

DYNET: DYNAMIC CONVOLUTION FOR ACCELERATING CONVOLUTIONAL NEURAL NETWORKS

Yikang Zhang, Jian Zhang, Qiang Wang, Zhao Zhong

HUAWEI

{zhangyikang7, zhangjian157, zorro.zhongzhao}@huawei.com

ABSTRACT

Convolution operator is the core of convolutional neural networks (CNNs) and occupies the most computation cost. To make CNNs more efficient, many methods have been proposed to either design lightweight networks or compress models. Although some efficient network structures have been proposed, such as MobileNet or ShuffleNet, we find that there still exists redundant information between convolution kernels. To address this issue, we propose a novel **dynamic convolution** method to adaptively generate convolution kernels based on image contents. To demonstrate the effectiveness, we apply dynamic convolution on multiple state-of-the-art CNNs. On one hand, we can reduce the computation cost remarkably while maintaining the performance. For ShuffleNetV2/MobileNetV2/ResNet18/ResNet50, DyNet can reduce 37.0/54.7/67.2/71.3% FLOPs without loss of accuracy. On the other hand, the performance can be largely boosted if the computation cost is maintained. Based on the architecture MobileNetV3-Small/Large, DyNet achieves 70.3/77.1% Top-1 accuracy on ImageNet with an improvement of 2.9/1.9%. To verify the scalability, we also apply DyNet on segmentation task, the results show that DyNet can reduce 69.3% FLOPs while maintaining Mean IoU on segmentation task.

1 INTRODUCTION

Convolutional neural networks (CNNs) have achieved state-of-the-art performance in many computer vision tasks (Krizhevsky et al., 2012; Szegedy et al., 2013), and the neural architectures of CNNs are evolving over the years (Krizhevsky et al., 2012; Simonyan & Zisserman, 2014; Szegedy et al., 2015; He et al., 2016; Hu et al., 2018; Zhong et al., 2018a;b). However, modern high-performance CNNs often require a lot of computation resources to execute a large amount of convolution kernel operations. Aside from the accuracy, to make CNNs applicable on mobile devices, building lightweight and efficient deep models has attracting much more attention recently (Howard et al., 2017; Sandler et al., 2018; Howard et al., 2019; Zhang et al., 2018; Ma et al., 2018). These methods can be roughly categorized into two types: efficient network design and model compression. Representative methods for the former category are MobileNet (Howard et al., 2017; Sandler et al., 2018; Howard et al., 2019) and ShuffleNet (Ma et al., 2018; Zhang et al., 2018), which use depth-wise separable convolution and channel-level shuffle techniques to reduce computation cost. On the other hand, model compression-based methods tend to obtain a smaller network by compressing a larger network via pruning, factorization, mimic and quantization (Chen et al., 2015; Han et al., 2015a; Jaderberg et al., 2014; Lebedev et al., 2014; Ba & Caruana, 2014; Zhu et al., 2016).

Although some handcrafted efficient network structures have been designed, we observe that the significant correlations still exist among convolutional kernels, and introduce a large amount of redundant calculations. Moreover, these small networks are hard to compress. For example, Liu et al. (2019) compress MobileNetV2 to 124M, but the accuracy drops by 5.4% on ImageNet compared with MobileNetV2 (1.0). This implies that traditional compression methods cannot solve the inherent redundancy problem in CNNs well. We theoretically analyze this phenomenon and find that it is caused by the nature of conventional convolution, where correlated kernels are cooperated to extract noise-irrelevant features. Thus it is hard to compress the conventional convolution kernels without information loss. We also find that if we linearly fuse several fixed convolution kernels to generate one dynamic kernel based on the input, we can obtain the noise-irrelevant features without the cooperation of multiple kernels.

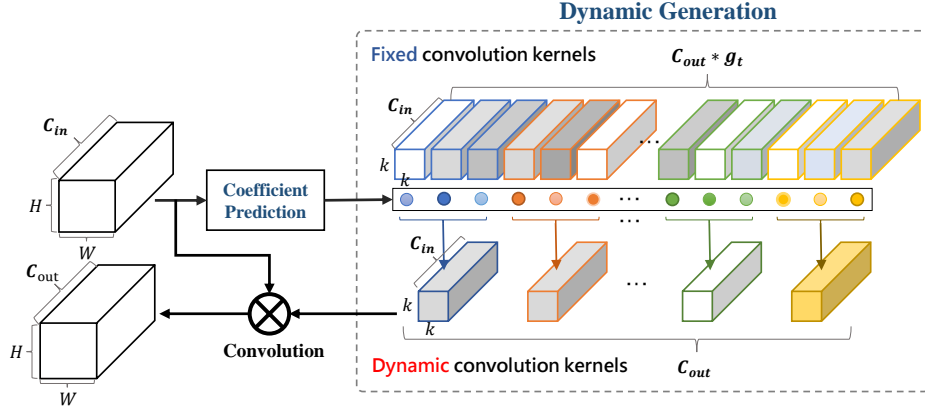


Figure 1: The overall framework of the dynamic convolution.

Based on the above observation and analysis, we propose dynamic convolution to address this issue, which adaptively generates convolution kernels based on image contents. The overall framework of dynamic convolution is shown in Figure 1, which consists of a *coefficient prediction module* and a *dynamic generation module*. The coefficient prediction module is trainable and designed to predict the coefficients of fixed convolution kernels. Then the dynamic generation module further generates a dynamic kernel based on the predicted coefficients.

Our proposed method is simple to implement and can be used as a drop-in plugin for any convolution layer to reduce redundancy. We evaluate the proposed dynamic convolution on state-of-the-art networks. On one hand, we can reduce the computation cost remarkably while maintaining the performance. For ShuffleNetV2 (1.0), MobileNetV2 (1.0), ResNet18 and ResNet50, DyNet reduces 37.0%, 54.7%, 67.2% and 71.3% FLOPs respectively while the Top-1 accuracy on ImageNet changes by +1.0%, -0.27%, -0.6% and -0.08%. On the other hand, the performance can be largely boosted if the computation cost is maintained. For MobileNetV3-Small(1.0) and MobileNetV3-Large(1.0), DyNet improve the Top-1 accuracy on ImageNet by 2.9% and 1.9% respectively while the FLOPs changes by +4.1% and +5.3%. Meanwhile, dynamic convolution further accelerates the inference speed of MobileNetV2 (1.0), ResNet18 and ResNet50 by $1.87\times$, $1.32\times$ and $1.48\times$ on CPU platform respectively.

2 RELATED WORK

We review related works from three aspects: efficient convolution neural network design, model compression, and dynamic convolution kernel.

