

# An Improved Structured Pruning Approach to Channel-level Pruning for Designing Compact Audio-Visual Wake Word Spotting System

Haotian Wang<sup>1,†</sup>, Gongzhen Zou<sup>1</sup>, Jun Du<sup>1,\*</sup>, Hengshun Zhou<sup>1</sup>, and Shifu Xiong<sup>2</sup>

<sup>1</sup> University of Science and Technology of China, Hefei, Anhui, P.R.China

<sup>2</sup> iFlytek Research, Hefei, Anhui, P. R. China

**Abstract.** In this paper, we propose a novel structured pruning approach based on channel-level pruning with learning-rate rewinding (CPLR) for designing the compact and low latency audio-visual wake word spotting (AVWWS) system. First, an efficient operation that aims at reducing network parameters is explored by integrating blueprint separable convolution (BSCConv). Next, the channel-level pruning with learning-rate rewinding strategy is applied to the improved network to prune network parameters and obtain a compact system. Finally, a binary regulation (BR) strategy is further proposed for reducing the inference time of the above compact system, called BS-CPLR. Tested on the MISP2021 AVWWS database, the results show that the proposed BS-CPLR approach achieves better system performance with fewer model parameters being used. We further tested the compact systems on the TB-RK3399ProD development board. The results show that the proposed BS-CPLR approach also achieves lower delay than BSCConv and CPLR.

**Keywords:** AVWWS, CPLR, BSCConv, binary regulation, development board

## 1 Introduction

Wake word spotting (WWS) plays a very important role in man-machine interaction. The goal of the WWS task is to recognize a predefined wake word [1]. Typical WWS systems are based on audio modality [2, 3], one of the main challenges is that the performance of these systems usually declines a lot under noisy conditions [4–6]. In the recent years, research efforts have been made in detection of wake word spotting under noisy and far-field conditions, including the use of speech enhancement (SE) module [6, 7], the design of new network structures, the design of novel architecture and new training strategy [8–10]. Recently, the authors demonstrate that visual information can improve the system performance based on clear and noisy audio signals [11, 12].

---

† Two authors contributed equally to this research

\* Correspondence: jundu@ustc.edu.cn

However, the introduction of visual information also greatly increases the total amount of network parameters and the computation cost, which will limit its application on lightweight devices.

Accordingly, effective pruning methods for neural network is essential for the feasibility of audio-visual systems. Generally, there are two efficient approaches to compress the complex network, including network pruning and network structure modifying. Network pruning includes structured pruning and unstructured pruning [13, 14]. Structured pruning can not only improve the computational efficiency, but also effectively reduces the number of network parameters, which has attracted more and more attention [15, 16]. Filter pruning is one type of structure pruning which can be realized by Taylor expansion, geometric median (FPGM) and other methods [17–19]. Channel pruning using the batch normalization (BN) layers has also been proven effective [20]. In our previous work, a more advanced compression method call CPLR [21] that based on channel-level pruning [20] and leaning-rate rewinding [22] was proposed. The network structure modifying method includes low-rank approximation and heterogeneous convolution. Low-rank approximation regards the weight matrix in the network as non-singular matrix that can be decomposed into the combination of several low-rank matrices [23]. The heterogeneous convolution method attempts to change the structure of the convolution layers in the network, transform the regular convolution into the combination of depth wise convolution and point wise convolution and MobileNets are successful applications of the heterogeneous convolution [24]. Recently, a new approach of heterogeneous convolution called blueprint separable convolution(BSConv) have been proposed [25]. BSConv explains the efficiency of heterogeneous convolution and performs better than MobileNets on several tasks.

Moreover, the training and pruning strategy is also very important. The classic strategy is the “weights initialization - pruning - fine-tuning” procedure [26]. Then the authors proposed a dynamic pruning approach, which uses sparse training in convenient pruning operation [20]. In 2019, the lottery ticket hypothesis (LTH) [27] was proposed by integrating existing pruning and training modes. In recent years, some researchers have compared various strategies based on LTH such as learning rate rewinding [22].

In this study, we proposed a new network compression approach called blueprint separable convolution channel-level pruning with learning rate rewinding (BS-CPLR) that integrates the channel-level pruning method, BSConv and LTH strategy. Firstly, the BSConv was adopted to transform all the convolution layers in the original network into BSConv layers. Next the CPLR strategy was used to prune and fine-tune the transformed network, slimming the network further. Then, binary regulation was introduced to the number of channels, proposed another version of BS-CPLR called BRBS-CPLR. Finally, we apply the proposed algorithm to neural network architectures of WWS. Evaluated on MISP-2021 AVWWS database and TB-RK3399ProD, the BS-CPLR approach yields consistent improvements, obtains better performance and reduces the time delay on development board.

