

# Trainable Nonlinear Reaction Diffusion: A Flexible Framework for Fast and Effective Image Restoration

Yunjin Chen and Thomas Pock

**Abstract**—Image restoration is a long-standing problem in low-level computer vision with many interesting applications. We describe a flexible learning framework based on the concept of nonlinear reaction diffusion models for various image restoration problems. By embodying recent improvements in nonlinear diffusion models, we propose a dynamic nonlinear reaction diffusion model with time-dependent parameters (*i.e.*, linear filters and influence functions). In contrast to previous nonlinear diffusion models, all the parameters, including the filters and the influence functions, are simultaneously learned from training data through a loss based approach. We call this approach TNRD – *Trainable Nonlinear Reaction Diffusion*. The TNRD approach is applicable for a variety of image restoration tasks by incorporating appropriate reaction force. We demonstrate its capabilities with three representative applications, Gaussian image denoising, single image super resolution and JPEG deblocking. Experiments show that our trained nonlinear diffusion models largely benefit from the training of the parameters and finally lead to the best reported performance on common test datasets for the tested applications. Our trained models preserve the structural simplicity of diffusion models and take only a small number of diffusion steps, thus are highly efficient. Moreover, they are also well-suited for parallel computation on GPUs, which makes the inference procedure extremely fast.

**Index Terms**—nonlinear reaction diffusion, loss specific training, image denoising, image super resolution, JPEG deblocking



## 1 INTRODUCTION

IMAGE restoration is the process of estimating uncorrupted images from noisy or blurred ones. It is one of the most fundamental operations in image processing, video processing, and low-level computer vision. For several decades, image restoration remains an active research topic and hence new approaches are constantly emerging. There exists a huge amount of literature addressing the topic of image restoration problems, see for example [41] for a survey.

In recent years, the predominant approaches for image restoration are non-local methods based on patch modeling, for example, image denoising with (i) Gaussian noise [40], [26], [15], [46], (ii) multiplicative noise [14], or (iii) Poisson noise [24], image interpolation [47], image deconvolution [20], etc. Most state-of-the-art techniques mainly concentrate on achieving utmost image restoration quality, with little consideration on the computational efficiency, *e.g.*, [40], [26], [47], despite the fact that it is a critical factor for real applications. However, there are a few exceptions. For example, there are two notable exceptions for the task of Gaussian denoising, BM3D [15] and the recently proposed Cascade of Shrinkage Fields (CSF) [52] model, which simultaneously offer high efficiency and high image restoration quality.

It is well-known that BM3D is a highly engineered Gaussian

image denoising algorithm. It involves a block matching process, which is challenging for parallel computation on GPUs, alluding to the fact that it is not straightforward to accelerate BM3D algorithm on parallel architectures. In contrast, the recently proposed CSF model offers high levels of parallelism, making it well suited for GPU implementation, thus owning high computational efficiency.

In this paper, we propose a flexible learning framework to generate fast and effective models for a variety of image restoration problems. Our approach is based on learning optimal nonlinear reaction diffusion models. The learned models preserve the structural simplicity of these models and hence it is straightforward to implement the corresponding algorithms on massive parallel hardware such as GPUs.

### 1.1 Nonlinear diffusion for image restoration

Partial differential equation (PDEs) have become a standard approach for various problems in image processing. On the one hand they come along with a sound mathematical framework that allow to make clear statements about the existence and regularity of the solutions. On the other hand, efficient numerical algorithms have been developed, that allow to compute the solution of PDEs in very short time [45], [55]. Although recent PDE approaches have shown good performance for a number of image processing task, they still fail to produce state-of-the-art quality for classical image restoration tasks.

In the seminal work [45], Perona and Malik (PM) proposed a nonlinear diffusion model, which is given as the following PDE

$$\begin{cases} \frac{\partial u}{\partial t} = \operatorname{div}(g(|\nabla u|)\nabla u) \\ u|_{t=0} = f, \end{cases} \quad (1)$$

- *Y.J. Chen is with the Institute for Computer Graphics and Vision, Graz University of Technology, 8010 Graz, Austria. E-mail: chenyunjin\_nudt@hotmail.com*
- *T. Pock is with the Institute for Computer Graphics and Vision, Graz University of Technology, 8010 Graz, Austria, as well as Digital Safety & Security Department, AIT Austrian Institute of Technology GmbH, 1220 Vienna, Austria. E-mail: pock@icg.tugraz.at*  
*This work was supported by the Austrian Science Fund (FWF) under the China Scholarship Council (CSC) Scholarship Program and the START project BIVISION, No. Y729.*

where  $\nabla$  is the gradient operator,  $t$  denotes the time,  $f$  is a initial image. The function  $g$  is known as edge-stopping function [7] or diffusivity function [55], and a typical  $g$ -function is given by  $g(z) = 1/(1+z^2)$ . The proposed PM diffusion model (1) leads to a nonlinear anisotropic<sup>1</sup> diffusion model which is able to preserve and enhance image edges. Hence, it is well suited for a number of image processing tasks such as image denoising and segmentation.

### 1.1.1 Improvements from the side of PDEs

A first variant of the PM model is the so-called biased anisotropic diffusion (also known as reaction diffusion) proposed by Nordström [43], which introduces a bias term (forcing term) to free the user from the difficulty of specifying an appropriate stopping time for the PM diffusion process. This additional term reacts against the strict smoothing effect of the pure PM diffusion, therefore resulting in a nontrivial steady-state.

Subsequent works consider modifications of the diffusion or the reaction term for the reaction diffusion model [21], [13], [1], *e.g.*, Acton *et al.* [1] exploited a more complicated reaction term to enhance oriented textures; [4] proposed to replace the ordinary diffusion term with a flow equation based on mean curvature. A notable work is the forward-backward diffusion model proposed by Gilboa *et al.* [23], which incorporates explicit inverse diffusion with negative diffusivity coefficient by carefully choosing the diffusivity function. The resulting diffusion processes can adaptively switch between forward and backward diffusion processes. In subsequent work [56], the theoretical framework for discrete forward-and-backward diffusion filtering has been investigated. Researchers also propose to exploit higher-order nonlinear diffusion filtering, which involves larger linear filters, *e.g.*, fourth-order diffusion models [28], [17], [27]. Meanwhile, theoretical properties about the stability and local feature enhancement of higher-order nonlinear diffusion filtering are established in [16].

It should be noted that the above mentioned diffusion models are handcrafted, including elaborate selections of the diffusivity functions, optimal stopping times and proper reaction forces. It is a generally difficult task to design a good-performing PDE for a specific image processing problem, as good insights into this problem and a deep understanding of the behavior of the PDEs are usually required. Therefore, an attempt to learn PDEs from training data via an optimal control approach was made in [39], where the PDEs to be trained have the form of

$$\begin{cases} \frac{\partial u}{\partial t} = \kappa(u) + a(t)^\top \mathcal{O}(u) \\ u|_{t=0} = f. \end{cases} \quad (2)$$

Coefficients  $a(t)$  are free parameters (*i.e.*, combination weights) to train.  $\kappa(u)$  is related to the TV regularization [49] and  $\mathcal{O}(u)$  denotes a set of operators (invariants) over  $u$ , *e.g.*,  $\|\nabla u\|_2^2 = u_x^2 + u_y^2$ .

### 1.1.2 Improvements from the side of image statistics/regularization

As shown in [51], [43], [34], [50], there exist a strong connection between anisotropic diffusion models and variational models

1. Anisotropic diffusion in this paper is understood in the sense that the smoothing induced by PDEs can be favored in some directions and prevented in others. The diffusivity is not necessary to be a tensor. It should be noted that this definition is different from Weickert's terminology [55], where anisotropic diffusion always involves a diffusion tensor, and the PM model is regarded as an isotropic model.

adopting image priors derived from the statistics of natural images. Let us consider the discrete version of the PM model (1), where images are represented as column vectors, *i.e.*,  $u \in \mathbb{R}^N$ . The discrete PM model is formulated as the following discrete PDE with an explicit finite difference scheme

$$\frac{u_{t+1} - u_t}{\Delta t} = - \sum_{i=\{x,y\}} \nabla_i^\top \Lambda(u_t) \nabla_i u_t \doteq - \sum_{i=\{x,y\}} \nabla_i^\top \phi(\nabla_i u_t), \quad (3)$$

where matrices  $\nabla_x$  and  $\nabla_y \in \mathbb{R}^{N \times N}$  are finite difference approximation of the gradient operators in  $x$ -direction and  $y$ -direction, respectively and  $\Delta t$  denotes the time step.  $\Lambda(u_t) \in \mathbb{R}^{N \times N}$  is defined as a diagonal matrix

$$\Lambda(u_t) = \text{diag} \left( g \left( \sqrt{(\nabla_x u_t)_p^2 + (\nabla_y u_t)_p^2} \right) \right)_{p=1, \dots, N},$$

where function  $g$  is the edge-stopping function mentioned before. If ignoring the coupled relation between  $\nabla_x u$  and  $\nabla_y u$ , the PM model can also be written in the form  $\phi(\nabla_i u) = (\phi(\nabla_i u)_1, \dots, \phi(\nabla_i u)_N)^\top \in \mathbb{R}^N$  with function  $\phi(z) = zg(z)$ , known as influence function [7] or flux function [55]. In this paper, we stick to this discrete and decoupled formulation, as it is the starting point of our approach.

As shown in previous works, *e.g.*, [51], [59], the diffusion step (3) corresponds to a gradient descent step to minimize the energy functional given as

$$\mathcal{R}(u) = \sum_{i \in \{x,y\}} \sum_{p=1}^N \rho((k_i * u)_p), \quad (4)$$

the functions  $\rho$  (*e.g.*,  $\rho(z) = \log(1+z^2)$ ) is the so-called penalty function.<sup>2</sup> It is worthwhile to mention that the matrix-vector product  $\nabla_x u$  can be interpreted as a 2D convolution of  $u$  with the linear filter  $k_x = [-1, 1]$  ( $\nabla_y$  corresponds to the linear filter  $k_y = [-1, 1]^\top$ ). The energy functional (4) can be also understood from the aspects of image statistics, image prior and image regularization. As a consequence, a lot of efforts listed below have been made to improve the capability of model (4).

- More filters of larger kernel size were considered in [59], [48], [10], [3], instead of relatively small kernel size, such as usually used pair-wise kernels. The resulting regularization model leads to the so-called fields of experts (FoE) [48] image prior, which works well for many image restoration problems.
- Instead of hand-crafted ones with fixed shape, more flexible penalty functions were exploited, and they were learned from data [59], [50], [34], [52]. Especially, as shown in [59] that those unusual penalties such as inverted penalties (*i.e.*,  $\rho(z)$  decreasing as a function of  $|z|$ ) were found to be necessary.
- In order accelerate the inference phase related to the model (4), in [3], [18], it was proposed to truncate the gradient descent procedure to fixed iterations/stages, and then train this truncated optimization algorithm based on the FoE prior model.
- In order to further increase the flexibility of multi-stage models, Schmidt and Roth [52] considered varying parameters per stage.

## 1.2 Motivations and Contributions

In this paper we concentrate on nonlinear diffusion process due to its high efficiency. Taking into consideration the improvements mentioned in Sec. 1.1.2, we propose a trainable nonlinear diffusion

2. Note that  $\rho'(z) = \phi(z)$ .

model with (1) fixed iterations (also referred to as stages), (2) more filters of larger kernel size, (3) flexible penalties in arbitrary shapes, (4) varying parameters for each iteration, *i.e.*, time varying linear filters and penalties. Then all the parameters (*i.e.*, linear filters and penalties) in the proposed model are simultaneously trained from training data in a supervised way. The proposed approach results in a novel learning framework to train effective image diffusion models. It turns out that the trained diffusion processes leads to state-of-the-art performance, while preserve the property of high efficiency of diffusion based approaches. In summary, our proposed nonlinear diffusion process offers the following advantages:

- 1) It is conceptually simple as it is merely a standard nonlinear diffusion model with trained filters and influence functions;
- 2) It has broad applicability to a variety of image restoration problems. In principle, all the diffusion based models can be revisited with appropriate training;
- 3) It yields excellent results for several tasks in image restoration, including Gaussian image denoising, single image super resolution and JPEG deblocking;
- 4) It is highly computationally efficient, and well suited for parallel computation on GPUs.

A shorter paper has been presented as a conference version [12]. In this paper, we incorporate additional contents listed as follows

- 1) We investigate more details of the training phase, such as the influence of (a) initialization, (b) the model capacity and (c) the number of training samples;
- 2) We consider more detailed analysis of trained models, such as how the trained models generate the patterns;
- 3) We exploit an additional application of single image super resolution to further illustrate the potential breadth of our proposed learning framework.

