# A LIGHTWEIGHT DYNAMIC FILTER FOR KEYWORD SPOTTING

*Donghyeon Kim[1], Kyungdeuk ko[1], Jeonggi Kwak[1], David K. Han[2], Hanseok Ko[1]*

[1]Korea University, South Korea
[2]Drexel university, USA

## ABSTRACT

Keyword Spotting (KWS) from speech signals is widely applied to perform fully hands-free speech recognition. The KWS network is designed as a small-footprint model so it can continuously be active. Recent efforts have explored dynamic filter-based models in deep learning frameworks to enhance the system's robustness or accuracy. However, as a dynamic filter framework requires high computational costs, the implementation is limited to the computational condition of the device. In this paper, we propose a lightweight dynamic filter to improve the performance of KWS. Our proposed model divides the dynamic filter into two branches to reduce computational complexity: pixel level and instance level. The proposed lightweight dynamic filter is applied to the front end of KWS to enhance the separability of the input data. The experimental results show that our model is robustly working on unseen noise and small training data environments by using a small computational resource.

**Index Terms**: keyword spotting, dynamic filter, dynamic weight, computational cost

## 1. INTRODUCTION

Recently, deep learning based speech applications have been applied in real world applications [1, 2, 3, 4]. For smart devices on standby for human user's commands, Keyword Spotting (KWS) is an essential capability. In addition, KWS system can be leveraged for Over-The-Top (OTT) media service or smart TV to enhance accessibility between user and computer. As these systems listen continuously to the audio stream in the environment, their KWS process should cost minimally to conserve battery power. For this reason, model parameters and subsequent computational cost are important elements in evaluating the KWS system. This Conventional KWS has been based on Hidden Markov Models (HMMs) by using a large vocabulary speech recognizer and a background model [5, 6]. With the advent of deep learning framework, Multi-Layer Perceptron (MLP) and Convolutional Neural Network (CNN) have been shown to outperform the conventional methods [7, 8, 9]. As CNN requires high computations in general to extract spectral and temporal features, some proposed models focused to reduce computational cost. Depthwise Separable Convolution (DSConv) extracts channel-wise features separately and later combines channel information by using 1D convolution [10]. As an alternative to the low computational model, temporal convolution-based networks [11, 12] have been proposed to extract spectral features by performing convolution on the temporal axis. Neural Architecture Search (NAS)

methods [13, 14] have been also proposed to find the best network structure with a reduced computational requirement for the KWS model. Although these networks would reduce memory footprints and computational operations, their performances have been shown limited, particularly with unknown speakers or unseen noise.

An alternative to improve performance is integrating a filter on the front end of the process [15, 16]. In the conventional dynamic filter process, dynamic weights are generated by each input patch. In this case, computations would be required for all the input pixels. Such an extra process of filtering, however, adds to the computational cost, which defeats the purpose of reducing computations. To balance out the requirement between high performance with minimal computation, we propose a Lightweight Dynamic (LDy) convolution for the front-end filtering. In our proposed network, two branches of dynamic filters are implemented at the front end of the network: pixel-level and instance-level. A pixel-level dynamic filter performs pixel attention for Time-Frequency (T-F) features while an instance-level dynamic filter produces a global representative weight vector from the temporal pooling of the acoustic features. These two ways of dynamic weights are merged to conduct CNN process. We also leverage dynamic weights to conduct Instance normalization to the output of dynamic convolution. This is because conventional normalization methods [17, 18, 19], which use static weights, have limitations in addressing the noise robustness problem. Our proposed model is applied to the convolution process in a computationally efficient manner to deliver dynamic filtering robust to noisy environments. The KWS experiments are carried out on Speech command data [20] v1 and v2. Additionally, we utilize three different noise data [21, 22, 23], to set up the unseen noise environments. The experimental results show that the proposed lightweight dynamic filter improves KWS performance with robustness over recently developed methods. Especially our lightweight dynamic convolution only utilizes 220K Flops., and 2K parameters for its implementation.

