



ALMA MATER STUDIORUM
UNIVERSITÀ DI BOLOGNA

ARCHIVIO ISTITUZIONALE
DELLA RICERCA

Alma Mater Studiorum Università di Bologna Archivio istituzionale della ricerca

CAS-CNN: A deep convolutional neural network for image compression artifact suppression

This is the final peer-reviewed author's accepted manuscript (postprint) of the following publication:

Published Version:

CAS-CNN: A deep convolutional neural network for image compression artifact suppression / Cavigelli, Lukas; Hager, Pascal; Benini, Luca. - ELETTRONICO. - (2017), pp. 7965927.752-7965927.759. (Intervento presentato al convegno 2017 International Joint Conference on Neural Networks (IJCNN) tenutosi a Anchorage, Alaska, USA nel May 14-19, 2017) [10.1109/IJCNN.2017.7965927].

Availability:

This version is available at: <https://hdl.handle.net/11585/624746> since: 2018-09-22

Published:

DOI: <http://doi.org/10.1109/IJCNN.2017.7965927>

Terms of use:

Some rights reserved. The terms and conditions for the reuse of this version of the manuscript are specified in the publishing policy. For all terms of use and more information see the publisher's website.

This item was downloaded from IRIS Università di Bologna (<https://cris.unibo.it/>).
When citing, please refer to the published version.

(Article begins on next page)

This is the post peer-review accepted manuscript of:

L. Cavigelli, P. Hager and L. Benini, "CAS-CNN: A deep convolutional neural network for image compression artifact suppression," 2017 International Joint Conference on Neural Networks (IJCNN), Anchorage, AK, 2017, pp. 752-759. doi: 10.1109/IJCNN.2017.7965927

The published version is available online at: <https://doi.org/10.1109/IJCNN.2017.7965927>

© 2017 IEEE. Personal use of this material is permitted. Permission from IEEE must be obtained for all other uses, in any current or future media, including reprinting/republishing this material for advertising or promotional purposes, creating new collective works, for resale or redistribution to servers or lists, or reuse of any copyrighted component of this work in other works

CAS-CNN: A Deep Convolutional Neural Network for Image Compression Artifact Suppression

Lukas Cavigelli, Pascal Hager, Luca Benini

Integrated Systems Laboratory, ETH Zurich, Zurich, Switzerland, Email: surname@iis.ee.ethz.ch

Abstract—Lossy image compression algorithms are pervasively used to reduce the size of images transmitted over the web and recorded on data storage media. However, we pay for their high compression rate with visual artifacts degrading the user experience. Deep convolutional neural networks have become a widespread tool to address high-level computer vision tasks very successfully. Recently, they have found their way into the areas of low-level computer vision and image processing to solve regression problems mostly with relatively shallow networks.

We present a novel 12-layer deep convolutional network for image compression artifact suppression with hierarchical skip connections and a multi-scale loss function. We achieve a boost of up to 1.79 dB in PSNR over ordinary JPEG and an improvement of up to 0.36 dB over the best previous ConvNet result. We show that a network trained for a specific quality factor (QF) is resilient to the QF used to compress the input image—a single network trained for QF 60 provides a PSNR gain of more than 1.5 dB over the wide QF range from 40 to 76.

I. INTRODUCTION

Compression methods can be split into two categories: lossless (e.g. PNG) and lossy (e.g. JPEG) [1]. While lossless methods provide the best visual experience to the user, lossy methods have an non-invertible compression function but can achieve a much higher compression ratio. They often come with a parameter to span the trade-off between file size and quality of the decompressed image. In practical uses, lossy compression schemes are often preferred on consumer devices for their much higher compression rate [1].

Particularly at high compression rates, the differences between the decompressed and the original image become visible with artifacts that are specific of the applied compression scheme. These are not only unpleasant to see, but also have a negative impact on many low-level vision algorithms [2]. Many compression algorithms rely on tiling the images into blocks, applying a sparsifying transform and re-quantization, followed by a generic loss-less data compression [3].

JPEG has become the most widely accepted standard in lossy image compression [4], with many efficient software transcoders publicly available and specialized hardware accelerators deployed in many cameras. Due to its popularity, JPEG-compressed images are also widely found on storage devices containing memories of moments experienced with family and friends, capturing the content of historic documents, and holding on to evidence in legal investigations.

Image compression is also used in wireless sensors systems to transfer visual information from sensor nodes to central storage and processing sites. In such systems, the transmitting node is often battery-powered and thus heavily power-

constrained [5]. Transmitting data is often the most expensive part in terms of energy, and strong compression can mitigate this by reducing the required transmit energy at the expense of introducing compression artifacts [3]. Similar challenges are also seen in mobile devices storing data: size and cost constraints limit the amount of memory for data storage, and the energy available on such devices is depleted rapidly when writing to flash memory—so much that it pays off to apply compression before writing to flash memory [6], [7]. On the processing site, these space and energy constraints are absent and much more computational power is available to decompress and possibly post-process the transmitted or stored images [3].

Deep convolutional neural networks (ConvNets) have become an essential tool for computer vision, even exceeding human performance in tasks such as image classification [8], object detection [9], and semantic segmentation [10], [11]. In addition, they have also started to gain relevance for regression tasks in low-level image and video processing, computing saliency maps [12], optical flow fields [13] and single-image super-resolution images [14] with state-of-the-art performance.