## 2 PROPOSED REACTION DIFFUSION PROCESS

In this section, we first describe our learning based reaction diffusion model for image restoration, and then we show the relations between the proposed model and existing image restoration models.

### 2.1 Proposed nonlinear diffusion model

The fundamental idea of our proposed learning based reaction diffusion model is described in Sec. 1.2. In addition, we incorporate a reaction term in order to apply our model for different image processing problems, as shown later. As a consequence, our proposed nonlinear reaction diffusion model is formulated as

$$\frac{u_t - u_{t-1}}{\Delta t} = - \underbrace{\sum_{i=1}^{N_k} K_i^t \phi_i^t(K_i^t u_{t-1})}_{\text{diffusion term}} - \underbrace{\psi^t(u_{t-1}, f)}_{\text{reaction term}}, \quad (5)$$

where  $K_i \in \mathbb{R}^{N \times N}$  is a highly sparse matrix, implemented as 2D convolution of the image  $u$  with the filter kernel  $k_i$ , *i.e.*,  $K_i u \Leftrightarrow k_i * u$ ,  $K_i$  is a set of linear filters and  $N_k$  is the number of filters. In practice, we set  $\Delta t = 1$ , as we can freely scale the functions  $\phi_i^t$  and  $\psi^t$  on the right hand side.

Note that our proposed diffusion process is truncated after a few stages, usually less than 10. Moreover, the linear filters and influence functions are adjustable and vary across stages. As

our proposed approach is inspired by nonlinear reaction diffusion model but with trainable filters and influence functions, we coin our method Trainable Nonlinear Reaction Diffusion (TNRD).

In practice, we train the proposed nonlinear diffusion model (5) for specific image restoration problem by exploiting application specific reaction terms  $\psi(u)$ . For classical image restoration problems, such as Gaussian denoising, image deblurring, image super resolution and image inpainting, we can set the reaction term to be the gradient of a data term, *i.e.*  $\psi(u) = \nabla_u \mathcal{D}(u)$ .

For example, if we choose  $\mathcal{D}^t(u, f) = \frac{\lambda^t}{2} \|Au - f\|_2^2$ , we have  $\psi^t(u) = \lambda^t A^\top (Au - f)$ , where  $f$  is the degraded input image,  $A$  is the associated linear operator, and  $\lambda^t$  is related to the strength of the reaction term. In the case of Gaussian denoising,  $A$  is the identity matrix; for image super resolution,  $A$  is related to the down sampling operation and for image deconvolution,  $A$  corresponds to the linear blur kernel.

### 2.2 A more general formulation of the proposed diffusion model

Note that the data term  $\mathcal{D}(u)$  related to (5) should be differentiable. In order to handle the problems involving a non-differentiable data term, *e.g.*, the JPEG deblocking problem investigated in Section 6, we consider a more general form of our proposed diffusion model as follows

$$u_t = \text{Prox}_{\mathcal{G}^t} \left( u_{t-1} - \left( \sum_{i=1}^{N_k} (K_i^t)^\top \phi_i^t(K_i^t u_{t-1}) + \psi^t(u_{t-1}, f) \right) \right), \quad (6)$$

where  $\text{Prox}_{\mathcal{G}^t}(\hat{u})$  is the proximal mapping operation [44] related to the non-differentiable function  $\mathcal{G}^t$ , given as

$$\text{Prox}_{\mathcal{G}^t}(\hat{u}) = \min_u \frac{\|u - \hat{u}\|_2^2}{2} + \mathcal{G}^t(u).$$

### 2.3 Related image restoration models

As mentioned in Sec. 1.1.2, there exist natural connection between anisotropic diffusion and image regularization based energy functional. Therefore, Eq. (6) can be interpreted as a forward-backward step<sup>3</sup> [37] at  $u_{t-1}$  for the energy functional given by

$$E^t(u, f) = \sum_{i=1}^{N_k} \mathcal{R}_i^t(u) + \mathcal{D}^t(u, f) + \mathcal{G}^t(u, f), \quad (7)$$

where  $\mathcal{R}_i^t(u) = \sum_{p=1}^N \rho_i^t((K_i^t u)_p)$  are the regularizers. Since the parameters  $\{K_i^t, \rho_i^t\}$  vary across the stages, (7) is a dynamic energy functional, which changes at each iteration.

If  $\{K_i^t, \rho_i^t\}$  keep the same across stages, the functional (7) with  $\mathcal{G} = 0$  corresponds to the FoE prior regularized variational model for image restoration [48], [11], [10]. In our work, we do not exactly solve this minimization problem anymore. In contrast, we run the gradient descent step for several iterations, and each gradient descent step is optimized by training. More importantly, we are allowed to investigate more generalized penalties.

To the best of our knowledge, Zhu and Mumford [59] were the first to consider learning reaction diffusion models from data. The linear filters appearing in the prior were chosen (not fully trained) from a general filter bank by minimizing the entropy of probability distributions of natural images. The associated penalty functions were learned based on the maximum entropy principle.

3. The forward step is a gradient descent step w.r.t the function  $\mathcal{D}$  and the backward step is the proximal operation w.r.t the function  $\mathcal{G}$ .

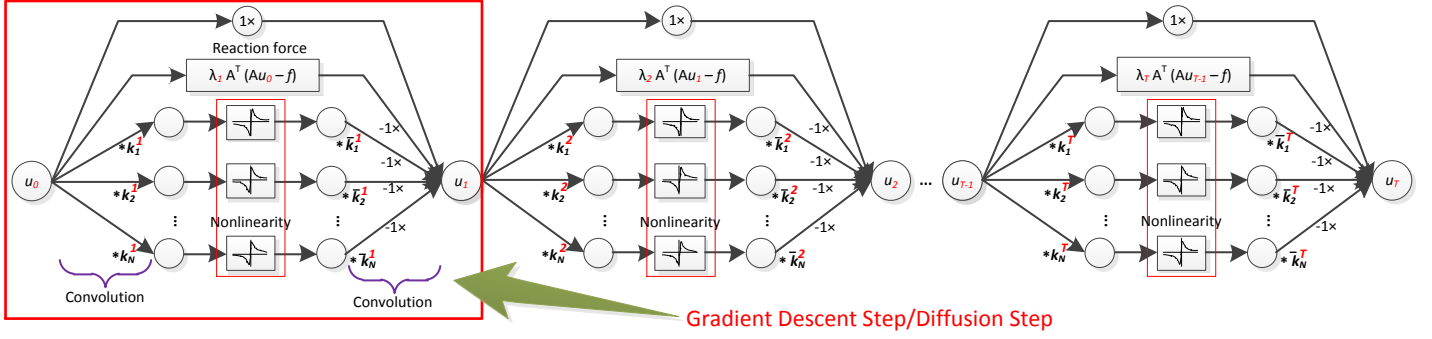


Fig. 1. The architecture of our proposed diffusion model with a reaction force  $\psi^t = \lambda^t A^T (Au_t - f)$  and  $\mathcal{G} = 0$ . It is represented as a feed-forward network. Note that the additional convolution step with the rotated kernels  $\bar{k}_i$  (cf. Equ. 15) does not appear in conventional feed-forward CNs.

However, our proposed diffusion model is a multi-stage diffusion process with multiple image priors, where the filters and penalties are fully trained from clean/degraded pairs.

Some more works have also been dedicated to train the penalties in a diffusion model. In [34], the authors trained the penalty functions in the way that they first computed the image statistics appropriate to certain chosen filters (e.g., horizontal and vertical image derivatives), and then considered mixture models of fixed shape having a peak near zero and two heavy tails in both sides to fit to image statistics. Analogously, in [50], potential functions of fixed shape are chosen to resemble zero mean filter responses.

In very recent work [52], Schmidt and Roth exploited an additive form of half-quadratic optimization to solve problem (7). The resulting model shares similarities with classical wavelet shrinkage and hence it is termed ‘‘cascade of shrinkage fields’’ (CSF). The CSF model relies on the assumption that the data term in (7) is quadratic and the operator  $A$  can be interpreted as a convolution operation, such that the corresponding subproblem can be solved in closed-form using the discrete Fourier transform (DFT). However, our proposed diffusion model does not have this restriction on the data term. In principle, any smooth data term is appropriate. Moreover, as shown in the following sections, we can even handle the case of non-smooth data terms.

As already mentioned, [3], [18] also proposed to train an optimized gradient descent algorithm for the energy functional similar to (7). However, their model is much more constrained, since they exploited the same filters for each gradient descent step and more importantly, they make use of a hand-selected influence function. This clearly restricts the model capability, as demonstrated in Sec. 4.

Comparing our model to the model of [39] (cf. (2)), one can see that this approach learns only a linear model based on pre-defined image invariants. Therefore, this model can be interpreted as a simplified version of our model (5) with fixed linear filters and influence functions, and only the weight of each term is optimized.

The proposed diffusion model also bears an interesting link to convolutional networks (CNs) applied to image restoration problems in [32]. One can see that each iteration (stage) of our proposed diffusion process involves convolution operations with a set of linear filters, and thus it can be treated as a convolutional network. The architecture of our proposed diffusion model is shown in Figure 1, where it is represented as a common feed-forward network. We refer to this network in the following as diffusion network.

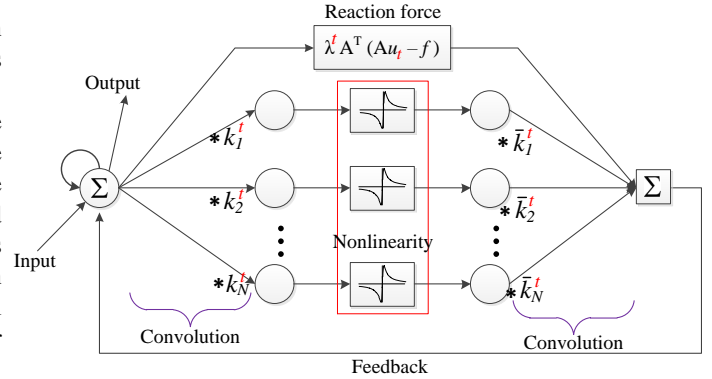


Fig. 2. Our diffusion network can also be interpreted as a CN with a feedback step, which makes it different from conventional feed-forward networks. Due to the feedback step, it can be categorized into recurrent networks [25].

However, we can introduce a feedback step to explicitly illustrate the special architecture of our diffusion network that we subtract ‘‘something’’ from the input image. Therefore, our diffusion model can be represented in a more compact way in Figure 2, where one can see that the structure of our CN model is different from conventional feed-forward networks. Due to this feedback step, it can be categorized into recurrent networks [25]. It should be noted that the nonlinearity (i.e., influence functions in the context of nonlinear diffusion) in our proposed network are trainable. However, conventional CNs make use of fixed activation function, e.g., the ReLU function [42] or sigmoid functions [32].

### 3 LEARNING FRAMEWORK

We train our diffusion networks in a supervised manner, namely we firstly prepare the input/output pairs for a certain image processing task, and then exploit a loss minimization scheme to learn the model parameters  $\Theta_t$  for each stage  $t$  of the diffusion process. The training dataset consists of  $S$  training samples  $\{f^s, u_{gt}^s\}_{s=1}^S$ , where  $f^s$  is a noisy observation and  $u_{gt}^s$  is the corresponding ground truth clean image. The model parameters  $\Theta_t$  of each stage include the parameters of (1) the reaction force weight  $\lambda$ , (2) linear filters and (3) influence functions, i.e.,  $\Theta_t = \{\lambda^t, \phi_i^t, k_i^t\}$ .

#### 3.1 Overall training model

In the supervised manner, a training cost function is required to measure the difference between the output of the diffusion network

and the ground-truth image. As our goal is to train a diffusion network with  $T$  stages, the cost function is formulated as

$$\mathcal{L}(\Theta_{1,\dots,T}) = \sum_{s=1}^S \ell(u_T^s, u_{gt}^s), \quad (8)$$

where  $u_T$  is the output of the final stage  $T$ . In our work we exploit the usual quadratic loss function<sup>4</sup>, defined as

$$\ell(u_T^s, u_{gt}^s) = \frac{1}{2} \|u_T^s - u_{gt}^s\|_2^2. \quad (9)$$

As a consequence, the training task is generally formulated as the following optimization problem

$$\begin{cases} \min_{\Theta} \mathcal{L}(\Theta) = \sum_{s=1}^S \frac{1}{2} \|u_T^s - u_{gt}^s\|_2^2 \\ \text{s.t.} \begin{cases} u_0^s = I_0 \\ u_t^s = \text{Prox}_{\mathcal{G}^t} \left( u_{t-1}^s - \left( \sum_{i=1}^{N_k} (K_i^t)^\top \phi_i^t(K_i^t u_{t-1}^s) + \psi^t(u_{t-1}^s, f^s) \right) \right) \\ t = 1 \dots T, \end{cases} \end{cases} \quad (10)$$

where  $\Theta = \{\Theta^t\}_{t=1}^T$  and  $I_0$  is the initial status of the diffusion process. Note that the above loss function only depends on the output of the final stage  $T$ , *i.e.*, the parameters in all stages are simultaneously trained such that the output of the diffusion process -  $u_T$  is optimized. We call this training scheme **joint training** similar to [52]. The joint training strategy is a minimization problem with respect to the parameters in all stages  $\{\Theta_1, \Theta_2, \dots, \Theta_T\}$ .

