AV	85.26%	81.13%	74.32%	60.21%	55.26%

**Table 2.** FOM performances of the audio -only, visual-only, our AV-KWS system (AV) and the AV-KWS method in [16] using noise-free training data under the test condition with babble noise.

SNR(dB)	20	15	10	5	0
[16]	78.94%	70.69%	58.34%	44.52%	41.67%
Audio	71.84%	62.43%	32.67%	21.38%	14.17%
Visual	47.89%	47.89%	47.89%	47.89%	47.89%
AV	83.26%	79.46%	68.25%	56.84%	51.34%

Tables 1 and 2 show the average FOMs for the audio-only, visual-only, our AV-KWS system using the networks trained with noise-free data and the method in [16] under white and babble noise test conditions, respectively. It is obvious that the visual-only system achieves constant performance for different noise intensities, because the visual condition is invariant and the variation of noise has no influence on the extraction of lip area. The performance of audio-only system is acceptable when SNR=20dB, but the performance decreases substantially with decreasing SNR. The method



**Fig. 3.** FOM performances of different features and methods for KWS based on training data with different noises under the test condition without noise.

in [16] achieves better performance than audio-only recognition with the benefit of using visual information, while it obtains lower performance than visual-only recognition under strong noisy conditions ( $SNR \leq 5dB$ ). By directly taking the lip area sequence as input, the visual network in our method can learn more robust spatiotemporal information for visual stream compared with HMM in [16]. The proposed audio-visual method obtains the best performance for both types of noise, mainly due to the effective combination of audio and visual information through the proposed audio-visual network based on MCNN.

The average FOMs for the four methods under noise-free test conditions are shown in Fig. 3. Here, our methods are trained and evaluated in three noisy conditions, respectively. The proposed AV-KWS method obtains the highest performance, because the audio and visual information is effectively learned and combined by the MCNN-based audio-visual neural network at decision level. The time-frequency information of audio stream can be better learned through the 2D convolutional operations with non-square kernels. By directly taking the lip area sequence as the input of visual network, the temporal and spatial information of lip motion can be synchronously learned through the 3D convolutional operations, which can avoid the error caused by the lip feature extraction in [16].

## 6. CONCLUSION

In order to effectively fuse audio and visual information, a novel AV-KWS method based on multidimensional CNN is proposed. The audio network effectively models the time-frequency information of audio stream by doing 2D convolutional operation on the log mel-spectrogram. The visual network effectively learns the spatiotemporal information of visual stream by doing 3D convolutional operation on the lip area sequence. By properly integrating the outputs of audio and visual networks at decision level, the proposed method achieves preferable and robust KWS performance under noisy conditions. Experiments based on the PKU-AV database demonstrate the effectiveness and adaptability of our method for various noisy environments.