In this work, we present 1) the construction of a new deep convolutional neural network architecture to remove compression artifacts in JPEG compressed image data, 2) a strategy to train this deep network, adaptable to other low-level vision tasks, and 3) extensive evaluations on the LIVE1 dataset, highlighting the properties of our network and showing that this is the current state-of-the-art performance ConvNet for compression artifact suppression (CAS).

II. RELATED WORK

Traditional approaches to suppress compression artifacts can be split into three categories. Various types of intelligent edge-aware denoising such as SA-DCT [15], [16], BM3D [17] have been proposed to address this task during the late 2000s. In recent years, dictionary-based sparse recovery algorithms such as DicTV [18], RTF [19], S-D2 [20], D^3 [21], DDCN [22] have achieved outstanding results by directly addressing the deficiencies such as ringing and blocking very specific to JPEG. These algorithms explicitly attempt to optimally reverse the effect of DCT-domain quantization using learned dictionaries very specific to the applied compressor and quantization tables.

This work was inspired by single-image super-resolution ConvNets, which are a special case of compression artifact removal, where the compression is a simple sub-sampling

operation. Several networks have shown to be very successful at this task, such as SRCNN [14] or DRCN [23]. They use different training procedures and approaches for network construction, but both ConvNets are a simple sequence of convolution and point-wise non-linearity layers.

Recently, two important works have been published, which apply ConvNets for compression artifact suppression: AR-CNN [2], [24] and the approach presented in [25]. The former starts from the architecture presented in SRCNN. In order to overcome convergence problems, they use transfer-learning from the 4-layer network retrained for artifact reduction to a deeper 5-layer network, as well as between networks trained for different JPEG quality factors (QFs) and datasets. In [25] a residual structure extends the simple stacking of convolutional, non-linearity and pooling layers, such that the network is only trained to produce an increment compensating for the distortions. Furthermore, skip elements where some feature maps are bypassing one or multiple layers and are then concatenated to the feature maps at a later stage were introduced. Additionally, they do not use a plain MSE loss function but also include an additional term to emphasize edges.

The networks of both works were trained on the 400 images contained in the BSDS500 train and test sets and evaluated on the remaining 100 images in the validation set. Testing of these networks was then performed on the LIVE1 dataset (29 images) [26] and, in case of AR-CNN, on the 5 test images of [15] and a self-collected dataset of 40 photographs from *twitter* as well. We will adopt their test datasets, procedures and quality measures. Our choice of the training dataset is discussed in Section III-D.

III. METHODOLOGY

We start from the basic concept of training a deep ConvNet for a regression problem, as has been done for the related task of superresolution [14], [23] or other low-level computer vision operations such as optical flow estimation [13]. The authors of [25] propose several new elements for artifact reduction ConvNets: A residual architecture, an edge-emphasized loss function, symmetric weight initialization, and skip connections. All these elements were introduced to alleviate the obstacles preventing the training of deep networks for regression tasks. Taking inspiration from deep neural networks such as FlowNet [13] and FCN [10] developed for optical flow estimation and semantic segmentation respectively, we propose a neural network with hierarchical skip connections (cf. Section III-A) and a multi-scale loss function (cf. Section III-C) for compression artifact suppression.

A. Network Architecture

An overview of our proposed network is shown in Figure 1. The blocks A, \dots, D each consist of two convolutional layers, increasing the number of channels from 1 to 128 and later to 256, the deeper they are in the network. At the same time the resolution is reduced by down-sampling (DS), which is implemented with 2×2 pixel average-pooling layers with 2×2 stride. The main path through the ConvNet (marked

TABLE I: Hyperparameters of the Layers

| name | type | #outp. ch. | #inp. ch. | filter size | #param. |
|-------------|----------|------------|-----------|----------------|---------|
| $A^{(1)}$ | conv | 128 | 1 | 3×3 | 1k |
| $A^{(2)}$ | conv | 128 | 128 | 3×3 | 147k |
| $B^{(1)}$ | conv | 128 | 128 | 3×3 | 147k |
| $B^{(2)}$ | conv | 128 | 128 | 3×3 | 147k |
| $C^{(1)}$ | conv | 128 | 256 | 3×3 | 295k |
| $C^{(2)}$ | conv | 256 | 256 | 3×3 | 590k |
| $D^{(1)}$ | conv | 256 | 256 | 3×3 | 590k |
| $D^{(2)}$ | conv | 256 | 256 | 3×3 | 590k |
| \tilde{D} | fullconv | 256 | 256 | $4 \times 4/2$ | 1049k |
| \hat{D} | conv | 1 | 256 | 3×3 | 2k |
| \tilde{C} | fullconv | 128 | 513 | $4 \times 4/2$ | 1051k |
| \hat{C} | conv | 1 | 513 | 3×3 | 5k |
| \tilde{B} | fullconv | 128 | 257 | $4 \times 4/2$ | 526k |
| \hat{B} | conv | 1 | 257 | 3×3 | 2k |
| \hat{A} | conv | 1 | 257 | 3×3 | 2k |
| Total | | | | | 5144k |

blue in Figure 1) then proceeds through the full-convolution¹ layers $\tilde{D}, \dots, \tilde{B}$ and the normal convolution layer \hat{A} . This way we obtain a 12-layer ConvNet, which however cannot be trained to achieve state-of-the-art accuracy using standard training methods. In the following, we list modifications to the network reducing the average path length, allowing it to converge to beyond state-of-the-art accuracy.

To reduce the path length, the higher-resolution intermediate results after each full-convolution layer are enhanced in the next layer by concatenating the lower-level features extracted earlier in the network natively at this resolution (marked red in Figure 1). We expect this to benefit the network two-fold: once through the additional information to help infer high-resolution outputs, and second to aid in training the early layers of the network by means of bypassing the middle layers.