One can see that our training model is also a deep model with many stages (layers). It is well-known that deep models are usually sensitive to initialization, and therefore training from scratch is prone to getting stuck at bad local minima. As a consequence, people usually consider a greedy layer-wise pre-training [5] to provide a good initialization for the joint training (fine tune).

In our work, we also consider a **greedy training** scheme similar to [52], to pre-train our diffusion network stage-by-stage, where each stage is greedily trained such that the output of each stage is optimized, *i.e.*, for stage  $t$ , we minimize the cost function

$$\mathcal{L}(\Theta_t) = \sum_{s=1}^S \ell(u_t^s, u_{gt}^s), \quad (11)$$

where  $u_t^s$  is the output of stage  $t$  of the diffusion process. Note that this is a minimization problem only with respect to the parameters  $\Theta_t$  in stage  $t$ .

### 3.2 Parameterizing the influence functions $\phi_i^t$ , linear filters $k_i^t$ and weights $\lambda^t$

In this paper, we aim to investigate arbitrary influence functions. In order to conduct a fast and accurate training, an effective function parameterization method is required. Following the work of [52], we parameterize the influence function via standard radial basis functions (RBFs), *i.e.*, each function  $\phi$  is represented as a weighted linear combination of a family of RBFs as follows

$$\phi_i^t(z) = \sum_{j=1}^M w_{ij}^t \varphi \left( \frac{|z - \mu_j|}{\gamma_j} \right), \quad (12)$$

4. This loss function is related to the PSNR quality measure. Note that as shown in [33], other quality measures, such as structural similarity (SSIM) and mean absolute error (MAE) can be chosen to define the loss function. At present we only consider the quadratic loss function due to its simplicity.

where  $\varphi$  represents different RBFs. In this paper, we exploit RBFs with equidistant centers  $\mu_j$  and unified scaling  $\gamma_j$ . We investigate two typical RBFs [31]: (1) Gaussian radial basis  $\varphi_g$  and (2) triangular-shaped radial basis  $\varphi_t$ , given as

$$\varphi_g(z) = \varphi \left( \frac{|z - \mu|}{\gamma} \right) = \exp \left( -\frac{(z - \mu)^2}{2\gamma^2} \right)$$

and

$$\varphi_t(z) = \varphi \left( \frac{|z - \mu|}{\gamma} \right) = \begin{cases} 1 - \frac{|z - \mu|}{\gamma} & |z - \mu| \leq \gamma \\ 0 & |z - \mu| > \gamma \end{cases}$$

respectively. The basis functions are shown in Figure 3, together with an example of the function approximation by using two different RBF methods. In our work, we have investigated both function approximation methods, and we find that they lead to similar results. We only present the results obtained by the Gaussian RBF in this paper.

In our work, the linear kernels  $k_i^t$  related to the linear operators  $K_i^t$  are defined as a linear combination of Discrete Cosine Transform (DCT) basis kernels  $b_r$ , *i.e.*,

$$k_i^t = \frac{\sum_r \omega_{i,r}^t b_r}{\|\omega_i^t\|_2},$$

where the kernels  $k_i^t$  are normalized to get rid of an ambiguity appearing in the proposed diffusion model. More details can be found in the *supplemental material*. The kernels are formed in this way in order to keep the expected property of zero-mean.

The weights  $\lambda^t$  in our model are constrained to be positive. To this end, we set  $\lambda^t \leftarrow e^{\lambda^t}$  in the training phase for our implementation.

### 3.3 Computing gradients

For both greedy training and joint training, we make use of gradient-based algorithms (*e.g.*, the L-BFGS algorithm [38]) for optimization. The key point is to compute the gradients of the loss function with respect to the training parameters. In greedy training, the gradient of the loss function at stage  $t$  with respect to the model parameters  $\Theta_t$  is computed using standard chain rule, given as

$$\frac{\partial \ell(u_t, u_{gt})}{\partial \Theta_t} = \frac{\partial u_t}{\partial \Theta_t} \cdot \frac{\partial \ell(u_t, u_{gt})}{\partial u_t}, \quad (13)$$

where  $\frac{\partial \ell(u_t, u_{gt})}{\partial u_t} = u_t - u_{gt}$  is directly derived from (9),  $\frac{\partial u_t}{\partial \Theta_t}$  is computed from the diffusion process for specific task. For the applications exploited in this paper, such as image denoising with Gaussian noise, single image super resolution and JPEG deblocking, we present the detailed derivations of  $\frac{\partial u_t}{\partial \Theta_t}$  in the *supplemental material*.

In the joint training, we compute the gradients of the loss function with respect to  $\Theta_t$  by using the standard back-propagation technique widely used in the neural networks learning [35], namely,  $\frac{\partial u_t}{\partial \Theta_t}$  is computed by using

$$\frac{\partial \ell(u_T, u_{gt})}{\partial \Theta_t} = \frac{\partial u_t}{\partial \Theta_t} \cdot \frac{\partial u_{t+1}}{\partial u_t} \dots \frac{\partial \ell(u_T, u_{gt})}{\partial u_T}.$$

Compared to the greedy training, we additionally need to calculate  $\frac{\partial u_{t+1}}{\partial u_t}$ . For the investigated image processing problems in this paper, we provide all necessary derivations in the *supplemental material*.

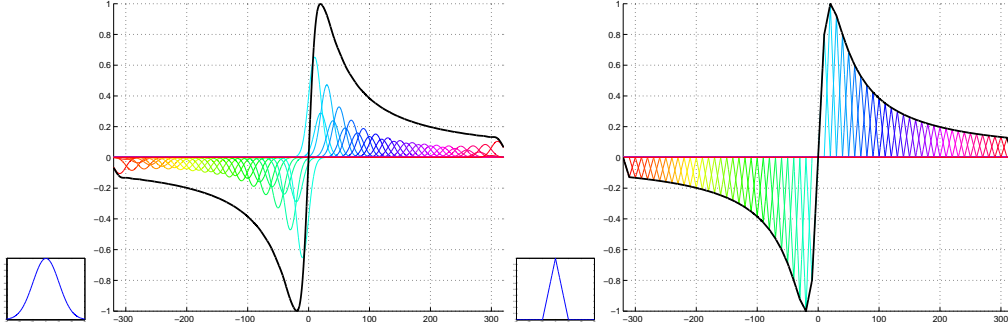


Fig. 3. Function approximation via Gaussian  $\varphi_g(z)$  or triangular-shaped  $\varphi_t(z)$  radial basis function, respectively for the function  $\phi(z) = \frac{2sz}{1+s^2z^2}$  with  $s = \frac{1}{20}$ . Both approximation methods use 63 basis functions equidistantly centered at  $[-310 : 10 : 310]$ .

### 3.4 Experimental setup and implementation details

#### 3.4.1 Boundary condition of the convolution operations

In our convolution based diffusion network, the image size stays the same when an image goes through the network, and we use the symmetric boundary condition for convolution calculation. In our original diffusion model (6), there is matrix transpose  $K^\top v$ , which exactly corresponds to the convolution operation  $\bar{k} * v$  ( $\bar{k}$  is obtained by rotating the kernel  $k$  180 degrees) in the cases of periodic and zero-padding boundary conditions. It should be noted that in the case of symmetric boundary condition used in this paper, this result holds only in the central image region. However, we still want to explicitly use the formulation  $\bar{k} * v$  to replace  $K^\top v$ , because the former can significantly simplify the derivation of the gradients required for training.

We find that the direct replacement introduces some artifacts at the image boundary. In order to avoid these artifacts, we symmetrically pad the input image before it is sent to the diffusion network, and then we discard those padding pixels in the final output. More details are found in the *supplemental material*.

#### 3.4.2 RBF kernels

Images exploited in this paper have the dynamic range in  $[0, 255]$ , and the filters have unit norm. In order to cover most of the filter response, we consider influence functions in the range  $[-310, 310]$ . We use 63 Gaussian RBFs with equidistant centers at  $[-310 : 10 : 310]$ , and set the scaling parameter  $\gamma = 10$ .

#### 3.4.3 Experimental setup

**Model capacity:** In our work, we train the proposed diffusion network with at most 8 stages to observe its saturation behavior after certain stages. We first greedily train  $T$  stages of our model with specific model capacity, then conduct a joint training to refine the parameters of the whole  $T$  stages.

In this paper, we mainly consider four different diffusion networks with increasing capacity:

- TNRD $_{3 \times 3}^T$ , Fully trained model with 8 filters of size  $3 \times 3$ ,
- TNRD $_{5 \times 5}^T$ , Fully trained model with 24 filters of size  $5 \times 5$ ,
- TNRD $_{7 \times 7}^T$ , Fully trained model with 48 filters of size  $7 \times 7$ ,
- TNRD $_{9 \times 9}^T$ , Fully trained model with 80 filters of size  $9 \times 9$ ,

where TNRD $_{m \times m}^T$  denotes a nonlinear diffusion process of stage  $T$  with filters of size  $m \times m$ . The filters number is  $m^2 - 1$ , if not specified. For example, TNRD $_{7 \times 7}^T$  model contains  $(48 \times 48 \text{ (filters)} + 48 \times 63 \text{ (penalties)} + 1 \text{ } (\lambda)) \cdot T = 5329 \cdot T$  free parameters.

**Training and test dataset:** In order to make a fair comparison to previous works, we make use of the same training datasets used in previous works for our training, and then evaluate the trained models on commonly used test datasets. For image processing problems investigated in this paper, *i.e.*, Gaussian denoising, single image super resolution and JPEG deblocking, we consider the following training and test datasets, respectively.

- a) Gaussian denoising. Following [52], we use the same 400 training images, and crop a  $180 \times 180$  region from each image, resulting in a total of 400 training samples of size  $180 \times 180$ , *i.e.*, roughly 13 million pixels. We then evaluate the denoising performance of a trained model on a standard test dataset of 68 natural images, which is suggested by [48], and later widely used for Gaussian denoising testing. Note that the test images are strictly separate from the training datasets.
- b) Single image super resolution. The publicly available framework of Timofte *et al.* [54] provides a perfect base to compare single image super resolution algorithms. It includes 91 training images and two test datasets **Set5** and **Set14**. Many recent state-of-the-art learning based image super resolution approaches [53], [19] accomplish their comparison based on this framework. Therefore, we also use the same 91 training images. We crop 4-5 sub-images of size  $150 \times 150$  from each training image, according to its size, and this finally gives us 421 training samples. We then evaluate the performance of the trained models on the **Set5** and **Set14** dataset.
- c) JPEG deblocking. We train the diffusion models using the same training images as in the case of Gaussian denoising. In the test phase, we follow the test procedure in [33] for performance evaluation. The test images are converted to gray-value, and scaled by a factor of 0.5, resulting 200 images of size  $240 \times 160$ .

#### 3.4.4 Approximate training time

Note that the calculation of the gradients of the loss function in (13) is very efficient even using a simple Matlab implementation, since it mainly involves 2D convolutions. The training time varies greatly for different configurations. Important factors include (1) model capacity, (2) number of training samples, (3) number of iterations taken by the L-BFGS algorithm, and (4) number of Gaussian RBF kernels used for function approximation. We report the most time consuming cases as follows.