2.1 EFFICIENT CONVOLUTION NEURAL NETWORK DESIGN

In many computer vision tasks (Krizhevsky et al., 2012; Szegedy et al., 2013), model design plays a key role. The increasing demands of high-quality networks on mobile/embedding devices have driven the study on efficient network design (He & Sun, 2015). For example, GoogleNet (Szegedy et al., 2015) increases the depth of networks with lower complexity compared to simply stacking convolution layers; SqueezeNet (Iandola et al., 2016) deploys a bottleneck approach to design a very small network; Xception (Chollet, 2017), MobileNet (Howard et al., 2017; Sandler et al., 2018) use depth-wise separable convolution to reduce computation and model size. ShuffleNet (Zhang et al., 2018; Ma et al., 2018) shuffle channels to reduce the computation of 1×1 convolution kernel and improve accuracy. MobileNetV3 Howard et al. (2019) are designed based on a combination of complementary search techniques. Despite the progress made by these efforts, we find that there still exists redundancy between convolution kernels and cause redundant computation. Dynamic convolution can reduce the redundant computation, thus complement those efficient networks.

2.2 MODEL COMPRESSION

Another trend to obtaining a small network is model compression. Factorization based methods (Jaderberg et al., 2014; Lebedev et al., 2014) try to speed up convolution operation by using tensor decomposition to approximate original convolution operation. Knowledge distillation based methods (Ba & Caruana, 2014; Romero et al., 2014; Hinton et al., 2015) learn a small network to mimic a larger teacher network. Pruning based methods (Han et al., 2015a;b; Wen et al., 2016; Liu et al., 2019) try to reduce computation by pruning the redundant connections or convolution channels. Compared with those methods, DyNet is more effective especially when the target network is already efficient enough. For example, in (Liu et al., 2019), they get a smaller model of 124M FLOPs by pruning the MobileNetV2, however, it drops the accuracy by 5.4% on ImageNet compared with the model with 291M FLOPs. Moreover, the pruned MobileNetV2 with 137M FLOPs in (Ye et al., 2020) drops the accuracy by 3.2% and the pruned ResNet50 with 2120M FLOPs in (Wang et al., 2019) drops the accuracy by 5.1%. While in DyNet, we can reduce the FLOPs of MobileNetV2 (1.0) from 298M to 129M with the accuracy drops only 0.27% and reduce the FLOPs of ResNet50 from 3980M to 1119M with the accuracy drops only 0.08%.

2.3 DYNAMIC CONVOLUTION KERNEL

Generating dynamic convolution kernel appears in both computer vision and natural language processing (NLP) tasks.

In computer vision domain, Klein et al. (Klein et al., 2015) and Brabandere et al. (Jia et al., 2016) directly generate convolution kernels via a linear layer based on the feature maps of previous layers. Because convolution kernels have a large number of parameters, the linear layer will be inefficient on the hardware. Our proposed method solves this problem by merely predicting the coefficients for linearly combining fixed kernels and achieve real speed up for CNN on hardware. This technique has been deployed in HUAWEI at the beginning of 2019 and the patent is filed in May 2019 as well. The attention paid by the academic community (Yang et al., 2019; Chen et al., 2019; 2020) demonstrates the great potential of this direction. In this paper, we derive insight into dynamic convolution from the perspective of 'noise-irrelevant feature' and conduct a correlation experiment to prove that the correlation among convolutional kernels can be largely reduced in DyNet.

In NLP domain, some works (Shen et al., 2018; Wu et al., 2019; Gong et al., 2018) incorporate context information to generate input-aware convolution filters which can be changed according to input sentences with various lengths. These methods also directly generate convolution kernels via a linear layer, etc. Because the size of CNN in NLP is smaller and the dimension of the convolution kernel is one, the inefficiency issue for the linear layer is alleviated. Moreover, Wu et al. (Wu et al., 2019) alleviate this issue utilizing the depthwise convolution and the strategy of sharing weight across layers. These methods are designed to improve the adaptivity and flexibility of language modeling, while our method aims to cut down the redundant computation cost.

3 DYNET: DYNAMIC CONVOLUTION IN CNNs

In this section, we first describe the motivation of dynamic convolution. Then we explain the proposed dynamic convolution in detail. Finally, we illustrate the architectures of our proposed DyNet.

3.1 MOTIVATION

As indicated in previous works (Han et al., 2015a;b; Wen et al., 2016; Liu et al., 2019), convolutional kernels are naturally correlated in deep models. For some of the well-known networks, we plot the distribution of the Pearson product-moment correlation coefficient between feature maps in Figure 2. Most existing works try to reduce correlations by compressing, however, it is hard to accomplish for efficient and small networks like MobileNets, even though the correlation is significant. **We think these correlations are vital for maintaining the performance because they are cooperated to obtain noise-irrelevant features.** We take face recognition as an example, where the pose or the illumination is not supposed to change the classification results. **Therefore, the feature maps will gradually become noise-irrelevant when they go deeper.** Based on the theoretical analysis in appendix A, we find this procedure needs the cooperation of multiple correlated kernels and we can get noise-irrelevant

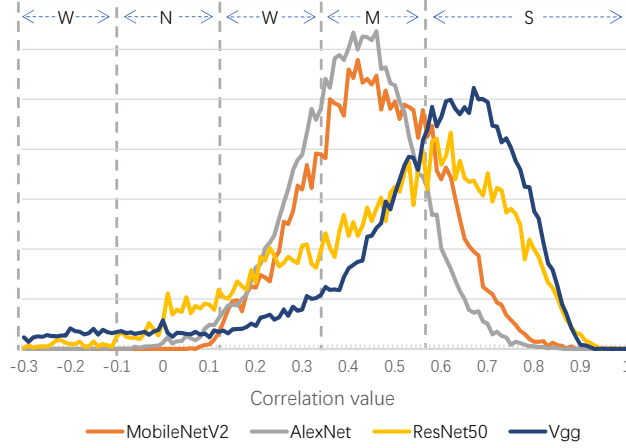


Figure 2: Pearson product-moment correlation coefficient between feature maps. S, M, W, N denote strong, middle, weak and no correlation respectively.

features without this cooperation if we dynamically fuse several kernels. In this paper, we propose a dynamic convolution method, which learns the coefficients to fuse multiple kernels into a dynamic one based on image contents. We give a more in-depth analysis of our motivation in appendix A.

3.2 DYNAMIC CONVOLUTION

The goal of dynamic convolution is to learn a group of kernel coefficients, which fuse multiple fixed kernels to a dynamic one. We illustrate the overall framework of dynamic convolution in Figure 1. We first utilize a trainable coefficient prediction module to predict coefficients. Then we further propose a dynamic generation module to fuse fixed kernels to a dynamic one. We will introduce the coefficient prediction module and dynamic generation module in detail in the following of this section.

