

Block Walsh-Hadamard Transform Based Binary Layers in Deep Neural Networks

HONGYI PAN, DIAA BADAWI, and AHMET ENIS CETIN, Department of Electrical and Computer Engineering, University of Illinois Chicago, USA

Convolution has been the core operation of modern deep neural networks. It is well-known that convolutions can be implemented in the Fourier Transform domain. In this paper, we propose to use binary block Walsh-Hadamard transform (WHT) instead of the Fourier transform. We use WHT-based binary layers to replace some of the regular convolution layers in deep neural networks. We utilize both one-dimensional (1-D) and two-dimensional (2-D) binary WHTs in this paper. In both 1-D and 2-D layers, we compute the binary WHT of the input feature map and denoise the WHT domain coefficients using a nonlinearity which is obtained by combining soft-thresholding with the tanh function. After denoising, we compute the inverse WHT. We use 1D-WHT to replace the 1×1 convolutional layers, and 2D-WHT layers can replace the 3×3 convolution layers and Squeeze-and-Excite layers. 2D-WHT layers with trainable weights can be also inserted before the Global Average Pooling (GAP) layers to assist the dense layers. In this way, we can reduce the number of trainable parameters significantly with a slight decrease in trainable parameters. In this paper, we implement the WHT layers into MobileNet-V2, MobileNet-V3-Large, and ResNet to reduce the number of parameters significantly with negligible accuracy loss. Moreover, according to our speed test, the 2D-FWHT layer runs about 24 times as fast as the regular 3×3 convolution with 19.51% less RAM usage in an NVIDIA Jetson Nano experiment.

Additional Key Words and Phrases: fast Walsh-Hadamard transform, block division, smooth-thresholding, image classification

1 INTRODUCTION

Recently, deep convolution neural networks (CNNs) have enjoyed a great success in many important applications such as image classification [1–10], object detection [11–14] and semantic segmentation [15–19]. While more and more parameters are needed in deep neural networks, deploying them on real-time resource-constrained environments such as embedded devices becomes a difficult task due to insufficient memory and limited computational capacity. To overcome the aforementioned limitations, smaller and computationally efficient neural networks are necessary for many practical applications.

Although 1×1 convolutions reduce the computational load, they are still computationally expensive and time-consuming in regular deep neural networks. In our previous paper [20], we proposed a binary layer based on Fast Walsh-Hadamard transform (FWHT) to replace the 1×1 convolution layer. We revised MobileNet-V2 [21] with the FWHT layer, and the new network is remarkably more slimmed and computationally efficient compared to the original structure according to our experiments. However, we can take advantage of the fast $O(m \log_2 m)$ Walsh-Hadamard Transform (WHT) algorithm only when m is an integer power of 2. As a result, we pad zeros to the end of the input vector in the FWHT layer to make the size of the vector an integer power of 2. We notice that padding zeros to increase the size of the Walsh-Hadamard Transform (WHT) brings some redundant parameters. In this paper, we introduce both one-dimensional (1D) and 2-Dimensional (2D) Walsh-Hadamard Transform layers. The new version of the 1D WHT computes the WHT in small blocks and avoids the zero-padding operations. As a result it needs fewer parameters compared to our earlier paper [20]. Our 2D FWHT layer can replace the 3×3 convolution layers and Squeeze-and-Excite

This work is funded by an award from the University of Illinois Chicago Discovery Partners Institute Seed Funding Program and NSF grants 1739396 and 1934915.

Authors' address: Hongyi Pan, hpan21@uic.edu; Diaa Badawi, dbadaw2@uic.edu; Ahmet Enis Cetin, aecyy@uic.edu, Department of Electrical and Computer Engineering, University of Illinois Chicago, 1200 W Harrison St, Chicago, Illinois, USA, 60607.

Manuscript submitted to ACM

layers. It can also be employed to assist the dense layers of a given network. Our contribution can be summarized as follows:

- We have a new 1D-FWHT layer. It improves the results of our previous work [20] by using Blocks of Walsh-Hadamard Transform (BWHT) instead of a single WHT computation to replace an entire 1×1 convolution layer. Compared to the old 1D-FWHT layer, the new 1D-BWHT layer retains the accuracy loss with a little more parameter reduction because it eliminates the zero-padding operation required to make the transform size a power of 2. For example, our version of the MobileNet-V2 network with the BWHT layers reaches a 0.08% higher accuracy with 11,919 fewer parameters compared to the MobileNet-V2 network with the FWHT layer on the CIFAR-10 dataset.
- We introduce the two-dimensional (2D) Walsh-Hadamard Transform layer in this paper. It can be used to replace the “Squeeze-and-Excite” layers and the regular 2D convolutional layers. For example, we reduce the number of trainable parameters by 48.62% with only a 0.76% accuracy loss in our MobileNet-V3-Large structure with 2D-WHT layers on the CIFAR-100 dataset. Our version of RESNET-34 with 2D-WHT layers has 53.57% fewer parameters than the regular RESNET-34 with only a 0.72% accuracy loss on the tiny ImageNet dataset.
- In addition, we introduce a weighted 2D-FWHT layer that can be easily inserted before the global average pooling (GAP) layer (or the flatten layer) to assist the dense layers. This novel layer improves the accuracy of the network with a slight (almost negligible) increase in parameters due to the additional weights. For example, our additional weighted 2D-FWHT layer improves the accuracy of the ResNet-20 network by 0.5% with only 256 additional parameters in the CIFAR-10 dataset.

This 2D-FWHT layer improves the accuracy of the MobileNet-V3-Large network by 0.15% in the CIFAR-100 dataset. In this case, only 704 more parameters are required.

Our novel binary layers do not increase the processing time even in conventional processors. For example, our 2D-FWHT layer runs about 24 times as fast as the regular 3×3 convolution layer on NVIDIA Jetson Nano.

2 RELATED WORK

Efficient neural network models include compressing a large neural network using quantization [22, 23], hashing [24], pruning [25], vector quantization [26] and Huffman encoding [27]. Another approach is the SqueezeNet [28], which is designed as a small network with 1×1 convolutional filters. Yet another approach is to use binary weights in neural networks [29–42]. Since the weights of the neurons are binary, they can be used to slim and accelerate networks in specialized hardware including compute-in-memory systems [43].

In 2020, Maneesh *et al.* proposed a trimmed version of MobileNet architecture called Reduced Mobilenet-V2 [44]. They replace bottleneck layers with heterogeneous kernel-based convolutions (HetConv) blocks. HetConv was first proposed by Pravendra *et al.* in 2019 [45]. Although HetConv reduced the parameters effectively, it is still based on the convolution, and we experimentally showed that our approach can outperform the network described in [44].

In 2021, James *et al.* proposed a novel layer based on the Fast Fourier Transform (FFT) called FNet for natural language processing applications. They insert the FFT in hidden layers and they can train the weights in the frequency domain (Fourier domain). They showed that their FNet has a faster speed with the guaranteed accuracy on BERT counterparts on the GLUE benchmark. However, the main weakness of the FNet is that their proposed method is only based on the real part of the FFT. After applying the FFT, they only keep the real part to avoid complex arithmetic. The imaginary

part of the FFT is simply ignored. Therefore, the information in the imaginary part is lost. Since the WHT is a real transform we do not suffer from information loss as a result of WHT computation.

Other methods using the Hadamard transform include [42, 46] but they did not perform any "convolutional" and non-linear filtering in the Hadamard transform domain. The "convolutional" filtering in the transform domain is possible by introducing multiplicative weights and non-linear filtering is possible with the use of two-sided smooth-thresholding which "denoises" small valued transform domain coefficients [20]. The use of one-sided RELU will lead to a significant information loss because transform domain coefficients can take both positive and negative values even if the input to the transform consists of all positive numbers. Our novel smooth thresholding nonlinearity is a two-sided version of the RELU and it can retain both positive and negative large amplitude coefficients while eliminating small amplitude ones. After performing the filtering operations in the transform domain we compute the inverse transform and continue processing the data in the feature map domain, which is a unique feature of our paper.

3 REVIEW OF THE FAST WALSH-HADAMARD (FWH) TRANSFORM

The Walsh-Hadamard Transform (WHT) is an example of a generalized class of the Fourier transforms. The transform matrix consists of +1 and -1 only. The convolution in the time (feature) domain leads to multiplication in the frequency domain in Fourier Transform. The Hadamard transform can be considered as a simplified version of the Fourier and wavelet transforms [47]. Therefore, we can approximately implement 1×1 and 3×3 convolutions in the WHT domain. We can not only train multiplicative weights but also train the threshold values of the soft and smooth-thresholds using the backpropagation algorithm [20]. We experimentally observed that this approach approximates 1×1 and 3×3 convolution operations very effectively while significantly reducing the number of network parameters. In addition, we used trainable weights in 2D Hadamard transform layers (Section 4.4) and they positively contributed to the accuracy as shown in Table 5 in Section 5.2 and Tables 9m 11 in Section 5.3.

Let $\mathbf{X}, \mathbf{Y} \in \mathbb{R}^m$ be the vectors in the "time" and transform domains, respectively, where $m = 2^k, k \in \mathbb{N}$. The WHT vector \mathbf{Y} is obtained from \mathbf{X} via a matrix multiplication as follows:

$$\mathbf{Y} = \mathbf{W}_k \mathbf{X} \quad (1)$$

where \mathbf{W}_k is called the 2^k -by- 2^k Walsh matrix, which can be generated using the following steps [48]:

- (1) Construct the Hadamard matrix \mathbf{H}_k :

$$\mathbf{H}_k = \begin{cases} 1, & k = 0, \\ \begin{bmatrix} \mathbf{H}_{k-1} & \mathbf{H}_{k-1} \\ \mathbf{H}_{k-1} & -\mathbf{H}_{k-1} \end{bmatrix}, & k > 0, \end{cases} \quad (2)$$

Alternatively, for $k > 1$, \mathbf{H}_k can also be computed using Kronecker product \otimes :

$$\mathbf{H}_k = \mathbf{H}_1 \otimes \mathbf{H}_{k-1}. \quad (3)$$

- (2) Shuffle the rows of \mathbf{H}_k to obtain \mathbf{W}_k by applying the bit-reversal permutation and the Gray-code permutation on row index.

For $k = 2$ and $m = 4$, the Hadamard matrix \mathbf{H}_2 and the corresponding Walsh matrix \mathbf{W}_2 are given by

$$\mathbf{H}_2 = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & -1 & 1 & -1 \\ 1 & 1 & -1 & -1 \\ 1 & -1 & -1 & 1 \end{bmatrix}, \quad \mathbf{W}_2 = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & -1 & -1 \\ 1 & -1 & -1 & 1 \\ 1 & -1 & 1 & -1 \end{bmatrix}. \quad (4)$$

which can be implemented using a wavelet filterbank with binary filters $h[n] = \{1, 1\}$ and $g[n] = \{-1, 1\}$ in two stages [47]. This process is the basis of the fast algorithm and the \mathbf{W}_2 matrix can be expressed as follows

$$\mathbf{W}_2 = \begin{bmatrix} 1 & 1 & 0 & 0 \\ 1 & -1 & 0 & 0 \\ 0 & 0 & 1 & -1 \\ 0 & 0 & 1 & 1 \end{bmatrix} \times \begin{bmatrix} 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 \\ 1 & -1 & 0 & 0 \\ 0 & 0 & 1 & -1 \end{bmatrix} \quad (5)$$

Similar to the Fast Fourier Transform (FFT) algorithm, the complexity of the Fast Walsh-Hadamard transform (FWHT) algorithm is also $O(m \log_2 m)$ when m is an integer power of 2 as it is completely based on the butterfly operations described in Eq. (1) in [49]. Because the Walsh matrix only contains +1 and -1, the implementation of FWHT can be achieved using only addition and subtraction operations via butterflies. It was shown that the WH transform is the same as the block Haar wavelet transform in [47]. As for the inverse transform, we have

$$\mathbf{X} = \frac{1}{m} \mathbf{W}_k \mathbf{Y}, \quad (6)$$

which implies that the inverse Walsh-Hadamard transform is itself with normalization by m .