In training, computing the gradients  $\frac{\partial \mathcal{L}}{\partial \Theta}$  with respect to the parameters of one stage for 400 images of size  $180 \times 180$  takes about 35s (TNRD $_{5 \times 5}$ ), 75s (TNRD $_{7 \times 7}$ ) or 165s (TNRD $_{9 \times 9}$ ) using Matlab implementation on a server with CPUs: Intel(R) Xeon E5-2680 @ 2.80GHz (eight parallel threads using *parfor*

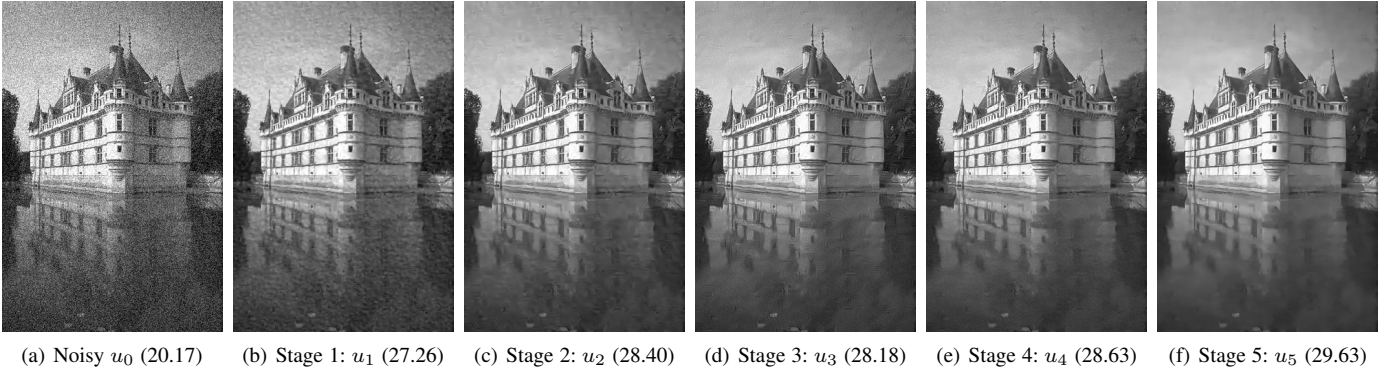


Fig. 4. An image denoising example for noise level  $\sigma = 25$  to illustrate how our learned  $\text{TNRD}_{5 \times 5}^5$  works. (b) - (e) are intermediate results at stage 1 - 4, and (f) is the output of stage 5, *i.e.*, the final denoising result.

in Matlab, 63 Gaussian RBF kernels for the influence function parameterization). We typically run 200 L-BFGS iterations for optimization. Therefore, the total training time, *e.g.*, for the  $\text{TNRD}_{7 \times 7}^5$  model is about  $5 \times (200 \times 75) / 3600 = 20.8h$ . Code for learning and inference is available on the authors' homepage [www.GPU4Vision.org](http://www.GPU4Vision.org)<sup>5</sup>. For the training of the Gaussian denoising task, we have also accomplished a GPU implementation, which is about 4-5 times faster than our CPU implementation.

#### 4 TRAINING FOR GAUSSIAN DENOISING

For the task of Gaussian denoising, we consider the following energy functional

$$\min_u E(u) = \sum_{i=1}^{N_k} \rho_i(k_i * u) + \frac{\lambda}{2} \|u - f\|_2^2.$$

By setting  $\mathcal{D}(u) = \frac{\lambda}{2} \|u - f\|_2^2$  and  $\mathcal{G}(u) = 0$ , we arrive at the following diffusion process with  $u_0 = f$

$$u_t = u_{t-1} - \left( \sum_{i=1}^{N_k} \bar{k}_i^t * \phi_i^t(k_i^t * u_{t-1}) + \lambda^t (u_{t-1} - f) \right), \quad (15)$$

where we explicitly use a convolution kernel  $\bar{k}_i$  (obtained by rotating the kernel  $k_i$  180 degrees) to replace the  $K_i^\top$  for the sake of model simplicity, but we have to pad the input image. The gradients  $\frac{\partial u_t}{\partial \theta_t}$  and  $\frac{\partial u_t}{\partial u_{t-1}}$  required in training are computed from this equation. Detailed derivations are presented in the *supplemental material*.

We started with the training for  $\text{TNRD}_{5 \times 5}^T$ . We first considered the greedy training phase to train a diffusion process up to 8 stages (*i.e.*,  $T \leq 8$ ), in order to observe the asymptotic behavior of the diffusion network. In the greedy training, only parameters in one stage were trained at a time. We exploited a plain initialization to start the training, namely linear filters and influence functions were initialized from the modified DCT filters and the function  $\phi(z) = 2z/(1+z^2)$ , respectively.

After the greedy training was completed, we conducted joint training for a diffusion model of certain stages (*e.g.*,  $T = 2, 5, 8$ ), to simultaneously tune the parameters in all stages. We initialized the joint training using parameters learned in greedy training, as this is guaranteed not to degrade the training performance.

We first trained our diffusion models for the Gaussian denoising problem with standard deviation  $\sigma = 25$ . The noisy training

images were generated by adding synthetic Gaussian noise with  $\sigma = 25$  to the clean images. Once we obtained a trained model, we evaluated its denoising performance on 68 natural images following the same test protocol as in [52], [10]. Figure 4 shows a denoising example for noise level  $\sigma = 25$  to illustrate the denoising process of our learned  $\text{TNRD}_{5 \times 5}^5$  diffusion network.

We present the final results of the joint training in Table 1, together with a selection of recent state-of-the-art denoising algorithms, namely BM3D [15], LSSC [40], EPLL-GMM [60], opt-MRF [10], RTF model [33], the CSF model [52] and WNNM [26], as well as two similar approaches ARF [3] and opt-GD [18], which also train an optimized gradient descent inference. We downloaded these algorithms from the corresponding author's homepage, and used them as is. Unfortunately, we are not able to present comparisons with [39], [32], as their codes are not available.

From Table 1, one can see that (1) the performance of the  $\text{TNRD}_{5 \times 5}^T$  model saturates after stage 5, *i.e.*, in practice, 5 stages are typically enough; (2) our  $\text{TNRD}_{5 \times 5}^5$  model has achieved significant improvement (28.78 vs. 28.60), compared to a similar model  $\text{CSF}_{5 \times 5}^5$ , which has the same model capacity and (3) moreover, our  $\text{TNRD}_{5 \times 5}^8$  model is on par with so far the best-reported algorithm - WNNM. It turns out that our trained models perform surprisingly well for image denoising. Then, a natural question arises: what is the critical factor for the effectiveness of the trained diffusion models?

#### 4.1 Understanding the proposed diffusion models

There are actually two main aspects in our training model: (1) the linear filters and (2) the influence functions. In order to have a better understanding of the trained models, we went through a series of experiments to investigate the impact of these two aspects.

Concentrating on the model capacity of 24 filters of size  $5 \times 5$ , we considered the training of a diffusion process with 10 steps, *i.e.*,  $T = 10$  for the Gaussian denoising of noise level  $\sigma = 25$ . We exploited two main classes of configurations: (A) the parameters of every stage are the same and (B) every diffusion stage is different from each other. In both configurations, we consider two cases: (I) only train the linear filters with fixed influence function  $\phi(z) = 2z/(1+z^2)$  and (II) simultaneously train the filters and influence functions.

Based on the same training dataset and test dataset, we obtained the following results: (A.I) every diffusion step is the same, and only the filters are optimized with fixed influence function.

5. <http://gpu4vision.icg.tugraz.at/binaries/nonlinear-diffusion.zip#pub94>

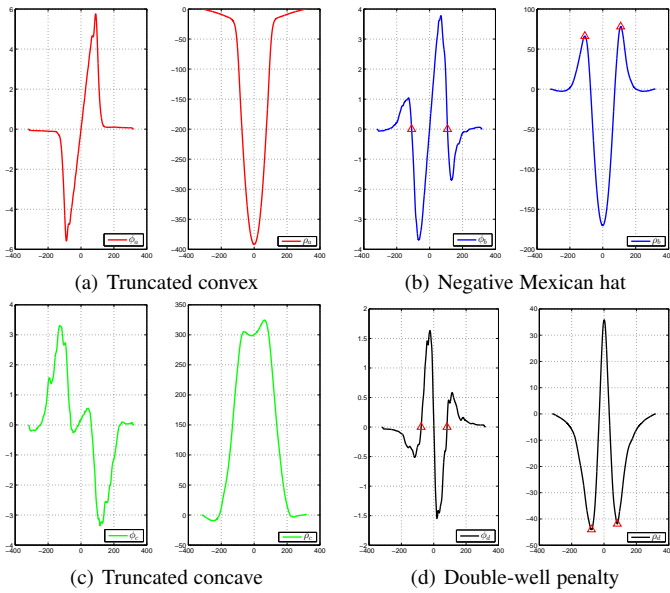


Fig. 5. The figure shows four characteristic influence functions (left plot in each subfigure) together with their corresponding penalty functions (right plot in each subfigure), learned by our proposed method in the  $\text{TNRD}_{5 \times 5}^3$  model. A major finding in this paper is that our learned penalty functions significantly differ from the usual penalty functions adopted in partial differential equations and energy minimization methods. In contrast to their usual robust smoothing properties which is caused by a single minimum around zero, most of our learned functions have multiple minima different from zero and hence are able to enhance certain image structures. See Sec. 4.3 for more information.

This is a similar configuration to previous works [3], [18]. The trained model achieves a test performance of 28.47dB. (A.II) with additional tuning of the influence functions, the resulting performance is boosted to 28.60dB. (B.I) every diffusion step can be different, but only the linear filters are trained with fixed influence functions. The corresponding model obtains a result of 28.56dB, which is equivalent to the variational model [10] with the same model capacity. Finally (B.II) with additional optimization of the influence functions, the trained model leads to a significant improvement with the result of 28.86dB.

The analytical experiments demonstrate that without the training of the influence functions, there is no chance to achieve significant improvements over previous works, no matter how hard we tune the linear filters. Therefore, we believe that the most critical factor of our proposed training model lies in the adjustable influence functions. A closer look at the learned influence functions of the  $\text{TNRD}_{5 \times 5}^3$  model in Sec.4.2 strengthens our argument.

Comparing our proposed TNRD model to the CSF model [52], one can see that the degree of freedom is in principle the same, since in both models the filters and the non-linear functions can be learned. Therefore, one would expect a similar performance of both models in practice. However, it turns out the performance of the CSF model is inferior to our TNRD model in the case of Gaussian denoising task. The reason for the performance gap is still unclear and we plan to investigate it in future work.

## 4.2 Learned influence functions

A close inspection of the learned 120 penalty functions<sup>6</sup>  $\rho$  in the  $\text{TNRD}_{5 \times 5}^3$  model demonstrated that most of the penalties resemble

6. The penalty function  $\rho(z)$  is integrated from the influence function  $\phi(z)$  according to the relation  $\phi(z) = \rho'(z)$

four representative shapes shown in Figure 5.

- Truncated convex penalty functions with low values around zero to promote smoothness.
- Negative Mexican hat functions, which have a local minimum at zero and two symmetric local maxima.
- Truncated concave functions with smaller values at the two tails.
- Double-well functions, which have a local maximum (not a minimum any more) at zero and two symmetric local minima.

At first glance, the learned penalty functions (except (a)) differ significantly from the usually adopted penalty functions used in PDE and energy minimization methods. However, it turns out that they have a clear meaning for image regularization.

Regarding the penalty function (b), there are two critical points (indicated by red triangles). When the magnitude of the filter response is relatively small (*i.e.*, less than the critical points), probably it is stimulated by the noise and therefore the penalty function encourages smoothing operation as it has a local minimum at zero. However, once the magnitude of the filter response is large enough (*i.e.*, across the critical points), the corresponding local patch probably contains a real image edge or certain structure. In this case, the penalty function encourages to increase the magnitude of the filter response, alluding to an image sharpening operation. Therefore, the diffusion process controlled by the influence function (b), can adaptively switch between image smoothing (forward diffusion) and sharpening (backward diffusion). We find that the learned influence function (b) is closely similar to an elaborately designed function in a previous work [23], which leads to an adaptive forward-and-backward diffusion process.

Similar forms of the learned penalty function in (c) with a concave shape are also observed in previous work on image prior learning [59]. This penalty function also encourages to sharpen the image edges. Concerning the learned penalty function (d), as it has local minima at two specific points, it prefers specific image structure, implying that it helps to form certain image structure. We also find that this penalty function is exactly the type of bimodal expert functions for texture synthesis employed in [30].

Therefore, our learned penalty functions confirmed existing robust penalties based prior models and many priors exploiting some unusual penalties, which can produce patterns and enhance preferred features. As a consequence, the diffusion process involving the learned influence functions does not perform pure image smoothing any more for image processing. In contrast, it leads to a diffusion process for adaptive image smoothing and sharpening, distinguishing itself from previous commonly used image regularization techniques.

