
Blind Image Restoration with Flow Based Priors

Leonhard Helminger^{1,*} Michael Bernasconi^{1,*} Abdelaziz Djelouah²

Markus Gross¹ Christopher Schroers²

¹Department of Computer Science
ETH Zurich, Switzerland

²DisneyResearch|Studios
Zurich, Switzerland

Abstract

Image restoration has seen great progress in the last years thanks to the advances in deep neural networks. Most of these existing techniques are trained using full supervision with suitable image pairs to tackle a specific degradation. However, in a blind setting with unknown degradations this is not possible and a good prior remains crucial. Recently, neural network based approaches have been proposed to model such priors by leveraging either denoising autoencoders or the implicit regularization captured by the neural network structure itself. In contrast to this, we propose using normalizing flows to model the distribution of the target content and to use this as a prior in a maximum a posteriori (MAP) formulation. By expressing the MAP optimization process in the latent space through the learned bijective mapping, we are able to obtain solutions through gradient descent. To the best of our knowledge, this is the first work that explores normalizing flows as prior in image enhancement problems. Furthermore, we present experimental results for a number of different degradations on data sets varying in complexity and show competitive results when comparing with the deep image prior approach.

1 Introduction

In today's digitized world, there is an increased demand to process existing older content. Examples are the archival of photo prints (Liu et al.) for more reliable long-term data storage, preparing heritage footage (Ame) for more engaging documentaries, and making classic films and existing catalog contents available to large new audiences through streaming services. This old content is however often in low quality and may be deteriorated in complex ways, which creates a need for *blind* image restoration methods that are generic and able to address a wide range of possibly combined degradations. Blind image restoration can be formulated as solving the following energy minimization problem:

$$\mathbf{x}^* = \arg \min_{\mathbf{x}} [\mathcal{L}_{\text{data}}(\hat{\mathbf{x}}, \mathbf{x}) + \mathcal{L}_{\text{reg}}(\mathbf{x})] , \quad (1)$$

where $\hat{\mathbf{x}}$ is the observed image and \mathbf{x}^* the restored image to be estimated. The first term, $\mathcal{L}_{\text{data}}$, is a data fidelity term which can be problem dependent and ensures that the solution agrees with the observation; the second term, $\mathcal{L}_{\text{reg}}(\mathbf{x})$, is a regularizer that typically encodes certain smoothness assumptions on the expected solution and thus pushes it to lie within a given space. From a Bayesian viewpoint, the posterior distribution of the restored image is $p(\mathbf{x}|\hat{\mathbf{x}}) \propto p(\hat{\mathbf{x}}|\mathbf{x})p(\mathbf{x})$. This allows rewriting the above restoration problem into the following equivalent maximum a posteriori (MAP) estimate:

$$\mathbf{x}^* = \arg \max_{\mathbf{x}} \log(p(\mathbf{x}|\hat{\mathbf{x}})) = \arg \max_{\mathbf{x}} \underbrace{\log(p(\hat{\mathbf{x}}|\mathbf{x}))}_{\text{data}} + \underbrace{\log p(\mathbf{x})}_{\text{reg}} , \quad (2)$$

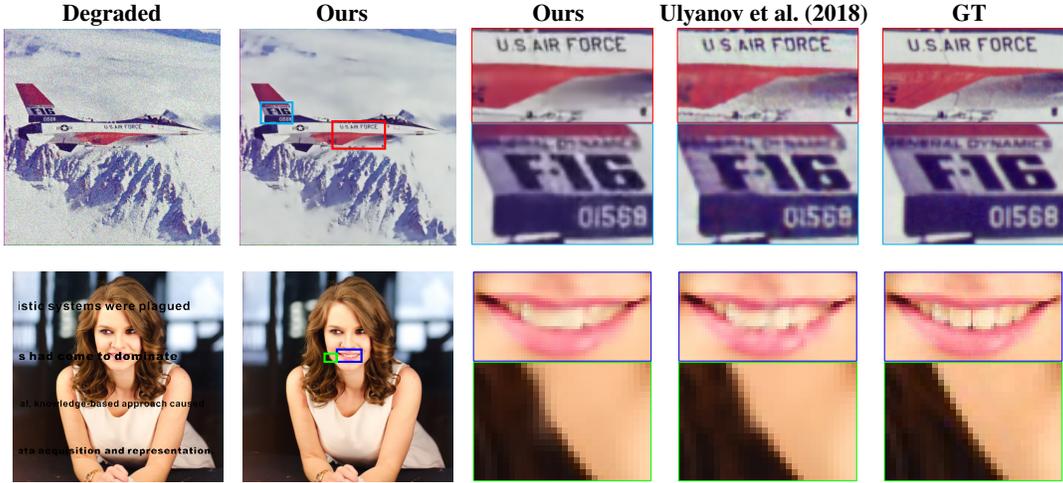


Figure 1: Comparative results with Deep Image Prior (Ulyanov et al., 2018) on different image restoration tasks. The first example corresponds to denoising whereas the second is image inpainting. Our approach is able to remove the degradation and produces visually more pleasing results in some regions like the text and the teeth.

which makes it more explicit that the regularizer should model prior knowledge about the unknown solution. Many handcrafted priors have been proposed reflecting desired properties based on total variation (Rudin et al., 1992), gradient sparsity (Fergus et al., 2006) or the dark pixel prior (He et al., 2010). More recently, learning based priors have been explored, in particular the usage of denoising autoencoders (DAEs) as regularizers for inverse imaging problems (Meinhardt et al., 2017). Building on DAEs, Bigdeli et al. (2017) propose to use a Gaussian smoothed natural image distribution as prior. In a different direction, Ulyanov et al. (2018) showed that an important part of the image statistics is captured by the structure of a convolutional image generator even independent of any learning.

