
Structural Compression of ResNet-like Convolutional Neural Networks

Nicholas Kastanos (nk569), Queens' College
Department of Computer Science and Technology
University of Cambridge
Cambridge, CB3 0FD
nk569@cam.ac.uk

Abstract

abstract

1 Introduction

2 Related Work

2.1 ResNet

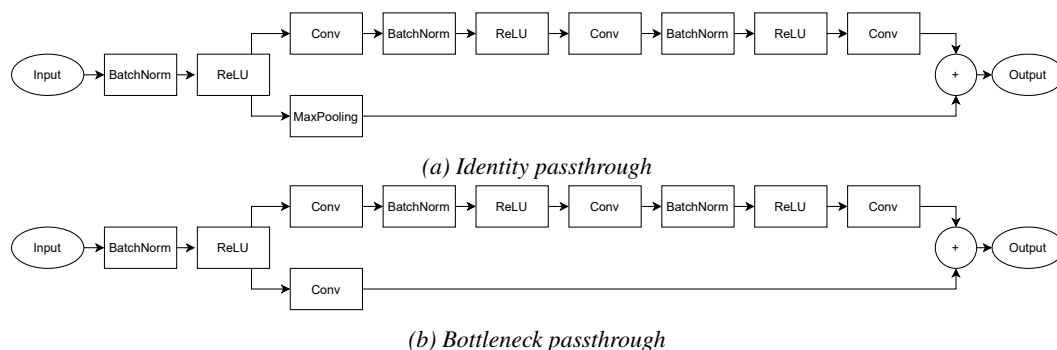
The residual block first postulated for use in ResNet has become a common-place feature in many subsequent networks. It was designed to combat network degradation in very deep networks, where they were learning more complicated features than were required. The skip connections introduced by the residual block allow features to bypass convolutional blocks, while also being affected by them. This can be summarised in the formula $g(x) = f(x) + x$ where x is the input, $g(x)$ is the desired output, and $f(x)$ is the input affected by the convolutional block.

This architecture was further developed until the form presented in ResNetV2, which includes two variants: identity passthrough which applies pooling to the skip connection x ; and bottleneck passthrough, which includes a single convolution in the skip connection. These two residual blocks can be seen in Figure 1.

2.2 Separable convolutions

Convolution layers contain a vast majority of the parameters in modern CNNs. By targeting parameter reductions to these layers, the compression can be spread throughout the network. Separable

L46 - Principles of Machine Learning Systems (2020)



convolutions reduce the number of parameters by decomposing the convolution into multiple stages through spatially and depthwise separable convolutions. While these convolutions reduce the memory and computation requirements of the system, the reduced parameters limits the number of possible kernels explored in training, and the resulting network may be suboptimal.

2.2.1 Spatially separable convolutions

A convolution kernel can be decomposed on its 2D spatial axis, i.e. height and width. Conceptually, the $n \times n$ kernel can be separated into two smaller kernels, a $n \times 1$ followed by a $1 \times n$ kernel. These kernels can be applied in sequential convolutions to obtain the same output shape as the single convolution. These decomposed kernels scale the parameters required by the convolution by a factor $P_s(n)$ (see Equation 1).

Similarly, the multiplication «MAC» operations of a spatially separated convolution are reduced. For a $M \times M$ input convolved with a $n \times n$ kernel, the number of multiplications are reduced by a factor of $M_s(n)$ (see Equation 2).

Equations 1) and 2 show that spatially separable convolutions show computational benefits when $n > 2$.

$$P_s(n) = \frac{2}{n} \quad (1) \quad \Rightarrow M_s(M, n) = \frac{2}{n} + \frac{2}{n(M-2)} \quad (2)$$

$$\Rightarrow M_s(n) = \frac{2}{n}, \text{ where } M \gg n$$

2.2.2 Depthwise separable convolutions

Depthwise separable convolutions separate the spatial convolution from the depth of the filters. This is accomplished by an initial depthwise convolution, followed by a pointwise convolution. The initial depthwise convolution separates the channels of the input and kernel, and convolves them independently. The pointwise convolution is a $N_F \times 1 \times 1 \times N_C$ convolution where N_F and N_C are the number of filters and channels respectively.

