

BP-DIP: A Backprojection based Deep Image Prior

Jenny Zukerman

*Department of Biomedical Engineering
Tel Aviv University
Tel Aviv, Israel
jennyz@mail.tau.ac.il*

Tom Tirer

*School of Electrical Engineering
Tel Aviv University
Tel Aviv, Israel
tomtirer@mail.tau.ac.il*

Raja Giryes

*School of Electrical Engineering
Tel Aviv University
Tel Aviv, Israel
raja@tauex.tau.ac.il*

Abstract—Deep neural networks are a very powerful tool for many computer vision tasks, including image restoration, exhibiting state-of-the-art results. However, the performance of deep learning methods tends to drop once the observation model used in training mismatches the one in test time. In addition, most deep learning methods require vast amounts of training data, which are not accessible in many applications. To mitigate these disadvantages, we propose to combine two image restoration approaches: (i) Deep Image Prior (DIP), which trains a convolutional neural network (CNN) from scratch in test time using the given degraded image. It does not require any training data and builds on the implicit prior imposed by the CNN architecture; and (ii) a backprojection (BP) fidelity term, which is an alternative to the standard least squares loss that is usually used in previous DIP works. We demonstrate the performance of the proposed method, termed BP-DIP, on the deblurring task and show its advantages over the plain DIP, with both higher PSNR values and better inference run-time.

Index Terms—Deep learning, loss functions, image deblurring

I. INTRODUCTION

Image restoration refers to the recovery of an original unknown image from its degraded version, which suffers from defects, such as blur, noise and low resolution. In a linear image restoration problem, the goal is to recover the original image $x^* \in \mathbb{R}^n$ from the degraded measurements

$$y = Ax^* + e, \quad (1)$$

where $e \in \mathbb{R}^m$ is an additive noise and $A \in \mathbb{R}^{m \times n}$ is a degradation operator. For example, in image deblurring $m = n$ and A is a square ill-conditioned matrix which represents a blur operator that filters the image by a blur kernel.

Image restoration tasks often involve minimization of a cost function, composed of a fidelity term and a prior term

$$\min_x \ell(x, y) + \beta s(x), \quad (2)$$

where ℓ is the fidelity term, s is the prior term and β is a positive parameter that controls the level of regularization. The fidelity term forces the output image to comply with the observation model, while the prior poses an underlying assumption about the latent image, e.g. that natural images are sharp and free of noise and holes.

Since inverse problems represented by (1) are usually ill-posed, a vast amount of research has been focused on the prior

term $s(x)$. Various natural image priors have been researched. Some of them can be described by explicit and interpretable functions, e.g. total-variation (TV) [1], and others, such as BM3D [2] and (pre-trained) deep generative models [3], are more implicit (i.e. cannot be associated with explicit prior functions). Many image priors (of both kinds) also exploit the non-local similarity of natural images [4], [5].

The fidelity term, however, has been less researched and is often chosen as the typical least squares (LS) objective. Recently, the authors of [6] have presented the backprojection (BP) fidelity term, which shows an advantage over the standard LS term on different image restoration tasks, such as deblurring and super-resolution, using priors such as TV, BM3D, and convolutional neural network (CNN) denoisers [7], [8]. In [6], the BP fidelity term is also mathematically analyzed and compared to LS for the case where $s(x)$ is the Tikhonov regularization (i.e. the ℓ_2 prior).

Nowadays, CNNs are a very powerful tool for many computer vision tasks, including image restoration. For a given reconstruction task, CNNs can perform the inverse mapping from the observations to the signal domain and achieve state-of-the-art performances due to their ability to learn from large datasets. Researchers, motivated by deep learning great results, are applying deep neural networks to solve imaging inverse problems such as denoising [9]–[12], super-resolution [13]–[15] and deblurring [16]–[18].

