

Deep Residual Learning in the JPEG Transform Domain

Max Ehrlich and Larry Davis

maxehr@umiacs.umd.edu lsd@umiacs.umd.edu

University of Maryland, College Park, MD, USA.

Abstract

We introduce a general method of performing Residual Network inference and learning in the JPEG transform domain that allows the network to consume compressed images as input. Our formulation leverages the linearity of the JPEG transform to redefine convolution and batch normalization with a tune-able numerical approximation for ReLu. The result is mathematically equivalent to the spatial domain network up to the ReLu approximation accuracy. A formulation for image classification and a model conversion algorithm for spatial domain networks are given as examples of the method. We show skipping the costly decompression step allows for faster processing of images with little to no penalty in the network accuracy.

1. Introduction

The popularization of deep learning since the 2012 AlexNet [15] architecture has led to unprecedented gains for the field. Many applications that were once academic are now seeing widespread use of machine learning with success. Although the performance of deep neural networks far exceeds classical methods, there are still some major problems with the algorithms from a computational standpoint. Deep networks require massive amounts of data to learn effectively, especially for complex problems [18]. Further, the computational and memory demands of deep networks mean that for many large problems, only large institutions with GPU clusters can afford to train from scratch, leaving the average scientist to fine tune pre-trained weights.

This problem has been addressed many times in the literature. Batch normalization [12] is ubiquitous in modern networks to accelerate their convergence. Residual learning [11] allows for much deeper networks to learn effective mappings without overfitting. Techniques such as pruning and weight compression [9] are becoming more commonplace. As problems become even larger and more complex, these techniques are increasingly being relied upon for efficient training and inference.

We approach this problem at the level of the image rep-

resentation. JPEG is the most widespread image file format. Traditionally, the first step in using JPEGs for machine learning is to decompress them. We propose to skip this step and instead reformulate the ResNet architecture to perform its operations directly on compressed images. The goal is to produce a new network that is mathematically equivalent to the spatial domain network, but which operates on compressed images by including the compression transform into the network weights, which can be done because they are both linear maps. Because the ReLu function is non-linear, we develop an approximation technique for it. This is a general method and, to our knowledge, is the first attempt at formulating a piecewise linear function in the transform domain.

The contributions of this work are as follows

1. The general method for expressing convolutional networks in the JPEG domain
2. Concrete formulation for residual blocks to perform classification
3. A model conversion algorithm to apply pretrained spatial domain networks to JPEG images
4. Approximated Spatial Masking: the first general technique for application of piecewise linear functions in the transform domain

By skipping the decompression step and by operating on the compressed format, we show a notable increase in speed for testing and a marginal speed for training.

2. Prior Work

We review prior work separated into three categories: compressed domain operations, machine learning in the compressed domain, and deep learning in the compressed domain.

2.1. Compressed Domain Operations

The expression of common operations in the compressed domain was an extremely active area of study in the late 80s and early 90s, motivated by the lack of computing power to

quickly decompress, process, and recompress images and video. For JPEG, Smith and Rowe [25] formulate fast JPEG compatible algorithms for performing scalar and pixelwise addition and multiplication. This was extended by Shen and Sethi [23] to general blockwise operations and by Smith [24] to arbitrary linear maps. Natarajan and Vasudev [19] additionally formulate an extremely fast approximate algorithm for scaling JPEG images. For MPEG, Chang *et al.* [2] introduce the basic algorithms for manipulating compressed video. Chang and Messerschmitt [3] give a fast algorithm for decoding motion compensation before DCT which allows arbitrary video compositing operations to be performed.

2.2. Machine Learning in the Compressed Domain

Compressed domain machine learning grew out of the work in the mid 90s. Arman *et al.* [1] give the basic framework for image processing of compressed images. Feng and Jiang [5] show how image retrieval can be performed directly on compressed JPEGs. He *et al.* [10] extend their work with a hypothesis testing technique. Wu *et al.* [30] formulate the popular SIFT feature extraction in the DCT domain.

2.3. Deep Learning in the Compressed Domain

Because deep networks are non-linear maps, deep learning has received limited study in the compressed domain. Ghosh and Chellappa [7] use a DCT as part of their network’s first layer and show that it speeds up convergence for training. This is extended by Ulicny *et al.* [26] to create separate filters for each DCT basis function. Wu *et al.* [29] formulate a deep network for video action recognition that uses a separate network for i-frames and p-frames. Since the p-frame network functions on raw motion vectors and error residuals it is considered compressed domain processing, although it works in the spatial domain and not the quantized frequency domain as in this work. Wu *et al.* show a significant efficiency advantage compared to traditional 3D convolution architectures, which they attribute to the p-frame data being a minimal representation of the video motion. Gueguen *et al.* [8] formulate a traditional ResNet that operates on DCT coefficients directly instead of pixels, *e.g.* the DCT coefficients are fed to the network. They show that learning converges faster on this input, further motivating the JPEG representation.

3. Background

We briefly review the JPEG compression/decompression algorithm [27] and introduce the multilinear method that we use to formulate our networks [24].

3.1. JPEG Compression

The JPEG compression algorithm is defined as the following steps.

1. Divide the image into 8×8 blocks
2. Compute the 2D forward Discrete Cosine Transform (DCT Type 2) of each block
3. Linearize the blocks using a zigzag order to produce a 64 component vector
4. Element-wise divide each vector by a quantization coefficient
5. Round the the vector elements to the nearest integer
6. Run-length code and entropy code the vectors

This process is repeated independently for each image plane. In most cases, the original image is transformed from the RGB color space to YUV and chroma subsampling is applied since the human visual system is less sensitive to small color changes than to small brightness changes [28]. The decompression algorithm is the inverse process. Note that the rounding step (step 5) must be skipped during decompression. This is the step in JPEG compression where information is lost and is the cause of artifacts in decompressed JPEG images.