Training deep networks for regression tasks is problematic and while we have reduced the path length for some paths (e.g. input $\rightarrow A \rightarrow \hat{A} \rightarrow$ output) using the aforementioned method, some very long paths remain. The gradients for adjusting the weights of D are propagated from the output through $\hat{A}, \tilde{B}, \tilde{C}, \tilde{D}, D$. To improve on this, we introduce a multi-scale optimization criterion: instead of optimizing input-to-output, we reconstruct low-resolution images already from deep within the network using a single convolutional layer (marked green in Figure 1), i.e. $\hat{D}, \hat{C}, \hat{B}$ for 1/64-th, 1/16-th, and 1/4-th of the resolution, respectively. We do not discard the output, but up-sample (US) it by a factor of $2 \times$ in each spatial dimension using nearest-neighbor interpolation and concatenate it to the feature maps generated by the full-convolution layer parallel to this path (marked yellow in Figure 1). Using this configuration, we have further shortened the deepest stack of layers significantly by reducing the distance from the middle layers to the output.

The parameters of the convolution and full-convolution

¹We use the definition of full-convolution (also known as up-convolution, deconvolution, backwards convolution, or fractional-strided convolution) as described in [10], [27].

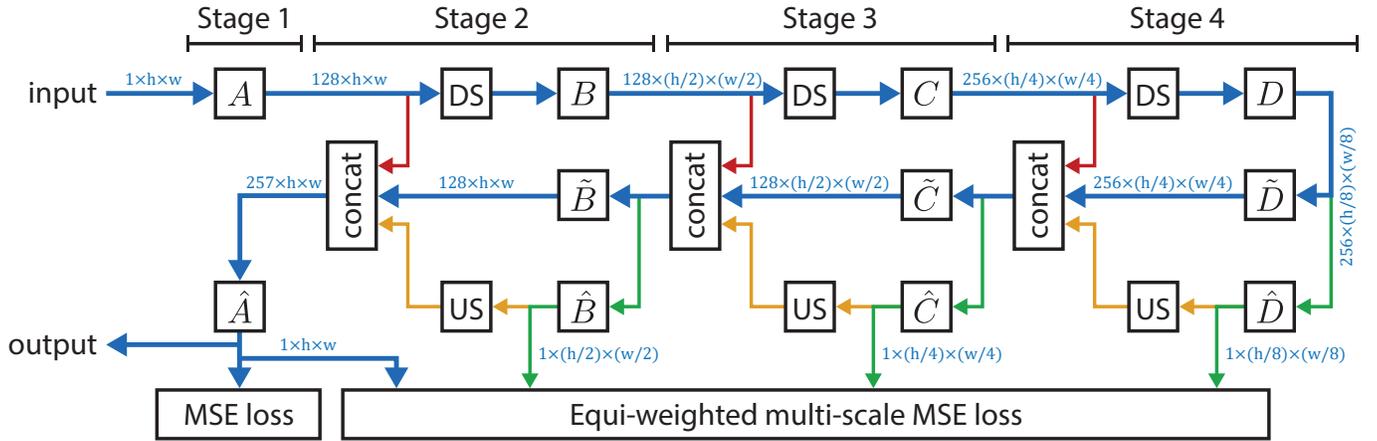


Fig. 1: Structure of the proposed ConvNet. The paths are color coded: **main path (bold)**, **concatenation of lower-level features**, **multi-scale output paths**, **re-use of multi-scale outputs**.

layers are listed in Table I. All these layers are followed by a Parametric Rectified Linear Unit (PReLU) [28] activation layer, where the slope for negative inputs is learned from data rather than pre-defined. These units have shown superior performance for ImageNet classification [28], reducing the issues of dead features [29].

We have found that learning a residual to the input image instead of the reconstructed image as suggested in previous work [25] did not improve the performance of the proposed ConvNet and thus do not include it in our network. The initial weight and bias values are drawn uniformly from the interval $(-n_{in}^{-1/2}, n_{in}^{-1/2})$, where n_{in} is the number of input channels into that layer.

Batch normalization has shown to reduce the achievable accuracy. The batch-wise normalization of means and variances introduces batch-to-batch jitter thereof into the system, preventing full convergence of the network to the maximum accuracy obtained otherwise.

B. Performance Metrics

The most wide-spread performance metrics to evaluate differences between images and many other signals are the mean-squared error (MSE) and the closely related peak signal-to-noise ratio (PSNR). The MSE is the pixel-wise average over the squared difference in intensity between the distorted and the reference image. The PSNR is the MSE normalized to the maximum possible signal values typically expressed in decibel (dB). Following [24], [25] with pixel values normalized to the range [0, 1], we use

$$\text{PSNR}(\mathbf{X}, \hat{\mathbf{X}}) = 10 \log_{10} \left(1/\text{MSE}(\mathbf{X}, \hat{\mathbf{X}}) \right), \quad (1)$$

$$\text{MSE}(\mathbf{X}, \hat{\mathbf{X}}) = \left(\sum_{p \in \mathcal{P}} e(\mathbf{x}_p, \hat{\mathbf{x}}_p)^2 \right) / |\mathcal{P}|, \quad (2)$$

where \mathcal{P} is the set of pixel indexes, \mathbf{X} is the reference image, $\hat{\mathbf{X}}$ is the image to evaluate, and e is the per-pixel error function (e.g. $|x_p - \hat{x}_p|$ for grayscale images).