However, the performance of deep learning methods tends to significantly drop once the observation model used in training mismatches the one in test time. This is the reason for the growing popularity of alternative methods, which are not biased to the observation model used in the offline training phase. One such example is the plug-and-play (P&P) denoisers approach [19] and its successors, which have been proposed in the last few years. In the original P&P paper [19], the authors minimize (2) using an optimization algorithm that decouples the fidelity term and the prior term. They propose to avoid explicit formulation of $s(x)$, and instead, handle the prior by an arbitrary denoising operation (e.g. BM3D or CNN denoisers). A related work is IDBP [7], which presents an alternative framework for solving inverse problems using off-the-shelf denoisers. This method requires less parameter tuning and implicitly uses the unique BP fidelity term [6]. Another recent variant of [19] modifies the way that the denoising engine is used to regularize the inverse problem [20], [21].

An additional disadvantage of deep learning originates from its requirement to often use a lot of data. Indeed, large datasets are used in the majority of works in this field. However, recently several deep learning methods, which are trained only using a single image, have been proposed and have demonstrated surprisingly good results. In Deep Image Prior (DIP) [22] the authors use the deep CNN itself as a prior for various restoration tasks, e.g. denoising, super-resolution and inpainting. As part of this approach, a CNN is trained from scratch during test time, using only the degraded image, with no requirement for a training dataset. Another work that trains a CNN super-resolver from scratch at test time is ZSSR [23], which presents a combination of deep architectures with internal-example learning. Besides the test image, no other images are used—all the pairs of training patches are extracted/synthesized from the test image. In SinGAN [24], a generative model is learned from a single image. The SinGAN contains a pyramid of fully convolutional generative adversarial networks (GANs) [25], where each of them learns the patch distribution at different scales of the image. Both ZSSR and SinGAN focus on a specific task (super-resolution and samples generation, respectively), in contrast to the DIP approach. There are also works that incorporate training on external data with image-adaptation (via fine-tuning using the test image), such as IDBP-CNN-IA [8] and IAGAN [26].

In this paper, we focus on the DIP strategy. Most of the papers that apply this approach use the typical LS loss function and achieve rather limited performance and/or require a very large number of backprop iterations at test time [22], [27], [28]. To mitigate these deficits, we propose to use a BP loss function instead of the LS loss. We demonstrate this approach for deblurring, using multiple kernels and noise levels, and present improved restoration results, which are reflected by higher PSNR and SSIM values and much better inference runtime.

II. BACKGROUND

Before we turn to describe our method, we first describe in more details the DIP and BP strategies.

A. Deep Image Prior (DIP)

In recent years, training deep convolution neural networks have become a common way to perform image restoration tasks. The popularity of this machine learning approach stems from the difficulty in accurately modelling natural images with explicit prior functions. Typically, the implicit prior that is associated with a CNN is achieved by supervised learning using a large dataset, where the degraded images serve as the networks inputs and the weights are optimized such that the networks outputs match the original images.

The Deep Image Prior (DIP) work [22] has demonstrated a remarkable phenomenon: CNNs can be used for solving image restoration problems without any offline training and external data. The DIP paper disputes the idea that supervised learning is mandatory for restoring images with CNNs. It shows that the network's architecture itself is sufficient for forcing a strong

prior that allows reconstructing an image from its degraded version. Therefore, there is no need in large datasets and offline learning, as the image restoration can be performed only by using the single degraded image.

In DIP, the estimated image $x \in \mathbb{R}^{3 \times H \times W}$ is parameterized by

$$x = f_\theta(z), \quad (3)$$

where $f_\theta(z)$ is (typically) a deep CNN with U-Net architecture [29], $z \in \mathbb{R}^{C' \times H' \times W'}$ is a fixed tensor filled with uniform noise, and θ are the network parameters. Currently, DIP works consider a loss function that is given by the LS objective

$$\min_{\theta} \|y - Af_\theta(z)\|_2^2, \quad (4)$$

and minimize it, with respect to θ , using first-order methods such as SGD and Adam [30]. Note that overfitting y may lead to $f_\theta(z)$ with artifacts. Therefore, in DIP the optimization process is terminated early.

