

CNN Project: Dog Breed Classifier

Project definition

Project overview

Ever since early days of computing, neural network architectures have been proposed for simple tasks, and their performances were often as good as clever designed algorithms such as Support Vector Machines, random forest, dictionary extraction, etc. For image processing though, neural networks promised high accuracy results, but it was easier to design algorithms using some sort of signal processing techniques such as Wavelet and Fourier transforms, Haar cascade filters or other types of custom designed filters depending on the problem to be solved. CNN methods were often overlooked because it was cumbersome to come up with an architecture that performed well at validation, and at the cost of spending valuable amount of time training and fine tuning it due the need of huge amounts of data and access to expensive hardware.

Since 2006 the CNN discussion started to gain momentum again thanks to the Internet, which facilitated collaboration, access to data, monetary investment and a gradual improvement in hardware for the Gaming industry. In 2012 AlexNet results on ImageNet challenge detonated the use of CNN for image classification tasks due to its impressive 10.8 percentage points reduction of the top-5 error versus the best performing algorithm at the time.

There have been a lot of new proposed CNN models for object detection and classification lately and their speeds and accuracies vary depending on the platform they are intended to be implemented.

The goal of this project is to build and train a custom model from scratch to classify dog breeds given a dog image. Since this a challenging area and will be cumbersome to come up with a model of my own design, I am going to base the CNN model on a MobileDet set of convolutions as proposed in [4]. The *dogImage* dataset will be used to train and test the CNN classifier [8]. The *lfw* human face dataset [9] will be used as input data to evaluate the trained model and see to which is the predicted dog breed for a given human face.

Problem statement

In ML, image classification is a well-known field of research where there have been multiple approaches proposed. But it is often proved to be a difficult domain because there are few models achieving good results in general datasets. Classifying dog breeds could be challenging even for humans, because some breeds can share some characteristics. Take Figures 1 and 2 as examples, they are different dogs but look alike. That's why for some applications it pays off to specialize a CNN model.



Figure 1 Welsh springer spaniel



Figure 2 Brittany

Metrics

The cross-entropy loss will be used to propagate backwards the error and adjust weights, since we hope the model will output a *softmax* probability vector indicating the confidence it has of an image being this or that breed of dog. Since we are using Pytorch as our default framework, a softmax layer is not necessary since the built-in cross entropy loss layer already computes the output probabilities.

The cross-entropy loss function is defined as:

$$l(x, y) = -\frac{1}{N} \sum_{n=1}^N x_n \log(s(y_n)), \quad s(y) = \frac{e^y}{\sum_{c=1}^C e^{y_c}}$$

Where x and y are the input target and predicted output respectively, N is the minibatch size and C is the number of classes in the dataset.

We are going to use an Adam optimizer due to its adaptative strategy, which consists in computing individual learning rates for each parameter, improving model convergence speed as explained in [6, 10, 11].

Analysis

Data exploration and visualizations

The dog image dataset contains 8351 images; 6680 images for training, 835 for validation and 836 for testing, distributed in 133 dog breed classes. Images vary in size, background, image quality, resolution, illumination, pose, some of them have humans in there, dogs often are cropped, and classes are imbalanced.

Figure 2 shows a sample of 9 dog images retrieved from the test folder.

