

Test-Time Adaptation for Super-Resolution: You Only Need to Overfit on a Few More Images

Mohammad Saeed Rad Thomas Yu Behzad Bozorgtabar Jean-Philippe Thiran
 Signal Processing Lab (LTS5), EPFL, Lausanne, Switzerland
 {saeed.rad, firstname.lastname}@epfl.ch

Abstract

Existing reference (RF)-based super-resolution (SR) models try to improve perceptual quality in SR under the assumption of the availability of high-resolution RF images paired with low-resolution (LR) inputs at testing. As the RF images should be similar in terms of content, colors, contrast, etc. to the test image, this hinders the applicability in a real scenario. Other approaches to increase the perceptual quality of images, including perceptual loss and adversarial losses, tend to dramatically decrease fidelity to the ground-truth through significant decreases in PSNR/SSIM. Addressing both issues, we propose a simple yet universal approach to improve the perceptual quality of the HR prediction from a pre-trained SR network on a given LR input by further fine-tuning the SR network on a subset of images from the training dataset with similar patterns of activation as the initial HR prediction, with respect to the filters of a feature extractor. In particular, we show the effects of fine-tuning on these images in terms of the perceptual quality and PSNR/SSIM values. Contrary to perceptually driven approaches, we demonstrate that the fine-tuned network produces a HR prediction with both greater perceptual quality and minimal changes to the PSNR/SSIM with respect to the initial HR prediction. Further, we present novel numerical experiments concerning the filters of SR networks, where we show through filter correlation, that the filters of the fine-tuned network from our method are closer to “ideal” filters, than those of the baseline network or a network fine-tuned on random images.

1. Introduction

Super-resolution (SR) is the ill-posed problem of transforming low-resolution (LR) images (I_{LR}) to their high-resolution (HR) counterparts (I_{HR}) [29, 34, 1, 10, 23, 28]. A common way to model the interaction between LR and HR images can be formulated as $I_{LR} = (I_{HR} * \mathbf{k}) \downarrow_s + N$, where $*$ denotes convolution, \mathbf{k} is the blur kernel, \downarrow_s de-

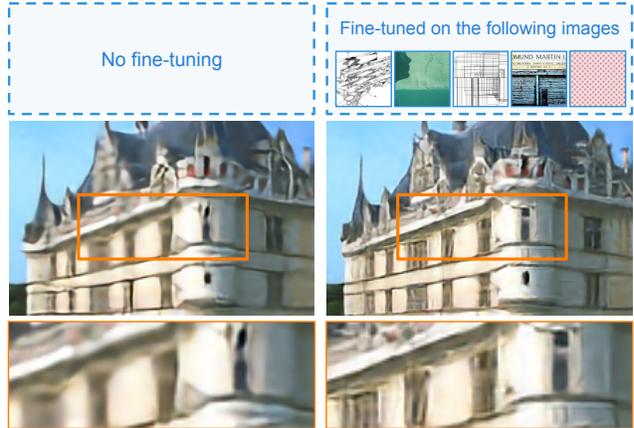


Figure 1. We demonstrate how we can improve the perceptual quality of Super-Resolution images produced by a generic SR network and a given LR image by fine-tuning the network on specific images which activate the same filters of a pre-trained feature extractor as those activated by the initial SR prediction. Left: Initial SR predictions from the baseline network, right: Predictions from the network after fine-tuning for a few iterations on selected images by our method. Zoom in for the best view.

notes downsampling by a factor s , and N is noise. In this paper, we focus on a common setting for SR, where the down-sampling kernel is known and is a bicubic downscaling kernel [29].

In this setting, deep learning algorithms [32, 8, 39, 21, 22] have made remarkable progress in image super-resolution that aim to obtain a I_{HR} output from one of its I_{LR} versions by leveraging the power of deep convolutional neural networks. Going even further, in the field of reference (RF)-based SR, an external high-resolution reference image is provided, where the reference image and I_{HR} share similar textures and qualities [41, 42, 14, 36]. In this way, the networks are trained to leverage additional information from the reference HR image. **This has the drawback of assuming the existence of and finding HR images similar to a given LR image -in terms of content, colors, contrast** as well as the increased size of the net-