Coefficient prediction module Coefficient prediction module is proposed to predict coefficients based on image contents. As shown in Figure 3, the coefficient prediction module can be composed of a global average pooling layer and a fully connected layer with Sigmoid as activation function. Global average pooling layer aggregates the input feature maps into a $1 \times 1 \times C_{in}$ vector, which serves as a feature extraction layer. Then the fully connected layer further maps the feature into a $1 \times 1 \times C$ vector, which are the coefficients for fixed convolution kernels of several dynamic convolution layers.

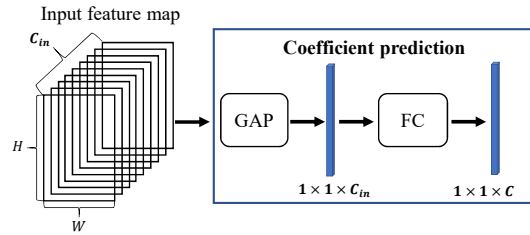


Figure 3: The coefficient prediction module.

Dynamic generation module For a dynamic convolution layer with weight $[C_{out} \times g_t, C_{in}, k, k]$, it corresponds with $C_{out} \times g_t$ fixed kernels and C_{out} dynamic kernels, the shape of each kernel is $[C_{in}, k, k]$. g_t denotes the group size, it is a hyperparameter. We denote the fixed kernels as w_t^i , the dynamic kernels as \tilde{w}_t , the coefficients as η_t^i , where $t = 0, \dots, C_{out}, i = 0, \dots, g_t$.

After the coefficients are obtained, we generate dynamic kernels as follows:

$$\tilde{w}_t = \sum_{i=1}^{g_t} \eta_t^i \cdot w_t^i \quad (1)$$

Training algorithm For the training of the proposed dynamic convolution, it is not suitable to use the batch-based training scheme. It is because the convolution kernel is different for different input images in the same mini-batch. Therefore, we fuse feature maps based on the coefficients rather than kernels during training. They are mathematically equivalent as shown in Eq. 2:

$$\begin{aligned}\tilde{O}_t &= \tilde{w}_t \otimes x = \sum_{i=1}^{g_t} (\eta_t^i \cdot w_t^i) \otimes x = \sum_{i=1}^{g_t} (\eta_t^i \cdot w_t^i \otimes x) \\ &= \sum_{i=1}^{g_t} (\eta_t^i \cdot (w_t^i \otimes x)) = \sum_{i=1}^{g_t} (\eta_t^i \cdot O_t^i),\end{aligned}\tag{2}$$

where x denotes the input, \tilde{O}_t denotes the output of dynamic kernel \tilde{w}_t , O_t^i denotes the output of fixed kernel w_t^i .

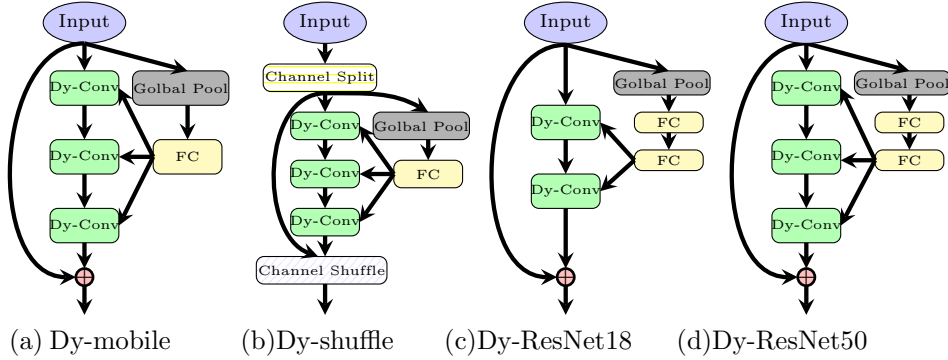


Figure 4: Basic building blocks for Dynamic Network variants of MobileNet, ShuffleNet, ResNet18, and ResNet50.

3.3 DYNAMIC CONVOLUTION NEURAL NETWORKS

We equip MobileNetV2, ShuffleNetV2, and ResNets with our proposed dynamic convolution, and propose Dy-mobile, Dy-shuffle, Dy-ResNet18, and Dy-ResNet50 respectively. The building blocks of these 4 networks are shown in Figure 4. Based on dynamic convolution, each dynamic kernel can get a noise-irrelevant feature without the cooperation of other kernels. Therefore we can reduce the channels for DyNets and remain the performance. We set the hyper-parameter g_t as 6 for all of them, and we give details of these dynamic CNNs below. To verify the performance can also be largely boosted if the number of channels is kept, we simply replace the convolution of MobileNetV3-Small(1.0) and MobileNetV3-Large(1.0) with the dynamic one to get the Dy-MobileNetV3-Small and Dy-MobileNetV3-Large.

Dy-mobile In our proposed Dy-mobile, we replace the original MobileNetV2 block with our dy-mobile block, which is shown in Figure 4 (a). The input of coefficient prediction module is the input of block, it produces the coefficients for all three dynamic convolution layers. Moreover, we further make two adjustments:

- We do not expand the channels in the middle layer like MobileNetV2. If we denote the output channels of the block as C_{out} , then the channels of all the three convolution layers will be C_{out} .
- Since the depth-wise convolution is efficient, we set $groups = \frac{C_{out}}{6}$ for the dynamic depth-wise convolution. We will enlarge C_{out} to make it becomes the multiple of 6 if needed.

After the aforementioned adjustments, the first dynamic convolution layer reduces the FLOPs from $6C^2HW$ to C^2HW . The second dynamic convolution layer keeps the FLOPs as $6CHW \times 3^2$ unchanged because we reduce the output channels by 6x while setting the groups of convolution

6x smaller, too. For the third dynamic convolution layer, we reduce the FLOPs from $6C^2HW$ to C^2HW as well. The ratio of FLOPs for the original block and our dy-mobile block is:

$$\frac{6C^2HW + 6CHW \times 3^2 + 6C^2HW}{C^2HW + 6CHW \times 3^2 + C^2HW} = \frac{6C + 27}{C + 27} = 6 - \frac{135}{C + 27} \quad (3)$$

Dy-shuffle In the original ShuffleNet V2, channel split operation will split feature maps to right-branch and left-branch, the right branch will go through one pointwise convolution, one depthwise convolution, and one pointwise convolution sequentially. We replace conventional convolution with dynamic convolution in the right branch as shown in Figure 4 (b). We feed the input of the right branch into coefficient prediction module to produce the coefficients. In our dy-shuffle block, we split channels into left-branch and right-branch with ratio 3 : 1, thus we reduce the 75% computation cost for two dynamic pointwise convolution. Similar to dy-mobile, we adjust the parameter "groups" in dynamic depthwise convolution to keep the FLOPs unchanged.

