

Solving Inverse Problems With Deep Neural Networks – Robustness Included?

Martin Genzel, Jan Macdonald, and Maximilian März

Abstract—In the past five years, deep learning methods have become state-of-the-art in solving various inverse problems. Before such approaches can find application in safety-critical fields, a verification of their reliability appears mandatory. Recent works have pointed out instabilities of deep neural networks for several image reconstruction tasks. In analogy to adversarial attacks in classification, it was shown that slight distortions in the input domain may cause severe artifacts. The present article sheds new light on this concern, by conducting an extensive study of the robustness of deep-learning-based algorithms for solving underdetermined inverse problems. This covers compressed sensing with Gaussian measurements as well as image recovery from Fourier and Radon measurements, including a real-world scenario for magnetic resonance imaging (using the NYU-fastMRI dataset). Our main focus is on computing adversarial perturbations of the measurements that maximize the reconstruction error. A distinctive feature of our approach is the quantitative and qualitative comparison with total-variation minimization, which serves as a provably robust reference method. In contrast to previous findings, our results reveal that standard end-to-end network architectures are not only resilient against statistical noise, but also against adversarial perturbations. All considered networks are trained by common deep learning techniques, without sophisticated defense strategies.

Index Terms—Inverse problems, image reconstruction, deep neural networks, adversarial robustness, medical imaging.



1 INTRODUCTION

SIGNAL reconstruction from indirect measurements plays a central role in a variety of applications, including medical imaging [1], communication theory [2], astronomy [3], and geophysics [4]. Such tasks are typically formulated as an inverse problem, which in its prototypical, finite-dimensional form reads as follows:

$$\left\{ \begin{array}{l} \text{Given a linear forward operator } \mathcal{A} \in \mathbb{R}^{m \times N} \\ \text{and corrupted measurements } \mathbf{y} = \mathcal{A}\mathbf{x}_0 + \mathbf{e} \\ \text{with } \|\mathbf{e}\|_2 \leq \eta, \text{ reconstruct the signal } \mathbf{x}_0. \end{array} \right\} \quad (1)$$

The ubiquitous presence of noise makes it indispensable that a reconstruction method has to be *robust* against additive perturbations \mathbf{e} . Furthermore, the measurement process is often costly and potentially harmful. Therefore, the underdetermined regime where $m \ll N$ has gained much attention during the last two decades. This restriction turns (1) into an *ill-posed inverse problem*, which does not possess a unique solution.

Under the additional assumption of sparsity, the methodology of *compressed sensing* has proven that accurate and robust reconstruction from incomplete measurements is still possible [5]. This means that a solution map $\text{Rec}: \mathbb{R}^m \rightarrow \mathbb{R}^N$ for (1) satisfies an error bound of the form

$$\|\mathbf{x}_0 - \text{Rec}(\mathbf{y})\|_2 \leq C \cdot \eta, \quad (2)$$

where $C > 0$ is a small constant. Although state-of-the-art in various real-world applications, the practicability of

the associated algorithms is often limited by computational costs, manual parameter tuning, and a mismatch between sparsity models and data.

Building on the recent success of artificial intelligence in computer vision [6]–[8], there has been a considerable effort to solve the inverse problem (1) by means of *deep learning*, e.g., see [9]–[18] and [19] for a recent survey. This advance is primarily based on fitting an artificial *neural network* (NN) model to a large set of data points in a supervised training procedure. It is fair to say that such data-driven approaches can significantly outperform classical methods in terms of reconstruction accuracy and speed. On the other hand, one may argue that the underlying mechanisms of NNs remain largely unclear [20]. Hence, in the absence of theoretical guarantees of the form (2), an empirical verification of their accuracy and robustness against measurement noise is crucial.