## 4.3 Pattern formation using the learned influence functions

In the previous work on Gibbs reaction diffusion<sup>7</sup> [59], it is shown that those unconventional penalty functions such as Figure 5(c) have significant meaning in visual computation, as they can produce patterns. We also find that those unconventional penalty functions learned in our models can produce some interesting image patterns.

7. The terminology of “reaction diffusion” in [59] is a bit different from ours. In our formulation, “reaction term” is related to the data term, while in [59], it means the diffusion term controlled by those downright penalty functions.



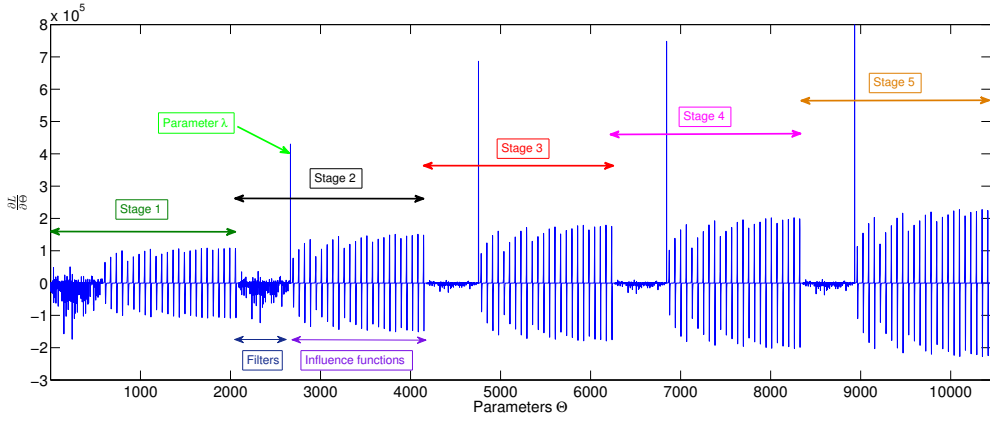


Fig. 6. Well distributed gradients  $\frac{\partial L}{\partial \Theta}$  over stages for the  $\text{TNRD}_{5 \times 5}^5$  model at the initialization point  $\Theta_0$  with a plain setting. One can see that the “vanishing gradient” phenomenon [6] in the back-propagation phase of a conventional deep model does not appear in our training model.

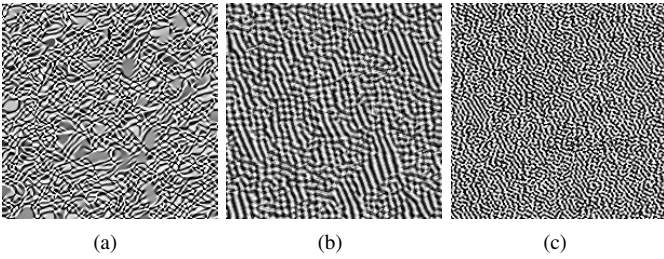


Fig. 7. Patterns synthesized from uniform noise using our learned diffusion models. (a) is generated by (16) using the parameters (linear filters and influence functions) in a stage of our learned  $\text{TNRD}_{5 \times 5}^5$  for image denoising, (b) is generated by (16) using the parameters in a stage of our learned  $\text{TNRD}_{7 \times 7}^5$  for image super resolution and (c) is also from a stage of our learned  $\text{TNRD}_{7 \times 7}^5$  for image super resolution.

We consider the following diffusion process involving our learned linear filters and the corresponding influence functions

$$\frac{u_t - u_{t-1}}{\Delta t} = - \sum_{i=1}^{N_k} \bar{k}_i * \phi_i(k_i * u_{t-1}), \quad (16)$$

where the filters  $k_i$  and influence functions  $\phi_i$  are chosen from a certain stage of the learned models. Note that we do not incorporate a reaction term in this diffusion model. We run (16) from starting points  $u_0$  (uniform noise images in the range  $[0, 255]$ ), and it converges to a local minimum<sup>8</sup>. Some synthesized patterns are shown in Figure 7. One can see that the diffusion model with our learned influence functions and filters can produce edge-like image structure and repeated patterns from fully random images. This kind of diffusion model is known as Gibbs reaction diffusion in [59]. We provide another example in Figure 13 to demonstrate how our learned diffusion models can generate meaningful patterns for image super resolution.

## 4.4 Important aspects of the training framework

### 4.4.1 Influence of initialization

Our training model is also a deep model with many stages (layers), but we find that it is not very sensitive to initialization. Based on the training for Gaussian denoising, we conducted experiments with fully random initializations and some plain settings.

<sup>8</sup>. The corresponding diffusion processes are unstable, and therefore we have to restrict the image dynamic range to  $[0, 255]$ .

We firstly investigated the case of greedy training where the model was trained stage by stage and the parameters of one stage were trained at a time. We initialized the parameters using fully random numbers in the range  $[-0.5, 0.5]$ . It turns out that the resulting models with different initializations lead to a deviation within 0.01dB in the test phase. That is to say, the greedy training strategy is not sensitive to initialization.

Then, we considered the case of joint training, where all the parameters in all stages were trained at a time. We also initialized the training with fully random numbers in the range  $[-0.5, 0.5]$ . In this case, it turns out that the resulting models lead to inferior results, e.g., in the case of  $\text{TNRD}_{5 \times 5}^5$  (28.61 vs.28.78). However, plain initializations can generate equivalent results. For example, we considered a plain initialization (all stages were initialized from the modified DCT filters and an unified influence function  $\phi(z) = 2z/(1+z^2)$ ), the resulting models performed almost the same as those models trained from some good initializations such as parameters obtained from greedy training, e.g.,  $\text{TNRD}_{5 \times 5}^5$  (28.75 vs.28.78) and  $\text{TNRD}_{7 \times 7}^5$  (28.91 vs.28.92).

We believe that this appealing property of our training framework is attributed to the well-distributed gradients across stages. We show in Figure 6 an example to illustrate the gradients of the training loss function with respect to the parameter of all stages. One can see that the well-known phenomenon of “vanishing gradient” [6] in the back-propagation phase of a usual deep model does not appear in our training model. We believe that the reason for the well-distributed gradients is that our training model is more constrained. For example, in a more general sense, the rotated kernel  $\bar{k}_i$  in our formulation is not necessary to be the rotated version of the kernel  $k_i$ , and it can be an arbitrary kernel. However, we stick to this form, as it has a clear meaning derived from energy minimization.

### 4.4.2 Influence of the number of training samples

In our training, we do not consider any regularization for the training parameters, and we finally reach good-performing models. A probable reason is that we have exploited sufficient training samples (400 samples of size  $180 \times 180$ ). Thus an interesting question arises: how many samples are sufficient for our training?

In order to answer this question, we re-train the  $\text{TNRD}_{5 \times 5}^5$  model using different size of training dataset, and then evaluate the test performance of trained models. We summarize the results in Figure 8. One can see that (1) too few training samples (e.g.,

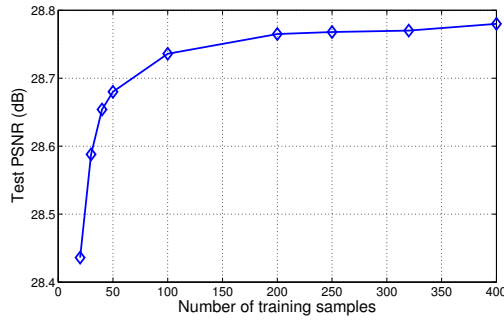


Fig. 8. Influence of the number of training examples for the training model  $\text{TNRD}_{5 \times 5}^{5 \times 5}$

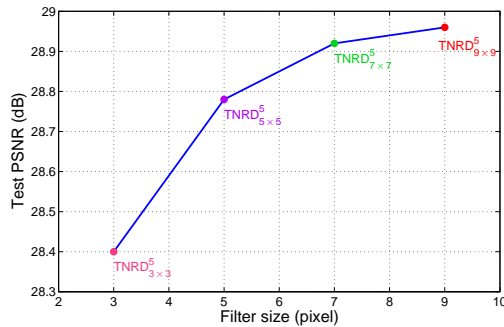


Fig. 9. Influence of the filter size (based on a relatively small training data set of 400 images of size  $180 \times 180$ )

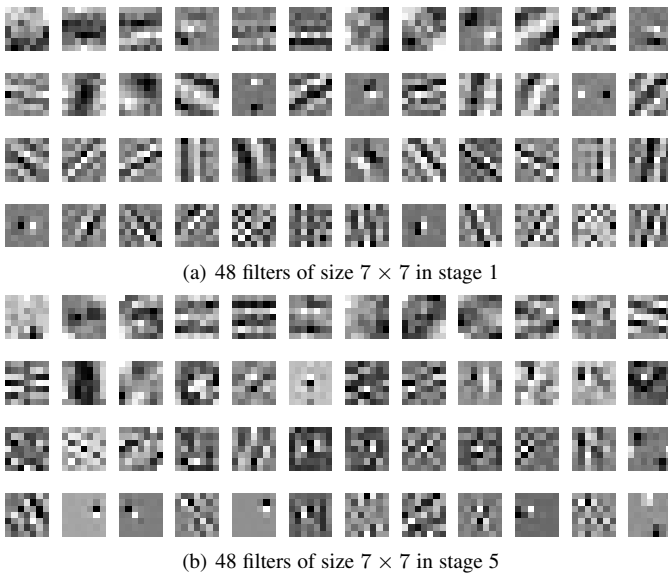


Fig. 10. Trained filters (in the first and last stage) of the  $\text{TNRD}_{7 \times 7}^{5 \times 7}$  model for the noise level  $\sigma = 25$ . We can find first, second and higher-order derivative filters, as well as rotated derivative filters along different directions. These filters are effective for image structure detection, such as image edge and texture.

40 images) will clearly lead to over-fitting, thus inferior test performance, and (2) 200 images are typically enough to prevent over-fitting.

#### 4.4.3 Influence of the filter size

In our model, the size of involved filters is a free parameter. In principle, we can exploit filters of any size, but in practice, we need to consider the trade-off between run time and accuracy.

In order to investigate the influence of the filter size, we increase the filter size to  $7 \times 7$  and  $9 \times 9$ . We find that increasing the filter size from  $5 \times 5$  to  $7 \times 7$  brings a significant improvement of 0.14dB ( $\text{TNRD}_{7 \times 7}^{5 \times 7}$  vs.  $\text{TNRD}_{5 \times 5}^{5 \times 5}$ ) as show in Table 1. However, when we further increase the filter size to  $9 \times 9$ , the resulting  $\text{TNRD}_{9 \times 9}^{5 \times 9}$  only leads to a performance of 28.96dB (a slight improvement of 0.05dB relative to the  $\text{TNRD}_{7 \times 7}^{5 \times 7}$  model). We can conjecture that further increasing the filter size to  $11 \times 11$  might bring negligible improvements.

Note that the above conclusion is drawn from a relatively small training data set of 400 images of size  $180 \times 180$ . It should be mentioned that when the size of the model increase, the size of training data set should also increase to avoid over-fitting. However, our current CPU implementation for training prevents us from training with larger model and large-scale data sets (millions). A faster implementation on GPUs together with the stochastic gradient descent optimization strategy is left to future work.

We also consider a model with smaller filters,  $3 \times 3$ . We summarize the results of different model capacities in Figure 9. In practice, we prefer the  $\text{TNRD}_{7 \times 7}^{5 \times 7}$  model as it provides the best trade-off between performance and computation time. Therefore, in later applications, we only consider  $\text{TNRD}_{7 \times 7}^{5 \times 7}$  models.

Fig. 10 shows the trained filters of the  $\text{TNRD}_{7 \times 7}^{5 \times 7}$  model in the first and last stage for the task of Gaussian denoising. One can find many edge and image structure detection filters along different directions and in different scales.

## 4.5 Training for different noise levels and comparison to recent state-of-the-arts

The above training experiments are based on Gaussian noise of level  $\sigma = 25$ . We also trained diffusion models for the noise levels  $\sigma = 15$  and  $\sigma = 50$ . The test performance is summarized in Table 1, together with comparison to very recent state-of-the-art denoising algorithms. In experiments, we observed that joint training can always gain an improvement of about 0.1dB over the greedy training for the cases of  $T \geq 5$ .

From Table 1, one can see that for all noise levels, the resulting  $\text{TNRD}_{7 \times 7}^{5 \times 7}$  model achieves the highest average PSNR. The  $\text{TNRD}_{7 \times 7}^{5 \times 7}$  model outperforms the benchmark - BM3D method by 0.35dB in average. This is a notable improvement as few methods can surpass BM3D more than 0.3dB in average [36]. Moreover, the  $\text{TNRD}_{7 \times 7}^{5 \times 7}$  model also surpasses the best-reported algorithm - WNNM method, which is quite slow as shown in Table 2.

