

Super-resolution

Sebastien Strban

Faculty of Computer and Information Science

Email: ss4774@student.uni-lj.si

I. ABSTRACT

Many breakthroughs in speed and accuracy of single image super-resolution (SISR) have been achieved. One of the biggest challenges is how to recover finer texture details when super-resolution is applied at large upscaling factors. A typical solution to SISR involves using a convolutional neural network (CNN), however new approaches using a generative adversarial network (GAN) are now also popular. The behavior of optimization-based super-resolution methods is principally driven by the choice of the objective function. In this work, we present an evaluation of SRResNet and SRGAN. SRResNet is a deep residual network and SRGAN is a generative adversarial network for image super-resolution (SR). SRResNet is able to recover reasonable quality photo-realistic textures from heavily downsampled images. SRGAN is capable of inferring photo-realistic natural images with high upscaling factors. This is achieved using a perceptual loss function which consists of an adversarial loss and a content loss. The adversarial loss pushes the solution to the natural image manifold using a discriminator network that is trained to differentiate between the super-resolved images and original photo-realistic images. The content loss is motivated by perceptual similarity instead of similarity in pixel space.

II. INTRODUCTION

The challenging task of estimating and predicting a high-resolution (HR) image from its low-resolution (LR) counterpart is referred to as super-resolution (SR). SR received substantial attention from within the computer vision research community and has a wide range of applications [1], [2], [3].

The optimization target of a supervised SR algorithms is typically the minimization of the mean squared error (MSE) between the recovered HR image and the ground truth. This is convenient as minimizing MSE also maximizes the peak signal-to-noise ratio (PSNR), which is a common measure used to evaluate and compare SR algorithms [4]. However, the ability of MSE to capture perceptually important differences, such as high level texture detail, is very limited as they are defined based on pixel-wise image differences [5], [6], [7].

In this work we present an evaluation of a super-resolution deep residual network with skip-connections (SRResNet) and super-resolution generative adversarial network (SRGAN) which also employs the use of SRResNet. SRResNet uses the MSE as the primary optimization target, while SRGAN is able to diverge from MSE as the sole optimization target and improve on performance and visual quality of results. New perceptual loss is also defined which uses high-level feature

maps of the VGG network [8], [9], [10] in combination with a discriminator that encourages solutions perceptually hard to distinguish from the HR reference images.

A. Related work

1) *Image super-resolution*: Some of the recent overview articles on image SR include Nasrollahi and Moeslund [3] or Yang et al. [4]. In our analysis we are focusing on single image super-resolution (SISR). However, there exist many approaches that recover HR images from multiple images [11], [12].

One of the first methods to tackle SISR were prediction-based methods. These are usually filtering approaches, such as linear, bicubic or Lanczos [13] filtering. Mentioned approaches can be very fast, however they oversimplify the SISR problem and usually yield solutions with overly smooth texture results. Methods that put particularly focus on edge-preservation have also been proposed [14], [15].

More powerful and advanced approaches aim to establish a complex mapping between low-resolution and high-resolution image information. These approaches usually rely on the training data. Many other methods, which are based on example-pairs, rely on LR training patches for which the corresponding HR counterparts are known. Early work was done and presented by Freeman et al. [16], [17].

Recently, convolutional neural network (CNN) based SR algorithms have shown excellent performance. In Wang et al. [18] the authors encode a sparse representation and use their feed-forward network architecture which is based on the learned iterative shrinkage and threshold algorithm [19]. In Dong et al. [20], [21] authors used bicubic interpolation to upscale an input image and trained a three layer deep fully convolutional network in an end-to-end fashion to achieve good performance in SR.

2) *Loss functions*: It is notable that pixel-wise loss functions (such as MSE) struggle to handle the uncertainty which is inherent in recovering high-frequency details (e.g. a texture). Minimizing MSE encourages finding of pixel-wise averages of plausible solutions which are generally smooth and thus give us poor perceptual quality [10], [22], [9]. The problem of minimizing MSE is illustrated in figure 1, where multiple potential solutions, with high texture details, are averaged to create a smooth reconstruction.

In Mathieu et al. [22] and Denton et al. [23] the authors tackled this problem by employing generative adversarial networks (GANs) for the application of image generation.

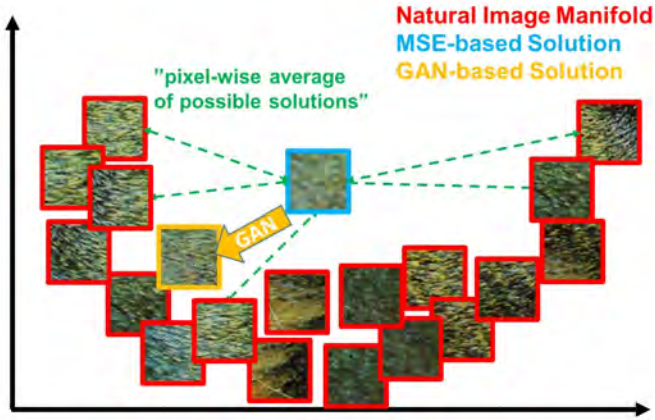


Fig. 1. Depiction of image patches from the natural image manifold (red), super-resolved image patches obtained with MSE (blue) and super-resolved image patches obtained with GAN (orange). The MSE-based solution appears overly smooth due to the pixel-wise average of possible solutions in the pixel space, while GAN-based solution drives the reconstruction towards the natural image manifold producing perceptually more convincing solutions.

Work of Dosovitskiy and Brox [24] uses loss functions based on Euclidean distances computed in the feature space of neural networks in combination with adversarial training. It is shown that the proposed loss allows visually superior image generation and can be used to solve the ill-posed inverse problem of decoding nonlinear feature representations. Similarly, Johnson et al. [9] and Bruna et al. [10] propose the use of features extracted from a pre-trained VGG network instead of low-level pixel-wise error measures. The authors formulate a loss function, which is based on the euclidean distance between feature maps that are extracted from the VGG19 network. Perceptually more convincing results were obtained for both super-resolution and artistic style-transfer [25], [26].

B. Contribution

GANs provide a powerful way for generating plausible-looking natural images which give us high perceptual quality. The GAN procedure encourages the reconstructions to move towards regions of the search space with high probability of containing photo-realistic images. This brings us closer to the natural image manifold as shown in figure 1.

In this work we describe and analyze the first very deep ResNet architecture using the concept of GANs to form a perceptual loss function for photo-realistic SISR [27]. The key aspects of this work are:

- a new state of the art for image SR with high upscaling factors ($4\times$) as measured by PSNR and structural similarity (SSIM) with 16 blocks deep ResNet (SRResNet) optimized for MSE;
- proposed SRGAN which is a GAN-based network optimized for a new perceptual loss. Here the MSE-based content loss was replaced with a loss calculated on feature maps of the VGG network, which are more invariant to changes in pixel space [28].

We describe the network architecture and the perceptual loss

in section III. An evaluation is provided in section IV. The work concludes with section V

III. METHOD

The aim of SISR is to create a super-resolved image I^{SR} , which is an approximation of a high-resolution image I^{HR} from a low-resolution input image I^{LR} . Here I^{LR} is the low-resolution image version of its high-resolution image counterpart I^{HR} . The high-resolution images are only available during training. In training, I^{LR} is obtained by applying a Gaussian filter to I^{HR} followed by a downsampling operation with downsampling factor r . For an image with C color channels, we describe I^{LR} by a real-valued tensor of size $W \times H \times C$, I^{HR} by $rW \times rH \times C$ and I^{SR} by $rW \times rH \times C$.

The ultimate goal is to train a generating function G that estimates for a given LR input image its corresponding HR counterpart. This is achieved by training a generator network as a feed-forward CNN G_{θ_G} parameterized by θ_G . Here $\theta_G = \{W_{1:L}; b_{1:L}\}$ denotes the weights and biases of a L -layer deep network and is obtained by optimizing a SR-specific loss function l^{SR} . For training images I_n^{HR} , $n = 1, \dots, N$ with corresponding I_n^{LR} , $n = 1, \dots, N$, we solve:

$$\hat{\theta}_G = \arg \min_{\theta_G} \frac{1}{N} \sum_{n=1}^N l^{SR}(G_{\theta_G}(I_n^{LR}), I_n^{HR}) \quad (1)$$

In this work we will use perceptual loss l^{SR} which is specifically designed as a weighted combination of several loss components that model distinct desirable characteristics of the recovered SR image.

1) *Adversarial network architecture:* A discriminator network D_{θ_D} is defined and based on Goodfellow et al. [29] which is optimized in an alternating manner along with G_{θ_G} to solve the adversarial min-max problem:

$$\min_{\theta_G} \min_{\theta_D} \mathbb{E}_{I^{HR} \sim p_{train}(I^{HR})} [\log D_{\theta_D}(I^{HR})] + \mathbb{E}_{I^{LR} \sim p_G(I^{LR})} [\log(1 - D_{\theta_D}(G_{\theta_G}(I^{LR})))] \quad (2)$$

The general idea behind this formulation is that it allows one to train a generative model G with the goal of fooling a differentiable discriminator D that is trained to distinguish super-resolved images from real images. With this approach this generator can learn to create solutions that are highly similar to real images and thus difficult to classify by discriminator D . This encourages the manifestation of perceptually superior solutions that reside in the subspace of natural images (the manifold). This approach is in contrast to SR solutions obtained by minimizing pixel-wise error measurements, such as the MSE.