The magnitude of the information loss can be tuned using the quantization coefficients. If a larger coefficient is applied in step 4, then the result will be closer to 0 which increases its likelihood of being dropped altogether during rounding. In this way, the JPEG transform forces sparsity on the representation, which is why it compresses image data so well. This is coupled with the tendency of the DCT to push the magnitude of the coefficients into the upper left corner (the DC coefficient and the lowest spatial frequency) resulting in high spatial frequencies being dropped. Not only do these high spatial frequencies contribute less response to the human visual system, but they are also the optimal set to drop for a least squares reconstruction of the original image:

Theorem 1 (DCT Least Squares Approximation Theorem). *Given a set of N samples of a signal $X = \{x_0, \dots, x_N\}$, let $Y = \{y_0, \dots, y_N\}$ be the DCT coefficients of X . Then, for any $1 \leq m \leq N$, the approximation*

$$p_m(t) = \frac{1}{\sqrt{n}}y_0 + \sqrt{\frac{2}{n}} \sum_{k=1}^m y_k \cos\left(\frac{k(2t+1)\pi}{2n}\right) \quad (1)$$

of X minimizes the least squared error

$$e_m = \sum_{i=0}^n (p_m(i) - x_i)^2 \quad (2)$$

Theorem 1 states that a reconstruction using the m lowest spatial frequencies is optimal with respect to any other set of m spatial frequencies. Proof of Theorem 1 is given in the supplementary material.

3.2. JPEG Linear Map

A key observation of the JPEG algorithm, and the foundation of most compressed domain processing methods [2, 3, 19, 23, 22, 21, 25, 24] is that steps 1-4 of the JPEG compression algorithm are linear maps, so they can be composed, along with other linear operations, into a single linear map which performs the operations on the compressed representation. Step 5, the rounding step, is irreversible and ignored by decompression. Step 6, the entropy coding, is a nonlinear map and its form is computed from the data directly, so it is difficult to work with this representation. We define the JPEG Transform Domain as the output of Step 4 in the JPEG encoding algorithm. This is a standard convention of compressed domain processing. Inputs to the algorithms described here will be JPEGs after reversing the entropy coding.

Formally, we model a single plane image as the type (0, 2) tensor $I \in H^* \otimes W^*$ where H and W are vector spaces and $*$ denotes the dual space. We always use the standard orthonormal basis for these vector spaces which allows the free raising and lowering of indices without the use of a metric tensor. We define the JPEG transform as the type (2, 3) tensor $J \in H \otimes W \otimes X^* \otimes Y^* \otimes K^*$. J represents a linear map $J : H^* \otimes W^* \rightarrow X^* \otimes Y^* \otimes K^*$ and is computed as (in Einstein notation)

$$I'_{xyk} = J_{xyk}^{hw} I_{hw} \quad (3)$$

We say that I' is the representation of I in the JPEG transform domain. The indices h, w give pixel position, x, y give block position, and k gives the offset into the block.

The form of J is constructed from the JPEG compression steps listed in the previous section. Let the linear map $B : H^* \otimes W^* \rightarrow X^* \otimes Y^* \otimes M^* \otimes N^*$ be defined as

$$B_{xymn}^{hw} = \begin{cases} 1 & h, w \text{ belongs in block } x, y \text{ at offset } m, n \\ 0 & \text{otherwise} \end{cases} \quad (4)$$

then B can be used to break the image represented by I into blocks of a given size such that the first two indices x, y index the block position and the last two indices m, n index the offset into the block.

Next, let the linear map $D : M^* \otimes N^* \rightarrow A^* \otimes B^*$ be defined as

$$D_{\alpha\beta}^{mn} = \frac{1}{4} V(\alpha) V(\beta) \cos\left(\frac{(2m+1)\alpha\pi}{16}\right) \cos\left(\frac{(2n+1)\beta\pi}{16}\right) \quad (5)$$

where $V(u)$ is a normalizing scale factor. Then D represents the 2D discrete forward (and inverse) DCT. Let $Z : A^* \otimes B^* \rightarrow \Gamma^*$ be defined as

$$Z_{\gamma}^{\alpha\beta} = \begin{cases} 1 & \alpha, \beta \text{ is at } \gamma \text{ under zigzag ordering} \\ 0 & \text{otherwise} \end{cases} \quad (6)$$

then Z creates the zigzag ordered vectors. Finally, let $S : \Gamma^* \rightarrow K^*$ be

$$S_k^{\gamma} = \frac{1}{q_k} \quad (7)$$

where q_k is a quantization coefficient. This scales the vector entries by the quantization coefficients.

With linear maps for each step of the JPEG transform, we can then create the J tensor described at the beginning of this section

$$J_{xyk}^{hw} = B_{xymn}^{hw} D_{\alpha\beta}^{mn} Z_{\gamma}^{\alpha\beta} S_k^{\gamma} \quad (8)$$

The inverse mapping also exists as a tensor \tilde{J} which can be defined using the same linear maps with the exception of S . Let \tilde{S} be

$$\tilde{S}_{\gamma}^k = q_k \quad (9)$$

Then

$$\tilde{J}_{hw}^{xyk} = B_{xymn}^{xymn} D_{\alpha\beta}^{\alpha\beta} Z_{\gamma}^{\alpha\beta} \tilde{S}_{\gamma}^k \quad (10)$$

Next consider a linear map $C : H^* \otimes W^* \rightarrow H^* \otimes W^*$ which performs an arbitrary pixel manipulation on an image plane I . To apply this mapping to a JPEG image I' , we first decompress the image, apply C to the result, then compress that result to get the final JPEG. Since compressing is an application of J and decompressing is an application of \tilde{J} , we can form a new linear map $\Xi : X^* \otimes Y^* \otimes K^* \rightarrow X^* \otimes Y^* \otimes K^*$ as

$$\Xi_{x'y'k'}^{xyk} = \tilde{J}_{hw}^{xyk} C_{h'w'}^{hw} J_{x'y'k'}^{h'w'} \quad (11)$$

which applies C in the JPEG transform domain. There are two important points to note about Ξ . The first is that, although it encapsulates decompression, applying C and compressing, it uses far fewer operations than doing these processes separately since the coefficients are multiplied out. The second is that it is mathematically equivalent to performing C on the decompressed image and compressing the result. It is not an approximation.