## 2 Related Work

In this section, we describe the two different types of neural network pruning approaches related to our work, namely blueprint separable convolution (BSConv) and channel-level pruning with learning-rate rewinding (CPLR).

### 2.1 Blueprint separable convolution (BSConv)

Blueprint separable convolution can be regarded as one kind of heterogeneous convolution [25]. All the original convolution layer in the network are splitted into two kinds of convolution kernels: One point-wise convolution layer followed by one depth-wise group convolution layer. Compared with depth-wise separable convolution (DSC) [24], BSConv is a reversed DSC, changes the order of the the point-wise and group-wise convolution because intra-kernel correlations are considered more important than cross-kernel correlations. BSConv outperforms better than DSC on several image classification tasks and can reduce the parameters in the network to about 70% without performance loss.

### 2.2 Channel-level pruning with learning-rate rewinding (CPLR)

Channel-level pruning is one kind of structured pruning method. Proposed in [20], channel-level pruning focus on reduce the number of channels in convolution layers. The channels in the convolution layers represent several feature filters to generate feature map. Some filters are considered less important than others, can not extract powerful features to offer information gain, can be pruned off without performance loss. The authors added a BN layer after each convolution layer, and considered that the importance of channel are highly correlated with the absolute value (ABS) of the scale factor  $\gamma$ . The channel with smaller  $\gamma$  is considered less important for the network and can be pruned freely. Learning-rate rewinding is a train and pruning strategy [22]. Based on the LTH [27], authors compare three training and adjustment strategies for LTH, namely weight rewinding strategy (lottery strategy), pruning and fine-tuning strategy, and learning rate rewinding strategy. The conclusion is the learning rate fine-tuning strategy is the best in most instances [22]. In our previous work [21], we combined the aforementioned two methods, proposed CPLR. Firstly, the original channel-level pruning method was improved from the one-shot pattern to an iterative pattern. Then the learning rate rewinding strategy was adopted to guide the channel-level pruning. CPLR can fully utilize the advantages of both channel-level and the LTH pruning. Evaluated on MISP-2021 AVWWS database, the CPLR approach can yield consistent improvements in terms of both system performance and model efficiency.

## 3 Proposed Methods

Fig. 1 shows the complete architecture of the proposed (BR)BS-CPL approach on AVWWS system.

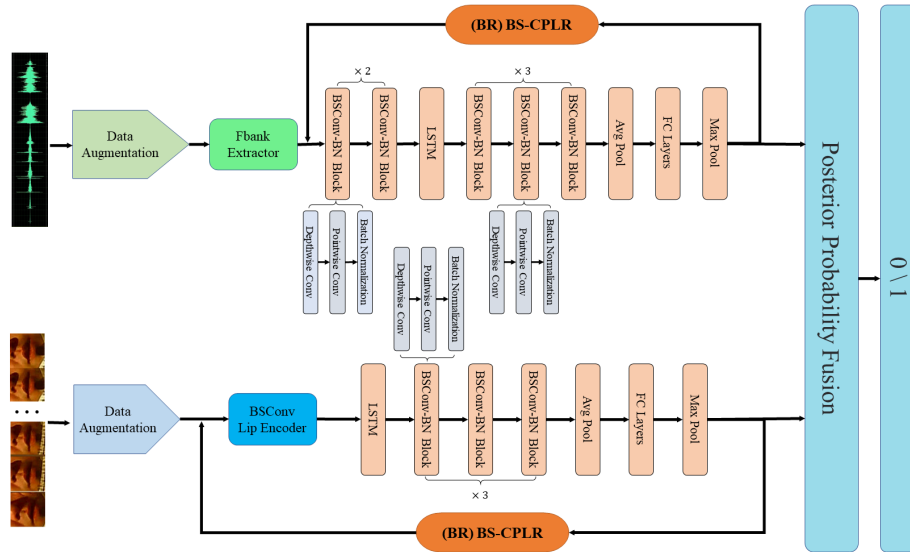


Fig. 1: The overall architecture of the audio-visual wake word spotting system using the proposed BS-CPLR.