## 4.6 Run time

The algorithm structure of our TNRD model is similar to the CSF model, which is well-suited for parallel computation on GPUs. We implemented our trained models on GPU using CUDA programming to speed up the inference procedure, and finally it leads to significantly improved performance, see Table 2. We make a run time comparison to other denoising algorithms based on strictly enforced single-threaded CPU computation (e.g., start Matlab with `-singleCompThread`) for a fair comparison, see Table 2. We only present the results of some selective algorithms, which either have the best denoising result or run time performance. We refer to [52] for a comprehensive run time comparison of various algorithms<sup>9</sup>.

9. LSSC, EPLL, opt-MRF and  $\text{RTF}^5$  methods are much slower than BM3D on the CPU, cf. [52].

Method	$\sigma$			St.	TNRD <sub>5×5</sub>	TNRD <sub>7×7</sub>
	15	25	50		$\sigma = 15$	
BM3D	31.08	28.56	25.62	2	31.14	31.30
LSSC	31.27	28.70	25.72	5	31.30	<b>31.42</b>
EPLL-GMM	31.19	28.68	25.67	8	31.34	<b>31.43</b>
opt-MRF <sub>7×7</sub>	31.18	28.66	25.70		$\sigma = 25$	
RTF <sup>5</sup>	–	28.75	–	2	28.58	28.77
WNNM	31.37	28.83	25.83	5	28.78	<b>28.92</b>
CSF <sub>5×5</sub> <sup>5</sup>	31.14	28.60	–	8	28.83	<b>28.95</b>
CSF <sub>7×7</sub> <sup>2</sup>	31.24	28.72	–		$\sigma = 50$	
ARF <sub>5×5</sub> <sup>4</sup>	30.70	28.20	–	2	25.54	25.78
opt-GD <sub>5×5</sub> <sup>10</sup>	–	28.39	–	5	25.80	<b>25.96</b>
				8	<b>25.87</b>	<b>26.01</b>

TABLE 1

Average PSNR (dB) on 68 images from [48] for image denoising with  $\sigma = 15, 25$  and  $50$ . All the TNRD models are jointly trained. Note that among those algorithms similar to our model, opt-MRF<sub>7×7</sub>, ARF<sub>5×5</sub><sup>4</sup> and opt-GD<sub>5×5</sub><sup>10</sup> only train the filters with fixed penalty function  $\log(1 + z^2)$ . In the opt-MRF<sub>7×7</sub> model, 48 filters of size  $7 \times 7$  (**2304 free parameters**), for ARF<sub>5×5</sub><sup>4</sup>, 13 filters of size  $5 \times 5$  (**325 free parameters**) and for the opt-GD<sub>5×5</sub><sup>10</sup> algorithm, 24 filters of size  $5 \times 5$  (**600 free parameters**) are trained. The CSF model and our approach train both the filters and nonlinearities, thus involving more parameters, e.g., the TNRD<sub>7×7</sub><sup>5</sup> model involves **26,645 free parameters** and the corresponding CSF<sub>7×7</sub><sup>2</sup> model involves **24,245 free parameters**.

Method	256 <sup>2</sup>	512 <sup>2</sup>	1024 <sup>2</sup>	2048 <sup>2</sup>	3072 <sup>2</sup>
BM3D [15]	1.1	4.0	17	76.4	176.0
CSF <sub>7×7</sub> <sup>2</sup> [52]	3.27	11.6	40.82	151.2	494.8
WNNM [26]	122.9	532.9	2094.6	–	–
	0.51	1.53	5.48	24.97	53.3
TNRD <sub>5×5</sub> <sup>5</sup>	<b>0.43</b>	<b>0.78</b>	<b>2.25</b>	<b>8.01</b>	<b>21.6</b>
	<b>0.005</b>	<b>0.015</b>	<b>0.054</b>	<b>0.18</b>	<b>0.39</b>
	1.21	3.72	14.0	62.2	135.9
TNRD <sub>7×7</sub> <sup>5</sup>	<b>0.56</b>	<b>1.17</b>	<b>3.64</b>	<b>13.01</b>	<b>30.1</b>
	<b>0.01</b>	<b>0.032</b>	<b>0.116</b>	<b>0.40</b>	<b>0.87</b>

TABLE 2

Run time comparison for image denoising (in seconds) with different implementations. (1) The run time results with gray background are evaluated with the single-threaded implementation on Intel(R) Xeon(R) CPU E5-2680 v2 @ 2.80GHz; (2) the blue colored run times are obtained with multi-threaded computation using Matlab *parfor* on the above CPUs; (3) the run time results colored in red are executed on a NVIDIA GeForce GTX 780Ti GPU. We do not count the memory transfer time between CPU/GPU for the GPU implementation (if counted, the run time will nearly double)

We see that our TNRD model is generally faster than the CSF model with the same model capacity. It is reasonable, because in each stage the CSF model involves additional DFT and inverse DTF operations, i.e., our model only requires a portion of the computation of the CSF model. Even though the BM3D is a non-local model, it still possesses high computational efficiency. In contrast, another non-local model - WNNM achieves compelling denoising results at the expense of huge computation time. Moreover, the WNNM algorithm is hardly applicable for high resolution images (e.g., 10 mega-pixels) due to its huge memory requirements. Note that our model can be also easily implemented with multi-threaded CPU computation.

In summary, our TNRD<sub>7×7</sub><sup>5</sup> model outperforms these recent state-of-the-arts, meanwhile it is the fastest method even with a CPU implementation. We present an illustrative denoising example in Figure 11 on an image from the test dataset. More denoising examples can be found in the *supplemental material* based on images from the test dataset and a megapixel-size natural image of size  $1050 \times 1680$ .

## 5 SINGLE IMAGE SUPER RESOLUTION (SISR)

As demonstrated in the last section that our trained diffusion model can lead to explicit backward diffusion process, which sharpens image structures like edges. This is the very property demanded for the task of image super resolution. Therefore, we are motivated to investigate the SISR problem with our proposed approach.

We start with the following energy functional

$$\min_u E(u) = \sum_{i=1}^{N_k} \rho_i(k_i * u) + \frac{\lambda}{2} \|Au - f\|_2^2, \quad (17)$$

where the linear operator  $A$  is a bicubic interpolation which links the high resolution (HR) image  $h$  to the low resolution (LR) image  $f$  via  $f = Ah$ . Casting  $\mathcal{D}(u) = \frac{\lambda}{2} \|Au - f\|_2^2$  and  $\mathcal{G}(u) = 0$ , the energy functional (17) suggests the following diffusion process

$$u_t = u_{t-1} - \left( \sum_{i=1}^{N_k} \bar{k}_i^t * \phi_i^t(k_i^t * u_{t-1}) + \lambda^t A^\top (Au_{t-1} - f) \right), \quad (18)$$

where the starting point  $u_0$  is given by the direct bicubic interpolation of the LR image  $f$ . Computing the gradients  $\frac{\partial u_t}{\partial \theta_t}$  and  $\frac{\partial u_t}{\partial u_{t-1}}$  with respect to (18) can be done with little modifications to the derivations for image denoising. Detailed derivations are presented in the *supplemental material*.

We considered the model capacity of TNRD<sub>7×7</sub><sup>5</sup>, and trained diffusion models for three upscaling factors  $\times 2$ ,  $\times 3$  and  $\times 4$ , using exactly the same 91 training images as in previous works [54], [53]. The trained models are evaluated on two different test data sets: **Set5** and **Set14**. Following previous works [54], [19], [53], the trained models are only applied to the luminance component of an image, and a regular bicubic upscaling method is applied to the color components.

The test results are summarized in Table 3 and Table 4. One can see that in terms of average PSNR, our trained diffusion model TNRD<sub>7×7</sub><sup>5</sup> leads to significant improvements over very recent state-of-the-arts in all cases, meanwhile it is still among the fast algorithms<sup>10</sup>. A SISR example is shown in Figure 12 to illustrate its effectiveness. One can see that our approach can obtain more accurate image edges, as shown in the zoom-in parts. More SISR examples can be found in the *supplemental material*.

We apply the learned diffusion parameters to the diffusion equation (16). It turns out that the diffusion process can also generate some interesting patterns from random images, as shown in Figure 7. We believe that this ability to generate image patterns from weak evidence is the main reason for the superiority of our trained model for the SISR task. In order to further validate our argument, we carry out a toy SISR experiment based on a synthesized image with repeated hexagons. The results are shown in Figure 13, where one can see that our trained model can better reconstruct those repeated image structures.

## 6 JPEG DEBLOCKING EXPERIMENTS

In order to further demonstrate the applicability of our proposed framework for those problems with a non-smooth data term, we investigate the JPEG deblocking problem - suppressing the block artifacts in the JPEG compressed images, which is formulated as a

<sup>10</sup>. Note that our approach is a Matlab implementation, while some of other algorithms are based on C++ implementations, such as SR-CNN.

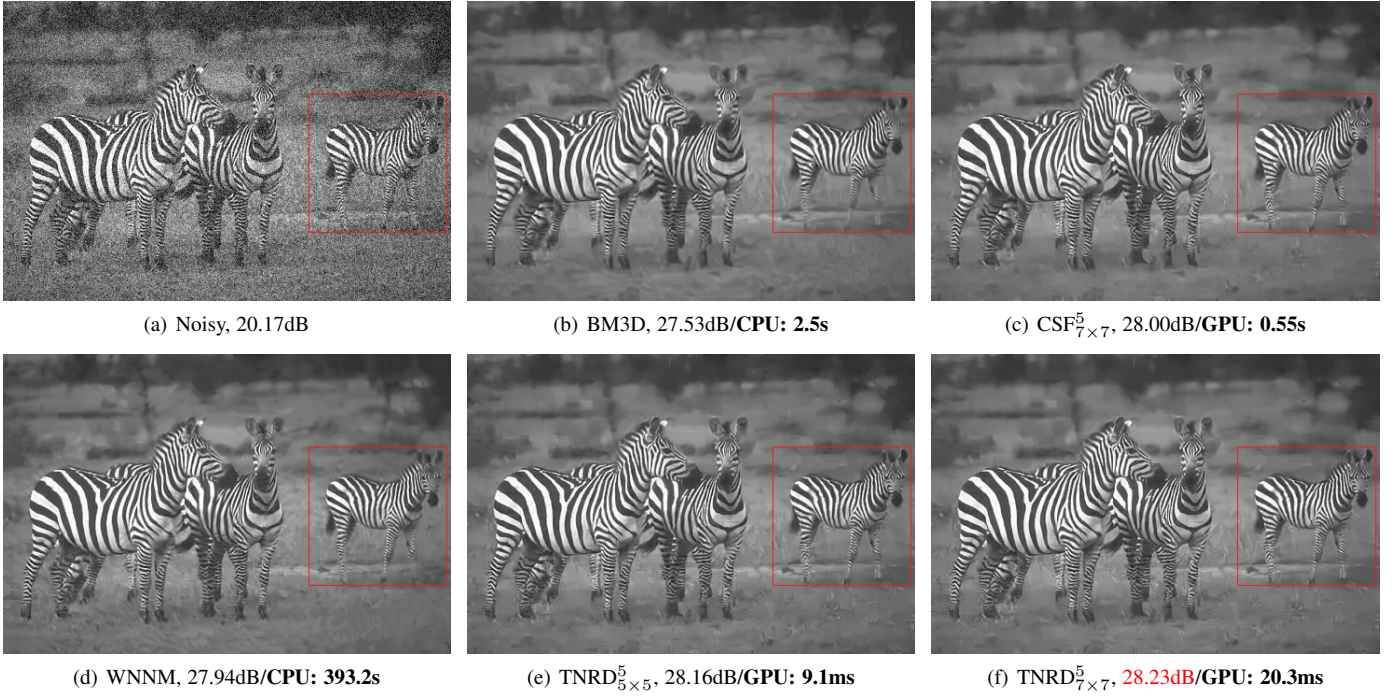


Fig. 11. Denoising results on a test image of size  $481 \times 321$  ( $\sigma = 25$ ) by different methods (compared with BM3D [15], WNNM [26] and CSF model [52]), together with the corresponding computation time either on CPU or GPU. Note the differences in the highlighted region.