## 2. RELATED WORKS

Dynamic Filter Network (DFN) is an adaptive deep neural network architecture [15] that generates filter parameters by neural networks. DFN consists of a filter generator and dynamic filter layers. The filter generator produces weights of the neural network, and the process of filtering occurs in the dynamic filter layer with the generated weights. This method conditionally adjusts weight vectors for given patch-level features. This learning method is simpler and works more efficiently over a self-attention network [24]. Kim *et al.* [16] proposed a convolution-based dynamic filter network to enhance salient features from the unseen noisy audio stream, and their experimental results showed that their approach outperformed conventional feature enhancement methods. Fujita *et al.* [25] also utilized a dynamic filter-based method for a lightweight ASR model. They confirmed that applying dynamic convolution to the decoder
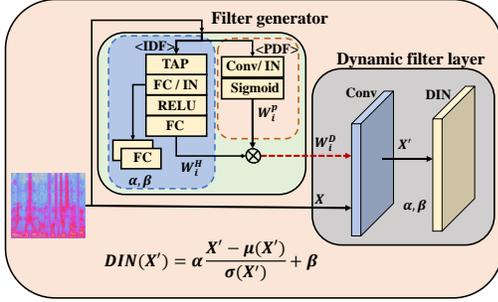
**Fig. 1**: Pipeline of Lightweight Dynamic Convolution process: PDF denotes Pixel Dynamic Filter; IDF denotes Instance-level Dynamic Filter. N denotes total pixels in the T-F feature

part in the encoder-decoder model improves accuracy and reduces computational load over transformer[26] based models. Although dynamic filter-based approaches have shown progressive performance improvements recently over other methods, they typically require high computational cost in filter implementation since the weights are produced by each patch basis. To alleviate this issue, we borrow the idea of Decoupled Dynamic Filter (DDF) [27] from computer vision. Instead of a single filter to take up channel-spatial features directly, DDF divides the filtering into two parts: channel and spatial. The channel-wise filter applies computationally efficient 1-D convolution while the spatial-wise side uses a simple global average pooling. These two filters are trained at separate branches of the network and they are then combined into forming a dynamic filter for input feature maps. By using simple computations, this method yields a significant reduction in memory and floating-point operations (FLOPs). In our adaptation of DDF, we split the dynamic filter process into Pixel Dynamic Filter (PDF) and Instance-level Dynamic Filter (IDF), and the process is described in Fig1. The PDF follows a spatial attention mechanism to capture pixel saliency and the IDF determines the direction of the dynamic weight vector adaptive to the input audio clip. The outputs of PDF and IDF are utilized to conduct the dynamic filter process.

## 3. LIGHTWEIGHT DYNAMIC FILTER

### 3.1. Pixel Dynamic Filter

The goal of PDF is to obtain an attention-based Time-Frequency (T-F) mask to measure a pixel level saliency. A single channel convolution is applied to T-F features ($x \in \mathcal{R}^{[T \times F]}$), and normalization with sigmoid activation is subsequently applied for mapping the feature scale from 0 to 1. The process of PDF is as follows:

$$w_i^P = \rho(IN(Conv(x_i, \mathbf{w}))), \quad (1)$$

where $Conv(x_i, \mathbf{w})$ denotes the convolution process with $x_i$ as the $i_{th}$ input data in mini-batch of the convolution and $\mathbf{w}$ represents the trainable weight vector ($\mathbf{w} \in \mathbb{R}^{[K \times K]}$, where $K$ denotes size of convolution kernel). Then Instance Normalization (IN) with sigmoid activation ($\rho(\cdot)$) are performed in series to obtain pixel level weights ($w_i^P \in \mathbb{R}^{[T \times F, 1]}$).

### 3.2. Instance-level Dynamic Filter

Conventional deep learning-based dynamic filters generate adaptive weights for the filter by learning them through the training process

of a Convolutional Neural Network (CNN). While this approach has shown to be flexible and adaptive in various domains of diverse speaker characteristics or noise, it comes at the cost of high computational loads.