It is interesting to note that the above strategy can be obtained from the general formulation in (2), for the LS fidelity term

$$\ell(x, y) = \frac{1}{2} \|y - Ax\|_2^2, \quad (5)$$

and the following indicator prior

$$s(x) = \begin{cases} 0, & x = f_\theta(z) \\ +\infty, & \text{otherwise} \end{cases}. \quad (6)$$

B. The backprojection (BP) fidelity term

Recently, the backprojection (BP) fidelity term has been proposed in [6] as an alternative to the widely used LS term (5). Under the practical assumptions that $m \leq n$ and $\text{rank}(A) = m$, the pseudoinverse of A is given by $A^\dagger = A^T(AA^T)^{-1}$, and the backprojection fidelity term is

$$\ell(x, y) = \frac{1}{2} \|A^\dagger(y - Ax)\|_2^2. \quad (7)$$

This fidelity term encourages agreement between the projection of the optimization variable onto the row space of the linear operator (i.e. $A^\dagger Ax$) and the pseudoinverse of the linear operator (back-projection) applied on the observations (i.e. $A^\dagger y$). Note that (7) can also be written as¹

$$\ell(x, y) = \frac{1}{2} \|(AA^T)^{-\frac{1}{2}}(y - Ax)\|_2^2. \quad (8)$$

It has been demonstrated in [6]–[8] that for different priors (e.g. TV, BM3D, and pre-trained CNNs) the BP fidelity term can yield better recoveries than LS for badly conditioned A and requires fewer iterations of optimization algorithms (the improved convergence rate is also analyzed in [31]). Yet, when the singular values of A are small, the performance advantage of BP is inversely proportional to the noise level.

¹The equivalence between (7) and (8) can be observed by expanding the two quadratic forms.