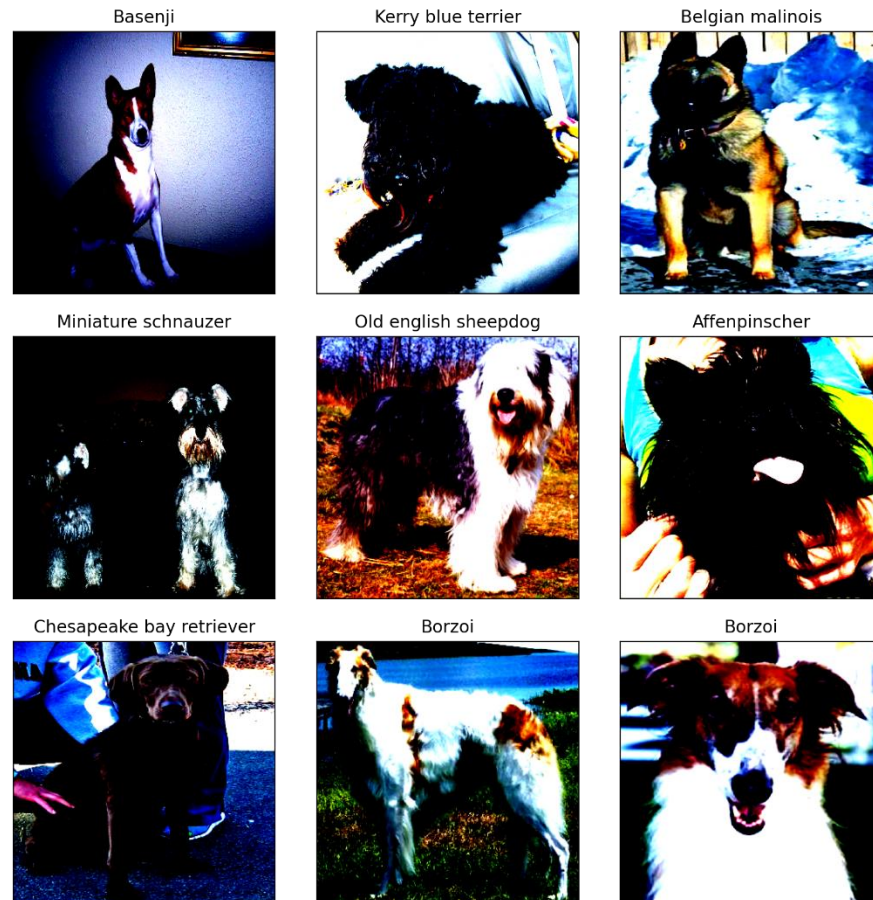


Figure 2 Dog images from test set. These images are preprocessed with resize and normalization transformations.

The human faces dataset contains 13233 human images; distributed in 5750 folders separated by name. Images have the same 250x250 size but vary in background, quality, pixel resolution, illumination, pose and are imbalanced.

Figure 3 shows a sample of 9 human face images from the dataset.

Data preprocessing

An image loader class was written to perform training, validation and testing steps. Since the class is the same for each step there is a parameter that is used only in training which is the transformation argument for data augmentation. The class accepts the batch size parameter to control how many images are used at each epoch step in training.

Thanks to the data augmentation transformation we can (as implied in the name) augment the data we have and increase the number of unique images, since each image is going to be different each time it is requested by the training algorithm.

I used two different sets of augmentation transformations; one was used when I trained the model from scratch the first time, and the second set of transformations was used in the second training step known as transfer learning, using the weights from the first training.



Figure 3 Human face images sample from the dataset.

Methodology

These are the transformations used to augment data in the first training:

- *Color jitter*: this transformation changes the hue, contrast and saturation with probabilities of 0.1, 0.5 and 0.5 respectively. I'm interested in teaching the model to distinguish the same dog shape without the proper color information, hence the change in hue. Hue changes are suppressed in the second set of transformation, since I'm expecting the model to distinguish color variations as much as possible.
- *Random mirror*: this transformation applies a horizontal flip to the image with a 0.5 probability of occurrence each time.
- *Random translation*: with this transformation we teach the model how a dog looks like in different positions in the image.
- *Random resize and crop*: this transformation crops the image at random location and resize it to a specific size, which in this case is the expected model input size, which is 320x320.
- *Normalization*: this transformation normalizes the image using the mean and standard deviation values recommended in Pytorch documentation [6]. For the R, G and B channels, these values are [0.485, 0.456, 0.406] for mean and [0.229, 0.224, 0.225] for standard deviation. These values were calculated using millions of color images.

