

Politechnika Warszawska

WYDZIAŁ ELEKTRONIKI
I TECHNIK INFORMACYJNYCH



Instytut Radioelektroniki i Technik Multimedialnych

Praca dyplomowa

na studiach: Studium Podyplomowe
Głębokie sieci neuronowe - zastosowania w mediach cyfrowych

Lightweight neural architectures by example - car model classification
using GhostNet

Piotr Chaberski

opiekun pracy:
Xin Chang

WARSZAWA 2020

Lightweight neural architectures by example - car model classification using GhostNet

Abstract. This is the abstract!

Keywords: XXX, XXX, XXX

Contents

1. Introduction	1
1.1. Problem background	1
1.2. Problem statement	2
2. Project description	3
2.1. Stanford Cars Dataset	3
2.2. GhostNet architecture	4
3. Experimentation setup	9
3.1. Project structure	9
3.2. Working environment	9
3.3. Configuration	9
3.4. Experiment tracking	9
4. Results	11
4.1. Best model	11
4.2. Experiments step-by-step	11
4.2.1. Loss function	13
4.2.2. Normalization	13
4.2.3. Augmentations	13
4.2.4. Grayscale conversion	13
4.2.5. Bounding boxes utilization	13
4.2.6. Optimizer change and L2 regularization	13
4.2.7. Dropout rate tests	13
4.2.8. Last layer size tests	13
4.2.9. Automatic learning rate scheduling	13
4.2.10. Controlled learning rate scheduling	13

4.2.11. Weight decay adjustment	13
4.2.12. Dropout rate verification	13
4.2.13. Additional augmentations tests	13
4.2.14. Learning rate scheduler adjustment	13
4.2.15. Last layer size sanity check	13
4.2.16. Learning rate annealing tests	13
References	15

1. Introduction

1.1. Problem background

One of the ongoing directions of deep learning research in computer vision and image recognition (but not only) is related to the reduction of neural network size and the number of operations needed for inference while preserving good level of performance in terms of classification accuracy or other metrics. The one important reason to do so is to reduce training times and costs and speed up the iterative process of hyperparameter optimization. Another drawback of large networks in some applications can be extensive overfitting. But the most important reason justifying the search for more efficient neural architectures is that in many practical applications models are needed to be deployed not on a multi-GPU servers or on cloud, but rather as a part of embedded systems on devices with very limited computational power and memory like smartphones, car systems or other devices with so-called intelligent modules.

At the time of completing this work there is already a significant number of different propositions of architectures aiming to reduce the number of parameters and FLOPS needed to efficiently perform image classification tasks [1]. Those architectures are most commonly trained on ImageNet (<http://www.image-net.org/>) and, among others, two metrics are reported on this dataset: accuracy and FLOPS, along with the total number of parameters. Those values give an impression about architecture efficiency in terms of trade-off between prediction quality, inference speed and required memory. Some, but definitely not all, of the successful implementations are:

- SqueezeNet (2016) [2]
- MobileNet (V1: 2017, V2: 2018, V3: 2019) [3][4][5]
- SqueezeNext (2018) [6]
- ShuffleNet (2018) [7][8]

- EfficientNet (2019) [9]
- HarDNet (2019) [10]
- GhostNet (2020) [11]

Most of these architectures come with different customizable variants. For example, EfficientNet has 8 different basic configurations (named b0 to b7) that differ in terms of complexity. Others, like MobileNet, were reworked and upgraded resulting in different versions (there are currently three versions of MobileNets, named simply V1, V2 and V3).

The above-mentioned architectures are capable of achieving good accuracy scores with very limited number of parameters and floating point operations required. For example, EfficientNet-b0 has 77.1% accuracy on ImageNet with 5.3 M parameters and 0.39 GFLOPS. GhostNet gets 73.98% accuracy with only 4.1 M parameters and 0.142 GFLOPS. On the contrasts, ResNet-50 to achieve 75.3% accuracy requires 25.6 M parameters and 4.1 GFLOPS to process the image of the same size (224x224 RGB).

