BINARIZED CONVOLUTIONAL NEURAL NETWORKS WITH CHANNEL QUADRUPLING AND SMOOTH DOWN SAMPLING

Anonymous authors

Paper under double-blind review

Abstract

This paper proposes novel binarized convolutional neural networks (BCNNs) named **OB-Net** and **OSB-Net**, specifically designed to **Q**uadruple the number of channels and incorporate a so-called Smooth downsampling in BCNNs for lowcost mobile environments. The proposed models combine FP32 depthwise separable (DS) convolutions with binarized 1×1 pointwise convolutions, offering reduced computational costs in the pointwise convolutions. To enhance the degraded performance of the above naive combination, the proposed models start with a small number of channels in shallow layers and expand them during downsampling by a factor of four, effectively managing model complexity in the downsampling. The proposed model structure maintains low computational costs in the shallow blocks and increases model complexity in the deep blocks, providing a wider dynamic range to manage information in the frequency domain. As a result, the proposed models overcome the limitations of existing BCNNs, delivering improved performance while reducing the total computational costs. For further performance enhancements, we propose a novel smooth downsampling with heightwise and widthwise sequential downsampling steps, doubling the number of channels at each step. Besides, we show that the channelwise self-attention (SE) is applicable with minimal additional computational costs in the proposed models. Besides, multiple binarized convolutions in the fully-connected (FC) layer reduce storage costs without requiring 8-bit quantized convolutions. Experimental results demonstrate the efficiency of the proposed models in terms of performance, computational costs, and inference latency on real hardware. Notably, the QSB-Net-Large with SE achieve 71.2% Top-1 accuracy on ImageNet-1K and 69.2 mean intersection over union (mIoU) in the semantic segmentation on the PASCAL VOC dataset, outperforming other counterparts.

036

038

006

008 009 010

011

013

014

015

016

017

018

019

021

023

024

025

026

027

028

029

031

032

034

1 INTRODUCTION

039 Although substantial parallelism in GPUs or specialized accelerators can achieve a considerable 040 speedup, edge devices lack sufficient parallelism for accommodating the increasing model com-041 plexity. BCNNs binarize both weights and activations into -1 and +1 in binarized convolutions. As 042 a result, multiply-accumulate operations are replaced by bitwise XNOR and bit-counting operations, 043 leveraging bit-level parallelism in CPU-based devices. While earlier BCNNs experienced significant 044 accuracy drops compared with their FP32 counterparts, recent BCNNs achieve good performance comparable to mobile-friendly CNNs. However, as computational costs increase within BCNNs, the expected inference speedup could not be satisfactory. For example, whereas ReActNetA (Liu et al., 046 2020) adopts binarized 3×3 spatial convolutions, its FP32 counterpart, MobileNetV1 (Howard 047 et al., 2017) utilizes FP32 DS spatial convolutions. Compared with FP32 DS convolutions, bina-048 rized 3×3 spatial convolutions require more computations, failing to provide faster inference on real hardware. Therefore, novel lightweight structures should be studied to minimize computational costs while avoiding significant performance degradation in BCNNs. 051

Conventional CNNs typically double the number of channels or adopt customized channel expansion during downsampling (He et al., 2016; Ridnik et al., 2021) to mitigate information loss. In BCNNs, however, we expect that the binarization error and low resolution of BCNNs require a dif-



Figure 1: Main idea and comparison in terms of Top-1 accuracy and OPs on ImageNet-1K (Russakovsky et al., 2015). In Figure 1 (b), the model in parentheses is the FP32 counterpart of its binarized model.

ferent strategy for performing the channel expansion and downsampling. Our hypothesis suggests that the performance of BCNNs can be significantly influenced by the complexity of deep blocks 071 and the configuration of downsampling. Although increasing the number of channels of all blocks 072 by the same ratio can mitigate the issue, computational costs dramatically increase. Therefore, we 073 think that developing a novel strategy for channel expansion and downsampling has the potential 074 to achieve significant benefits in BCNNs. The approach in the conceptual illustration of Figure 1 075 (a) quadruply increases the number of channels in deeper blocks during downsampling. By hav-076 ing a small number of channels in the shallow block, the proposed strategy can reduce the total 077 computations compared with existing BCNNs.

This paper proposes novel BCNNs named QB-Net and QSB-Net, quadrupling the number of channels in channel expansion and configuring smooth downsampling in the new structures of BCNNs. The newly developed designs and main contributions are as follows:

- 1. New Structure for Channel Quadrupling: To reduce the computations in pointwise convolutions, the proposed models utilize 3×3 DS convolutions and binarized pointwise convolutions. However, the existing channel expansion with the above structure shows significant performance degradation. To overcome the problem, the proposed models have a small number of channels in the shallow blocks and increase the number of channels by a factor of four during downsampling. The small number of channels in the shallow blocks reduces the total computations, as shown in Figure 1 (b). The channel quadrupling can provide a wider dynamic range in deep blocks, enhancing the ability to manage the information in the frequency domain and achieving good performance in BCNNs.
- 2. Smooth Downsampling and Techniques for Better Models: In QSB-Net, the proposed smooth downsampling processes two sequential steps involving heightwise and widthwise 1-D downsampling, producing noticeable performance enhancements. Each 1-D downsampling layer can double the number of channels, resulting in quadrupled channels after the two steps of downsampling in Figure 1 (a). Experimental results show that channelwise self-attention in DS convolutions is effective, as shown in Figure 1 (b). Besides, 8-bit quantization can be valid for the FP32 DS convolutional and FC layers. Notably, the usage of multiple binarized convolutions in the FC layer can reduce storage costs, which can solve the issue of the increasing number of channels in the FC layer.
- 3. Experiments on Various Datasets and Real Hardware: We evaluated the proposed models for image classification and semantic segmentation. The experiments on a Raspberry Pi 4B (RPi 4B) and a Samsung Exynos processor using Larq Compute Engine (LCE) (Bannink et al., 2021) show that the proposed models have significant inference speedup compared with the baseline model and mobile-friendly CNNs.
- 103 104

102

068

069

082

084

090

092

095

096

098

099

The proposed model reaches up to 71.2% Top-1 accuracy on ImageNet-1K. and 69.2 mIoU in semantic segmentation on the PASCAL VOC dataset, outperforming other counterparts. The above results prove that the proposed channel quadrupling and smooth downsampling can be effective in BCNNs.

108 2 RELATED WORKS

110 Courbariaux et al. (2016) showed that BCNNs can reduce memory usage and energy consumption by 111 nearly $32 \times$ compared with baselined FP32 models. Many existing works have focused on improv-112 ing the performance of BCNNs due to their significant performance drops. XNOR-Net (Rastegari 113 et al., 2016) adopted a scaling factor for both weights and input features of binarized convolutional layers, achieving 51.4% Top-1 accuracy with binarized ResNet18 (He et al., 2016) on ImageNet-1K. 114 Bi-RealNet (Liu et al., 2018) utilized the single skipped shortcut, enhancing the Top-1 accuracy of 115 binarized ResNet18 up to 56.4%. XNOR-Net++ (Bulat et al., 2019) advanced BCNNs by introduc-116 ing heightwise and widthwise scaling factors for the output features. Real-to-Bin (Martinez et al., 117 2019) enhanced the performance up to 65.4% by incorporating a self-attention block. 118

119 Whereas the above BCNNs are based on ResNet (He et al., 2016), MobiNet (Phan et al., 2020a) and ReActNetA (Liu et al., 2020) followed the structure of MobileNetV1 (Howard et al., 2017). 120 However, MobiNet produced only 53.5% Top-1 accuracy on ImageNet-1K. ReActNetA (Liu et al., 121 2020) achieved 69.4% Top-1 accuracy by deploying binarized 3×3 convolutions instead of us-122 ing FP32 DS convolutions. The outstanding performance of Real-to-Bin and ReActNetA showed 123 the effectiveness of teacher-student training (Hinton et al., 2015) in BCNNs. There have been sev-124 eral works to develop specific activation functions and training methods for BCNNs. IR-Net (Qin 125 et al., 2020) introduced a binarization method and so-called error decaying estimator for the train-126 ing. SI-BNN (Wang et al., 2020) proposed a binarization function with trainable thresholds. BNSC-127 Net (Wu et al., 2021) decomposed 2-D convolutions to employ additional skip connections. RB-128 Net (Liu et al., 2022) reshaped pointwise convolutions and incorporated a balanced activation func-129 tion. ReCU (Xu et al., 2021) employed a rectified clamp unit to improve its training results. SA-130 BNN (Liu et al., 2021) mitigated the weight flip issues in BCNNs. AdaBin (Tu et al., 2022) adap-131 tively optimized weights and features to follow the value distributions of its FP32 baseline model. ReBNN (Xu et al., 2023) introduced weighted reconstruction loss to reduce the weight oscillation 132 during training. Whereas PokeBNN (Zhang et al., 2022) achieved impressive results by adopting 133 ResNet50 as its baseline and teacher models, its hyperparameter for the clipping bound and stride 134 configuration in the first convolutional layer introduced a different optimization strategy. Different 135 from Zhang et al. (2022), the proposed QSB-Net shows a novel strategy for channel quadrupling and 136 smooth downsampling. Therefore, we conclude that the fundamental techniques and contributions 137 of the proposed models are totally different from those of PokeBNN. 138

139 140

3 BACKGROUNDS

For a binarized convolution in BCNNs, both filter weights and input features are quantized into binary values. For given binarized input features I_b and binarized $K \times K$ filters F_b , the conventional binarized convolution, denoted as $I_b * F_b$, can be formulated as:

$$(I_b * F_b)(i,j) = \gamma \cdot \sum_{c=0}^{C-1} \sum_{m=0}^{K-1} \sum_{n=0}^{K-1} I_b(c,i+m,j+n) \cdot F_b(c,m,n),$$
(1)

where i, j are the spatial indices of the output feature map in a channel. The term γ is the scaling factor of the binarized convolutional outputs. During inference, whereas pre-binarized weights F_b can be pre-stored, input features I are binarized into I_b during feedforwarding. While so-called binarization-aware training updates real-valued weights in the backpropagation during training, the updated real-valued weights are binarized during feedforwording. Let $\epsilon(a, b)$ be the error between a and b. Binarization errors of $\epsilon(F(c, m, n), F_b(c, m, n))$ exist between the real-valued updated weights and the binarized weights used during feedforwarding.

154 Several works studied better value distributions for the binarized convolution (Helwegen et al., 2019; 155 Liu et al., 2021; Tu et al., 2022; Xu et al., 2023; Zhang et al., 2022). However, the calibration of 156 value distributions requires specific hyperparameters, so we think that there could be no consistently 157 valid rules. A customized sign function (Liu et al., 2018; 2020) requires a long model training time 158 due to the complex derivatives of the customized sign function. We do not focus on new methods for calibrating value distributions and customized activation functions. Without using the above new 159 methods and hyperparameters, we adopt a straight-through estimator (STE) (Bengio et al., 2013) to 160 approximate the derivatives of the sign conversion. Besides, we assume that the output features of 161 the binarized convolution are scaled with channelwise learnable scaling factors.

162 In order to show the computational complexity, the number of floating-point operations (FLOPs) has 163 been generally used. In BCNNs, when BOPs denote the number of binarized multiply-accumulate 164 operations, $OPs = FLOPs + \frac{BOPs}{64}$ are used to estimate the total computations. However, the 165 latencies on real hardware platforms were not reduced by a factor of $\frac{1}{64}$ (Bannink et al., 2021). On 166 the other hand, entropy is commonly used to quantify the average amount of information produced 167 by a stochastic data source. Although the entropy could not be directly related to the model perfor-168 mance, we believe that it could demonstrate the information characteristics in the proposed models. To explain the information characteristics of the proposed model, the entropy of the output features 169 170 in each block is analyzed in Appendix A.3.