While a number of works report a remarkable resilience against noise [17], [21], [22], several alarming findings indicate that deep-learning-based reconstruction schemes are typically unstable [23]–[26]. In particular, the recent study of Antun et al. [24] suggests that deep learning for inverse problems comes at the cost of instabilities, in the sense that “[...] *certain tiny, almost undetectable perturbations, both in the image and sampling domain, may result in severe artifacts in the reconstruction [...]*”. In machine learning research on classification, such a sensitivity of NNs is a well-established phenomenon. Initiated by Szegedy et al. [27], a substantial body of literature is devoted to *adversarial attacks* (and their defenses), i.e., the computation of a visually imperceptible change to the input that fools the NN. Typically, an “attacker” exploits gradient-based information in order to cross the discontinuous decision boundary of a classifier. This can be a serious issue for sensitive applications where wrong

- M. Genzel is with the Mathematical Institute of Utrecht University, Utrecht, Netherlands.
- J. Macdonald and M. März are with the Institute of Mathematics of Technical University of Berlin, Berlin, Germany.

All authors contributed equally. Correspondence may be addressed to maerz@math.tu-berlin.de.

predictions impose a security risk—imagine a misclassified stop sign in autonomous driving [28], [29].

Despite these findings, it appears peculiar that solving inverse problems by deep-learning-based schemes might become unstable. Learning a reconstruction algorithm can be seen as a regression task, where measurements are mapped to a high-dimensional signal manifold (e.g., medical images). In contrast, a NN classifier maps to a low-dimensional, discrete output domain, resulting in a “vulnerable” decision boundary. Moreover, it is well known that robust and accurate algorithms exist for many inverse problems. Since these are often used as templates for NN architectures, it seems surprising that the latter should suffer from severe instabilities. Clearly, the robustness against noise is quintessential for an application of deep learning in practice, especially in sensitive fields such as biomedical imaging. Therefore, we believe that a profound study of this topic is indispensable.

1.1 Contributions

This article is dedicated to a comprehensive numerical study of the robustness of NN-based methods for solving underdetermined inverse problems. The primary objective of our experiments is to analyze how much the reconstruction error grows with the noise level η . We investigate this relationship in terms of statistical and adversarial noise: the former means that measurement noise is drawn from an appropriate probability distribution, while the latter explores worst-case perturbations that maximize the reconstruction error for fixed η . Similar to adversarial attacks in classification, computing worst-case noise is based on a non-convex formulation that is addressed by automatic differentiation and a gradient descent scheme. In the absence of an empirical certificate of robustness, a central and distinctive component of our analysis is the systematic comparison with a classical benchmark method with provable guarantees, namely total-variation (TV) minimization. In this case, evaluating the gradient is non-trivial and carried out by unrolling the underlying optimization problem.

Our experiments consider several prototypical inverse problems as use cases. This includes classical compressed sensing with Gaussian measurements as well as the reconstruction of phantom images from Radon and Fourier measurements. Furthermore, a real-world scenario for magnetic resonance imaging (MRI) is investigated, based on the NYU-fastMRI dataset [30], [31]. We examine a representative selection of learned reconstruction architectures, reaching from simple post-processing NNs to iterative schemes. In total, this work presents a robustness analysis of more than 25 NNs, each of them trained in-house with publicly available code.¹

Our main findings may be summarized as follows:

(i) In every considered scenario, we find deep-learning-based methods that are at least as robust as TV minimization with respect to adversarial noise. This does not require sophisticated architectures or defense strategies.

(ii) All trained NNs are remarkably robust against statistical noise. Although TV minimization may yield exact

recovery for noiseless measurements, it is still outperformed by learned methods in mid- to high-noise regimes.

(iii) The reconstruction performance is affected by the underlying NN architecture. For instance, promoting data consistency in iterative schemes may improve both accuracy and robustness.

(iv) One should not commit the “inverse crime” of training a NN with *noiseless* data, which may cause an unstable behavior for higher noise levels. We demonstrate that simply adding white Gaussian noise to the training measurements is an effective remedy—a regularization technique that is commonly known as *jittering* in machine learning research.