Set5 images	scale	Bicubic		K-SVD [58]		ANR [54]		SR-CNN [19]		RFL [53]		TNRD <sub>7x7</sub> <sup>5</sup>	
		PSNR	Time	PSNR	Time	PSNR	Time	PSNR	Time	PSNR	Time	PSNR	Time
baby	2	37.07	-	38.25	8.21	38.44	1.39	38.30	0.38	38.39	1.31	<b>38.51</b>	1.52
bird	2	36.81	-	39.93	2.67	40.04	0.44	40.64	0.14	40.99	0.52	<b>41.29</b>	0.59
butterfly	2	27.43	-	30.65	2.14	30.48	0.38	32.20	0.10	32.46	0.41	<b>33.16</b>	0.56
head	2	34.86	-	35.59	2.46	35.66	0.41	35.64	0.13	<b>35.70</b>	0.48	<b>35.71</b>	0.60
woman	2	32.14	-	34.49	2.45	34.55	0.43	34.94	0.13	35.19	0.46	<b>35.50</b>	0.57
<b>average</b>	2	33.66	-	35.78	3.59	35.83	0.61	36.34	0.18	36.55	0.64	<b>36.83</b>	0.77
baby	3	33.91	-	35.08	3.77	35.13	0.79	35.01	0.38	35.04	0.79	<b>35.28</b>	1.52
bird	3	32.58	-	34.57	1.34	34.60	0.27	34.91	0.14	35.15	0.31	<b>36.11</b>	0.59
butterfly	3	24.04	-	25.94	1.08	25.90	0.24	27.58	0.10	27.18	0.25	<b>28.90</b>	0.56
head	3	32.88	-	33.56	1.35	33.63	0.24	33.55	0.13	33.68	0.29	<b>33.78</b>	0.60
woman	3	28.56	-	30.37	1.14	30.33	0.24	30.92	0.13	30.92	0.28	<b>31.77</b>	0.57
<b>average</b>	3	30.39	-	31.90	1.74	31.92	0.35	32.39	0.18	32.39	0.39	<b>33.17</b>	0.77
baby	4	31.78	-	33.06	2.63	33.03	0.59	32.98	0.38	33.05	0.60	<b>33.29</b>	1.52
bird	4	30.18	-	31.71	0.70	31.82	0.18	31.98	0.14	32.14	0.23	<b>32.98</b>	0.59
butterfly	4	22.10	-	23.57	0.54	23.52	0.14	25.07	0.10	24.44	0.19	<b>26.22</b>	0.56
head	4	31.59	-	32.21	0.66	32.27	0.16	32.19	0.13	32.31	0.22	<b>32.57</b>	0.60
woman	4	26.46	-	27.89	0.72	27.80	0.23	28.21	0.13	28.31	0.23	<b>29.17</b>	0.57
<b>average</b>	4	28.42	-	29.69	1.05	29.69	0.26	30.09	0.18	30.05	0.29	<b>30.85</b>	0.77

TABLE 3

PSNR (dB) and run time (s) performance for upscaling factors  $\times 2$ ,  $\times 3$  and  $\times 4$  on the Set5 dataset. All the methods use the same 91 training images as in [54].

non-smooth optimization problem. Motivated by [8], we consider the following variational model based on the FoE image prior

$$\min_u E(u) = \sum_{i=1}^{N_k} \rho_i(k_i * u) + \mathcal{I}_Q(Du), \quad (19)$$

where  $\mathcal{I}_Q$  is a indicator function over the set  $Q$  (quantization constraint set). In JPEG compression, information loss happens in the quantization step, where all possible values in the range  $[d - 0.5, d + 0.5]$  ( $d$  is an integer) are quantized to a single number  $d$ . Given a compressed data, we only know  $d$ . Therefore, all possible values in the interval  $[d - 0.5, d + 0.5]$  define a convex set  $Q$

which is a box constraint. The sparse matrix  $D \in \mathbb{R}^{N \times N}$  denotes the block DCT transform. We refer to [8] for more details.

By setting  $\mathcal{D}(u) = 0$  and  $\mathcal{G}(u) = \mathcal{I}_Q(Du)$ , we obtain the following diffusion process

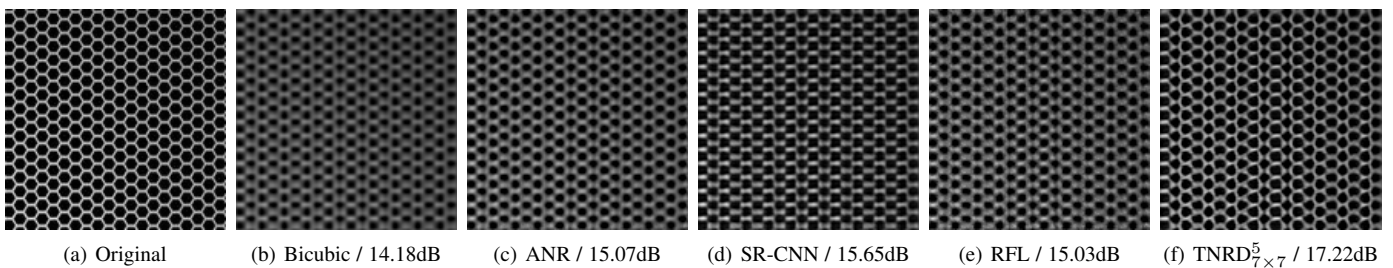
$$u_t = D^\top \text{proj}_Q \left( D \left( u_{t-1} - \sum_{i=1}^{N_k} \bar{k}_i^t * \phi_i^t(k_i^t * u_{t-1}) \right) \right), \quad (20)$$

where  $\text{proj}_Q(\cdot)$  denotes the orthogonal projection onto  $Q$ . More details can be found in the *supplemental material*.

We also trained diffusion models for the JPEG deblocking problem. We followed the test procedure in [33] for performance evaluation. We distorted the images by JPEG block-

Set14 images	Bicubic		K-SVD [58]		ANR [54]		SR-CNN [19]		RFL [53]		TNRD $_{7 \times 7}^5$	
	PSNR	Time	PSNR	Time	PSNR	Time	PSNR	Time	PSNR	Time	PSNR	Time
baboon	23.21	-	23.52	3.54	23.56	0.77	23.60	0.40	23.57	0.75	<b>23.62</b>	1.30
barbara	26.25	-	<b>26.76</b>	6.24	26.69	1.23	26.66	0.70	26.63	1.18	26.25	1.75
bridge	24.40	-	25.02	3.98	25.01	0.80	25.07	0.44	25.11	0.81	<b>25.29</b>	1.19
coastguard	26.55	-	27.15	1.54	27.08	0.36	<b>27.20</b>	0.17	27.16	0.35	27.12	0.65
comic	23.12	-	23.96	1.37	24.04	0.27	24.39	0.15	24.27	0.34	<b>24.67</b>	0.65
face	32.82	-	33.53	1.10	33.62	0.24	33.58	0.13	33.65	0.29	<b>33.82</b>	0.57
flowers	27.23	-	28.43	2.66	28.49	0.57	28.97	0.30	28.86	0.61	<b>29.55</b>	0.90
foreman	31.18	-	33.19	1.54	33.23	0.30	33.35	0.17	33.87	0.36	<b>34.65</b>	0.65
lenna	31.68	-	33.00	3.89	33.08	0.79	33.39	0.44	33.38	0.77	<b>33.77</b>	1.19
man	27.01	-	27.90	3.81	27.92	0.76	28.18	0.44	28.20	0.80	<b>28.52</b>	1.17
monarch	29.43	-	31.10	6.13	31.09	1.13	32.39	0.66	32.10	1.12	<b>33.61</b>	1.66
pepper	32.39	-	34.07	3.84	33.82	0.80	34.35	0.44	34.55	0.82	<b>35.06</b>	1.20
ppt3	23.71	-	25.23	4.53	25.03	1.01	26.02	0.58	25.84	0.98	<b>27.08</b>	1.48
zebra	26.63	-	28.49	3.36	28.43	0.69	28.87	0.38	29.03	0.72	<b>29.40</b>	1.04
<b>average performance</b>	27.54	-	28.67	3.40	28.65	0.69	29.00	0.39	29.02	0.71	<b>29.46</b>	1.10

TABLE 4

Upscaling factor  $\times 3$  performance in terms of PSNR(dB) and runtime (s) per image on the Set14 dataset.Fig. 12. A super resolution example for the “Monarch” image from Set14 with an upscaling factor  $\times 3$ . Note the differences in the highlighted region that our model achieves more clean and sharp image edges. **Best viewed on screen and zoom in.**Fig. 13. A toy experiment on a synthesized image with repeated hexagons for the upscaling factor  $\times 3$ .

ing artifacts. We considered three compression quality settings  $q = 10, 20$  and  $30$  for the JPEG encoder.

We trained three nonlinear diffusion TNRD $_{7 \times 7}$  models for different compression parameter  $q$ . We found that for JPEG deblocking, 4 stages are already enough. Results of the trained models are shown in Table 5, compared with several representative deblocking approaches. We see that our trained TNRD $_{7 \times 7}^4$  outperforms all the competing approaches in terms of PSNR. Concerning the

run time, our model takes about 11.2s to handle an image of size  $1024 \times 1024$  with CPU computation, while the strongest competitor (in terms of run time) - SADCT consumes about 56.5s<sup>11</sup>. Furthermore, our model is extremely fast on GPUs. For the same image size the GPU implementation takes about 0.095s. See the *supplemental material* for JPEG deblocking examples.

11. RTF is slower than SADCT, as it depends on the output of SADCT.

$q$	JPEG decoder	TGV [8]	Dic-SR[9]	SADCT [22]	RTF[33]	TNRD $_{7 \times 7}^4$
10	26.59	26.96	27.15	27.43	27.68	<b>27.85</b>
20	28.77	29.01	29.03	29.46	29.83	<b>30.06</b>
30	30.05	30.25	30.13	30.67	31.14	<b>31.41</b>

TABLE 5

JPEG deblocking results for natural images, reported with average PSNR values.

## 7 DISCUSSION, SUMMARY AND FUTURE WORK

### 7.1 Summary

We have proposed a trainable nonlinear reaction diffusion framework for effective image restoration. Its critical point lies in the additional training of the influence functions. We have trained our models for the problem of Gaussian denoising, single image super resolution and JPEG deblocking. Based on standard test datasets, the trained models result in the best-reported results. We believe that the effectiveness of the trained diffusion models is attributed to the following desired properties of the models

- *Anisotropy*. In the trained filters, we can find rotated derivative filters in different directions, *cf.* Fig 10, which will make the diffusion happen in some special directions.
- *Higher order*. The learned filters contain first, second and higher-order derivative filters, *cf.* Fig 10.
- *Adaptive forward/backward diffusion through the learned nonlinear functions*. Nonlinear functions corresponding to explicit backward diffusion appear in the learned nonlinearity, *cf.* Fig 5.

Meanwhile, the structure of trained models is very simple and well-suited for parallel computation on GPUs. As a consequence, the resulting algorithms are significantly faster than all competing algorithms and hence are also applicable to the restoration of high resolution images.

### 7.2 Discussion

One possible limitation of the proposed TNRD approach is that one has to define the ground truth - the expected output of the diffusion network during training. For image restoration applications in this paper, this is not a problem as we have a clear choice for the ground truth. However, for those applications with ambiguous ground truth, *e.g.*, image structure extraction [57], we will have to make efforts to define the ground truth.

Furthermore, the trained diffusion networks will only perform well in the way they are trained. For example, the trained model based on noise level  $\sigma = 25$  will break for an input image with noise  $\sigma = 50$ , and the trained model for upscaling factor  $\times 3$  will also lead to inferior performance when it is applied to the SISR problem of upscaling factor  $\times 2$ . It is generally hard to train a universal diffusion model to handle all the noise levels or all upscaling factors.

Our approach is to optimize a time-discrete PDE, which is inspired by FoE based model, but we do not aim to minimize a series of FoE based energies. Our model directly learns an optimal trajectory for a certain possibly unknown energy functional, the minimizer of which provides a good estimate of the demanded solution. Probably, such a functional can not be modeled by a single FoE energy, while our learned gradient descent/forward-backward steps provide good approximation to the local gradients of this unknown functional.

### 7.3 Future work

From an application point of view, we think that it will be interesting to consider learned nonlinear reaction diffusion based models also for other image processing tasks such as image inpainting, blind image deconvolution, optical flow. Moreover, since learning the influence functions turned out to be crucial, we believe that learning optimal nonlinearities (*i.e.*, activation functions) in standard CNs could lead to a similar performance increase. There are actually two recent works [2], [29] to investigate a parameterized ReLU function in standard deep convolutional networks, which indeed brings improvements even with little freedom to tune the activation functions. Finally, it will also be interesting to investigate the unconventional penalty functions learned by our approach in usual energy minimization approaches.