As an efficient alternative to the convolution-based dynamic filter generation, we developed IDF which guides in capturing directional features of dynamic filter. IDF is aimed to produce an audio representative weight vector, which is similar to speaker embedding in speaker verification. To this end, the temporally averaged feature ($H \in \mathbb{R}^{[F]}$) that preserves the frequency trait gets fed to an MLP-based model, and the output is used as vector direction of dynamic weight. The process of IDF is as follows:

$$\mathbf{w}^H = max(0, IN(H \cdot \mathbf{w_1} + \mathbf{b_1})) \cdot \mathbf{w_2} + \mathbf{b_2}, \quad (2)$$

where $\mathbf{w_i}$ and $\mathbf{b_i}$ are learning parameters. Two layers of MLP with IN and relu activation extract a weight vector ($\mathbf{w}^H \in \mathbb{R}^{[1, K \times K]}$). In this way, the weight vector is generated uniquely per input audio signal, and the dynamic filters ($\mathbf{w}^H$ and $w_i^P$) are multiplied to perform dynamic convolution.

### 3.3. Dynamic Convolution

The process of dynamic convolution is as follows:

$$\mathbf{w}_i^D = w_i^P \odot \mathbf{w}^H, \quad (3)$$

$$x_i^D = IN(Conv(x_i, \mathbf{w}_i^D)), \quad (4)$$

$$x_i^{'} = x_i^D + x_i, \quad (5)$$

where $\mathbf{w}_i^D \in \mathbb{R}^{[T \times F, K \times K]}$ denotes dynamic weights for CNN computation. From equation 3, weight vector ($\mathbf{w}^H$) and scalar weight of pixel ($w_i^P$) perform element-wise product ($\odot$) to obtain the dynamic weights. Then, the convolution process using the dynamic weights is performed on the input and the result is normalized as in Equation 4. Through the skip connection, the normalized output ($x_i^D \in \mathbb{R}^{[T \times F]}$) is added to the input $x_i$ to be fed to the KWS model as shown in Equation 5. As such, the dynamic convolution-based filter generation adapts to each input effectively by enhancing features relevant to the KWS task. In conventional dynamic filters, a patch-level weight vector gets generated by a filter generator. Thus, the weight vector has different directions and magnitudes per patch. Our method, however, uniquely generates a unit directional vector for all the patches while the magnitude gets conditionally adjusted depending on each input patch. By decomposing the weight vector in terms of its direction and magnitude, the dimension of the weights becomes reduced significantly. This would significantly reduce the computational complexity of the dynamic filter.

### 3.4. Dynamic Instance Normalization

Normalization technique is an important part to achieve promising KWS performance [19]. However, leveraging a conventional normalization process, where statistic weights (scale and bias) are used, might not efficiently handle blind environments. From the motivation of the previous study [28], to address this, we leverage dynamic weights to implement feature normalization for the output of dynamic convolution. We firstly normalize the output of dynamic convolution by zero mean unit variance ($\mu_x$ and $\sigma_x$) and two different FC layers produce dynamic weights ($\alpha$ and $\beta \in \mathbb{R}^{[F]}$) respectively to adjust scale and bias. Here, an input of the linear layer is the output of the first linear layer in IDF. The process of Dynamic Instance Normalization (DIN) is as follows:

$$DIN(x_i^D) = \alpha(x_i^D - \mu_{x_i^D})/\sigma_{x_i^D} + \beta, \quad (6)$$

**Table 1**: Comparison of the parameter numbers, computation time, and memory cost. Conv, DyConv, and LDyConv stand for static filter-based convolution, Dynamic filter convolution, and Lightweight Dynamic convolution respectively.

| Filter | Conv | DyConv | LDyConv |
|--------|------|--------|---------|
| Parameter | $K^2$ | $K^4$ | $K^2(F+1)$ |
| Time | $O(K^2N)$ | $O(K^4N)$ | $O(K^2N)$ |
| Space | – | $O(K^2N)$ | $O(N^2+K)$ |

where $\mu_{x_i^D}$ and $\sigma_{x_i^D}$ denote mean and standard variation of the $x_i^D$. Instead of IN, DIN is applied to equation (4).