Apart from these observations, our work is, to the best of our knowledge, the first to empirically characterize the performance gap between adversarial and statistical noise in the context of (1). In particular, this gap is not exclusive to deep-learning-based schemes but also appears for classical methods such as TV minimization. Our central conclusion is:

The existence of adversarial examples in classification tasks does not naturally carry over to NN-based solvers for inverse problems. Such reconstruction schemes may not only supersede classical methods as state-of-the-art, but they can also exhibit a similar degree of robustness.

Since our study as it is has required massive computational resources (>2 years of GPU computation time), some aspects have to remain unexplored, see Section 6 for a discussion. In particular, given the sheer number of NN architectures, we explicitly do *not* claim that every deep-learning-based method is stable (cf. Section 5.1). Nevertheless, our findings suggest that fairly standard workflows allow for surprisingly robust reconstruction schemes. This offers an alternative and novel perspective on the reliability of deep learning strategies in inverse problems. Therefore, we believe that the present work takes an important step towards their safe use in practice.

1.2 Organization of This Article

Section 2 is devoted to relevant previous works, followed by a conceptual overview of our approach in Section 3. The latter introduces all considered reconstruction methods, the associated NN architectures as well as our attack strategy to analyze their adversarial robustness. The main results are then presented in Section 4, complemented by several additional experiments in Section 5. We conclude with a general discussion of our findings in Section 6.

2 RELATED WORK

Initiated by Szegedy et al. [27], the vulnerability of deep NNs to adversarial examples has been the subject of more than 2500 publications [33]. We refer to [34], [35] for recent surveys of the field and further references. The vast majority of existing articles is concerned with classification and related tasks, such as image segmentation [36]. On the other hand, only few works have explicitly addressed the adversarial robustness of learned solvers for inverse problems.

1. Our Python implementation, based on the *PyTorch* package [32], can be found under <https://github.com/jmaces/robust-nets>.

To the best of our knowledge, Huang et al. [23] have made the first effort to transfer adversarial attacks to NN-based reconstruction methods. They demonstrate that a distortion of the network’s input may result in the loss of small image features. However, their initial findings are restricted to the specific problem of limited angle computed tomography, where the robust recovery of certain parts of the image is provably impossible [37]. Moreover, the proposed perturbation model is non-standard and does not correspond to noise in the measurements.

More recently, the topic was brought to attention by the inspiring article of Antun et al. [24]. Their numerical experiments show instabilities of existing deep NNs with respect to adversarial noise, out-of-distribution features, and changes in the number of measurements. An important difference to our work is that adversarial noise is only computed for learned schemes. We believe that a comparative “attack” of a classical benchmark method is crucial for a fair assessment of robustness. Furthermore, the results of [24] are reported qualitatively by visualizing reconstructed images, as it is common in adversarial machine learning. We argue that the mathematical setup of the inverse problem (1) calls for a quantitative error analysis that is in line with the bound of (2). Finally, the training stage of the networks in [24] does not seem to account for noise, which we have identified as a potential source of instability, see Section 5.1. Note that our study also analyzes the FBPCvNet architecture [13], a relative of AUTOMAP [17], and an iterative scheme similar to DeepMRI [38]. Nevertheless, a one-to-one comparison to [24] is subtle due to task-specific architectures and data processing. A follow-up work of [24] presents a theoretical characterization of instabilities in terms of the kernel of the forward operator [25]. Our results provide empirical evidence that the considered deep-learning-based schemes could be kernel aware (cf. Section 5.3).

As a countermeasure to the outcome of [24], Raj et al. [26] suggest a sophisticated defense strategy resulting in robust networks. This work also addresses shortcomings of the attack strategy in [24], see Section 3.4 for details. Finally, in line with our findings, Kobler et al. [39] propose the data-driven *total deep variation* regularizer and demonstrate its adversarial robustness for image denoising.

3 METHODS AND PRELIMINARIES

In this section, we briefly introduce the considered reconstruction schemes for solving the inverse problem (1). This includes a representative selection of NN-based methods and total-variation minimization as a classical benchmark. Furthermore, our attack strategy to analyze their adversarial robustness is presented.