Both metrics are fully referenced, comparing individual pixels to the original image and converging to zero for a perfect reconstruction. They are known to differ from perceived visual quality [1], [30]–[32] but find wide-spread use due to their simplicity. A variation of the PSNR measure is the IPSNR (increase in PSNR), which is the PSNR difference to the baseline distorted image and thus measures quality improvement. It is also more stable across different datasets.

A popular alternative is to use the structural similarity index (SSIM) [30], which is the mean of the product of three terms assessing similarity in luminance, contrast and structure over multiple localized windows. We use the Matlab implementation provided with [30] for evaluation and use the same parameters as related work [2], [24], [25]: $K_1 = 0.01$, $K_2 = 0.03$, and a 8×8 local statistics window \mathbf{w} of ones. A third measure used in related work is the PSNR-B [33], which adds a (non-referenced) blocking effect factor (BEF) term to the MSE measure. The BEF measures luminance discontinuities at the horizontally and vertically oriented block boundaries. We define the IPSNR-B analogous to the IPSNR.

C. Loss Function

During the training of the ConvNets we minimize the MSE criterion, penalizing deviations from the reference image by the squared distance. However, as mentioned in Section III-A, in order to improve the training procedure we include not only the full-resolution output, but also the low-resolution outputs from within the network. The reference for these is computed by down-sampling the input image, averaging across 4, 16 and 64 pixels, respectively. Each of these outputs' MSE contributes equally to the overall *multi-scale (MS) loss* function.

We run the training until convergence with this objective, before removing the lower resolution images from the loss function and continue the training for several epochs to minimize the MSE of only the full-resolution output image (*output loss*), fine-tuning (FT) the network with this optimization objective. For QF 20, training is performed with MS loss until saturation of the testing MSE (Epoch 200), and then fine-tuned

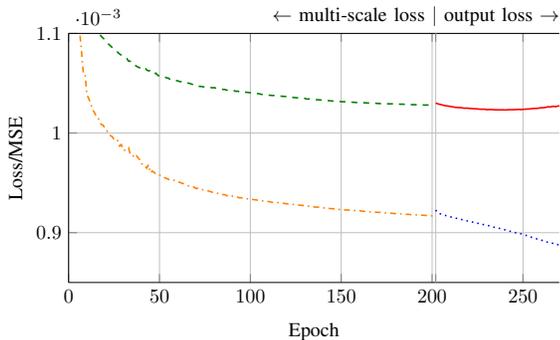


Fig. 2: Loss improvement by number of training epochs for compression with quality factor 20. The multi-scale MSE loss is minimized (--- , scaled up by $3\times$ for readability) until the testing MSE (---) is saturated at Epoch 200, after which the network is fine-tuned using the MSE loss on the output image (⋯). Then the parameters for which the testing MSE (—) is minimal are selected. Note: An epoch during the fine-tuning phase contains 150k instead of 50k images.

with the output loss/MSE criterion before selecting the weights with the lowest testing MSE (cf. Figure 2).

In previous work, including an edge-emphasized term into the loss function has been proposed [25]. We decided not to introduce such a loss term because it leads to an additional hyperparameter to adjust the weight and because we consider it inconsistent to train the network with a loss function different from the quality measure used to benchmark the results. Tuning the hyperparameters for the best PSNR would result in choosing the weight value of the edge-emphasized loss term to be zero.

As such, it prevents further improvement in terms of PSNR and SSIM beyond some limit, and the factor with which it is weighted can be used to trade-off overall reconstruction quality and deblocking. We do not include such a term in our setup because our main objective is to maximize the overall reconstruction, which already implies a high-quality deblocking. By training on a large dataset we do not require such a regularization term.

D. Dataset

Previous networks for compression artifact reduction were trained on the 400 train and test images of the BSDS500 dataset and tested on the 100 remaining validation images [2], [24], [25]. The authors of [25] show that this is the limiting factor for further improvement of their larger L8 network with 220k learned parameters. We do not want to constrain the size of our network by the amount of available training data, particularly since we do not need hard-to-obtain labels for it. We thus use the large, widely-known and publicly available ImageNet 2013 detection dataset [35], which consists of 396k training and 20k validation color images of various sizes. From each image we take cut-outs of 120×120 pixels to generate our dataset.

The color images are transformed to YCbCr space and only the luminance channel is used further. The input to the network

is then generated by compressing the resulting image using the Matlab JPEG compressor² with a bit depth of 8.

For training our network we take 50k images of the 120×120 pixel cut-outs from the training set and 10k cut-outs for the validation set. We increase the size of the training set to 150k for fine-tuning with the output loss function. Testing is performed on the 29 images of the LIVE1 dataset.

We use the Torch framework [36] with cuDNN v5.1.3 [37] for our evaluations. We use the Adam optimizer [38] starting with a learning rate of 10^{-4} . A minibatch size of 20 images was chosen and training was parallelized over two Nvidia Titan X Maxwell GPUs. We have not applied any preprocessing to the images before feeding them into the network. Our main training was conducted for QF 20 compressed input data and we have trained the networks for other quality factors starting from this one to reduce training time. For the forward pass, a throughput of 1.01 Mpixel/s has been measured with a Nvidia GTX1080 using single-precision floating-point operations.