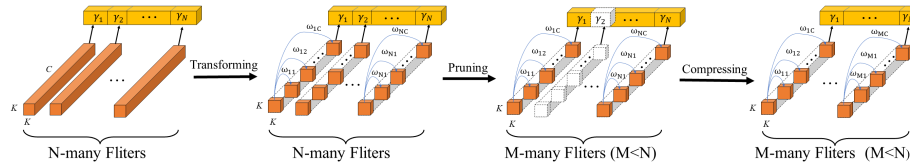


Fig. 2: The application process of BS-CPLR on BSCConv-BN block.

### 3.1 Blueprint separable convolution CPLR (BS-CPLR)

When the original convolution layers are transformed into BSCConv layers, the number of the channels remains the same. So we can use CPLR to prune the transformed network. The main process of applying BSCConv on convolution and BN (ConvBN) block is shown in Fig. 2. Firstly, we transform all the convolution layers in the network into blueprint separable convolution (BSCConv). Then we add one BN layer after each depth-wise group convolution layers. Afterwards, we design a iterative channel-level pruning method to prune the point-wise and depth-wise convolution layers in BSCConv based on scale factor  $\gamma$  of BN layers. We adopt the same training and pruning process as CPLR, using the same learning-rate to retrain the pruned network. Sparse-training method is also adopted to optimize the pruning process, making the unimportant channels easier to be identified.

### 3.2 Binary Regulation for BS-CPLR (BRBS-CPLR)

To enhance the efficiency of our compressed system working on board, we improved the channel-pruning method used in CPLR. The memory space and threads in standard board are always the power of 2. Similarly, we always set our channel numbers and batch size to be the power of 2 to fully occupy the memory space and threads on board, so the network will work more efficiently on board. Inspired by this, we constraint the channel numbers in the obtained pruning mask to be  $2^N$  ( $N$  represents natural number). Then we arrange the channels and prune off the unimportant channels based on the pruning mask. The full algorithm is shown in Algorithm 1.

---

#### Algorithm 1 BIBS-CPLR algorithm

---

- 1** : Transform all the convolution layers in network into BSConv layers, generate the initial network  $f(\theta; \gamma)$ .
  - 2** : Pre-train the initial network to the early-stop point  $f(x; \theta_0; \gamma_0)$  by using sparse-training.
  - 3** : Set the pruning rate  $k = 30\%$  per round.
  - 4** : Create the initial mask  $m = \{m_1, m_2, \dots, m_i, \dots, m_n\}$  based on BN scale factor  $\gamma$  and the pruning rate  $k$ .
  - 5** : Traverse the initial mask, constraint every element  $m_i$  to be the nearest  $2^{N_i}$ , generate a new mask  $\tilde{m} = \{2^{N_1}, 2^{N_2}, \dots, 2^{N_i}, \dots, 2^{N_n}\}$
  - 6** : Arrange and prune the channels using the mask  $\tilde{m}$ , obtain the reinit network  $f(x; \theta_0 \odot \tilde{m}; \gamma_0 \odot \tilde{m})$ .
  - 7** : Use learning rate rewinding strategy to retrain the network to obtain network  $f(x; \theta_1; \gamma_1)$ .
  - 8** : Repeat **3** to **7** until the desired weight sparsity ratio is reached.
  - 9** : Fine-tune the pruned network using the same learning rate as the last retraining process, return fine-tuned network  $f(x; \theta_N; \gamma_N)$ .
- 

We decide the nearest  $2^{N_i}$  in step4 using the following decision criterion:

$$|2^{N_i} - m_i| \leq |2^{N_i \pm 1} - m_i| \quad (1)$$

## 4 Applications on AVWWS System

Our proposed methods are evaluated on an audio-visual wake word spotting task based on the MISP2021 challenge dataset [28]. Our model architecture mainly consists of two parts: audio-only system, audio-visual system. Each system employs three different methods: CPLR, BS-CPLR, and BRBS-CPLR to compare performance. The details are elaborated in the following subsections.

### 4.1 Audio-only WWS system

The baseline AVWWS system provided by the challenge was selected as the original audio-only system, which consists of two 2D-convolution layers, one

long short term memory (LSTM) layer, and three 2D-convolution layers. We add a BN layer after each convolution layer. Then, we replace all convolutions into BSConv and adopted the CPLR, BS-CPLR, and BRBS-CPLR to compress the system.