171 In BCNNs, increasing model complexity has been proven to produce better model performance, 172 suppressing negative effects of the binarization error (Lin et al., 2017). In other words, when the 173 number of values to be summed increases in the binarized convolution in Eq.(1), the output repre-174 sentation capacity can be enhanced, offsetting the binarization errors. The enhanced representation capacity can capture different aspects of features that exist in a wide range of the frequency domain. 175 Although the information loss due to the low resolution of the binarized convolutions can be mit-176 igated by capturing the features in the high-frequency domain, additional computational costs are 177 required. By increasing the depth multiplier (Howard et al., 2017), the number of channels in all 178 blocks is multiplied, which can enhance the performance of BCNNs. However, the total computa-179 tions quadruply increase, which can be a burden in mobile-friendly BCNNs. 180

181

183

185

4 QSB-NET WITH CHANNEL QUADRUPLING & SMOOTH DOWNSAMPLING

4.1 MOTIVATIONS

186 The depth multiplier for multiplying the number of channels in all blocks dramatically increases the 187 total computations, having long inference latency in mobile-friendly BCNNs. Whereas most exist-188 ing BCNNs have focused on training techniques and specific blocks on the baseline model structure, 189 the efficient BCNN structures were not sufficiently studied. To reduce the total computations, the structural development of BCNNs should mainly increase the model complexity of the efficient part 190 that significantly contributes to model performance. On the other hand, we hypothesize that informa-191 tion loss during downsampling can impact model performance. Although channel expansion during 192 downsampling can improve representation capacity, it is questionable whether the doubled or cus-193 tomized channel expansions in conventional CNNs are the most effective in BCNNs. Besides, the 194 binarization error could compound the effects of information loss during downsampling. Therefore, 195 we develop a novel strategy for configuring channel expansion and downsampling in BCNNs.

196 197 198

4.2 BLOCKS FOR CHANNEL QUADRUPLING AND SMOOTH DOWNSAMPLING

199 Depending on the configurations of the number of channels and smooth downsampling, we develop 200 four models named QB-Net-Small, QB-Net-Large, QSB-Net-Small, and QSB-Net-Large. Figure 2 201 illustrates the block structures used in QB-Net. The structure of MobileNetV1 (Howard et al., 2017) 202 inspires the usage of DS convolutions in BCNNs. MobileNetV1 (Howard et al., 2017) deploys 203 both DS convolutions and 1×1 pointwise convolutions. Besides, DS convolutions account for 204 only 3% of the total computations in MobileNetV1. Several existing BCNNs were motivated by 205 the structure of MobileNet (Phan et al., 2020a;b; Liu et al., 2020). However, these works adopted 206 binarized 3 × 3 convolution (Liu et al., 2020) or grouped binarized convolution (Phan et al., 2020a;a) 207 instead of using FP32 DS convolution. In Liu et al. (2020), as the number of channels increased, the binarized 3×3 convolutions dramatically increase the total computations compared with the usage 208 of DS convolutions. In Phan et al. (2020a;b), the grouped binarized DS convolutions did not achieve 209 acceptable accuracy compared with MobileNetV1. 210

Considering the small computational portion of the DS convolutions in MobileNetV1, we do not
prioritize the binarization of DS convolutions. Therefore, we determine that the proposed models
deploy FP32 3 × 3 DS convolutional layer (DSCONV) in a block. However, a naive modification
of the baseline model can show significant performance degradation. When applying FP32 DS convolutions instead of binarized 3 × 3 convolutions in ReActNetA (Liu et al., 2020), the modification
of ReActNetA only achieved 64.2% Top-1 accuracy on ImageNet-1K using teacher-student train-

239 240

241

242

243

244

245

246

247

248

249

250

255



Figure 2: Illustration of blocks used in QS-Net. Terms DSCONV and BCONV refer to the DS convolutional 231 and binarized convolutional layers, respectively. Suffixes N, D, and Q mean normal, double in downsampling, 232 and quadruple in downsampling, respectively. While block TypeN does not change the shape of output features, blocks TypeD and TypeQ expand channels along with downsampling feature maps. When the number 233 of input channels C is small (e.g., 16 or 32), 3×3 BCONV is used instead of 1×1 BCONV. The values in 234 parentheses are the heightwise and widthwise strides, respectively. Terms H and W denote the height and 235 width of feature maps. During downsampling, the shortcut connection utilizes the 2-D average pooling layer 236 with a stride of (2, 2). The learnable bias layer (Liu et al., 2020) is deployed just before its convolutional layer. 237 BPReLU layers are used after the shortcut is added.



Figure 3: Illustration of blocks used in QSB-Net. Suffixes DS and QS denote double in smooth downsampling and quadruple in smooth downsampling, respectively. For the smooth downsampling, DS convolutions are sequentially performed with strides of (2, 1) and (1, 2). In the second DSCONV, whereas the number of channels is doubled in TypeQS, TypeDS does not expand the channels.

ing (Hinton et al., 2015). Therefore, a novel structural breakthrough is required to achieve better
 performance while maintaining small total computational costs.

Figure 2 illustrates **TypeN**, **TypeD**, and **TypeQ** blocks deployed in QB-Net. In the TypeD and TypeQ blocks, after performing downsampling with a stride of (2, 2) in DSCONVs, the output channels from DSCONVs are concatenated to expand the number of channels. After performing the convolutions, batch normalization (BN) is performed in both DSCONV and BCONV, which is not shown in Figure 2 for clarity.

263 While QB-Net deploys TypeN, TypeD, and TypeQ blocks, the proposed QSB-Net can utilize other 264 blocks named **TypeDS** and **TypeQS** during downsampling. Figure 3 illustrates the block structures 265 in QSB-Net, which support so-called smooth downsampling. Whereas QB-Net performs conven-266 tional downsampling with a stride of (2, 2), QSB-Net adopts two-stage downsampling to provide 267 additional model complexity during downsampling. In TypeQS, two DSCONVs are performed se-268 quentially with strides of (2, 1) and (1, 2). A BCONV is deployed after performing each DSCONV. 269 In a block, after performing the first DSCONV and channel concatenation, the shape of feature maps 261 can be $2 \times C \times \frac{H}{2} \times W$. Then, the downsampling in the second DSCONV produces the feature

Index	ReActNet, MobileNetV1			QB-Net-Small ¹			QSB-Net-Small ¹			
	Block	Input	Output	Block	Input	Output	Block	Input	Output	
1	Conv ²	$3 \times 224 \times 224$	$32 \times 112 \times 112$	Conv ²	$3 \times 224 \times 224$	$16\!\times\!112\!\times\!112$	Conv ²	$3 \times 224 \times 224$	$ 16 \times 112 \times 11 $	
2	Reduction ³	$32 \times 112 \times 112$	$64 \times 112 \times 112$	TypeN	$16 \times 112 \times 112$	$16\times112\times112$	TypeN	$16 \times 112 \times 112$	$16 \times 112 \times 1$	
3	Reduction	$64 \times 112 \times 112$	$128 \times 56 \times 56$	TypeD	$16 \times 112 \times 112$	$32 \times 56 \times 56$	TypeDS	$16 \times 112 \times 112$	$32 \times 56 \times 50$	
4	Normal	$128 \times 56 \times 56$	$128 \times 56 \times 56$	TypeN	$32 \times 56 \times 56$	$32 \times 56 \times 56$	TypeN	$32 \times 56 \times 56$	$32 \times 56 \times 50$	
5	Reduction	$128 \times 56 \times 56$	$256 \times 28 \times 28$	TypeQ	$32 \times 56 \times 56$	$128 \times 28 \times 28$	TypeQS	$32 \times 56 \times 56$	$128 \times 28 \times 2$	
6	Normal	$256 \times 28 \times 28$	$256 \times 28 \times 28$	TypeN	$128 \times 28 \times 28$	$128 \times 28 \times 28$	TypeN	$128 \times 28 \times 28$	$128 \times 28 \times 2$	
7	Reduction	$256 \times 28 \times 28$	$512 \times 14 \times 14$	TypeQ	$128 \times 28 \times 28$	$512 \times 14 \times 14$	TypeQS	$128 \times 28 \times 28$	$512 \times 14 \times 12$	
8	Normal	$512 \times 14 \times 14$	$512 \times 14 \times 14$	TypeN	$512 \times 14 \times 14$	$512\!\times\!14\!\times\!14$	TypeN	$512 \times 14 \times 14$	$512 \times 14 \times 10^{-10}$	
9	Normal	$512 \times 14 \times 14$	$512 \times 14 \times 14$	TypeN	$512 \times 14 \times 14$	$512\!\times\!14\!\times\!14$	TypeN	$512 \times 14 \times 14$	$512 \times 14 \times 1$	
10	Normal	$512 \times 14 \times 14$	$512 \times 14 \times 14$	TypeN	$512 \times 14 \times 14$	$512\!\times\!14\!\times\!14$	TypeN	$512 \times 14 \times 14$	$512 \times 14 \times 1$	
11	Normal	$512 \times 14 \times 14$	$512 \times 14 \times 14$	TypeN	$512 \times 14 \times 14$	$512\!\times\!14\!\times\!14$	TypeN	$512 \times 14 \times 14$	$512 \times 14 \times 1$	
12	Normal	$512 \times 14 \times 14$	$512 \times 14 \times 14$	TypeN	$512 \times 14 \times 14$	$512 \times 14 \times 14$	TypeN	$512 \times 14 \times 14$	$512 \times 14 \times 1$	
13	Reduction	$512 \times 14 \times 14$	$1024 \times 7 \times 7$	TypeQ	$512 \times 14 \times 14$	$2048 \times 7 \times 7$	TypeQS	$512 \times 14 \times 14$	$2048 \times 7 \times$	
14	Normal	$1024 \times 7 \times 7$	$1024 \times 7 \times 7$	TypeN	$2048 \times 7 \times 7$	$2048 \times 7 \times 7$	TypeN	$2048 \times 7 \times 7$	$2048 \times 7 \times$	
15	AvgPool	$1024 \times 7 \times 7$	$1024 \times 1 \times 1$	AvgPool	$2048 \times 7 \times 7$	$2048 \times 1 \times 1$	AvgPool	$2048 \times 7 \times 7$	$2048 \times 1 \times$	
16	FC	$1024 \times 1 \times 1$	1000	FC	$2048 \times 1 \times 1$	1000	FC	$2048 \times 1 \times 1$	1000	

Table 1: Comparison of model structures. QB-Net-Small and QSB-Net-Small quadruple the number of chan-nels in TypeQ and TypeQS blocks. In the Input and Output columns, each term denotes $C \times H \times W$ in Figure 2.

¹ Whereas QB-Net-Small and QSB-Net-Small adopt TypeQ and TypeQS in the 13-th block, QB-Net-Large and QSB-Net-Large deploy them in the 9-th

block and set the number of channels as 2048 from the 10-th block The first block denoted as Conv consists of the conventional FP32 convolution, batch normalization, and ReLU layers





Figure 4: OPs of convolutions in each block. Yellow regions indicate the blocks that perform downsampling. Whereas Figure 4 (a) shows the OPs of the DS convolution in the proposed models, Figure 4 (b) visualizes the OPs of the next binarized pointwise convolution. In Figure 4, the proposed QB-Net-Small, QB-Net-Large, QSB-Net-Small, and QSB-Net-Large do not decrease OPs in the downsampling blocks. Other counterparts such as MobileNetV1 (Howard et al., 2017) and ReActNetA (Liu et al., 2020) have small OPs in the downsam-pling blocks. It is noted that only QB-Net-Large and QSB-Net-Large perform downsampling in the 9-th block.

maps of $4 \times C \times \frac{H}{2} \times \frac{W}{2}$ after performing the second channel concatenation. On the other hand, in the second DSCONV of TypeDS, the number of channels is not doubled.