4 METHODOLOGY

In this section, we will first review our previous work on Fast Walsh-Hadamard Transform (FWHT) layer [20], then we describe the Block Walsh-Hadamard Transform (BWHT) layer. Next, we will introduce a novel 2D-Walsh-Hadamard Transform layer which can be easily inserted into any deep network before the Global Average Pooling (GAP) layer and/or the flatten layer. We will also describe how it can replace the 3×3 convolution layers and the Squeeze-and-Excite layers to reduce the number of parameters significantly. Finally, we describe how we introduce multiplicative weights in the 2D WHT domain. Addition of weights in the 2D WHT transform domain causes a minor increase in the number of parameters but returns a higher accuracy than the regular 2D WHT layer.

4.1 Fast Walsh-Hadamard Transform Layer

The 1×1 convolution layer provides amazing benefit in modern deep convolution layers [5, 6, 21, 50, 51]. It is widely used to change the dimensions of channels. For example, in ResNet [5], a 1×1 convolution layer is applied to the residual blocks when the input and output sizes are different. The 1×1 convolution layer named “conv_expand” is applied to increase the number of channels in each bottleneck layer of MobileNet-V2 [21]. Then, a depthwise convolution layer is employed as the main component for feature extraction. After this step, another 1×1 convolution layer named “conv_projection” reduces the number of channels.

However, there are as many parameters as the number of input channels in each 1×1 convolution layer. Computing these operations is very time-consuming during the inference. Since the main duty of the 1×1 convolution layer is just to adjust the number of channels we proposed a novel FWHT layer to replace the 1×1 convolution layer in [20].

The FWHT layer is summarized in Algorithms 1 and 2. In brief, an FWHT layer consists of an FWHT operation to change the input tensor to the WHT domain, a smooth-thresholding operation as the non-linearity function in the WHT domain, and another inverse FWHT to change the tensor back to the feature-map domain. Each FWHT is applied among the channel axis, which implies that completing FWHT on a tensor $\mathbf{X} \in \mathbb{R}^{n \times w \times h \times m}$ means performing $n \times w \times h$ m -length FWHTs in parallel.

We apply "denoising" in the WH transform domain to eliminate small-amplitude coefficients. This is a nonlinear filtering operation in the transform domain. Inspired by soft-thresholding [52, 53] which is defined as

$$y = S_T(x) = \text{sign}(x)(|x| - T)_+ = \begin{cases} x + T, & x < -T \\ 0, & |x| \leq T \\ x - T, & x > T \end{cases} \quad (7)$$

we have proposed the variant called smooth-thresholding [20]

$$y = S'_T(x) = \tanh(x)(|x| - T)_+, \quad (8)$$

where T is the thresholding parameter and is trainable in our FWHT layer.

Due to its definition given in Eq. (9), the denoising parameter T in soft-thresholding can only be updated by either +1 or -1. On the other hand, the derivative in Eq. (10) is the derivative in Eq. (9) multiplied by $\tanh(x)$. As a result, the convergence of the smooth-thresholding operator is smooth and steady in the back-propagation algorithm.

$$\frac{\partial(\text{sign}(x)(|x| - T)_+)}{\partial T} = \begin{cases} 1, & x < -T \\ 0, & |x| \leq T \\ -1, & x > T \end{cases} \quad (9)$$

$$\frac{\partial(\tanh(x)(|x| - T)_+)}{\partial T} = \begin{cases} -\tanh(x), & |x| > T \\ 0, & |x| \leq T \end{cases} \quad (10)$$

We learn a different threshold T value for each WHT domain coefficient.

We do not use the ReLU function in the transform domain because the WHT domain coefficients can take both positive and negative values, and large positive and negative transform domain coefficients are equally important. In our MobileNet-V2 experiments in [20] we have verified this observation, and both soft-thresholding and smooth-thresholding improve the recognition accuracy compared to the ReLU.

Two types of the FWHT layer are summarized in Algorithms 1 and 2. Because the Direct Current (DC) channel $\mathbf{Y}[:, :, :, 0]$ usually contains essential information about the mean value of the input feature map, we do not apply any thresholding on it. If we perform smooth-thresholding on tensor \mathbf{Y} , each slice among the channel axis will share a common threshold value. Therefore, there are total $(2^d - 1)$ trainable thresholding parameters in the 2^d -channel FWHT layer.

To expand the number of channels, we first compute the 2^d point WHTs and perform smooth-thresholding in the transform domain. We pad $(2^d - c)$ zeros to the end of each input vector before the 2^d -by- 2^d WHT to increase the dimension. After smooth-thresholding in the WH domain, we calculate the inverse WH transform.

To project channels by a factor of $r = \frac{2^p}{2^q}$, we first compute the 2^p point WHTs and perform smooth-thresholding in the transform domain as described in Algorithm 2. After this step, we compute the 2^q point WH transforms to reduce

Algorithm 1 The FWHT layer for channel expansion [20]**Input:** Input tensor $\mathbf{X} \in \mathbb{R}^{n \times w \times h \times c}$ **Output:** Output tensor $\mathbf{Z} \in \mathbb{R}^{n \times w \times h \times tc}$

- 1: Find minimum $d \in \mathbb{N}$, s.t. $2^d \geq tc$
- 2: $\hat{\mathbf{X}} = \text{pad}(\mathbf{X}, 2^d - c) \in \mathbb{R}^{n \times w \times h \times 2^d}$
- 3: $\mathbf{Y} = \text{FWHT}(\hat{\mathbf{X}}) \in \mathbb{R}^{n \times w \times h \times 2^d}$
- 4: $\hat{\mathbf{Y}} = \text{concat}(\mathbf{Y}[:, :, :, 0], \text{ST}(\mathbf{Y}[:, :, :, 1 :]))$
- 5: $\hat{\mathbf{Z}} = \text{FWHT}(\hat{\mathbf{Y}}) \in \mathbb{R}^{n \times w \times h \times 2^d}$
- 6: $\mathbf{Z} = \hat{\mathbf{Z}}[:, :, :, tc]$
- 7: **return** \mathbf{Z} .

Comments: Function $\text{pad}(\mathbf{A}, b)$ pads b zeros on the channel axis of tensor \mathbf{A} . $\text{FWHT}(\cdot)$ is the normalized fast Walsh-Hadamard transform on the last axis. Function $\text{concat}(\cdot, \cdot)$ concatenates two tensors along the last axis. $\text{ST}(\cdot)$ performs smooth-thresholding. Index follows Python's rule.

Algorithm 2 The FWHT layer for channel projection [20]**Input:** Input tensor $\mathbf{X} \in \mathbb{R}^{n \times w \times h \times tc}$ **Output:** Output tensor $\mathbf{Z} \in \mathbb{R}^{n \times w \times h \times c}$

- 1: Find minimum $p, q \in \mathbb{N}$, s.t. $2^p \geq tc$, $2^q \geq c$
- 2: $r = 2^{p-q}$
- 3: $\hat{\mathbf{X}} = \text{pad}(\mathbf{X}, 2^p - tc) \in \mathbb{R}^{n \times w \times h \times 2^p}$
- 4: $\mathbf{Y} = \text{FWHT}(\hat{\mathbf{X}}) \in \mathbb{R}^{n \times w \times h \times 2^p}$
- 5: $\hat{\mathbf{Y}} = \text{concat}(\mathbf{Y}[:, :, :, 0]/r, \text{avgpool}(\text{ST}(\mathbf{Y}[:, :, :, 1 : 2^p - r + 1]), r)) \in \mathbb{R}^{n \times w \times h \times 2^q}$
- 6: $\hat{\mathbf{Z}} = \text{FWHT}(\hat{\mathbf{Y}}) \in \mathbb{R}^{n \times w \times h \times 2^q}$
- 7: $\mathbf{Z} = \hat{\mathbf{Z}}[:, :, :, c]$
- 8: **return** \mathbf{Z} .

Comments: Function $\text{pad}(\mathbf{A}, b)$ pads b zeros on the channel axis of tensor \mathbf{A} . $\text{FWHT}(\cdot)$ is the normalized fast Walsh-Hadamard transform on the last axis. Function $\text{concat}(\cdot, \cdot)$ concatenates two tensors along the last axis. Function $\text{avgpool}(\mathbf{A}, b)$ is the average pooling on \mathbf{A} with pooling size and strides are b . $\text{ST}(\cdot)$ performs smooth-thresholding. Index follows Python's rule.

the dimension of the feature map. We divide the DC channel values by r to keep the energy at the same level as other channels after pooling. In Step 5 of Algorithm 2, we average pool the transform domain coefficients to reduce the dimension of the WHT and discard the last $(r - 1)$ transform domain coefficients of \mathbf{Y} to make the dimension equal to 2^q . The last $(r - 1)$ coefficients are high-frequency coefficients, and usually, their amplitudes are negligible compared to other WHT coefficients.

Therefore, the dimension change operation from m dimensions to n dimensions can be summarized as follows:

$$\mathbf{Z} = \begin{cases} \frac{1}{2^q} \mathbf{U} \mathbf{W}_q \mathbf{S}'_{\mathbf{T}} \mathbf{W}_q \mathbf{P} \mathbf{X}, & m \leq n \\ \frac{1}{\sqrt{2^{p+q}}} \mathbf{U} \mathbf{W}_q \mathbf{A}_{\text{vg}} \mathbf{S}'_{\mathbf{T}} \mathbf{W}_p \mathbf{P} \mathbf{X}, & m > n \end{cases} \quad (11)$$

where p is the minimum integer such that $2^p \geq m$, q is the minimum integer such that $2^q \geq n$, \mathbf{P} describes the zero-padding operation to make \mathbf{X} multipliable by \mathbf{W}_p or \mathbf{W}_q . $\mathbf{S}'_{\mathbf{T}}$ represents the smooth-thresholding layer with DC channel excluded. \mathbf{U} is unpadding function to make the dimension the same as \mathbf{Z} , and \mathbf{A}_{vg} is average pooling on the channel axis without the DC channel. Figure 1 (a) shows an example to increase the number of channels and (b) another example to decrease the number of channels using FWHT layers, respectively.