Figure 4 shows how the same dog image is shown in training at different steps.

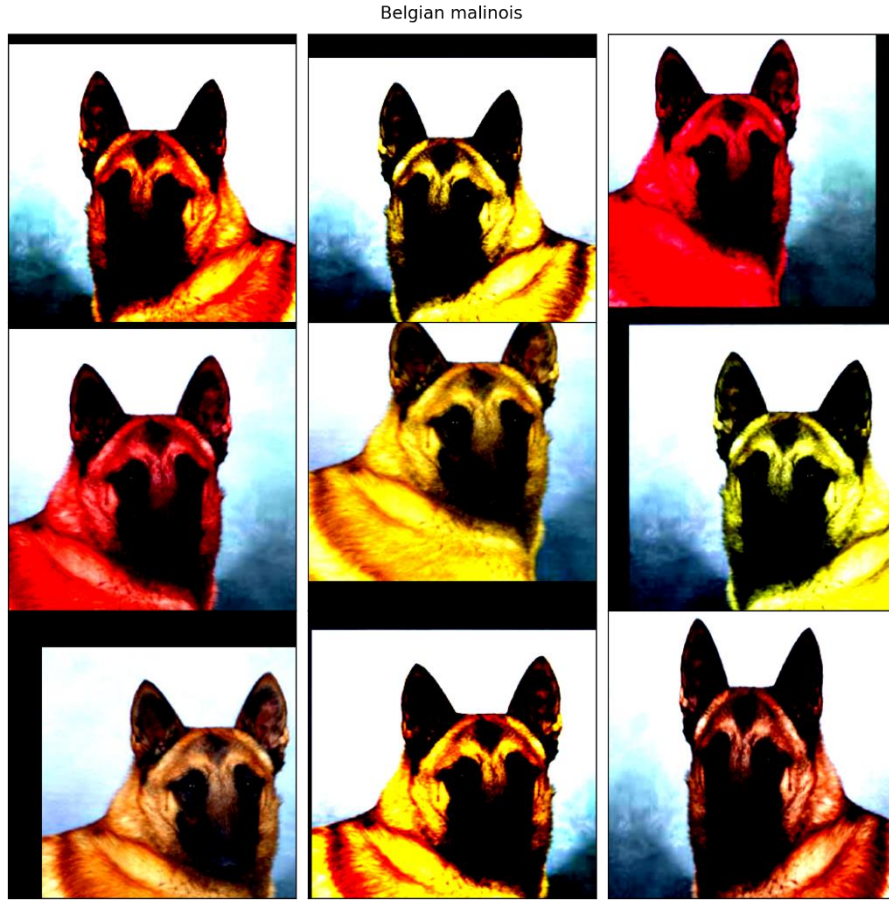


Figure 4 Nine different images of the same dog.

Implementation

The proposed model is a MobileDet architecture [4]. The authors use a neural architecture search (NAS) to find the best architecture configuration based on the target hardware the model is going to be deployed. Instead of using only inverted bottlenecks (IBN) for convolution layers, which are the default choice for state-of-the-art architectures due their low latency and high performance, they add two convolution layer families to the search space; Fused (expansion) and Tucker (compression) convolution layers. These layers are shown in Figure 5.

The reason behind these additions is that an IBN-only architecture, with less floating-point operations per second (FLOPS) is going to be fast and optimized in CPU but will be slower than full convolutions with orders of magnitude more FLOPS in TPU and GPU, since these hardware architectures are designed to accelerate these types of operations.