IV. RESULTS & DISCUSSION

We have evaluated the mean PSNR, PSNR-B and SSIM across the LIVE1 dataset for the JPEG quality factors 10, 20, 40, 60 and 80, and compare them to related work in Table II. We use the same JPEG compressor as in AR-CNN [24] and Svoboda *et al.* [25] (i.e. Matlab), with which we obtain the identical baseline PSNR of 30.07 dB for QF 20 and 27.77 dB for QF 10 for the JPEG compressed image with respect to the uncompressed reference.

²We have used this compressor to remain comparable with related work. Other implementations such as libjpeg or libjpeg-turbo use different quantization tables and, in case of these two libraries, result in a significantly larger file size and as a consequence also a better PSNR for the same quality factor.

TABLE II: Restoration Quality Comparison on LIVE1

| QF | Algorithm | | PSNR [dB] | PSNR-B [dB] | SSIM |
|------------------|------------------|--------------|--------------|--------------|--------------|
| 10 | JPEG | [34] | 27.77 | 25.33 | 0.791 |
| | SA-DCT | [15] | 28.65 | 28.01 | 0.809 |
| | AR-CNN | [2] | 29.13 | 28.74 | 0.823 |
| | L4 | [25] | 29.08 | 28.71 | 0.824 |
| | ours, MS loss | | 29.36 | 28.92 | 0.830 |
| | ours, w/ loss FT | | 29.44 | 29.19 | 0.833 |
| 20 | JPEG | [34] | 30.07 | 27.57 | 0.868 |
| | SA-DCT | [15] | 30.81 | 29.82 | 0.878 |
| | AR-CNN | [2] | 31.40 | 30.69 | 0.890 |
| | L4 | [25] | 31.42 | 30.83 | 0.890 |
| | L8 | [25] | 31.51 | 30.92 | 0.891 |
| | ours, MS loss | | 31.67 | 30.84 | 0.894 |
| ours, w/ loss FT | | 31.70 | 30.88 | 0.895 | |
| 40 | JPEG | [34] | 32.35 | 29.96 | 0.917 |
| | SA-DCT | [15] | 32.99 | 31.79 | 0.924 |
| | AR-CNN | [2] | 33.63 | 33.12 | 0.931 |
| | L4 | [25] | 33.77 | — | — |
| | ours, MS loss | | 33.98 | 32.83 | 0.935 |
| | ours, w/ loss FT | | 34.10 | 33.68 | 0.937 |
| 60 | JPEG | [34] | 33.99 | 31.89 | 0.940 |
| | ours, w/ loss FT | | 35.78 | 35.10 | 0.954 |
| 80 | JPEG | [34] | 36.88 | 35.47 | 0.964 |
| | ours, w/ loss FT | | 38.55 | 37.73 | 0.973 |

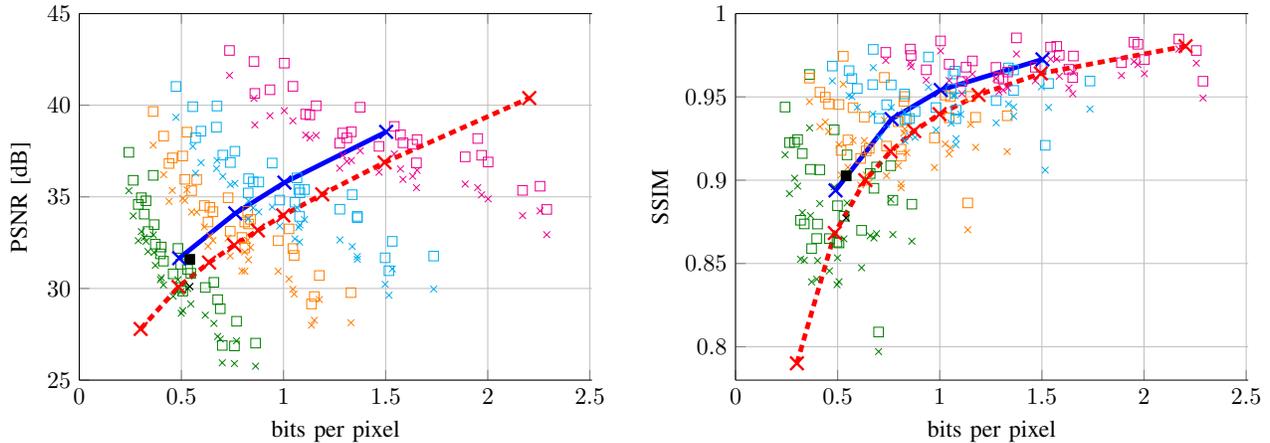


Fig. 3: PSNR (left) and SSIM (right) evaluated on the LIVE1 dataset with respect to the number of bits per pixel required to store the compressed image. The ordinary JPEG performance is shown as (---) for QF 10 to 90 in steps of 10, averaged over all images in the dataset. Individual images are shown with markers: ordinary JPEG (\times), after CAS-CNN (\square). The image depicted in Figure 5 is marked with (\blacksquare and \times). The different quality factors are color coded: QF 20 (\square), QF 40 (\square), QF 60 (\square), QF 80 (\square). The CAS-CNN output quality averaged over the dataset is shown as (—).