1.2. Problem statement

This project is a part of a broader conception to create a mobile application to recognize car models from pictures taken by the users. The initial idea was to:

1. Pick some of the efficient mobile architectures (the project was intended to be carried out in a group), train them on a open dataset of car images and compare in terms of accuracy, model size and FLOPS.
2. Prepare custom dataset of images taken and labelled personally, finetune the best model from step 1 to reflect car models distribution on the streets of Poland.
3. Prepare model for deployment, create a simple Android application that allow to take a picture and recognize a car model.

This work focuses only on step one with selected architecture. Specifically, it describes the process of training and optimizing hyperparameters of GhostNet [11] model using Stanford Cars Dataset [12] (https://ai.stanford.edu/~jkrause/cars/car_dataset.html) to check the performance of this particular novel mobile architecture in a car model recognition task.

2. Project description

2.1. Stanford Cars Dataset

Stanford Cars Dataset [12] is a dataset published by Jonathan Krause of Stanford University and is publicly available at https://ai.stanford.edu/~jkrause/cars/car_dataset.html.



Figure 1: Example images from Stanford Cars Dataset

The dataset contains 16,185 images of 196 classes of car models (precisely, class label contains information about make, model and production year of a car). Dataset has been splitted with stratification into two parts:

- 8,144 images as a training set
- 8,041 images as a test set

In addition to class labels, both subsets have also bounding boxes (as 4 coordinates in metadata files).

Images are originally of different sizes, mostly in RGB, but there are some grayscale images which has to be taken into account during preprocessing. Another thing to be aware of is

that the dataset has been updated at some point - the images and the split did not change, but the file names were reordered and metadata was reorganized for the ease of use.

2.2. GhostNet architecture

GhostNet [11] is the architecture designed and first implemented by the research team at Huawei Noah's Ark Lab (<http://www.noahlab.com.hk/>). It is based on the observation, that standard convolutional layers with many filters are large in terms of number of parameters and computationally expensive, while often producing redundant feature maps that are very much alike each other (they might be considered as “ghosts” of the original feature map). The goal of the GhostNet design is not to get rid of those redundant feature maps, because they often help the network to comprehensively understand all the features of the input data. Instead of that, the focus is on obtaining those redundant feature maps in a cost-efficient way.