### 3.5. Computational Complexity

Our proposed dynamic convolution is applied to the front end of the KWS network for the dynamic feature extraction. A single-channel audio input gets fed to dynamic convolution and a single-channel output is employed as input for the KWS network. For this single-channel configuration, we compare computational costs (model parameters, time complexity, and space complexity) as summarized in Table 1. $N$ denotes the size of the pixel in the T-F feature which is equal to $[T \times F]$. For simplicity, we only compare the convolution process of the conventional Dynamic Convolution (DyConv) and our Lightweight Dynamic Convolution (LDyConv).

**Model parameters.** In DyConv, the $K^2$ dimensional weight vector is driven by the patch size $K^2$. Thus, $K^4$ of parameters are used for DyConv implementation. In LDyConv, the PDF produces $K^2$ scalar weights from the patch dimension and the IDF generates a weight vector from $F$ by using two FC layers. Thus, LDyConv requires $K^2(F+1)$ parameters. Since $K^2$ and $F$ are of the same order, the number of parameters required for DyConv and LDyConv are similar.

**Time complexity.** DyConv computes $K^4$ for every sample (pixel). With the total number of samples ($N$), DyConv has $O(K^4N)$ of time complexity. In LDyConv, the PDF computes $K^2$ for every sample while the IDF is only performed per audio clip. Thus, time complexity of LDyConv is $O(K^2(N+F))$ or it is approximately $O(K^2N)$ since $N >> F$. Our method would result in similar time complexity over the static filter.

**Space complexity** As the DyConv generates patch-level dynamic weight, it has $O(K^2N)$ of space complexity for saving dynamic weight vector. On the other hand, our dynamic filter only requires $O(K^2+N)$ of space complexity.

In summary, our proposed lightweight dynamic convolution model has similar model complexity compared to a static convolution filter model. We confirm that our proposed model consumes 2K parameters and 220K Flops when $K=3$, $F=40$, and $T=98$.

## 4. EXPERIMENT AND DISCUSSION

### 4.1. Experimental Setup

**Dataset.** We used speech command datasets v1 and v2 [20] for evaluating the KWS performance. By following the DB guideline of the dataset, we utilized 10 keywords with two extra classes (unknown or silent) for model training, injected background noise, and added random time-shifting.

**Noise setup.** For evaluating robustness against noise, we utilized DCASE [21], Urbansound8K [22] and WHAM [23] datasets. These three datasets contain background noise of urban locations. For data augmentation, we randomly selected an audio sample from the noise data and mixed it with the test audio of speech command with 5 different Signal-to-Noise ratios (SNRs) [20dB, 15dB, 10dB, 5dB, and 0dB].

**Acoustic feature extraction.** The acoustic feature we used is Mel Frequency Cepstral Coefficients (MFCC) constructed with 30ms of windows with 10ms overlap from an audio clip sampled at 16kHz. 64 Mel filters are employed to extract a Mel-spectrogram and 40 MFCC coefficients are extracted. This process gives a [40,98] size of audio features.

**Computation setup.** All our experiments are done by using the Tensorflow deep learning package with RTX-2080 ti GPU. In the training process, we used a batch size of 100, 30K iterations, and an ADAM optimizer with a 0.001 initial learning rate. For every 10K iteration, the learning rate is decreased by 0.1.

**Implementation detail.** In the PDF and the dynamic convolution process, we used $3 \times 3$ CNN kernel ($k=3$) dilated by (2,2) with a stride of 1. In the IDF, the first FC and second FC follow $40 \times 40$ and $40 \times k$ layer dimensions respectively. For DIN, two layers of FC which have $40 \times 40$ filter size respectively are utilized to produce $\alpha$ and $\beta$.

### 4.2. Baselines

Four different baseline architectures are used for comparisons. We implemented the proposed LDyConv on TENet architecture [12] and compared its performance with the following baseline models.

**TCNet.** TCNet [11] (or TC-Resnet) utilized temporal convolution and skip-connection for a fast and low computational cost model. TCNet8 contains 3 convolution blocks and 1 FC layer. Each convolution block has two layers of temporal convolution with a skip connection. Similarly, TCNet14 contains 6 convolution blocks and 1 FC layer.

**TENet.** TENet [12] utilizes a depth-separable convolution framework. A convolution block contains three convolutions with batch normalization. TENet6 has 6 convolution blocks and 1 FC layer. Each convolution block has two 1D convolutions and 1 temporal convolution. TENet12 contains 12 convolution blocks and 1 FC layer. TENet has 32 output channels for each convolution block and TENet-n has 16 output channels for each convolution block.