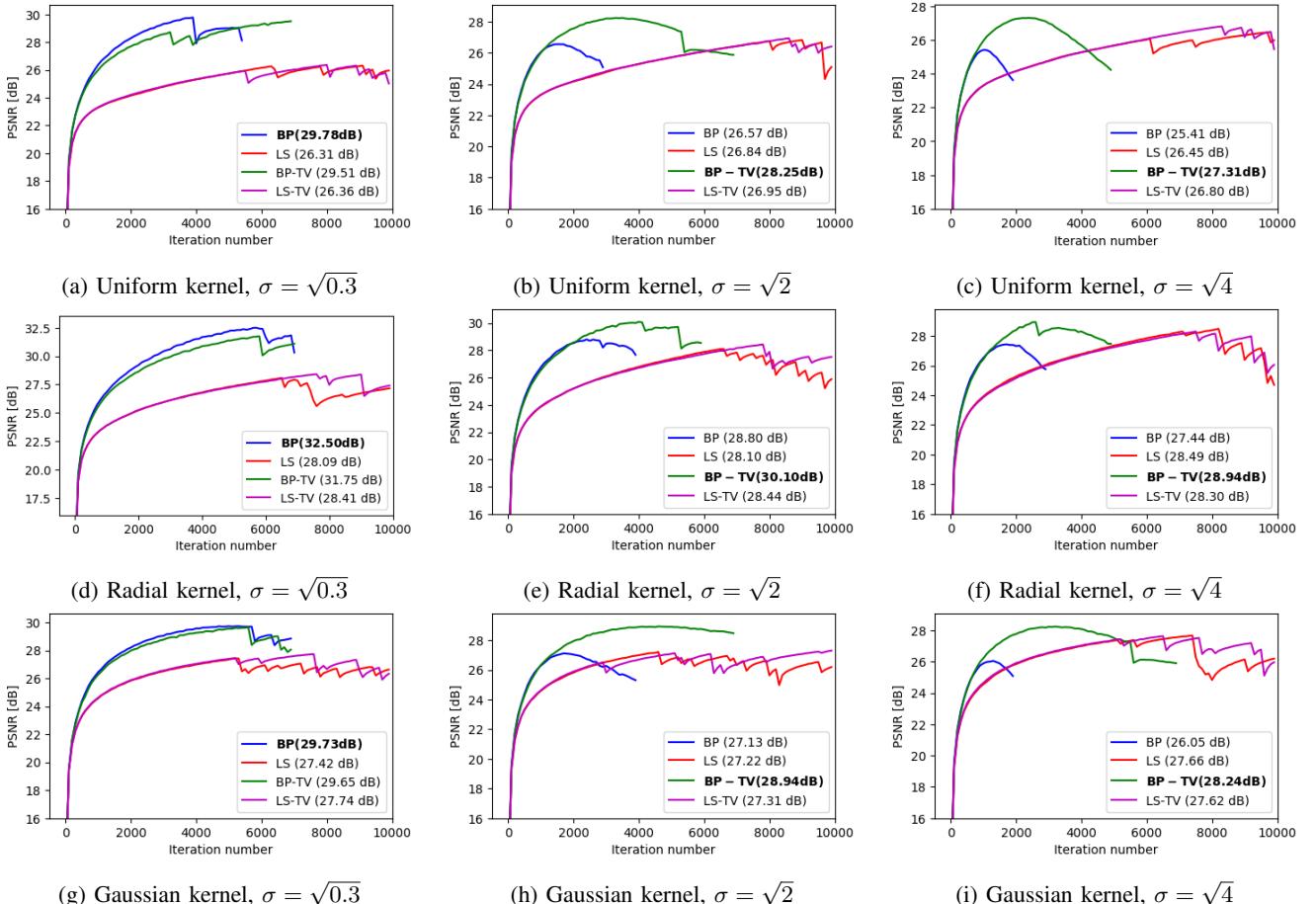


Fig. 1: Deblurring results (averaged over Set14) for: Uniform kernel in (a)-(c), Radial kernel in (d)-(f) and Gaussian kernel in (g)-(i), for noise levels $\sigma = \sqrt{0.3}, \sqrt{2}$ and $\sqrt{4}$. Note that BP achieves best PSNR in all kernels with $\sigma = \sqrt{0.3}$, and BP-TV achieves best PSNR in all kernels with higher noise levels ($\sigma = \sqrt{2}, \sqrt{4}$).

III. BACK-PROJECTION BASED DEEP IMAGE PRIOR

In this paper, we demonstrate that using the BP fidelity term improves the performance of standard DIP, which uses the LS fidelity term as the loss function. The use of BP fidelity term (8) in DIP leads to the following cost function

$$\min_{\theta} \left\| (AA^T)^{-\frac{1}{2}} (y - Af_{\theta}(z)) \right\|_2^2. \quad (9)$$

For the image deblurring task, A represents convolution with a blur kernel h . Therefore, in this case, AA^T represents convolution with the filter $h * \text{flip}(h)$. This operator, as well as its square root inverse $(AA^T)^{-\frac{1}{2}}$, has a very fast implementation using the Fast Fourier Transform (FFT) [32]. To conclude, the loss function (9) can be efficiently implemented by²

$$\min_{\theta} \left\| \mathcal{F}^* \left(\frac{1}{\sqrt{|\mathcal{F}(h)|^2 + \epsilon_1 \sigma^2 + \epsilon_2}} \mathcal{F}(y - h * f_{\theta}) \right) \right\|_2^2, \quad (10)$$

²Note that this formulation also allows to easily obtain the loss gradients using popular software packages, such as TensorFlow and PyTorch.

where h is a blur kernel, σ is the noise level in (1), and \mathcal{F} , \mathcal{F}^* stand for Fourier transform and inverse Fourier transform, respectively. ϵ_1 and ϵ_2 are regularization parameters, which are required since A is ill-conditioned in case of deblurring, and the scenarios include noise. Note that FFT implementation of the operator AA^T and its inverse can be done in more restoration tasks, such as in super-resolution (e.g. see [33]).

As explained in [6], in scenarios such as deblurring, where A has many singular values that are much smaller than 1, the BP term is more sensitive to noise than the LS term. Therefore, the inversion of AA^T has to be regularized. The higher the noise level in y is, the more sensitive the expression $\frac{1}{\sqrt{|\mathcal{F}(h)|^2 + \epsilon_1 \sigma^2 + \epsilon_2}} \mathcal{F}(y)$ in (10).