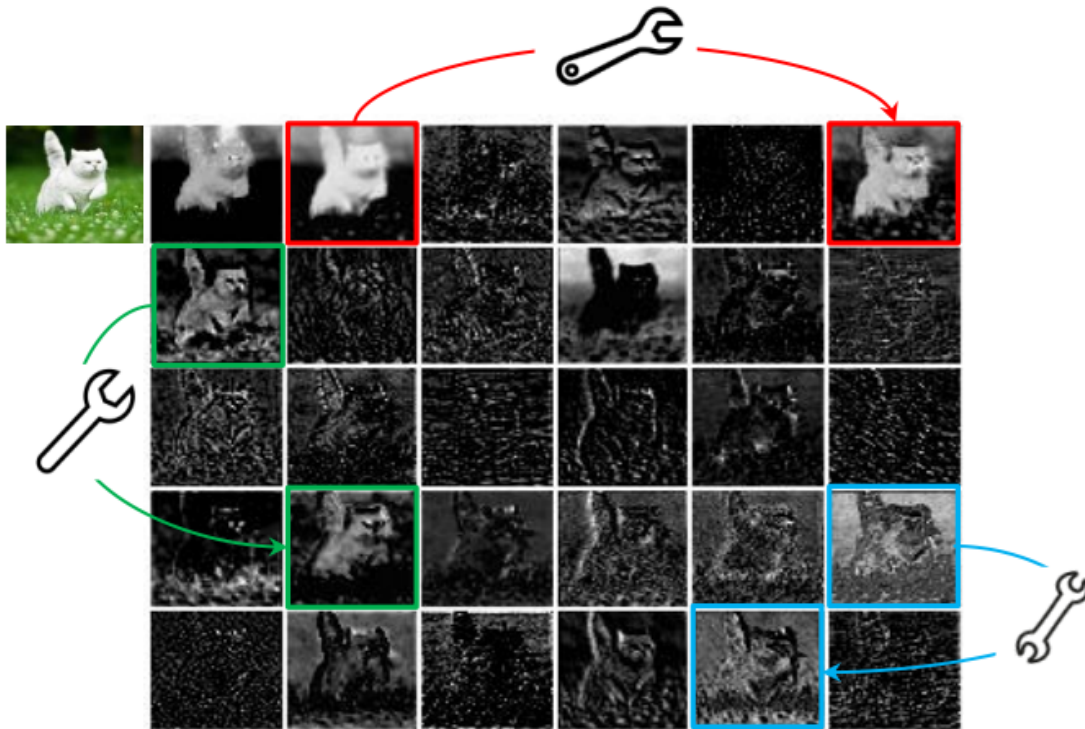


Figure 2: Redundant feature maps from ResNet-50 (picture from paper)

This cost-efficiency in creating feature maps is achieved by introducing GhostModule, namely splitting standard convolutional layer with many filters into two parts. The first

part, still being a standard convolutional layer but with less filters, produces a set of base feature maps. Then the second part, by applying cheap linear operations, produces redundant feature maps from the original set (so-called “ghosts”). In the end, the outputs of the first and the second part are concatenated.

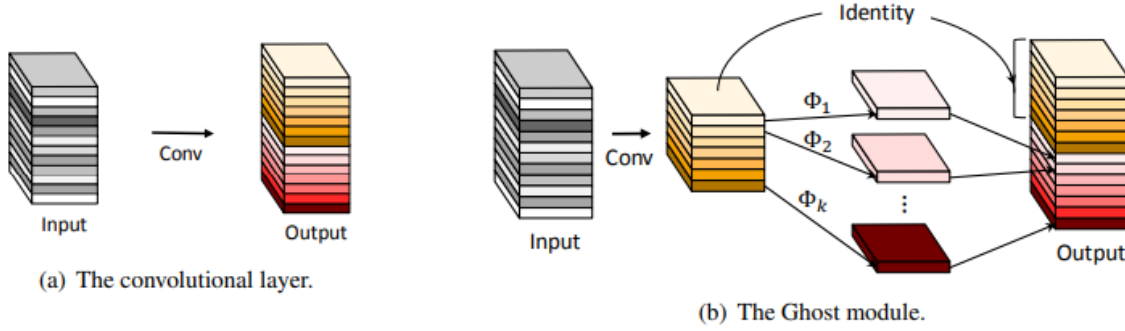


Figure 3: Comparison of standard convolution (a) and GhostModule (b) (picture from paper)

The above mentioned cheap linear operations are implemented using depthwise convolutions [13] (although other options like affine or wavelet transforms were also tested by the authors). With this assumption, GhostModule can be implemented in PyTorch as follows:

```
class GhostModule(nn.Module):
    def __init__(
        self, inp, oup,
        kernel_size=1, ratio=2, dw_size=3, stride=1,
        relu=True
    ):
        super().__init__()
        self.oup = oup
        init_channels = math.ceil(oup / ratio)
        new_channels = init_channels*(ratio-1)

        self.primary_conv = nn.Sequential(
            nn.Conv2d(
                inp, init_channels, kernel_size, stride,
                kernel_size//2, bias=False
            ),
            nn.BatchNorm2d(init_channels),
```

```

        nn.ReLU(inplace=True) if relu else nn.Sequential(),
    )

    self.cheap_operation = nn.Sequential(
        nn.Conv2d(init_channels, new_channels, dw_size, 1,
            dw_size//2, groups=init_channels, bias=False
        ),
        nn.BatchNorm2d(new_channels),
        nn.ReLU(inplace=True) if relu else nn.Sequential(),
    )

    def forward(self, input):
        output_1 = self.primary_conv(input)
        output_2 = self.cheap_operation(output_1)
        output = torch.cat([output_1, output_2], dim=1)
        return output[:, :self.oup, :, :]