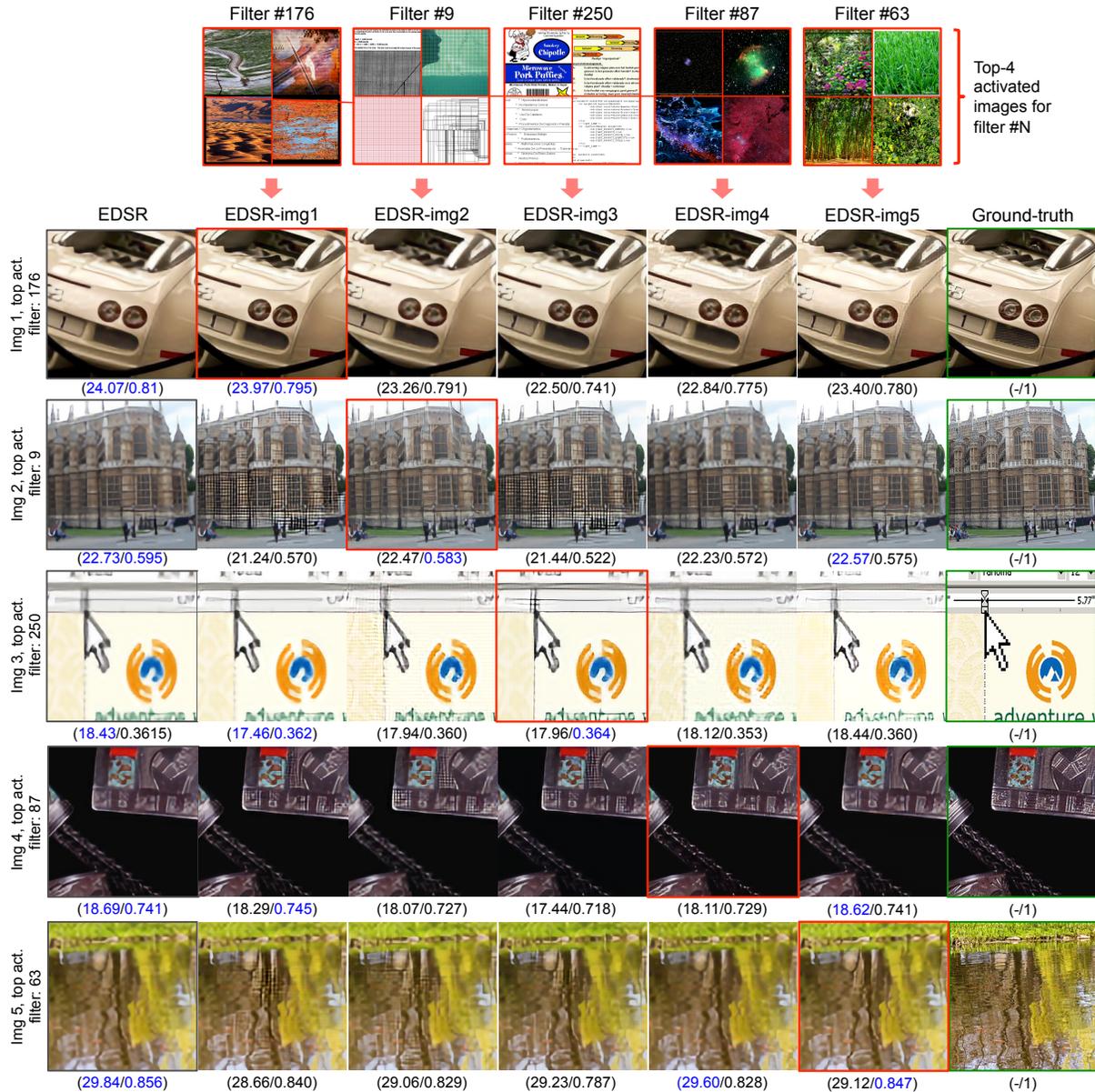


Figure 2. We demonstrate the effect of fine-tuning on images, which maximally activate specific filters in a pre-trained classification network with respect to perceptual quality and PSNR/SSIM values. The first column shows the initial HR predictions from the baseline network while subsequent columns show predictions from the network after fine-tuning on the images bordered by red at the top. Note that in each row, **the network fine-tuned on the image set which shares the filter of maximal activation with the initial HR prediction gives the best perceptual quality without affecting the PSNR or SSIM significantly**. Fine-tuning on image sets which maximally activate different filters results in oversmoothing or image artifacts as compared to the ground truth. Two best values are in blue. **Please zoom in on the screen.**

works trained to incorporate the additional HR input.

In the SR literature, pixel-based metrics, which compare predicted HR images to the ground truth HR image such as the peak signal to noise ratio (PSNR) or structural similarity index (SSIM) are commonly used to judge the performance of SR methods [29]. However, it is known

that optimizing neural networks for PSNR, SSIM, or other pixel-based metrics generally result in over-smoothed, perceptually unappealing HR images [4, 15, 18]. In fact, [4] shows that there is a mathematical tradeoff between performance on these pixel-based metrics and perceptual quality. However, **we note that in theory, a perfect reconstruct-**

tion would have the highest performance on both pixel-based metrics and perceptual quality. Strategies to increase perceptual quality include training networks with a perceptual loss [15], which computes the distance between predicted and ground truth images in the feature space using a pre-trained classification network. Generative adversarial networks (GANs) [11] are also used to improve perceptual quality [13, 18, 32, 5, 6, 24, 30]. **However, these approaches significantly decrease PSNR, SSIM and other pixel-based metrics with respect to trained networks using only the pixel-wise losses [15, 18, 4].**

In this paper, inspired by RF-based SR and previous analysis of learned filters of classification networks, we propose a novel method to increase the perceptual quality of the output of a generic PSNR-based SR network on a given LR image without significantly affecting the PSNR or SSIM. This is done through test-time adaptation of the generic SR network, to tailor it to a given LR image used for testing. Concretely, given an input LR image and the SR network pre-trained using only pixel-wise losses, e.g., L_1 , we first obtain the initial HR prediction from the network. We then fine-tune the network on a few pairs of LR/HR images from the training dataset, where the images are chosen by the similarity of their activations of filters from a pre-trained classification network with respect to the corresponding activations of the initial HR prediction. We show that the perceptual quality of the HR image from the fine-tuned network increases without significantly decreasing the PSNR or SSIM values. Further, we demonstrate that this does not contradict past studies on the trade-off between PSNR and perceptual quality [25, 17], as this results from fine-tuning on images that activate the same filters as the initial LR input. The fine-tuned SR network performs worse on images dissimilar to the LR input; hence, overall performance is in conformity with the trade-off. As shown in Fig. 2, our method can improve perceptual quality with minimal impact on PSNR/SSIM with fine-tuning on images with similar activations as the LR input.

Our contributions are as follows:

- We propose a novel, test-time adaptation method to improve SR, which guides PSNR-based SR networks toward perceptually more compelling images by fine-tuning on selected images at the test-time, **without significant impact on the PSNR or SSIM.**
- To our knowledge, we are the first to investigate how overfitting/fine-tuning on selected images, which differ by what filters in a pre-trained classification network they maximally activate, can change SR reconstructions for better or worse.
- We also show, to our knowledge, novel numerical experiments in the field of SR, where we quantitatively

relate the filters of the pre-trained SR network, the fine-tuned network, and an “ideal” SR network (ideal with respect to the given LR input) to show that our method moves the filters of the pre-trained SR network closer to the “ideal” filters.