All existing methods proposed alternatives and approximations to the true image prior $p(\mathbf{x})$ in Equation 2. However, with deep normalizing flows, we have an approach for a tractable *and* exact log-likelihood computation (Dinh et al., 2017). Therefore, we propose to use normalizing flows for capturing the distribution of target high quality content to serve as a prior in the MAP formulation. In addition to this, the inference of the latent value that corresponds to a data point can be done exactly without any approximation since our generative model is invertible. We use this learned bijective mapping to express the MAP optimization process in the latent space and are able to obtain solutions through gradient descent. In a number of experiments, we explore our approach for different degradations on data sets of varying complexity and we show that we can achieve competitive results as illustrated in Figure 1.

The contribution of this paper is three fold: 1) to the best of our knowledge, our work is the first using normalizing flows to learn a prior for blind image restoration; 2) we take advantage of the bijective mapping learned by our model to express the MAP problem of image reconstruction in latent space, where gradient descent can be used to estimate the solution; 3) we propose using new loss terms during model training for regularizing the latent space which yields a better behavior during the MAP inference.

Our paper is organized as follows. In Section 2, we recap important background regarding normalizing flow before describing our method in Section 3. Section 4 covers important related work and Section 5 discusses our experimental results. We give our conclusions in Section 6.

2 Normalizing Flow

Borrowing the notation from Papamakarios et al. (2019), let's consider two random variables X and U that are related through the reversible transformation $T : \mathbb{R}^d \rightarrow \mathbb{R}^d$, $\mathbf{x} = T(\mathbf{u})$. In this case, the distribution of the two variables are related as follows:

$$p_X(\mathbf{x}) = p_U(\mathbf{u}) |\det J_T(\mathbf{u})|^{-1}, \quad (3)$$

where $\mathbf{u} = T^{-1}(\mathbf{x})$ and $J_T(\mathbf{u})$ is the Jacobian of T . Here, the determinant preserves total probability and can be understood as the *amount* of squeezing and stretching of the space induced by the transformer T . The objective of normalizing flows (Rezende & Mohamed, 2015) is to map a base distribution to an arbitrary distribution through a change of variable. In practice, a series T_1, \dots, T_K of such mappings are applied to transform the base distribution into a more complex multi-modal one

$$\mathbf{x} \xleftrightarrow{T_K^{-1}} \mathbf{h}_{K-1} \xleftrightarrow{T_{K-1}^{-1}} \mathbf{h}_{K-2} \cdots \mathbf{h}_1 \xleftrightarrow{T_1^{-1}} \mathbf{u}, \quad (4)$$

$$p_X(\mathbf{x}) = p_U(T^{-1}(\mathbf{x})) \prod_{k=1}^K \left| \det \frac{d\mathbf{h}_{k-1}}{d\mathbf{h}_k} \right|, \quad (5)$$

where we define $\mathbf{h}_K \triangleq \mathbf{x}$ and $\mathbf{h}_0 \triangleq \mathbf{u}$. It is clear that computing the determinant of these Jacobian matrices, as well as the function inverses, must remain easy to allow their integration as part of a neural network. This is not the case for arbitrary Jacobians and recent successes in normalizing flow are due to the proposition of invertible transformations with easy to compute determinants.

Normalizing flows as generative model. Recent works (Kingma & Dhariwal, 2018; Dinh et al., 2017) have shown the great potential of using normalizing flow as generative model where an image observation \mathbf{x} is generated from a latent representation \mathbf{u}

$$\mathbf{x} = T_\theta(\mathbf{u}) \quad \text{with} \quad \mathbf{u} \sim p(\mathbf{u}). \quad (6)$$

Here $\mathbf{x} \in \mathcal{X}$ is a high-dimensional vector, T_θ denotes a composition of invertible transformations, and $p(\mathbf{u})$ is the base distribution e.g. a normal distribution. Considering a discrete set \mathcal{X} of N natural images, the flow based model is trained by minimizing the following log-likelihood objective:

$$\mathcal{L} = \frac{1}{N} \sum_{i=1}^N -\log p_\theta(\mathbf{x}^{(i)}). \quad (7)$$

In the next section, we will describe our approach for leveraging flow based models for various image restoration applications.

3 Blind Image Restoration with Flow Based Priors

By training a generative flow model as described in the previous section, we learn a mapping T_θ from a latent space \mathcal{U} , with a known base distribution $p(\mathbf{u})$, to the complex image space \mathcal{X} . In this work, we propose to use the capacity of normalizing flows to compute the exact likelihood of images $p_\theta(\mathbf{x})$, as prior in the image restoration problem

$$\mathbf{x}^* = \arg \min_{\mathbf{x}} -\log p(\hat{\mathbf{x}}|\mathbf{x}) - \log p_\theta(\mathbf{x}). \quad (8)$$

In addition to the prior, we also take advantage of the bijective mapping in normalizing flows to rewrite the optimization with respect to the latent \mathbf{u}

$$\mathbf{u}^* = \arg \min_{\mathbf{u}} [-\log p(\hat{\mathbf{x}}|T_\theta(\mathbf{u})) - \log p_\theta(T_\theta(\mathbf{u}))]. \quad (9)$$

With this new formulation, we are leveraging the learned mapping between the complex input space (the image space) and the base space (the latent space) that follows a simpler distribution. This new space is more adapted for such an optimization problem. In this work, we solve it through an iterative gradient descent, where each step is applied on the latents according to

$$\mathbf{u}^{t+1} = \mathbf{u}^t - \eta \nabla_{\mathbf{u}} L(\theta, \mathbf{u}, \hat{\mathbf{x}}). \quad (10)$$

Here $L(\theta, \mathbf{u}, \hat{\mathbf{x}})$ abbreviates the objective defined in equation 9 and η is the weighting applied to the gradient. We used the Adam optimizer (Kingma & Ba, 2015) to compute the gradient steps. The model is generic and once trained on *target quality* images, different applications can be considered by adapting the data loss term. In this work we use a generic data fidelity term between the input image $\hat{\mathbf{x}}$ and the restored result $\mathbf{x} = T_\theta(\mathbf{u})$:

$$\mathcal{L}_{\text{data}}(\hat{\mathbf{x}}, T_\theta(\mathbf{u})) = -\log p(\hat{\mathbf{x}}|T_\theta(\mathbf{u})) = \mathbf{m} \odot \lambda \|\hat{\mathbf{x}} - T_\theta(\mathbf{u})\|_2, \quad (11)$$

where \odot is the Hadamard product. The mask \mathbf{m} is a binary mask that indicates pixel locations with valid color values and allows to handle the inpainting scenario. The parameter λ controls the deviation tolerance from the original degraded input $\hat{\mathbf{x}}$. Next we provide details on the normalizing flow architecture used, the training losses, and our coarse to fine optimization procedure.