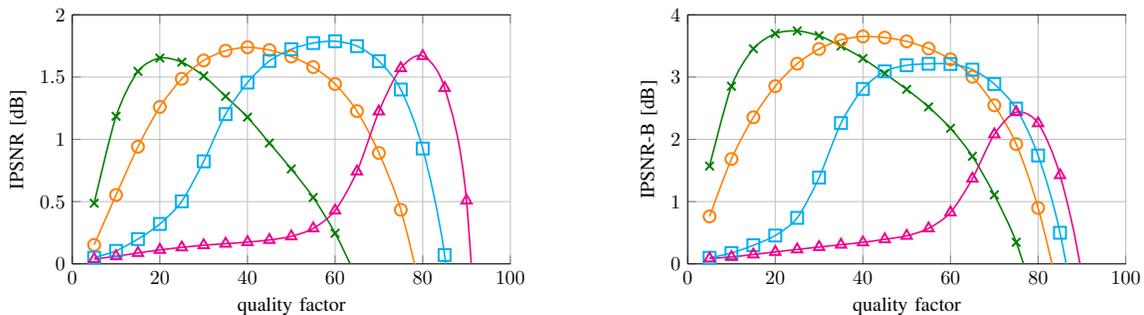


Fig. 4: PSNR and PSNR-B improvement for various compression quality factors for networks trained with images compressed with a single quality factor: QF 20 (\rightarrow), QF 40 (\rightarrow), QF 60 (\rightarrow), QF 80 (\rightarrow), evaluated on the LIVE1 dataset.

For our network, we list results directly after training with the multi-scale loss function as well as after fine-tuning with the output loss. The already state-of-the-art results are further improved by this two-step learning procedure. Overall, we can see a significant improvement in PSNR of 0.19 dB over the L8 network [25], 0.30 dB over AR-CNN and 1.63 dB over ordinary JPEG for QF 20. The SSIM is also improved to 0.895. For QF 10 we see a gain of 1.67 dB over ordinary JPEG, 0.36 dB over the L4 network and 0.31 dB over AR-CNN, the state-of-the-art ConvNet for this configuration.

For QF 10, we improve the PSNR-B by 0.45 dB over previous work. However, for a lower compression rate, we do not exceed the PSNR-B value achieved by the L8 network. As described in the next paragraph, there are no visible blocking artifacts after applying our ConvNet. PSNR-B has been introduced for benchmarking deblocking algorithms, and by its definition the blocking artifact-penalizing term measuring the differences between pixels along the block boundary does not vanish even for a perfect reconstruction. An image with higher reconstruction quality might thus suffer from a lower PSNR-B value because of clearer edges all over the image including at

the block boundaries.

In Figure 3 we show the distribution of the individual images of the LIVE1 dataset in terms of PSNR and SSIM with respect to the used number of bits per pixel for several QFs. The average PSNR and SSIM for each QF is also shown, visualizing that this method works for strong as well as for weak compression. Looking at the individual images, it becomes visible that our method improves not only the mean PSNR and SSIM, but enhances each individual image.

As discussed in Section III-B, the visual perception can differ from quantitative evaluations using classical quality measures. To give a visual impression as well, we provide a qualitative visual comparison in Figure 5. The *lighthouse3* image serves as a basis for this comparison and is the same one used in [25]. It is shown with black markers in Figure 3, indicating that this image is not a particularly well-working outlier. A clear improvement is visible, there are no perceptible blocking artifacts anymore and the ringing artifacts are strongly suppressed without blurring the railing depicted in the image. For completeness, we also provide the results for the 5 classical test images used throughout many compression

papers (cf. Figure 6). The trained models and scripts required to reproduce these images are available online³.

In Figure 4, we show that the networks trained for a specific quality factor do not need to be retrained for the specific quality factor with which the image was compressed to achieve a high improvement in PSNR or PSNR-B. The network trained for QF 60 already boosts the PSNR by more than 1.5 dB for quality factors ranging from 25 to almost 60. This resilience to variations in quantization has not been shown for approaches focusing on DCT-domain recovery.

V. CONCLUSION

We have presented a 12-layer deep convolutional neural network for compression artifact suppression in JPEG images with hierarchical skip connections and trained with a multi-scale loss function. The result is a new state-of-the-art ConvNet achieving a boost of up to 1.79 dB in PSNR over ordinary JPEG and showing an improvement of up to 0.36 dB over the best previous ConvNet result. We have shown that a network trained for a specific quality factor is resilient to the QF used to compress the input image—a single network trained for QF 60 provides a PSNR gain of more than 1.5 dB over the wide QF range from 40 to 76. The obtained results are also qualitatively superior to those of existing ConvNets. The network is not tailored to the JPEG-specific compression procedure, and can thus potentially be applied to a wide range of image compression algorithms.

ACKNOWLEDGMENTS

The authors would like to thank Thilo Weber and Jonas Wiesendanger for their preliminary explorations on this topic. This work has received funding from *armasuisse Science & Technology* and the European Union’s Horizon 2020 research and innovation programme under grant agreement no. 732631.

REFERENCES

- [1] Z. Wang, A. C. Bovik, and L. Lu, “Why is image quality assessment so difficult?” in *IEEE Int. Conf. Acoust. Speech Signal Process.*, 2002.
- [2] K. Yu, C. Dong, C. C. Loy, and X. Tang, “Deep Convolutional Networks for Compression Artifacts Reduction,” *arXiv:1608.02778*, 2016.
- [3] L. Chew and L. Ang, “Survey of image compression algorithms in wireless sensor networks,” *2008 Int. Symp. Inf. Technol.*, pp. 1–9, 2008.
- [4] S. Souders, “HTTP Archive - Interesting Stats,” 2016. [Online]. Available: <http://httparchive.org/interesting.php>
- [5] A. Kerhet, M. Magno, F. Leonardi, A. Boni, and L. Benini, “A low-power wireless video sensor node for distributed object detection,” *J. Real-Time Image Process.*, vol. 2, no. 4, pp. 331–342, 2007.
- [6] K. Barr and K. Asanovic, “Energy Aware Lossless Data Compression,” *Proc. of MobiSys*, no. May, 2003.
- [7] Y. Joo, Y. Cho, D. Shin, and N. Chang, “Energy-aware data compression for Multi-Level Cell (MLC) flash memory,” in *Proc. ACM/IEEE Des. Autom. Conf.*, 2007, pp. 716–719.
- [8] K. He, X. Zhang, S. Ren, and J. Sun, “Deep Residual Learning for Image Recognition,” *arXiv:1512.03385*, dec 2015.
- [9] S. Ren, K. He, R. Girshick, and J. Sun, “Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks,” *arXiv:1506.01497*, 2015.
- [10] J. Long, E. Shelhamer, and T. Darrell, “Fully Convolutional Networks for Semantic Segmentation,” in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2015.