Several layers in the blocks are based on previous works as follows: the single skipped shortcut (Liu et al., 2018) is used for each DSCONV and BCONV, where symbol \oplus denotes the elementwise addition. A learnable bias (Liu et al., 2020) is used to calibrate the distribution of input features in each channel. The activation layer with PReLU-BN, denoted as BPReLU, is deployed after the addition with a shortcut. In Phan et al. (2020a), the activation layer with PReLU-BN is used. However, its effects were not compared with other works such as RPReLU (Liu et al., 2020). In our experiments, BPReLU slightly outperformed RPReLU by 0.1%-0.3% in the image classification on ImageNet-1K.

4.3 MODEL STRUCTURE WITH CHANNEL QUADRUPLING

In order to achieve better performance, the number of channels can be expanded to increase the rep-resentative capacity. However, when the numbers of channels in all blocks are naively multiplied, its computational costs dramatically increase. Therefore, the model structure is important to deter-mine which block has more channels, considering both computational costs and model performance. In existing FP32 models, the channel expansion during downsampling is either doubled (He et al., 2016) or customized (Sandler et al., 2018; Huang et al., 2019; Tan & Le, 2019). However, many existing BCNNs are based on the structure of ResNet (He et al., 2016) or MobileNetV1 (Howard et al., 2017), so previous studies about new structures seem to be insufficient. Table 1 shows the model structures of the proposed QS-Net-Small and QSB-Net-Small, comparing with MobileNetV1



Figure 5: Frequency domain histograms of proposed models and counterparts in the outputs of the third (left) and final (right) downsampling blocks. The output features of a goose image are hooked and used to obtain the frequency domain histograms. In FP32 models, wide dynamic ranges are shown in the output of the final downsampling blocks. Whereas ReActNetA (Liu et al., 2020) shows a narrow range in the frequency domain in the output of the final downsampling block, the proposed models have wide dynamic ranges. The proposed Large models have wider dynamic ranges than the Small models in the frequency domain.

and ReActNet (Liu et al., 2020). Because QB-Net-Small and QSB-Net-Small adopt TypeQ and TypeQS in the 13-th block, TypeQ and TypeQS blocks have $2048 \times 7 \times 7$ output feature maps. Compared with MobileNetV1 and ReActNet, when H and W of feature maps are large, the number of channels C is small. On the other hand, when H and W of feature maps are small, C increases by quadrupling channels during downsampling. QB-Net-Large and QSB-Net-Large adopt TypeQ and TypeQS in the 9-th block and set the number of output channels as 2048. The model structures of QB-Net-Large and QSB-Net-Large are summarized in Appendix A.1.

347 348 349

350

351

340 341

342

343

344

345

346

4.4 INCREASED REPRESENTATION CAPACITY IN DEEP BLOCKS

352 Figure 4 illustrates the OPs of the convolutions in each block, representing the main idea of the 353 OB-Net and OSB-Net. Whereas the counterparts such as MobileNetV1 (Howard et al., 2017) and 354 ReActNetA (Liu et al., 2020) reduce OPs in the downsampling blocks, the channel expansion in 355 QB-Net does not decrease OPs in the downsampling block. Moreover, additional computations are 356 performed in the TypeDS and TypeQS blocks. We expect that the representation capacity in a block 357 of BCNNs during downsampling can have a significant impact on the model performance. The 358 low resolution of binarized convolutions could degrade representation capacity, so it is thought that the low resolution can be compensated with additional model complexity in the proposed blocks. 359 The proposed model structures increase model complexity for deep blocks during downsampling. 360 Several existing works show that increasing representation capacity during downsampling is helpful 361 for achieving better performance in BCNNs. For example, XNOR-Net (Rastegari et al., 2016) and 362 Bi-RealNet (Liu et al., 2018) deploy FP32 1×1 convolutional layer for the shortcut connection. However, the shortcut connection with FP32 1×1 convolutions significantly increases the total 364 OPs, showing long latencies on real mobile hardware (Bannink et al., 2021). We note that the channel quadrupling can increase the representation capacity during downsampling without using 366 the FP32 1×1 convolutions. It is expected that the proposed structure using the channel quadrupling 367 can have a wider dynamic range in the frequency domain from the increased complexity in deep 368 blocks. Although the number of channels in the deep blocks increases, the total computations are significantly reduced by having a small number of channels in shallow blocks. 369

370 To show the effects of the increased representation capacity and wider dynamic range in deep blocks, 371 the frequency domain histograms are compared. Figure 5 shows the frequency domain histograms 372 of proposed models and counterparts in the outputs of the third (left) and final (right) downsampling 373 blocks. Notably, the proposed models have wider dynamic ranges than ReActNetA (Liu et al., 374 2020). The proposed QSB-Net-Large and QSB-Net-Large(SE1) have dynamic ranges comparable 375 to those of FP32 models. However, the similarity between Figures 5 (g) and (h) indicates that the self-attention block does not widen the dynamic range in the frequency domain. Therefore, we 376 conclude that the role of the channel quadrupling is critical to obtain wide dynamic ranges in the 377 frequency domain. Besides, the feature maps of deep blocks are visualized in Appendix A.8.

381	Model	Top-1	FLOPs	BOPs	OPs	Model	Top-1	FLOPs	BOPs	OPs
382	WIOUCI	(%)	$(\times 10^{8})$	$(\times 10^8)$	$(\times 10^8)$	Widdel	(%)	$(\times 10^8)$	$(\times 10^8)$	$(\times 10^8)$
383	ResNet18	69.6	18.2	-	18.2	MobileNetV1	70.6	5.75	-	5.75
384	MobileNetV2	71.8	3.0	-	3.0	SuffleNetV2	69.4	1.46	-	1.46
007	GhostNet	73.9	1.41	-	1.41	AlphaNet-A0	77.8	2.03	-	2.03
385	XNOR-Net(ResNet18)	51.2	1 4 1	17.0	1.67	MobiNet-Mid(MobileNetV1) ¹	54.4	0.52	_	0.52
386	Bi-RealNet(ResNet18)	56.4	1.54	16.8	1.63	SI-BNN(Bi-RealNet)	59.7	1.54	16.8	1.63
387	XNOR-Net++(ResNet18)	57.1	1.41	17.0	1.67	BNSC-Net(ResNet18)	59.9	-	-	1.40
000	RB-Net(ResNet18)	66.8	-	-	0.52	RB-Net(ResNet34)	70.2	-	-	0.71
388	IR-Net(ResNet18)	58.1	-	-	1.63	SA-BNN(ResNet18)	61.7	-	-	1.69
389	Real-to-Bin(ResNet18)	65.4	1.56	1.68	1.83	AdaBin(ResNet18)	66.4	-	-	0.88
200	ReBNN(ResNet18)	66.9	-	-	1.63	ReActNetA(MobileNetV1) ¹	69.4	0.12	48.2	0.87
390	PokeBNN×0.75(ResNet50) ²	70.5	0.10	20.3	0.42	ReCU(ResNet18)	66.4	-	-	1.63
391	OB-NFT-Small	66.9	0.28	93	0.43	OSB-NET-Small	68.8	0.36	10.8	0.53
392	OB-NET-Large	69.8	0.20	15.5	0.53	OSB-NET-Large	70.6	0.36	16.9	0.62
393	QSB-NET-Large(SE1)	71.1	0.37	16.9	0.64	QSB-NET-Large(SE2)	71.2	0.42	16.9	0.69

378 Table 2: Comparison with existing BCNNs on ImageNet-1K. The names of the baseline FP32 models are in 379 parentheses.

¹ Instead of DS convolutions in MobileNetV1 (Howard et al., 2017), MobiNet-Mid (Phan et al., 2020a) and ReActNetA (Liu et al., 2020) adopted grouped convolutions and binarized 3×3 convolutions, respectively.

² Both the baseline and teacher models are ResNet50 (He et al., 2016). In PokeBNN (Zhang et al., 2022), the first convolution

adopted a stride of 4, which made PokeBNN have small FLOPs. We did not consider n-bit quantized operations in this comparison, counting n-bit quantized operations as FLOPs.

³ While term SE2 denotes that self-attention (SE) blocks are deployed after all DS convolutions, term SE1 indicates that SE blocks are deployed only after the DS convolutions during downsampling

APPLICABILITY OF TECHNIQUES FOR BETTER MODELS 4.5

The applicability of the techniques for achieving enhanced performance can be considered with 402 small additional computations. Whereas DSCONVs during downsampling reduce the computa-403 tional costs by removing channelwise operations, existing BCNNs (Phan et al., 2020a; Liu et al., 404 2020) adopted the channelwise spatial convolutions. However, the channelwise operations signifi-405 cantly increase computational costs, so it could be better to deploy the self-attention (SE) block (Hu 406 et al., 2018) after the DSCONVs, applying the attention mechanism to different channels with small 407 computational costs. Because the number of channels is 2048 in the FC layer, its parameters are 408 doubled compared with MobileNetV1 (Howard et al., 2017) and ReActNetA (Liu et al., 2020). Any 409 fixed-point format in the FC layer can reduce storage costs, so we performed an experiment to know 410 whether the FC layer with 8-bit quantized weights is applicable, which is explained in Ablations 411 Studies. Besides, the FC layer with 8-bit quantized weights can be implemented with 8 binarized 412 convolutions, which is described in Appendix A.2. The usage of binarized convolutions in the FC layer reduces storage costs of the FC layer by $\times 4$ without using 8-bit fixed-point operations. 413

414 415

416 417

418

380 38

396

397

399 400

401

5 EXPERIMENTAL RESULTS AND ANALYSIS

5.1 IMAGE CLASSIFICATION AND LATENCY ON REAL HARDWARE

419 We experimented with the proposed models on ImageNet-1K. Like ReActNetA (Liu et al., 420 2020), experiments adopted the two-stage teacher-student training using pretrained ResNet34 as 421 a teacher (Hinton et al., 2015). In the first stage with 256 epochs, whereas the input features for 422 BCONVs were binarized, weights were FP32 values. In the second stage, the pretrained weights 423 from the first stage were used in the initialization. Both input features and weights for BCONVs were binarized during 256 training epochs. The detailed training process is described in Appendix A.4. 424

425 Table 2 shows the comparison with several mobile-friendly CNNs and BCNNs on ImageNet-1K, 426 where a dash (-) denotes that the value was not reported in the reference. The FP32 models and their 427 accuracies are listed above the first midline. The proposed models can significantly reduce OPs 428 compared with ShuffleNetV2 (Ma et al., 2018) and GhostNet (Han et al., 2020). Mostly, the proposed models outperformed existing BCNNs above the second midline in terms of OPs. Although 429 MobiNet-Mid (Phan et al., 2020a) had small OPs, Top-1 accuracy was only 54.4%. ReActNetA (Liu 430 et al., 2020) needed additional OPs due to the usage of binarized 3×3 convolutions instead of DS 431 convolutions. PokeBNN $\times 0.75$ (Zhang et al., 2022) only had 0.42 $\times 10^8$ OPs. Unlike the proposed

432 Table 3: Comparison of parameters and latency for ImageNet- Table 4: Comparison of semantic segmentations 433 1K images.

on PASCAL VOC dataset.