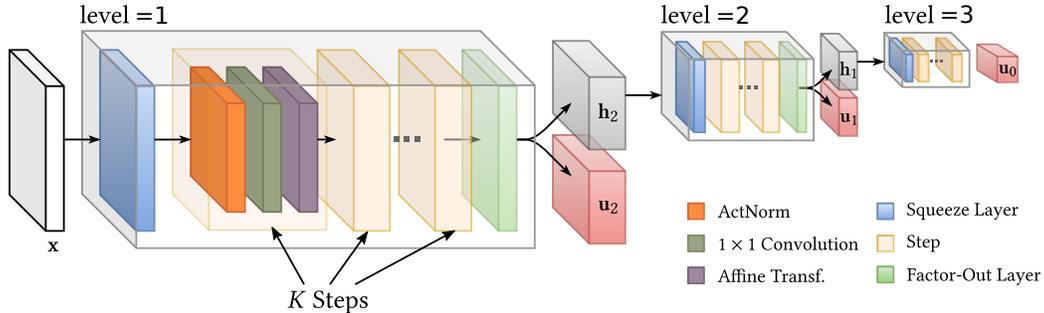


Figure 2: Overview of the normalizing flow architecture. The input image \mathbf{x} is processed by an $L = 3$ level network, where each level consists of a squeeze operation followed by a series of K steps. Each step is a succession of *ActNorm*, 1×1 convolution and an *affine layer*. The image latent representation is $(\mathbf{u}_0, \mathbf{u}_1, \mathbf{u}_2)$. The number of levels and steps can be adapted to the complexity of the data.

3.1 Generative Flow Architecture

The proposed generative model is based on the architecture described by Kingma & Dhariwal (2018). We first present the individual building layers

- **Activation normalization.** Proposed by Kingma & Dhariwal (2018), this is an alternative to batch normalization. It performs an affine transformation on the activations using a learned scale and bias parameter per channel.
- **Invertible 1×1 convolution.** Kingma & Dhariwal (2018) also proposed to replace the random permutation of channels, in coupling layers between the transformations, with a learned invertible 1×1 convolution.
- **Affine transformation.** This layer is a coupling introduced by Dinh et al. (2015). The input is split into two partitions, where one is the input for the conditioner, a neural network to modify the channels of the second partition. Here, the transformation is affine.
- **Factor-out layers.** The objective of factoring-out parts of the base distribution (Dinh et al., 2017) is to allow a coarse to fine modeling by introducing conditional distributions and dependencies on deeper levels.

Using these layers, we propose the model illustrated in Figure 2. It consists of L levels, each one is a succession of K steps, where a step defined as the composition of the layers: *ActNorm*, 1×1 convolution and *Affine*. At the end of each intermediate level, the transformed values (*latents*) are split in two parts \mathbf{h}_i and \mathbf{u}_i , with the factor-out layer. The parameters (μ_i, σ_i) of the conditional distribution $p(\mathbf{u}_i | \mathbf{h}_i)$ are predicted by a neural network. In our case, this is a zero initialized 2D convolution as proposed in (Kingma & Dhariwal, 2018). In the experimental part and in supplementary material, we provide more details about the architecture used for each dataset.

3.2 Training and Latent Space Regularization

When using normalizing flows to learn a continuous distribution, the input images have to be *dequantized*. Following common practices in generative flows, we redefine the negative log-likelihood objective (*nll*) of equation 7

$$\mathcal{L}_{nll} = \frac{1}{N} \sum_{i=1}^N -\log p_{\theta}(\mathbf{x}^{(i)} + \epsilon) . \quad (12)$$

Here ϵ is uniformly sampled from $[0, 1]$. This model is sufficient for simple datasets as we show in the experimental section with the MNIST examples (see Figure 3). However for more complex data, a regularization of the learned latent space is needed. The main objective is to structure this space in a beneficial way for the optimization.

Latent-Noise loss. In order to enforce some regularization of the latent space, we add uniform noise to the latents $\mathbf{u}_\xi = \mathbf{u} + \xi$ where $\xi \sim \mathcal{U}(-0.5, 0.5)$. The proposed loss term

$$\mathcal{L}_{ln} = \|T_\theta(\mathbf{u}_\xi) - \mathbf{x}\|_2 \quad (13)$$

penalizes parameters θ that would map back \mathbf{u}_ξ far from the initial input image \mathbf{x} . It is interesting to note that this loss does not make any assumption regarding the degraded images, but it still results in a latent space better suited for our optimization problem.

Auto-Encoder loss. If we consider the model illustrated in Figure 2, the image \mathbf{x} is mapped to its representation $(\mathbf{u}_0, \mathbf{u}_1, \mathbf{u}_2)$. From only the latent value \mathbf{u}_0 , we compute $\tilde{\mathbf{x}}$ by sampling the most likely intermediate values $\tilde{\mathbf{u}}_l \sim p(\mathbf{u}_l | \mathbf{h}_l)$. Since we use a Gaussian distribution, this corresponds to the mean value of the predicted distribution. The proposed loss

$$\mathcal{L}_{ae} = \|\tilde{\mathbf{x}} - \mathbf{x}\|_2 \quad (14)$$

forces the model to store sufficient information in the deepest level to reconstruct the image. This allows a more robust coarse-to-fine strategy during the optimization.