In order to mitigate the sensitivity to noise of BP for DIP (which seems somewhat increased compared to other priors [6]–[8]), we propose to add TV regularization [1] to the loss function (10). The TV term encourages piecewise smoothness in the image and has led to a more balanced restoration in our experiments. The (anisotropic) TV loss, used in this work, is



(a) Ground truth (b) Blurred image (c) LS (d) LS-TV (e) BP (f) BP-TV

Fig. 2: Deblurring using BP-DIP and LS-DIP (with and without TV). Uniform kernel, $\sigma = \sqrt{0.3}$



(a) Ground truth (b) Blurred image (c) LS (d) LS-TV (e) BP (f) BP-TV

Fig. 3: Deblurring using BP-DIP and LS-DIP (with and without TV). Radial kernel, $\sigma = \sqrt{2}$



(a) Ground truth (b) Blurred image (c) LS (d) LS-TV (e) BP (f) BP-TV

Fig. 4: Deblurring using BP-DIP and LS-DIP (with and without TV). Gaussian kernel, $\sigma = \sqrt{2}$



(a) Ground truth (b) Blurred image (c) LS (d) LS-TV (e) BP (f) BP-TV

Fig. 5: Deblurring using BP-DIP and LS-DIP (with and without TV). Gaussian kernel, $\sigma = \sqrt{4}$

given by

$$\sum_{i,j} |x_{i+1,j} - x_{i,j}| + |x_{i,j+1} - x_{i,j}| \quad (11)$$

for a two-dimensional signal x .

IV. EXPERIMENTS

We demonstrate our approach on the deblurring task, which aims at recovering the original, sharp image from a blurred image. The performance of the deblurring process is compared between two cases: DIP with backprojection loss and DIP with least squares loss, denoted by BP-DIP and LS-DIP,

respectively. We use 3 different blur kernels; radial, Gaussian and uniform, with 3 different noise levels: $\sqrt{0.3}$, $\sqrt{2}$ and $\sqrt{4}$. The radial kernel is of size 15 x 15 and can be written as $\frac{1}{1+x_1^2+x_2^2}, x_1, x_2 = -7, \dots, 7$, the Gaussian kernel is of size 15 x 15 with STD 1.6, and the uniform kernel is of size 9 x 9. All three kernels are normalized to the sum of 1.

The hyper-parameters and CNN are chosen to be the same as in the original DIP paper [22], such that all of the experiments are performed using a U-Net architecture [29] with skip-connections where the input and the output have the same spatial size, the optimiser is Adam [30] and the learning rate

TABLE I: Deblurring results (PSNR [dB] / SSIM averaged over Set14) of the different methods

	Uniform kernel			Radial kernel			Gaussian kernel		
	$\sigma = \sqrt{0.3}$	$\sigma = \sqrt{2}$	$\sigma = \sqrt{4}$	$\sigma = \sqrt{0.3}$	$\sigma = \sqrt{2}$	$\sigma = \sqrt{4}$	$\sigma = \sqrt{0.3}$	$\sigma = \sqrt{2}$	$\sigma = \sqrt{4}$
LS	26.31 / 0.83	26.84 / 0.84	26.45 / 0.83	28.09 / 0.88	28.10 / 0.88	28.49 / 0.87	27.42 / 0.87	27.22 / 0.86	27.66 / 0.86
LS-TV	26.36 / 0.83	26.95 / 0.84	26.80 / 0.84	28.41 / 0.88	28.44 / 0.88	28.30 / 0.88	27.74 / 0.86	27.31 / 0.85	27.62 / 0.86
BP	29.78 / 0.91	26.57 / 0.84	25.41 / 0.80	32.50 / 0.95	28.80 / 0.89	27.44 / 0.86	29.73 / 0.91	27.13 / 0.85	26.05 / 0.82
BP-TV	29.51 / 0.90	28.25 / 0.88	27.31 / 0.86	31.75 / 0.93	30.10 / 0.91	28.94 / 0.89	29.65 / 0.91	28.94 / 0.90	28.24 / 0.88

is 0.01. The deblurring performance is evaluated on Set14 dataset.