## 4.2 Audio-visual fusion system

The audio-visual fusion system is consistent with [28], which includes two systems: audio-only and video-only. The audio-only system is described in the paragraph audio-only WWS system. Similar to audio-only system, the corresponding original video-only system provided by the challenge was selected, which consists of a lip feature extractor, one LSTM layer, and three 2D-convolution layers, where the lip feature extractor consists of a 3D-convolution layer, a BN layer, a 3D-pooling layer, and a ResNet-18. Due to the development board does not support 3D operators, all 3D operators in the lip feature extractor are replaced into 2D operators. Similar to the audio-only system, all convolution layers are transformed into BSConv layers, and then CPLR, BS-CPLR, and BRBS-CPLR were adopted to compress the system.

The final predicted results are determined by the posterior probability of the audio-only and video-only WWS systems, as shown in the following formula:

$$P_{AV} = \alpha \times P_A(y_a|\mathbf{f}_a) + \beta \times P_V(y_v|\mathbf{f}_v) \quad (2)$$

where  $P_A(y_a|\mathbf{f}_a)$  and  $P_V(y_v|\mathbf{f}_v)$  is posterior of wake word presence ( $y_a/y_v$ ) generated by feeding audio features  $\mathbf{f}_a$  and video features  $\mathbf{f}_v$  into the audio-only and video-only models, respectively.  $\alpha$  and  $\beta$  are the weights of audio-only and video-only systems. The output value of these models is compared with the preset threshold ( $th_A, th_V, th_{AV}$ ) after the sigmoid operation.

## 5 Experiments and Results

### 5.1 Experimental setup

In this section, we introduce the database we used, the metric to measure the performance of different systems and the development board used for testing.

**Database and metric** We conducted experiments on the MISP2021 AVWWS database [28], which includes about 125 hours of the same amount of near-field, mid-field, and far-field data and has been divided into training, development, and evaluation sets without overlapping speakers.

Following the requirements of the challenge committee, the evaluation metric is determined by the combination of False Reject Rate (FRR) and False Alarm Rate (FAR). The calculation is as follows:

$$Score = FRR + FAR = \frac{N_{FR}}{N_{wake}} + \frac{N_{FA}}{N_{non-wake}} \quad (3)$$

where  $N_{\text{wake}}$  and  $N_{\text{non-wake}}$  denotes the number of samples with and without wake words, respectively.  $N_{\text{FR}}$  denotes the number of samples containing the wake word while not recognized by the system.  $N_{\text{FA}}$  denotes the number of samples that do not contain the wake word but are predicted to be positive by the system. The lower *Score*, the better the system performance. For a more detailed description of the dataset, please refer to [28].

**Development board** The TB-RK3399ProD development board was adopted in our experiments. TB-RK3399ProD is a hardware development board developed for Rockchip RK3399Pro chip, which integrates chip debugging and simulation testing, etc. It adopts high-performance AI processing chip RK3399Pro, integrated with AI neural network processor NPU, supports 8Bit/16Bit operation, and its computing power reaches 3.0Tops, meeting various AI applications such as vision and audio. It is compatible with mainstream AI frameworks TensorFlow, PyTorch, Caffe, etc., and supports both Android and Linux systems. Since its complete interface, the design has strong expansibility, it can apply different use scenarios and full function verification.

## 5.2 Results on standard metric

This section presents the experimental results and analysis of our proposed methods. We first evaluated the performance of the audio-only system. Then, we train the audio-only and video-only systems of the audio-visual fusion system respectively and evaluated the audio-visual fusion system. The *Score* is used as an indicator of the generalization ability of the system. Lower *Score*, better performance.

**Results for audio-only system** We adopted the same data augmentation methods as [28] for audio data to improve the generalization ability of the model. Considering the complexity and challenges of the far-field environment, the proposed methods are evaluated on the far-field audio database.