The final training loss for the normalizing flows is

$$\mathcal{L} = \mathcal{L}_{nll} + \beta_{ln}\mathcal{L}_{ln} + \beta_{ae}\mathcal{L}_{ae}, \quad (15)$$

where β_{ln} and β_{ae} are the weightings for each loss term. We used $\beta_{ln} = 100$ and $\beta_{ae} = 1$. The ablation study in the experimental section shows the necessity of training the generative flow model with all these loss terms.

3.3 Coarse-To-Fine Optimization

The optimization procedure described in Equation 10 is iterative and we need to set its initial value \mathbf{u}^0 . In order to choose a good starting point, we leverage the introduced multi-scale architecture. Our starting point is

$$\mathbf{u}^0 = (\hat{\mathbf{u}}_0, \tilde{\mathbf{u}}_1, \tilde{\mathbf{u}}_2) \quad \text{with } \hat{\mathbf{u}}_0 \text{ defined by } T_\theta^{-1}(\hat{\mathbf{x}}). \quad (16)$$

The values of the other components, $\tilde{\mathbf{u}}_1$ and $\tilde{\mathbf{u}}_2$, are sampled as the mean values of the respective predicted distributions. $p(\mathbf{u}_1 | \mathbf{h}_1)$ and $p(\mathbf{u}_2 | \mathbf{h}_2)$. As our auto-encoder loss enforces the possibility to reconstruct the image from $\hat{\mathbf{u}}_0$ only, this lowest level contains coarse image information while details are stored in the upper levels. This is advantageous for image restoration tasks where the degradation often affects the *detail* of an image.

Given this starting point, the optimization is done in a coarse-to-fine fashion. First, only the lowest level variables are optimized while the upper levels are respectively sampled from the predicted means. These are then progressively included in the optimization

$$\mathbf{u}_0^{t+1} = \mathbf{u}_0^t - \eta \nabla_{\mathbf{u}_0} L(\theta, \mathbf{u}, \hat{\mathbf{x}}), \quad (17)$$

$$(\mathbf{u}_0, \mathbf{u}_1)^{t+1} = (\mathbf{u}_0, \mathbf{u}_1)^t - \eta \nabla_{(\mathbf{u}_0, \mathbf{u}_1)} L(\theta, \mathbf{u}, \hat{\mathbf{x}}), \quad (18)$$

$$(\mathbf{u}_0, \mathbf{u}_1, \mathbf{u}_2)^{t+1} = (\mathbf{u}_0, \mathbf{u}_1, \mathbf{u}_2)^t - \eta \nabla_{(\mathbf{u}_0, \mathbf{u}_1, \mathbf{u}_2)} L(\theta, \mathbf{u}, \hat{\mathbf{x}}). \quad (19)$$

With this coarse-to-fine scheme, we are able to incrementally refine the reconstructed images by making sure that the lower level information is correct first.

4 Related Work and Discussion

Despite the success of supervised deep learning approaches for dedicated image restoration problems such as super-resolution (Wang et al., 2018; Zhang et al., 2018), denoising (Zhang et al., 2017a), inpainting (Pathak et al., 2016) or a combination of them (Park & Mu Lee, 2017), one important drawback is the need for retraining whenever the specific degradation or its parameters change. Some recent works (Cornillère et al., 2019; Bell-Kligler et al., 2019) have investigated the blind setting for super-resolution. However that concerns the parameters of the degradation only and such solutions are not applicable to an unknown degradation.

When addressing the blind restoration problem, the common approach is to consider the Bayesian perspective where recovering the original image is expressed as solving a maximum a posteriori (MAP)

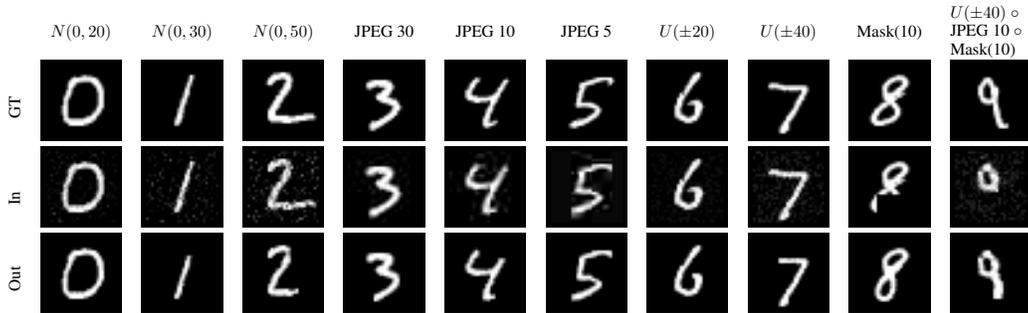


Figure 3: Results produced by a single-level normalizing flow trained on the MNIST dataset. Each column corresponds to a different type of degradation. From top to bottom the ground truth, the degraded image and the reconstructed image are shown.

problem. The objective function consists of a fidelity term and a regularization term. The fidelity term can be problem specific and easier to express than the prior that is supposed to reflect desired properties of the reconstructed image. Existing handcrafted priors are based on total variation (Rudin et al., 1992), gradient sparsity (Fergus et al., 2006) or the dark pixel prior (He et al., 2010). Recently, several works have investigated the usage of CNNs as priors. For example, (Rick Chang et al., 2017; Zhang et al., 2017b) show how a deep CNN trained for image denoising can effectively be used as prior in various image restoration tasks. Additionally, Meinhardt et al. (2017) provide new insights on how the denoising strength of the neural network relates to the weight on the data fidelity term. Bigdeli et al. (2017) define a utility function that includes the smoothed natural image distribution and relate this to denoising autoencoders. In a different direction, Ulyanov et al. (2018) showed that an important part of the image statistics is already captured by the structure of a convolutional image generator itself, independent of any learning. This work was further analyzed from a Bayesian perspective (Cheng et al., 2019) and combined with a denoising autoencoder prior (Mataev et al., 2019).