Model	Parameters (Mbyets)	$\begin{array}{c} \text{OPs} \\ (\times 10^8) \end{array}$	Latency(ms) RPi 4B	Latency(ms Exynos
MobileNetV1	16.9	5.75	160.8	34
MobileNetV2	14.0	3.0	117.4	23
XNOR-Net(ResNet18)	4.2	1.67	87.0	21
Bi-RealNet(ResNet18)	4.2	1.63	80.2	19
Real-to-Bin(ResNet18)	5.4	1.83	100.9	-
ReActNetA	7.7	0.87	120.4	30
QuickNet-Small	4.0	-	17.5	9
QuickNet	4.2	-	27.4	13
QB-Net-Small	9.8(3.5) ²	0.43	55.5	14
QB-Net-Large	12.0(5.8) ²	0.53	65.5	14
QSB-Net-Small	10.0(3.9) ²	0.53	76.1	18
OSB-Net-Large	12.3(6.2) ²	0.62	86.2	21
OSB-Not-Large(SE1)	12 4(6 3) ²	0.69	89.7	22

¹ In Bannink et al. (2021), QuickNet-Small and QuickNet had 59.4% and 63.3% Top-1 accuracies on ImageNet-1K.

² The value in parentheses denotes the storage costs when 8-bit quantization was applied to the FC layer. 8 binarized convolutions can replace the 8-bit quantized FC layer. In the ablation studies, the quantization showed no significant

and other existing models, PokeBNN adopted a stride of 4 in the first FP32 convolutional layer, 453 which was the main reason for having the small OPs. 454

455 Table 3 summarizes the comparison in terms of parameters and latency using Larg Compute Engine 456 (LCE) (Bannink et al., 2021) on a single thread of RPi 4B and a Samsung Exynos-9820 processor. 457 We note that the supported layers in LCE were limited, so only several mobile-friendly CNNs and BCNNs based on ResNet18 and MobileNetV1 were compared in Table 3. All proposed models 458 were faster than FP32 models in Table 3. QB-Net showed good efficiency in terms of Top-1 accu-459 racy and latency. For example, QB-Net-Large can achieve 69.8% Top-1 accuracy on ImageNet-1K, 460 having 0.53×10^8 OPs and 65.5 ms latency on the RPi 4B. Although QSB-Net-Large can enhance 461 Top-1 accuracy by 0.8%, its latency increased by 20.7 ms. Because QSB-Net-Large(SE1) only 462 adopted SE blocks during downsampling, the increasing latency was small. Compared with Re-463 ActNetA (Liu et al., 2020), the proposed Large models provided better performances, having faster 464 inference speed. QuickNet-Small showed fast inference speed because QuickNet models were op-465 timized considering the mechanism of LCE (Bannink et al., 2021). However, its performance was 466 only 59.4% Top-1 accuracy on ImageNet-1K.

On the other hand, the experiments with the Samsung Exynos-9820 processor adopted TensorFlow 468 Lite interpreter via Android app. The latencies of the models in Table 3 were measured by averaging 469 the results of 300 runs. In the experiments, all proposed models were faster than the listed FP32 470 models. Furthermore, the proposed models achieved significant speedups compared to the baseline 471 ReActNetA. Besides, QSB-Net-Large(SE1) provided the highest Top-1 accuracy of 71.1%, having 472 22 ms inference latency. As shown in Table 3, the main weakness of the proposed models is the increasing storage costs for the final FC layer. The final FC layer consumed about 8 Mbytes storage, 473 474 which was more than $\frac{2}{3}$ storage costs of the proposed models. As explained in subsection 4.5, the quantization of the FC layer can mitigate the problem, so the values in parentheses in Table 3 prove 475 the reduction of storage costs using the quantization of the FC layer. The effects of quantization on 476 the FC layer in terms of model performance will be discussed in Ablation Studies. 477

478

467

479 5.2 SEMANTIC SEGMENTATION 480

481 We evaluated the proposed models on the PASCAL VOC 2012 dataset (Everingham et al., 2010) 482 for semantic segmentation. The detailed information about the dataset is in Appendix A.10. We 483 measured mIoU on 20 object classes and 1 background class, following the training recipe described in DeepLabv3 (Chen et al., 2017). ReActNetA (Liu et al., 2020) had 0.5 mIoU lower result 484 compared with QB-Net-Small. Besides, QSB-Net-Large was the best-performing binarized seg-485 mentation model, outperforming FP32 ResNet18 (He et al., 2016) and CBNN (Zhou et al., 2023).

447

448

449

450 451 452

Model	W/F ¹	mIoU	Model	W/F	mIoU
ResNet18	32/32	64.9	LQ-Net	3/3	62.5
GroupNet	1/1	60.5	GroupNet+BPAC ²	1/1	65.1
CBNN(Sum) ³	1/1	66.2	CBNN(Cat) ³	1/1	66.5
ReActNetA	1/1	61.8			
OR-Net-Small	1/1	62.3	OSB-Net-Small	1/1	66.2

OB-Net -Net-Small 1/1 68.1 QSB-Net-Large 1/169.2 **OB-Net-Large**

¹ Terms W and F indicate *n*-bit quantization of weights and features

² Term BPAC represents binary parallel atrous convolution. ³ Terms Sum and Cat denote the summation and concatenation

accuracy drop.

Furthermore, the QSB-Net-Small and Large models significantly enhanced performance compared with LQ-Net (Zhang et al., 2018) using 3-bit quantization.

489 5.3 ABLATION STUDIES

Deployments of 1×1 **convolutions:** In a modification of QB-Net-Large, when the number of input 491 492 channels was 16 or 32, binarized pointwise convolutions were adopted after DS convolutions. In this case, the Top-1 accuracy on ImageNet-1K was 68.0%, which was significantly degraded by 493 1.8%. In another modification of QB-Net-Large, when the number of input channels was 16, 32, 494 or 128, binarized 3×3 convolutions were adopted after DS convolutions. In this case, its Top-1 495 accuracy was 70.0%, which was only increased by 0.2%. Based on the above empirical findings, 496 we concluded that when the number of channels in shallow blocks was small, the accuracy can be 497 significantly degraded. Therefore, we determined that instead of binarized pointwise convolutions, 498 binarized 3×3 convolutions were used when the number of input channels was 16 or 32. 499

Effects of $\times 8$ **channel expansion in deep blocks:** A modification of QB-Net-Small changed the number of output channels of the last 13-th and 14-th blocks to 4048 (C = 4048). The modification produced 69.5% Top-1 accuracy, having 2.6% performance enhancement. This study shows that the increasing complexity of deep blocks can significantly enhance model performance. However, we thought that the storage costs of the FC layer and computational costs of binarized pointwise convolutions were significant, so we did not consider the case with C = 4048.

Learnable bias deployed before DS convolutions: In a modification of QB-Net-Large, when the
 learnable bias layer was not deployed just before DS convolutions, Top-1 accuracy was degraded by
 2.1% on ImageNet-1K. The above result indicates that the calibration using the learnable bias can
 be effective in the DS convolutions of the proposed models.

Quantization of DS convolutions and FC layer: In a modification of QSB-Net-Large(SE2), we applied 8-bit quantization to the weights and input features for DS convolutions in each block. In the two-stage training, its training result surprisingly achieved up to 71.5% Top-1 accuracy on ImageNet-1K. Then, we applied 8-bit quantization to the weights and input features for the DS convolution and FC layers, where the tuning adopted the learning rate of 10^{-4} and 25 epochs on the pretrained model. In the evaluation, Top-1 accuracy was 70.8%, demonstrating that 8-bit quantization was applicable to the proposed models.

Training from scratch without teacher: To know the model performance without using the teacher-student training, the proposed models were trained from scratch with the scheduler of co-sine annealing with warmup during 600 epochs, which was based on the recipe of QuickNet (Bannik et al., 2021). QB-Net-Large, QSB-Net-Large, and QSB-Net-Large(SE1) had 65.8%, 67.0%, and 67.5% Top-1 accuracies on ImageNet-1K, which indicates that the teacher-student training was critical for enhancing model performance. Whereas the proposed models were slower than Quick-Net (Bannink et al., 2021), Top-1 accuracies of the proposed models without using the teacher-student training significantly outperformed 63.3% Top-1 accuracy of QuickNet in Table 3.

524 525 526

6 CONCLUSION

527 This paper proposes new BCNNs having low computational costs and high performance using chan-528 nel quadrupling and smooth downsampling. The proposed structure using the channel quadrupling 529 has low-cost computations in shallow blocks and a wider dynamic range in the frequency domain 530 due to the increased model complexity in deep blocks. The techniques for achieving better models, 531 such as the self-attention (Hu et al., 2018) and the quantization of DS convolutional and FC layers, 532 are applicable to the proposed models. Using multiple binarized convolutions in the FC layer, the 533 proposed models can significantly reduce storage costs without using an 8-bit quantized FC layer. 534 When the self-attention is applied to the DS convolutions during downsampling, its performance on 535 ImageNet-1K reaches up to 71.1% Top-1 accuracy, only requiring 0.63×10^8 OPs. The small laten-536 cies on an RPi 4B and a Samsung Exynos processor prove that the proposed models are suitable for implementing high-speed inference on real hardware. Considering the above performance enhance-537 ments and structural benefits, it is concluded that the proposed models are efficient for mobile-based 538 applications. Additional explanations and visualizations of the proposed models are included in Appendixes A.8, A.9, and A.10.

540 ETHIC STATEMENT

This paper proposes new binarized convolutional neural network (BCNN) models, introducing novel architectures to enhance model performance and inference speed on real hardware. In this study, we have kept the following ethical principles of ICLR 2025 as:

- 1. **Contribute to Society and Human Well-being:** The aim of this study is to develop efficient BCNNs by reducing computational costs in shallow blocks and increasing representation capacity in deep blocks. The development of the proposed mobile-friendly models has the potential to benefit a wide range of societal applications, including healthcare, education, and environmental monitoring, particularly in resource-constrained applications.
- 2. Uphold High Standards of Scientific Excellence: We have intensively performed experiments to validate our proposed models. The motivations, ideas, and conclusions are presented to contribute to the scientific community.
 - 3. **Avoid Harm:** The study does not include any human subjects or sensitive personal data. We strongly discourage any misuse of our work that could harm individuals, although there is no explicit information about the misuse in the manuscript.
- 4. **Be Honest, Trustworthy, and Transparent:** We have honestly reported our research findings, including both strengths and limitations. We sincerely reported the weak points of the proposed models and the method to achieve better models. All data sources, model structures, and experimental environments are fully disclosed to ensure transparency.
 - 5. **Be Fair and Take Action to Avoid Discrimination:** Because we adopt public datasets such as ImageNet-1K and public Pytorch library, experiments can be fair without any discrimination.
 - 6. **Respect the Work Required to Produce New Ideas and Artefacts:** We cite all relevant references in the manuscript to respect existing works. This paper is written considering the previous works and knowledge.
 - 7. **Respect Privacy:** The datasets adopted in our experiments, such as ImageNet-1K and PASCAL VOC, are publicly available and do not contain critical personal information.
 - 8. Honour Confidentiality: This paper does not have any confidentiality agreements.

Reproducibility Statement

We adopt conventional ImageNet-1K and PASCAL VOC datasets for easy reproduction. The attached code as supplementary materials can run when Pytorch dataset formats are prepared. The detailed model structures are described in the main body of this paper and Appendix A.1. Detailed explanations of experimental environments and training processes are included in Appendix A.4. Besides, the environments for evaluating inference speed on real hardware are described in Appendix A.7.