**MHA-RNN.** MHA-RNN[29] utilizes CRNN and self-attention mechanism for the KWS model training. The output of CRNN feeds to a dot product-based Multi-Head Attention (MHA) model, and two layers of FC produce probability values for KWS.

**Neural Architecture Search.** NAS is a network architecture designing method for deep learning applications and Differentiable Architecture Search (DARTS) is a variant of NAS that reduces search costs by weight sharing. We compared various state-of-art NAS methods in keyword spotting. Please see details of the model in [13, 14].

**BC-Resnet.** BroadCasted Resnet[30] uses BC block which contains frequency and temporal depth-wise convolution with a SubSpectral-Norm. In the BC block, the output of a 2D convolution is fed to the pooling layer and FC layer to represent audio features.

**Lightweight convolution.** Lightweight convolution [24] (Lconv) is a separable convolution method by using weight sharing and weight normalization. A single block of Lconv contains two linear layers, Gated Linear Unit (GLU) activation, and lightweight convolution. we apply the single Lconv block at the front end of the TENet12 model. In the first linear layer, frequency is expended by 80, and GLU is carried out to the frequency dimension. Then, lightweight convolution ($H=10$) and the other linear layer are computed by

**Table 2**: Comparison with lightweight models on Speech Command v1 and v2: Par. and Flops. denote Model parameters and computational cost respectively. Notation of † denotes the application of Spec-Augmentation [31]. For an accurate experiment, 8 times averaging accuracy and best performance are presented.

| Model | (Par.,Flops.) | V1 Acc | V1 Best | V2 Acc | V2 Best |
|---|---|---|---|---|---|
| TCNet8[11] | (145K,4.40M) | - | 96.2 | - | - |
| TCNet14[11] | (305K,8.26M) | - | 96.6 | 96.53 | 96.8 |
| TENet12[12] | (100K,6.42M) | - | 96.6 | 97.10 | 97.3 |
| TENet12†[12] | (100K,6.42M) | 97.19 | 97.3 | 97.43 | 97.6 |
| MHA-RNN†[29] | (743K,87.2M) | - | 97.2 | - | 98.0 |
| BC-ResNet3†[30] | (54.2K,32.4M) | 97.6 | - | 98.2 | - |
| BC-ResNet6†[30] | (188K,106.2M) | 97.9 | - | 98.6 | - |
| BC-ResNet8†[30] | (321K,178.2M) | 98.0 | - | 98.7 | - |
| NAS2[13] | (886K,-) | - | 97.2 | - | - |
| Random[14] | (196K,8.8M) | 96.58 | 96.8 | - | - |
| DARTS[14] | (93K,4.9M) | 96.63 | 96.9 | 96.92 | 97.1 |
| F-DARTS[14] | (188K,10.6M) | 96.70 | 96.9 | 97.1 | 97.4 |
| N-DARTS[14] | (109K,6.3M) | 96.79 | 97.2 | 97.18 | 97.4 |
| Lconv[24] | (105K,7.40M) | 96.88 | 97.0 | 97.24 | 97.3 |
| Dconv[24] | (107K,7.69M) | 96.89 | 97.1 | 96.26 | 97.4 |
| LDy-TENet12 (w/o DIN) | (102K,6.64M) | 96.95 | 97.1 | 97.35 | 97.6 |
| LDy-TENet12 (w/o DIN)† | - | 97.42 | 97.6 | 97.66 | 97.7 |
| LDy-TENet12 | (105K,6.97M) | 96.94 | 97.1 | 97.40 | 97.6 |
| LDy-TENet12† | - | 97.43 | 97.6 | 97.67 | 97.7 |

preserving temporal and spectral features. Additionally, by following [24], we perform Dynamic convolution (Dyconv). Instead of employing the static weight in the lightweight convolution, a single linear layer is used to produce weights for Lconv. In our implementation, the Lconv block requires 982K of Flops. and 4.9K of parameters. In the Dconv, 1373K of Flops. and 6.6K of parameters are used.