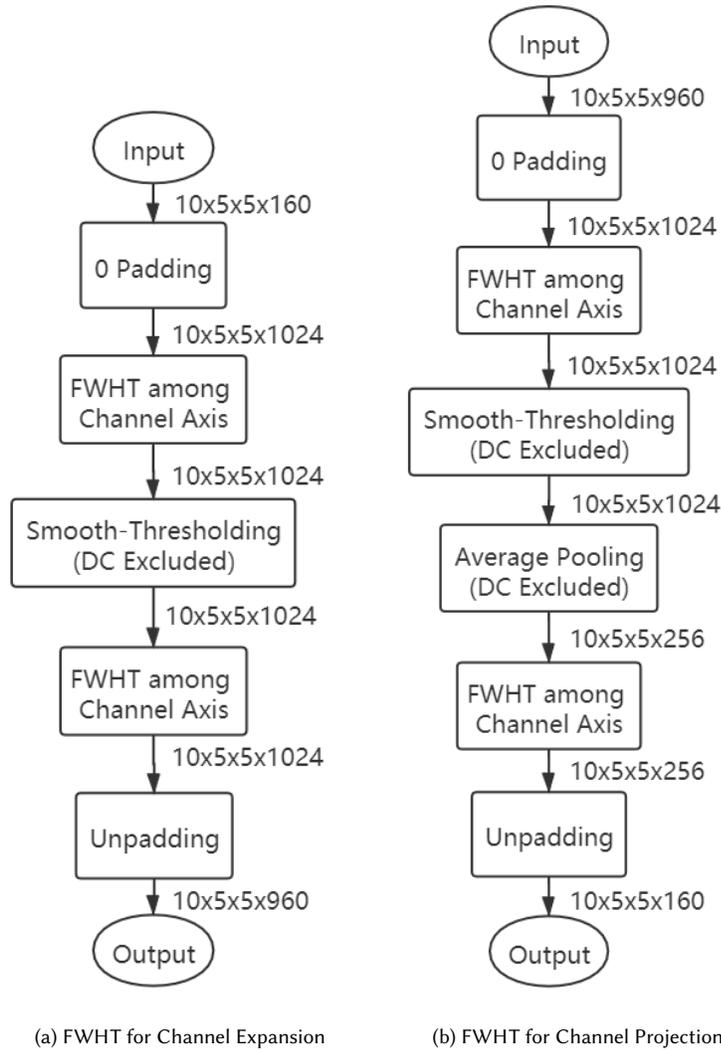


Fig. 1. The FWHT layer (left) increases the number of channels by 6, and the FWHT layer (right) decreases the number of channels of a tensor by $\frac{1}{6}$, respectively. We pad zeros to the tensor to 1024 channels because 1024 is the minimum power of 2 greater than 960. There are 1023 trainable parameters in each FWHT layer.

In consequence, if d is the minimum integer such that 2^d is no less than the number of input channels, the trainable number of parameters in FWHT layers is no more than the $(2^d - 1)$. The trainable parameters are only the threshold values in the smooth-thresholding operator. Hence, it is clear that FWHT layer requires significantly fewer parameters than the regular 1×1 convolution layer, which requires a different set of filter coefficients for each 1×1 convolution.

4.2 Block Walsh-Hadamard Transform layer

Although the FWHT layer is very efficient and can save a huge number of parameters, it requires zero padding before computing the FWHT when the size of the vector is not an integer power of 2. For example, if the input tensor has 384 channels, we have to pad 128 zeros to make it 512. These zeros contain no information but increase the number of trainable parameters significantly. Inspired by the block-division strategy of the Discrete Cosine Transform (DCT)-based JPEG image compression [54] method, we propose a layer in which we divide the input to small blocks and compute the Walsh-Hadamard transforms of small blocks of data. In this way, we do not need to pad a large number of zeros to the end of the input vector. For example, we divide the feature map into blocks of size 32 and compute 12 WHTs in an input tensor which has 384 channels. If necessary, we pad zeros to the end of the last block.

The BWHT layer for channel expansion is described in Algorithm 3. When we want to increase the number of channels by t using the WHT of size s , we overlap data blocks for a c -channel tensor. We change c channels to $\frac{tc}{s}$ blocks and each block has s channels. This is described in Algorithm 4. To achieve overlapping, we first create a $\lfloor \frac{tc}{s} \rfloor$ -length arithmetic sequence $K = [k_0, k_1, \dots, k_{\lfloor \frac{tc}{s} \rfloor - 1}]$ from 0 to $(c - s)$, then we use channels from index $[k_i]$ to index $([k_i] + s - 1)$ to build the i -th block. After overlapping, we take the FWHT of each block along the last axis, perform smooth-thresholding with DC channel excluded, and compute the inverse FWHT. Finally, we reshape the result to a tensor with tc channels.

Algorithm 3 The BWHT Layer for Channel Expansion

Input: Input tensor $\mathbf{X} \in \mathbb{R}^{n \times w \times h \times c}$, Hadamard size s

Output: Output tensor $\mathbf{Z} \in \mathbb{R}^{n \times w \times h \times tc}$

- 1: Resample \mathbf{X} to get $\hat{\mathbf{X}} \in \mathbb{R}^{n \times w \times h \times \frac{tc}{s} \times s}$
- 2: $\mathbf{Y} = \text{FWHT}(\hat{\mathbf{X}}) \in \mathbb{R}^{n \times w \times h \times \frac{tc}{s} \times s}$
- 3: $\hat{\mathbf{Y}} = \text{concat}(\mathbf{Y}[:, :, :, 0], \text{ST}(\mathbf{Y}[:, :, :, 1 :])) \in \mathbb{R}^{n \times w \times h \times \frac{tc}{s} \times s}$
- 4: $\hat{\mathbf{Z}} = \text{FWHT}(\hat{\mathbf{Y}}) \in \mathbb{R}^{n \times w \times h \times \frac{tc}{s} \times s}$
- 5: Reshape $\hat{\mathbf{Z}}$ to get $\mathbf{Z} \in \mathbb{R}^{n \times w \times h \times tc}$
- 6: **return** \mathbf{Z} .

Comments: FWHT(\cdot) is the normalized fast Walsh-Hadamard transform on the last axis. Function $\text{concat}(\cdot, \cdot)$ concatenates two tensors along the last axis. ST(\cdot) performs smooth-thresholding. Index follows Python's rule.

Algorithm 4 Resampling

Input: Input tensor $\mathbf{X} \in \mathbb{R}^{n \times w \times h \times c}$, Hadamard size s

Output: Output tensor $\hat{\mathbf{X}} \in \mathbb{R}^{n \times w \times h \times \frac{tc}{s} \times s}$

- 1: $K = \lfloor \text{linspace}(0, c - s, \lfloor \frac{tc}{s} \rfloor) \rfloor$
- 2: **for** i in $\text{range}(\frac{tc}{s} - 1)$ **do**
- 3: $\hat{\mathbf{X}}[:, :, :, i, :] = \mathbf{X}[:, :, :, K[i] : K[i] + s]$
- 4: **end for**
- 5: **return** $\hat{\mathbf{X}}$.

Comments: $\lfloor \cdot \rfloor$ denotes the floor function. Function $\text{linspace}(a, b, c)$ creates a c -length arithmetic sequence from a to b . Index follows Python's rule.

The BWHT layer for channel projection is described in Algorithm 5. We first divide the tensor into $\frac{tc}{s}$ blocks. Then, in each block, we perform an FWHT, smooth-thresholding, and compute the inverse FWHT. Finally, we reshape the

output tensor to tc channels and take an average pooling to c channels. Unlike the FWHT for channel projection, we take the average pooling in the feature map domain instead of doing it in the WHT domain because the transform size is much smaller in this case, and we want to keep as much information as possible.

Algorithm 5 The BWHT Layer for Channel Projection

Input: Input tensor $\mathbf{X} \in \mathbb{R}^{n \times w \times h \times tc}$, Hadamard size s

Output: Output tensor $\mathbf{Z} \in \mathbb{R}^{n \times w \times h \times c}$

- 1: Reshape \mathbf{X} to get $\hat{\mathbf{X}} \in \mathbb{R}^{n \times w \times h \times \frac{tc}{s} \times s}$
- 2: $\mathbf{Y} = \text{FWHT}(\hat{\mathbf{X}}) \in \mathbb{R}^{n \times w \times h \times \frac{tc}{s} \times s}$
- 3: $\hat{\mathbf{Y}} = \text{concat}(\mathbf{Y}[:, :, :, 0], \text{ST}(\mathbf{Y}[:, :, :, 1 :])) \in \mathbb{R}^{n \times w \times h \times \frac{tc}{s} \times s}$
- 4: $\hat{\mathbf{Z}} = \text{FWHT}(\hat{\mathbf{Y}}) \in \mathbb{R}^{n \times w \times h \times \frac{tc}{s} \times s}$
- 5: Reshape $\hat{\mathbf{Z}}$ to get $\tilde{\mathbf{Z}} \in \mathbb{R}^{n \times w \times h \times tc}$
- 6: $\mathbf{Z} = \text{avgpool}(\tilde{\mathbf{Z}}, t) \in \mathbb{R}^{n \times w \times h \times c}$
- 7: **return** \mathbf{Z} .

Comments: FWHT(\cdot) is the normalized fast Walsh-Hadamard transform on the last axis. Function $\text{concat}(\cdot, \cdot)$ concatenates two tensors along the last axis. Function $\text{avgpool}(\mathbf{A}, b)$ is the average pooling on \mathbf{A} with pooling size and strides are b . ST(\cdot) performs smooth-thresholding. Index follows Python's rule.

Figure 2 shows an example to increase the number of channels and an example to decrease the number of channels using BWHT layers. The dimension change operation from m dimensions to n dimensions can be summarized as follows:

$$\mathbf{Z} = \begin{cases} \frac{1}{2^s} \mathbf{R} \mathbf{W}_q \mathbf{S}'_{\mathbf{T}} \mathbf{W}_q \mathbf{G} \mathbf{X}, & m \leq n \\ \frac{1}{\sqrt{2^s}} \mathbf{A}_{\text{vg}} \mathbf{R} \mathbf{W}_q \mathbf{S}'_{\mathbf{T}} \mathbf{W}_p \mathbf{G} \mathbf{X}, & m > n \end{cases} \quad (12)$$

where \mathbf{G} describes resampling or reshaping operations to divide \mathbf{X} into blocks, $\mathbf{S}'_{\mathbf{T}}$ is the smooth-thresholding layer with DC channel excluded, \mathbf{R} represents the overlapping function to make the dimension the same as \mathbf{Z} , and \mathbf{A}_{vg} represents the average pooling on the channel axis.

Consequently, compared to the FWHT layer with 2^d channels, we can reduce the trainable parameters from $(2^d - 1)$ to s by using the BWHT layer with a block size of s . For example, if there are 384 channels, a single FWHT layer contains 511 parameters while a BWHT layer only contains 31 parameters for a block size of $s = 32$.

We replace the two 1×1 convolution layers in each bottleneck layer of MobileNet-V2 with BWHT layers.

4.3 2D Walsh-Hadamard Transform layer

In this section, we will introduce the two-dimensional (2D) Walsh-Hadamard Transform layers. In convolution neural networks, the global average pooling (GAP) and flatten layers are universally employed before the dense layers to reduce the dimension of the output tensor. James Lee-Thorp *et al.* proposed a Fourier Transform-based layer before the dense layer in their natural language processing network called FNet [55]. Inspired by their work, we design a 2D WHT layer that can be inserted before the GAP layer or the flatten layer as shown in Figure 3b and Equation 13.

$$\mathbf{y} = \mathbf{W} \mathbf{S}'_{\mathbf{T}} \mathbf{W} \mathbf{x}. \quad (13)$$

where, \mathbf{W} denotes 2D-FWHT operation and $\mathbf{S}'_{\mathbf{T}}$ represents the 2D smooth thresholding.