The weight of the TV regularizer was chosen as the best value for BP and LS separately: $1e - 3$ for BP (all kernels), $1e - 5$ and $1e - 6$ for LS with Gaussian/radial kernels and uniform kernel, respectively. The ϵ_1 and ϵ_2 values in equation (10) were chosen as 0.01 and $1e - 3$ respectively.

Curves of PNSR vs. iteration number for all 9 experiments are displayed in Figure 1, and the final PNSR/SSIM are also displayed in Table I. Several visual examples are presented in Figures 2-5. It is apparent that: (a) with BP loss, DIP yields higher PSNR than with LS loss; (b) BP-DIP reaches its peak PSNR faster than LS-DIP; (c) In most cases, adding the TV term improves the PSNR for both LS and BP, however, the bigger improvement is seen with BP-DIP; (d) As in the original DIP paper, early stopping is needed, otherwise the images are corrupted with noise.

Notice that when $\sigma = \sqrt{0.3}$, BP presents the best PSNR out of the 4 methods (BP, BP-TV, LS, LS-TV) and there is, in fact, no need in adding the TV term. However, when σ is higher, BP-TV presents the best PSNR out of the 4 methods. In general, when the noise level rises, the gap between BP and LS results reduces and adding TV to the loss term usually boosts the results. Also notice that the PSNR of LS-DIP starts descending at some iteration, similarly to the PSNR of BP-DIP (i.e. both of them require early stopping). However, for LS-DIP it happens at a much later iteration. This behavior shows that when the number of iterations is tuned for best performance, BP-DIP (even with TV) is faster than the LS-DIP.

V. CONCLUSION

In this work, we examined the influence of the backprojection (BP) fidelity term on Deep Image Prior (DIP). We conducted multiple deblurring experiments, using various blur kernels and noise levels and achieved significant improvement over the DIP work, both in PSNR and in the required number of optimization iterations (and thus in the inference runtime). Our approach presents another empirical evidence that untrained CNNs can reconstruct a clean and sharp image using only its degraded version. Yet, it demonstrates that very good results can be obtained even after a relatively small number of iterations. Future work includes examining the same concept with other restoration tasks, e.g. super-resolution.