## REFERENCES

- [1] S. T. Acton, D. Prasad Mukherjee, J. P. Havlicek, and A. Conrad Bovik. Oriented texture completion by AM-FM reaction-diffusion. *IEEE Trans. Image Processing*, 10(6):885–896, 2001.
- [2] F. Agostinelli, M. Hoffman, P. Sadowski, and P. Baldi. Learning activation functions to improve deep neural networks. *arXiv preprint arXiv:1412.6830*, 2014.
- [3] A. Barbu. Training an active random field for real-time image denoising. *IEEE Trans. Image Processing*, 18(11):2451–2462, 2009.
- [4] C. A. Z. Barcelos, M. Boaventura, and E. C. Silva Jr. A well-balanced flow equation for noise removal and edge detection. *IEEE Trans. Image Processing*, 12(7):751–763, 2003.
- [5] Y. Bengio, P. Lamblin, D. Popovici, and H. Larochelle. Greedy layer-wise training of deep networks. In *Advances in Neural Information Processing Systems*, pages 153–160, 2006.
- [6] Y. Bengio, P. Simard, and P. Frasconi. Learning long-term dependencies with gradient descent is difficult. *IEEE Trans. Neural Networks*, 5(2):157–166, 1994.
- [7] M. Black, G. Sapiro, D. Marimont, and D. Heeger. Robust anisotropic diffusion and sharpening of scalar and vector images. In *Proc. IEEE Int'l Conf. Image Processing*, pages 263–266, 1997.
- [8] K. Bredies and M. Holler. Artifact-free decompression and zooming of JPEG compressed images with total generalized variation. In *Computer Vision, Imaging and Computer Graphics. Theory and Application*, pages 242–258. Springer, 2013.
- [9] H. Chang, M. K. Ng, and T. Zeng. Reducing artifact in JPEG decompression via a learned dictionary. *IEEE Trans. Signal Processing*, 62(3):718–728, 2014.
- [10] Y. Chen, T. Pock, R. Ranftl, and H. Bischof. Revisiting loss-specific training of filter-based MRFs for image restoration. In *Proc. German Conf. Pattern Recognition*, pages 271–281, 2013.
- [11] Y. Chen, R. Ranftl, and T. Pock. Insights into analysis operator learning: From patch-based sparse models to higher order MRFs. *IEEE Trans. Image Processing*, 23(3):1060–1072, 2014.
- [12] Y. Chen, W. Yu, and T. Pock. On learning optimized reaction diffusion processes for effective image restoration. In *Proc. IEEE Int'l Conf. Computer Vision and Pattern Recognition*, 2015.
- [13] G.-H. Cottet and L. Germain. Image processing through reaction combined with nonlinear diffusion. *Mathematics of Computation*, pages 659–673, 1993.
- [14] D. Cozzolino, S. Parrilli, G. Scarpa, G. Poggi, and L. Verdoliva. Fast adaptive nonlocal SAR despeckling. *IEEE Trans. Geosci. Remote Sensing Lett.*, 11(2):524–528, 2014.
- [15] K. Dabov, A. Foi, V. Katkovnik, and K. O. Egiazarian. Image denoising by sparse 3-d transform-domain collaborative filtering. *IEEE Trans. Image Processing*, 16(8):2080–2095, 2007.
- [16] S. Didas, J. Weickert, and B. Burgeth. Stability and local feature enhancement of higher order nonlinear diffusion filtering. In *Proc. German Conf. Pattern Recognition*, pages 451–458, 2005.
- [17] S. Didas, J. Weickert, and B. Burgeth. Properties of higher order nonlinear diffusion filtering. *J. of Mathematical Imaging and Vision*, 35(3):208–226, 2009.
- [18] J. Domke. Generic methods for optimization-based modeling. In *Int'l Conf. Artificial Intelligence and Statistics*, pages 318–326, 2012.
- [19] C. Dong, C. C. Loy, K. He, and X. Tang. Learning a deep convolutional network for image super-resolution. In *Proc. European Conf. Computer Vision*, pages 184–199. Springer, 2014.
- [20] W. Dong, L. Zhang, G. Shi, and X. Li. Nonlocally centralized sparse representation for image restoration. *IEEE Trans. Image Processing*, 22(4):1620–1630, 2013.

- [21] J. Esclarín and L. Alvarez. Image quantization using reaction-diffusion equations. *SIAM J. Applied Mathematics*, 57(1):153–175, 1997.
- [22] A. Foi, V. Katkovnik, and K. Egiazarian. Pointwise shape-adaptive DCT for high-quality denoising and deblurring of grayscale and color images. *IEEE Trans. Image Processing*, 16(5):1395–1411, 2007.
- [23] G. Gilboa, N. Sochen, and Y. Y. Zeevi. Forward-and-backward diffusion processes for adaptive image enhancement and denoising. *IEEE Trans. Image Processing*, 11(7):689–703, 2002.
- [24] R. Giryes and M. Elad. Sparsity-based Poisson denoising with dictionary learning. *IEEE Trans. Image Processing*, 23(12):5057–5069, 2014.
- [25] A. Graves and J. Schmidhuber. Offline handwriting recognition with multidimensional recurrent neural networks. In *Advances in Neural Information Processing Systems*, pages 545–552, 2009.
- [26] S. Gu, L. Zhang, W. Zuo, and X. Feng. Weighted nuclear norm minimization with application to image denoising. In *Proc. IEEE Int'l Conf. Computer Vision and Pattern Recognition*, 2014.
- [27] P. Guidotti and K. Longo. Two enhanced fourth order diffusion models for image denoising. *J. of Mathematical Imaging and Vision*, 40(2):188–198, 2011.
- [28] M. R. Hajiaboli. An anisotropic fourth-order diffusion filter for image noise removal. *Int'l J. Computer Vision*, 92(2):177–191, 2011.
- [29] K. He, X. Zhang, S. Ren, and J. Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. *arXiv preprint arXiv:1502.01852*, 2015.
- [30] N. Heess, C. K. I. Williams, and G. E. Hinton. Learning generative texture models with extended Fields-of-Experts. In *British Machine Vision Conference*, pages 1–11, 2009.
- [31] Y. H. Hu and J.-N. Hwang. *Handbook of neural network signal processing*. CRC press, 2010.
- [32] V. Jain and S. Seung. Natural image denoising with convolutional networks. In *Advances in Neural Information Processing Systems*, pages 769–776, 2009.
- [33] J. Jancsary, S. Nowozin, and C. Rother. Loss-specific training of non-parametric image restoration models: A new state of the art. In *Proc. European Conf. Computer Vision*, pages 112–125, 2012.
- [34] K. Krajssek and H. Schar. Diffusion filtering without parameter tuning: Models and inference tools. In *Proc. IEEE Int'l Conf. Computer Vision and Pattern Recognition*, pages 2536–2543, 2010.
- [35] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- [36] A. Levin and B. Nadler. Natural image denoising: Optimality and inherent bounds. In *Proc. IEEE Int'l Conf. Computer Vision and Pattern Recognition*, pages 2833–2840, 2011.
- [37] P.-L. Lions and B. Mercier. Splitting algorithms for the sum of two nonlinear operators. *SIAM J. Numerical Analysis*, 16(6):964–979, 1979.
- [38] D. C. Liu and J. Nocedal. On the limited memory BFGS method for large scale optimization. *Mathematical Programming*, 45(1):503–528, 1989.
- [39] R. Liu, Z. Lin, W. Zhang, and Z. Su. Learning PDEs for image restoration via optimal control. In *Proc. European Conf. Computer Vision*, pages 115–128, 2010.
- [40] J. Mairal, F. Bach, J. Ponce, G. Sapiro, and A. Zisserman. Non-local sparse models for image restoration. In *Proc. IEEE Int'l Conf. Computer Vision*, pages 2272–2279, 2009.
- [41] P. Milanfar. A tour of modern image filtering: new insights and methods, both practical and theoretical. *IEEE Signal Processing Magazine*, 30(1):106–128, 2013.
- [42] V. Nair and G. E. Hinton. Rectified linear units improve restricted boltzmann machines. In *Proc. IEEE Int'l Conf. Machine Learning*, pages 807–814, 2010.
- [43] K. Niklas Nordström. Biased anisotropic diffusion: a unified regularization and diffusion approach to edge detection. *Image and Vision Computing*, 8(4):318–327, 1990.
- [44] P. Ochs, Y. Chen, T. Brox, and T. Pock. iPiano: Inertial Proximal Algorithm for Nonconvex Optimization. *SIAM J. Imaging Sciences*, 7(2):1388–1419, 2014.
- [45] P. Perona and J. Malik. Scale-space and edge detection using anisotropic diffusion. *IEEE Trans. Pattern Analysis and Machine Intelligence*, 12(7):629–639, 1990.
- [46] A. Rajwade, A. Rangarajan, and A. Banerjee. Image denoising using the higher order singular value decomposition. *IEEE Trans. Pattern Analysis and Machine Intelligence*, 35(4):849–862, 2013.
- [47] Y. Romano, M. Protter, and M. Elad. Single image interpolation via adaptive nonlocal sparsity-based modeling. *IEEE Trans. Image Processing*, 23(7):3085–3098, 2014.
- [48] S. Roth and M. J. Black. Fields of Experts. *Int'l J. Computer Vision*, 82(2):205–229, 2009.
- [49] L. I. Rudin, S. Osher, and E. Fatemi. Nonlinear total variation based noise removal algorithms. *Physica D: Nonlinear Phenomena*, 60(1):259–268, 1992.
- [50] H. Schar, M. J. Black, and H. W. Haussecker. Image statistics and anisotropic diffusion. In *Proc. IEEE Int'l Conf. Computer Vision*, pages 840–847, 2003.
- [51] O. Scherzer and J. Weickert. Relations between regularization and diffusion filtering. *J. of Mathematical Imaging and Vision*, 12(1):43–63, 2000.
- [52] U. Schmidt and S. Roth. Shrinkage fields for effective image restoration. In *Proc. IEEE Int'l Conf. Computer Vision and Pattern Recognition*, 2014.
- [53] S. Schuler, C. Leistner, and H. Bischof. Fast and accurate image upscaling with super-resolution forests. In *Proc. IEEE Int'l Conf. Computer Vision and Pattern Recognition*, 2015.
- [54] R. Timofte, V. De, and L. V. Gool. Anchored neighborhood regression for fast example-based super-resolution. In *Proc. IEEE Int'l Conf. Computer Vision*, pages 1920–1927, 2013.
- [55] J. Weickert. *Anisotropic diffusion in image processing*. Teubner Stuttgart, 2008.
- [56] M. Welk, G. Gilboa, and J. Weickert. Theoretical foundations for discrete forward-and-backward diffusion filtering. In *Proc. Int'l Conf. Scale Space and Variational Methods in Computer Vision*, pages 527–538, 2009.
- [57] L. Xu, Q. Yan, Y. Xia, and J. Jia. Structure extraction from texture via relative total variation. *ACM Trans. Graphics*, 31(6):139, 2012.
- [58] R. Zeyde, M. Elad, and M. Protter. On single image scale-up using sparse-representations. In *Curves and Surfaces*, pages 711–730. Springer, 2012.
- [59] S. C. Zhu and D. Mumford. Prior learning and Gibbs reaction-diffusion. *IEEE Trans. Pattern Analysis and Machine Intelligence*, 19(11):1236–1250, 1997.
- [60] D. Zoran and Y. Weiss. From learning models of natural image patches to whole image restoration. In *Proc. IEEE Int'l Conf. Computer Vision*, pages 479–486, 2011.



lems and convex optimization.

**Yunjin Chen** received a B.Sc. degree in applied physics from Nanjing University of Aeronautics and Astronautics, China, a M.Sc. degree in optical engineering from National University of Defense Technology, China, and a Ph.D degree in computer science from Graz University of Technology, Austria in 2007, 2010 and 2015, respectively. Since July 2015, he serves as a scientific researcher in the military of China. His current research interests are learning image prior model for low-level computer vision prob-



lems and convex optimization.

**Thomas Pock** received a M.Sc and a Ph.D degree in Computer Engineering from Graz University of Technology in 2004 and 2008, respectively. In 2013 he received the START price of the Austrian Science Fund (FWF) and the German Pattern recognition award of the German association for pattern recognition (DAGM). Since June 2014, he is a Professor of Computer Science at Graz University of Technology (AIT Stiftungsprofessur "Mobile Computer Vision") and a principal scientist at the Department of Safety and Security at the Austrian Institute of Technology (AIT). The focus of his research is the development of mathematical models for computer vision and image processing in mobile scenarios as well as the development of efficient algorithms to compute these models.