## 7. REFERENCES

- [1] Martin Wollmer, Florian Eyben, Joseph Keshet, Alex Graves, Bjorn Schuller, and Gerhard Rigoll, “Robust discriminative keyword spotting for emotionally colored spontaneous speech using bidirectional lstm networks,” in *IEEE International Conference on Acoustics, Speech and Signal Processing*, 2009, pp. 3949–3952.
- [2] Berlin Chen, “Word topic models for spoken document retrieval and transcription,” *ACM Transactions on Asian Language Information Processing*, vol. 8, no. 1, pp. 1–27, 2009.
- [3] Damianos Karakos, Richard Schwartz, Stavros Tsakalidis, Le Zhang, Shivesh Ranjan, Tim Tim Ng, Roger Hsiao, Guruprasad Saikumar, Ivan Bulyko, Long Nguyen, et al., “Score normalization and system combination for improved keyword spotting,” in *IEEE Workshop on Automatic Speech Recognition and Understanding*, 2013, pp. 210–215.
- [4] Cheng Pang, Hong Liu, Jie Zhang, and Xiaofei Li, “Binaural sound localization based on reverberation weighting and generalized parametric mapping,” *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 25, no. 8, pp. 1618–1632, 2017.
- [5] George Saon and Jen-Tzung Chien, “Large-vocabulary continuous speech recognition systems: A look at some recent advances,” *IEEE Signal Processing Magazine*, vol. 29, no. 6, pp. 18–33, 2012.
- [6] Jay G Wilpon, Lawrence R Rabiner, C-H Lee, and ER Goldman, “Automatic recognition of keywords in unconstrained speech using hidden markov models,” *IEEE Transactions on Acoustics, Speech, and Signal Processing*, vol. 38, no. 11, pp. 1870–1878, 1990.
- [7] Anupam Mandal, KR Prasanna Kumar, and Pabitra Mitra, “Recent developments in spoken term detection: a survey,” *International Journal of Speech Technology*, vol. 17, no. 2, pp. 183–198, 2014.
- [8] Guoguo Chen, Carolina Parada, and Georg Heigold, “Small-footprint keyword spotting using deep neural networks,” in *IEEE International Conference on Acoustics, Speech and Signal Processing*, 2014, pp. 4087–4091.
- [9] Tara N Sainath and Carolina Parada, “Convolutional neural networks for small-footprint keyword spotting,” in *Interspeech*, 2015, pp. 1478–1482.
- [10] Sercan . Ark, Markus Kliegl, Rewon Child, Joel Hestness, Andrew Gibiansky, Chris Fougner, Ryan Prenger, and Adam Coates, “Convolutional recurrent neural networks for small-footprint keyword spotting,” in *Interspeech*, 2017, pp. 1606–1610.
- [11] Darryl Stewart, Rowan Seymour, Adrian Pass, and Ji Ming, “Robust audio-visual speech recognition under noisy audio-video conditions,” *IEEE Transactions on Cybernetics*, vol. 44, no. 2, pp. 175–184, 2014.
- [12] Ahmed Hussen Abdelaziz, Steffen Zeiler, and Dorothea Kolossa, “Learning dynamic stream weights for coupled-hmm-based audio-visual speech recognition,” *IEEE/ACM Transactions on Audio, Speech and Language Processing*, vol. 23, no. 5, pp. 863–876, 2015.
- [13] Haomain Zheng, Meng Wang, and Zhu Li, “Audio-visual speaker identification with multi-view distance metric learning,” in *IEEE International Conference on Image Processing*. IEEE, 2010, pp. 4561–4564.
- [14] Kun Lu and Yunde Jia, “Audio-visual emotion recognition using boltzmann zippers,” in *IEEE International Conference on Image Processing*. IEEE, 2012, pp. 2589–2592.
- [15] Ming Liu, Ziyou Xiong, Stephen M Chu, Zhenqiu Zhang, and Thomas S Huang, “Audio visual word spotting,” in *IEEE International Conference on Acoustics, Speech, and Signal Processing*, 2004, vol. 3, pp. 785–788.
- [16] Pingping Wu, Hong Liu, Xiaofei Li, Ting Fan, and Xuewu Zhang, “A novel lip descriptor for audio-visual keyword spotting based on adaptive decision fusion,” *IEEE Transactions on Multimedia*, vol. 18, no. 3, pp. 326–338, 2016.
- [17] Vahid Kazemi and Josephine Sullivan, “One millisecond face alignment with an ensemble of regression trees,” in *IEEE Conference on Computer Vision and Pattern Recognition*, 2014, pp. 1867–1874.
- [18] Nitish Srivastava, Geoffrey E Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov, “Dropout: a simple way to prevent neural networks from overfitting,” *Journal of machine learning research*, vol. 15, no. 1, pp. 1929–1958, 2014.
- [19] Andrew Varga and Herman JM Steeneken, “Assessment for automatic speech recognition: Ii. noisex-92: A database and an experiment to study the effect of additive noise on speech recognition systems,” *Speech communication*, vol. 12, no. 3, pp. 247–251, 1993.
- [20] Igor Szöke, Petr Schwarz, Pavel Matejka, Lukás Burget, Martin Karafiát, Michal Fapso, and Jan Cernocký, “Comparison of keyword spotting approaches for informal continuous speech,” in *Interspeech*, 2005, pp. 633–636.