The generator architecture is illustrated in figure 2 and the discriminator architecture in figure 3. At the core of the very deep generator network G are B residual blocks with identical layout. Two convolutional layers are used with small 3×3 kernels and 64 feature maps followed by batch-normalization layers and ParametricReLU as the activation function. The resolution of the input image is increased with two trained sub-pixel convolution layers.

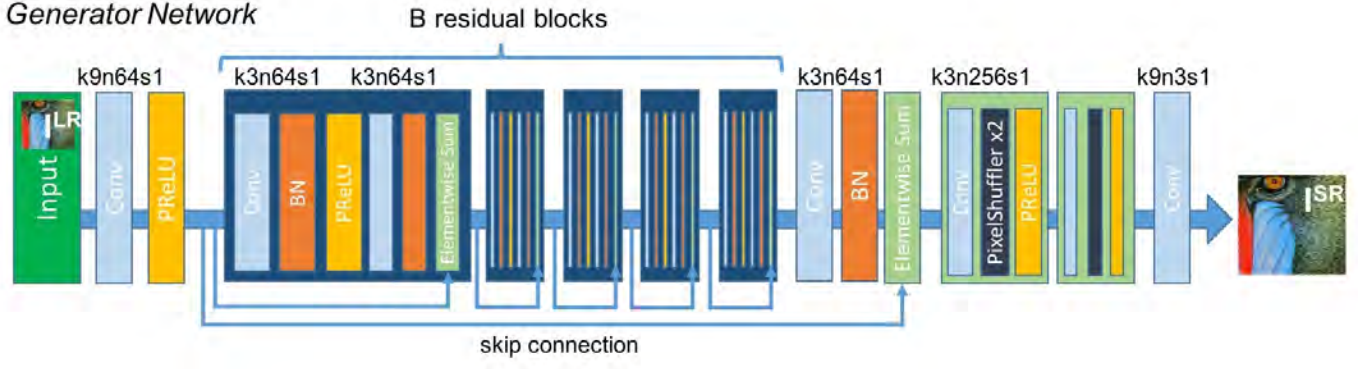


Fig. 2. Architecture of generator network with corresponding kernel size (k), number of feature maps (n) and stride (s) indicated for each convolutional layer.

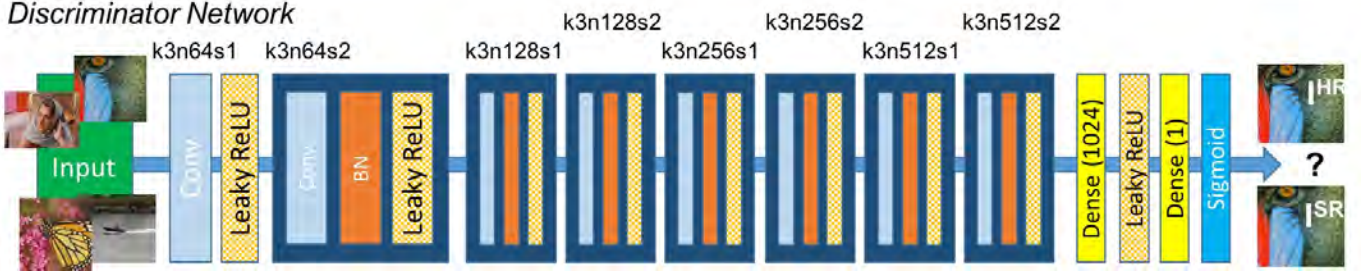


Fig. 3. Architecture of discriminator network with corresponding kernel size (k), number of feature maps (n) and stride (s) indicated for each convolutional layer.

To discriminate real HR images from generated SR samples a discriminator network is trained. LeakyReLU activation is used ($\alpha = 0.2$) and max-pooling is avoided throughout the network. The discriminator network is trained to solve the maximization problem in Equation III-1. It contains 8 convolutional layers with an increasing number of 3×3 filter kernels, increasing by a factor of 2 from 64 to 512 kernels. Strided convolutions are used to reduce the image resolution each time the number of features is doubled. The resulting 512 feature maps are followed by two dense layers and a final sigmoid activation function to obtain a probability for sample classification.

A. Perceptual loss function

The definition the perceptual loss function l^{SR} is critical for the performance of the generator network. This formulation of the perceptual loss is as the weighted sum of a content loss (l_X^{SR}) and an adversarial loss (l_{Gen}^{SR}) component as:

$$l^{SR} = l_X^{SR} + 10^{-3} l_{Gen}^{SR} \quad (3)$$

1) *Content loss*: The pixel-wise MSE loss is calculated as:

$$l_{MSE}^{SR} = \frac{1}{r^2 W H} \sum_{x=1}^{rW} \sum_{y=1}^{rH} (I_{x,y}^{HR} - G_{\theta_G}(I^{LR})_{x,y})^2 \quad (4)$$

This is one of the most widely used optimization target for image SR on which many state-of-the-art approaches

rely [20], [30]. However, while achieving particularly high PSNR, solutions of MSE optimization problems often lack high-frequency content which results in perceptually unsatisfying and unpleasing solutions with overly smooth texture details.

Better choice would be to use a loss function that is closer to perceptual similarity. The VGG loss is defined based on the ReLU activation layers of the pre-trained 19 layer VGG network.

With $\phi_{i,j}$ we indicate the feature map obtained by the j -th convolution (after activation) before the i -th maxpooling layer within the VGG19 network. The VGG loss is then defined as the euclidean distance between the feature representations of a reconstructed image $G_{\theta_G}(I^{LR})$ and the reference image I^{HR} :

$$l_{VGG/i,j}^{SR} = \frac{1}{W_{i,j} H_{i,j}} \sum_{x=1}^{W_{i,j}} \sum_{y=1}^{H_{i,j}} (\phi_{i,j}(I^{HR})_{x,y} - \phi_{i,j}(G_{\theta_G}(I^{LR}))_{x,y})^2 \quad (5)$$

Here $W_{i,j}$ and $H_{i,j}$ describe the dimensions of the respective feature maps within the VGG network.

2) *Adversarial loss*: In addition to the content losses described so far, the generative component of this GAN is also added to the perceptual loss. This encourages our network to favor solutions that reside on the manifold of natural images, when trying to fool the discriminator network.

The generative loss l_{Gen}^{SR} is defined based on the probabil-

ities of the discriminator $D_{\theta_D}(G_{\theta_G}(I^{LR}))$ over all training samples as:

$$l_{Gen}^{SR} = \sum_{n=1}^N -\log D_{\theta_D}(G_{\theta_G}(I^{LR})) \quad (6)$$

Here, $D_{\theta_D}(G_{\theta_G}(I^{LR}))$ is the probability that the reconstructed image $G_{\theta_G}(I^{LR})$ is a natural HR image. For better gradient behavior we minimize $-\log D_{\theta_D}(G_{\theta_G}(I^{LR}))$ instead of $\log(1 - D_{\theta_D}(G_{\theta_G}(I^{LR})))$ [29].

IV. EXPERIMENTS AND EVALUATION

Training is done on the COCO training set 2014 and validation is performed with COCO validation set 2014. Testing and evaluation of experiments is performed on four widely used public benchmark datasets: COCO testing set 2014, Set5, Set14 and BSD100, the testing set of BSD300.

All experiments are performed with a scale factor of $4\times$ between low-resolution and high-resolution images. This corresponds to a $16\times$ reduction in number of image pixels.

A. Base model

The base model follows the architectural implementation of original model as described in [27] and is used as a base for other evaluations.

The base model uses Adam as an optimizer with $\beta_1 = 0.9$. The SRResNet network is trained with a learning rate of 10^{-4} and 10^6 update iterations, which yield 200 epochs in total on our training dataset. The trained SRResNet network is based on MSE loss and is used as initialization for generator when training the actual GAN to avoid undesired local optima. The SRGAN network is based on perceptual loss and trained with 10^{-5} update iterations at a learning rate of 10^{-4} and another 10^5 iterations at a lower learning rate of 10^{-5} , which yield 40 epochs in total on our training dataset. Updates to the generator and discriminator network are alternating ($k = 1$). The generator network has 16 identical residual blocks ($B = 16$).