4. JPEG Domain Residual Networks

The ResNet architecture, consists of blocks of four basic operations: Convolution (potentially strided), ReLU, Batch Normalization, and Component-wise addition, with

the blocks terminating with a global average pooling operation [11] before a fully connected layer performs the final classification. Our goal will be to develop JPEG domain equivalents to these five operations. Network activations are given as a single tensor holding a batch of multi-channel images, that is $I \in N^* \otimes P^* \otimes H^* \otimes W^*$.

4.1. Convolution

The convolution operation follows directly from the discussion in Section 3.2. The convolution operation is a shorthand notation for a linear map $C : N^* \otimes P^* \otimes H^* \otimes W^* \rightarrow N^* \otimes P^* \otimes H^* \otimes W^*$. Since the same operation is applied to each image in the batch, we can represent C with a type (3, 3) tensor. The entries of this tensor give the coefficient for a given pixel in a given input channel for each pixel in each output channel. We now develop the algorithm for representing discrete convolutional filters using this data structure.

A naive algorithm can simply copy randomly initialized convolution weights into this larger structure following the formula for convolution and then apply the JPEG compression and decompression tensors to the result. However, this is difficult to parallelize and incurs additional memory overhead to store the spatial domain operation. A more efficient algorithm would produce the JPEG domain operation directly and be easy to express as a compute kernel for a GPU. Start by considering the JPEG decompression tensor \tilde{J} . Note that since $\tilde{J} \in X \otimes Y \otimes K \otimes H^* \otimes W^*$ the last two indices of \tilde{J} form single channel image under our image model (e.g. the last two indices are in $H^* \otimes W^*$). If the convolution can be applied to this "image", then the resulting map would decompress and convolve simultaneously. We can formulate a new tensor $\hat{J} \in N \otimes H^* \otimes W^*$ by reshaping \tilde{J} and treating this as a batch of images¹. Then, given randomly initialized filter weights, K computing

$$\hat{C}^b = K \star \hat{J}^b \tag{12}$$

where \star indicates the convolution operation and \hat{J}^b indexes \hat{J} in the batch dimension, gives us the desired map. After reshaping \hat{C} back to the original shape of \tilde{J} to give \tilde{C} , the full compressed domain operation can be expressed as

$$\Xi_{p'x'y'k'}^{pxyk} = \tilde{C}_{p'hw}^{pxyk} J_{x'y'k'}^{hw} \tag{13}$$

where p and p' index the input and output channels of the image respectively. This algorithm skips the overhead of computing the spatial domain map explicitly and depends only on the batch convolution operation which is available in all GPU accelerated deep learning libraries. Further, the

¹Consider as a concrete example using 32×32 images. Then \tilde{J} is of shape $4 \times 4 \times 64 \times 32 \times 32$ and the described reshaping gives \hat{J} of shape $1024 \times 1 \times 32 \times 32$ which can be treated as a batch of size 1024 of 32×32 images for convolution.

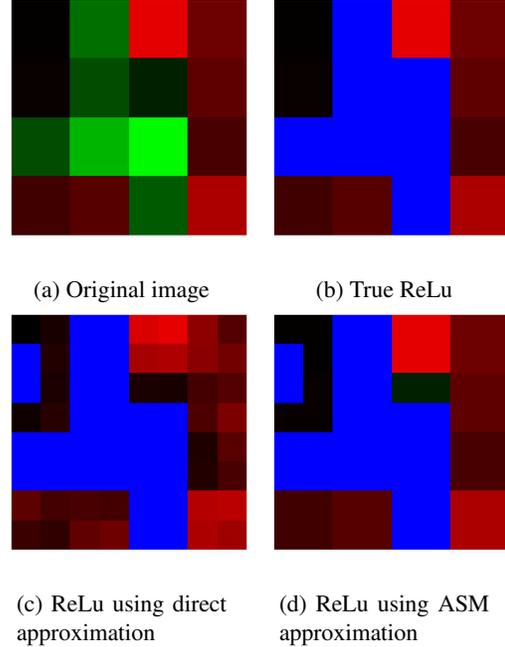


Figure 1: Example of ASM ReLU on an 8×8 block. Green pixels are negative, red pixels are positive, and blue pixels are zero. 6 spatial frequencies are used for both approximations. Note that the direct approximation fails to preserve positive pixel values.

map can be precomputed to speed up inference by avoiding repeated applications of the convolution. At training time, the gradient of the compression and decompression operators is computed and used to find the gradient of the original convolution filter with respect to the previous layers error, then the map Ξ is updated using the new filter. So, while inference efficiency of the convolution operation is greatly improved, training efficiency is limited by the more complex update. We show in Section 5.4 that the training throughput is still higher than the equivalent spatial domain model.

4.2. ReLU

Computing ReLU in the JPEG domain is not as straightforward since ReLU is a non-linear function. Recall that the ReLU function is given by

$$r(x) = \begin{cases} x & x > 0 \\ 0 & x \leq 0 \end{cases} \tag{14}$$

We begin by defining the ReLU in the DCT domain and show how it can be trivially extended to the JPEG transform domain. To do this, we develop a general approximation technique called Approximated Spatial Masking that can apply any piecewise linear function to JPEG compressed

images.

To develop this technique we must balance two seemingly competing criteria. The first is that we want to use the JPEG transform domain, since it has a computational advantage over the spatial domain. The second is that we want to compute a non-linear function which is incompatible with the JPEG transform. Can we balance these two constraints by sacrificing a third criterion? Consider an approximation of the spatial domain image that uses only a subset of the DCT coefficients. Computing this is fast, since it does not use the full set of coefficients, and gives us a spatial domain representation which is compatible with the non-linearity. What we sacrifice is accuracy. The accuracy-speed tradeoff is tunable to the problem by changing the size of the set of coefficients.