Dy-ResNet18/50 In Dy-ResNet18 and DyResNet50, we simply reduce half of the output channels for dynamic convolution layers of each residual block. Because the input channels of each block are large compared with dy-mobile and dy-shuffle, we use two linear layers as shown in Figure 4 (c) and Figure 4 (d) to reduce the number of parameters. If the input channel is C_{in} , the output channels of the first linear layer will be $\frac{C_{in}}{4}$ for Dy-ResNet18/50.

4 EXPERIMENTS

4.1 IMPLEMENTATION DETAILS

For the training of the proposed dynamic neural networks. Each image has data augmentation of randomly cropping and flipping, and is optimized with SGD strategy with cosine learning rate decay. We set batch size, initial learning rate, weight decay and momentum as 2048, 0.8, 5e-5 and 0.9 respectively. We also use the label smoothing with a rate of 0.1. We evaluate the accuracy on the test images with center crop.

4.2 EXPERIMENT SETTINGS AND COMPARED METHODS

We evaluate DyNet on ImageNet (Russakovsky et al., 2015), which contains 1.28 million training images and 50K validation images collected from 1000 different classes. We train the proposed networks on the training set and report the top-1 error on the validation set. To demonstrate the effectiveness, we compare the proposed dynamic convolution with state-of-the-art networks under mobile setting, including MobileNetV1 (Howard et al., 2017), MobileNetV2 (Sandler et al., 2018), ShuffleNet (Zhang et al., 2018), ShuffleNet V2 (Ma et al., 2018), Xception (Chollet, 2017), DenseNet (Huang et al., 2017), IGCv2 (Xie et al., 2018) and IGCv3 (Sun et al., 2018).

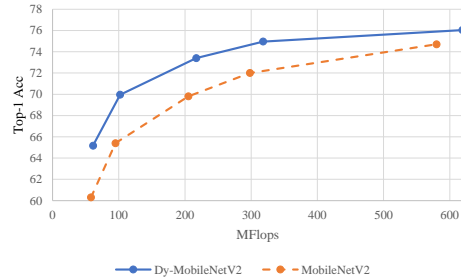


Figure 5: Compare with MobileNetV2 under the similar Flops constraint.

4.3 EXPERIMENT RESULTS AND ANALYSIS

Analysis of accuracy and computation cost We demonstrate the results in Table 1, where the number in the brackets indicates the channel number controller (Sandler et al., 2018). We partitioned the result table into three parts: (1) The original networks corresponding to the implemented dynamic networks; (2) Compared state-of-the-art networks under mobile settings; (3) The proposed dynamic networks.

Table 1 provides several valuable observations: (1) Compared with these well-known models under mobile setting, the proposed Dy-mobile, Dy-shuffle, and Dy-MobileNetV3 achieves the best

Table 1: Comparison of different network architectures over classification error and computation cost. The number in the brackets denotes the channel number controller (Sandler et al., 2018).

Methods	MFLOPs	Top-1 err. (%)
MobileNetV3-Small(1.0) (Howard et al., 2019)	56	32.60
ShuffleNet V2 (1.0) (Ma et al., 2018)	146	30.60
MobileNetV2 (1.0) (Sandler et al., 2018)	298	28.00
MobileNetV3-Large(1.0) (Howard et al., 2019)	219	24.8
ResNet18	1730	30.41
ResNet50	3890	23.67
ShuffleNet v1 (1.0) (Zhang et al., 2018)	140	32.60
MobileNet v2 (0.75) (Sandler et al., 2018)	145	32.10
MobileNet v2 (0.6) (Sandler et al., 2018)	141	33.30
MobileNet v1 (0.5)(Howard et al., 2017)	149	36.30
DenseNet (1.0) (Huang et al., 2017)	142	45.20
Xception (1.0) (Chollet, 2017)	145	34.10
IGCV2 (0.5) (Xie et al., 2018)	156	34.50
IGCV3-D (0.7) (Sun et al., 2018)	210	31.50
Dy-MobileNetV3-Small	59	29.7
Dy-shuffle (1.0)	92	29.6
Dy-mobile (1.0)	135	28.27
Dy-MobileNetV3-Large	228	22.9
Dy-ResNet18	567	31.01
Dy-ResNet50	1119	23.75

classification error with lowest computation cost. This demonstrates that the proposed dynamic convolution is a simple yet effective way to reduce computation cost. (2) Compared with the corresponding basic neural structures, the proposed Dy-shuffle (1.0), Dy-mobile (1.0), Dy-ResNet18 and Dy-ResNet50 reduce 37.0%, 54.7%, 67.2% and 71.3% computation cost respectively with little drop on Top-1 accuracy. This shows that even though the proposed network significantly reduces the convolution computation cost, the generated dynamic kernel can still capture sufficient information from image contents. (3) Compared with MobileNetV3-Small(1.0) and MobileNetV3-Large(1.0), Dy-MobileNetV3-Small and Dy-MobileNetV3-Large improve the Top-1 accuracy on ImageNet by 2.9% and 1.9% respectively with the FLOPs only increasing by 3M and 9M. The results also indicate that the performance can be largely boosted if the computation cost is maintained

Furthermore, we conduct detailed experiments on MobileNetV2. We replace the conventional convolution with the proposed dynamic one and get Dy-MobileNetV2. The accuracy of classification for models with different numbers of channels is shown in Figure 5. It is observed that Dy-MobileNetV2 consistently outperforms MobileNetV2 but the ascendancy is weakened with the increase of the number of channels.

Analysis of the dynamic kernel Aside from the quantitative analysis, we also demonstrate the redundancy of the generated dynamic kernels compared with conventional kernels in Figure 6. We calculate the correlation between 160 feature maps output by the 7th stage for the original MobileNetV2(1.0) and Dy-MobileNetV2 (1.0) based on the validation set. Note that Dy-MobileNetV2 (1.0) is different with Dy-mobile(1.0). Dy-MobileNetV2(1.0) keeps the channels of each layer the same as the original one while replacing the conventional convolution with dynamic convolution. As shown in Figure 6, we can observe that the correlation distribution of dynamic kernels have more values distributed between -0.1 and 0.2 compared with conventional convolution kernels, which indicates that the redundancy between dynamic convolution kernels are much smaller than the conventional convolution kernels.