```

Two GhostModules combine for a basic building block of GhostNet - the GhostBottleneck, which is based on the concept taken from MobileNet-V3 design [5] (additionally, in some GhostBottlenecks, similarly to MobileNet-V3, Squeeze-and-Excitation modules are used [14]). The first GhostModule in a GhostBottleneck expands the number of channels, while the second one, after ReLU, reduces them again. There is also a residual connection over the two GhostModules. GhostBottleneck has also strided version (with `stride=2` depthwise convolution between GhostModules) which is applied at the end of each stage of GhostNet.

To form up the entire GhostNet architecture several GhostBottlenecks are combined in a sequence which is followed by global average pooling and a convolution which transforms feature maps to the feature vector of length 1280. This feature vector, after dropout layer, is then transformed with a fully connected layer to the size of output number of classes.

GhostNet architecture based on paper:

Input	Operator	#exp	#out	SE	Stride
2242 x 3	Conv2d 3x3	-	16	-	2
1122 x 16	G-bneck	16	16	-	1
1122 x 16	G-bneck	48	24	-	2

Input	Operator	#exp	#out	SE	Stride
562 x 24	G-bneck	72	24	-	1
562 x 24	G-bneck	72	40	1	2
282 x 40	G-bneck	120	40	1	1
282 x 40	G-bneck	240	80	-	2
142 x 80	G-bneck	200	80	-	1
142 x 80	G-bneck	184	80	-	1
142 x 80	G-bneck	184	80	-	1
142 x 80	G-bneck	480	112	1	1
142 x 112	G-bneck	672	112	1	1
142 x 112	G-bneck	672	160	1	2
72 x 160	G-bneck	960	160	-	1
72 x 160	G-bneck	960	160	1	1
72 x 160	G-bneck	960	160	-	1
72 x 160	G-bneck	960	160	1	1
72 x 160	Conv2d 1x1	-	960	-	1
72 x 960	AvgPool 7x7	-	-	-	-
12 x 960	Conv2d 1x1	-	1280	-	1
12 x 1280	FC	-	1000	-	-

GhostNet architecture described above (and in original paper as well) is the basic setup which can be modified by structuring GhostBottlenecks in different sequences. This basic setup, as mentioned before, gets 73.98% accuracy on ImageNet with 4.1 M parameters and requires only 0.142 GFLOPS to process 224x224 RGB image. Other more complex variations, as presented in paper, show superiority over previous designs like MobileNet or ShuffleNet getting better accuracy with less FLOPS and latency.

Full PyTorch implementation of GhostNet that was used in this work is available at GitHub repository of the project.