2. Overview of the approach

The overview of the proposed method is shown in Fig. 3; the task is to predict an HR image from a given LR input by benefiting from a few more essential images with respect to the pre-trained model. The pipeline can be split into three main steps: First, we construct a reference dataset, namely the Activation dataset, containing essential images for further fine-tuning. Second, we use a novel technique to choose relevant images from the Activation dataset. Finally, we fine-tune the pre-trained SR network on these images and produce the final reconstruction. In what follows, let G , \mathcal{D} denote the baseline SR network and the dataset of paired LR and HR images used to train the SR network, respectively. We present each step in detail as follows:

Construction of Activation dataset We first construct a reference dataset from the HR images of \mathcal{D} by extracting their corresponding activations from the third layer of the VGG classification network [27]. For each channel in the third layer (*conv3*), we order (descending) the images by the channel’s corresponding activation and take the top K images. As there are 256 channels in the third layer, we form a reference dataset of $256 \times K$ HR images. We choose the third layer as the features from this layer have been shown to be more discriminative [7, 38]. As an example, in Fig. 4, we show for different filters in different layers of VGG19 [27], the top nine images by filter activation from a subset of 50 thousand images from ImageNet [9]. We further investigate the effectiveness of using other layers (*conv2*, 4, and 5, in our supplementary material).

Test-Time Adaptation of the SR network We obtain an initial HR prediction from passing LR to G , which we call SR . We pass SR through the third layer of the VGG classification network [27] and note the top M filters with the highest activations. From this list of filters, we can use our reference dataset to define a set of $M \times K$ images where for each of the M filters, we take the top K images in our dataset in terms of activation of the filter. We then fine-tune G on this set of images for a set number of epochs determined by performance on the validation set.

Prediction After fine-tuning G , we again pass the LR image to G to obtain our final HR prediction, which we call activated SR. **The activated SR image is perceptually more convincing than the initial SR, without significant decreases in its PSNR and SSIM values.**

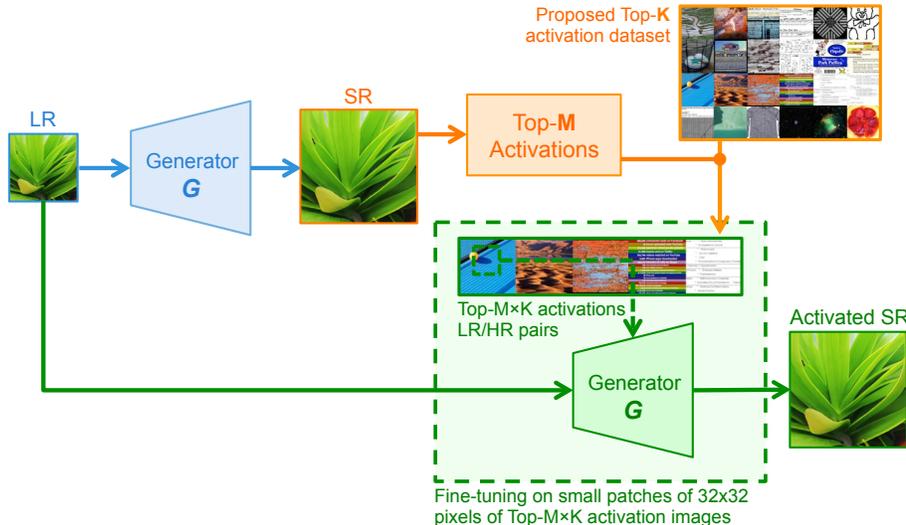


Figure 3. The overview of the proposed method: First, the LR input is passed to the SR network to generate an initial SR prediction. We then find the top M filters of the third layer of the VGG [27] network which are activated by the initial SR prediction. Then, we fine-tune the SR network on a set of $M * K$ images chosen from the training data, which maximally activate the same M filters. Finally, we pass the LR input to the fine-tuned SR network for the final SR prediction. In this example, $K = 1$ and $M = 5$.

3. Image Activations in SR

In the machine learning/computer vision literature, analysis of the activations of neural networks with respect to different inputs is often used for the purposes of understanding/interpretability [2] and extraction of relevant features for downstream processing, for instance, in unsupervised learning [7, 26]. In terms of SR, only perceptual loss uses this analysis by matching the activations, with respect to a layer of a pre-trained classification network, of the HR prediction and the ground-truth, showing the efficacy and importance of these features. We go further by explicitly analyzing the activations of the third layer of VGG19 [27] with respect to a large dataset of 50 thousand images. Then for each filter, we can assign a group of images with the highest activations. As perceptual loss shows that constructing images based on activations can improve perceptual quality, it stands to reason that fine-tuning a network on images that also triggers specific filters can enhance SR reconstructions on images that have similar activations with respect to those filters. Hence, in contrast to perceptual loss, we are able to exploit the analysis of activations by enhancing the perceptual quality of SR on a given LR input by using a set of images **which are visually different from the LR input**, but similar in terms of activation. To the best of our knowledge, we are the first to create and benefit from such a dataset for SR. Fig. 4 shows a few example images from the Activated dataset; the detailed procedure of generating this dataset is presented in section 2. *This dataset is available in supplementary material and we will release this dataset upon acceptance of the paper.*