In most convolution neural networks, the input tensor of the GAP and the flatten layer is in $\mathbb{R}^{N \times w \times h \times c}$, where N is the batch size. Usually, $w = h \ll c$. On the other hand, the output tensor of the GAP layer is in $\mathbb{R}^{N \times c}$ and the output

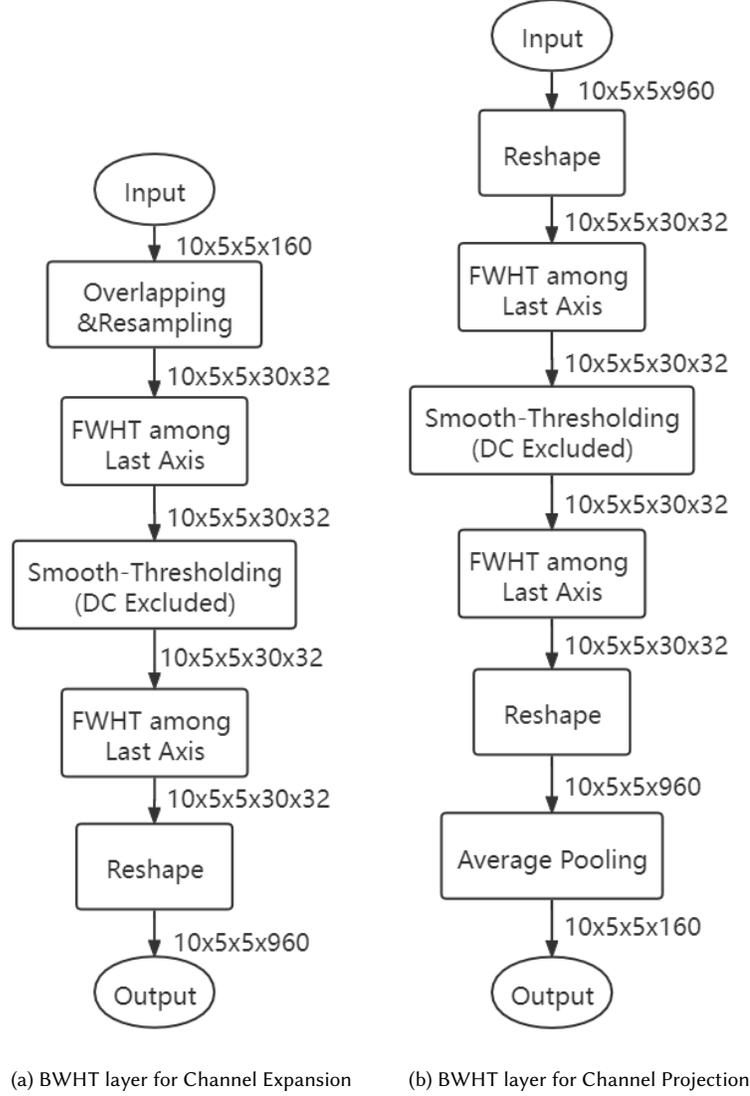


Fig. 2. A BWHT layer (left) increases channel number of a tensor by 6 and a BWHT layer (right) decreases channel number of a tensor by $\frac{1}{6}$, respectively. There are only 31 trainable parameters in each BWHT layer because the WHT size is 32.

tensor of the flatten layer is in $\mathbb{R}^{N \times whc}$. Therefore, the input of the GAP layer contains more information than the output of the GAP layer, and flatten layer reduces the 3D-spatial structure of the tensor. Due to these reasons, we insert the 2D-FWHT layer before the GAP layer or the flatten layer, and the 2D-FWHT layer aims to reinforce analysis among the width and height axes.

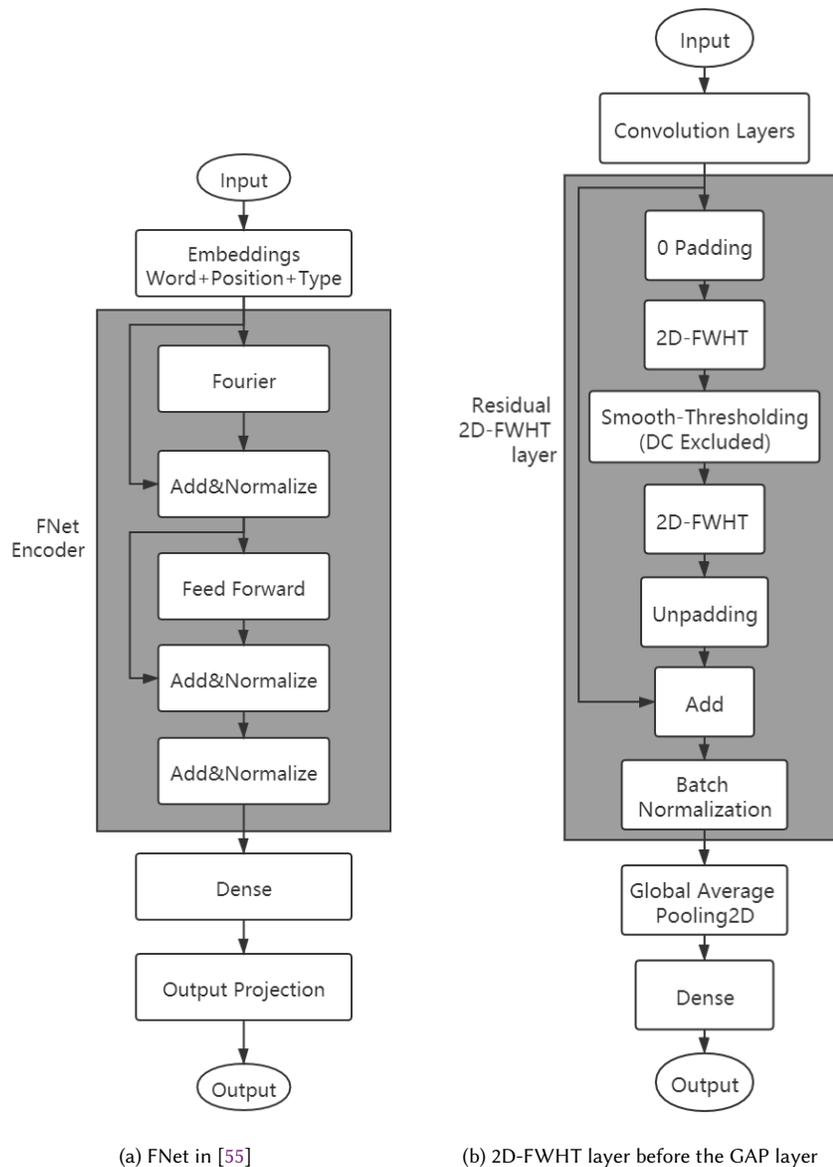


Fig. 3. FNet and CNN with a 2D-FWHT layer before the GAP layer.

As it is described in Algorithm 6 and Figure 4, we first pad zeros among the width and height axes to make the sizes of the transform powers of 2. Usually, w and h are very small integers if the convolution part is well-designed. Hence, the WHT size is not very large in this case (2, 4, or 8 is usually sufficient). Therefore, we do not apply block FWHT at this stage because we do not need to pad too many zeros. Then, we apply two 1D-FWHTs among the width and

height axes separately to convert the tensor into the WH domain. Next, we perform 2D-smooth-thresholding. There are 2^{p+q} trainable thresholding parameters, while the one at the DC channel is redundant but needed to perform Python's broadcasting. Each slice among the width and height axes will share a common threshold value. Afterwards, we reset the DC value to its original value and apply two 1D-FWHTs along the width and height axes separately to convert the tensor back to the feature domain. Finally, we un-pad the zeros from the result with the input tensor to get the output tensor. If it is with residual design, we will add the output with the input tensor to get the final output tensor.

Algorithm 6 The 2D-FWHT Layer

Input: Input tensor $\mathbf{X} \in \mathbb{R}^{n \times w \times h \times c}$, Hadamard size s

Output: Output tensor $\mathbf{Z} \in \mathbb{R}^{n \times w \times h \times c}$

- 1: find minimum $p, q \in \mathbb{N}$, s.t. $2^p \geq w, 2^q \geq h$
- 2: $\hat{\mathbf{X}} = \text{pad2}(\mathbf{X}, 2^p - w, 2^q - h) \in \mathbb{R}^{n \times 2^p \times 2^q \times c}$
- 3: $\mathbf{Y} = \text{FWHT2}(\hat{\mathbf{X}}) \in \mathbb{R}^{n \times 2^p \times 2^q \times c}$
- 4: $\hat{\mathbf{Y}} = \text{ST2}(\mathbf{Y}) \in \mathbb{R}^{n \times 2^p \times 2^q \times c}$
- 5: $\hat{\mathbf{Y}}[:, 0, 0, :] = \mathbf{Y}[:, 0, 0, :]$
- 6: $\hat{\mathbf{Z}} = \text{FWHT2}(\hat{\mathbf{Y}}) \in \mathbb{R}^{n \times 2^p \times 2^q \times c}$
- 7: **if** Residual is applied **then**
- 8: $\mathbf{Z} = \hat{\mathbf{Z}}[:, : w, : h, :] + \mathbf{X}$
- 9: **else**
- 10: $\mathbf{Z} = \hat{\mathbf{Z}}[:, : w, : h, :]$
- 11: **end if**

Comments: Function $\text{pad2}(A, b, c)$ pads b zeros among the width axis and c zeros among the height axis of tensor A . $\text{FWHT2}(\cdot)$ is the normalized 2D-fast Walsh-Hadamard transform on the width and height axes (applying two 1D-FWHTs separately for efficiency). $\text{ST2}(\cdot)$ performs 2D-smooth-thresholding. Index follows Python's rule.

We also use the 2D-FWHT layers to replace the regular 2D-convolution layers. In this case, it does not change the number of channels due to the fact that it takes no transformation among the channel axis. Therefore, if we want to change the number of channels, we can apply a 1D-BWHT first then apply a 1D-BWHT layer to change the number of channels, then apply a 2D-FWHT layer. If the weight and height of the input tensor are not very small, we can also apply the block division as the BWHT layer to avoid padding too many zeros.

4.4 2D Weighted Walsh-Hadamard Transform layer

Convolution in "time" domain can be implemented using multiplicative weights in the Fourier domain. However, Fourier transform requires complex arithmetic. In this paper, we replaced the Fourier Transform with the binary WHT. The weighted 2D WHT layer uses multiplicative weights in the transform domain. In the weighted-WHT layer we perform both regular "filtering" and nonlinear "denoising" in the transform domain.

To improve the accuracy, we use the weighted smooth-thresholding function, in which a trainable weight is applied:

$$y = S'_{WT}(x) = \tanh(x)(|vx| - T)_+ \quad (14)$$

with constraint $v \geq 0$, where v is a trainable multiplicative weight representing a WHT domain multiplication. The corresponding partial derivatives are

$$\frac{\partial(\tanh(x)(|vx| - T)_+)}{\partial T} = \begin{cases} \tanh(x)|x|, & |vx| > T \\ 0, & |vx| \leq T \end{cases} \quad (15)$$

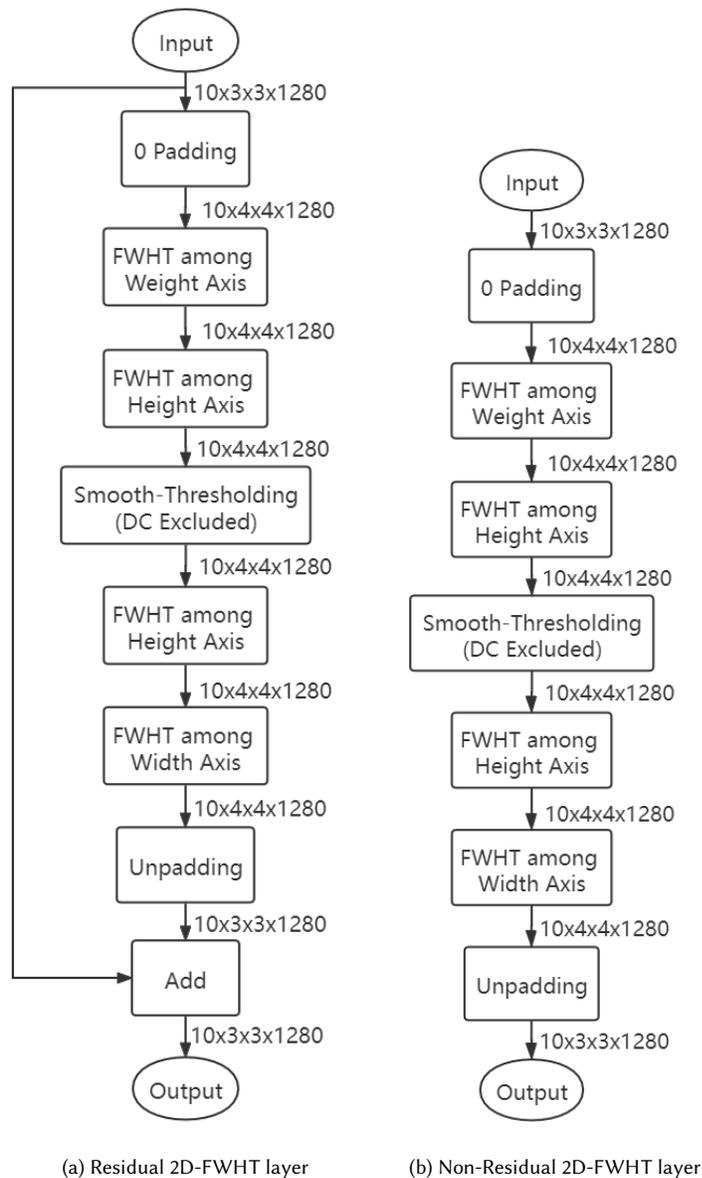


Fig. 4. 2D-FWHT layer examples. In each example, there are $4 \times 4 = 16$ trainable parameters while the one at DC is redundant. If weights are applied in the smooth-thresholding (details are in Section 4.4, there are $4 \times 4 = 16$ extra trainable parameters. The residual one can be inserted before the GAP layer or the flatten layer to assist the dense layers, while the non-residual one can be employed to replace the 2D convolution layers and Squeeze-and-Excite layers.