By Theorem 1 we use the lowest m frequencies for an optimal reconstruction. For the 8×8 DCT used in the JPEG algorithm, this gives 15 spatial frequencies total (numbered 0 to 14). We can then fix a maximum number of spatial frequencies k and use all coefficients ϕ such that $\phi \leq k$ as our approximation.

If we now compute the piecewise linear function on this approximation directly there are two major problems. The first is that, although the form of the approximation is motivated by a least squares minimization, it is by no means guaranteed to reproduce the original values of *any* of the pixels. The second is that this gives the value of the function in the spatial domain, and to continue using a JPEG domain network we would need to compress the result which adds computational overhead.

To solve the first problem we examine the intervals that the linear pieces fall into. The larger these intervals are, the more likely we are to have produced a value in the correct interval ² in our approximation. Further, since the lowest k frequencies minimize the least squared error, the higher the frequency, the less likely it is to push a pixel value out of the correct range. With this motivation, we can produce a binary mask for each piece of the function. The linear pieces can then be applied directly to the DCT coefficients, and then multiplied by their masks and summed to give the final result. This preserves all pixel values. The only errors would be in the mask which would result in the wrong linear piece being applied. This is the fundamental idea behind the Approximated Spatial Masking (ASM) technique.

The final problem is that we now have a mask in the spatial domain, but the original image is in the DCT domain. There is a well known algorithm for pixelwise multiplication of two DCT images [25], but it would require the mask to also be in the DCT domain. Fortunately, there is a straightforward solution that is a result of the multilinear

²For example if the original pixel value was 0.7 and the approximate value is 0.5, then the approximation is in the correct interval for ReLU (≥ 0) but its value is incorrect.

analysis given in Section 3.2. Consider the bilinear map

$$H : A^* \otimes B^* \times M^* \otimes N^* \rightarrow A^* \otimes B^* \quad (15)$$

that takes a DCT block, F , and a spatial mask G , and produces the masked DCT block by pixelwise multiplication. Our task will be to derive the form of H . We proceed by construction. The steps of such an algorithm naively would be

1. Take the inverse DCT of F : $I_{mn} = D_{mn}^{\alpha\beta} F_{\alpha\beta}$
2. Pixelwise multiply: $I'_{mn} = I_{mn} G_{mn}$
3. Take the DCT of I' : $F'_{\alpha'\beta'} = D_{\alpha'\beta'}^{mn} I'_{mn}$.

Since these three steps are linear or bilinear maps, they can be combined

$$F'_{\alpha'\beta'} = F^{\alpha\beta} [D_{\alpha\beta}^{mn} D_{\alpha'\beta'}^{mn}] G_{mn} \quad (16)$$

Giving the final bilinear map H as

$$H_{\alpha'\beta'}^{\alpha\beta mn} = D^{\alpha\beta mn} D_{\alpha'\beta'}^{mn} \quad (17)$$

We call H the Harmonic Mixing Tensor since it gives all the spatial frequency permutations that we need. H can be precomputed to speed up computation.

To use this technique to compute the ReLU function, consider this alternative formulation

$$\text{nnm}(x) = \begin{cases} 1 & x > 0 \\ 0 & x \leq 0 \end{cases} \quad (18)$$

We call the function $\text{nnm}(x)$ the nonnegative mask of x . This is our binary mask for ASM. We express the ReLU function as

$$r(x) = \text{nnm}(x)x \quad (19)$$

This new function can be computed efficiently from fewer spatial frequencies with much higher accuracy since only the sign of the original function needs to be correct. Figure 1 gives an example of this algorithm on a random block and compares it to computing ReLU on the approximation directly. Note that in the ASM image the pixel values of all positive pixels are preserved, the only errors are in the mask. In the direct approximation, however, none of the pixel values are preserved and it suffers from masking errors. The magnitude of the error is tested in Section 5.3 and pseudocode for the ASM algorithm is given in the supplementary material.

To extend this method from the DCT domain to the JPEG transform domain, the rest of the missing JPEG tensor can simply be applied as follows:

$$H_{k'}^{kmn} = Z_{\gamma'}^k \tilde{S}_{\alpha\beta}^{\gamma} D^{\alpha\beta mn} D_{\alpha'\beta'}^{mn} S_{\gamma'}^{\alpha'\beta'} Z_{k'}^{\gamma'} \quad (20)$$

Since the operation is the same for each block, and there are no interactions between blocks, the blocking tensor B can be skipped.

4.3. Batch Normalization

Batch normalization [12] has a simple and efficient formulation in the JPEG domain. Recall that batch normalization defines two learnable parameters: γ and β . A given feature map I is first centered and then normalized over the batch, then scaled by γ and translated by β . The full formula is

$$\text{BN}(I) = \gamma \frac{I - \mathbb{E}[I]}{\sqrt{\text{Var}[I]}} + \beta \quad (21)$$

So to define the batch normalization operation in the JPEG domain, we need four parts: the mean, the variance, scalar multiplication and scalar addition. Again, we first derive the result in the DCT domain and trivially extend to the JPEG transform domain.

We start with the sample mean. Observe, from the definition of the DCT, the first DCT coefficient is given by

$$D_{00} = \frac{1}{2\sqrt{2N}} \sum_{x=0}^N \sum_{y=0}^N I_{xy} \quad (22)$$

In other words, the (0,0) DCT coefficient is proportional to the mean of the block. Further, since the DCT basis is orthonormal, we can be sure that the remaining DCT coefficients do not depend on the mean. This means that to center the image we need only set the (0,0) DCT coefficient to 0. For tracking the running mean, we simply read this value. Note that this is a much more efficient operation than the mean computation in the spatial domain.

Next, to get the variance, we use the following theorem:

Theorem 2 (The DCT Mean-Variance Theorem). *Given a set of samples of a signal X such that $\mathbb{E}[X] = 0$, let Y be the DCT coefficients of X . Then*

$$\text{Var}[X] = \mathbb{E}[Y^2] \quad (23)$$

Intuitively this makes sense because the (0,0) coefficient represents the mean, the remaining DCT coefficients are essentially spatial oscillations around the mean, which should define the variance. Proof of this theorem is given in the supplementary material.