4. Overfitting: the good, the bad, and the ugly

Throughout this paper, we have used the word “fine-tuning” for continuing the training of a pre-trained SR network on a small set of images. Implicitly, this assumes that such training has a beneficial effect for the purpose of the network, which is to perform SR on a given LR image (“**The good**”). However, as seen in Fig. 2, such fine-tuning could also be labeled as overfitting, since our method only improves reconstructions on images with similar patterns of filter activation as the given LR image; other inputs can result in image artifacts and over smoothing (“**The bad**”). That is, the fine-tuned network no longer generalizes to all image classes. This can be understood in terms of the tradeoff between perceptual quality, and PSNR established in [4]. We conjecture that we are able to gain perceptual quality with minimal changes to PSNR/SSIM precisely because this gain occurs only on images similar in filter activation to those used in the fine-tuning. As both PSNR and perceptual quality can decrease in other images, the overall performance does not contravene the tradeoff. Thus, for a given LR image, overfitting is actually good for improving SR reconstructions. However, we note that the outcome of fine-tuning is dependent on the number of epochs of additional training (“**The ugly**”). Further, while generalization of the network performance is clearly compromised, it is possible for the fine-tuning to have no effect, good or bad, on different classes of images. In Fig. 5, we show the effects of fine-tuning on visual quality, PSNR, and SSIM values as a function of the number of epochs as well as how it can dramatically increase the perceptual quality of some



Figure 4. Top 9 activated images from a subset of 50 thousand images from ImageNet [9] for different filters in the conv1, conv3 and conv5 layers of VGG19 [27], respectively.

images while not affecting others.

It remains to address how overfitting using only a pixel-wise loss can improve perceptual quality. We emphasize that the fine-tuning is done with only L_1 loss; in contrast to perceptual loss or adversarial losses used to improve perceptual quality, only pixel-wise metrics are used in our approach. In Fig. 6, we show a diagram of our hypothesis that **overfitting guides the SR network to a local minimum, where the pixel-wise error is only slightly different, while the perceptual quality is dramatically improved**. As evidence, note that almost the same PSNR is achieved on image b (during the pretraining of the network, before fine-tuning by our approach) and image c (after fine-tuning), but image c is much sharper and realistic.

5. Experiments and results

5.1. Experimental settings

5.1.1 Generator architecture

While our method and experiments can generalize to arbitrary SR networks, we use an EDSR [19] as our baseline generator, which we denote as G . EDSR performs better than other conventional residual SR networks by eliminating some unnecessary modules e.g., batch normalization. This makes it a good candidate to investigate the effectiveness of our proposed approach as many other SR networks incorporate components designed for specific contributions/improvements that may not strictly be necessary. The architecture consists of 32 residual blocks and 256 filters per convolutional layer (more details in supplementary material). We train this network in a single step for 50 epochs, using the L_1 loss function. For the training data, we use a subset of 50 thousand images taken from Imagenet [9]. The Adam optimizer was used for the optimization. The learning rate was set to $1e-3$ and then decayed by a factor of ten every 20 epochs.

5.1.2 Fine-tuning/overfitting

Parameters: In order to force the fine-tuning to make changes to the filters of the network’s feature extractor rather than changing the last layers of the network, we freeze the convolutional layers related to up-sampling, more specifically, the filters coming after the pixel-shuffle layers. The images for fine-tuning are the random crops of 32×32 pixels from our constructed dataset. We choose a relatively low learning rate of $1e-4$ for gradual change.

K and M: We conduct sensitivity analysis to choose the best values for the number of images per filter K and the number of filter M used for our test image. We tune these parameters based on the perceptual quality of the generated images. The results of this work are produced by setting the values of K and M to two and five, respectively (10 images in total). a more detailed study can be found in the supplementary material.

Stoppage condition: The criteria to stop the fine-tuning was basically defined based on qualitative comparison of reconstructed images at different epochs where we could see at epoch 30, the vast majority of the images from our validation set were perceptually more convincing as compared to other epochs. However, considering Fig. 6, we can see this choice can also be justified as this epoch also coincides with the beginning of a significant drop in SSIM and PSNR values over all images on the test set.

5.1.3 Test-set

For our test-set, we randomly chose 100 images from the ImageNet dataset (non-overlapping between activation and training datasets), as both our baseline network and the Activation dataset are trained on/using a subset of 50,000 ImageNet images. As it is shown [12, 35] that SR network quality drops when doing cross-dataset tests, therefore, we focus on showing a proof of concept of improving a generic