and

$$\frac{\partial(\tanh(x)(|vx| - T)_+)}{\partial T} = \begin{cases} -\tanh(x), & |vx| > T \\ 0, & |vx| \leq T \end{cases} \quad (16)$$

When $|vx| > T$, Eq. (15) holds because when $v \geq 0$, $\tanh(x)x\text{sign}(vx) = \tanh(x)x\text{sign}(x) = \tanh(x)|x|$. We don't multiply v and x in $\tanh(\cdot)$ in Eq. (14) to avoid quadratic term of v in the derivatives. In this way, Eq. (13) can be rewritten as:

$$y = \mathbf{W}S'_{\mathbf{W}\mathbf{T}}\mathbf{W}\mathbf{x}. \quad (17)$$

where, \mathbf{W} denotes 2D-FWHT and $S'_{\mathbf{W}\mathbf{T}}$ denotes 2D weighted smooth thresholding. We initialize all weights $v = 1$ in the 2D weighted smooth-thresholding layer. We initialize the threshold value $T = 0$ at the DC channel. The threshold values at other AC channels are initialized as positive values. In weighted-WHT layer, the number of trainable parameters are doubled but the number of parameters in the 2D-FWHT layer is still significantly smaller than a regular convolution layer. For example, if the input and the output are in $\mathbb{R}^{3 \times 3 \times 1280}$ and the WHT size is 4, there are only $4 \times 4 + 4 \times 4 = 32$ trainable parameters in weighted 2D-FWHT layer, but 3×3 2D convolution layer requires $3 \times 3 \times 1280 = 11520$ trainable parameters.

5 EXPERIMENTAL RESULTS

Our training and testing experiments are carried on an HP-Z820 workstation with 2 Intel Xeon E5-2695 v2 CPUs, 2 NVIDIA RTX A4000 GPUs, and 128GB RAM. Speed tests are also carried out using a NVIDIA Jetson Nano board. The code is written in TensorFlow-Keras in Python 3.

First, we will compare the accuracy of MobileNet-V2 based neural networks on the Fashion MNIST dataset and CIFAR-10 dataset. Then, we will further investigate 2D-FWHT in MobileNet-V3 and ResNet. Finally, we will perform a speed test on our 2D-FWHT layer and the regular 3×3 convolution layer.

5.1 BWHT in MobileNet-V2

In MobileNet-V2 experiments, we use ImageNet-pretrained MobileNet-V2 model with 1.0 depth [56]. As it is shown in Table 1, to build the fine-tuned MobileNet-V2, we replace layers after the Global Average Pooling (GAP) layer with a dropout layer and a dense layer as TensorFlow official transfer learning demo [57]. Since the minimum input of MobileNet-V2 ImageNet-pretrained models is $96 \times 96 \times 3$, we resize the images to this resolution in our experiments. Because images in Fashion MNIST dataset are in gray-scale, after interpolating, we copy the values two times to convert them to $96 \times 96 \times 3$ format.

Results of MobileNet-V2 on Fashion MNIST and CIFAR-10 are shown in Tables 2 and 3 respectively. The number of parameters is counted by TensorFlow API "model.summary()", and the non-trainable parameters are from the batch normalization layers. We first investigate the 2D-FWHT layer before the GAP layer. We name the models as "2D-FWHT before GAP". We add one residual 2D-FWHT layer followed by a batch normalization layer. They bring an extra 2576 more trainable parameters (only 16 from the 2D-FWHT layer) and only 2560 non-trainable parameters but improve the accuracy by 0.27% and 0.47% on Fashion MNIST and CIFAR-10, respectively. The number of trainable parameters increases only by 0.12%.

We then beat the FWHT results in [20] which do not apply block division. As Figure 5 shows, we first revise the 1×1 convolution layers. The models are named as "1x1 conv. changed". The fraction at the beginning of the model's

Table 1. Structure of fine-tuned MobileNet-V2 (baseline). t , c , n and s represent expansion factor, channel, repeat time and stride [21]. Initial weights before dropout are from ImageNet checkpoint float_v2_1.0_96 in [56].

Input	Operator	t	c	n	s
$96^2 \times 3$	Conv2D	-	32	1	2
$48^2 \times 32$	Bottleneck	1	16	1	1
$48^2 \times 16$	Bottleneck	6	24	2	2
$24^2 \times 24$	Bottleneck	6	32	3	2
$12^2 \times 32$	Bottleneck	6	64	4	2
$6^2 \times 64$	Bottleneck	6	96	3	1
$6^2 \times 96$	Bottleneck	6	160	3	2
$3^2 \times 160$	Bottleneck	6	320	1	1
$3^2 \times 320$	Conv2D	-	1280	1	1
$3^2 \times 1280$	AvgPool	-	-	1	-
1280	Dropout (rate=0.2)	-	-	1	-
1280	Dense (units=10)	-	-	1	-

Table 2. MobileNet-V2 Fashion MNIST Result

Model	Trainable Parameters	Non-Trainable Parameters	Trainable Parameters Reduction Ratio	Accuracy
Fine-tuned MobileNet-V2 (baseline)	2,236,682	34,112	-	95.37%
2D-FWHT before GAP	2,239,258	36,672	-	95.64%
$\frac{1}{2}$ 1x1 conv. changed (FWHT) [20]	288,983	34,112	87.08%	94.50%
$\frac{1}{2}$ 1x1 conv. changed (BWHT)	273,177	34,112	87.79%	94.45%
$\frac{1}{2}$ 1x1 conv. changed (BWHT) +2D-FWHT before GAP	275,753	36,372	87.67%	94.75%

Table 3. MobileNet-V2 CIFAR-10 Result

Model	Trainable Parameters	Non-Trainable Parameters	Trainable Parameters Reduction Ratio	Accuracy
Baseline MobileNet-V2 model ^a in [44]	2.2378M	-	-	94.3%
RMNv2 ^b [44]	1.0691M	-	52.22%	92.4%
Fine-tuned MobileNet-V2 (baseline)	2,236,682	34,112	-	95.21%
Fourier Layer [55] before GAP	2,239,242	36,672	-	95.38%
2D-FWHT before GAP	2,239,258	36,672	-	95.68%
$\frac{1}{3}$ 1x1 conv. changed (FWHT) [20]	506,094	34,112	77.37%	93.14%
$\frac{1}{3}$ 1x1 conv. changed (BWHT)	494,175	34,112	77.91%	93.22%
$\frac{1}{3}$ 1x1 conv. changed (BWHT) +2D-FWHT before GAP	496,751	36,672	77.79%	93.46%

Parameters and accuracy of baseline model^a and RMNv2^b are from Table 3 in [44].

name denotes how many bottleneck layers are changed. For instance, $\frac{1}{2}$ means we change the last half bottleneck layers. Since the final convolution layer is also 1×1 , we replace it with a BWHT layer for channel expansion.

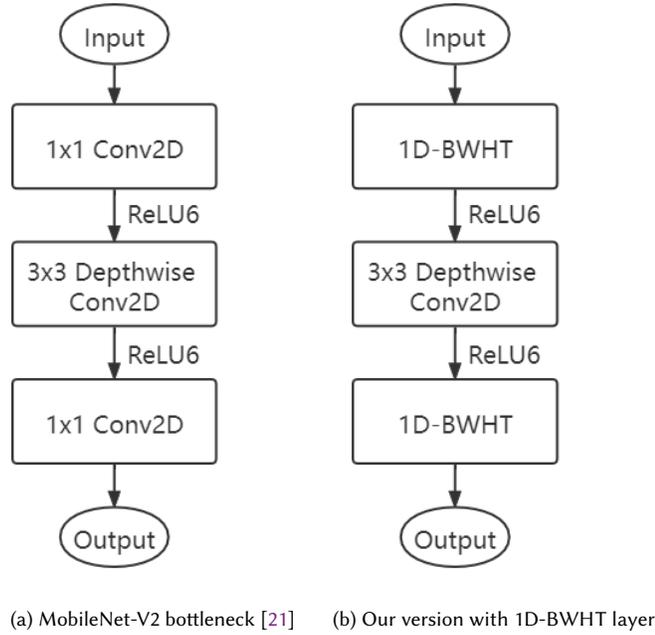


Fig. 5. MobileNet-V2 bottleneck and our version with 1D-BWHT layer.

We also compared our results with Reduced Mobilenet-V2 (RMNV2) [44] in CIFAR-10 experiment. They slim MobileNet-V2 by replacing bottleneck layers with a novel block called HetConv blocks. Their trimming method is still based on the convolution. When we change $\frac{1}{3}$ convolution by BWHT layer and add an extra 2D-FWHT layer before the GAP layer, we reduce more parameters than (77.79% VS 52.22%) [44] with a smaller accuracy loss (1.75% VS 1.9%).

Moreover, we compare our results with the FNet Fourier layer [55] in CIFAR-10 experiment. The FNet Fourier layer can be computed as

$$Y = X + \Re(\mathcal{F}_w(\mathcal{F}_h(X))) \quad (18)$$

where the input tensor is X and the output tensor is Y . $\mathcal{F}_w(\cdot)$ and $\mathcal{F}_h(\cdot)$ are the Fast Fourier Transform (FFT) among the weight and the height. $\Re(\cdot)$ is the real part of the tensor. We try adding one FNet layer before the GAP layer to compare the result of adding one 2D-FWHT layer before the GAP layer, and our accuracy is 0.30% higher than the accuracy of the FNet. This is because the FNet can only keep the real part of the tensor. In another words the information in the imaginary part is simply discarded. On the other hand, the Walsh-Hadamard Transform is a binary transform. Therefore, there is no information loss due to the Walsh-Hadamard Transform.

According to the above experiments, we have the following observations:

- In general, the weighted 2D-FWHT before the GAP layer improves the accuracy of an image recognition network. For example, it improves the accuracy of the fine-tuned MobileNet-V2 by 0.27% in Fashion MNIST and 0.47% in CIFAR-10. Furthermore, it increases the accuracy of the model “ $\frac{1}{2}$ 1x1 conv. changed BWHT” by 0.4% on Fashion MNIST and accuracy of model “ $\frac{1}{3}$ 1x1 conv. changed BWHT” on CIFAR-10 by 0.24%, respectively.