The number of parameters $P_d(n)$ and multiplications $M_d(n)$ are reduced by the same factor, which can be seen in Equation 3. Many CNNs have $N_F \gg 1$, therefore depthwise convolutions show compression kernel sizes greater than 1.

$$P_d(n) = M_d(n) = \frac{1}{n^2} + \frac{1}{N_F} \approx \frac{1}{n^2} \quad (3)$$

2.3 Quantization and datatype compression

Many resource-constrained devices do not have sufficient memory to use large neural networks, or may not have access to floating-point arithmetic units. Both of these factors can be mitigated by using low-precision integer datatypes, such as 8-bit integers. This effectively reduces the memory required for each parameter by $1/4$.

However, by reducing the precision of the datatype, the learned parameters are not represented fully, reducing the performance of the network. This has been shown to have significant effect on the network when reducing 32-bit floating point numbers to 8-bit integers, however quantization aware training can limit the impact.

TensorFlow, by default, uses 32-bit floating point precision for its network layers and training. TensorFlow Lite provides functionality to convert trained models to quantized, 8-bit integer models, however it does not have native 8-bit implementations for all layer types. In the event an incompatible layer is used, the engine upscales the activations to the full precision for that layer.

3 Methodology

In order to assess the benefits and costs of the compression, a baseline architecture is established and used as a comparison reference point. The ResNet50V2 image classification CNN-based neural

network is used as the baseline architecture. The baseline architecture and subsequent networks are trained using the same parameters using TensorFlow.

3.1 Training

The CNNs are trained using the CIFAR-10 dataset, with random horizontal flips, rotations, and cropping. The Adam optimizer is used with a learning rate of 0.001, and categorical crossentropy loss. The networks are trained for 30 epochs with a batch size of 32. The CIFAR-10’s training split is separated into train and validation splits. The splits are not shuffled to ensure consistency across models, and the validation split contains 10% of the total training data. This keeps the dataset’s test set free to evaluate the models on new data.

3.2 Modified residual blocks

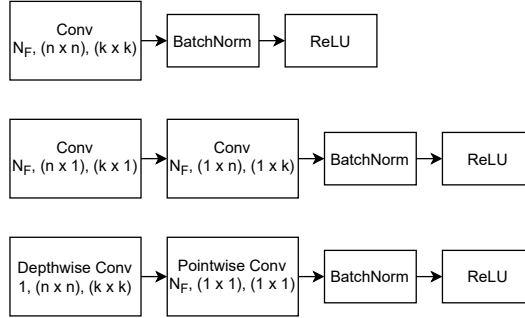


Figure 2: Example of replaced convolutions with separable convolutions with N_F filters, kernel size n , and stride k . From top to bottom: Traditional Convolution; Spatially Separable Convolution; Depthwise Separable Convolution

In order to modify the ResNetV2 network, the convolutions in the residual blocks are chosen to be substituted with separable convolutions. Additional batch normalization and activation layers are not added between these new layers. An example of this expansion can be seen in Figure 2. Only convolutions which would provide parameter and compute reduction are replaced with separable convolutions. For ResNet50V2, these are any convolutions where the kernel size $n > 1$.

These two additional CNNs are called ResNet50V2-Spatial and ResNet50V2-Depthwise, where convolutional layers are replaced with spatial and depthwise convolutions respectively.

3.3 Post-training quantization and compression

The fully trained models are compressed to 8-bit precision using TensorFlow Lite. The framework recommends the provision of the representative dataset to enhance the optimization process. This allows for quantization of both weights and activations for quantizable operations. The representative dataset is compiled using 1000 samples from the training dataset.

4 Evaluation and results

Table 1: Results

Name	Loss	Accuracy (%)	Compression Factor
ResNet50V2	0	100	1
ResNet50V2-Spatial	0	100	2000
ResNet50V2-Depthwise	0	100	2000
Compressed ResNet50V2	0	100	20000
Compressed ResNet50V2-Spatial	0	100	2000
Compressed ResNet50V2-Depthwise	0	100	2000

5 Conclusion