The evaluation results of the base model are presented in table I. We can observe that SRResNet gives better numerical

	COCO Loss	COCO PSNR	COCO SSIM	Set5 Loss	Set5 PSNR	Set5 SSIM
SRResNet	0.0219	27.3166	0.7766	0.0084	30.6233	0.8899
SRGAN	0.1938	25.3691	0.6998	0.2755	27.6378	0.8098
	Set14 Loss	Set14 PSNR	Set14 SSIM	BSD100 Loss	BSD100 PSNR	BSD100 SSIM
SRResNet	0.0116	28.5642	0.7977	0.0140	27.3658	0.7509
SRGAN	0.1935	25.9921	0.7149	0.1802	25.1075	0.6653

TABLE I
EVALUATION RESULTS OF THE BASE MODEL.

results across the all testing datasets. This phenomenon is also observed in the original paper [27].

Loss, PSNR and SSIM with respect to the number of epochs for SRResNet can be seen on figure 4. Loss, PSNR and SSIM with respect to the number of epochs for SRGAN can be seen on figure 5. We can observe a steady decrease in loss for both SRResNet and SRGAN. There is also a slight gap between validation and train set. For SRResNet it looks like a case of overfitting.

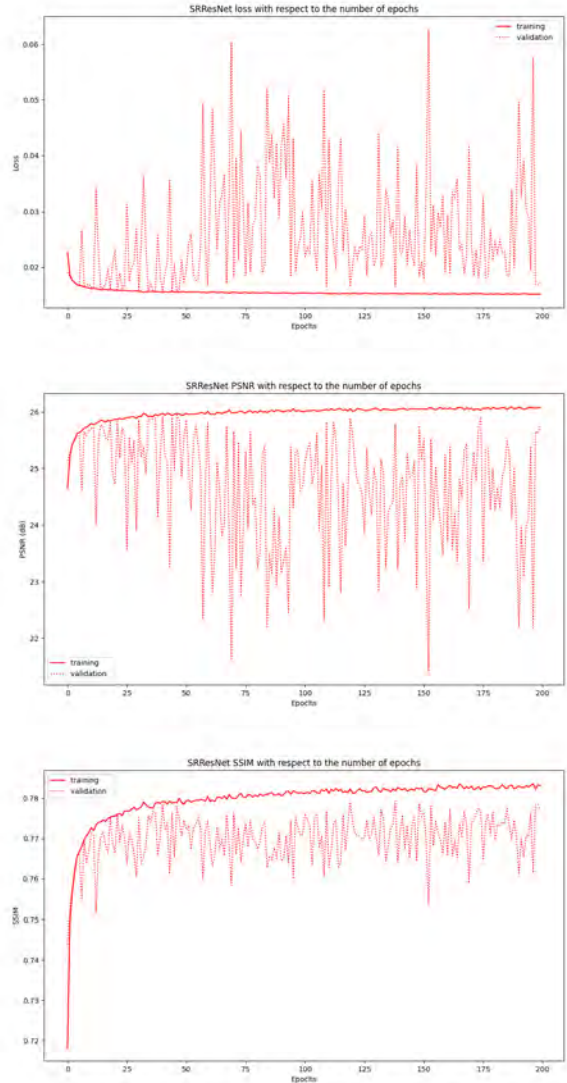


Fig. 4. Loss, PSNR and SSIM with respect to the number of epochs for SRResNet for the base model.

On figure 6 we present an example of super-resolved image. We can see that SRResNet and SRGAN, give superior results compared to bicubic upsampling. Although SRResNet gives better results than plain bicubic upsampling, the final result is overly smooth. SRGAN give much more natural results and is also able to create fine texture details.

Similar results can be seen in example on figure 7. SRResNet gives overly smooth yet better result than plain bicubic upsampling. SRGAN on the other hand presents us with much finer detail (note the grass on the ground and skin of the elephants).

Sometimes there is little difference between SRResNet and SRGAN. This can be partially observed on figure 8. Due to the smooth nature of the texture on the bananas, there is little difference between SRResNet and SRGAN. Note that the perceptually convincing reconstruction of text on the sticker is not well realized.

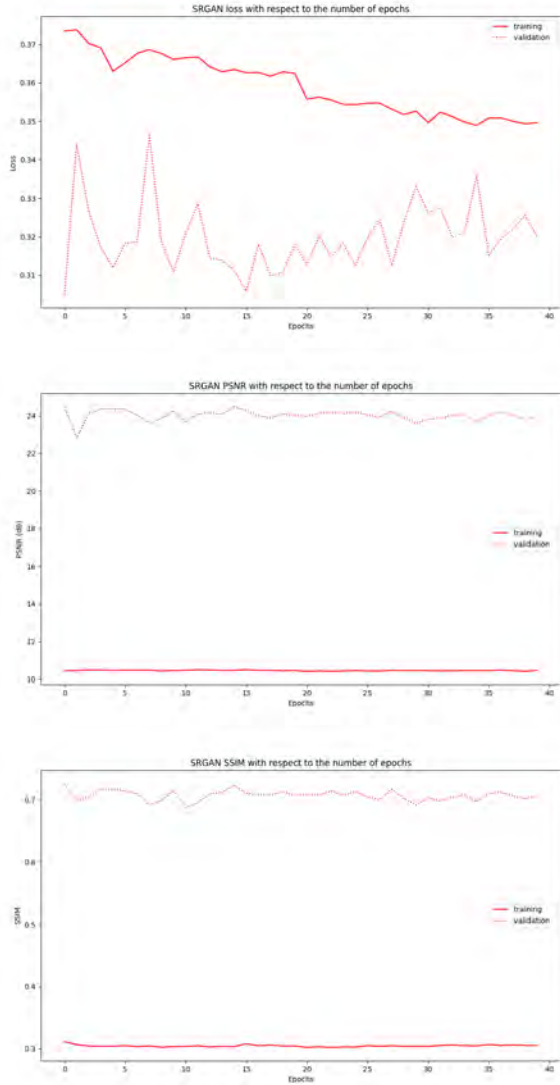


Fig. 5. Loss, PSNR and SSIM with respect to the number of epochs for SRGAN for the base model.

Final example is presented on figure 9. Similarly to previous examples SRResNet gives slightly smoother results than SRGAN. This time one might argue that SRResNet gives more pleasing results when compared to SRGAN.

B. Epochs

Here we explore the effect of increasing number of epochs during the training phase.

The evaluation results of the increasing number epochs are presented in table II. We can observe that SRResNet gives better numerical results across the all testing datasets.

Loss, PSNR and SSIM with respect to the number of epochs for SRResNet can be seen on figure 10. Loss, PSNR and SSIM with respect to the number of epochs for SRGAN can be seen on figure 11. We can observe a steady decrease in loss for both SRResNet and SRGAN. We also see that increasing



Fig. 6. An example of super-resolved image from Set14. Top left is LR image, top right is SR image using SRResNet, bottom left is SR image using SRGAN, bottom right is original image in HR.

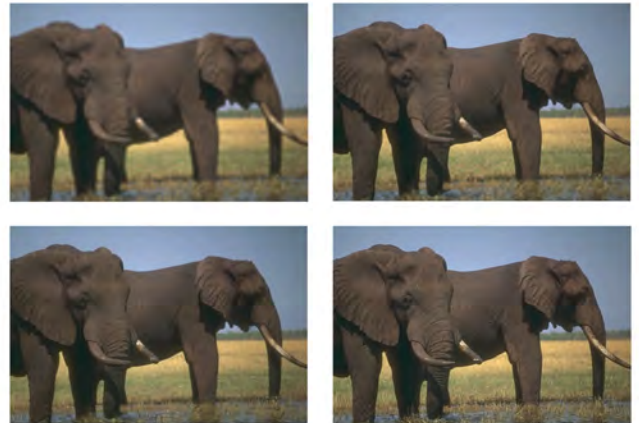


Fig. 7. An example of super-resolved image from BSD100. Top left is LR image, top right is SR image using SRResNet, bottom left is SR image using SRGAN, bottom right is original image in HR.

	COCO Loss	COCO PSNR	COCO SSIM	Set5 Loss	Set5 PSNR	Set5 SSIM
SRResNet	0.0183	27.2895	0.7784	0.0088	30.1063	0.8914
SRGAN	0.1945	25.2257	0.7048	0.2776	27.0301	0.8204
	Set14 Loss	Set14 PSNR	Set14 SSIM	BSD100 Loss	BSD100 PSNR	BSD100 SSIM
SRResNet	0.0114	28.5675	0.7990	0.0153	27.1376	0.7491
SRGAN	0.1834	25.9390	0.7180	0.1787	25.1277	0.6666

TABLE II
EVALUATION RESULTS OF INCREASED NUMBER OF EPOCHS.

the number of epochs helped in improving the model and diminishing the gap between training and validation losses.



Fig. 8. An example of super-resolved image from COCO test 2014. Top left is LR image, top right is SR image using SRResNet, bottom left is SR image using SRGAN, bottom right is original image in HR.



Fig. 9. An example of super-resolved image from Set5. Top left is LR image, top right is SR image using SRResNet, bottom left is SR image using SRGAN, bottom right is original image in HR.

C. Batch size

Here we explore the effect of the different batch sizes used during the training phase.