The idea presented in our work stems from recent developments in normalizing flows (Dinh et al., 2015, 2017; Kingma & Dhariwal, 2018) and their promising capacity of learning a bijective mapping from a space with a prescribed distribution to the complex space of images, additionally providing exact log-likelihood tractability. Using a learned prior that only depends on properties of high quality images is an exciting direction, as this removes the need to rely on other assumptions that are either explicit, in the case of handcrafted solutions, or implicit in the case of denoising autoencoders. This work is a first step demonstrating the potential of normalizing flows in image restoration tasks. We believe this is an exciting new direction that furthermore is expected to benefit from improvements and research that generally explores normalizing flow as a generative model.

5 Experiments

In this section we explore the usage of our proposed solution for different blind image restoration tasks. We show results on two synthetic datasets, the MNIST and the self generated Sprites, and on real images. We also include comparisons with the Deep Image Prior (DIP) (Ulyanov et al., 2018).

Since we do not focus on a specific degradation during training, our proposed approach can be applied on various types of restoration problems. In this work we present results on 3 different types of image degradation: noise (uniform and normal), JPEG compression artifacts, and missing regions. The noisy images are generated by adding i.i.d. samples of noise to the pixel values, with noise distributed according to $\mathcal{U}(\min, \max)$ or $\mathcal{N}(0, \sigma)$. The varying degrees of JPEG artifacts are generated by using different levels (10 to 70) for the JPEG compression. For the inpainting task, we masked multiple regions of size 10×10 pixels. An overview of the used degradations is visualized in Figure 3.

MNIST results. As a first step we tested our flow based image prior on the well studied MNIST dataset (LeCun et al., 1998). Given the simplicity of this dataset, the model used for this experiment consists of a single-level $L = 1$ with $K = 16$ steps. We choose the base distribution $p(\mathbf{u})$ to be a

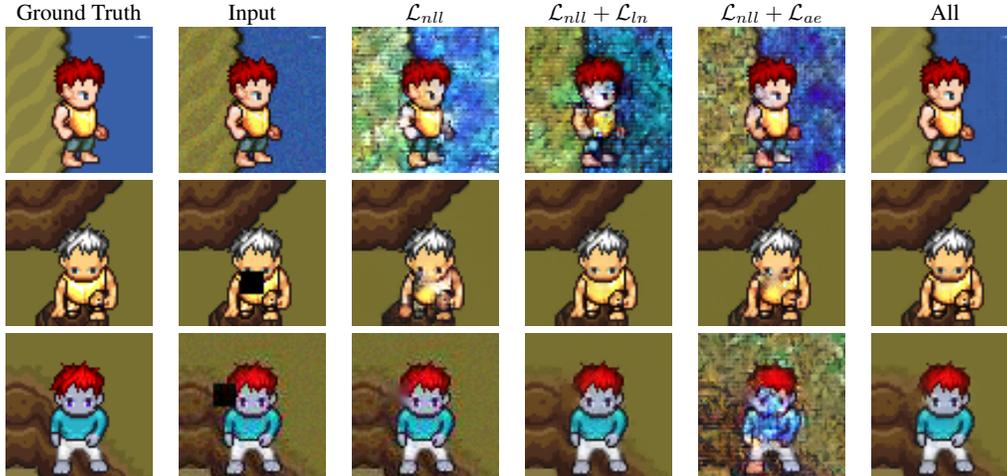


Figure 4: Restoration of degraded Sprites: first row corresponds to a Gaussian noise with $\sigma = 5$, second row is inpainting and last row combines denoising, inpainting and JPEG artifact removal. Columns correspond to different normalizing flow models, each one trained with the indicated loss term. Results show the importance of using all the proposed loss terms (see text for details).

Gaussian with unit variance and a trainable mean. Further, a ResNet (He et al., 2016) with 2 blocks and $C = 128$ intermediate channels, was used to learn the parameters for the affine transformations.

Given a degraded image \hat{x} the goal is to find the most likely image x^* by solving the optimization problem of Equation 9. Given the simplicity of the data set, we use the mean of the base distribution $p(\mathbf{u})$ as starting point \mathbf{u}^0 . It can be seen in Figure 3 that this is sufficient to enhance the binary digits for any degradation. A related experiment was conducted by Dinh et al. (2015), where the degraded digits were enhanced by maximizing the probability of the image trough back propagation to the pixel values. This is equivalent to only considering the prior term in Equation 9.

Sprites results. To handle this larger and more complex data set, we increased the capacity of our flow based prior. We use $L = 3$ levels, with $K = 8$ steps each. In the optimization, the learning rate η and the data weighting term λ are set to 1 and 99, respectively. The gradient descent is done in a coarse to fine way (see section 3.3), each time with 50 update steps per level before including the next one. When all latent levels are included, an additional 150 optimization steps are performed.

Figure 4 shows image restoration results on this data set: The first row corresponds to a denoising task, the second is image inpainting and the last combines both in addition to compression artifact removal. Note that these images were not observed during training. As the data becomes more complex, we can see the importance of the regularization losses proposed in Section 3.2. Using the negative-log-likelihood loss (\mathcal{L}_{nll}) is clearly not sufficient, and a prior trained only with this term is not suited for the latent space optimization. The most important improvement comes from using the latent-noise loss (\mathcal{L}_{ln}). This regularization enforces neighboring elements in latent space to be mapped back to similar images. This is highly beneficial to the gradient descent procedure in latent space and a prior trained with this loss already leads to some good restoration results. Finally, a coarse-to-fine approach is able to handle most cases, in particular high intensity noise levels. This requires training the normalizing flow model with the additional auto-encoder loss (\mathcal{L}_{ae}).