Our results are shown in Figure. 3. We can observe that after converting all convolution layers of the models into BSConv layers, the parameters of the model are reduced to 1.81M, which is equivalent to the WSR of BS-CPLR and BRBS-CPLL methods is 0.33. With the increase of WSR, the score of the three compression methods showed a trend of decreasing first and then increasing. When WSR is between 0.40 and 0.75, BS-CPLR and BRBS-CPLR achieve better performance than CPLR. Specifically, BRBS-CPLR outperforms BS-CPLR when WSR is between 0.4 and 0.63, but when WSR is between 0.63 and 0.75, BS-CPLR method obtains better results. When WSR exceeds 0.75, the CPLR method obtains better results. Throughout the pruning process, BS-CPLR method achieves the best *Score* (0.2497) on the audio far-field evaluation set when the WSR is 0.72. And BRBS-CPLR method reached the highest WSR of 0.82. More interestingly, we notice that compared to the other two methods, the curve of BRBS-CPLR method is more smooth and steady. In general, the BRBS-CPLR curve

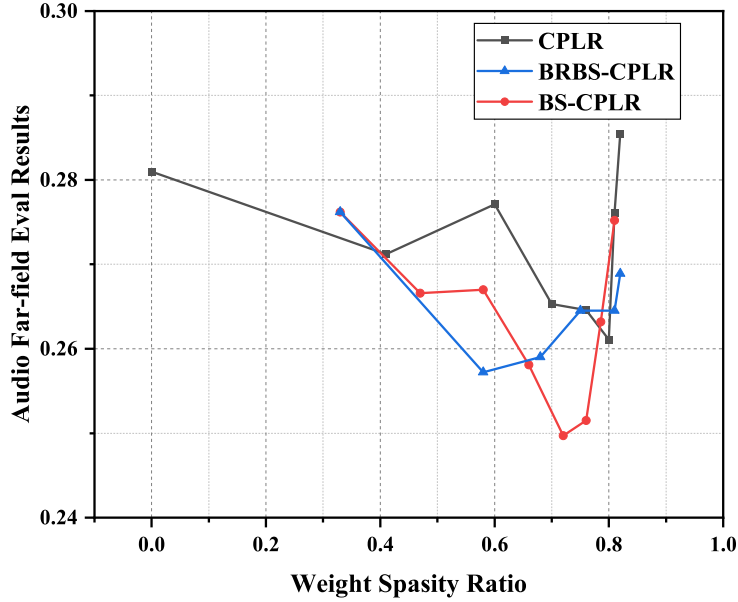


Fig. 3: Performance curves of CPLR, BS-CPLR and BRBS-CPLR on the far-field evaluation set for audio-only system.

is more stable than the other two methods when WSR changes, and BS-CPLR can achieve the best score.

**Results for audio-visual fusion system** For the video-only system, we train and test the trained model on the mid-field data, and the same RandAugment method as [28] was used for data augmentation.

As shown in the Fig. 4, with the decreasing of parameters, the *Score* of the AVWWS system compressed by BS-CPLR shows a pattern of first falling and then rising. When the parameters are about 3.2M, the minimum *Score* is about 0.24. When parameters are under 3.2M, the generalization ability of the system begins to deteriorate rapidly.

The curve of the performance of the AVWWS system compressed by BRBS-CPLR is more stable and obtains the best *Score* (0.2494) when the parameters are about 2.9M. Throughout the pruning process, BRBS-CPLR performs better than BS-CPLR generally and can obtain a higher WSR with no performance loss, but BS-CPLR can achieve a better *Score*.



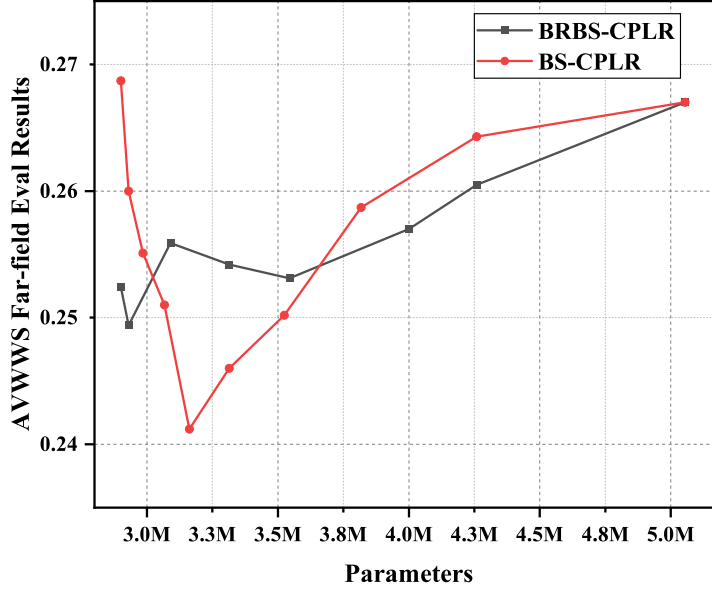


Fig. 4: Performance curves of BS-CPLR and BRBS-CPLR on the far-field evaluation set for AVVWS system.