The evaluation results of the different batch sizes are presented in table III. We can observe that SRResNet, with a batch size of 16 gives overall best numerical results across the all testing datasets.

Loss, PSNR and SSIM with respect to the number of epochs for SRResNet can be seen on figure 12. Loss, PSNR and SSIM

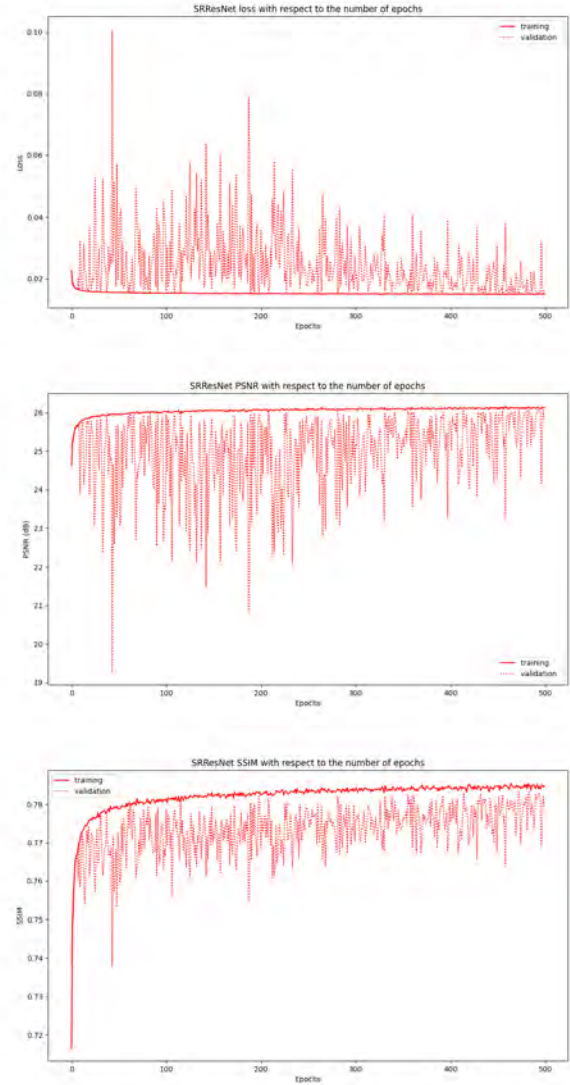


Fig. 10. Loss, PSNR and SSIM with respect to the number of epochs for SRResNet for increased number of epochs.

with respect to the number of epochs for SRGAN can be seen on figure 13. We can observe a steady decrease in loss for both SRResNet and SRGAN. The lower number of batch sizes seem to be beneficial for SRResNet loss and PSNR, while SRResNet SSIM seems best at around 16. Best SRGAN loss and SSIM is achieved with larger batch sizes, while the PSNR looks best with smaller batch sizes.

D. Optimizer learning rate

Here we explore the effect of the different optimizers and different learning rates used during the training phase.

The evaluation results of the Adam optimizer and different learning rates are presented in table IV. We can observe that using Adam optimizer with learning rate 10^{-4} yields best results.

Loss, PSNR and SSIM with respect to the number of

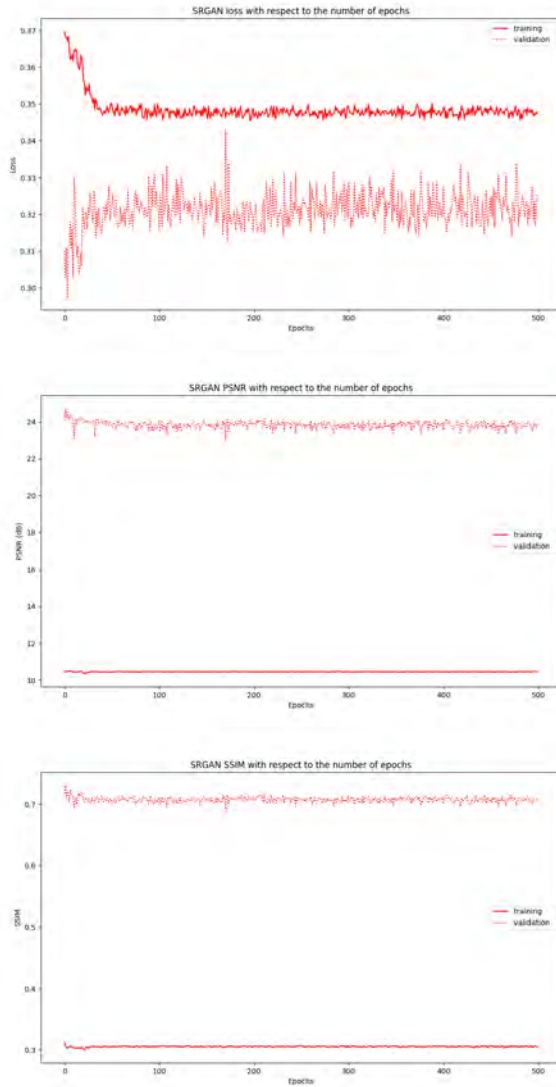


Fig. 11. Loss, PSNR and SSIM with respect to the number of epochs for SRGAN for increased number of epochs.

epochs for SRResNet using Adam optimizer can be seen on figure 14. Loss, PSNR and SSIM with respect to the number of epochs for SRGAN using Adam optimizer can be seen on figure 15. We can confirm that using Adam optimizer with learning rate 10^{-4} yields best results for both SRResNet and SRGAN. Increasing the learning rate has a negative impact on performance, while decreasing the learning rate has a positive impact on performance.

The evaluation results of the SGD optimizer and different learning rates are presented in table V. We can observe that using SGD optimizer with learning rate 10^{-1} yields best results on all testing datasets.

Loss, PSNR and SSIM with respect to the number of epochs for SRResNet using SGD optimizer can be seen on figure 16. Loss, PSNR and SSIM with respect to the number of epochs for SRGAN using SGD optimizer can be seen

	COCO Loss	COCO PSNR	COCO SSIM	Set5 Loss	Set5 PSNR	Set5 SSIM
SRResNet, 1	0.0143	27.8191	0.7782	0.0050	31.6714	0.8969
SRGAN, 1	0.2420	19.6681	0.6479	0.2208	28.5091	0.8620
SRResNet, 2	0.0139	27.9936	0.7858	0.0048	31.7809	0.9004
SRGAN, 2	0.2365	16.1535	0.5405	0.2465	24.8239	0.8378
SRResNet, 4	0.0137	28.0711	0.7868	0.0046	31.9375	0.9015
SRGAN, 4	0.3593	15.7794	0.5945	0.4470	13.3654	0.6416
SRResNet, 8	0.0137	28.0657	0.7877	0.0046	31.9493	0.9017
SRGAN, 8	0.2156	22.9547	0.6767	0.3313	20.6688	0.7466
SRResNet, 16	0.0136	28.1025	0.7883	0.0046	31.9345	0.9023
SRGAN, 16	0.1905	25.8082	0.7057	0.2529	29.1558	0.8324
SRResNet, 32	0.0276	27.0769	0.7799	0.0046	31.9011	0.9007
SRGAN, 32	0.1945	25.3633	0.7006	0.2484	29.0209	0.8382
SRResNet, 64	0.0137	28.0807	0.7873	0.0047	31.9276	0.9015
SRGAN, 64	0.1850	25.7564	0.7086	0.2374	29.1578	0.8449
SRResNet, 128	0.0137	28.0664	0.7872	0.0047	31.9168	0.9015
SRResNet, 256	0.1925	25.5932	0.7023	0.2397	29.1346	0.8411
SRResNet, 512	0.0139	28.0023	0.7852	0.0048	31.8391	0.9005
SRResNet, 1024	0.1862	25.4428	0.7020	0.2419	28.9599	0.8410
SRResNet, 2048	0.0141	27.9223	0.7832	0.0049	31.7255	0.8993
SRResNet, 4096	0.1833	25.4843	0.7029	0.2449	29.1373	0.8444
SRResNet, 1	0.0118	28.4080	0.7926	0.0135	27.3937	0.7475
SRGAN, 1	0.2027	22.6563	0.7173	0.2393	20.1547	0.6132
SRResNet, 2	0.0115	28.5317	0.7977	0.0132	27.5089	0.7539
SRGAN, 2	0.2389	20.2460	0.6603	0.2649	18.7091	0.5400
SRResNet, 4	0.0115	28.5848	0.7980	0.0130	27.5693	0.7546
SRGAN, 4	0.2629	20.4511	0.6664	0.3252	19.7112	0.6073
SRResNet, 8	0.0114	28.6114	0.7989	0.0130	27.5384	0.7553
SRGAN, 8	0.1932	24.6153	0.7041	0.1901	24.2988	0.6561
SRResNet, 16	0.0114	28.6075	0.7993	0.0130	27.5998	0.7559
SRGAN, 16	0.1793	26.2713	0.7173	0.1781	25.3709	0.6676
SRResNet, 32	0.0119	28.5291	0.7982	0.0132	27.5172	0.7550
SRResNet, 64	0.0115	28.5926	0.7964	0.0131	27.5150	0.7551
SRResNet, 128	0.0114	28.6219	0.7989	0.0130	27.5743	0.7551
SRGAN, 128	0.1781	28.8796	0.7043	0.1775	24.9667	0.6583
SRResNet, 256	0.0115	28.6023	0.7987	0.0131	27.5668	0.7551
SRResNet, 512	0.1808	25.7765	0.6992	0.1882	24.8520	0.6498
SRResNet, 1024	0.0116	28.5298	0.7967	0.0131	27.5262	0.7536
SRResNet, 2048	0.1782	25.6421	0.7005	0.1818	24.5470	0.6457
SRResNet, 4096	0.0118	28.4319	0.7953	0.0133	27.4700	0.7523
SRResNet, 8192	0.1780	25.7628	0.7017	0.1787	24.4982	0.6435

TABLE III
EVALUATION RESULTS OF THE DIFFERENT BATCH SIZES.