### 4.3. Result discussion

**Table 3**: Comparison with Unseen noise environment on Speech Command v1: experiment is performed on LDy-TENet12, TENet12, Lconv and Dconv models.

| Noise | SNR (dB) | LDy. w DIN | LDy. w/o DIN | Dconv | Lconv | TENet |
|---|---|---|---|---|---|---|
| DCASE | 20 | 97.08 | **97.17** | 96.90 | 97.07 | 97.10 |
| | 15 | 96.87 | **96.91** | 96.76 | 96.80 | 96.84 |
| | 10 | 96.02 | **95.98** | 95.73 | 95.92 | 95.77 |
| | 5 | **94.38** | 94.15 | 93.79 | 94.11 | 93.85 |
| | 0 | **90.64** | 90.42 | 89.18 | 89.53 | 89.20 |
| Urban | 20 | **96.36** | 96.34 | 96.27 | 96.25 | 96.34 |
| | 15 | **95.50** | 95.50 | 95.42 | 95.42 | 95.49 |
| | 10 | 93.99 | 94.15 | 93.74 | 93.76 | 93.49 |
| | 5 | **91.11** | 90.80 | 90.02 | 90.03 | 90.17 |
| | 0 | **82.32** | 81.39 | 79.73 | 80.00 | 80.26 |
| WHAM | 20 | **96.67** | 96.66 | 96.51 | 96.61 | 96.60 |
| | 15 | 95.97 | 95.93 | 95.86 | 95.94 | **95.98** |
| | 10 | 93.64 | **93.75** | 93.33 | 93.73 | 93.51 |
| | 5 | 89.74 | **89.93** | 89.49 | 90.11 | 89.32 |
| | 0 | **79.73** | 78.86 | 78.32 | 79.13 | 78.39 |

Tables 2, 3, and 4 summarize the KWS results on the Speech Command dataset v1 and v2. Our proposed method, Lconv block, and Dconv block are applied to the front end of the TENet12 model. For a more thorough model evaluation, we repeated the experiment 8 times.

**Small footprint KWS.** Table 2 shows small-footprint KWS performances over state-of-the-art methods. We compare the learning parameters, Flops., averaging accuracy, and the highest accuracy over the 8 repeated experiments. For the fair comparison, we perform

**Table 4**: Comparison with Unseen noise environment on Speech Command v2.

| Noise | SNR (dB) | LDy. w DIN | LDy. w/o DIN | Dconv | Lconv | TENet |
|---|---|---|---|---|---|---|
| DCASE | 20 | **97.24** | 97.01 | 96.80 | 96.94 | 96.79 |
| | 15 | **96.74** | 96.35 | 96.24 | 96.29 | 96.05 |
| | 10 | **96.06** | 95.49 | 95.28 | 95.23 | 95.18 |
| | 5 | **94.22** | 93.74 | 92.95 | 92.85 | 92.82 |
| | 0 | **90.08** | 89.05 | 87.71 | 87.66 | 87.36 |
| Urban | 20 | **96.55** | 96.22 | 96.04 | 96.16 | 96.06 |
| | 15 | **95.20** | 94.99 | 94.44 | 94.55 | 94.65 |
| | 10 | **93.80** | 93.45 | 92.61 | 92.41 | 92.67 |
| | 5 | **89.28** | 88.45 | 87.60 | 87.18 | 87.51 |
| | 0 | **81.77** | 80.25 | 78.39 | 78.16 | 78.24 |
| WHAM | 20 | **96.37** | 96.11 | 96.00 | 96.19 | 96.06 |
| | 15 | **95.72** | 95.17 | 94.87 | 94.88 | 94.99 |
| | 10 | **93.57** | 93.38 | 93.00 | 93.15 | 92.89 |
| | 5 | **88.91** | 88.46 | 87.99 | 87.83 | 87.41 |
| | 0 | **78.66** | 77.43 | 76.31 | 76.65 | 75.93 |

Spec-Augmentation [31] during the training, and the results are indicated by †. From the results, we confirm that our methods (LDy-TENet) show improved performance over TENet based model. Particularly, the LDy-TENet6-n shows similar performance over the TENet12 which is a 3 times larger model. Compared with the NAS-based models, our LDy-TENet12 achieves the best-averaging accuracy with low computational costs. The performance improvement is not significant between the TENet12 and the lightweight convolution models (Lcon and Donv). Additionally, the parameters and Flops. of Lconv and Dconv are higher than our method, while they take degraded performance. Although BC-Resnet-based models show improved performance over our method, these models require high computational resources (FLOPS.) for their implementation. On the other hand, our method can be worked by using small FLOPS.