Since my motivation is to study the MobileDet family, I have implemented the Bottleneck, Tucker and Fused layer classes in the jupyter notebook. As support material, I have used previous

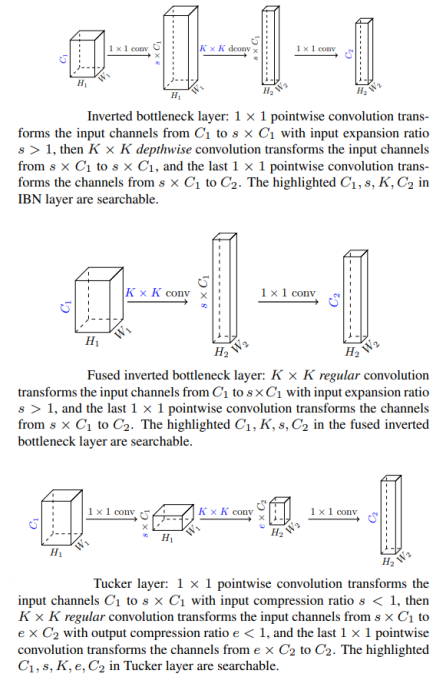


Figure 5 Bottleneck, Fused and Tucker layers explanation [4].

integrations using bottleneck architectures such as MobileNetV1 to V3, EfficientNetV1Bx [1, 2, 3], the MobileDet paper [4] and the Tensorflow research repository [7].

Bottleneck layer implemented in pytorch:

```
class Bottleneck(nn.Module):
    Params = namedtuple('BottleneckParams', 'filters kernel_size stride expansion use_residual use_se activation', defaults=(3, 1, 4, True, False, nn.ReLU6))

    def __init__(self, inputs, filters, kernel_size=3, stride=1, expansion=4, use_residual=True, use_se=False, activation=nn.ReLU6):
        super(Bottleneck, self).__init__()
        expand_filters = int(inputs) * expansion
        layers = []
        if expansion > 1:
            layers.append(RegularConv2d(inputs, expand_filters, 1, activation=activation))
        layers.append(DepthwiseConv2d(expand_filters, expand_filters, kernel_size, stride, activation))
        if use_se:
            hidden_dim = scale(expand_filters, 0.25)
            layers.append(SqueezeAndExcite(expand_filters, hidden_dim, activation))
        layers.append(RegularConv2d(expand_filters, filters, 1, activation=nn.Identity))
        self.net = nn.Sequential(*layers)
        self.use_residual = stride == 1 and inputs == filters and use_residual

    def forward(self, x):
        out = self.net(x)
        return out + x if self.use_residual else out
```

Fused layer implemented in pytorch:

```
class Fused(nn.Module):
    Params = namedtuple('FusedParams', 'filters kernel_size stride expansion use_residual use_se activation', defaults=(3, 1, 8, True, False, nn.ReLU6))

    def __init__(self, inputs, filters, kernel_size=3, stride=1, expansion=8, use_residual=True, use_se=False, activation=nn.ReLU6):
        super(Fused, self).__init__()
        expand_filters = int(inputs) * expansion
        layers = [ RegularConv2d(inputs, expand_filters, kernel_size, stride, activation=activation) ]
        if use_se:
            hidden_dim = scale(expand_filters, 0.25)
            layers.append(SqueezeAndExcite(expand_filters, hidden_dim, activation))
        layers.append(RegularConv2d(expand_filters, filters, 1, activation=nn.Identity))
        self.net = nn.Sequential(*layers)
        self.use_residual = stride == 1 and inputs == filters and use_residual

    def forward(self, x):
        out = self.net(x)
        return out + x if self.use_residual else out
```

Tucker layer implemented in pytorch:

```
class Tucker(nn.Module):
    Params = namedtuple('TuckerParams', 'filters kernel_size stride in_ratio out_ratio use_residual activation', defaults=(3, 1, 0.25, 0.25, True, nn.ReLU6))

    def __init__(self, inputs, filters, kernel_size=3, stride=1, in_ratio=0.25, out_ratio=0.25, use_residual=True, activation=nn.ReLU6):
        super(Tucker, self).__init__()
        self.use_residual = stride == 1 and inputs == filters and use_residual
        hidden_in = scale(inputs, in_ratio)
        hidden_out = scale(filters, out_ratio)
        self.net = nn.Sequential(
            RegularConv2d(inputs, hidden_in, 1, activation=activation),
            RegularConv2d(hidden_in, hidden_out, kernel_size, stride, activation=activation),
            RegularConv2d(hidden_out, filters, 1, activation=nn.Identity),
```

```
)  
def forward(self, x):  
    out = self.net(x)  
    return out + x if self.use_residual else out
```

I have used the proposed backbone architecture targeting the Pixel-1 CPU accelerator (Figure 6) because it is already the best configuration for that hardware according to the NAS algorithm.