	COCO Loss	COCO PSNR	COCO SSIM	Set5 Loss	Set5 PSNR	Set5 SSIM
SRResNet, 1e-1	1.4228	13.1536	0.4244	1.5209	12.6921	0.3888
SRGAN, 1e-1	0.2787	12.1388	0.4128	0.3264	12.1870	0.3806
SRResNet, 1e-2	0.4239	12.7827	0.4007	0.3967	13.1758	0.4630
SRGAN, 1e-2	0.2787	12.1388	0.4128	0.3264	12.1870	0.3806
SRResNet, 1e-3	0.0138	28.0172	0.7867	0.0047	31.8344	0.9008
SRGAN, 1e-3	0.1995	25.3966	0.6964	0.2556	28.8420	0.8385
SRResNet, 1e-4	0.0153	28.0184	0.7862	0.0047	31.9811	0.9012
SRGAN, 1e-4	0.1969	24.8958	0.6969	0.2872	26.3803	0.8078
SRResNet, 1e-5	0.0141	27.9295	0.7833	0.0049	31.7169	0.8990
SRGAN, 1e-5	0.2128	24.4837	0.7093	0.3560	24.9329	0.8307
SRResNet, 1e-6	0.0155	27.4038	0.7657	0.0060	30.9854	0.8848
SRGAN, 1e-6	0.2115	24.3678	0.7325	0.3368	23.5866	0.8350
SRResNet, 1e-7	0.0179	26.6381	0.7394	0.0095	29.5271	0.8458
SRGAN, 1e-7	0.1977	25.7059	0.7565	0.2653	26.5974	0.8601
SRResNet, 1e-1	1.3285	13.8049	0.3816	1.4595	14.6831	0.3827
SRGAN, 1e-1	0.2395	12.4859	0.3669	0.2602	13.7321	0.3799
SRResNet, 1e-2	0.4709	12.1461	0.3631	0.4628	12.3999	0.3143
SRGAN, 1e-2	0.2395	12.4859	0.3669	0.2602	13.7321	0.3799
SRResNet, 1e-3	0.0114	28.6002	0.7989	0.0131	27.5206	0.7552
SRGAN, 1e-3	0.1925	25.5104	0.6952	0.1921	24.9167	0.6557
SRResNet, 1e-4	0.0114	28.6433	0.7995	0.0130	27.5729	0.7557
SRGAN, 1e-4	0.1941	25.6344	0.7101	0.1829	24.9255	0.6595
SRResNet, 1e-5	0.0120	28.3930	0.7942	0.0132	27.4897	0.7525
SRGAN, 1e-5	0.2106	25.9821	0.7258	0.1926	24.9634	0.6649
SRResNet, 1e-6	0.0129	27.9709	0.7810	0.0141	27.1546	0.7394
SRResNet, 1e-7	0.2056	26.2074	0.7590	0.1858	25.2086	0.6983
SRResNet, 1e-8	0.0159	26.9259	0.7522	0.0159	26.5549	0.7179
SRResNet, 1e-9	0.1812	27.3106	0.7838	0.1804	26.2667	0.7324

TABLE IV
EVALUATION RESULTS OF ADAM OPTIMIZER AT DIFFERENT LEARNING RATES.

on figure 17. We can confirm that using SGD optimizer with learning rate 10^{-1} yields best results for SRResNet. Here decreasing the learning rate has a negative impact on performance, while increasing the learning rate has a positive impact on performance. Similarly we can say about SRGAN, with an exception of learning rate 10^{-2} , which gives poor results.

The evaluation results of the ADAGRAD optimizer and different learning rates are presented in table VI. We can observe that using ADAGRAD we minimize loss when using learning rate 10^{-3} . PSNR is minimized when using learning rate 10^{-6} and SSIM is minimized when using learning rate 10^{-7} .

Loss, PSNR and SSIM with respect to the number of epochs

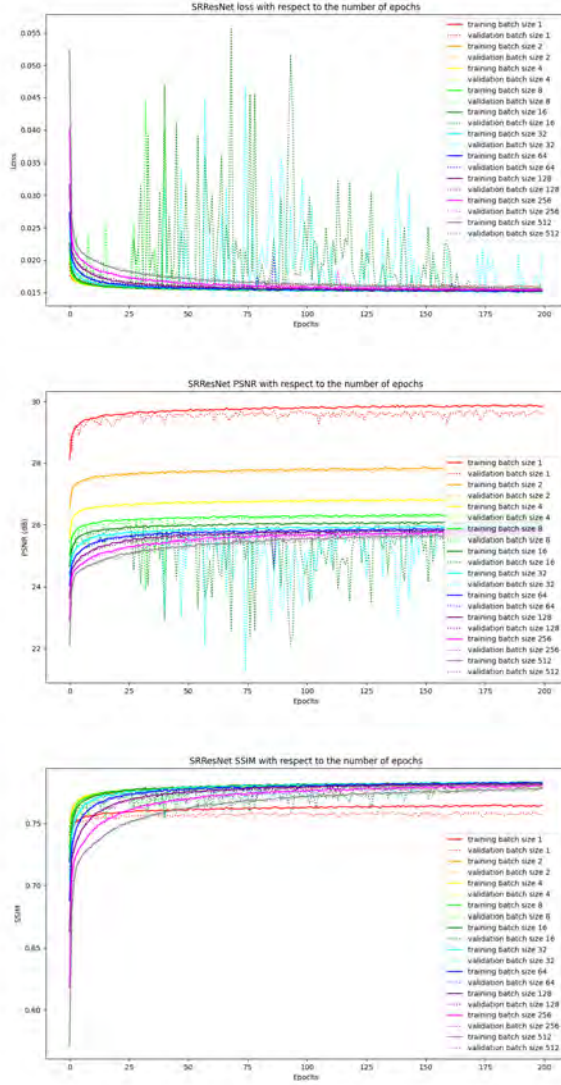


Fig. 12. Loss, PSNR and SSIM with respect to the number of epochs for SRResNet for different batch sizes.

for SRResNet using ADAGRAD optimizer can be seen on figure 18. Loss, PSNR and SSIM with respect to the number of epochs for SRGAN using ADAGRAD optimizer can be seen on figure 19. We can confirm that using ADAGRAD optimizer with learning rate 10^{-3} yields best results for SRResNet and SRGAN. Using learning rate 10^{-1} has negative impact on performance.

The evaluation results of the RMSPROP optimizer and different learning rates are presented in table VII. We can observe that using RMSPROP optimizer with learning rate 10^{-4} yields best results across all testing datasets.

Loss, PSNR and SSIM with respect to the number of epochs for SRResNet using RMSPROP optimizer can be seen on figure 20. Loss, PSNR and SSIM with respect to the number of epochs for SRGAN using RMSPROP optimizer can be seen on figure 21. We can confirm that using RMSPROP optimizer

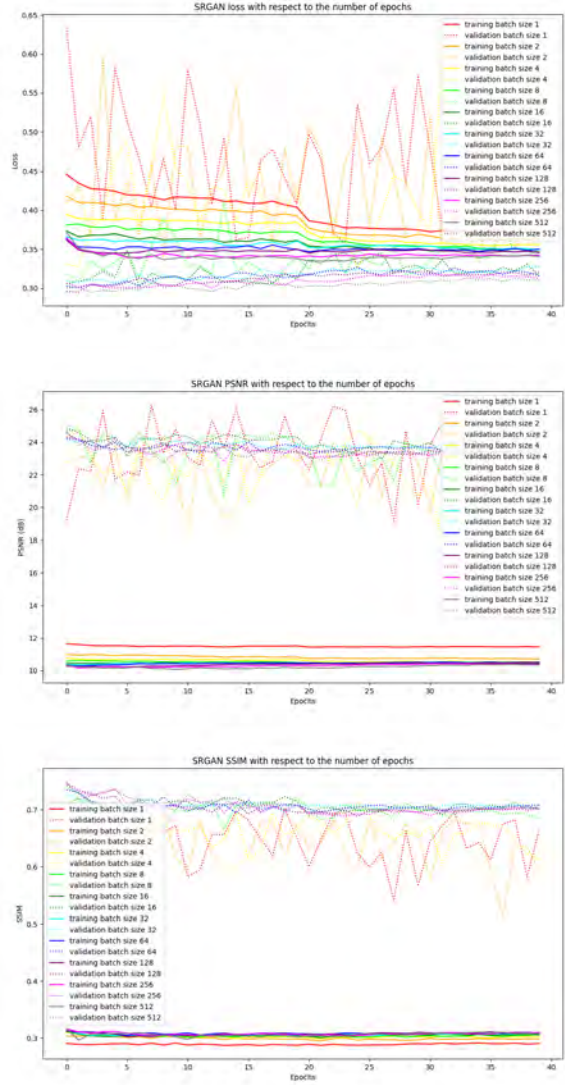


Fig. 13. Loss, PSNR and SSIM with respect to the number of epochs for SRGAN for different batch sizes.

with learning rate 10^{-4} yields best results for SRResNet and SRGAN. Using learning rate 10^{-1} or 10^{-2} has negative impact on performance and causes unstable results.