594	REFERENCES
595	

596	Tensorflow lite. https://www.tensorflow.org/lite, 2023. Accessed: 2023-09-20.
597	Tom Bannink Adam Hillier Lukas Geiger Tim de Bruin Leon Overweel Jelmer Neeven and
598	Koen Helwegen. Larg compute engine: Design, benchmark and deploy state-of-the-art binarized
599	neural networks. Proceedings of Machine Learning and Systems, 3:680–695, 2021.
600	
601	Yoshua Bengio, Nicholas Léonard, and Aaron Courville. Estimating or propagating gradients
602	through stochastic neurons for conditional computation. arXiv preprint arXiv:1308.3432, 2013.
603	Adrian Bulat Georgios Tzimiropoulos and Samsung AI Center Xnor-net++: Improved hinary
604	neural networks. In <i>The British Machine Vision Conf.</i> 2019.
605	
606	Holger Caesar, Jasper Uijlings, and Vittorio Ferrari. Coco-stuff: Thing and stuff classes in context.
609	In Computer vision and pattern recognition (CVPR), 2018 IEEE conference on. IEEE, 2018.
600	Bo Chen, Golnaz Ghiasi, Hanviao Liu, Tsung Vi Lin, Dmitry Kalenichenko, Hartwig Adam, and
610	Ouoc V Le, Mnasfpn: Learning latency-aware pyramid architecture for object detection on mobile
611	devices. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition.
612	pp. 13607–13616, 2020a.
613	
614	Liang-Chieh Chen, George Papandreou, Florian Schroff, and Hartwig Adam. Rethinking atrous
615	convolution for semantic image segmentation. ArXiv, abs/1706.05587, 2017. URL https:
616	//api.semanticscholar.org/CorpusiD:22655199.
617	Tse-Wen Chen, Pangfeng Liu, and Jan-Jan Wu. Exploiting data entropy for neural network com-
618	pression. In 2020 IEEE International Conference on Big Data (Big Data), pp. 5007–5016. IEEE,
619	2020b.
620	
621	Matthieu Courbariaux, Itay Hubara, Daniel Soudry, Ran El-Yaniv, and Yoshua Bengio. Binarized
622	or 1 arXiv preprint arXiv: 1602.02830, 2016
623	01-1. <i>urxiv preprint urxiv</i> .1002.02050, 2010.
624	M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The pascal visual object
625	classes (voc) challenge. International Journal of Computer Vision, 88(2):303–338, June 2010.
627	LiuS GeZ, F WANG, et al. Yolox: Exceedingyoloseriesin2021. arXiv preprint arXiv:2107.08430,
628	2021.
629	D Circhick Fost a opp on Viv proprint on Viv 1504 08082 2015
630	K Onsnick. Fast I-chin. <i>urxiv preprint urxiv.1504.08085</i> , 2015.
631	Kai Han, Yunhe Wang, Qi Tian, Jianyuan Guo, Chunjing Xu, and Chang Xu. Ghostnet: More
632	features from cheap operations. In Proceedings of the IEEE/CVF conference on computer vision
633	and pattern recognition, pp. 1580–1589, 2020.
634	Kaiming Ha Viangun Zhang, Shaaging Dan, and Jian Sun. Dean racidual laarning for imaga racag
635	nition In Proceedings of the IFFF conference on computer vision and pattern recognition pr
636	770–778 2016
637	
638	Koen Helwegen, James Widdicombe, Lukas Geiger, Zechun Liu, Kwang-Ting Cheng, and Roeland
639	Nusselder. Latent weights do not exist: Rethinking binarized neural network optimization. In
640	Advances in neural information processing systems, pp. 7531–7542, 2019.
641	Geoffrey Hinton Oriol Vinyals Jeff Dean et al Distilling the knowledge in a neural network arXiv
642	preprint arXiv:1503.02531, 2(7), 2015.
643	
644	Lilian Hollard, Lucas Mohimont, Nathalie Gaveau, and Luiz-Angelo Steffenel. Leyolo, new scalable
645	and efficient cnn architecture for object detection. arXiv preprint arXiv:2406.14239, 2024.
040 647	Yuanduo Hong Huihui Pan Weichao Sun and Yisong Jia Deen dual-resolution networks for real-
047	time and accurate semantic segmentation of road scenes. <i>arXiv preprint arXiv:2101.06085</i> , 2021.

- Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, 649 Marco Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural networks for 650 mobile vision applications. arXiv preprint arXiv:1704.04861, 2017. 651 Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In *Proceedings of the IEEE* 652 conference on computer vision and pattern recognition, pp. 7132–7141, 2018. 653 654 Gao Huang, Zhuang Liu, Geoff Pleiss, Laurens Van Der Maaten, and Kilian Weinberger. Convo-655 lutional networks with dense connectivity. IEEE Transactions on Pattern Analysis and Machine Intelligence, 2019. 656 657 Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. 658 CoRR, abs/1412.6980, 2014. URL https://api.semanticscholar.org/CorpusID: 659 6628106. 660 Xiaofan Lin, Cong Zhao, and Wei Pan. Towards accurate binary convolutional neural network. In 661 Advances in Neural Information Processing Systems, pp. 345–353, 2017. 662 663 Chunlei Liu, Peng Chen, Bohan Zhuang, Chunhua Shen, Baochang Zhang, and Wenrui Ding. Sa-664 bnn state-aware binary neural network. In Proceedings of the AAAI Conference on Artificial 665 Intelligence, volume 35, pp. 2091–2099, 2021. 666 Chunlei Liu, Wenrui Ding, Peng Chen, Bohan Zhuang, Yufeng Wang, Yang Zhao, Baochang Zhang, 667 and Yuqi Han. Rb-net: Training highly accurate and efficient binary neural networks with re-668 shaped point-wise convolution and balanced activation. IEEE Transactions on Circuits and Sys-669 tems for Video Technology, 32(9):6414-6424, 2022. 670 671 Zechun Liu, Baoyuan Wu, Wenhan Luo, Xin Yang, Wei Liu, and Kwang-Ting Cheng. Bi-real net: Enhancing the performance of 1-bit cnns with improved representational capability and advanced 672 training algorithm. In Proceedings of the European conference on computer vision (ECCV), pp. 673 722-737, 2018. 674 675 Zechun Liu, Zhiqiang Shen, Marios Savvides, and Kwang-Ting Cheng. Reactnet: Towards pre-676 cise binary neural network with generalized activation functions. In European Conference on 677 Computer Vision (ECCV), pp. 143–159, 2020. 678 Ningning Ma, Xiangyu Zhang, Hai-Tao Zheng, and Jian Sun. Shufflenet v2: Practical guidelines for 679 efficient cnn architecture design. In Proceedings of the European conference on computer vision 680 (ECCV), pp. 116–131, 2018. 681 682 Brais Martinez, Jing Yang, Adrian Bulat, and Georgios Tzimiropoulos. Training binary neural 683 networks with real-to-binary convolutions. In International Conference on Learning Representations, 2019. 684 685 Sachin Mehta, Mohammad Rastegari, Linda Shapiro, and Hannaneh Hajishirzi. Espnetv2: A light-686 weight, power efficient, and general purpose convolutional neural network. In Proceedings of the 687 IEEE/CVF conference on computer vision and pattern recognition, pp. 9190–9200, 2019. 688 Hai Phan, Yihui He, Marios Savvides, Zhiqiang Shen, et al. Mobinet: A mobile binary network for 689 image classification. In The IEEE Winter Conference on Applications of Computer Vision, pp. 690 3453-3462, 2020a. 691 692 Hai Phan, Zechun Liu, Dang Huynh, Marios Savvides, Kwang-Ting Cheng, and Zhiqiang Shen. 693 Binarizing mobilenet via evolution-based searching. In Proceedings of the IEEE/CVF conference 694 on computer vision and pattern recognition, pp. 13420–13429, 2020b. Haotong Qin, Ruihao Gong, Xianglong Liu, Mingzhu Shen, Ziran Wei, Fengwei Yu, and Jingkuan 696 Song. Forward and backward information retention for accurate binary neural networks. In 697 Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pp. 2250– 2259, 2020. 699 Mohammad Rastegari, Vicente Ordonez, Joseph Redmon, and Ali Farhadi. Xnor-net: Imagenet 700
- classification using binary convolutional neural networks. In *European conference on computer vision*, pp. 525–542. Springer, 2016.

- Tal Ridnik, Hussam Lawen, Asaf Noy, Emanuel Ben Baruch, Gilad Sharir, and Itamar Friedman. Tresnet: High performance gpu-dedicated architecture. In *proceedings of the IEEE/CVF winter conference on applications of computer vision*, pp. 1400–1409, 2021.
- Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)*, 115(3):211–252, 2015. doi: 10.1007/s11263-015-0816-y.
- Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mo bilenetv2: Inverted residuals and linear bottlenecks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 4510–4520, 2018.
- Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural net works. In *International conference on machine learning*, pp. 6105–6114. PMLR, 2019.
- Zhijun Tu, Xinghao Chen, Pengju Ren, and Yunhe Wang. Adabin: Improving binary neural networks with adaptive binary sets. In *European conference on computer vision*, pp. 379–395. Springer, 2022.
- Weitao Wan, Jiansheng Chen, Tianpeng Li, Yiqing Huang, Jingqi Tian, Cheng Yu, and Youze Xue.
 Information entropy based feature pooling for convolutional neural networks. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 3405–3414, 2019.
- Jian Wang, Chenhui Gou, Qiman Wu, Haocheng Feng, Junyu Han, Errui Ding, and Jingdong Wang.
 Rtformer: Efficient design for real-time semantic segmentation with transformer. *Advances in Neural Information Processing Systems*, 35:7423–7436, 2022.
- Peisong Wang, Xiangyu He, Gang Li, Tianli Zhao, and Jian Cheng. Sparsity-inducing binarized neural networks. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pp. 12192–12199, 2020.
- Lijun Wu, Xu Lin, Zhicong Chen, Jingchang Huang, Huawei Liu, and Yang Yang. An efficient binary convolutional neural network with numerous skip connections for fog computing. *IEEE Internet of Things Journal*, 8(14):11357–11367, 2021.
- Sheng Xu, Yanjing Li, Teli Ma, Mingbao Lin, Hao Dong, Baochang Zhang, Peng Gao, and Jinhu Lu.
 Resilient binary neural network. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pp. 10620–10628, 2023.
- Zihan Xu, Mingbao Lin, Jianzhuang Liu, Jie Chen, Ling Shao, Yue Gao, Yonghong Tian, and
 Rongrong Ji. Recu: Reviving the dead weights in binary neural networks. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 5198–5208, 2021.
- Changqian Yu, Changxin Gao, Jingbo Wang, Gang Yu, Chunhua Shen, and Nong Sang. Bisenet
 v2: Bilateral network with guided aggregation for real-time semantic segmentation. *International Journal of Computer Vision*, 129:3051–3068, 2021.
- Dongqing Zhang, Jiaolong Yang, Dongqiangzi Ye, and Gang Hua. Lq-nets: Learned quantization
 for highly accurate and compact deep neural networks. In *European Conference on Computer Vision (ECCV)*, 2018.
- Yichi Zhang, Zhiru Zhang, and Lukasz Lew. Pokebnn: A binary pursuit of lightweight accuracy. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 12475–12485, 2022.
- Hengshuang Zhao, Jianping Shi, Xiaojuan Qi, Xiaogang Wang, and Jiaya Jia. Pyramid scene parsing network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 2881–2890, 2017.
- Yang Zhao and Hao Zhang. Quantitative performance assessment of cnn units via topological entropy calculation. *arXiv preprint arXiv:2103.09716*, 2021.
- Xichuan Zhou, Rui Ding, Yuxiao Wang, Wenjia Wei, and Haijun Liu. Cellular binary neural network
 for accurate image classification and semantic segmentation. *IEEE Transactions on Multimedia*, 25:8064–8075, 2023. doi: 10.1109/TMM.2022.3233255.