Implementation of the backbone settings and the Net class.

```
LAYERS = {  
    'conv': RegularConv2d,  
    'fused': Fused,  
    'tucker': Tucker,  
    'bottleneck': Bottleneck  
}  
  
SETTINGS = [  
    ('conv', RegularConv2d.Params(16, stride=2)),  
    ('bottleneck', Bottleneck.Params(8, expansion=1)),  
    ('fused', Fused.Params(16, stride=2, expansion=4)),  
    ('bottleneck', Bottleneck.Params(16)),  
    ('tucker', Tucker.Params(16, out_ratio=0.75)),  
    ('tucker', Tucker.Params(16, out_ratio=0.75)),  
    ('bottleneck', Bottleneck.Params(32, 5, 2, 8)),  
    ('bottleneck', Bottleneck.Params(32)),  
    ('tucker', Tucker.Params(32, in_ratio=0.75, out_ratio=0.75)),  
    ('tucker', Tucker.Params(32, in_ratio=0.75, out_ratio=0.75)),  
    ('bottleneck', Bottleneck.Params(72, 5, 2, 8)),  
    ('bottleneck', Bottleneck.Params(72)),  
    ('bottleneck', Bottleneck.Params(72)),  
    ('tucker', Tucker.Params(72, out_ratio=0.75)),  
    ('bottleneck', Bottleneck.Params(96, kernel_size=5, expansion=8)),  
    ('bottleneck', Bottleneck.Params(96, expansion=8)),  
    ('bottleneck', Bottleneck.Params(96, kernel_size=5, expansion=8)),  
    ('bottleneck', Bottleneck.Params(96, kernel_size=5)),  
    ('bottleneck', Bottleneck.Params(104, kernel_size=5, stride=2)),  
    ('tucker', Tucker.Params(104, 3, 1, 0.25, 0.75)),  
    ('tucker', Tucker.Params(104, 3, 1, 0.75, 0.75)),  
    ('tucker', Tucker.Params(104, 3, 1, 0.25, 0.75)),  
    ('bottleneck', Bottleneck.Params(192, kernel_size=5, expansion=8))  
]  
  
class Net(nn.Module):  
    def __init__(self, n_classes=133, settings=SETTINGS, hidden_top=256, dropout=0.25, use_se=false, depth_multiplier=1.0):  
        super(Net, self).__init__()  
        inputs = 3  
        layers = []  
        for layer, params in settings:  
            if layer in ['fused', 'bottleneck']:  
                params = params._replace(use_se=use_se)  
            params = params._replace(filters=scale(params[0], depth_multiplier), activation=Swish6)  
            layers.append(LAYERS[layer](inputs, **params._asdict()))  
            inputs = params[0]  
        layers.extend([  
            GlobalAvgPool2d(True),  
            nn.Dropout(dropout),  
            nn.Linear(inputs, hidden_top),  
            nn.Linear(hidden_top, n_classes)  
        ])
```

```
self.net = nn.Sequential(*layers)
def forward(self, x):
    return self.net(x)

model = Net(hidden_top=512, use_se=True, depth_multiplier=1.5)
```

A global average pool is used in the classification head to map the 288 feature maps in the last convolution to a fully connected block with 512 hidden connections in the hidden layer and use a dropout rate of 0.25 in between. A depth multiplier of 1.5 is used to expand the number of original filters used in the Bottleneck, Fused and Tucker convolution layers. This depth multiplier will help me achieve a good accuracy in the first epochs at the expense of having a resulting model with more parameters, but since the convolution layers are well designed, the model is still going to be light.