Blind image restoration. We show that the proposed model is applicable to the restoration of generic images. In order to do so, the model must generalize to patches of high resolution good quality images. For this we use the DIV2K dataset (Agustsson & Timofte, 2017) that serves as training and test set for most image super-resolution works. We use the same train/test split with 800 images in the training set and 100 in the test set. Training is done on random image patches of size 64×64 . The normalizing flow architecture used here is very similar to the one described for the Sprites (see supplementary material for details). The restoration of full images of arbitrary size can be done by reconstructing each patch individually. A margin is used to avoid boundary artifacts between patches. Restoration results are presented in Figure 5 for different image degradations. For

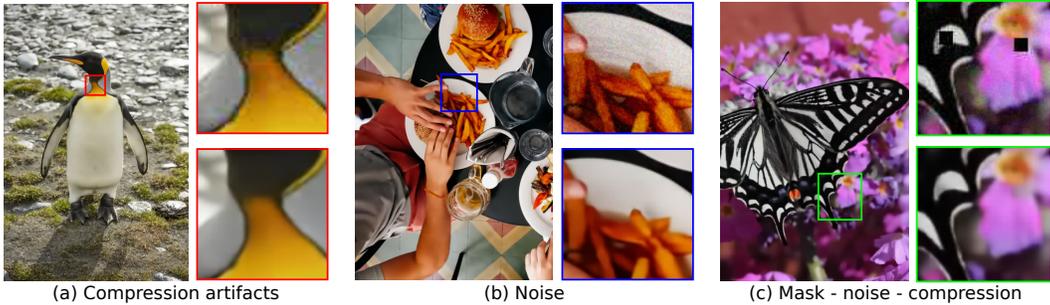


Figure 5: Results on DIV2K dataset. The proposed prior is used to restore arbitrary size images. Degradations include: (a) JPEG compression artifacts; (b) denoising; and (c) a combination of masked regions, noise and compression artifacts.



Figure 6: Compared to DIP, restoring large missing regions is not possible (green), but on this example it produced better denoising results (red).

Type of degradation	DIP	Ours
JPEG artifacts	27.91	30.29
Noise	29.45	28.99
Multiple degradations	25.96	29.87

Figure 7: Quantitative evaluation on DIV2K using PSNR (see text for details).

each example, we show the full resolution result, then focus on a part of the image, illustrating the change.

Comparison with Deep Image Prior (DIP). We first compare the two methods on the images presented in the original DIP paper (Ulyanov et al., 2018). We use our same model trained on the DIV2K dataset. We show competitive restoration results (Figure 1), producing even visually more pleasing reconstruction than DIP on some regions (such as the text and the mouth). The main limit in our case is the patch size used during training. Because of this, it is not possible to inpaint large masked regions such as in the library image (Figure 6). Interestingly however, in this case background regions are better denoised. We also conduct a quantitative evaluation with results presented in Figure 7. Using the test set from DIV2K, we try to restore different degradations: Noise ($\mathcal{N}(0, 5)$), JPEG artifacts and a combination of artifact removal, denoising and inpainting. For this comparison it is unclear how to best set the number of iterations for the DIP. To handle this, we started from the observation that our method converges to the result in approximately 1 hour of computation. Using the DIP online implementation, this corresponds to around 10k optimization steps on the denoising task. We used this maximum number of steps as the threshold for all images and degradations of the test set. The evaluation using PSNR as error metric (Figure 7), demonstrates that our approach is able to achieve competitive results and even outperform DIP on some of the restoration tasks.

6 Conclusion

In this paper, we explored using normalizing flows for capturing the distribution of target high quality content to serve as a prior in blind image restoration. To the best of our knowledge, this is the first time such a direction is explored. One advantage of this formulation is the learned bijective mapping from image to latent space that we use to express the MAP problem of image reconstruction in latent space. We also show the importance of using regularizing losses during training. Finally, we present experimental results illustrating the capacity of the proposed solution to handle different degradations on data sets of varying complexity. We believe this is an exciting new direction as there is still a lot of potential for improvement.

References

- America In Color. <https://www.smithsonianchannel.com/shows/america-in-color/1004516>. Accessed: 2018-03-12.
- Agustsson, E. and Timofte, R. NTIRE 2017 challenge on single image super-resolution: Dataset and study. In *2017 IEEE Conference on Computer Vision and Pattern Recognition Workshops, CVPR Workshops 2017, Honolulu, HI, USA, July 21-26, 2017*, pp. 1122–1131. IEEE Computer Society, 2017. doi: 10.1109/CVPRW.2017.150. URL <https://doi.org/10.1109/CVPRW.2017.150>.
- Bell-Kligler, S., Shocher, A., and Irani, M. Blind super-resolution kernel estimation using an internal-gan. In *Advances in Neural Information Processing Systems*, pp. 284–293, 2019.
- Bigdeli, S. A., Zwicker, M., Favaro, P., and Jin, M. Deep mean-shift priors for image restoration. In *Advances in Neural Information Processing Systems*, pp. 763–772, 2017.
- Cheng, Z., Gadelha, M., Maji, S., and Sheldon, D. A bayesian perspective on the deep image prior. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 5443–5451, 2019.
- Cornillère, V., Djelouah, A., Yifan, W., Sorkine-Hornung, O., and Schroers, C. Blind image super resolution with spatially variant degradations. *ACM Transactions on Graphics (SIGGRAPH Asia Conference Proceedings)*, 2019.
- Dinh, L., Krueger, D., and Bengio, Y. NICE: non-linear independent components estimation. In Bengio, Y. and LeCun, Y. (eds.), *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Workshop Track Proceedings*, 2015. URL <http://arxiv.org/abs/1410.8516>.
- Dinh, L., Sohl-Dickstein, J., and Bengio, S. Density estimation using real NVP. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*, 2017. URL <https://openreview.net/forum?id=HkpbmH91x>.
- Fergus, R., Singh, B., Hertzmann, A., Roweis, S. T., and Freeman, W. T. Removing camera shake from a single photograph. In *ACM SIGGRAPH 2006 Papers*, pp. 787–794. 2006.
- He, K., Sun, J., and Tang, X. Single image haze removal using dark channel prior. *IEEE transactions on pattern analysis and machine intelligence*, 33(12):2341–2353, 2010.
- He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. In Bengio, Y. and LeCun, Y. (eds.), *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015. URL <http://arxiv.org/abs/1412.6980>.
- Kingma, D. P. and Dhariwal, P. Glow: Generative flow with invertible 1x1 convolutions. In *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, 3-8 December 2018, Montréal, Canada.*, pp. 10236–10245, 2018. URL <http://papers.nips.cc/paper/8224-glow-generative-flow-with-invertible-1x1-convolutions>.
- LeCun, Y., Bottou, L., Bengio, Y., and Haffner, P. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- Liu, C., Rubinstein, M., Krainin, M., and Freeman, B. PhotoScan: Taking Glare-Free Pictures of Pictures. <https://ai.googleblog.com/2017/04/photoscan-taking-glare-free-pictures-of.html>. Accessed: 2020-05-25.
- Mataev, G., Milanfar, P., and Elad, M. Deepred: Deep image prior powered by red. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*, pp. 0–0, 2019.