To apply γ and the variance, we use scalar multiplication. Since JPEG is linear, this is unchanged

$$J(\gamma I) = \gamma J(I) \quad (24)$$

For scalar addition to apply β , note that since the (0,0) coefficient is the mean, and adding β to every pixel in the image is equivalent to raising the mean by β , we can simply add β to the (0,0) coefficient.

To extend this to JPEG is simple. The proportionality constant for the (0,0) coefficient is $\frac{1}{2\sqrt{2 \times 8}} = \frac{1}{8}$. For this reason, many quantization matrices use 8 as the (0,0) quantization coefficient. This means that the 0th block entry for

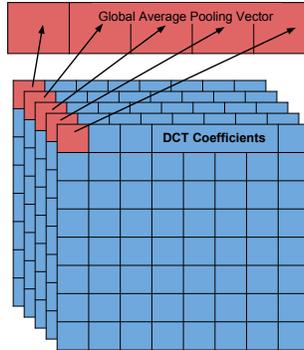


Figure 2: Global average pooling. The 0th coefficient of each block can be used directly with no computation.

a block does not need any proportionality constant, it stores exactly the mean. So for adding β , we can simply set the 0th position to β without performing additional operations. The other operations are unaffected.

4.4. Component-wise Addition

Component-wise addition is the simplest formulation in our network. This is a well known result detailed in [2, 21, 23, 25] among others. Since the JPEG transform, J , is a linear map, for two images F and G , we have

$$J(F + G) = J(F) + J(G) \quad (25)$$

meaning that we can simply perform a component-wise addition of the JPEG compressed results with no need for further processing.

4.5. Global Average Pooling

Global average pooling also has a simple formulation in JPEG domain. Recall from the discussion of Batch Normalization (Section 4.3) that the 0th element of the block after quantization is equal to the mean of the block. Then this element can be extracted channel-wise from each block and the global average pooling result is the channel-wise mean of these elements.

Furthermore, our network architecture for classification will always reduce the input images to a single block, which can then have its mean extracted and reported as the global average pooling result directly. Note the efficiency of this process: rather than channel-wise averaging in a spatial domain network, we simply have an unconditional read operation, one per channel. This is illustrated in Figure 2.

4.6. Model Conversion

The previous sections described how to build the ResNet component operations in the JPEG domain. While this implies straightforward algorithms for both inference and learning on JPEGs, we can also convert pre-trained models

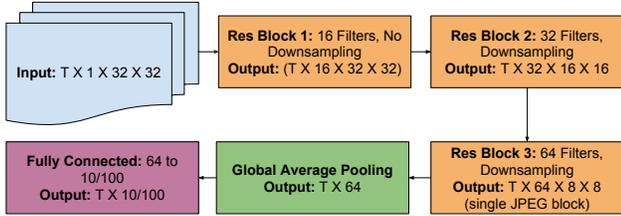


Figure 3: Simple network architecture. T indicates the batch size.

for JPEG inference. This allows any model that was trained on spatial domain images to benefit from our algorithms at inference time. Consider Equation 12. In this equation, K holds the randomly initialized convolution filter. By instead using pretrained spatial weights for K , the convolution will work as expected on JPEGs. Similarly, pretrained $\alpha, \beta, \mu, \sigma$ for batch normalization can be provided. By doing this for each layer in a pretrained network, the network will operate on JPEGs. The only caveat is that the ReLU approximation accuracy can effect the final performance of the network since the weights were not trained to cope with it. This is tested in Section 5.3.

5. Experiments

We give experimental evidence for the efficacy of our method, starting with a discussion of the architectures we use and the datasets. We use model conversion as a sanity check, ensuring that the JPEG model with exact ReLU matches exactly the testing accuracy of a spatial domain model. Next we show how the ReLU approximation accuracy effects overall network performance. We conclude by showing the training and testing time advantage of our method.

5.1. Network Architectures and Datasets

Since we are concerned with reproducing the inference results of spatial domain networks, we choose the MNIST [16] and CIFAR-10/100 [14] datasets since they are easy to work with. The MNIST images are padded to 32×32 to ensure an even number of JPEG blocks. Our network architecture is shown in Figure 3. The classification network consists of three residual blocks with the final two performing downsampling so that the final feature map consists of a single JPEG block. The goal of this architecture is not to get high accuracy, but rather to serve as a point of comparison for the spatial and JPEG algorithms.

5.2. Model Conversion

For this first experiment, we show empirically that the JPEG formulation is mathematically equivalent to the spatial domain network. To show this, we train 100 spatial

domain models on each of the three datasets and give their mean testing accuracies. We then use model conversion to transform the pretrained models to the JPEG domain and give the mean testing accuracies of the JPEG models. The images are losslessly JPEG compressed for input to the JPEG networks and the exact (15 spatial frequency) ReLU formulation is used. The result of this test is given in Table 1. Since the accuracy difference between the networks is extremely small, the deviation is also included.

Dataset	Spatial	JPEG	Deviation
MNIST	0.988	0.988	2.999e-06
CIFAR10	0.725	0.725	9e-06
CIFAR100	0.385	0.385	1e-06

Table 1: Model conversion accuracies. Spatial and JPEG testing accuracies are the same to within floating point error.