REFERENCES

- [1] L. I. Rudin, S. Osher, and E. Fatemi, “Nonlinear total variation based noise removal algorithms,” 1992.
- [2] K. Dabov, A. Foi, V. Katkovnik, and K. O. Egiazarian, “Image denoising by sparse 3-d transform-domain collaborative filtering,” *IEEE Transactions on Image Processing*, vol. 16, pp. 2080–2095, 2007.
- [3] A. Bora, A. Jalal, E. Price, and A. G. Dimakis, “Compressed sensing using generative models,” in *ICML*, 2017.
- [4] A. Buades, B. Coll, and J.-M. Morel, “A non-local algorithm for image denoising,” *2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, vol. 2, pp. 60–65 vol. 2, 2005.
- [5] D. Glasner, S. Bagon, and M. Irani, “Super-resolution from a single image,” *2009 IEEE 12th International Conference on Computer Vision*, pp. 349–356, 2009.
- [6] T. TIREZ and R. Giry, “Back-projection based fidelity term for ill-posed linear inverse problems,” *IEEE Transactions on Image Processing*, vol. 29, no. 1, pp. 6164–6179, 2020.
- [7] ———, “Image restoration by iterative denoising and backward projections,” *IEEE Transactions on Image Processing*, vol. 28, pp. 1220–1234, 2017.
- [8] ———, “Super-resolution via image-adapted denoising CNNs: Incorporating external and internal learning,” *IEEE Signal Processing Letters*, vol. 26, pp. 1080–1084, 2019.
- [9] K. Zhang, W. Zuo, Y. Chen, D. Meng, and L. Zhang, “Beyond a gaussian denoiser: Residual learning of deep cnn for image denoising,” *IEEE Transactions on Image Processing*, vol. 26, pp. 3142–3155, 2017.
- [10] X.-J. Mao, C. Shen, and Y.-B. Yang, “Image restoration using very deep convolutional encoder-decoder networks with symmetric skip connections,” in *NIPS*, 2016.
- [11] V. Jain and H. S. Seung, “Natural image denoising with convolutional networks,” in *NIPS*, 2008.
- [12] J. Xie, L. Xu, and E. Chen, “Image denoising and inpainting with deep neural networks,” in *NIPS*, 2012.
- [13] B. Lim, S. Son, H. Kim, S. Nah, and K. M. Lee, “Enhanced deep residual networks for single image super-resolution,” *2017 IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pp. 1132–1140, 2017.
- [14] C. Dong, C. C. Loy, K. He, and X. Tang, “Image super-resolution using deep convolutional networks,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 38, pp. 295–307, 2014.
- [15] J. Kim, J. K. Lee, and K. M. Lee, “Accurate image super-resolution using very deep convolutional networks,” *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 1646–1654, 2015.
- [16] J. Sun, W. Cao, Z. Xu, and J. Ponce, “Learning a convolutional neural network for non-uniform motion blur removal,” *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 769–777, 2015.
- [17] S. Nah, T. H. Kim, and K. M. Lee, “Deep multi-scale convolutional neural network for dynamic scene deblurring,” *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 257–265, 2016.
- [18] A. Chakrabarti, “A neural approach to blind motion deblurring,” *ArXiv*, vol. abs/1603.04771, 2016.
- [19] S. V. Venkatakrishnan, C. A. Bouman, and B. Wohlberg, “Plug-and-play priors for model based reconstruction,” *2013 IEEE Global Conference on Signal and Information Processing*, pp. 945–948, 2013.
- [20] Y. Romano, M. Elad, and P. Milanfar, “The little engine that could: Regularization by denoising (red),” *ArXiv*, vol. abs/1611.02862, 2016.
- [21] S. A. Bigdeli, M. Zwicker, P. Favaro, and M. Jin, “Deep mean-shift priors for image restoration,” in *Advances in Neural Information Processing Systems*, 2017, pp. 763–772.
- [22] D. Ulyanov, A. Vedaldi, and V. S. Lempitsky, “Deep image prior,” *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9446–9454, 2018.
- [23] A. Shocher, N. Cohen, and M. Irani, “Zero-shot super-resolution using deep internal learning,” *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 3118–3126, 2017.

- [24] T. R. Shaham, T. Dekel, and T. Michaeli, “Singan: Learning a generative model from a single natural image,” *ArXiv*, vol. abs/1905.01164, 2019.
- [25] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, “Generative adversarial nets,” in *NIPS*, 2014, pp. 2672–2680.
- [26] S. Abu Hussein, T. Tirer, and R. Giryes, “Image-adaptive GAN based reconstruction,” *AAAI Conference on Artificial Intelligence*, 2020.
- [27] D. Van Veen, A. Jalal, M. Soltanolkotabi, E. Price, S. Vishwanath, and A. G. Dimakis, “Compressed sensing with deep image prior and learned regularization,” *arXiv preprint arXiv:1806.06438*, 2018.
- [28] G. Mataev, P. Milanfar, and M. Elad, “Deepred: Deep image prior powered by red,” in *Proceedings of the IEEE International Conference on Computer Vision Workshops*, 2019.
- [29] O. Ronneberger, P. Fischer, and T. Brox, “U-net: Convolutional networks for biomedical image segmentation,” *ArXiv*, vol. abs/1505.04597, 2015.
- [30] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” *CoRR*, vol. abs/1412.6980, 2014.
- [31] T. Tirer and R. Giryes, “On the convergence rate of projected gradient descent for a back-projection based objective,” *arXiv preprint arXiv:2005.00959*, 2020.
- [32] J. W. Cooley and J. W. Tukey, “An algorithm for the machine calculation of complex Fourier series,” *Mathematics of computation*, vol. 19, no. 90, pp. 297–301, 1965.
- [33] S. Abu Hussein, T. Tirer, and R. Giryes, “Correction filter for single image super-resolution: Robustifying off-the-shelf deep super-resolvers,” in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.