- The 1D-BWHT layer provides better results than the single block 1D-FWHT layer in terms of both the accuracy and the number of parameters when the number of channels is not the power of 2. Although on Fashion MNIST, the 1D-FWHT model reaches higher accuracy, the BWHT models reach higher accuracy on CIFAR-10. CIFAR-10 is a more difficult dataset so it is more convincing than the Fashion MNIST. Plus, in each case, the BWHT models have fewer parameters than the 1D FWHT models because the BWHT layer avoids padding 0s and it uses a smaller WHT size.

5.2 2D-FWHT in MobileNet-V3

In this section, we use MobileNet-V3-Large [58] on CIFAR-100 dataset and Tiny ImageNet to show the advantage of the 2D-FWHT layer. As Figure 6 shows, compared to the MobileNet-V2 bottleneck block, the MobileNet-V3 bottleneck block contains one Squeeze-and-Excite layer after the depthwise convolution layer. In this section, we replace each Squeeze-and-Excite layer in the final $\frac{1}{3}$ bottleneck blocks with the 2D-FWHT layer. As Table 4 shows, the base model is MobileNet-V3-large with the weights in the early layers are pretrained on the ImageNet dataset. Due to the input size of public ImageNet-pretrained MobileNet-V3 models are all 224×224 but CIFAR image size is 32×32 , we upsample the images to 224×224 .

Table 4. Structure of the fine-tuned MobileNet-V3-Large. The table follows Table 1 in [58]. Initial weights before dropout are from ImageNet checkpoint “Large dm=1 (float)” in [56]. SE denotes whether there is a Squeeze-And-Excite in that block. “NL” denotes the type of nonlinearity used. Here, “HS” denotes h-swish and “RE” denotes ReLU. “s” denotes stride.

Input	Operator	exp size	#out	SE	NL	s
$224^2 \times 3$	Conv2d	-	16	-	HS	2
$112^2 \times 16$	bneck, 3x3	16	16	N	RE	1
$112^2 \times 16$	bneck, 3x3	64	24	N	RE	2
$56^2 \times 24$	bneck, 3x3	72	24	N	RE	1
$56^2 \times 24$	bneck, 5x5	72	40	Y	RE	2
$28^2 \times 40$	bneck, 5x5	120	40	Y	RE	1
$28^2 \times 40$	bneck, 5x5	120	40	Y	RE	1
$28^2 \times 40$	bneck, 3x3	240	80	N	HS	2
$14^2 \times 80$	bneck, 3x3	200	80	N	HS	1
$14^2 \times 80$	bneck, 3x3	184	80	N	HS	1
$14^2 \times 80$	bneck, 3x3	184	80	N	HS	1
$14^2 \times 80$	bneck, 3x3	480	112	Y	HS	1
$14^2 \times 112$	bneck, 3x3	672	112	Y	HS	1
$14^2 \times 112$	bneck, 5x5	672	160	Y	HS	2
$7^2 \times 160$	bneck, 5x5	960	160	Y	HS	1
$7^2 \times 160$	bneck, 5x5	960	160	Y	HS	1
$7^2 \times 160$	conv2d, 1x1	-	960	-	HS	1
$7^2 \times 960$	AvgPool	-	-	-	-	1
960	Dropout(rate=0.2)	-	-	-	-	-
960	Dense(units=10 for CIFAR-10, 100 for CIFAR-100, 200 for Tiny ImageNet)	-	-	-	-	-

As Table 5 shows, we reduce 48.55% trainable parameters with only a 0.33% accuracy loss by changing one-third of ($\frac{1}{3}$) Squeeze-and-Excite layers and applying one 2D-FWHT layer before the GAP layer. Compared to the cases without applying weights in the 2D smooth-thresholding, all cases with weights get improved accuracy. Specifically, when we

insert a 2D-FWHT layer in the 2D-smooth-threshold before the GAP layer, the network with weights in the 2D smooth-thresholding reaches 0.51% accuracy higher than the network without weights in the 2D smooth-thresholding. When we change $\frac{1}{3}$ Squeeze-and-Excite layers, the network with weights in the 2D smooth-thresholding reaches 0.15% accuracy higher than the network without weights in the 2D smooth-thresholding. In each pair, the parameters amount is almost the same respectively. Therefore, weighted smooth-thresholding is superior to non-weighted smooth-thresholding in accuracy.

We also compare the WHT layer with FNet Fourier layer [55]. When we insert a layer before the GAP layer, the model with our 2D-FWHT reaches a 0.62% accuracy higher than the model with FNet Fourier layer. In addition, when we change $\frac{1}{3}$ Squeeze-and-Excite layers, the model with 2D-FWHT reaches a 0.29% accuracy higher than the model with FNet Fourier layer.

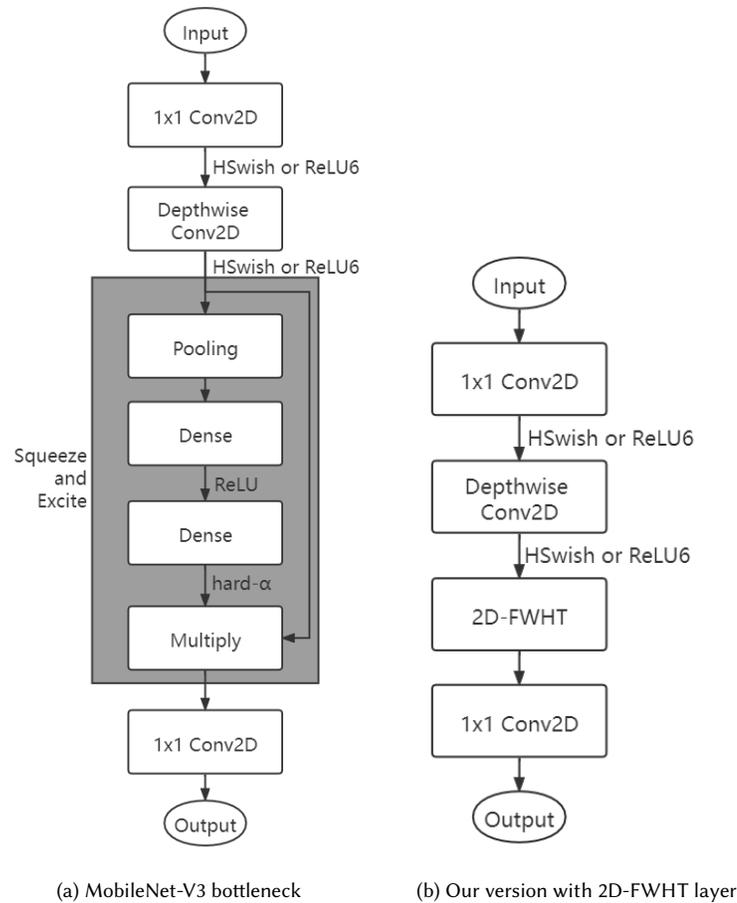


Fig. 6. MobileNet-V3 bottleneck [58] and our version with 2D-FWHT layer.

Furthermore, we also try our best model on CIFAR-10 and Tiny ImageNet. We replace each Squeeze-and-Excite layer in the final $\frac{1}{3}$ bottleneck blocks by the 2D-FWHT layer and insert one 2D-FWHT layer before the GAP layer. As

Table 5. MobileNet-V3-Large CIFAR-100 Result

Model	Weights in 2D Thresholding	Trainable Parameters	Non-Trainable Parameters	Trainable Parameters Reduction Ratio	Accuracy
Fine-tuned model (Baseline)	-	3,068,052	24,400	-	78.25%
Fourier layer [55] before GAP	-	3,069,972	26,320	-	80.01%
2D-FWHT before GAP	N	3,070,036	26,320	-	80.22%
2D-FWHT before GAP	Y	3,070,100	26,320	-	80.73%
$\frac{1}{3}$ S&E changed (Fourier [55])	-	1,574,988	24,400	-	77.20%
$\frac{1}{3}$ S&E changed (2D-FWHT)	N	1,575,692	24,400	48.64%	77.34%
$\frac{1}{3}$ S&E changed (2D-FWHT)	Y	1,576,396	24,400	48.62%	77.49%
$\frac{1}{3}$ S&E changed+2D-FWHT before GAP	Y	1,578,444	26,320	48.55%	77.92%

S&E is the Squeeze-and-Excite layer.

Table 6 shows, our model has 49.96% less trainable parameters than the fine-tuned MobileNet-V3-Large model and the accuracy is only 1.20% lower on the CIFAR-10. As Table 7 shows, our model has 47.08% less trainable parameters than the fine-tuned MobileNet-V3-Large model and the accuracy is only 1.59% lower on the Tiny ImageNet.

Table 6. MobileNet-V3-Large CIFAR-10 Result

Model	Weights in 2D Thresholding	Trainable Parameters	Non-Trainable Parameters	Trainable Parameters Reduction Ratio	Accuracy
Fine-tuned model (Baseline)	-	2,981,562	24,400	-	95.52%
$\frac{1}{3}$ S&E changed+2D-FWHT before GAP	Y	1,491,954	26,320	49.96%	94.32%

S&E is the Squeeze-and-Excite layer.

Table 7. MobileNet-V3-Large Tiny ImageNet Result

Model	Weights in 2D Thresholding	Trainable Parameters	Non-Trainable Parameters	Trainable Parameters Reduction Ratio	Accuracy
Fine-tuned model (Baseline)	-	3,164,152	24,400	-	64.37%
$\frac{1}{3}$ S&E changed+2D-FWHT before GAP	Y	1,674,544	26,320	47.08%	62.78%

S&E is the Squeeze-and-Excite layer.

5.3 2D-FWHT in ResNet

In this section, we will investigate 2D-FWHT in ResNet. ResNet is built with standard convolution layers [5]. It does not use any depthwise convolution layer. Therefore, we can employ the ResNet to show the advantage of our 2D-FWHT layer. We do not use any pre-trained weights but we initialize the weights with Kaiming He’s initialization [59]. We first build ResNet-20 as shown in Table 8 and Figure 7a. As Table 9 shows, it reaches an accuracy of 91.60% on CIFAR-10 with 273,066 trainable parameters. Then, we revise all residual blocks by replacing 3×3 convolution layers with 2D-FWHT

layers and replacing 1×1 convolution layers with 1D-BWHT layers as shown in Figure 7b. Here we apply 2D-FWHT layers without residual design because the blocks already contain the residual design. Since the dimension numbers are already in the power of 2, we do not need to pad any zeros before computing the Walsh-Hadamard transforms. We do not implement an additional 2D-FWHT layer before the GAP layer because the input of the GAP layer is already from a 2D-FWHT layer. In this way, we save 95.76% parameters. In other words, there are virtually no parameters left before the GAP layers. Note that the dense layer contains 650 parameters, the batch normalization layers have 2,752 parameters, and the first convolution layer holds 448 parameters. There are only 7,740 parameters in the 2D-FWHT layers and the 1D-BWHT layers. With an extremely low number of trainable parameters, the network still reaches 60.47% accuracy. We then revise the residual partially as shown in Figure 7c. We retain the first 3×3 convolution layer in each residual block and change other convolution layers. In this way, we save 52.76% trainable parameters with only a 1.72% accuracy loss. Finally, we add weights in the 2D smooth-threshold, and we reduce 51.26% trainable parameters with only a 1.48% accuracy loss.

Table 8. Structure of ResNet-20 for CIFAR-10. Building blocks are shown in brackets, with the numbers of blocks stacked. Down-sampling is performed by conv3_1 and conv4_1 with a stride of 2. Batch normalization layer is applied after each convolution layer.