E. Number of residual blocks

Here we explore the effect of different number of residual blocks in SRResNet and generator part of SRGAN.

The evaluation results of the different number of residual blocks are presented in table VIII. We can observe that using 16 residual blocks yields best results on all testing datasets.

Loss, PSNR and SSIM with respect to the number of epochs for SRResNet can be seen on figure 22. Loss, PSNR and SSIM with respect to the number of epochs for SRGAN can be seen on figure 23. We observe that using 8 or 16 residual blocks gives the most stable results on both SRResNet and SRGAN. We can also see that increasing the residual blocks further increases the loss on the validation set.

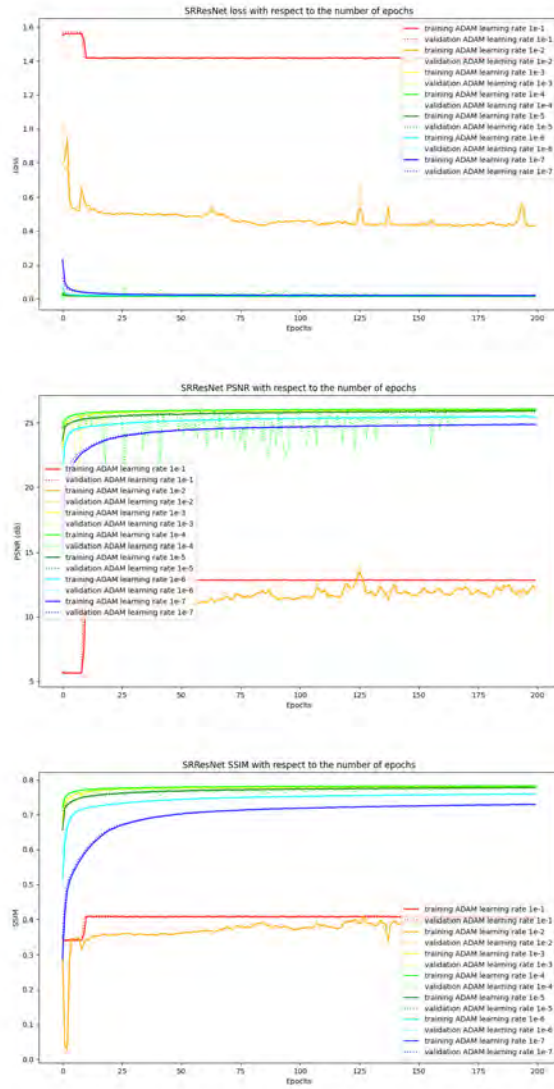


Fig. 14. Loss, PSNR and SSIM with respect to the number of epochs for SRResNet for using Adam optimizer at different learning rates.

V. CONCLUSION

We have described a deep residual network SRResNet and a generative adversarial network SRGAN. SRResNet can give us plausible results in certain circumstances, where an image in question is naturally smooth. SRGAN uses combination of content loss and adversarial loss to give us much more detail results.

We can note that the ideal loss function depends on the application. Approaches that try to create finer texture detail might be generally desirable, but less suited for medical applications. The perceptually convincing reconstruction of text is also very challenging to realize. The development of content loss functions that describe image spatial content, that are more invariant to changes in pixel space, will further improve photo-realistic image SR results.

It is also worth mentioning that standard quantitative mea-

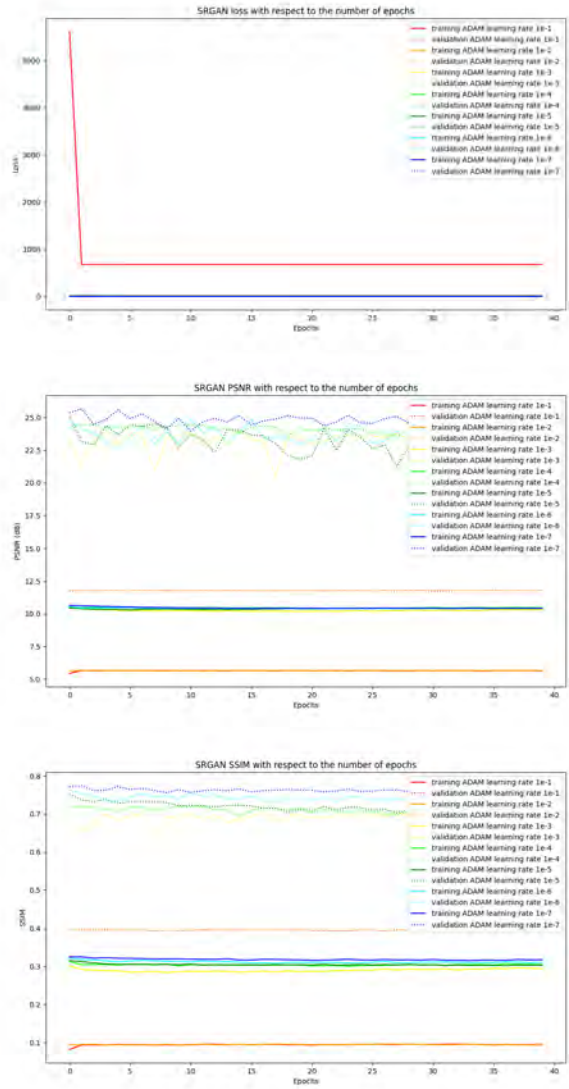


Fig. 15. Loss, PSNR and SSIM with respect to the number of epochs for SRGAN for using Adam optimizer at different learning rates.

asures (such as PSNR and SSIM) fail to capture and accurately assess the image quality with respect to the human visual system perception.

Of particular importance when aiming for photo-realistic solutions to the SR problem is the choice of the content loss

REFERENCES

- [1] Q. Yang, R. Yang, J. Davis, and D. Nistér, "Spatial-depth super resolution for range images," *2007 IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1–8, 2007. [Online]. Available: <https://api.semanticscholar.org/CorpusID:6788025>
- [2] W. Zou and P. C. Yuen, "Very low resolution face recognition problem," vol. 21, 10 2010, pp. 1 – 6.
- [3] K. Nasrollahi and T. Moeslund, "Super-resolution: A comprehensive survey," *Machine Vision and Applications*, vol. 25, pp. 1423–1468, 08 2014.
- [4] C.-Y. Yang, C. Ma, and M.-H. Yang, "Single-image super-resolution: A benchmark," 09 2014, pp. 372–386.

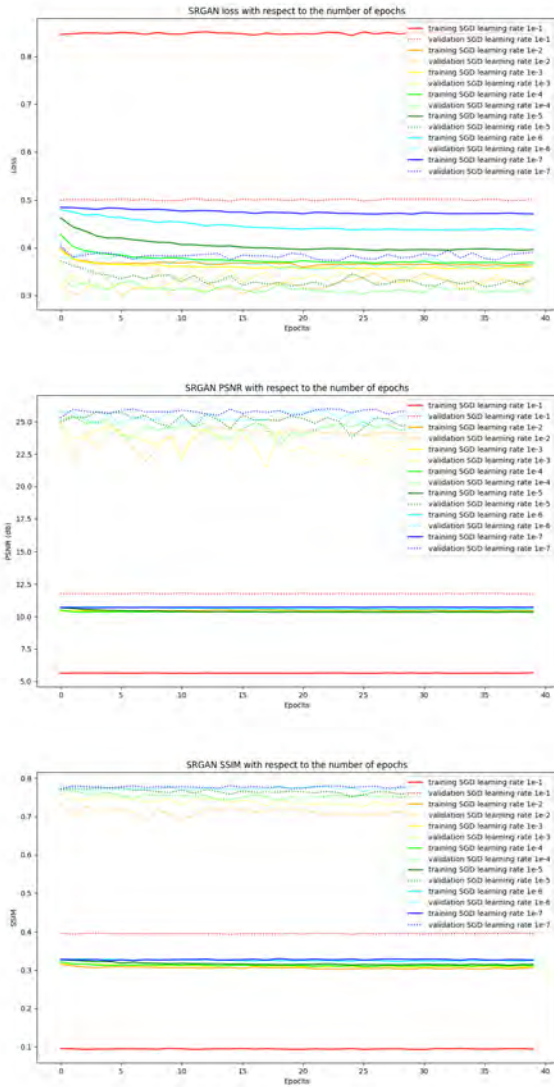


Fig. 17. Loss, PSNR and SSIM with respect to the number of epochs for SRGAN for using SGD optimizer at different learning rates.