**Unseen noise environment.** Tables 3 and 4 summarize KWS results on the unseen noisy environment with 3 different noise datasets and 5 different SNRs. The results of the Lconv and Dconv are less significant since they show similar performance with the TENet. On the other hand, our proposed model shows more robust performance over the baselines, and particularly the performance is improved when the SNR is low. Compared with the TENet, 2.1% (v1) and 3.5% (v2) performance improvements are shown in the Urban 0dB condition. Especially, leveraging DIN enhances the robustness of the model in the blind environment. DIN delivers performance improvements when 5 and 0 dB SNR environments

As a result, our proposed dynamic filter in the front of the network would enhance the performance of KWS in unseen noisy environments. Especially, as the model takes two main parts (PDF and IDF), it can be implemented with a small computational cost.

## 5. CONCLUSION

The main focus of this study was to develop a lightweight dynamic filter model for an acoustic feature extractor in Keyword spotting. We proposed the Lightweight Dynamic Convolution model which decomposes a dynamic filter into two parts (pixel and kernel) for alleviating the issues of computational cost and noise robustness. In addition, Dynamic Instance Normalization delivers performance improvement over noisy environments. The process has a small footprint and through the relevant experiments, it is shown fast compared to the conventional dynamic convolution method while retaining the adaptability of a dynamic filter. The experiments confirmed that our proposed lightweight model is robust on unseen noise over lightweight models.

# 6. REFERENCES

[1] Y. Lee, J. Min, D. K. Han, and H. Ko, "Spectro-temporal attention-based voice activity detection," *IEEE Signal Processing Letters*, vol. 27, pp. 131–135, 2020.

[2] S. Lee, D. K. Han, and H. Ko, "Multimodal emotion recognition fusion analysis adapting bert with heterogeneous feature unification," *IEEE Access*, vol. 9, pp. 94 557–94 572, 2021.

[3] G. Kim, D. K. Han, and H. Ko, "Specmix: A mixed sample data augmentation method for training withtime-frequency domain features," *arXiv preprint arXiv:2108.03020*, 2021.

[4] D. Kim, S. Park, D. K. Han, and H. Ko, "Multi-band cnn architecture using adaptive frequency filter for acoustic event classification," *Applied Acoustics*, vol. 172, p. 107579, 2021.

[5] Y. Benayed, D. Fohr, J. P. Haton, and G. Chollet, "Confidence measures for keyword spotting using support vector machines," in *2003 IEEE International Conference on Acoustics, Speech, and Signal Processing, 2003. Proceedings.(ICASSP'03).*, vol. 1. IEEE, 2003, pp. I–I.

[6] H. Ketabdar, J. Vepa, S. Bengio, and H. Bourlard, "Posterior based keyword spotting with a priori thresholds," in *Ninth International Conference on Spoken Language Processing*, 2006.

[7] G. Chen, C. Parada, and G. Heigold, "Small-footprint keyword spotting using deep neural networks," in *2014 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2014, pp. 4087–4091.

[8] T. N. Sainath and C. Parada, "Convolutional neural networks for small-footprint keyword spotting," in *Sixteenth Annual Conference of the International Speech Communication Association*, 2015.

[9] S. O. Arik, M. Kliegl, R. Child, J. Hestness, A. Gibiansky, C. Fougner, R. Prenger, and A. Coates, "Convolutional recurrent neural networks for small-footprint keyword spotting," *arXiv preprint arXiv:1703.05390*, 2017.

[10] Y. Zhang, N. Suda, L. Lai, and V. Chandra, "Hello edge: Keyword spotting on microcontrollers," *arXiv preprint arXiv:1711.07128*, 2017.