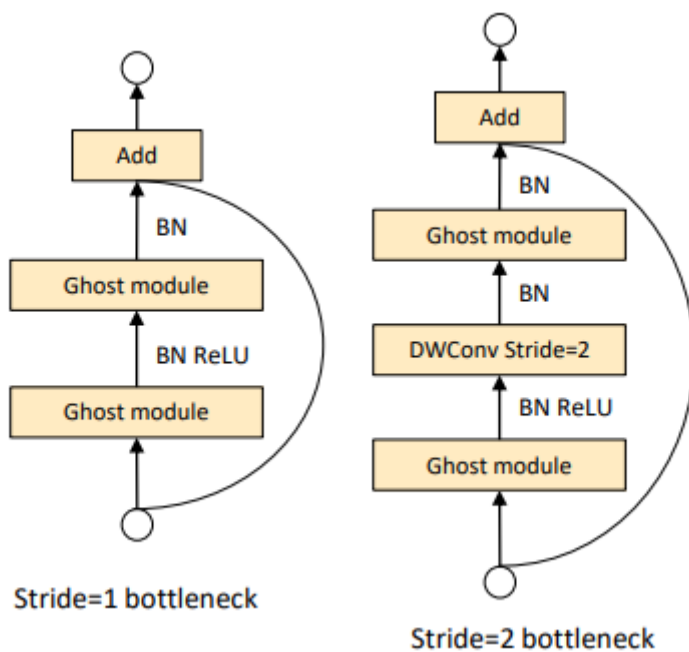


Figure 4: GhostBottleneck (picture from paper)

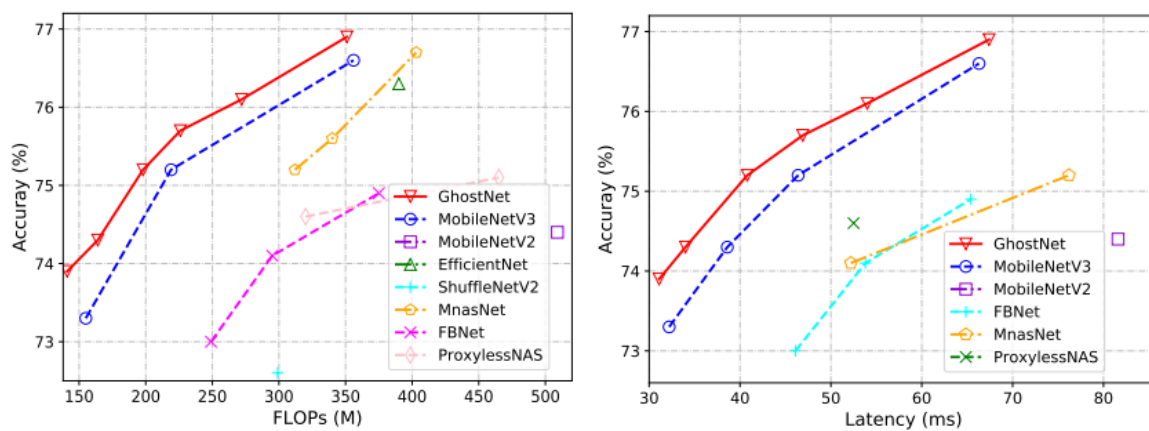


Figure 5: GhostNet comparison with some other mobile architectures (pictures from paper)

3. Experimentation setup

The experimentation setup is entirely based on Python. GhostNet (and some other networks, which also can be used) implementation is written in PyTorch. Training process is orchestrated using `pytorch-lightning` package and controlled by parameters passed through YAML config file. Neptune.ai experiment management tool (<https://neptune.ai/>) was used for experiment tracking. To build an environment for data preparation and model training Python `virtual env` utility was used. In addition to local training setup there is also a possibility to recreate the project environment and run training on Google Colab platform using a prepared Jupyter Notebook.

3.1. Project structure

Source code for the project is available in GitHub repository: <https://github.com/pchaberski/cars>. Repository contains the following elements:

- `documentation` folder containing markdown files with project documentation, images, bibliography as a `bib` file and some tools for bibtex and PDF conversion