These decisions helped me to build a model that when trained from scratch, with randomly initialized weights achieves acceptable results when trained in fewer epochs.

The model was trained for 30 epochs from scratch. The Adam optimizer was used with learning rate of 0.001. The image loader feeds 20 images per training step augmenting the data, and the validation loader feeds 8 images per validation step. I decided to use these values due to hardware limitations.

Refinement

In the transfer learning step, the following changes were made to the augment transformation:

- Remove hue variations in the *color jitter* transform.
- Random rotations of -20^0 to 20^0 were added to the images.
- Random gaussian blur is applied with 0.5 probability, using a kernel size of 3x3, sigma x and y of 0.1 and 2 respectively.

A dropout rate of 0.1 is used to let the network rely on more connections. The model was trained for 30 epochs with an early stop of 10 epochs if the validation loss doesn't improve. The learning rate of the optimizer was initialized with 0.0001 to let the model improve at a much slower step than before.

A reduce-on-plateau learning rate scheduler to limit the maximum lambda to be used for the Adam optimizer was used in both trainings. This scheduler doesn't do much since the Adam adapts the learning rate for each parameter accordingly.

Results

Model evaluation and validation

The validation loss stopped improving after 25 epochs, with validation loss of 1.6539 and a validation accuracy of 0.5599 when training from scratch. The test loss and test accuracy were 1.9752 and 0.48 respectively, classifying 409 out of 836 images correctly.

With the changes made in the augment transformation in the transfer learning step, the validation loss stopped improving after 3 epochs, with validation loss of 1.3228 and validation accuracy of 0.6206. The test loss and test accuracy were 1.3385 and 0.62 respectively, classifying 523 out of 836 images correctly.

Target: Pixel-1 CPU
(23.7 mAP @ 122 ms)

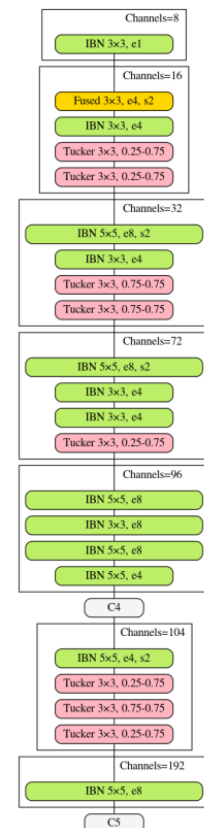


Figure 6 Proposed Pixel-1 CPU accelerator architecture.

The following table summarizes the results.

<i>Training step</i>	<i>Epochs</i>	<i>Test loss</i>	<i>Test accuracy</i>
<i>From scratch</i>	25	1.9752	48% (409/836)
<i>Transfer learning</i>	3	1.3385	62% (523/836)

Justification

The results are not as impressive as if the model would have been trained using transfer learning from a previous long running training in the ImageNet dataset but given the fact it was trained from literally scratch, with randomly initialized weights, achieving 62% accuracy after 25 + 3 epochs is quite remarkable.

The following images show the model in action in the application:

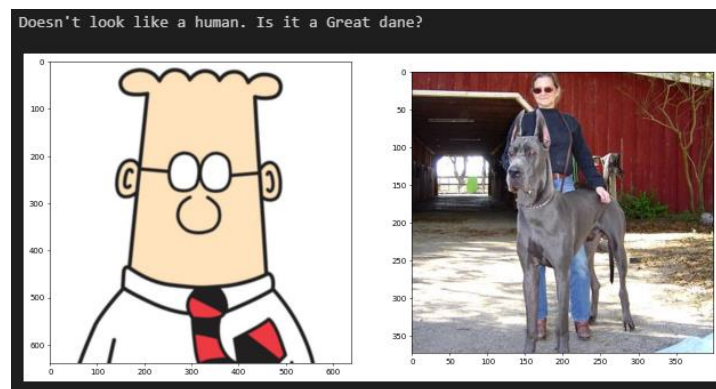


Figure 7 Dilbert being confused with a great dane.