A APPENDIX / SUPPLEMENTAL MATERIAL

A.1 MODEL STRUCTURES OF QB-NET-LARGE AND QSB-NET-LARGE

Table 5: Model structures of QB-Net-Large and QSB-Net-Large for ImageNet-1K.

Index		QB-Net-Lar	ge	QSB-Net-Large				
	Block	Input	Output	Block	Input	Output		
1	Conv	$3 \times 224 \times 224$	$16 \times 112 \times 112$	Conv	$3 \times 224 \times 224$	$ 16 \times 112 \times 112 $		
2	TypeN	$16 \times 112 \times 112$	$16\!\times\!112\!\times\!112$	TypeN	$16 \times 112 \times 112$	$16 \times 112 \times 112$		
3	TypeD	$16 \times 112 \times 112$	$32 \times 56 \times 56$	TypeDS	$16 \times 112 \times 112$	$32 \times 56 \times 56$		
4	TypeN	$32 \times 56 \times 56$	$32 \times 56 \times 56$	TypeN	$32 \times 56 \times 56$	$32 \times 56 \times 56$		
5	TypeQ	$32 \times 56 \times 56$	$128 \times 28 \times 28$	TypeQS	$32 \times 56 \times 56$	$128 \times 28 \times 28$		
6	TypeN	$128 \times 28 \times 28$	$128\!\times\!28\!\times\!28$	TypeN	$128 \times 28 \times 28$	$128 \times 28 \times 28$		
7	TypeQ	$128 \times 28 \times 28$	$512 \times 14 \times 14$	TypeQS	$128 \times 28 \times 28$	$512 \times 14 \times 14$		
8	TypeN	$512 \times 14 \times 14$	$512 \times 14 \times 14$	TypeN	$512 \times 14 \times 14$	$512 \times 14 \times 14$		
9	TypeQ	$512 \times 14 \times 14$	$2048 \times 7 \times 7$	TypeQS	$512 \times 14 \times 14$	$2048 \times 7 \times 7$		
10	TypeN	$2048 \times 7 \times 7$	$2048 \times 7 \times 7$	TypeN	$2048 \times 7 \times 7$	$2048 \times 7 \times 7$		
11	TypeN	$2048 \times 7 \times 7$	$2048 \times 7 \times 7$	TypeN	$2048 \times 7 \times 7$	$2048 \times 7 \times 7$		
12	TypeN	$2048 \times 7 \times 7$	$2048 \times 7 \times 7$	TypeN	$2048 \times 7 \times 7$	$2048 \times 7 \times 7$		
13	TypeN	$2048 \times 7 \times 7$	$2048 \times 7 \times 7$	TypeN	$2048 \times 7 \times 7$	$2048 \times 7 \times 7$		
14	TypeN	$2048 \times 7 \times 7$	$2048 \times 7 \times 7$	TypeN	$2048 \times 7 \times 7$	$2048 \times 7 \times 7$		
15	AvgPool	$2048 \times 7 \times 7$	$2048 \times 1 \times 1$	AvgPool	$2048 \times 7 \times 7$	$2048 \times 1 \times 1$		
16	FC	$2048 \times 1 \times 1$	1000	FC	$2048 \times 1 \times 1$	1000		

A.2 BINARIZED CONVOLUTIONS IN FINAL FC LAYER

Let us assume that x_i denotes a feature in input channel *i*. A *n*-bit quantized weight w_i for input channel *i* is denoted as $w_i = w_i^{n-1}2^{n-1} + w_i^{n-2}2^{n-2} + \cdots + w_i^{1}2^1 + w_i^{0}2^0$. When $w_i^7, w_i^6, \cdots, w_i^1, w_i^0 \in \{-1, +1\}, w_i \in \{-255, -253, \cdots, +253, +255\}$, showing the uniform interval between quantized weights of w_i . When we assume n = 8 for 8-bit quantized operations, a multiply operation is formulated as $x_i \cdot w_i = w_i^7 x_i 2^7 + w_i^6 x_i 2^6 + \cdots + w_i^1 x_i 2^1 + w_i^0 x_i 2^0$.

787 When the number of input channels is C, an output feature x_{out} of a binarized pointwise convolution can be calculated as:

$$x_{out} = scale \cdot \sum_{i=1}^{C} (w_i^7 x_i 2^7 + w_i^6 x_i 2^6 + \dots + w_i^1 x_i 2^1 + w_i^0 x_i 2^0),$$
(2)

where the index of an output channel is not shown for clarity. Term *scale* is the scaling factor for a binarized convolution. Eq.(2) can be represented as:

$$x_{out} = scale \cdot \left(2^7 \sum_{\substack{i=1\\1\times1 \text{ BCONV}}}^C w_i^7 x_i + 2^6 \sum_{\substack{i=1\\1\times1 \text{ BCONV}}}^C w_i^6 x_i + \dots + 2^1 \sum_{\substack{i=1\\1\times1 \text{ BCONV}}}^C w_i^1 x_i + 2^0 \sum_{\substack{i=1\\1\times1 \text{ BCONV}}}^C w_i^0 x_i\right).$$
(3)

While the binarized convolution in Rastegari et al. (2016); Bulat et al. (2019) is formulated as $x_{out} =$ $scale \cdot \sum_{i=1}^{C} w_i^0 x_i^0, x_i \in \{-1, +1\}$, 8 binarized convolutions denoted as 1×1 BCONV are required in Eq.(3). Considering Eq.(2), the FC layer with 8-bit quantized weights can be decomposed into 8 binarized convolutions. The change of the latency using the decomposed binarized convolutions was negligible on real hardware because the amount of computations in the FC layer was very small. However, the costs for storing the parameters for the FC layer can be significantly reduced. In the proposed QSB-Net, the number of input channels for the FC layer was 2048. When ImageNet-1K was used as the dataset, the FP32 FC layer consumed 8.19 Mbytes in the proposed models. However, when 8 binarized convolutions were deployed in the FC layer, the storage costs were only 2.05 Mbytes, having $\times 4$ memory efficiency. When the mixed precision using both FP32 and 8-bit fixed-point operations was not available in the hardware platform for BCNNs, the decomposed binarized convolutions in the FC layer are effective in reducing the total storage costs.

822

823

824 825

826 827

828

829

830

834 835 836



Figure 6: Visualized entropies of block outputs from ReActNetA and proposed models.

A.3 DETAILED PRESENTATION OF ENTROPY IN PROPOSED MODELS

We think that the downsampling of feature maps and binarization error of BCNNs could affect the information loss. Entropy can be exploited to design the behavior of a specific layer (Wan et al., 2019; Chen et al., 2020b; Zhao & Zhang, 2021). The above existing works motivate us to believe that entropy can be used to explain the information characteristics of output features in BCNNs.

The formula for entropy H(X) is specifically for discrete variable $x_i \in X = \{x_1, x_2, \dots, x_{n-1}, x_n\}$, which is expressed as follows:

$$H(X) = -\sum_{i=1}^{n} p(x_i) \log_2 p(x_i),$$
(4)

where $p(x_i)$ is the probability of each possible value $x_i \in X$. With Eq.(4), the entropy for a block was calculated as: firstly, because the values were continuous, the values of output feature maps in a layer were discretized by binning the values using a histogram. The probability distribution was determined by the relative frequency of each bin from the histogram. Then, the above probability distribution was used to calculate the entropy using the discrete entropy formula in Eq.(4).

We adopted five images from the validation dataset of ImageNet-1K. The images were used as the 842 inputs for the pretrained models. In the inference, the output features were hooked and used to cal-843 culate the entropy of each block. Figure 6 (a) illustrates that the entropy of each block significantly 844 fluctuated in baselined ReActNetA, QB-Net-Small, and QSB-Net-Small. On the other hand, Fig-845 ure 6 (b) shows that the fluctuation of entropy is reduced. We think that the above trends show that 846 feature distributions can change smoothly during feedforwarding. Besides, it is concluded that the 847 increasing representation capacity in deep blocks of QB-Net-Large and QSB-Net-Large can make 848 the change of feature distributions more stable. 849

A.4 DETAILED DESCRIPTION OF TRAINING PROCESS IN IMAGE CLASSIFICATION

ImageNet-1K contains 1.3M training and 50K validation images with 1,000 classes. During training on the ImageNet-1K dataset, 224×224 images based on the augmentations in Liu et al. (2020) were adopted. In inference, 224×224 center-cropped images from the validation dataset were adopted without any specific image augmentation.

856 Like other BCNN models, the first convolutional and final FC layers adopted FP32 weights and 857 activations. In order to reduce the storage costs of the final FC layer, 8-bit quantization can be 858 applied. For an apple-to-apple comparison, we adopted ADAM (Kingma & Ba, 2014) optimizer 859 in all cases, having $\beta_1 = 0.9$ and $\beta_2 = 0.999$. When training during E_{epochs} epochs, the initial learning rate $lr_{base} = 0.001$ was assigned. During training, the learning rate lr in the e_{epochs} -th epoch decreased based on **poly** policy, which limited the maximum learning rate of the ADAM 861 optimizer (Kingma & Ba, 2014) by $lr_{base} \times (1 - e_{epochs}/E_{epochs})$. In the two-stage training, a 862 teacher-student training method (Hinton et al., 2015) was adopted using the pretrained ResNet34 (He 863 et al., 2016) from Pytorch official site as a teacher.

In order to know the effects of structural benefits, the training recipe of ReActNetA (Liu et al., 2020), which was a well-known BCNN model motivated by MobileNetV1 (Phan et al., 2020a), was adopted. By employing the same training recipe, our analysis demonstrates that the proposed model structures achieve comparable or enhanced performance with reduced computational costs. Based on the above apple-to-apple comparison, we mainly focused on the demonstration of enhancements by adapting model structures. In the first stage with 256 epochs, whereas the input features for BCONVs were binarized, weights were FP32 values. The weight decay in the first stage was set as 10^{-5} . In the second stage, the pretrained weights from the first stage were used in the initialization. Both input features and weights for BCONVs were binarized during 256 training epochs. The weight decay in the second stage was set as zero. All experiments were conducted on a machine having an AMD Ryzen Threadripper PRO 5955WX 16-Core CPU, 2 NVIDIA RTX 4090 GPUs, and 256 GB RAM. Although the exact training time depended on the status of computing resource usage, the total training times of QSB-Net-Large and QB-Net-Small on our machine were up to about 6 and 7.5 days, respectively. For the approximation of the gradient of sign function

In order to know the effects of the teacher-student training, the proposed models were trained from scratch with the scheduler of cosine annealing with warmup during 600 epochs. In this experiments, we adopted ADAM (Kingma & Ba, 2014) optimizer, having $\beta_1 = 0.9$ and $\beta_2 = 0.999$. The initial learning rate of $lr_{base} = 0.001$ was assigned, having 5 warmup steps with 0.001 maximum learning rate and zero minimum learning rate. After performing warmup steps, the learning rate decreased. The total training times of QB-Net-Large and QSB-Net-Large(SE1) in our machine were up to about 8 and 9 days, respectively.

918 A.5 DETAIL DISCUSSION FOR ABLATION STUDIES

Effects of Channel Expansion and Computational Complexity in Shallow and Deep Layers:
 Firstly, by deploying TypeN and TypeD blocks and considering the channel expansion of Mo bileNetV1 and ReActNetA in Table 1, the model performance without channel quadrupling was
 investigated. In the above ablation model denoted as Ablation 1, the classification on ImageNet
 achieved only 64.2% Top-1 accuracy.