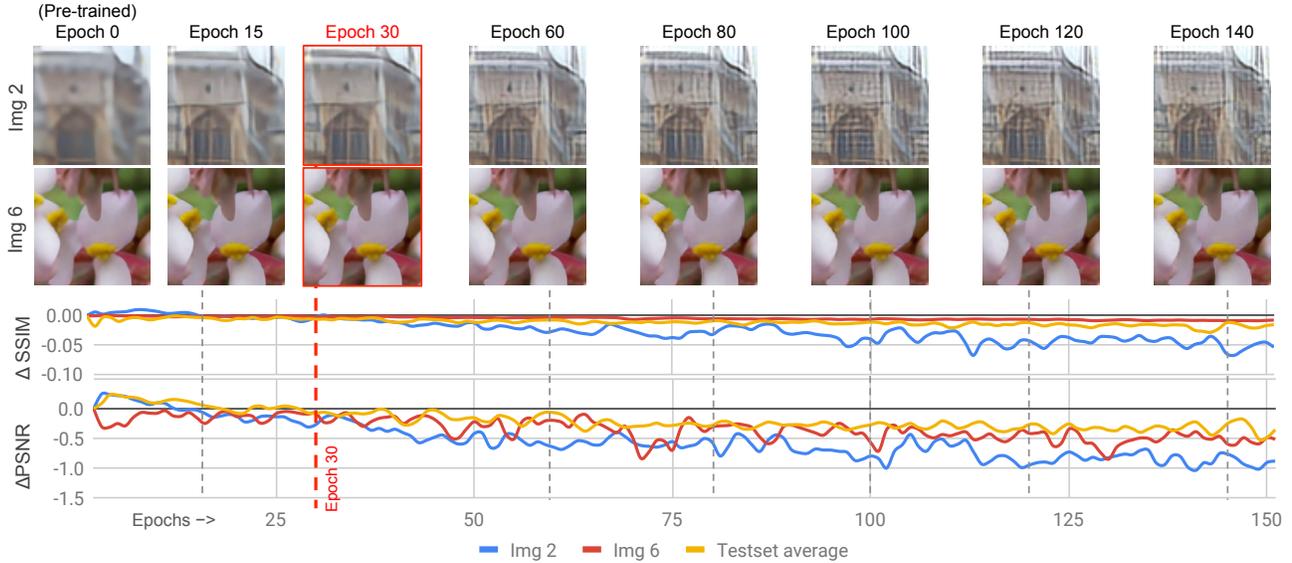


Figure 5. The effects of fine-tuning as a function of the number of epochs. We show the average change of PSNR and SSIM values over the test set, as well as explicit examples of visual, PSNR, and SSIM evolution on two images. We see in image 2 that perceptual quality can dramatically increase with fine-tuning, while image 6 is not affected significantly. **Please zoom in on the screen.**

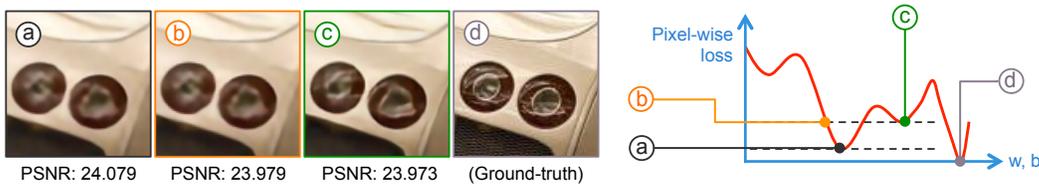


Figure 6. Image (a) is obtained by a pre-trained baseline with pixel-wise optimization on a large dataset. Images (b,c) are obtained during the fine-tuning by our proposed method, reaching almost the same PSNR. Image (d) is the ground truth. We see from comparing images (b) and (c) that our method is guiding the SR network to a different local minimum with a better perceptual quality, as the same loss is achieved but with dramatically different quality.

SR network on a generic dataset and do not add an additional variable of different datasets to the mix.

5.2. Filter selection analysis

In the following, we provide, to our knowledge, novel experiments and investigations into SR networks, where we examine, at the level of the network’s filters, how the SR network changes in response to our selective overfitting. For our experiments, we draw on [31], where authors found that two networks trained from scratch for the same task can have different filter orders and different filter patterns; however, fine-tuning a network to perform a different, but related task preserved the filter orders and patterns of the original network. They further show that the changes in filters by doing fine-tuning are gradual, by proposing to quantitatively assess the similarities between the filters of two different instances of the same network through correlation;

concretely, given filter F_i, F_j ,

$$\rho_{ij} = \frac{(F_i - \bar{F}_i)(F_j - \bar{F}_j)}{\sqrt{\|F_i - \bar{F}_i\|_2} \sqrt{\|F_j - \bar{F}_j\|_2}} \quad (1)$$

where ρ_{ij} is the correlation index. We use this correlation index to quantitatively study the changes in the filters of the SR network after fine-tuning. Given an LR image with HR ground truth, let G_{per} denote the EDSR baseline which is fine-tuned on solely this LR image to produce a perfect reconstruction. We can, in some sense, assume that G_{per} possesses the ideal or optimal set of filters for super-resolving this LR image, as we overfit it on this image; further, we verified that, consistent with [31], the overall structure/filter orders are preserved from the baseline network, indicating that G_{per} is not simply memorizing the image within its parameters.

Let G' denote the fine-tuned network produced from our



Figure 7. The average correlations over the test-set images of the filters of the final layer of feature extractor of G' and G_{rand} to the filters of G_{per} as a function of the number of epochs of fine-tuning in red/blue respectively, with the correlation of the baseline as a dotted black line. We see that the correlation of G' to G_{per} is higher than G_{rand} : This is consistent with our hypothesis that the proposed method of fine-tuning transforms the filters of the baseline to be closer to the “ideal” filters for a particular image.

method on this LR image. Let G_{rand} denote the EDSR network fine-tuned on a set of random images. In Fig 7, we show the average correlations of the filters of the final layer of G' and G_{rand} to the filters of G_{per} as a function of the number of epochs of fine-tuning. The average was computed by constructing G' , G_{rand} , G_{per} for each image in the test set, then taking the average correlation over the images. We also show the correlation of the filters of the baseline G with G_{per} . We see that the correlation of G' to G_{per} is generally higher than those of G_{rand} and G , including at 30 epochs, which is the number that we use for our method. This provides evidence that our method of fine-tuning in some sense brings the baseline closer to the “ideal” set of filters for a given LR image.

5.3. Comparison to PSNR-based approaches

From the qualitative results in Fig. 2, we can observe that when we fine-tune the pre-trained EDSR network using the images chosen through our method, namely activated-SR approach, **the perceptual quality increases with minimal impact on the PSNR/SSIM**. This minimal impact on the PSNR/SSIM has been also shown in Fig. 5, where we can see that over a test set of 100 images, the mean changes in PSNR/SSIM are minimal.