Analysis of speed on the hardware We also analyze the inference speed of DyNet. We carry out experiments on the CPU platform (Intel(R) Core(TM) i7-7700 CPU @ 3.60GHz) with Caffe (Jia et al., 2014). We set the size of input as 224 and report the average inference time of 50 iterations. It is

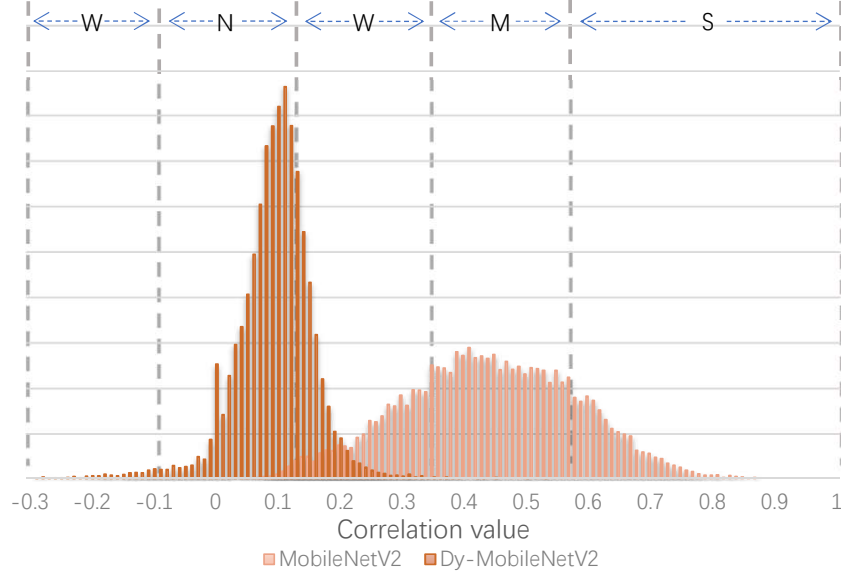


Figure 6: Pearson product-moment correlation coefficient between feature maps, S, M, W, N denote strong, middle, weak and no correlation respectively. We can observe that compared with conventional kernels, the generated dynamic kernels have small correlation values.

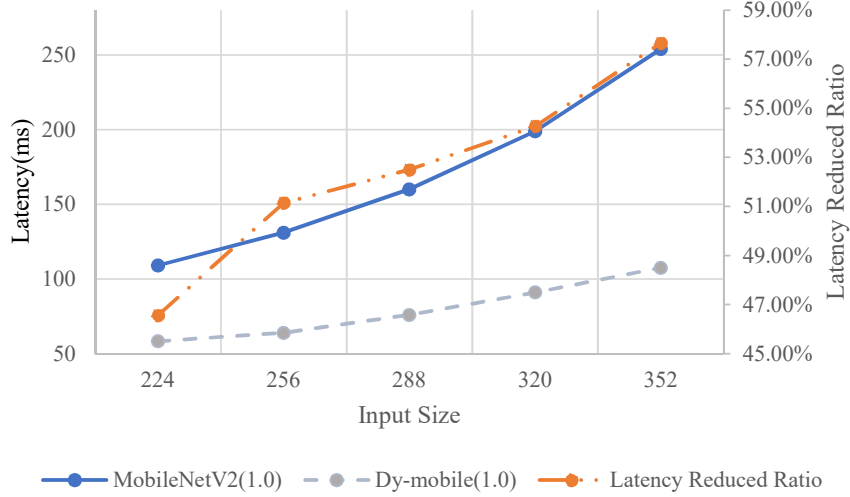


Figure 7: Latency for different input size. If we denote the latency of MobileNetV2(1.0), Dy-mobile as L_{Fix} and L_{Dym} , then Latency Reduced Ratio is defined as $100\% - \frac{L_{Dym}}{L_{Fix}}$.

reasonable to set mini-batch size as 1, which is consistent with most inference scenarios. The results are shown in Table 2. Moreover, the latency of fusing fixed kernels is independent with the input size, thus we expect to achieve a bigger acceleration ratio when the input size of networks becomes larger. We conduct experiments to verify this assumption, the results are shown in Figure 7. We can observe that the ratio of reduced latency achieved by DyNet gets bigger as the input size becomes larger. As shown in (Tan & Le, 2019), a larger input size can make networks perform significantly better, thus DyNet is more effective in this scenario.

We also analyze the training speed on the GPU platform. The model is trained with 32 NVIDIA Tesla V100 GPUs and the batch size is 2048. We report the average training time of one iteration in Table

2. It is observed that the training speed of DyNet is slower, it is reasonable because we fuse feature maps rather than kernels according to Eq. 2 in the training stage.

Table 2: Speed on the hardware.

Methods	Top-1 err. (%)	Inference Time	Training Time
MobileNetV2(1.0)	28.00	109.1ms	173ms
ResNet18	30.41	90.7ms	170ms
ResNet50	23.67	199.6ms	308ms
Dy-mobile(1.0)	28.27	58.3ms	250ms
Dy-ResNet18	31.01	68.7ms	213ms
Dy-ResNet50	23.75	135.1ms	510ms

4.4 EXPERIMENTS ON SEGMENTATION

To verify the scalability of DyNet on other tasks, we conduct experiments on segmentation. Compared to the method Dilated FCN with ResNet50 as backbone (Fu et al., 2018), Dilated FCN with Dy-ResNet50 reduces 69.3% FLOPs while maintaining the MIoU on Cityscapes validation set. The result is shown in Table 3.

Table 3: Experiments of segmentation on Cityscapes val set.

Methods	BaseNet	GFLOPs	Mean IoU%
Dilated FCN(Fu et al., 2018)	ResNet50	310.8	70.03
Dilated FCN(Fu et al., 2018)	Dy-ResNet50	95.6	70.48

4.5 ABLATION STUDY

Comparison between convolution with conventional kernel and dynamic kernel We correspondingly design two baseline networks for Dy-mobile (1.0) and Dy-shuffle (1.5), denoted as Fix-mobile(1.0) and Fix-shuffle (1.5) respectively. Specifically, we remove the coefficient prediction module and dynamic generation module, using fixed convolution kernels directly, the channel number is kept changeless. The results are shown in Table 4, compare with baseline networks Fix-mobile (1.0) and Fix-shuffle (1.5), the proposed Dy-mobile (1.0) and Dy-shuffle (1.5) achieve absolute classification improvements by 5.19% and 2.82% respectively. This shows that directly decreasing the channel number to reduce computation cost influences the classification performance a lot. While the proposed dynamic kernel can retain the representation ability as much as possible.

Table 4: Ablation experiments results of convolution with conventional kernel and dynamic kernel.