- Meinhardt, T., Moller, M., Hazirbas, C., and Cremers, D. Learning proximal operators: Using denoising networks for regularizing inverse imaging problems. In *Proceedings of the IEEE International Conference on Computer Vision*, pp. 1781–1790, 2017.
- Papamakarios, G., Nalisnick, E., Rezende, D. J., Mohamed, S., and Lakshminarayanan, B. Normalizing flows for probabilistic modeling and inference. *arXiv preprint arXiv:1912.02762*, 2019.
- Park, H. and Mu Lee, K. Joint estimation of camera pose, depth, deblurring, and super-resolution from a blurred image sequence. In *Proceedings of the IEEE International Conference on Computer Vision*, pp. 4613–4621, 2017.
- Pathak, D., Krahenbuhl, P., Donahue, J., Darrell, T., and Efros, A. A. Context encoders: Feature learning by inpainting. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 2536–2544, 2016.
- Rezende, D. J. and Mohamed, S. Variational inference with normalizing flows. In *Proceedings of the 32nd International Conference on Machine Learning, ICML 2015, Lille, France, 6-11 July 2015*, pp. 1530–1538, 2015. URL <http://proceedings.mlr.press/v37/rezende15.html>.
- Rick Chang, J., Li, C.-L., Póczos, B., Vijaya Kumar, B., and Sankaranarayanan, A. C. One network to solve them all—solving linear inverse problems using deep projection models. In *Proceedings of the IEEE International Conference on Computer Vision*, pp. 5888–5897, 2017.
- Rudin, L. I., Osher, S., and Fatemi, E. Nonlinear total variation based noise removal algorithms. *Physica D: nonlinear phenomena*, 60(1-4):259–268, 1992.
- Ulyanov, D., Vedaldi, A., and Lempitsky, V. Deep image prior. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 9446–9454, 2018.
- Wang, Y., Perazzi, F., McWilliams, B., Sorkine-Hornung, A., Sorkine-Hornung, O., and Schroers, C. A fully progressive approach to single-image super-resolution. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pp. 864–873, 2018.
- Zhang, K., Zuo, W., Chen, Y., Meng, D., and Zhang, L. Beyond a gaussian denoiser: Residual learning of deep cnn for image denoising. *IEEE Transactions on Image Processing*, 26(7):3142–3155, 2017a.
- Zhang, K., Zuo, W., Gu, S., and Zhang, L. Learning deep cnn denoiser prior for image restoration. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 3929–3938, 2017b.
- Zhang, Y., Li, K., Li, K., Wang, L., Zhong, B., and Fu, Y. Image super-resolution using very deep residual channel attention networks. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 286–301, 2018.

A Supplementary Material

A.1 Additional Comparison with Deep Image Prior

We provide an additional comparison with Deep Image Prior for the task of compression artifact removal.



Figure 8: JPG artifact removal. We can observe that our results are sharper around the eyes.

A.2 MNIST

For MNIST the network architecture is kept simple, only consisting of a single level. We use $K = 16$ steps in our model. Due to the fact that squeezing layers require the input’s height and width to be divisible by two the input images are zero-padded to size 32×32 .

As coupling transform we use the one depicted in Figure 9 with two blocks ($N = 2$) and 128 intermediate channels ($C_{inter} = 128$). Finally, we choose a Gaussian with unit variance as our base distribution. The Gaussian’s mean is set to a trainable parameter. All other parameters are listed in Table 1.

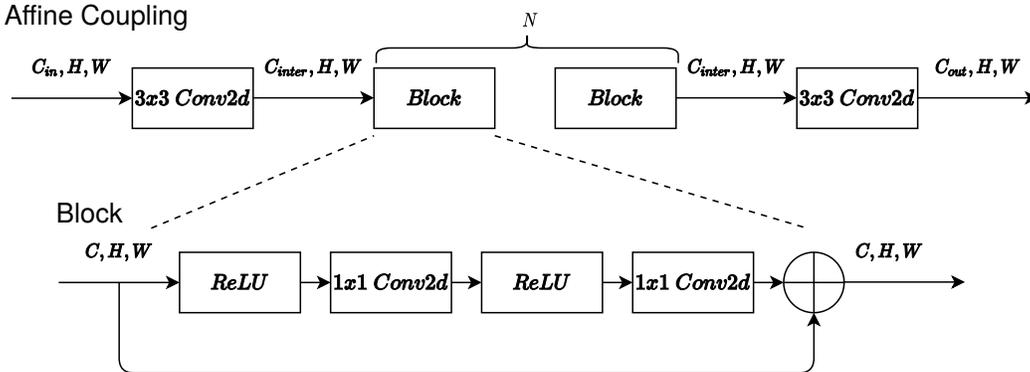


Figure 9: Details of the affine coupling transform. $3 \times 3 \text{ Conv}2d$ and $1 \times 1 \text{ Conv}2d$ refer to standard 2D convolutions using a kernel size of 3×3 and 1×1 respectively. The “+” at the end of the block is an element wise addition.

A.3 Sprites

Each image in the Sprites dataset consists of a figure performing some pose in front of a random background. Figures are centered in the image and have varying color for hair and clothing. Each image is of size 64×64 . Dataset will be made available upon acceptance.