- [11] L. Cavigelli, M. Magno, and L. Benini, “Accelerating Real-Time Embedded Scene Labeling with Convolutional Networks,” in *Proc. ACM/IEEE Des. Autom. Conf.*, 2015.
- [12] R. Zhao, W. Ouyang, H. Li, and X. Wang, “Saliency detection by multi-context deep learning,” *Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit.*, vol. 07-12-June, pp. 1265–1274, 2015.
- [13] P. Fischer, A. Dosovitskiy, E. Ilg, P. Haeusser, C. Hazirbas, V. Golkov, P. Van der Smagt, D. Cremers, and T. Brox, “FlowNet: Learning Optical Flow with Convolutional Networks,” in *arXiv:15047.06852*, 2015.
- [14] C. Dong, C. C. Loy, K. He, and X. Tang, “Learning a deep convolutional network for image super-resolution,” *Proc. Eur. Conf. Comput. Vis.*, pp. 184–199, 2014.
- [15] A. Foi, V. Katkovnik, and K. Egiazarian, “Pointwise shape-adaptive DCT for high-quality denoising and deblocking of grayscale and color images,” *IEEE Trans. Image Process.*, vol. 16, no. 5, pp. 1395–1411, 2007.
- [16] —, “Pointwise shape-adaptive DCT for high-quality deblocking of compressed color images,” in *Eur. Signal Process. Conf.*, 2006.
- [17] K. Dabov, A. Foi, and V. Katkovnik, “Image denoising by sparse 3D transformation-domain collaborative filtering,” *IEEE Trans. Image Process.*, vol. 16, no. 8, pp. 1–16, 2007.
- [18] H. Chang, M. K. Ng, and T. Zeng, “Reducing artifacts in JPEG decomposition via a learned dictionary,” *IEEE Trans. Signal Process.*, vol. 62, no. 3, pp. 718–728, 2014.
- [19] J. Jancsary, S. Nowozin, and C. Rother, “Loss-Specific Training of Non-Parametric Image Restoration Models: A New State of the Art,” pp. 112–125, 2012.
- [20] X. Liu, X. Wu, J. Zhou, and D. Zhao, “Data-driven sparsity-based restoration of JPEG-compressed images in dual transform-pixel domain,” pp. 5171–5178, 2015.
- [21] Z. Wang, D. Liu, S. Chang, Q. Ling, and T. S. Huang, “D3: Deep Dual-Domain Based Fast Restoration of JPEG-Compressed Images,” *IEEE Conf. Comput. Vis. Pattern Recognit.*, 2016.
- [22] J. Guo and H. Chao, “Building Dual-Domain Representations for Compression Artifacts Reduction,” in *ECCV*, 2016, pp. 628–644.
- [23] J. Kim, J. K. Lee, and K. M. Lee, “Deeply-Recursive Convolutional Network for Image Super-Resolution,” in *arXiv:1511.04491*, 2015.
- [24] C. Dong, Y. Deng, C. C. Loy, and X. Tang, “Compression Artifacts Reduction by a Deep Convolutional Network,” in *2015 IEEE Int. Conf. Comput. Vis.* IEEE, dec 2015, pp. 576–584.
- [25] P. Svoboda, M. Hradis, D. Barina, and P. Zemcik, “Compression Artifacts Removal Using Convolutional Neural Networks,” *J. WSCG*, vol. 24, no. 2, pp. 63–72, 2016.
- [26] H. R. Sheikh, Z. Wang, L. Cormack, and A. C. Bovik, “LIVE image quality assessment database release 2,” 2005.
- [27] H. Noh, S. Hong, and B. Han, “Learning Deconvolution Network for Semantic Segmentation,” *arXiv:1505.04366*, vol. 1, 2015.
- [28] K. He, X. Zhang, S. Ren, and J. Sun, “Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification,” in *arXiv:1502.01852*, 2015.
- [29] M. D. Zeiler and R. Fergus, “Visualizing and Understanding Convolutional Networks,” in *ECCV 2014, LNCS 8689*, nov 2014, pp. 818–833.
- [30] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli, “Image quality assessment: From error visibility to structural similarity,” *IEEE Trans. Image Process.*, vol. 13, no. 4, pp. 600–612, 2004.
- [31] H. Sheikh and A. Bovik, “Image information and visual quality,” *IEEE Trans. Image Process.*, vol. 15, no. 2, pp. 430–444, 2006.
- [32] B. Girod, “Digital Images and Human Vision,” A. B. Watson, Ed. Cambridge, MA, USA: MIT Press, 1993, ch. What’s wro, pp. 207–220.
- [33] C. Yim and A. C. Bovik, “Quality Assessment of Deblocked Images,” *IEEE Trans. Image Process.*, vol. 20, no. 1, pp. 88–98, jan 2011.
- [34] *MATLAB version 8.5 (R2015a)*, The Mathworks, Inc., Natick, Massachusetts, 2015.
- [35] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, “ImageNet: A large-scale hierarchical image database,” in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2009.
- [36] R. Collobert, “Torch7: A Matlab-like Environment for Machine Learning,” *Adv. Neural Inf. Process. Syst. Work.*, 2011.
- [37] S. Chetlur, C. Woolley, P. Vandermersch, J. Cohen, J. Tran, B. Catanzaro, and E. Shelhamer, “cuDNN: Efficient Primitives for Deep Learning,” in *arXiv:1410.0759*, oct 2014.
- [38] D. Kingma and J. Ba, “Adam: A Method for Stochastic Optimization,” in *Proc. Int. Conf. Learn. Represent.*, dec 2015.