[11] S. Choi, S. Seo, B. Shin, H. Byun, M. Kersner, B. Kim, D. Kim, and S. Ha, "Temporal convolution for real-time keyword spotting on mobile devices," *arXiv preprint arXiv:1904.03814*, 2019.

[12] X. Li, X. Wei, and X. Qin, "Small-footprint keyword spotting with multi-scale temporal convolution," *arXiv preprint arXiv:2010.09960*, 2020.

[13] T. Mo, Y. Yu, M. Salameh, D. Niu, and S. Jui, "Neural architecture search for keyword spotting," *arXiv preprint arXiv:2009.00165*, 2020.

[14] B. Zhang, W. Li, Q. Li, W. Zhuang, X. Chu, and Y. Wang, "Autokws: Keyword spotting with differentiable architecture search," in *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2021, pp. 2830–2834.

[15] X. Jia, B. De Brabandere, T. Tuytelaars, and L. V. Gool, "Dynamic filter networks," in *Advances in neural information processing systems*, 2016, pp. 667–675.

[16] D. Kim, J. Park, D. K. Han, and H. Ko, "Dual stage learning based dynamic time-frequency mask generation for audio event classification," *Proc. Interspeech 2020*, pp. 836–840, 2020.

[17] J. L. Ba, J. R. Kiros, and G. E. Hinton, "Layer normalization," *arXiv preprint arXiv:1607.06450*, 2016.

[18] D. Ulyanov, A. Vedaldi, and V. Lempitsky, "Instance normalization: The missing ingredient for fast stylization," *arXiv preprint arXiv:1607.08022*, 2016.

[19] S. Chang, H. Park, J. Cho, H. Park, S. Yun, and K. Hwang, "Subspectral normalization for neural audio data processing," in *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2021, pp. 850–854.

[20] P. Warden, "Speech commands: A dataset for limited-vocabulary speech recognition," *arXiv preprint arXiv:1804.03209*, 2018.

[21] A. Mesaros, T. Heittola, and T. Virtanen, "Acoustic scene classification in dcase 2019 challenge: Closed and open set classification and data mismatch setups," 2019.

[22] J. Salamon, C. Jacoby, and J. P. Bello, "A dataset and taxonomy for urban sound research," in *Proceedings of the 22nd ACM international conference on Multimedia*, 2014, pp. 1041–1044.

[23] G. Wichern, J. Antognini, M. Flynn, L. R. Zhu, E. McQuinn, D. Crow, E. Manilow, and J. Le Roux, "Wham!: Extending speech separation to noisy environments," in *Proc. Interspeech*, Sep. 2019.

[24] F. Wu, A. Fan, A. Baevski, Y. N. Dauphin, and M. Auli, "Pay less attention with lightweight and dynamic convolutions," *arXiv preprint arXiv:1901.10430*, 2019.

[25] Y. Fujita, A. S. Subramanian, M. Omachi, and S. Watanabe, "Attention-based asr with lightweight and dynamic convolutions," in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 7034–7038.

[26] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," in *Advances in neural information processing systems*, 2017, pp. 5998–6008.

[27] J. Zhou, V. Jampani, Z. Pi, Q. Liu, and M.-H. Yang, "Decoupled dynamic filter networks," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 6647–6656.

[28] Y. Jing, X. Liu, Y. Ding, X. Wang, E. Ding, M. Song, and S. Wen, "Dynamic instance normalization for arbitrary style transfer," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 04, 2020, pp. 4369–4376.

[29] O. Rybakov, N. Kononenko, N. Subrahmanya, M. Visontai, and S. Laurenzo, "Streaming keyword spotting on mobile devices," *arXiv preprint arXiv:2005.06720*, 2020.

[30] B. Kim, S. Chang, J. Lee, and D. Sung, "Broadcasted residual learning for efficient keyword spotting," *arXiv preprint arXiv:2106.04140*, 2021.

[31] D. S. Park, W. Chan, Y. Zhang, C.-C. Chiu, B. Zoph, E. D. Cubuk, and Q. V. Le, "Specaugment: A simple data augmentation method for automatic speech recognition," *arXiv preprint arXiv:1904.08779*, 2019.