3.2. Working environment

3.3. Configuration

3.4. Experiment tracking

4. Results

4.1. Best model

4.2. Experiments step-by-step

	experiment description	train_acc	valid_acc
C-1	Baseline (Cross Entropy Loss)	92.49%	8.15%
C-2	Loss function change (Label Smoothing Cross Entropy)	98.89%	9.12%
C-3	Augmentations: horizontal flip, affine	99.45%	11.96%
C-4	Augmentations: horizontal flip, affine, erasing	99.76%	51.92%
C-5	Augmentations: horizontal flip, erasing, color jitter	98.12%	38.08%
C-6	Augmentations: horizontal flip, affine, erasing, color jitter	93.68%	38.68%
C-7	Augmentations: horizontal flip, affine, color jitter	99.73%	54.28%
C-8	Grayscale: no normalization, no augmentations	99.49%	6.58%
C-9	Grayscale: with normalization, no augmentations	97.13%	8.68%
C-10	Grayscale: normalization, best RGB augmentations	7.58%	3.91%
C-11	Training set cropping with bounding boxes	4.36%	3.07%
C-12	Training set cropping + background erasing	99.67%	50.51%
C-13	L2 regularization with AdamW: weight decay = 0.1	99.44%	63.39%
C-14	L2 regularization with AdamW: weight decay = 0.2	98.84%	68.50%
C-15	L2 regularization with AdamW: weight decay = 0.3	95.83%	61.84%
C-16	L2 regularization with AdamW: weight decay = 0.4	95.95%	65.14%
C-17	L2 regularization with AdamW: weight decay = 0.5	90.38%	59.95%
C-18	Dropout rate tests: dropout = 0.1	99.11%	66.90%
C-19	Dropout rate tests: dropout = 0.3	98.62%	67.81%
C-20	Dropout rate tests: dropout = 0.4	96.52%	64.88%

	experiment description	train_acc	valid_acc
C-21	Dropout rate tests: dropout = 0.5	96.28%	66.75%
C-22	Last layer size tests: out channels = 320	97.13%	68.93%
C-23	Last layer size tests: out channels = 640	96.13%	63.13%
C-24	Last layer size tests: out channels = 960	98.23%	64.96%
C-25	Last layer size tests: out channels = 1600	98.99%	63.11%
C-26	Automatic LR scheduling: take #1	99.82%	74.60%
C-27	Automatic LR scheduling: take #2	99.78%	76.20%
C-28	Automatic LR scheduling: take #3	99.83%	75.14%
C-29	Automatic LR scheduling: take #4	99.78%	74.82%
C-30	Controlled LR scheduling: milestones = [28, 48, 68, 88]	80.66%	57.82%
C-31	Controlled LR scheduling: milestones = [36, 56, 76, 96]	95.03%	64.93%
C-32	Controlled LR scheduling: milestones = [44, 64, 84, 104]	98.68%	68.79%
C-33	Controlled LR scheduling: milestones = [52, 72, 92, 112]	99.60%	71.59%
C-36	Weight decay adjustment: weight decay = 0.5	98.84%	79.40%
C-37	Weight decay adjustment: weight decay = 0.3	99.57%	74.44%
C-38	Weight decay adjustment: weight decay = 0.4	99.37%	78.82%
C-39	Weight decay adjustment: weight decay = 0.6	98.67%	82.55%
C-40	Weight decay adjustment: weight decay = 0.7	99.24%	75.12%
C-41	Dropout rate verification: dropout = 0.3	98.49%	82.08%
C-42	Dropout rate verification: dropout = 0.4	95.34%	79.57%
C-43	Dropout rate verification: dropout = 0.5	96.08%	77.87%
C-44	Dropout rate verification: dropout = 0.25	98.79%	82.45%
C-45	Additional augmentations test: resized crop	97.56%	78.73%
C-46	Additional augmentations test: rotation	97.03%	78.25%
C-47	Additional augmentations test: perspective	97.42%	80.22%
C-48	Additional augmentations test: erasing	93.68%	80.56%
C-50	LR scheduler adjustment: milestones = [67, 82, 95, 107]	98.94%	83.79%
C-51	LR scheduler adjustment: milestones = [63, 78, 91, 103]	98.86%	82.54%
C-53	LR scheduler adjustment: milestones = [66, 81, 94, 106]	98.96%	83.02%
C-55	LR scheduler adjustment: milestones = [68, 83, 96, 108]	98.78%	83.72%
C-56	LR scheduler adjustment: milestones = [64, 79, 92, 104]	98.99%	82.79%
C-58	Last layer size sanity check: out channels = 1280	99.44%	78.83%
C-63	LR annealing test: LR geometric sequence	99.80%	70.51%
C-64	LR annealing test: exponentiation base = 0.955	98.49%	60.70%

experiment description		train_acc	valid_acc
C-65	LR annealing test: exponentiation base = 0.975	99.66%	73.07%
C-66	LR annealing test: exponentiation base = 0.98	98.72%	70.46%