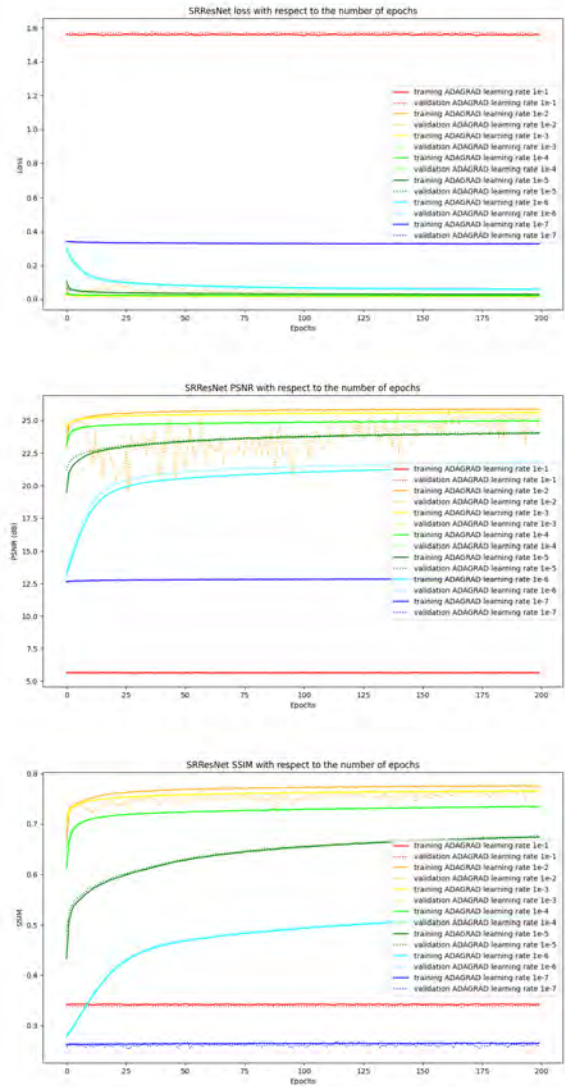


Fig. 18. Loss, PSNR and SSIM with respect to the number of epochs for SRResNet for using ADAGRAD optimizer at different learning rates.

[6] P. Gupta, P. Srivastava, S. Bhardwaj, and V. Bhateja, "A modified psnr metric based on hvs for quality assessment of color images," 12 2011.

[7] Z. Wang, A. Bovik, H. Sheikh, and E. Simoncelli, "Image quality assessment: From error visibility to structural similarity," *Image Processing, IEEE Transactions on*, vol. 13, pp. 600 – 612, 05 2004.

[8] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv 1409.1556*, 09 2014.

[9] J. Johnson, A. Alahi, and L. Fei-Fei, "Perceptual losses for real-time style transfer and super-resolution," vol. 9906, 10 2016, pp. 694–711.

[10] J. Bruna, P. Sprechmann, and Y. Lecun, "Super-resolution with deep convolutional sufficient statistics," 11 2015.

[11] S. Borman and R. Stevenson, "Super-resolution from image sequences - a review," 09 1998, pp. 374 – 378.

[12] S. Farsiu, M. Robinson, M. Elad, and P. Milanfar, "Fast and robust multiframe super resolution," *IEEE transactions on image processing : a publication of the IEEE Signal Processing Society*, vol. 13, pp. 1327–44, 11 2004.

[13] C. Duchon, "Lanczos filtering in one and two dimensions," *Journal of Applied Meteorology - J APPL METEOROL*, vol. 18, pp. 1016–1022, 08 1979.

[14] J. P. Allebach and P. W. Wong, "Edge-directed interpolation," *Proceedings of 3rd IEEE International Conference on Image Processing*, vol. 3, pp. 707–710 vol.3, 1996. [Online]. Available: <https://api.semanticscholar.org/CorpusID:16375380>

[15] X. Li and M. Orchard, "New edge-directed interpolation," *Image Processing, IEEE Transactions on*, vol. 10, pp. 1521 – 1527, 11 2001.

[16] W. Freeman, E. Pasztor, and O. Carmichael, "Learning low-level vision," *International Journal of Computer Vision*, vol. 40, pp. 25–47, 10 2000.

[17] W. Freeman, T. Jones, and E. Pasztor, "Example-based super-resolution," *IEEE Computer Graphics and Applications*, vol. 22, no. 2, pp. 56–65, 2002.

[18] Z. Wang, D. Liu, J. Yang, W. Han, and T. Huang, "Deep networks for image super-resolution with sparse prior," 2015.

[19] K. Gregor and Y. LeCun, "Learning fast approximations of sparse coding," in *International Conference on Machine Learning*, 2010. [Online]. Available: <https://api.semanticscholar.org/CorpusID:13333501>

[20] C. Dong, C. C. Loy, K. He, and X. Tang, "Image super-resolution using deep convolutional networks," 2015.

[21] —, "Image super-resolution using deep convolutional networks," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 38, no. 2, pp. 295–307, 2016.

[22] M. Mathieu, C. Couprie, and Y. Lecun, "Deep multi-scale video predic-

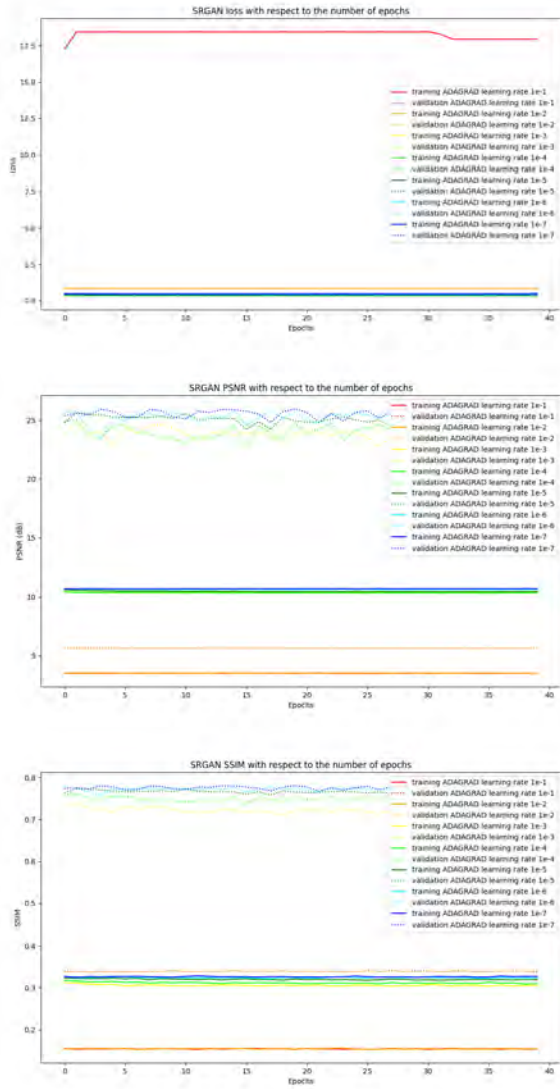


Fig. 19. Loss, PSNR and SSIM with respect to the number of epochs for SRGAN for using ADAGRAD optimizer at different learning rates.