In Fig. 8, we additionally compare our method to LapSRN [16], RCAN [40] and EDSR [19] methods and by using test images from Set5 [3], Set14 [37] and BSD100 [20] standard datasets. **For a fair comparison, in this section, we only considered PSNR-based approaches as our methods still relies only on minimizing the pixel-wise distance of the SR and ground-truth images and does not benefit from any perceptual losses**. This figures shows that activated-SR images produced by our method have superior perceptual quality, while Table 1 confirms that this increase had a minimal impact on the PSNR/SSIM over the whole test set.

5.4. Comparison to perceptual-based approaches

Finally, in Fig 9, we provide a comparison between SR network trained using our proposed method and using perceptual losses (pixel-wise loss + vgg loss + adversarial loss,

Dataset	Metric	LapSRN	RCAN	EDSR	Ours
Set5	SSIM	0.887	0.918	0.893	0.891
	PSNR	31.56	32.61	32.41	32.40
Set14	SSIM	0.772	0.773	0.774	0.776
	PSNR	28.20	28.86	28.81	28.70
BSD100	SSIM	0.742	0.815	0.802	0.819
	PSNR	27.41	29.32	29.24	29.15

Table 1. Comparison LapSRN [16], RCAN [40], EDSR [19], and activated-SR (ours) on various test sets. **We emphasize that our method is EDSR using our test-time adaption method**. We show the results from other methods for comparison. Considering Fig. 8 the proposed method improves the perceptual quality of EDSR with minimal impact on the PSNR/SSIM.

with the same setting and discriminator as described in ESRGAN [33] work). We note that the perceptual loss adds more sharpness than that of our method, but can also provide highly distorted textures. In all cases, the images from our method are sharper/more detailed than those of the EDSR baseline, without distorting the texture. This can be explained by the fact that optimizing SR networks with only perceptual loss sometimes leads to the incitement of high frequency details in image e.g., sharp edges, entailing over-sharpened images. Therefore, they do not conform with the distortion based metrics.

On average, the decrease in PSNR and SSIM using perceptual loss is 628 and 355 percent larger, respectively, than the corresponding decreases using our method. Hence, our method provides images with much greater fidelity to the ground truth, while increasing the perceptual quality without distorted textures.

5.5. Inference time

We note that as our method fine-tunes the baseline network for every test image, this is computationally more expensive than simply using the baseline network. However, we note that relatively small patches of 32×32 pixels, and a small number of images (10 in our case) used for fine-tuning still keeps the computation time practical for single image SR tasks; the additional fine-tuning takes ~ 13 seconds by using a GeForce GTX 1080Ti GPU, which results in a total

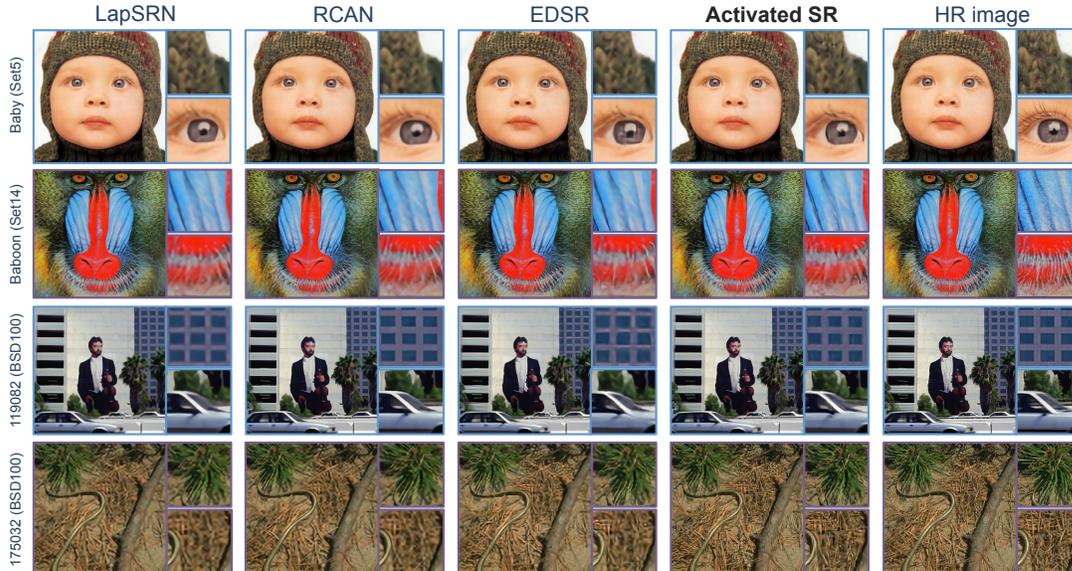


Figure 8. Qualitative comparison to **PSNR-based approaches**. From left to right: Bicubic, LapSRN [16], RCAN [40], EDSR [19], Activated-SR (ours), and HR image, tested on images from Set 5 [3], Set14 [37] and BSD100 [20] testsets. **We emphasize that our method is EDSR using our test-time adaption method.** We show results from other networks for comparison. Zoom in for the best view.