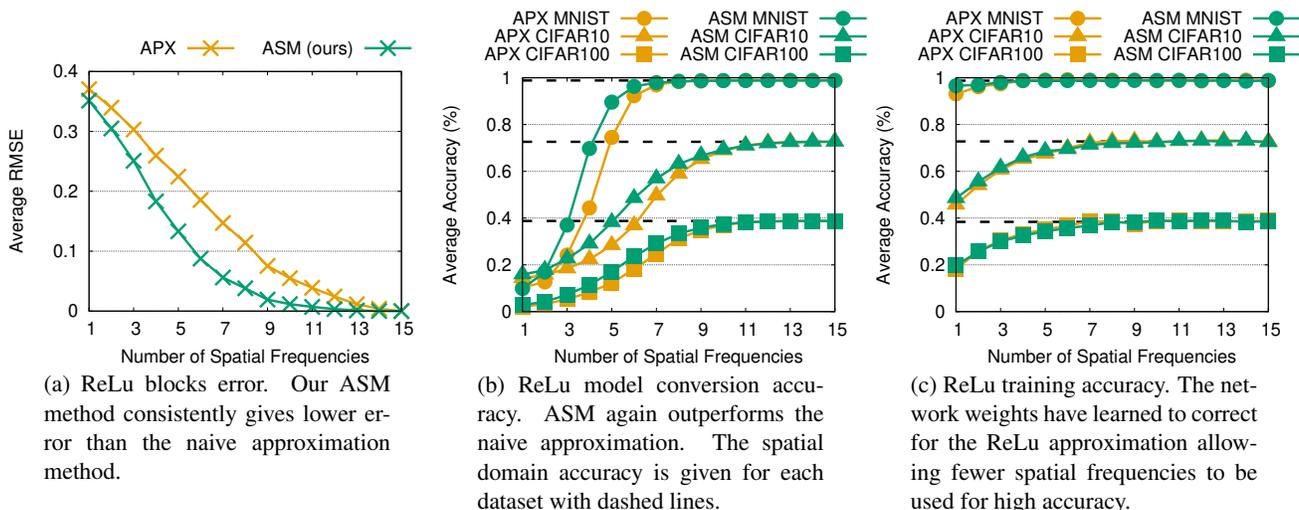
5.3. ReLU Approximation Accuracy

Next, we examine the impact of the ReLU approximation. We start by examining the raw error on individual 8×8 blocks. For this test, we take random 4×4 pixel blocks in the range $[-1, 1]$ and scale them to 8×8 using a box filter. Fully random 8×8 blocks do not accurately represent the statistics of real images and are known to be a worst case for the DCT transform. The 4×4 blocks allow for a large random sample size while still approximating real image statistics. We take 10 million blocks and compute the average RMSE of our ASM technique and compare it to computing ReLU directly on the approximation (APX). This test is repeated for all one to fifteen spatial frequencies. The result, shown in Figure 4a shows that our ASM method gives a better approximation (lower RMSE) through the range of spatial frequencies.

This test provides a strong motivation for the ASM method, so we move on to testing it in the model conversion setting. For this test, we again train 100 spatial domain models and then perform model conversion with the ReLU layers ranging from 1-15 spatial frequencies. We again compare our ASM method with the APX method. The result is given in Figure 4b. Again the ASM method outperforms the APX method.

As a final test, we show that if the models are trained in the JPEG domain, the CNN weights will actually learn to cope with the approximation and fewer spatial frequencies are required for good accuracy. The result in Figure 4c shows that the ASM method again outperforms the APX method and that the network weights have learned to cope with the approximation.

Figure 4: ReLu accuracy results.



5.4. Efficiency of Training and Testing

Finally, we show the throughput for training and testing. For this we test on all three datasets by training and testing a spatial model and training and testing a JPEG model and measuring the time taken. This is then converted to an average throughput measurement. The experiment is performed on an NVIDIA Pascal GPU with a batch size of 40 images. The results, shown in Figure 5, show that the JPEG model is able to outperform the spatial model in all cases, but that the performance on training is still limited. This is caused by the more complex gradient created by the convolution and ReLu operations. At inference time, however, performance

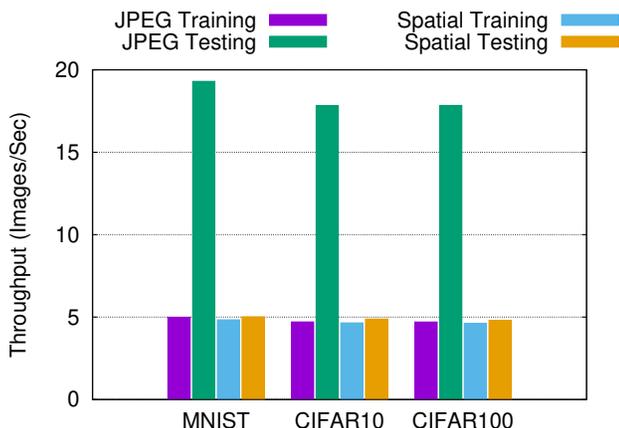


Figure 5: Throughput. The JPEG model has a more complex gradient which limits speed improvement during training. Inference, however, sees considerably higher throughput.

is greatly improved over the spatial model.

6. Conclusion and Future Work

In this work we showed how to formulate deep residual learning in the JPEG transform domain, and that it provides a notable performance benefit in terms of processing time per image. Our method expresses convolutions as linear maps [24] and introduces a novel approximation technique for ReLu. We showed that the approximation can achieve highly performant results with little impact on classification accuracy.

Future work should focus on two main points. The first is efficiency of representation. Our linear maps take up more space than spatial domain convolutions. This makes it hard to scale the networks to datasets with large image sizes. Secondly, library support in commodity deep learning libraries for some of the features required by this algorithm are lacking. As of this writing, true sparse tensor support is missing in all of PyTorch [20], TensorFlow [17], and Caffe [13], with these tensors being represented as coordinate lists which are known to be highly non-performant. Additionally, the `einsum` function for evaluating multilinear expressions is not fully optimized in these libraries when compared to the speed of convolutions in libraries like CuDNN [4], though we make use of the `opt_einsum` [6] tool to partially mitigate this.

7. Acknowledgment

This research was partially funded by Facebook AI Research. We especially thank Dr. Ser-Nam Lim and his team at Facebook for their continued support of our work.