Layer Name	Output Shape	Implementation Details
Input	$32 \times 32 \times 3$	-
Conv1	$32 \times 32 \times 16$	$3 \times 3, 16$
Conv2_x	$32 \times 32 \times 16$	$\left[\begin{array}{l} 3 \times 3, 16 \\ 3 \times 3, 16 \end{array} \right] \times 3$
Conv3_x	$16 \times 16 \times 32$	$\left[\begin{array}{l} 3 \times 3, 32 \\ 3 \times 3, 32 \end{array} \right] \times 3$
Conv4_x	$8 \times 8 \times 64$	$\left[\begin{array}{l} 3 \times 3, 32 \\ 3 \times 3, 64 \end{array} \right] \times 3$
GAP	64	Global Average Pooling2D
Output	10	Dense(unit = 10)

Table 9. ResNet-20 CIFAR-10 Result

Model	Weights in 2D Thresholding	Trainable Parameters	Non-Trainable Parameters	Trainable Param Reduction Ratio	Accuracy
ResNet-20 in [5]	-	0.27M	-	-	91.25%
ResNet-20 in our trial (baseline)	-	273,066	1,376	-	91.60%
2D-FWHT before GAP	N	273,258	1,504	-	91.74%
2D-FWHT before GAP	Y	273,322	1,504	-	91.75%
Residual blocks completely revised	N	11,590	1,376	95.76%	60.47%
Residual blocks partially revised	N	129,000	1,376	52.76%	89.88%
Residual blocks partially revised	Y	133,082	1,376	51.26%	90.12%

We further explore ResNet-34 on the Tiny ImageNet. The base model is defined as Table 10. As it is shown in Table 11, we reduce 53.66% trainable parameters with only 0.89% accuracy loss without weights in 2D smooth-thresholding and 53.37% trainable parameters with only 0.72% accuracy loss with weights in 2D smooth-thresholding.

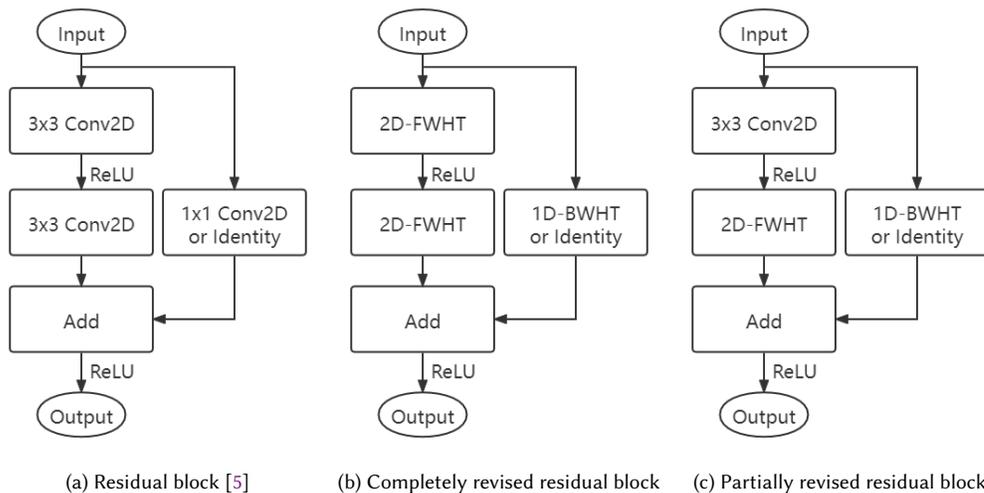


Fig. 7. Residual Block and our revised versions. 1×1 convolution or 1D-BWHT layer is applied if the number of channels increases.

Table 10. Structure of ResNet-34 for Tiny ImageNet. Building blocks are shown in brackets, with the numbers of blocks stacked. Downsampling is performed by conv3_1, conv4_1 and conv5_1 with a stride of 2. Batch normalization layer is applied after each convolution layer.

Layer Name	Output Shape	Implementation Details
Input	$64 \times 64 \times 3$	-
Conv1	$64 \times 64 \times 64$	$3 \times 3, 64$
Conv2_x	$64 \times 64 \times 64$	$\left[\begin{array}{l} 3 \times 3, 64 \\ 3 \times 3, 64 \end{array} \right] \times 3$
Conv3_x	$32 \times 32 \times 128$	$\left[\begin{array}{l} 3 \times 3, 128 \\ 3 \times 3, 128 \end{array} \right] \times 4$
Conv4_x	$16 \times 16 \times 256$	$\left[\begin{array}{l} 3 \times 3, 256 \\ 3 \times 3, 256 \end{array} \right] \times 6$
Conv5_x	$8 \times 8 \times 512$	$\left[\begin{array}{l} 3 \times 3, 512 \\ 3 \times 3, 512 \end{array} \right] \times 3$
GAP	512	Global Average Pooling2D
Output	200	Dense(unit = 200)

Table 11. ResNet-34 Tiny ImageNet Result

Model	Weights in 2D Thresholding	Trainable Parameters	Non-Trainable Parameters	Trainable Param Reduction Ratio	Accuracy
ResNet-34 (baseline)	-	21,386,312	15,232	-	53.06%
Residual blocks partially revised	N	9,910,565	15,232	53.66%	52.17%
Residual blocks partially revised	Y	9,928,677	15,232	53.57%	52.34%

5.4 Speed and Memory Tests

In our early work [20], we have shown that the 1D-FWHT layer runs about 2 times as fast as the 1×1 convolution layer on the NVIDIA Jetson Nano (4GB version) when the input and the output are both in $\mathbb{R}^{10 \times 32 \times 32 \times 1024}$ and the input tensor is initialized randomly. Now we will compare the 2D-FWHT layer versus the 3×3 convolution layer to show the speed advantage of our 2D-FWHT layer for the ResNet and other deep CNNs with regular convolution layers.

In this experiment, we set the input and the output are both in $\mathbb{R}^{10 \times 8 \times 8 \times 1024}$. We use the same laptop (with Intel Core i7-7700HQ CPU, NVIDIA GTX-1060 GPU with Max-Q design, and 16GB DDR4 RAM) as [20] to run the TensorFlow PB model, and we deploy the TFLite model to the same NVIDIA Jetson Nano. As it is stated in Table 12, the 2D-FWHT layer runs about 24 times as fast as the 3×3 convolution layer on the NVIDIA Jetson Nano. There are two reasons why the time difference on NVIDIA Jetson Nano is much more significant than on the laptop. One reason is that the TFLite model on the NVIDIA Jetson Nano is optimized by TensorFlow for ARM embedded devices, while the PB model on the laptop is not optimized as much as the TFLite model. The other reason is that the Windows backend apps slow down the system, while the NVIDIA Jetson Nano system is “clean”.

We also compare the 2D-FWHT layer with the Squeeze-and-Excite layer. The squeeze ratio is $\frac{1}{4}$ as the Squeeze-and-Excite layer in the MobileNet-V3. We achieve comparable times in squeeze-and-excite layers. This is because current processors cannot distinguish between multiplication by ± 1 and a real number.

Table 12. Speed test

Device	3×3 Conv2D	2D-FWHT	Squeeze-and-Excite (Ratio = $\frac{1}{4}$)
Laptop (GPU)	0.0508 S	0.0459 S	0.0449 S
Laptop (CPU)	0.0768 S	0.0479 S	0.0424 S
NVIDIA Jetson Nano	1.8861 S	0.0776 S	0.0519 S

Input tensor and output tensor are in $\mathbb{R}^{10 \times 8 \times 8 \times 1024}$. code is available at [60].

In addition, We record the RAM usage of the three layers on the NVIDIA Jetson Nano. According to Figure 8, the 2D-FWHT layer and the Squeeze-and-Excite layer have close RAM usages, and they require about 40MB (19.51%) less memory than the 3×3 Conv2D in the inference.

6 CONCLUSIONS

In this paper, we proposed 1D and 2D Walsh-Hadamard Transform (WHT)-based binary layers to replace 1×1 and 3×3 convolution layers in deep neural networks. The 2D-FWHT layer can also be used to replace the Squeeze-and-Excite layers common to many publicly available networks including the MobileNet-V3 [58].

We implement WHT in blocks of data using its $O(m \log_2 m)$ Fast-WHT (FWHT) algorithm, where m is the number of elements in each block. Compared to the FWHT layer in [20] which does not apply block by block computation, the Block-WHT (BWHT) layer avoids zero paddings, if the transform size m is chosen wisely as a power of 2, and our new approach can save more trainable parameters compared to our earlier paper [20].

In addition, we propose a residual 2D-FWHT layer that can be easily inserted before the global average pooling (GAP) layer (or the flatten layer) to assist the dense layers. By applying the 1D-BWHT layers and a 2D-FWHT layer before the GAP layer on the MobileNet-V2, we save 87.67% trainable parameters with only a 0.62% accuracy loss on the Fashion MNIST and 77.79% trainable parameters with a 1.75% accuracy decrease on the CIFAR-10. Similarly, by

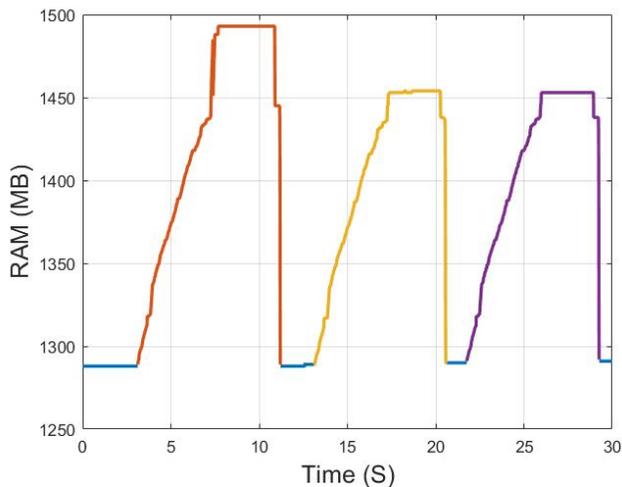


Fig. 8. Memory test on an NVIDIA Jetson Nano. From left to right, three peaks denote the inference memory of the 3×3 Conv2D (Max: 1493 MB), 2D-FWHT (Max: 1454 MB), and Squeeze-and-Excite layers (Max: 1453 MB). In each peak, after the layer is loaded, we run the inference code in a while-loop until we reach a steady state. RAM usage is 1288 MB when the device is vacant. Input and output tensors are in $\mathbb{R}^{10 \times 8 \times 8 \times 1024}$.

inserting one 2D-FWHT layer before the GAP layer and replacing Squeeze-and-Excite layers with the 2D-FWHT layers on the MobileNet-V3-Large, we save 49.96% trainable parameters with a 1.20% accuracy loss on the CIFAR-10, 48.55% trainable parameters with a 0.33% accuracy loss on the CIFAR-100, and 47.08% trainable parameters with a 1.59% accuracy decrease on the Tiny ImageNet.

By inserting a 2D-FWHT layer before the GAP layer in MobileNet-V3-Large, we achieve an increase of 2.48% in accuracy on CIFAR-100 with a slight increase (0.07%) in trainable parameters.

Furthermore, the non-residual 2D-FWHT layer can be also used as a replacement for the 3×3 convolution layers and the Squeeze-and-Excite layers. We can reduce nearly half trainable parameters for MobileNet-V3-Large with a negligible accuracy loss (0.33% accuracy loss on the CIFAR-100 dataset).

As a final point, we compare the speed of our 2D-FWHT layer with the regular 3×3 2D convolution layer. In an NVIDIA Jetson Nano board, our 2D-FWHT layer runs about 24 times as fast as the regular 2D convolution layer with 19.51% less RAM usage.

REFERENCES

- [1] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems*, 25:1097–1105, 2012.
- [2] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- [3] 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*, pages 1–9, 2015.
- [4] Fei Wang, Mengqing Jiang, Chen Qian, Shuo Yang, Cheng Li, Honggang Zhang, Xiaogang Wang, and Xiaoou Tang. Residual attention network for image classification. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3156–3164, 2017.
- [5] 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*, pages 770–778, 2016.
- [6] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Identity mappings in deep residual networks. In *European conference on computer vision*, pages 630–645. Springer, 2016.