Methods	MParams	MFLOPs	Top-1 err. (%)
Fix-mobile (1.0)	2.16	129	33.57
Fix-shuffle (1.5)	2.47	171	30.30
Dy-mobile (1.0)	7.36	135	28.27
Dy-shuffle (1.5)	11.0	180	27.48

Effectiveness of g_t for dynamic kernel The group size g_t in Eq. 1 does not change the computation cost of DyNet but affects the performance of the network. Thus we provide an ablative study on g_t . We set g_t as 2,4,6 for dy-mobile(1.0) respectively and the results are shown in Table 5. The performance of dy-mobile(1.0) becomes better when g_t gets larger. It is reasonable because a larger g_t means the number of kernels cooperated for obtaining one noise-irrelevant feature becomes larger. When $g_t = 1$, the coefficient prediction module can be regarded as merely learning the attention for different channels, which can improve the performance of networks as well (Hu et al., 2018). Therefore we provide ablative study for comparing $g_t = 1$ and $g_t = 6$ on Dy-mobile(1.0) and

Table 5: Ablation experiments on g_t .

Methods	MParams	MFLOPs	Top-1 err. (%)
Fix-mobile(1.0)	2.16	129	33.57
Dy-mobile(1.0, $g_t = 2$)	3.58	131	29.43
Dy-mobile(1.0, $g_t = 4$)	5.47	133	28.69
Dy-mobile(1.0, $g_t = 6$)	7.36	135	28.27

Table 6: Comparison for $g_t = 1$ and $g_t = 6$.

Methods	MParams	MFLOPs	Top-1 err. (%)
Dy-mobile (1.0, $g_t = 1$)	2.64	131	30.85
Dy-mobile (1.0, $g_t = 6$)	7.36	135	28.27
Dy-ResNet18 ($g_t = 1$)	3.04	553	33.8
Dy-ResNet18 ($g_t = 6$)	16.6	567	31.01

Dy-ResNet18. The results are shown in Table 6. From the table we can see that, setting $g_t = 1$ will reduce the Top-1 accuracy on ImageNet for Dy-mobile(1.0) and Dy-ResNet18 by 2.58% and 2.79% respectively. It proves that the improvement of our proposed dynamic networks does not only come from the attention mechanism.

5 CONCLUSION

In this paper, we propose a novel dynamic convolution method to adaptively generate convolution kernels based on image content, which reduces the redundant computation cost existed in conventional convolution kernels. Based on the proposed dynamic convolution, we design several dynamic convolution neural networks based on well-known architectures. The experiment results show that DyNet can reduce FLOPs remarkably while maintaining the performance or boost the performance while maintaining the computation cost. As future work, we want to further explore the redundancy phenomenon existed in convolution kernels, and find other ways to reduce computation cost, such as dynamically aggregate different kernels for different images other than fixed groups used in this paper.

REFERENCES

- Jimmy Ba and Rich Caruana. Do deep nets really need to be deep? In *Advances in neural information processing systems*, pp. 2654–2662, 2014.
- Jin Chen, Xijun Wang, Zichao Guo, Xiangyu Zhang, and Jian Sun. Dynamic region-aware convolution. *arXiv preprint arXiv:2003.12243*, 2020.
- Wenlin Chen, James Wilson, Stephen Tyree, Kilian Weinberger, and Yixin Chen. Compressing neural networks with the hashing trick. In *International Conference on Machine Learning*, pp. 2285–2294, 2015.
- Yinpeng Chen, Xiyang Dai, Mengchen Liu, Dongdong Chen, Lu Yuan, and Zicheng Liu. Dynamic convolution: Attention over convolution kernels. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2019.
- François Chollet. Xception: Deep learning with depthwise separable convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1251–1258, 2017.
- Jun Fu, Jing Liu, Haijie Tian, Zhiwei Fang, and Hanqing Lu. Dual attention network for scene segmentation. 2018.
- Jingjing Gong, Xipeng Qiu, Xinchu Chen, Dong Liang, and Xuanjing Huang. Convolutional interaction network for natural language inference. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pp. 1576–1585, 2018.
- Song Han, Huizi Mao, and William J Dally. Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding. *arXiv preprint arXiv:1510.00149*, 2015a.
- Song Han, Jeff Pool, John Tran, and William Dally. Learning both weights and connections for efficient neural network. In *Advances in neural information processing systems*, pp. 1135–1143, 2015b.
- Kaiming He and Jian Sun. Convolutional neural networks at constrained time cost. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 5353–5360, 2015.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015.
- Andrew Howard, Mark Sandler, Grace Chu, Liang-Chieh Chen, Bo Chen, Mingxing Tan, Weijun Wang, Yukun Zhu, Ruoming Pang, Vijay Vasudevan, et al. Searching for mobilenetv3. In *Proceedings of the IEEE International Conference on Computer Vision*, pp. 1314–1324, 2019.
- Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861*, 2017.
- Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 7132–7141, 2018.
- Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 4700–4708, 2017.
- Forrest N Iandola, Song Han, Matthew W Moskewicz, Khalid Ashraf, William J Dally, and Kurt Keutzer. Squeezenet: Alexnet-level accuracy with 50x fewer parameters and 0.5 mb model size. *arXiv preprint arXiv:1602.07360*, 2016.
- Max Jaderberg, Andrea Vedaldi, and Andrew Zisserman. Speeding up convolutional neural networks with low rank expansions. *arXiv preprint arXiv:1405.3866*, 2014.