References

- [1] Farshid Arman, Arding Hsu, and Ming-Yee Chiu. “Image processing on compressed data for large video databases”. In: *Proceedings of the first ACM international conference on Multimedia*. ACM. 1993, pp. 267–272.
- [2] Shih-Fu Chang. “Video Compositing in the DCT domain”. In: *IEEE Workshop on Visual Signal Processing and Communications, Raleigh, NC, Sep. 1992*. 1992.
- [3] Shih-Fu Chang and David G Messerschmitt. “A new approach to decoding and compositing motion-compensated DCT-based images”. In: *icassp*. IEEE. 1993, 421–424vol.
- [4] Sharan Chetlur et al. “cudnn: Efficient primitives for deep learning”. In: *arXiv preprint arXiv:1410.0759* (2014).
- [5] Guocan Feng and Jianmin Jiang. “JPEG image retrieval based on features from DCT domain”. In: *International Conference on Image and Video Retrieval*. Springer. 2002, pp. 120–128.
- [6] Daniel G. A. Smith and Johnnie Gray. “opt_einsum - A Python package for optimizing contraction order for einsum-like expressions”. In: *Journal of Open Source Software* 3.26 (June 29, 2018), p. 753. ISSN: 2475-9066. DOI: 10.21105/joss.00753. URL: <http://dx.doi.org/10.21105/joss.00753>.
- [7] Arthita Ghosh and Rama Chellappa. “Deep feature extraction in the DCT domain”. In: *Pattern Recognition (ICPR), 2016 23rd International Conference on*. IEEE. 2016, pp. 3536–3541.
- [8] Lionel Gueguen et al. “Faster Neural Networks Straight from JPEG”. In: *International Conference on Learning Representations*. 2018.
- [9] Song Han, Huizi Mao, and William J Dally. “Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding”. In: *arXiv preprint arXiv:1510.00149* (2015).
- [10] Daan He, Zhenmei Gu, and Nick Cercone. “Efficient image retrieval in DCT domain by hypothesis testing”. In: *Image Processing (ICIP), 2009 16th IEEE International Conference on*. IEEE. 2009, pp. 225–228.
- [11] Kaiming He et al. “Deep residual learning for image recognition”. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2016, pp. 770–778.
- [12] Sergey Ioffe and Christian Szegedy. “Batch normalization: Accelerating deep network training by reducing internal covariate shift”. In: *arXiv preprint arXiv:1502.03167* (2015).
- [13] Yangqing Jia et al. “Caffe: Convolutional Architecture for Fast Feature Embedding”. In: *arXiv preprint arXiv:1408.5093* (2014).
- [14] Alex Krizhevsky and Geoffrey Hinton. *Learning multiple layers of features from tiny images*. Tech. rep. Citeseer, 2009.
- [15] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. “Imagenet classification with deep convolutional neural networks”. In: *Advances in neural information processing systems*. 2012, pp. 1097–1105.
- [16] Yann LeCun. “The MNIST database of handwritten digits”. In: <http://yann.lecun.com/exdb/mnist/> ().
- [17] Martín Abadi et al. *TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems*. Software available from tensorflow.org. 2015. URL: <http://tensorflow.org/>.
- [18] Maryam M Najafabadi et al. “Deep learning applications and challenges in big data analytics”. In: *Journal of Big Data* 2.1 (2015), p. 1.
- [19] Balas K Natarajan and Bhaskaran Vasudev. “A fast approximate algorithm for scaling down digital images in the DCT domain”. In: *Image Processing, 1995. Proceedings., International Conference on*. Vol. 2. IEEE. 1995, pp. 241–243.
- [20] Adam Paszke et al. “Automatic differentiation in PyTorch”. In: *NIPS-W*. 2017.
- [21] Bo Shen and Ishwar K Sethi. “Block-based manipulations on transform-compressed images and videos”. In: *Multimedia Systems* 6.2 (1998), pp. 113–124.
- [22] Bo Shen and Ishwar K Sethi. “Direct feature extraction from compressed images”. In: *Storage and Retrieval for Still Image and Video Databases IV*. Vol. 2670. International Society for Optics and Photonics. 1996, pp. 404–415.
- [23] Bo Shen and Ishwar K Sethi. “Inner-block operations on compressed images”. In: *Proceedings of the third ACM international conference on Multimedia*. ACM. 1995, pp. 489–498.
- [24] Brian C Smith. “Fast software processing of motion JPEG video”. In: *Proceedings of the second ACM international conference on Multimedia*. ACM. 1994, pp. 77–88.

- [25] Brian C Smith and Lawrence A Rowe. “Algorithms for manipulating compressed images”. In: *IEEE Computer Graphics and Applications* 13.5 (1993), pp. 34–42.
- [26] Matej Ulicny, Vladimir A Krylov, and Rozenn Dahyot. “Harmonic Networks: Integrating Spectral Information into CNNs”. In: *arXiv preprint arXiv:1812.03205* (2018).
- [27] Gregory K Wallace. “The JPEG still picture compression standard”. In: *IEEE transactions on consumer electronics* 38.1 (1992), pp. xviii–xxxiv.
- [28] Stefan Winkler, Murat Kunt, and Christian J van den Branden Lambrecht. “Vision and video: models and applications”. In: *Vision Models and Applications to Image and Video Processing*. Springer, 2001, pp. 201–229.
- [29] Chao-Yuan Wu et al. “Compressed video action recognition”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2018, pp. 6026–6035.
- [30] Zhen Wu et al. “SIFT Feature Extraction Algorithm for Image in DCT Domain”. In: *Applied Mechanics and Materials*. Vol. 347. Trans Tech Publ. 2013, pp. 2963–2967.