In order to know the effects of model complexity in the shallow layer, the following experiments were performed as: Firstly, in order to know the effects of 3×3 binarized convolutions in shallow layers, when the number of output channels 16 and 32, we adopted binarized 1×1 pointwise convolutions in QB-Net-Small instead of 3×3 binarized convolutions. In the above ablation model denoted as Ablation 2, the classification on ImageNet achieved only 66.7% Top-1 accuracy. Com-pared with 66.9% Top-1 accuracy of QB-Net-Small in Table 1, there was only 0.2% Top-1 accuracy drop. Secondly, as shown in Table 6, we evaluated an ablation model by decreasing the number of channels in the shallow layer, which is denoted as Ablation 3. In Ablation 3, the classification on ImageNet-1K achieved 66.3% Top-1 accuracy. The evaluation in Ablations 2 and 3 indicated that the model complexity in the shallow layer could not be critical.

In another ablation study, we increased the model complexity in deep layers as: Firstly, instead of binarized 1×1 pointwise convolutions, binarized 3×3 convolutions were used in a modification of QB-Net-Small when the numbers of output channels were 16, 32, and 128. In the above ablation model denoted as Ablation 4, the classification on ImageNet-1K achieved 67.9% Top-1 accuracy, which outperformed QB-Net-Small by 1%. On the other hand, when the numbers of output channels were 16, 32, and 128, a modification of QB-Net-Large adopted binarized 3×3 convolutions instead of binarized 1×1 pointwise convolutions. In the above ablation model denoted as **Ablation 5**, the classification on ImageNet-1K achieved 70.0% Top-1 accuracy, which outperformed QB-Net-Large only by 0.2%. Moreover, when a modification of QB-Net-Small denoted as Ablation 6 had 4096 output channels in the last two depthwise separable convolutions and binarized pointwise convolu-tions in Table 6, it achieved 69.5% Top-1 accuracy on ImageNet-1K. In the comparison between Ablations 4, 5, and 6, as the computational complexity in the deeper layers increased, it was con-cluded that the increasing computational complexity in shallow layers has no significant impact on the model performance.

a	Δ	2		
9		0		
~		~		
u	21	u.		

Table 6: Model structures from ablation studies on ImageNet-1	Κ.
---	----

Index		Ablation 3	3	Ablation 6				
	Block	Input	Output	Block	Input	Output		
1	Conv	$3 \times 224 \times 224$	$8 \times 112 \times 112$	Conv	$3 \times 224 \times 224$	$ 16 \times 112 \times 112 $		
2	-	$8 \times 112 \times 112$	$16 \times 112 \times 112$	TypeN	$16 \times 112 \times 112$	$16 \times 112 \times 112$		
3	TypeD	$16 \times 112 \times 112$	$32 \times 56 \times 56$	TypeD	$16 \times 112 \times 112$	$32 \times 56 \times 56$		
4	TypeN	$32 \times 56 \times 56$	$32 \times 56 \times 56$	TypeN	$32 \times 56 \times 56$	$32 \times 56 \times 56$		
5	TypeQ	$32 \times 56 \times 56$	$128 \times 28 \times 28$	TypeQ	$32 \times 56 \times 56$	$128 \times 28 \times 28$		
6	TypeN	$128 \times 28 \times 28$	$128 \times 28 \times 28$	TypeN	$128 \times 28 \times 28$	$128 \times 28 \times 28$		
7	TypeQ	$128 \times 28 \times 28$	$512 \times 14 \times 14$	TypeQ	$128 \times 28 \times 28$	$512 \times 14 \times 14$		
8	TypeN	$512 \times 14 \times 14$	$512 \times 14 \times 14$	TypeN	$512 \times 14 \times 14$	$512 \times 14 \times 14$		
9	TypeN	$512 \times 14 \times 14$	$512 \times 14 \times 14$	TypeN	$512 \times 14 \times 14$	$512 \times 14 \times 14$		
10	TypeN	$512 \times 14 \times 14$	$512 \times 14 \times 14$	TypeN	$512 \times 14 \times 14$	$512 \times 14 \times 14$		
11	TypeN	$512 \times 14 \times 14$	$512 \times 14 \times 14$	TypeN	$512 \times 14 \times 14$	$512 \times 14 \times 14$		
12	TypeN	$512 \times 14 \times 14$	$512 \times 14 \times 14$	TypeN	$512 \times 14 \times 14$	$512 \times 14 \times 14$		
13	TypeQ	$512 \times 14 \times 14$	$2048 \times 7 \times 7$	-	$512 \times 14 \times 14$	$4096 \times 7 \times 7$		
14	TypeN	$2048 \times 7 \times 7$	$2048 \times 7 \times 7$	TypeN	$4096 \times 7 \times 7$	$4096 \times 7 \times 7$		
15	AvgPool	$2048 \times 7 \times 7$	$2048 \times 1 \times 1$	AvgPool	$4096 \times 7 \times 7$	$4096 \times 1 \times 1$		
16	FC	$2048 \times 1 \times 1$	1000	FC	$4096 \times 1 \times 1$	1000		

Effects of Smooth Downsampling: As shown in Figure 3, the smooth downsampling was only performed in TypeDS and TypeQS, so that the effects of smooth downsampling can be discussed with the performance gap between QB-Net and QSB-Net. In a modification of Ablation 1 having the channel expansion of MobileNetV1 and ReActNetA, the TypeD block was applied. In the above modification denoted as Ablation 7, Top-1 accuracy was enhanced by 70.5%, which showed that

972 the smooth downsampling is effective in both the existing baseline and proposed models. However, 973 the computational costs were greater than those of ReActNetA due to the additional computation 974 from the smooth downsampling. In the comparison between QB-Net and QSB-Net models, QSB-975 Net-Small and QSB-Net-Large enhanced Top-1 accuracy by 1.9% and 0.8% over QB-Net-Small 976 and QB-Net-Large, respectively. The above experimental results indicated that as the computational complexity in deep layers increased, the effects of the smooth downsampling were limited. The 977 below Table 7 summarizes the performance of the ablation models in the image classifications on 978 ImageNet-1K. 979

979

981

 Table 7: Summary of ablation studies with additional experiments on ImageNet-1K.

To June 1941, 216 and an
1Ablation 1 64.2 Channel expansion from MobileNetV1 and ReActivet without smooth downsampling2Ablation 2 66.7 Binarized 1×1 convolutions for # channels with 16 and 32(QB-Net-Small)
3 Ablation 3 66.3 8 output channels for the first convolution(QB-Net-Small) 4 Ablation 4 67.9 Binarized 3 × 3 convolutions when # output channels are 16 32 and 128 (QB-Net-Small)
5 Ablation 5 70.0 Binarized 3 × 3 convolutions when # output channels are 16, 32, and 128 (QB-Net-Large)
6 Ablation 6 69.5 4096 output channels in the last two DS layers(QB-Net-Small) 7 Ablation 7 70.5 Modification of Ablation 1 with smooth downsampling
8 QB-Net-Small 66.9 Small computational complexity in deep layers without smooth downsampling 9 OSB-Net-Small 68.8 Small computational complexity in deep layers with smooth downsampling
10 QB-Net-Large 69.8 Large computational complexity in deep layers without smooth downsampling 00 CRD N-term 69.8 Large computational complexity in deep layers without smooth downsampling
11 QSB-Net-Large 70.6 Large computational complexity in deep layers with smooth downsampling

A.6 EXPERIMENTAL RESULTS ON COCO-STUFF

To demonstrate the generalization ability of the proposed models, we conducted semantic segmentation experiments on the large-scale dataset COCO-Stuff Caesar et al. (2018). In Table 8, QSB-Net-Large and QB-Net-Large achieved higher performance than the FP32 precision models with a similar number of parameters while using 1-bit weights and features.

Table 8: Comparison of semantic segmentation on COCO-Stuff dataset.

Model	Params	W/F	mIoU	
BiSeNetv2-L (Yu et al., 2021)	-	32/32	28.7	
DDRNet23 (Hong et al., 2021)	20.1M	32/32	32.1	
PSPNet50 (Zhao et al., 2017)	-	32/32	32.6	
RTFormer-B (Wang et al., 2022)	16.8M	32/32	35.3	
QB-Net-Large	12.0M	1/1	36.4	
QSB-Net-Large	12.3M	1/1	37.5	

Notably, QSB-Net-Large outperformed the Transformer-based RTFormer-B (Wang et al. (2022)) by 2.2 mIoU. It can be seen that the proposed model has good generalization performance even on large datasets and has a high potential for practical applications.

Besides, object detection was performed on the COCO-Stuff dataset. Table 9 lists the comparison of object detection, where QSB-Net-Large has achieved higher performance than the FP32 precision lightweight models, using 1-bit weights and features. The results show that for vision tasks such as classification, semantic segmentation, and object detection, the proposed model outperformed the accuracy of FP32 precision models.

Table 9: Comparison of objection detection on COCO-Stuff dataset.

1053	Model	W/F	mAP
1054	Fast-RCNN (Girshick, 2015)	32/32	19.7
1055	LeYOLO-Nano (Hollard et al., 2024)	32/32	25.2
1057	MnasFPN (MobileNetv3) (Chen et al., 2020a)	32/32	25.5
1058	YOLOX-Nano (GeZ et al., 2021)	32/32	25.8
1059	ESPNetv2 (Mehta et al., 2019)	32/32	26.0
1060	QSB-Net-Large	1/1	20.4

Considering the above semantic segmentation and object detection on the CoCo-Stuff dataset, it is concluded that the proposed models can show an important achievement that overcame the limitations of BCNNs and enhanced its applicability in real applications.

1080 A.7 EXPERIMENTAL ENVIRONMENTS ON REAL HARDWARE: LARQ

Target models were prepared using TensorFlow Keras framework. XNOR-Net (Rastegari et al., 2016), Real-to-Bin (Martinez et al., 2019), Bi-RealNet (Liu et al., 2018), and QuickNet (Bannink et al., 2021) Keras models were from Larq Zoo (Bannink et al., 2021). The proposed models and other counterparts were coded based on the original Keras layers. The models can be converted into TFLite (TensorFlow Lite) filebuffer files. When checking the inference speed of a model, we adopted Larq Compute Engine (LCE) (Bannink et al., 2021), which provided a benchmark evalua-tion program based on TensorFlow Lite (ten, 2023) and customized binarized convolutional layers. The benchmark evaluation program ran on Manjaro 64-bit GNOME Desktop for an RPi 4B and An-droid app for a Samsung Exynos-9820 processor. It was known that LCE provided a collection of hand-optimized TFLite custom operators. Along with the full support of existing TFLite operators, each binarized convolutional layer can be performed using its custom binarized convolution. In our evaluations, the program showed the averaged latencies of 50 runs on the RPi 4B and 300 runs on the Samsung Exynos-9820 processor with randomly generated inputs.