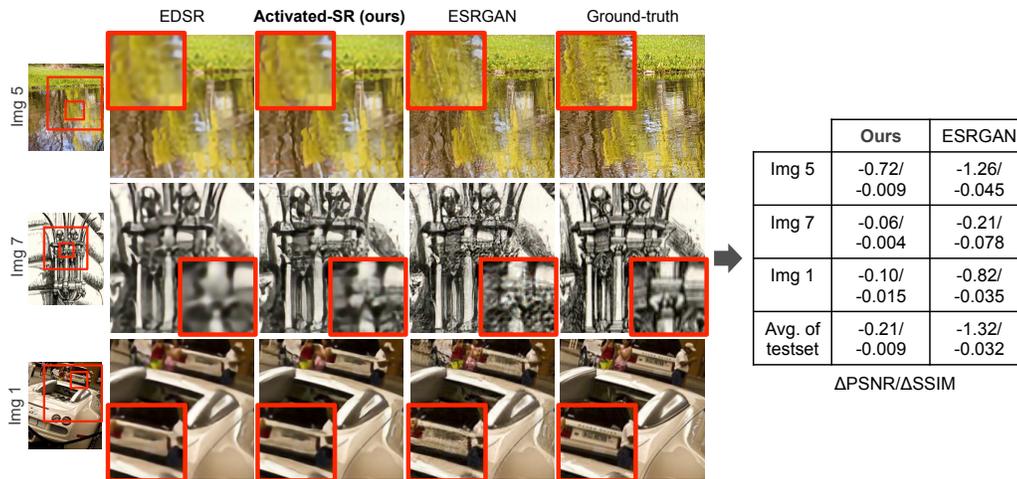


Figure 9. Comparing the proposed method and a perceptual-based approach [33]. In general, the perceptual loss provides sharper edges but also more distorted textures, whereas the proposed method provides images which are sharper and contain more details than the baseline without distortion. In the table, we show that this is reflected in the decrease in the PSNR/SSIM; **using perceptual loss decreases the PSNR/SSIM relative to the baseline far more than using our method.**

time of ~ 14 seconds for a 2560×1920 pixel output.

6. Conclusion

In this paper, we propose a novel approach to improve the perceptual quality of PSNR-based SR methods. In our approach, given a pre-trained SR network and LR input, we use test-time adaptation by fine-tuning the SR network on a subset of images from the training dataset with similar activation patterns as the initial HR prediction, with respect to the filters of a feature extractor. We show that the fine-tuned network produces an HR prediction with both greater

perceptual quality and minimal changes to the PSNR/SSIM, in contrast to perceptually driven approaches. Further, in contrast to reference-based SR, we use only images from our proposed activation dataset for fine-tuning, eliminating the issue with the availability of HR reference images close to the input image. Finally, through numerical experiments novel to the field of SR, we show that our fine-tuning can be interpreted as within the test-time adaptation paradigm, where we update the model parameters to be closer to the parameters of an "ideal" SR network, which is overfitted on the given LR input.

References

- [1] Saeed Anwar, Salman Khan, and Nick Barnes. A deep journey into super-resolution: A survey. *arXiv preprint arXiv:1904.07523*, 2019. 1
- [2] David Bau, Bolei Zhou, Aditya Khosla, Aude Oliva, and Antonio Torralba. Network dissection: Quantifying interpretability of deep visual representations. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6541–6549, 2017. 4
- [3] Marco Bevilacqua, Aline Roumy, Christine Guillemot, and Marie-Line Alberi Morel. Low-Complexity Single-Image Super-Resolution based on Nonnegative Neighbor Embedding. In *British Machine Vision Conference (BMVC)*, Guildford, Surrey, United Kingdom, Sept. 2012. 7, 8
- [4] Yochai Blau and Tomer Michaeli. The perception-distortion tradeoff. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6228–6237, 2018. 2, 3, 4
- [5] Behzad Bozorgtabar, Mohammad Saeed Rad, Hazim Kemal Ekenel, and Jean-Philippe Thiran. Learn to synthesize and synthesize to learn. *Computer Vision and Image Understanding*, 185:1–11, 2019. 3
- [6] Behzad Bozorgtabar, Mohammad Saeed Rad, Hazim Kemal Ekenel, and Jean-Philippe Thiran. Using photorealistic face synthesis and domain adaptation to improve facial expression analysis. In *2019 14th IEEE International Conference on Automatic Face & Gesture Recognition (FG 2019)*, pages 1–8. IEEE, 2019. 3
- [7] Mathilde Caron, Piotr Bojanowski, Armand Joulin, and Matthijs Douze. Deep clustering for unsupervised learning of visual features. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 132–149, 2018. 3, 4
- [8] Tao Dai, Jianrui Cai, Yongbing Zhang, Shu-Tao Xia, and Lei Zhang. Second-order attention network for single image super-resolution. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 11065–11074, 2019. 1
- [9] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009. 3, 5
- [10] Chao Dong, Chen Change Loy, Kaiming He, and Xiaoou Tang. Learning a deep convolutional network for image super-resolution. In *European conference on computer vision*, pages 184–199. Springer, 2014. 1
- [11] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680, 2014. 3
- [12] Zhen Han, Enyan Dai, Xu Jia, Shuaijun Chen, Chunjing Xu, Jianzhuang Liu, and Qi Tian. Unsupervised image super-resolution with an indirect supervised path. *arXiv preprint arXiv:1910.02593*, 2019. 5
- [13] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1125–1134, 2017. 3
- [14] Junjun Jiang, Yi Yu, Zheng Wang, Suhua Tang, Ruimin Hu, and Jiayi Ma. Ensemble super-resolution with a reference dataset. *IEEE transactions on cybernetics*, 2019. 1
- [15] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *European conference on computer vision*, pages 694–711. Springer, 2016. 2, 3
- [16] Wei-Sheng Lai, Jia-Bin Huang, Narendra Ahuja, and Ming-Hsuan Yang. Deep laplacian pyramid networks for fast and accurate super-resolution. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5835–5843, 2017. 7, 8
- [17] Christian Ledig, Lucas Theis, Ferenc Huszár, Jose Caballero, Andrew P. Aitken, Alykhan Tejani, Johannes Totz, Zehan Wang, and Wenzhe Shi. Photo-realistic single image super-resolution using a generative adversarial network. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 105–114, 2016. 3
- [18] Christian Ledig, Lucas Theis, Ferenc Huszár, Jose Caballero, Andrew Cunningham, Alejandro Acosta, Andrew Aitken, Alykhan Tejani, Johannes Totz, Zehan Wang, et al. Photo-realistic single image super-resolution using a generative adversarial network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4681–4690, 2017. 2, 3
- [19] Bee Lim, Sanghyun Son, Heewon Kim, Seungjun Nah, and Kyoung Mu Lee. Enhanced deep residual networks for single image super-resolution. *2017 IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pages 1132–1140, 2017. 5, 7, 8
- [20] D. Martin, C. Fowlkes, D. Tal, and J. Malik. A database of human segmented natural images and its application to evaluating segmentation algorithms and measuring ecological statistics. In *Proc. 8th Int’l Conf. Computer Vision*, volume 2, pages 416–423, July 2001. 7, 8
- [21] Mohammad Saeed Rad, Behzad Bozorgtabar, Urs-Viktor Marti, Max Basler, Hazim Kemal Ekenel, and Jean-Philippe Thiran. Srobb: Targeted perceptual loss for single image super-resolution. In *The IEEE International Conference on Computer Vision (ICCV)*, October 2019. 1
- [22] Mohammad Saeed Rad, Behzad Bozorgtabar, Claudiu Musat, Urs-Viktor Marti, Max Basler, Hazim Kemal Ekenel, and Jean-Philippe Thiran. Benefiting from multitask learning to improve single image super-resolution. *Neurocomputing*, 398:304–313, 2020. 1
- [23] Mohammad Saeed Rad, Thomas Yu, Claudiu Musat, Hazim Kemal Ekenel, Behzad Bozorgtabar, and Jean-Philippe Thiran. Benefiting from bicubically down-sampled images for learning real-world image super-resolution. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pages 1590–1599, January 2021. 1
- [24] Mehdi SM Sajjadi, Bernhard Scholkopf, and Michael Hirsch. Enhancenet: Single image super-resolution through