³<http://iis.ee.ethz.ch/~lukasc/cascnn/>



(a) uncompressed



(b) compressed (JPEG QF 20)



(c) SA-DCT



(d) AR-CNN



(e) L8



(f) CAS-CNN (ours)

Fig. 5: Qualitative comparison of reconstruction quality on the *lighthouse3* image of the LIVE1 dataset for JPEG quality factor 20. Images (a),(b),(d),(e) reprinted with permission from [25].

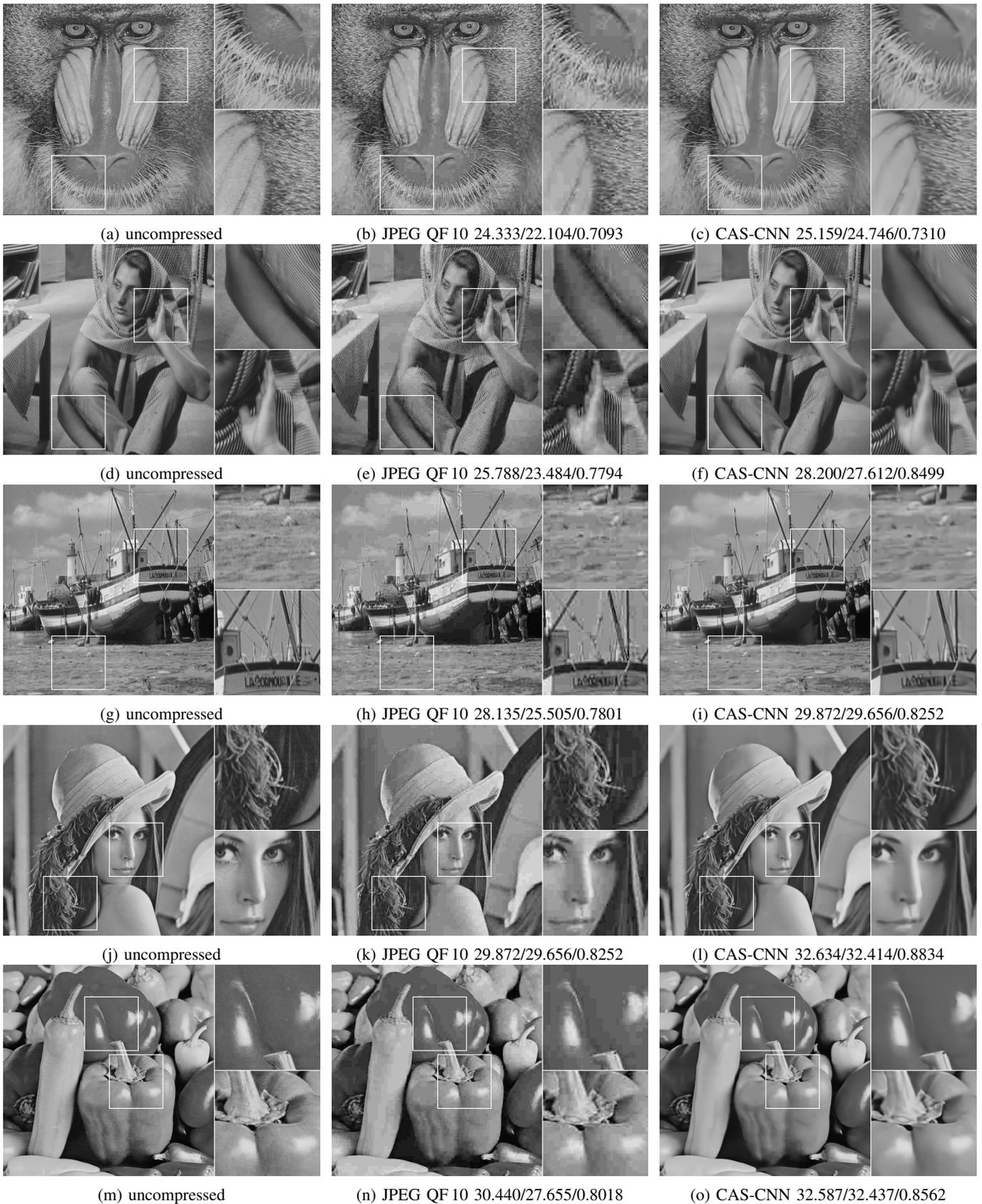


Fig. 6: Evaluation on the 5 classical test images. We show the uncompressed images (left), the Matlab JPEG QF 10 compressed images (center), and the result of applying our CAS-CNN to the compressed images. The PSNR/PSNR-B/SSIM with respect to the uncompressed images is indicated below the images.