1134 A.8 VISUALIZATION OF FEATURE MAPS

1136	Figures 7 and 8 visualize the first five feature maps from the downsampling blocks. The pretrained
1137	model of the baseline ReActNetA (Liu et al., 2020) was downloaded from its official GitHub. Com-
1138	pared with ReActNetA, the visualization shows a significant difference in the first and last downsam-
1139	pling blocks. In QSB-Net models, the output feature maps after performing heightwise downsam-
1140	pling are illustrated. The feature maps are denoted as the output of $stride = (2, 1)$. Compared with
1141	the feature maps after widthwise downsampling, the output feature maps after heightwise downsam-
1142	pling also show the diversity of features after heightwise downsampling, which could indicate that
1143	smooth downsampling in QSB-Net can enhance the representation capacity. In general, the visu-
1144	representation capacity of deep blocks during downsampling. Also, the visualization in Figures 7
1145	and 8 proves the wide dynamic range in the frequency domain of deep blocks illustrated in Figure 5
1146	and o proves the write dynamic range in the nequency domain of deep brocks must deed in Figure 5.
1147	
1148	
1149	
1150	
1151	
1152	
1153	
1154	
1155	
1156	
1157	
1158	
1159	
1160	
1161	
1162	
1163	
1164	
1165	
1166	
1167	
1168	
1169	
1170	
1171	
1172	
1173	
1174	
1175	
1176	
1177	
1178	
1179	
1180	
1181	
1182	
1183	
1184	
1185	
1186	
1187	

1189											
1190											
1191											
1192										The Mar	
1193	Output of 1-st Block	1				Output of 1-st Block	=		- 49, A		
1194											
1195	Output of 3-rd Block					Output of 3-rd Block			歌		
1196										S. A. See	
1197	Output of 5-th Block	Ser.				Output of 5-th Block				<u>.</u>	
1198	的。在 中的	1850	The de		5.24				100.0		1.500 (m
1199	Output of 7-th Block	12	1.0		1.3	Output of 7-th Block		13	23	10	1.1
1200			100 A						10.01		5.00
1201	Output of 13-th Block			.		Output of 13-th Block		120		.	
1202											
1203	(a). Feature	maps of	ReActNo	etA for a	goose	(b). l	Feature r	naps of I	ReActNe	tA for a v	vacuum
1204	-								liage	1411	
1205	Output of I-st Block	EL.				Output of 1-st Block					32.6
1206			1 100								
1207	Output of 3-rd Block	54			-25	Output of 3-rd Block		3%	4	11	
1208		56 C.	100.00		1205				11 × 12 ×		
1209	Output of 5-th Block	EL.	Sec.		- (9)	Output of 5-th Block	1. C.	34	25	王子说	
1210				100	1000			100.00	-1.5	-	1000
1211	Output of 7-th Block		10	10	2.20	Output of 7-th Block		を売	624	- N.	12.00
1212								2.01			100
1213	13-th Block	5.5	. 10		1.1	13-th Block	2.5	22			50 B
1214			D N C			(J) E			D Net C		
1215	(c). Feature f	naps of C	nage	mail for	a goose	(a). Fe	eature ma	aps of QI	nage	nall for a	vacuum
1210	A MAR	16	16		15	Cutowa of				-2.	
1217	1-st Block	-12				1-st Block				-14 - 14 -	
1219	Output of					Output of					
1220	stride=(2,1) in 3-rd Block			Constant State		<i>stride=(2,1)</i> in 3-rd Block				19.6	一個名
1221	Output of	20	and the	1	15	Output of	Ser Mar		- X -)		
1222	3-rd Block		$\lambda_{i,j} = 0$	1-1-2-1-1-1-1-1-1-1-1-1-1-1-1-1-1-1-1-1		3-rd Block			國王自		
1223	Output of stride=(2,1)	14			18	Output of stride=(2,1)				-10-11	
1224	in 5-th Block	C.C.S.				in S-th Block	第一492 22				
1225	Output of		1		12	Output of 5-th Block	1	Sec. 1	$-\frac{1}{2} \phi(t)$		
1226	5-01 DAGK	484	1944				学校包	1.	1.16		
1227	Output of stride=(2,1)	88 X.	1.57	1.10	36.	Output of stride=(2,1)	1		29%	10.00	-328.4
1228	in 7-th Block					In 7-th block	and a case				1.04
1229	Output of 7-th Block	38:		14	14	Output of 7-th Block	1.50	20	32.3	1.0	100
1230	1.000	1997	100	100			2010		1 0767	8. S	
1231	Output of stride=(2,1) in 13-th	3-11		120	100	Output of stride=(2,1) in 13-th		\mathbb{C}^{2}			34 C.
1232	Block					BIOCK					
1233	Output of 13-th Block					Output of 13-th Block	5		6 D		107
1234											
1235	(e). Feature m	naps of Q	SB-Net-S	Small for	a goose	(f)	. Feature	e maps of	f QSB-N	et-Small	for a
1236		11	nage					vacuu	in inage		
1237	Figure 7: Vis	ualizati	ions of	feature	maps fo	or ReA	.ctNet/	A and p	oropose	d Sma	ll models.
1238											

1294 1295	(e). Feature n	aps of Q goos	SB-Net-	Large(SE	E1) for a	(f). Fe	eature m	aps of Q vacuu	SB-Net-l m image	Large(SE	(1) for a
1293											
1292	Output of 9-th Block		5.0		2.2	Output of 9-th Block		10			
1291	in 9-th Block					in 9-th Block					
1290	Output of stride=(2,1)	•	122	1.1	283	Output of stride=(2,1)		6.25	100	\mathcal{L}	2.00
1289		100				. JI DIJEK	10	1.12	-	12	1.12
1288	Output of 7-th Block	1.	22	10	3.25	Output of 7-th Block	3	14			\mathbb{R}^{n}
1287	m /-th Block					in r-th Block					
1286	Output of stride=(2,1)	As =	12			Output of stride=(2,1)	Sec.	18	AL.	1.10	
1285	J-BI DOCK	120	了我们			o-uri Block	No.		一次是		1.4
1284	Output of S-th Block	36			11.1	Output of	13	18	All	292	
1283	in 5-th Block		S. A. C			in 5-th Block	132		1944E		
1282	Output of stride=(2,1)	$(\mathcal{A}_{\mathcal{A}})$	Sec. 1			Output of stride=(2,1)	Sta /)	20	1. Carl	- Aller	
1281	3-rd Block		See.	and the	-	3-rd Block		Sec.	副公告		
1280	Output of	445	- 16.	1925		Output of	-	Sec.			10
1279	stride=(2,1) in 3-rd Block	1. 3				stride=(2,1) in 3-rd Block		S. C.	· 松阳 P		
1278	Output of		and street			Output of			1.00 100		
1277	1-st Block					Output of 1-st Block		省	1 Ar		
1273	Output of	18	l.			Output of	-				
1274	(c). Feature m	aps of Q ii	SB-Net-l nage	Large for	a goose	(d).	Feature	e maps of vacuu	t QSB-N m image	et-Large	tor a
12/3			CD N - 1			7 I N	Ent				6
1272	Output of 9-th Block	-				Output of 9-th Block					
1271											
1270	Output of stride=(2,1) in 9-th Block	\mathcal{D}_{i}	8 <u>-</u>	.		Output of stride=(2,1) in 9-th Block	100	10	400	9 . %	
1269	1.000				3000				10.00	1.595	
1268	Output of 7-th Block	10	1	24	24	Output of 7-th Block	3.9			36	
1267	in /-th Block					in 7-th Block					
1266	Output of stride=(2,1)	Sec.	1	1250	1	Output of stride=(2,1)	19	3.4		:静子	10
1265	5-th Block			120	15-10	5-th Block	14.42		-16 V	1.1	
1264	Output of			5	×1.	Output of	14			24.5	- 8
1263	in S-th Block	1.543	12			sandte=(2,1) in S-th Block	2 5 5,		-65/4	32.2	S. C.
1262	Output of stride=(2.1)		· 76.	675X		Output of		17.0		12.77	14 12 2
1200	Output of 3-rd Block		-12	-2	1894 (C)	Output of 3-rd Block			1 · ·		14
1209			N PS						28.1	- dear	all a
1258	Output of stride=(2,1) in 3-rd Block		22			Output of stride=(2,1) in 3-rd Block		All a			A Co
1257						0					
1256	Output of 1-st Block	1	-		-1.11	Output of 1-st Block		= Y			
1255		11	nage					111	lage	Contraction of the second	1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1
1254	(a). Feature r	naps of C	QB-Net-L	arge for	a goose	(b). Fe	ature ma	aps of QI	B-Net-La	urge for a	vacuum
1253	100				- - -						
1252	Output of 9-th Block					Output of 9-th Block			12		1. A.
1251	2.88	123	26	40	19 P 10	1300Ck	1.0	-92	1.0	4.2	100.0
1250	Output of 7-th Block	2i	1	9 6	5. e	Output of 7-th Block		538		2.5	
1249	. Nik	1.20	19 10 2	No fee	1.0	2000 h.e.	-112P	用可能	四月(13	नदम	1.4
1248	Output of 5-th Block				10	Output of 5-th Block	19			14.19	-260
1247	- 4		A Strange		The second	ə-rd Block	132	-31			
1246	Output of 3-rd Block	18	1 m	1	-1 - T	Output of			-Rep.	-	=140%
1245	i st Block	NA.				1-st Block		14-	二個日		
1243	Output of	12		26		Output of	-			40	4
1242											
1040											

Figure 8: Visualizations of feature maps for proposed Large models.

A.9 T-SNE AS A TOOL FOR VISUALIZING AND INTERPRETING IMAGE CLASSIFICATION RESULTS 1298

The evaluations of models using t-distributed stochastic neighbor embedding (t-SNE) are illustrated in Figure 9. In FP32 ResNet18 and MobileNetV2, the visual patterns show that the clustered features are well separated. The visual patterns of ReActNetA (Liu et al., 2020) and proposed models indicate that BCNNs suffer from difficulties in distinguishing between several classes. Compared with the Small models, the Large models slightly reduced the intersection of clustered features. Therefore, we conclude that the visualized patterns indirectly prove the performance enhancements in the Large models.



1350 A.10 VISUALIZED RESULTS OF SEMANTIC SEGMENTATION

We demonstrated that QSB-Net-Large could learn a good representation of objects in semantic segmentation, where QSB-Net-Large combined with DeepLabv3+ was trained. The PASCAL VOC 2012 (Everingham et al., 2010) segmentation dataset contains 1,465 training, 1,449 validation, and 1,456 test images having pixel-level annotations. The dataset was augmented by the extra annotations from the PASCAL VOC 2011, resulting in 10,582 augmented training images. Figures 10 and 12 show the visualized results of semantic segmentation using QSB-Net-Large on the VOC PASCAL 2012 dataset (Everingham et al., 2010).



Figure 10: Semantic segmentation results of QSB-Net-Large. On the left, input images are illustrated. In the middle, the segmentation masks overlayed on its input image are shown. On the right, the predicted segmentation masks are given.



Figure 11: Semantic segmentation results of QSB-Net-Large. On the left, input images are illustrated. In the middle, the segmentation masks overlayed on its input image are shown. On the right, the predicted segmentation masks are given.

1445

1446 1447

A.11 ESTIMATION OF ENERGY CONSUMPTION ON REAL HARDWARE

1448 We performed an energy consumption analysis on real hardware using an RPi 4B on Manjaro OS. 1449 Due to the limitations of using the power management ICs (PMIC), we cannot use power profiling tools. Therefore, we employed an alternative method, where we measured the current drawn from 1450 the USB power line using a current probe over 500 test runs. To assess the impact of running a 1451 benchmark evaluation program in Appendix A.7, we compared the current consumption under two 1452 conditions: with and without executing the program on the RPi 4B. When the models were not 1453 being evaluated, the current consumption was approximately 713 mA @ 5V. In contrast, during the 1454 execution of the models listed in Table 3, the current consumption ranged from 912 mA to 936 mA. 1455 Considering the resolution of the current probe, we thought that these variations were not significant. 1456

1457 While other system activities such as parameter uploading and log dumping were also performed during the measurements, their impact on the comparison was not great, considering the negligible

relative difference. Based on these observations, we conclude that the systematic energy consump-tion is largely proportional to inference latency because no substantial differences in system power consumption were detected during the evaluations.

		4
	Allowet 1990 Allowet 1990 Allow	A CONCOLORIZA

Figure 12: A realistic environment for estimating energy consumption. Using a current probe with \pm 1.5% accuracy resolution, the amount of current in a USB power cable was measured.