Architecture. For this experiment, the number of levels is set to $L = 3$ and each level has $K = 8$ steps. The distributions $p(\mathbf{u}_1 | \mathbf{h}_1)$ and $p(\mathbf{u}_2 | \mathbf{h}_2)$ depend on a function which computes mean $\mu(\mathbf{h}_i)$ and variance $\sigma(\mathbf{h}_i)$. We call this function the context encoder. A single 2D convolution with kernel

Parameter	Value
# levels	1
# flow blocks per level N_f	16
Affine coupling C_{inter}	128
Base distribution $p(\mathbf{u}_0)$	$\mathcal{N}(\mu, 1)$
optimizer	Adam
learning rate	10^{-4}
batch size	50
# steps	10^5
max gradient value	10^5
max gradient L_2 -norm	10^4

Table 1: Details of architecture and training for the MNIST experiments

size 3×3 and twice the number of output dimension as input dimensions is used as the context encoder. The context encoder’s output is then split in half along the channel dimension. One half is used as $\mu(\mathbf{h}_i)$, the other as $\sigma(\mathbf{h}_i)$. The convolutions weight and bias are initialized to zero for stability reasons. The other parameters for the Sprites dataset are listed in Table 2.

Parameter	Value
# levels (L)	3
# flow steps per level (K)	8
Affine coupling C_{inter}	128
Base distribution $p(\mathbf{u}_1 \mathbf{h}_1), p(\mathbf{u}_2 \mathbf{h}_2)$	$\mathcal{N}(\mu(\mathbf{h}_i), \text{Diag}(\sigma(\mathbf{h}_i)))$
Base distribution $p(\mathbf{u}_0)$	$\mathcal{N}(\mu, \text{Diag}(\sigma))$
Context Encoder $p(\mathbf{u}_1 \mathbf{h}_1), p(\mathbf{u}_2 \mathbf{h}_2)$	zero initialized 2D Convolution, kernel size 3x3
optimizer	Adam
learning rate	10^{-4}
batch size	20
# steps	10^5
max gradient value	10^5
max gradient L_2 -norm	10^4
latent noise magnitude	± 0.5
latent noise loss (β_{ln})	100
autoencoder loss (β_{ae})	1

Table 2: Sprites training specification.

A.4 DIV2K

The number of levels in the architecture is set to $L = 8$ with $K = 4$ steps per level. The number of intermediate channels in the coupling transforms is 256. The context encoder architecture is deepened from 1 to 5 convolutional layer as is illustrated in Figure 10 and a dropout layer is added to the beginning. All the architecture parameters are listed in Table 3.

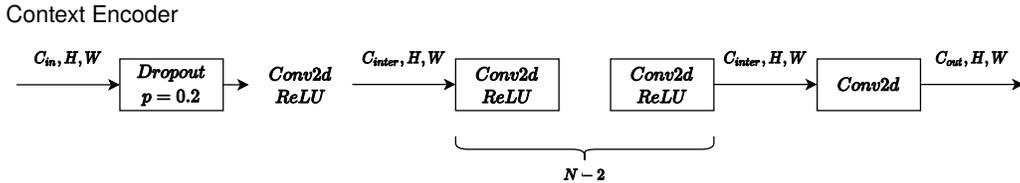


Figure 10: Architecture of the context encoder used for the DIV2K example. A dropout layer with $p = 0.2$ is used as the first layer to prevent overfitting. The last convolution’s weight and bias are initialized to zero for stability reasons.

In addition to this we found that at test time the optimization was faster when the model was trained with additional noise on the images.

The *Image-Noise-loss* \mathcal{L}_{in} works similarly to the *Latent-Noise-loss* \mathcal{L}_{ln} (see Equation 13 in the main paper) except the noise is added to the image \mathbf{x} and distortion is measured on the encoding $\mathbf{u} = T_{\theta}^{-1}(\mathbf{x})$.

$$\mathcal{L}_{in} = \|T_{\theta}^{-1}(\mathbf{x}) - T_{\theta}^{-1}(\mathbf{x} + \eta)\|_2 \quad (20)$$

Parameter	Value
# levels	3
# flow blocks per level N_f	4
Coupling transform C_{inter}	256
Base distribution $p(\mathbf{u}_1 \mathbf{h}_1), p(\mathbf{u}_2 \mathbf{h}_2)$	$\mathcal{N}(\mu(\mathbf{h}_i), \text{Diag}(\sigma(\mathbf{h}_i)))$
Base distribution $p(\mathbf{u}_0)$	$\mathcal{N}(\mu, \text{Diag}(\sigma))$
Context Encoder $p(\mathbf{u}_1 \mathbf{h}_1), p(\mathbf{u}_2 \mathbf{h}_2)$	$N = 5$
optimizer	Adam
learning rate	10^{-4}
batch size	15
# steps	20^5
max gradient value	10^5
max gradient L_2 -norm	10^4
latent noise magnitude	± 0.5
latent noise loss (β_{ln})	100
autoencoder loss (β_{ae})	1
image noise loss (β_{in})	100
image noise magnitude	± 10

Table 3: DIV2K training specification.

Patch-wise Reconstruction. A full image of arbitrary size can be reconstructed by reconstructing each patch individually. To avoid boundary artifacts between patches a margin is used as illustrated in Figure 11. The margin causes overlap between adjacent patches yielding more consistent results in boundary regions.

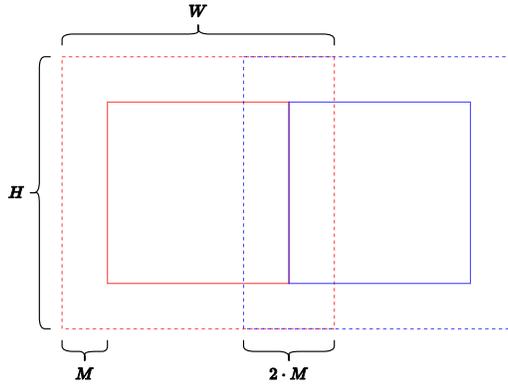


Figure 11: Illustrations of the tiles used for patch-wise reconstruction. H and W refer to the patches height and with respectively. M refers to the margin. Neighboring patches overlap in a region of width $2 \cdot M$. Analogously the same pattern extends in the vertical direction. In our work we use $H = W = 64$ and $M = 4$.