Supplementary Material

1. Proof of the DCT Least Squares Approximation Theorem

Theorem 1 (DCT Least Squares Approximation Theorem). *Given a set of N samples of a signal $X = \{x_0, \dots, x_N\}$, let $Y = \{y_0, \dots, y_N\}$ be the DCT coefficients of X . Then, for any $1 \leq m \leq N$, the approximation*

$$p_m(t) = \frac{1}{\sqrt{n}} y_0 + \sqrt{\frac{2}{n}} \sum_{k=1}^m y_k \cos\left(\frac{k(2t+1)\pi}{2n}\right) \quad (1)$$

of X minimizes the least squared error

$$e_m = \sum_{i=0}^n (p_m(i) - x_i)^2 \quad (2)$$

Proof. First consider that since Equation 1 represents the Discrete Cosine Transform, which is a Linear map, we can write rewrite it as

$$D_m^T y = x \quad (3)$$

where D_m is formed from the first m rows of the DCT matrix, y is a row vector of the DCT coefficients, and x is a row vector of the original samples.

To solve for the least squares solution, we use the the normal equations, that is we solve

$$D_m D_m^T y = D_m x \quad (4)$$

and since the DCT is an orthonormal transformation, the rows of D_m are orthogonal, so $D_m D_m^T = I$. Therefore

$$y = D_m x \quad (5)$$

Since there is no contradiction, the least squares solution must use the first m DCT coefficients. \square

2. Proof of the DCT Mean-Variance Theorem

Theorem 2 (DCT Mean-Variance Theorem). *Given a set of samples of a signal X such that $E[X] = 0$, let Y be the DCT coefficients of X . Then*

$$\text{Var}[X] = E[Y^2] \quad (6)$$

Proof. Start by considering $\text{Var}[X]$. We can rewrite this as

$$\text{Var}[X] = E[X^2] - E[X]^2 \quad (7)$$

Since we are given $E[X] = 0$, this simplifies to

$$\text{Var}[X] = E[X^2] \quad (8)$$

Next, we express the DCT as a linear map such that $X = DY$ and rewrite the previous equation as

$$\text{Var}[X] = E[(DY)^2] \quad (9)$$

Squaring gives

$$E[(DY)^2] = E[(D^T D)Y^2] \quad (10)$$

Since D is orthogonal this simplifies to

$$E[(D^T D)Y^2] = E[(D^{-1} D)Y^2] = E[Y^2] \quad (11)$$

\square

3. Algorithms

We conclude by outlining in pseudocode the algorithms for the three layer operations described in the paper. Algorithm 1 gives the code for convolution explosion, Algorithm 2 gives the code for the ASM ReLu approximation, and Algorithm 3 gives the code for Batch Normalization.

Algorithm 1 Convolution Explosion. K is an initial filter, p, p' are the input and output channels, h, w are the image height and width, s is the stride, \star_s denotes the discrete convolution with stride s . J and \tilde{J} are constants of shape (x, y, k, h, w) with $y = h/8, x = w/8, k = 64$.

```

function EXPLODE( $K, p, p', h, w, s$ )
   $d_j \leftarrow \text{shape}(\tilde{J})$ 
   $d_b \leftarrow (d_j[0], d_j[1], d_j[2], 1, h, w)$ 
   $\hat{J} \leftarrow \text{reshape}(\tilde{J}, d_b)$ 
   $\hat{C} \leftarrow \hat{J} \star_s K$ 
   $d_c \leftarrow (p, p', d_j[0], d_j[1], d_j[2], h/s, h/s)$ 
   $\tilde{C} \leftarrow \text{reshape}(\hat{C}, d_c)$ 
  return  $\tilde{C}_{p'hw}^{pxyk} J_{x'y'k'}$ 

```

Algorithm 2 Approximated Spatial Masking for ReLu. F is a DCT domain block, ϕ is the desired maximum spatial frequencies, N is the block size.

```

function RELU( $F, \phi, N$ )
   $M \leftarrow$  ANNM( $F, \phi, N$ )
  return APPLYMASK( $F, M$ )

function ANNM( $F, \phi, N$ )
   $I \leftarrow$  zeros( $N, N$ )
  for  $i \in [0, N)$  do
    for  $j \in [0, N)$  do
      for  $\alpha \in [0, N)$  do
        for  $\beta \in [0, N)$  do
          if  $\alpha + \beta \leq \phi$  then
             $I_{ij} \leftarrow I_{ij} + F_{ij} D_{ij}^{\alpha\beta}$ 

   $M \leftarrow$  zeros( $N, N$ )
   $M[I > 0] \leftarrow 1$ 
  return  $M$ 

function APPLYMASK( $F, M$ )
  return  $H_{\alpha'\beta'}^{\alpha\beta ij} F_{\alpha\beta} M_{ij}$ 

```

Algorithm 3 Batch Normalization. F is a batch of JPEG blocks (dimensions $N \times 64$), S is the inverse quantization matrix, m is the momentum for updating running statistics, t is a flag that denotes training or testing mode. The parameters γ and β are stored externally to the function. $\hat{\cdot}$ is used to denote a batch statistic and $\tilde{\cdot}$ is used to denote a running statistic.

```

function BATCHNORM( $F, S, m, t$ )
  if  $t$  then
     $\mu \leftarrow$  mean( $F[:, 0]$ )
     $\hat{\mu} \leftarrow F[:, 0]$ 
     $F[:, 0] = 0$ 
     $D_g \leftarrow F_k S_k$ 
     $\hat{\sigma}^2 \leftarrow$  mean( $F^2, 1$ )
     $\sigma^2 \leftarrow$  mean( $\hat{\sigma}^2 + \hat{\mu}^2$ ) -  $\mu^2$ 
     $\tilde{\mu} \leftarrow \hat{\mu}(1 - m) + \mu m$ 
     $\sigma^2 \leftarrow \hat{\sigma}^2(1 - m) + \mu m$ 
     $F[:, 0] \leftarrow F[:, 0] - \mu$ 
     $F \leftarrow \frac{\gamma F}{\sigma}$ 
     $F[:, 0] \leftarrow F[:, 0] + \beta$ 
  else
     $F[:, 0] \leftarrow F[:, 0] - \tilde{\mu}$ 
     $F \leftarrow \frac{\gamma F}{\tilde{\sigma}}$ 
     $F[:, 0] \leftarrow F[:, 0] + \beta$ 
  return  $F$ 

```