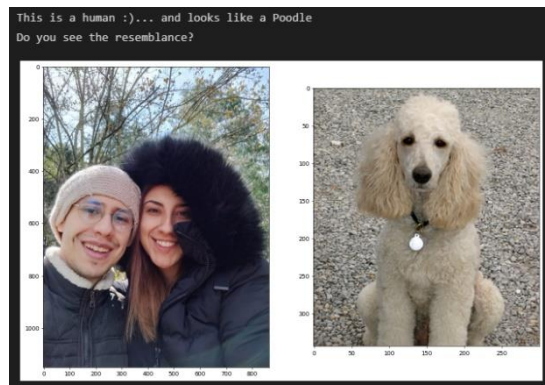


Figure 8 A couple resembling a poodle.

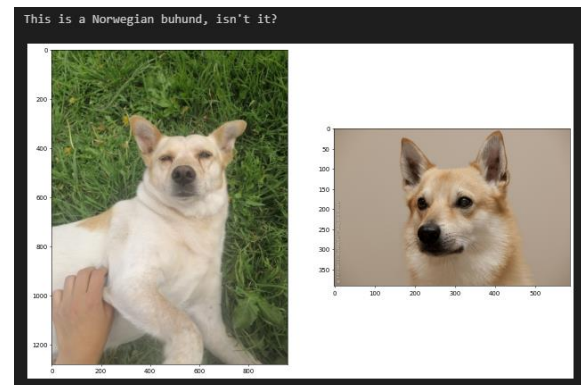


Figure 1 A possible mixed breed dog is classified as a norwegian buhund.

Conclusions

1. Further work is required to implement a proper transfer learning approach, training the model in the ImageNet dataset.
2. Further experimentation is needed to implement a good data augmentation strategy.

3. The NAS approach offers a good starting point to design hardware and possible problem aware CNN architectures, where we could have a model specially designed to maximize dog breed classification accuracy.
4. If properly implemented in a suitable CNN framework, this model could be used in mobile devices or in the edge for industrial applications.
5. I tried to export this model using ONNX to develop a PWA that could do funny things, but the ReLU6 operation was not supported. This might be a good starting point to deploy the model in a web application. With this many parameters a TFLite model would weight as much as 7MB.

References

1. A. G. Howard, W. Wang et al., "MobileNets: Efficient Convolutional Neural Networks for Mobile Vision Applications", <https://arxiv.org/pdf/1704.04861.pdf>
2. M. Sandler, A. Howard et al., "MobileNetV2: Inverted Residual and Linear Bottlenecks", <https://arxiv.org/pdf/1801.04381.pdf>
3. A. Howard, M. Sandler et al., "Searching for MobileNetV3", <https://arxiv.org/pdf/1905.02244.pdf>
4. Y. Xiong, H. Liu et al., "MobileDets: Searching for Object Detection Architectures for Mobile Accelerators", <https://arxiv.org/pdf/2004.14525.pdf>
5. Udacity's dog breed classifier repository, <https://github.com/udacity/deep-learning-v2-pytorch/tree/master/project-dog-classification>
6. Pytorch documentation, <https://pytorch.org/docs/stable/index.html>
7. Tensorflow models research repository, <https://github.com/tensorflow/models>
8. Dog images dataset, <https://s3-us-west-1.amazonaws.com/udacity-aind/dog-project/dogImages.zip>
9. Human face dataset, <https://s3-us-west-1.amazonaws.com/udacity-aind/dog-project/lfw.zip>
10. D. P. Kingman and J. L. Ba, "Adam: A Method for Stochastic Optimization", <https://arxiv.org/pdf/1412.6980.pdf>
11. Adam – latest trends in deep learning optimization, <https://towardsdatascience.com/adam-latest-trends-in-deep-learning-optimization-6be9a291375c>