- automated texture synthesis. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 4491–4500, 2017. 3
- [25] Mehdi S. M. Sajjadi, Bernhard Schölkopf, and Michael Hirsch. Enhancenet: Single image super-resolution through automated texture synthesis. *2017 IEEE International Conference on Computer Vision (ICCV)*, pages 4501–4510, 2016. 3
- [26] Marcel Simon and Erik Rodner. Neural activation constellations: Unsupervised part model discovery with convolutional networks. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1143–1151, 2015. 4
- [27] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *CoRR*, abs/1409.1556, 2014. 3, 4, 5
- [28] Ying Tai, Jian Yang, Xiaoming Liu, and Chunyan Xu. Memnet: A persistent memory network for image restoration. In *Proceedings of the IEEE international conference on computer vision*, pages 4539–4547, 2017. 1
- [29] JD Van Ouwerkerk. Image super-resolution survey. *Image and vision Computing*, 24(10):1039–1052, 2006. 1, 2
- [30] Xintao Wang, Ke Yu, Chao Dong, and Chen Change Loy. Recovering realistic texture in image super-resolution by deep spatial feature transform. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 606–615, 2018. 3
- [31] Xintao Wang, Ke Yu, Chao Dong, Xiaoou Tang, and Chen Change Loy. Deep network interpolation for continuous imagery effect transition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1692–1701, 2019. 6
- [32] Xintao Wang, Ke Yu, Shixiang Wu, Jinjin Gu, Yihao Liu, Chao Dong, Yu Qiao, and Chen Change Loy. Esrgan: Enhanced super-resolution generative adversarial networks. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 0–0, 2018. 1, 3
- [33] Xintao Wang, Ke Yu, Shixiang Wu, Jinjin Gu, Yihao Liu, Chao Dong, Yu Qiao, and Chen Change Loy. Esrgan: Enhanced super-resolution generative adversarial networks. In *ECCV Workshops*, 2018. 7, 8
- [34] Zhihao Wang, Jian Chen, and Steven CH Hoi. Deep learning for image super-resolution: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020. 1
- [35] Yunxuan Wei, Shuhang Gu, Yawei Li, and Longcun Jin. Unsupervised real-world image super resolution via domain-distance aware training. *arXiv preprint arXiv:2004.01178*, 2020. 5
- [36] Wenhan Yang, Sifeng Xia, Jiaying Liu, and Zongming Guo. Reference-guided deep super-resolution via manifold localized external compensation. *IEEE Transactions on Circuits and Systems for Video Technology*, 29(5):1270–1283, 2018. 1
- [37] Roman Zeyde, Michael Elad, and Matan Protter. On single image scale-up using sparse-representations. In *Proceedings of the 7th International Conference on Curves and Surfaces*, pages 711–730, Berlin, Heidelberg, 2012. Springer-Verlag. 7, 8
- [38] Richard Zhang, Phillip Isola, and Alexei A Efros. Split-brain autoencoders: Unsupervised learning by cross-channel prediction. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1058–1067, 2017. 3
- [39] Yulun Zhang, Kunpeng Li, Kai Li, Lichen Wang, Bineng Zhong, and Yun Fu. Image super-resolution using very deep residual channel attention networks. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 286–301, 2018. 1
- [40] Yulun Zhang, Kunpeng Li, Kai Li, Lichen Wang, Bineng Zhong, and Yun Fu. Image super-resolution using very deep residual channel attention networks. In *ECCV*, 2018. 7, 8
- [41] Zhifei Zhang, Zhaowen Wang, Zhe Lin, and Hairong Qi. Image super-resolution by neural texture transfer. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7982–7991, 2019. 1
- [42] Haitian Zheng, Mengqi Ji, Haoqian Wang, Yebin Liu, and Lu Fang. Crossnet: An end-to-end reference-based super resolution network using cross-scale warping. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 88–104, 2018. 1