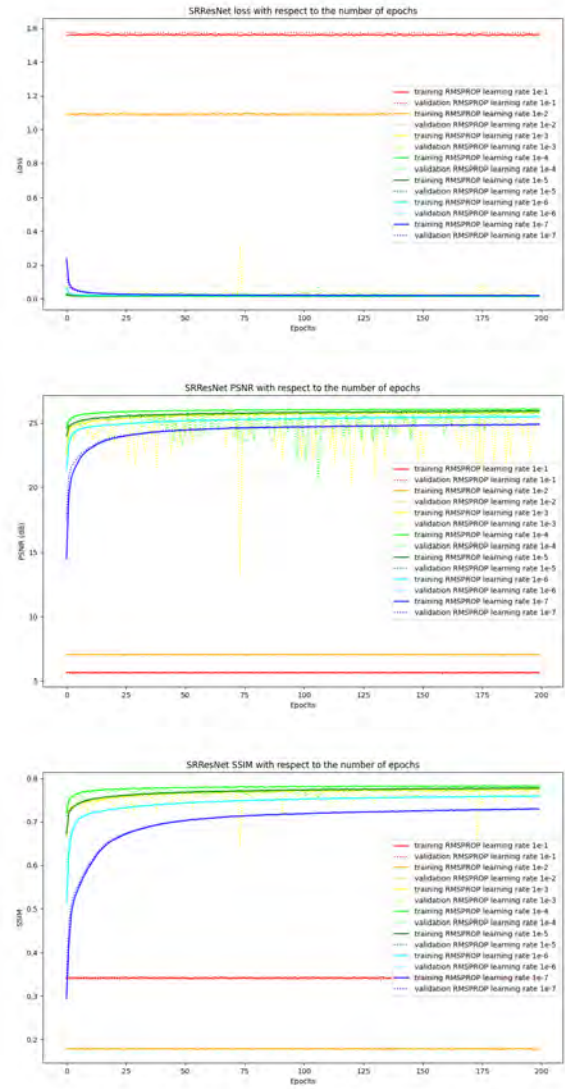


Fig. 20. Loss, PSNR and SSIM with respect to the number of epochs for SRResNet for using RMSPROP optimizer at different learning rates.

tion beyond mean square error,” 11 2016.

- [23] E. Denton, S. Chintala, A. Szlam, and R. Fergus, “Deep generative image models using a laplacian pyramid of adversarial networks,” 2015.
- [24] A. Dosovitskiy and T. Brox, “Generating images with perceptual similarity metrics based on deep networks,” 2016.
- [25] L. A. Gatys, A. S. Ecker, and M. Bethge, “Texture synthesis using convolutional neural networks,” 2015.
- [26] L. Gatys, A. Ecker, and M. Bethge, “Image style transfer using convolutional neural networks,” 06 2016, pp. 2414–2423.
- [27] C. Ledig, L. Theis, F. Huszar, J. Caballero, A. Cunningham, A. Acosta, A. Aitken, A. Tejani, J. Totz, Z. Wang, and W. Shi, “Photo-realistic single image super-resolution using a generative adversarial network,” 2017.
- [28] C. Li and M. Wand, “Combining markov random fields and convolutional neural networks for image synthesis,” 2016.
- [29] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, “Generative adversarial networks,” *Advances in Neural Information Processing Systems*, vol. 3, 06 2014.
- [30] W. Shi, J. Caballero, F. Huszár, J. Totz, A. P. Aitken, R. Bishop, D. Rueckert, and Z. Wang, “Real-time single image and video super-resolution using an efficient sub-pixel convolutional neural network,” in

2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2016, pp. 1874–1883.

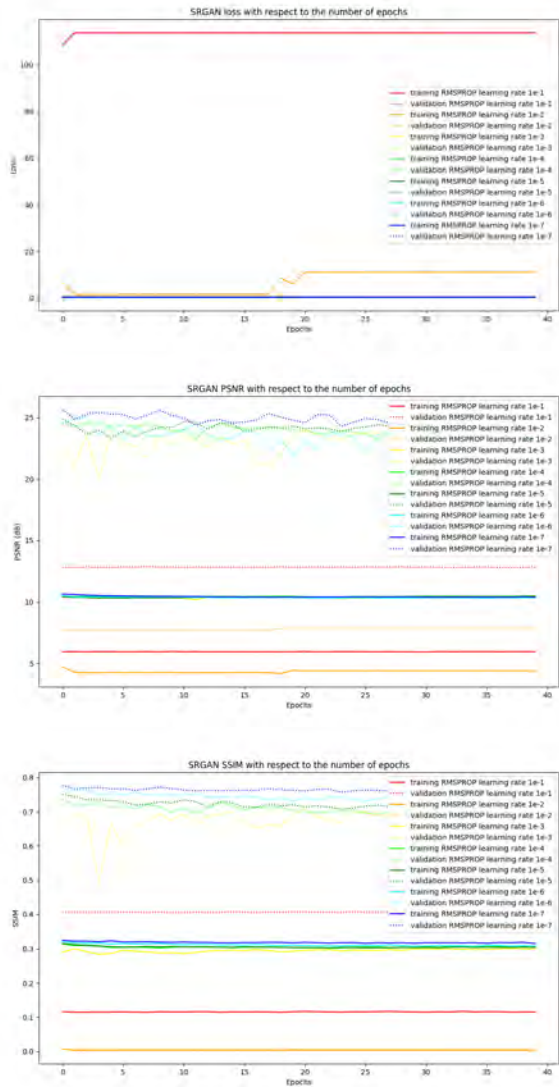


Fig. 21. Loss, PSNR and SSIM with respect to the number of epochs for RMSPROP for using Adam optimizer at different learning rates.

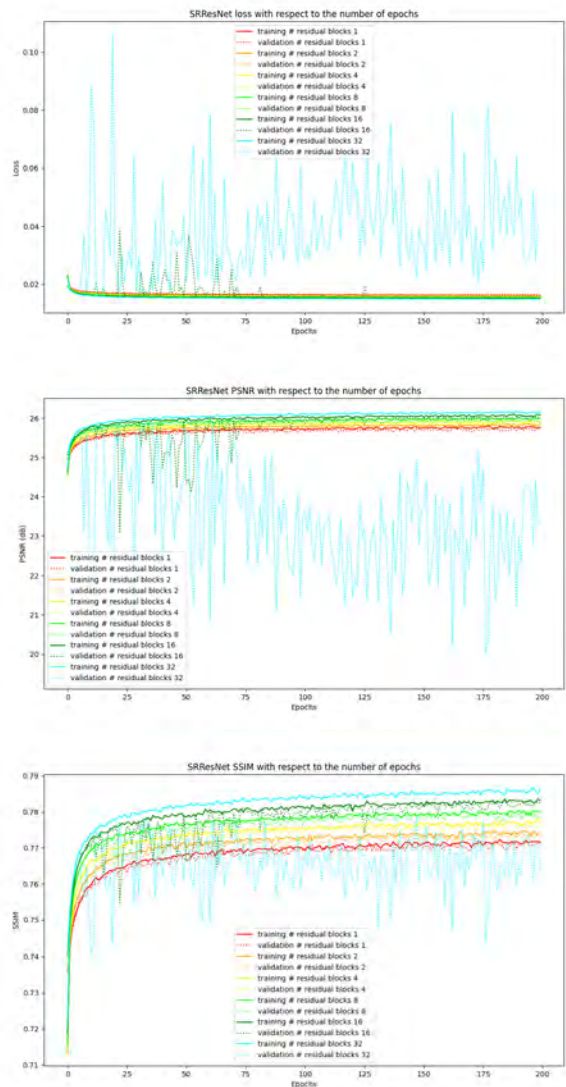


Fig. 22. Loss, PSNR and SSIM with respect to the number of epochs for SRResNet for the different number of residual blocks.

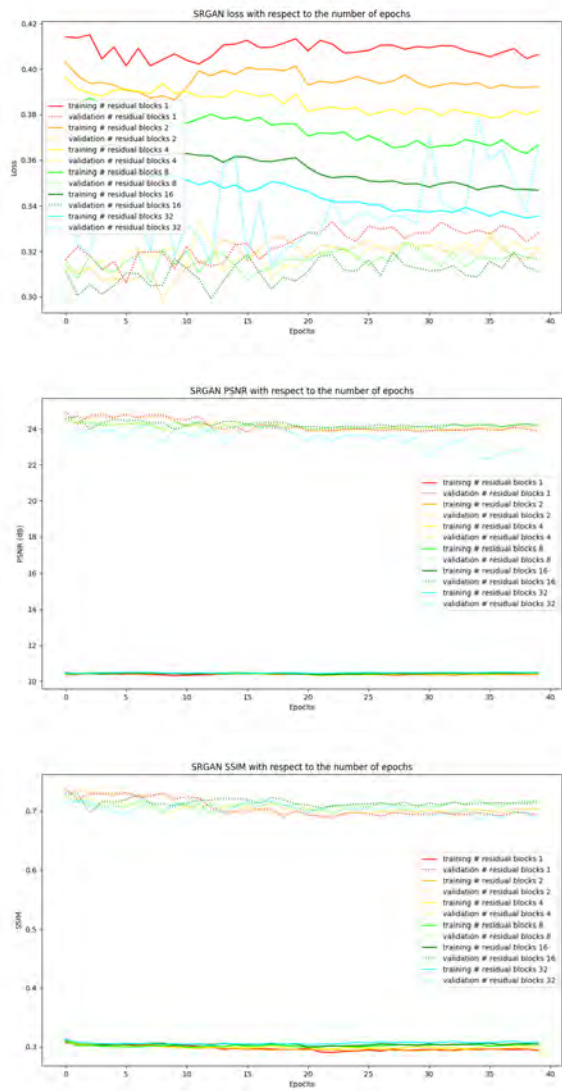


Fig. 23. Loss, PSNR and SSIM with respect to the number of epochs for SRGAN for the different number of residual blocks.