4.2.1. Loss function**4.2.2. Normalization****4.2.3. Augmentations****4.2.4. Grayscale conversion****4.2.5. Bounding boxes utilization****4.2.6. Optimizer change and L2 regularization****4.2.7. Dropout rate tests****4.2.8. Last layer size tests****4.2.9. Automatic learning rate scheduling****4.2.10. Controlled learning rate scheduling****4.2.11. Weight decay adjustment****4.2.12. Dropout rate verification****4.2.13. Additional augmentations tests****4.2.14. Learning rate scheduler adjustment****4.2.15. Last layer size sanity check****4.2.16. Learning rate annealing tests**

References

[1]

M. Hollemans. New mobile neural network architectures. Accessed 2020-09-26: <https://machinethink.net/blog/mobile-architectures>, 2020.

[2]

F. N. Iandola, S. Han, M. W. Moskewicz, K. Ashraf, W. J. Dally, and K. Keutzer. Squeezenet: Alexnet-level accuracy with 50x fewer parameters and <0.5mb model size, 2016. [arXiv]

[3]

A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, and H. Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications, 2017. [arXiv]

[4]

M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, and L.-C. Chen. Mobilenetv2: Inverted residuals and linear bottlenecks. 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition, Jun 2018. [DOI | [http](#)]

[5]

A. Howard, M. Sandler, B. Chen, W. Wang, L.-C. Chen, M. Tan, G. Chu, V. Vasudevan, Y. Zhu, R. Pang, and et al. Searching for mobilenetv3. 2019 IEEE/CVF International Conference on Computer Vision (ICCV), Oct 2019. [DOI | [http](#)]

[6]

A. Gholami, K. Kwon, B. Wu, Z. Tai, X. Yue, P. Jin, S. Zhao, and K. Keutzer. Squeezenext: Hardware-aware neural network design. 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW), Jun 2018. [DOI | [http](#)]

[7]

X. Zhang, X. Zhou, M. Lin, and J. Sun. Shufflenet: An extremely efficient convolutional neural network for mobile devices. 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition, Jun 2018. [DOI | [http](#)]

[8]

N. Ma, X. Zhang, H.-T. Zheng, and J. Sun. Shufflenet v2: Practical guidelines for efficient cnn architecture design. Lecture Notes in Computer Science, page 122–138, 2018. [DOI | [http](#)]

[9]

M. Tan and Q. V. Le. Efficientnet: Rethinking model scaling for convolutional neural networks, 2019. [[arXiv](#)]

[10]

P. Chao, C.-Y. Kao, Y. Ruan, C.-H. Huang, and Y.-L. Lin. Hardnet: A low memory traffic network. 2019 IEEE/CVF International Conference on Computer Vision (ICCV), Oct 2019. [DOI | [http](#)]

[11]

K. Han, Y. Wang, Q. Tian, J. Guo, C. Xu, and C. Xu. Ghostnet: More features from cheap operations. 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), Jun 2020. [DOI | [http](#)]

[12]

J. Krause, M. Stark, J. Deng, and L. Fei-Fei. 3d object representations for fine-grained categorization. In 4th International IEEE Workshop on 3D Representation and Recognition (3dRRR-13), Sydney, Australia, 2013.

[13]

A. Pandey. Depth-wise convolution and depth-wise separable convolution. Accessed 2020-09-26: <https://medium.com/@zurister/depth-wise-convolution-and-depth-wise-separable-convolution-37346565d4ec>, 2018.

[14]

J. Hu, L. Shen, and G. Sun. Squeeze-and-excitation networks. 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition, Jun 2018.