### 5.3 Results on development board

In this section, we deployed the audio-only and AVVWS system on TB-RK3399ProD and tested the time delay to process single sample on board. The shorter the time delay, the more efficient the system is. We use the systems compressed by BSconv as our baseline. The results are shown in Table 1.

Based on the above results. For both two systems, the processing time delay of the system compressed by BS-CPLR and BRBS-CPLR is apparently lower than the baseline system. And the delay of BRBS-CPLR is lower than that of BS-CPLR, with a reduction of 1.18ms on audio-only system and 58.41ms on AVVWS system. The system compressed by BS-CPLR has less channels and parameters compared with the baseline system, which reduces the amount of computation. So the proposed BS-CPLR obtains better model efficiency. Moreover, the introduction of binary regulation makes the system fully occupy the memory space and threads of the board, so the model efficiency can get a further promotion.

On the other hand, although the introduction of visual modality can improve the adaptability of the WWS in complex environments, the complexity of the video-only system also makes the computational cost of the fusion system higher,

Table 1: Time delay on board of different pruning methods.

Method used to compress	Storage	Time delay
<b>Audio-only system</b>		
BSCnv(baseline)	23.67 MB	29.33 ms
BS-CPLR	5.91 MB	15.85 ms
BRBS-CPLR	5.63 MB	14.67 ms
<b>AVWWS system</b>		
BSCnv(baseline)	51.43 MB	203.62 ms
BS-CPLR	33.67 MB	198.89 ms
BRBS-CPLR	33.39 MB	145.21 ms

which is the main reason for the increasing time delay. The increase in time delay greatly limits the application of the system, and we will focus on how to solve the problem of high delay in the future.

## 6 Conclusions

In this paper, we propose a new strategy based on blueprint separable convolution (BSCnv) and the channel-level pruning with learning-rate rewinding, called BS-CPLR. We first explored an efficient convolution operation which is used to reduce network parameters. Furthermore, a binary conditioning (BR) strategy is proposed to reduce the inference time for the network compressed by BS-CPLR. The system performance and time delay are tested on the MISP2021 AVWWS dataset and the TB-RK3399ProD development board respectively. The experimental results show consistent improvement, indicating the advantages of the proposed method.

## 7 Acknowledgements

This work was supported by the National Natural Science Foundation of China under Grant No. 62171427 and the Strategic Priority Research Program of Chinese Academy of Sciences under Grant No. XDC08050200.

## References

1. D. Leroy, A. Coucke, T. Lavril, T. Gisselbrecht, and J. Dureau, “Federated learning for keyword spotting,” in *Proc. ICASSP 2019*, 2019, pp. 6341–6345.
2. I. López-Espejo, Z.-H. Tan, and J. Jensen, “Improved external speaker-robust keyword spotting for hearing assistive devices,” *IEEE/ACM Transactions on Audio, Speech and Language Processing*, vol. 28, pp. 1233–1247, 2020.