-
- Xu Jia, Bert De Brabandere, Tinne Tuytelaars, and Luc V Gool. Dynamic filter networks. In *Advances in Neural Information Processing Systems*, pp. 667–675, 2016.
- Yangqing Jia, Evan Shelhamer, Jeff Donahue, Sergey Karayev, Jonathan Long, Ross Girshick, Sergio Guadarrama, and Trevor Darrell. Caffe: Convolutional architecture for fast feature embedding. *arXiv preprint arXiv:1408.5093*, 2014.
- Benjamin Klein, Lior Wolf, and Yehuda Afek. A dynamic convolutional layer for short range weather prediction. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 4840–4848, 2015.
- Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pp. 1097–1105, 2012.
- Vadim Lebedev, Yaroslav Ganin, Maksim Rakhuba, Ivan Oseledets, and Victor Lempitsky. Speeding-up convolutional neural networks using fine-tuned cp-decomposition. *arXiv preprint arXiv:1412.6553*, 2014.
- Zechun Liu, Haoyuan Mu, Xiangyu Zhang, Zichao Guo, Xin Yang, Tim Kwang-Ting Cheng, and Jian Sun. Metapruning: Meta learning for automatic neural network channel pruning. *arXiv preprint arXiv:1903.10258*, 2019.
- Ningning Ma, Xiangyu Zhang, Hai-Tao Zheng, and Jian Sun. Shufflenet v2: Practical guidelines for efficient cnn architecture design. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 116–131, 2018.
- Adriana Romero, Nicolas Ballas, Samira Ebrahimi Kahou, Antoine Chassang, Carlo Gatta, and Yoshua Bengio. Fitnets: Hints for thin deep nets. *arXiv preprint arXiv:1412.6550*, 2014.
- Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115(3):211–252, 2015.
- Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mobilenetv2: Inverted residuals and linear bottlenecks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 4510–4520, 2018.
- Dinghan Shen, Martin Renqiang Min, Yitong Li, and Lawrence Carin. Learning context-sensitive convolutional filters for text processing. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pp. 1839–1848, 2018.
- Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- Ke Sun, Mingjie Li, Dong Liu, and Jingdong Wang. Igcv3: Interleaved low-rank group convolutions for efficient deep neural networks. *arXiv preprint arXiv:1806.00178*, 2018.
- Christian Szegedy, Alexander Toshev, and Dumitru Erhan. Deep neural networks for object detection. In *Advances in neural information processing systems*, pp. 2553–2561, 2013.
- Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1–9, 2015.
- Mingxing Tan and Quoc V. Le. Efficientnet: Rethinking model scaling for convolutional neural networks. 2019.
- Zongyue Wang, Shaohui Lin, Jiao Xie, and Yangbin Lin. Pruning blocks for cnn compression and acceleration via online ensemble distillation. *IEEE Access*, 7:175703–175716, 2019.
- Wei Wen, Chunpeng Wu, Yandan Wang, Yiran Chen, and Hai Li. Learning structured sparsity in deep neural networks. In *Advances in neural information processing systems*, pp. 2074–2082, 2016.

-
- Felix Wu, Angela Fan, Alexei Baevski, Yann N Dauphin, and Michael Auli. Pay less attention with lightweight and dynamic convolutions. *arXiv preprint arXiv:1901.10430*, 2019.
- Guotian Xie, Jingdong Wang, Ting Zhang, Jianhuang Lai, Richang Hong, and Guo-Jun Qi. Interleaved structured sparse convolutional neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 8847–8856, 2018.
- Brandon Yang, Gabriel Bender, Quoc V Le, and Jiquan Ngiam. Condconv: Conditionally parameterized convolutions for efficient inference. In *Advances in Neural Information Processing Systems*, pp. 1305–1316, 2019.
- Mao Ye, Chengyue Gong, Lizhen Nie, Denny Zhou, Adam Klivans, and Qiang Liu. Good subnetworks provably exist: Pruning via greedy forward selection. *arXiv preprint arXiv:2003.01794*, 2020.
- Xiangyu Zhang, Xinyu Zhou, Mengxiao Lin, and Jian Sun. Shufflenet: An extremely efficient convolutional neural network for mobile devices. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 6848–6856, 2018.
- Zhao Zhong, Junjie Yan, Wei Wu, Jing Shao, and Cheng-Lin Liu. Practical block-wise neural network architecture generation. In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 2423–2432, 2018a.
- Zhao Zhong, Zichen Yang, Boyang Deng, Junjie Yan, Wei Wu, Jing Shao, and Cheng-Lin Liu. Block-qnn: Efficient block-wise neural network architecture generation. *arXiv preprint arXiv:1808.05584*, 2018b.
- Chenzhuo Zhu, Song Han, Huizi Mao, and William J Dally. Trained ternary quantization. *arXiv preprint arXiv:1612.01064*, 2016.

A APPENDIX

A.1 DETAILED ANALYSIS OF OUR MOTIVATION

We illustrate our motivation from a convolution with output $f(x)$, i.e.,

$$f(x) = x \otimes w, \quad (4)$$

where \otimes denotes the convolutional operator, $x \in R^n$ is a vectorized input and $w \in R^n$ means the filter. Specifically, the i_{th} element of the convolution output $f(x)$ is calculated as:

$$f_i(x) = \langle x_{(i)}, w \rangle, \quad (5)$$

where $\langle \cdot, \cdot \rangle$ provides an inner product and $x_{(i)}$ is the circular shift of x by i elements. We define the index i started from 0.

We denote the noises in $x_{(i)}$ as $\sum_{j=0}^{d-1} \alpha_j y_j$, where $\alpha_j \in R$ and $\{y_0, y_1, \dots, y_{d-1}\}$ are the base vectors of noise space Ψ . Then the kernels in one convolutional layer can be represented as $\{w_0, w_1, \dots, w_c\}$. The space expanded by $\{w_0, w_1, \dots, w_c\}$ is Ω . We can prove if the kernels are trained until $\Psi \subset \Omega$, then for each $w_k \notin \Psi$, we can get the noise-irrelevant $f_i(x^{white}) = \langle x_{(i)}^{white}, w_k \rangle$ by the cooperation of other kernels w_0, w_1, \dots .

Firstly $x_{(i)}$ can be decomposed as:

$$x_{(i)} = \bar{x}_{(i)} + \beta w_k + \sum_{j=0}^{d-1} \alpha_j y_j, \quad (6)$$

where $\beta \in R$ and $\bar{x} \in R^n$ is vertical to w_k and y_j .

For concision we assume the norm of w_k and y_j is 1. Then,

$$f_i(x) = \langle x_{(i)}, w_k \rangle = \langle \bar{x}_{(i)} + \beta w_k + \sum_{j=0}^{d-1} \alpha_j y_j, w_k \rangle = \beta \langle w_k, w_k \rangle + \sum_{j=0}^{d-1} \alpha_j \langle y_j, w_k \rangle \quad (7)$$

When there is no noise, i.e. $\alpha_j = 0$ for $j = 0, 1, \dots, d-1$, the white output $f_i(x^{white})$ becomes:

$$f_i(x^{white}) = \langle x_{(i)}^{white}, w_k \rangle = \langle \bar{x}_{(i)} + \beta w_k, w_k \rangle = \beta \langle w_k, w_k \rangle = \beta. \quad (8)$$

It is proved in the Appendix A.2 that:

$$f_i(x^{white}) = \langle a_{00} w_k + \sum_t \beta_t w_t, x_{(i)} \rangle = (a_{00} + \beta_k) \langle w_k, x_{(i)} \rangle + \sum_{t \neq k} \beta_t \langle w_t, x_{(i)} \rangle, \quad (9)$$

where β_0, \dots, β_c is determined by the input image.

Eq. 9 is fulfilled by linearly combine convolution output $\langle w_k, x_{(i)} \rangle$ and $\langle w_t, x_{(i)} \rangle$ for those $\beta_t \neq 0$ in the following layers. Thus if there are N coefficients in Eq. 9 that are not 0, then we need to carry out N times convolution operation to get the noise-irrelevant output of kernel w_t , this causes redundant calculation.