- [7] Diaa Badawi, Hongyi Pan, Sinan Cem Cetin, and A Enis Cetin. Computationally efficient spatio-temporal dynamic texture recognition for volatile organic compound (voc) leakage detection in industrial plants. *IEEE Journal of Selected Topics in Signal Processing*, 14(4):676–687, 2020.
- [8] Chirag Agarwal, Shahin Khobahi, Dan Schonfeld, and Mojtaba Soltanalian. Coronet: a deep network architecture for enhanced identification of covid-19 from chest x-ray images. In *Medical Imaging 2021: Computer-Aided Diagnosis*, volume 11597, page 1159722. International Society for Optics and Photonics, 2021.
- [9] Harris Partaourides, Kostantinos Papadamou, Nicolas Kourtellis, Ilias Leontiadis, and Sotirios Chatzis. A self-attentive emotion recognition network. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 7199–7203. IEEE, 2020.
- [10] Dimitrios Stamoulis, Ting-Wu Chin, Anand Krishnan Prakash, Haocheng Fang, Sribhuvan Sajja, Mitchell Bogner, and Diana Marculescu. Designing adaptive neural networks for energy-constrained image classification. In *Proceedings of the International Conference on Computer-Aided Design*, pages 1–8, 2018.
- [11] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You only look once: Unified, real-time object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 779–788, 2016.
- [12] Süleyman Aslan, Uğur Gündükbay, B Uğur Töreyn, and A Enis Cetin. Deep convolutional generative adversarial networks for flame detection in video. In *International Conference on Computational Collective Intelligence*, pages 807–815. Springer, 2020.
- [13] Guglielmo Menchetti, Zhanli Chen, Diana J Wilkie, Rashid Ansari, Yasemin Yardimci, and A Enis Cetin. Pain detection from facial videos using two-stage deep learning. In *2019 IEEE Global Conference on Signal and Information Processing (GlobalSIP)*, pages 1–5. IEEE, 2019.
- [14] Süleyman Aslan, Uğur Gündükbay, B Uğur Töreyn, and A Enis Cetin. Early wildfire smoke detection based on motion-based geometric image transformation and deep convolutional generative adversarial networks. In *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 8315–8319. IEEE, 2019.
- [15] Changqian Yu, Jingbo Wang, Chao Peng, Changxin Gao, Gang Yu, and Nong Sang. Bisenet: Bilateral segmentation network for real-time semantic segmentation. In *Proceedings of the European conference on computer vision (ECCV)*, pages 325–341, 2018.
- [16] Zilong Huang, Xinggang Wang, Lichao Huang, Chang Huang, Yunchao Wei, and Wenyu Liu. Ccnet: Criss-cross attention for semantic segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 603–612, 2019.
- [17] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3431–3440, 2015.
- [18] Rudra PK Poudel, Stephan Liwicki, and Roberto Cipolla. Fast-scn: fast semantic segmentation network. *arXiv preprint arXiv:1902.04502*, 2019.
- [19] Yinli Jin, Wenbang Hao, Ping Wang, and Jun Wang. Fast detection of traffic congestion from ultra-low frame rate image based on semantic segmentation. In *2019 14th IEEE Conference on Industrial Electronics and Applications (ICIEA)*, pages 528–532. IEEE, 2019.
- [20] Hongyi Pan, Diaa Badawi, and Ahmet Enis Cetin. Fast walsh-hadamard transform and smooth-thresholding based binary layers in deep neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4650–4659, 2021.
- [21] 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*, pages 4510–4520, 2018.
- [22] Jiaying Wu, Cong Leng, Yuhang Wang, Qinghao Hu, and Jian Cheng. Quantized convolutional neural networks for mobile devices. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4820–4828, 2016.
- [23] Usama Muneeb, Erdem Koyuncu, Yasaman Keshkarjahromd, Hulya Seferoglu, Mehmet Fatih Erden, and A Enis Cetin. Robust and computationally-efficient anomaly detection using powers-of-two networks. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 2992–2996. IEEE, 2020.
- [24] Wenlin Chen, James Wilson, Stephen Tyree, Kilian Weinberger, and Yixin Chen. Compressing neural networks with the hashing trick. In *International conference on machine learning*, pages 2285–2294. PMLR, 2015.
- [25] Hongyi Pan, Diaa Badawi, and Ahmet Enis Cetin. Computationally efficient wildfire detection method using a deep convolutional network pruned via fourier analysis. *Sensors*, 20(10):2891, 2020.
- [26] Mingchao Yu, Zhifeng Lin, Krishna Narra, Songze Li, Youjie Li, Nam Sung Kim, Alexander Schwing, Murali Annamaram, and Salman Avestimehr. Gradiveq: Vector quantization for bandwidth-efficient gradient aggregation in distributed cnn training. *arXiv preprint arXiv:1811.03617*, 2018.
- [27] 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*, 2019.
- [28] 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.
- [29] Matthieu Courbariaux, Itay Hubara, Daniel Soudry, Ran El-Yaniv, and Yoshua Bengio. Binarized neural networks: Training deep neural networks with weights and activations constrained to +1 or -1. *arXiv preprint arXiv:1602.02830*, 2016.
- [30] Adrian Bulat and Georgios Tzimiropoulos. Hierarchical binary cnns for landmark localization with limited resources. *IEEE transactions on pattern analysis and machine intelligence*, 42(2):343–356, 2018.
- [31] Mohammad Rastegari, Vicente Ordonez, Joseph Redmon, and Ali Farhadi. Xnor-net: Imagenet classification using binary convolutional neural networks. In *European conference on computer vision*, pages 525–542. Springer, 2016.
- [32] Zhiqiang Shen, Zechun Liu, Jie Qin, Lei Huang, Kwang-Ting Cheng, and Marios Savvides. S2-bnn: Bridging the gap between self-supervised real and 1-bit neural networks via guided distribution calibration. *arXiv preprint arXiv:2102.08946*, 2021.

- [33] Zechun Liu, Zhiqiang Shen, Marios Savvides, and Kwang-Ting Cheng. Reactnet: Towards precise binary neural network with generalized activation functions. In *European Conference on Computer Vision*, pages 143–159. Springer, 2020.
- [34] Brais Martinez, Jing Yang, Adrian Bulat, and Georgios Tzimiropoulos. Training binary neural networks with real-to-binary convolutions. *arXiv preprint arXiv:2003.11535*, 2020.
- [35] Adrian Bulat, Brais Martinez, and Georgios Tzimiropoulos. Bats: Binary architecture search. *arXiv preprint arXiv:2003.01711*, 2020.
- [36] Itay Hubara, Matthieu Courbariaux, Daniel Soudry, Ran El-Yaniv, and Yoshua Bengio. Binarized neural networks. In *Proceedings of the 30th International Conference on Neural Information Processing Systems*, pages 4114–4122, 2016.
- [37] Milad Alizadeh, Javier Fernández-Marqués, Nicholas D Lane, and Yarín Gal. An empirical study of binary neural networks’ optimisation. In *International Conference on Learning Representations*, 2018.
- [38] Tom Bannink, Arash Bakhtiari, Adam Hillier, Lukas Geiger, Tim de Bruin, Leon Overweel, Jelmer Neeven, and Koen Helwegen. Larq compute engine: Design, benchmark, and deploy state-of-the-art binarized neural networks. *arXiv preprint arXiv:2011.09398*, 2020.
- [39] Felix Juefei-Xu, Vishnu Naresh Boddeti, and Marios Savvides. Local binary convolutional neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 19–28, 2017.
- [40] Xiaofan Lin, Cong Zhao, and Wei Pan. Towards accurate binary convolutional neural network. *arXiv preprint arXiv:1711.11294*, 2017.
- [41] Ziwei Wang, Jiwen Lu, Chenxin Tao, Jie Zhou, and Qi Tian. Learning channel-wise interactions for binary convolutional neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 568–577, 2019.
- [42] Ritchie Zhao, Yuwei Hu, Jordan Dotzel, Christopher De Sa, and Zhiru Zhang. Building efficient deep neural networks with unitary group convolutions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11303–11312, 2019.
- [43] Shamma Nasrin, Diaa Badawi, Ahmet Enis Cetin, Wilfred Gomes, and Amit Ranjan Trivedi. Mf-net: Compute-in-memory sram for multibit precision inference using memory-immersed data conversion and multiplication-free operators. *IEEE Transactions on Circuits and Systems I: Regular Papers*, 68(5):1966–1978, 2021.
- [44] Maneesh Ayi and Mohamed El-Sharkawy. Rmnv2: Reduced mobilenet v2 for cifar10. In *2020 10th Annual Computing and Communication Workshop and Conference (CCWC)*, pages 0287–0292. IEEE, 2020.
- [45] Pravendra Singh, Vinay Kumar Verma, Piyush Rai, and Vinay P Nambodiri. Hetconv: Heterogeneous kernel-based convolutions for deep cnns. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4835–4844, 2019.
- [46] T Ceren Deveci, Serdar Cakir, and A Enis Cetin. Energy efficient hadamard neural networks. *arXiv preprint arXiv:1805.05421*, 2018.
- [47] A Enis Cetin, Omer N Gerek, and Sennur Ulukus. Block wavelet transforms for image coding. *IEEE Transactions on Circuits and Systems for Video Technology*, 3(6):433–435, 1993.
- [48] Joseph L Walsh. A closed set of normal orthogonal functions. *American Journal of Mathematics*, 45(1):5–24, 1923.
- [49] Bernard J. Fino and V. Ralph Algazi. Unified matrix treatment of the fast walsh-hadamard transform. *IEEE Transactions on Computers*, 25(11):1142–1146, 1976.
- [50] Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. Aggregated residual transformations for deep neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1492–1500, 2017.
- [51] Christian Szegedy, Sergey Ioffe, Vincent Vanhoucke, and Alexander A Alemi. Inception-v4, inception-resnet and the impact of residual connections on learning. In *Thirty-first AAAI conference on artificial intelligence*, 2017.
- [52] PM Agante and JP Marques De Sá. Ecg noise filtering using wavelets with soft-thresholding methods. In *Computers in Cardiology 1999. Vol. 26 (Cat. No. 99CH37004)*, pages 535–538. IEEE, 1999.
- [53] David L Donoho. De-noising by soft-thresholding. *IEEE transactions on information theory*, 41(3):613–627, 1995.
- [54] Michael W Marcellin, Michael J Gormish, Ali Bilgin, and Martin P Boliek. An overview of jpeg-2000. In *Proceedings DCC 2000. Data Compression Conference*, pages 523–541. IEEE, 2000.
- [55] James Lee-Thorp, Joshua Ainslie, Ilya Eckstein, and Santiago Ontanon. Fnet: Mixing tokens with fourier transforms. *arXiv preprint arXiv:2105.03824*, 2021.
- [56] Mobilenet in tensorflow’s official github. <https://github.com/tensorflow/models/tree/master/research/slim/nets/mobilenet>, 2019. Accessed: 2021-03-01.
- [57] Transfer learning and fine-tuning. https://www.tensorflow.org/tutorials/images/transfer_learning, 2020. Accessed: 2021-03-01.
- [58] 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/CVF International Conference on Computer Vision*, pages 1314–1324, 2019.
- [59] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *Proceedings of the IEEE international conference on computer vision*, pages 1026–1034, 2015.
- [60] Block walsh-hadamard transform layer speed test code. <https://github.com/phy710/Block-Walsh-Hadamard-Transform-Layer-Speed-Test>, 2021. Accessed: 2021-12-31.