3. Y. Wang, H. Lv, D. Povey, L. Xie, and S. Khudanpur, "Wake word detection with streaming transformers," in *ICASSP 2021 - 2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2021, pp. 5864–5868.
4. X. Ji, M. Yu, J. Chen, J. Zheng, D. Su, and D. Yu, "Integration of multi-look beamformers for multi-channel keyword spotting," in *Proc. ICASSP 2020*, 2020, pp. 7464–7468.
5. I. López-Espejo, Z.-H. Tan, and J. Jensen, "A novel loss function and training strategy for noise-robust keyword spotting," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 29, pp. 2254–2266, 2021.
6. Y. A. Huang, T. Z. Shabestary, and A. Gruenstein, "Hotword cleaner: Dual-microphone adaptive noise cancellation with deferred filter coefficients for robust keyword spotting," in *ICASSP*, 2019, pp. 6346–6350.
7. Y. Gao, Y. Mishchenko, A. Shah, S. Matsoukas, and S. Vitaladevuni, "Towards data-efficient modeling for wake word spotting," in *Proc. ICASSP 2020*, 2020, pp. 7479–7483.
8. Y. Gao, N. D. Stein, C. Kao, Y. Cai, M. Sun, and T. Z. et al., "On front-end gain invariant modeling for wake word spotting," in *Interspeech*, 2020, pp. 991–995.
9. H. Park, P. Zhu, I. Lopez-Moreno, and N. Subrahmanya, "Noisy student-teacher training for robust keyword spotting," in *Interspeech*, 2021, pp. 331–335.
10. A. Hard, K. Partridge, C. Nguyen, N. Subrahmanya, A. Shah, P. Zhu, I. Lopez-Moreno, and R. Mathews, "Training keyword spotting models on non-iid data with federated learning," in *Interspeech 2020*, 2020, pp. 4343–4347.
11. D. Stewart, R. Seymour, A. Pass, and J. Ming, "Robust audio-visual speech recognition under noisy audio-video conditions," *IEEE transactions on cybernetics*, vol. 44, no. 2, pp. 175–184, 2013.
12. L. Momeni, T. Afouras, T. Stafylakis, S. Albanie, and A. Zisserman, "Seeing wake words: Audio-visual keyword spotting," in *31st British Machine Vision Conference 2020, BMVC 2020*, 2020.
13. J. Cheng, P.-s. Wang, G. Li, Q.-h. Hu, and H.-q. Lu, "Recent advances in efficient computation of deep convolutional neural networks," *Frontiers of Information Technology & Electronic Engineering*, vol. 19, no. 1, pp. 64–77, 2018.
14. L. Deng, G. Li, S. Han, L. Shi, and Y. Xie, "Model compression and hardware acceleration for neural networks: A comprehensive survey," *Proceedings of the IEEE*, vol. 108, no. 4, pp. 485–532, 2020.
15. H. Mao, S. Han, J. Pool, W. Li, X. Liu, Y. Wang, and W. J. Dally, "Exploring the regularity of sparse structure in convolutional neural networks," *CoRR*, vol. abs/1705.08922, 2017.
16. Z. Huang and N. Wang, "Data-driven sparse structure selection for deep neural networks," in *ECCV*, 2018, pp. 304–320.
17. P. Molchanov, S. Tyree, T. Karras, T. Aila, and J. Kautz, "Pruning convolutional neural networks for resource efficient inference," in *ICLR*. OpenReview.net, 2017.
18. Y. He, P. Liu, Z. Wang, Z. Hu, and Y. Yang, "Filter pruning via geometric median for deep convolutional neural networks acceleration," in *CVPR*, 2019, pp. 4335–4344.
19. C. Gamanayake, L. Jayasinghe, B. K. K. Ng, and C. Yuen, "Cluster pruning: An efficient filter pruning method for edge ai vision applications," *IEEE Journal of Selected Topics in Signal Processing*, vol. 14, no. 4, pp. 802–816, 2020.
20. Z. Liu, J. Li, Z. Shen, G. Huang, S. Yan, and C. Zhang, "Learning efficient convolutional networks through network slimming," in *ICCV*, 2017, pp. 2755–2763.

21. H. Wang, J. Du, H. Zhou, H. Lu, and Y. Cao, "A novel approach to structured pruning of neural network for designing compact audio-visual wake word spotting system," in *The Asia-Pacific Signal and Information Processing Association (APSIPA)*, 2022.
22. A. Renda, J. Frankle, and M. Carbin, "Comparing rewinding and fine-tuning in neural network pruning," in *ICLR*, 2020.
23. A. Tulloch and Y. Jia, "High performance ultra-low-precision convolutions on mobile devices," *arXiv preprint arXiv:1712.02427*, 2017.
24. A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, and H. Adam, "Mobilenets: Efficient convolutional neural networks for mobile vision applications," *arXiv preprint arXiv:1704.04861*, 2017.
25. D. Haase and M. Amthor, "Rethinking depthwise separable convolutions: How intra-kernel correlations lead to improved mobilenets," in *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 14 588–14 597.
26. S. Han, J. Pool, J. Tran, and W. Dally, "Learning both weights and connections for efficient neural network," *Advances in neural information processing systems (NIPS)*, vol. 28, 2015.
27. J. Frankle and M. Carbin, "The lottery ticket hypothesis: Finding sparse, trainable neural networks," in *ICLR*, 2019.
28. H. Zhou, J. Du, G. Zou, Z. Nian, C. Lee, S. M. Siniscalchi, S. Watanabe, O. Scharenborg, J. Chen, S. Xiong, and J. Gao, "Audio-visual wake word spotting in misp2021 challenge: Dataset release and deep analysis," in *Interspeech*, 2022.