In Eq. 9, we can observe that the computation cost can be reduced to one convolution operation by linearly fusing those kernels to a dynamic one:

$$\begin{aligned} \tilde{w} &= (a_{00} + \beta_k) w_k + \sum_{t \neq k, \beta_t \neq 0} \beta_t w_t \\ f_i(x^{white}) &= \langle \tilde{w}, x_{(i)} \rangle. \end{aligned} \quad (10)$$

In Eq. 10, the coefficients β_0, β_1, \dots is determined by $\alpha_0, \alpha_1, \dots$, thus they should be generated based on the input of network. *This is the motivation of our proposed dynamic convolution.*

A.2 PROVING OF EQ. 9

We denote $g_{ij}(x)$ as $\langle x_{(i)}, y_j \rangle$, $j = 0, 1, \dots, d-1$. Then,

$$g_{ij}(x) = \langle x_{(i)}, y_j \rangle = \langle \bar{x}_{(i)} + \beta w_k + \sum_{t=0}^{d-1} \alpha_t y_t, y_j \rangle = \beta \langle w_k, y_j \rangle + \sum_{t=0}^{d-1} \alpha_t \langle y_t, y_j \rangle. \quad (11)$$

By summarize Eq. 7 and Eq. 11, we get the following equation:

$$\begin{bmatrix} \langle w_k, w_k \rangle & \langle y_0, w_k \rangle & \langle y_1, w_k \rangle & \dots & \langle y_{d-1}, w_k \rangle \\ \langle w_k, y_0 \rangle & \langle y_0, y_0 \rangle & \langle y_1, y_0 \rangle & \dots & \langle y_{d-1}, y_0 \rangle \\ \langle w_k, y_1 \rangle & \langle y_0, y_1 \rangle & \langle y_1, y_1 \rangle & \dots & \langle y_{d-1}, y_1 \rangle \\ \vdots & \vdots & \vdots & \dots & \vdots \\ \langle w_k, y_{d-1} \rangle & \langle y_0, y_{d-1} \rangle & \dots & \dots & \langle y_{d-1}, y_{d-1} \rangle \end{bmatrix} \begin{bmatrix} \beta \\ \alpha_0 \\ \alpha_1 \\ \vdots \\ \alpha_{d-1} \end{bmatrix} = \begin{bmatrix} f_i(x) \\ g_{i0}(x) \\ g_{i1}(x) \\ \vdots \\ g_{i(d-1)}(x) \end{bmatrix}, \quad (12)$$

We simplify this equation as:

$$A\vec{x} = \vec{b}. \quad (13)$$

Because $w_k \notin \Psi$, we can denote w_k as:

$$w_k = \gamma_{\perp} w_{\perp} + \sum_{j=0}^{d-1} \gamma_j y_j, \quad (14)$$

where w_{\perp} is vertical to y_0, \dots, y_{d-1} and $\gamma_{\perp} \neq 0$.

moreover because $|w_k| = 1$, thus

$$|\gamma_{\perp}|^2 + \sum_{j=0}^{d-1} |\gamma_j|^2 = 1. \quad (15)$$

It can be easily proved that:

$$A = \begin{bmatrix} 1 & \gamma_0 & \gamma_1 & \dots & \gamma_{d-1} \\ \gamma_0 & 1 & 0 & \dots & 0 \\ \gamma_1 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \dots & \vdots \\ \gamma_{d-1} & 0 & \dots & \dots & 1 \end{bmatrix}. \quad (16)$$

thus,

$$\begin{aligned} |A| &= \begin{vmatrix} 1 & \gamma_0 & \gamma_1 & \dots & \gamma_{d-1} \\ \gamma_0 & 1 & 0 & \dots & 0 \\ \gamma_1 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \dots & \vdots \\ \gamma_{d-1} & 0 & \dots & \dots & 1 \end{vmatrix} \\ &= \begin{vmatrix} 1 - \gamma_0^2 & 0 & \gamma_1 & \dots & \gamma_{d-1} \\ \gamma_0 & 1 & 0 & \dots & 0 \\ \gamma_1 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \dots & \vdots \\ \gamma_{d-1} & 0 & \dots & \dots & 1 \end{vmatrix} \\ &= \begin{vmatrix} 1 - \gamma_0^2 - \gamma_1^2 & 0 & 0 & \dots & \gamma_{d-1} \\ \gamma_0 & 1 & 0 & \dots & 0 \\ \gamma_1 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \dots & \vdots \\ \gamma_{d-1} & 0 & \dots & \dots & 1 \end{vmatrix} \\ &= \begin{vmatrix} 1 - \gamma_0^2 - \gamma_1^2 - \dots - \gamma_{d-1}^2 & 0 & 0 & \dots & 0 \\ \gamma_0 & 1 & 0 & \dots & 0 \\ \gamma_1 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \dots & \vdots \\ \gamma_{d-1} & 0 & \dots & \dots & 1 \end{vmatrix} \\ &= \begin{vmatrix} \gamma_{\perp}^2 & 0 & 0 & \dots & 0 \\ \gamma_0 & 1 & 0 & \dots & 0 \\ \gamma_1 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \dots & \vdots \\ \gamma_{d-1} & 0 & \dots & \dots & 1 \end{vmatrix} \\ &= \gamma_{\perp}^2 \neq 0. \end{aligned} \quad (17)$$

thus,

$$\vec{x} = A^{-1}\vec{b}. \quad (18)$$

If we denote the elements of the first row of A^{-1} as $a_{00}, a_{01}, \dots, a_{0d}$, then

$$\begin{aligned} f_i(x^{white}) &= \beta = a_{00}f_i(x) + \sum_{j=0}^{d-1} a_{0(j+1)}g_{i,j}(x) \\ &= a_{00}\langle w_k, x_{(i)} \rangle + \sum_{j=0}^{d-1} a_{0(j+1)}\langle y_j, x_{(i)} \rangle \\ &= \langle a_{00}w_k + \sum_{j=0}^{d-1} a_{0(j+1)}y_j, x_{(i)} \rangle. \end{aligned} \quad (19)$$

Because $\Psi \subset \Omega$, there exists $\{\beta_t \in R | t = 0, 1, \dots, c\}$ that

$$\sum_{j=0}^{d-1} a_{0(j+1)}y_j = \sum_t \beta_t w_t. \quad (20)$$

Then,

$$f_i(x^{white}) = \langle a_{00}w_k + \sum_t \beta_t w_t, x_{(i)} \rangle = (a_{00} + \beta_k)\langle w_k, x_{(i)} \rangle + \sum_{t \neq k} \beta_t \langle w_t, x_{(i)} \rangle, \quad (21)$$