3.1 Neural Network Architectures

In the past five years, numerous deep-learning-based approaches for solving inverse problems have been developed; see [19], [40] for overviews. The present work focuses on a selection of widely used *end-to-end network schemes* that define an explicit reconstruction map from \mathbb{R}^m to \mathbb{R}^N , see also Fig. 1.

The first considered method is a *post-processing network*:

$$\text{UNet}: \mathbb{R}^m \rightarrow \mathbb{R}^N, \mathbf{y} \mapsto [\mathcal{U} \circ \mathcal{A}^\dagger](\mathbf{y}).$$

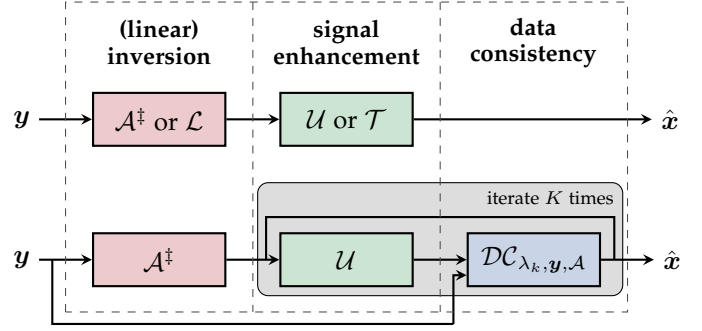


Fig. 1. Schematic network reconstruction pipelines of UNet, TiraFL (top), and ItNet (bottom).

It employs the U-Net architecture $\mathcal{U}: \mathbb{R}^N \rightarrow \mathbb{R}^N$ [41] as a residual network [42] to enhance an initial, model-based reconstruction $\mathcal{A}^\dagger(\mathbf{y})$. Here, $\mathcal{A}^\dagger: \mathbb{R}^m \rightarrow \mathbb{R}^N$ is an approximate inversion of the forward operator \mathcal{A} , e.g., the filtered back-projection for Radon measurements. Despite its simplicity, it has been demonstrated in [13] that UNet is an effective solution method for (1); see also [12], [15], [21], [43], [44] for related approaches.

Our second reconstruction scheme is a *fully-learned network*:

$$\text{TiraFL}: \mathbb{R}^m \rightarrow \mathbb{R}^N, \mathbf{y} \mapsto [\mathcal{T} \circ \mathcal{L}](\mathbf{y}),$$

which is closely related to UNet, but differs in two aspects: It is based on the Tiramisu architecture $\mathcal{T}: \mathbb{R}^N \rightarrow \mathbb{R}^N$ [45] as a residual network, which can be seen as a refinement of the U-Net. While \mathcal{T} shares the same multi-level structure, it is built from fully-convolutional dense-blocks [46] instead of standard convolutional blocks. More importantly, the fixed inversion \mathcal{A}^\dagger is replaced by a learnable linear layer $\mathcal{L} \in \mathbb{R}^{N \times m}$, so that TiraFL does not contain fixed model-based components anymore. The approach of TiraFL is similar to [17], [47], which makes use of a fully-learned reconstruction map for MRI. For the sake of completeness, we have also conducted experiments for Tira, a Tiramisu-based post-processing network, as well as for UNetFL, a U-Net-based fully-learned network, see Section S1–S3 in the supplementary material for results.

Finally, we also analyze an *iterative network*:

$$\text{ItNet}: \mathbb{R}^m \rightarrow \mathbb{R}^N, \mathbf{y} \mapsto \left[\left(\bigcirc_{k=1}^K [\mathcal{DC}_{\lambda_k, \mathbf{y}, \mathcal{A}} \circ \mathcal{U}] \right) \circ \mathcal{A}^\dagger \right](\mathbf{y})$$

where

$$\mathcal{DC}_{\lambda_k, \mathbf{y}, \mathcal{A}}: \mathbb{R}^N \rightarrow \mathbb{R}^N, \mathbf{x} \mapsto \mathbf{x} - \lambda_k \cdot \mathcal{A}^* (\mathcal{A}\mathbf{x} - \mathbf{y}).$$

The scalar parameters λ_k are learnable and \mathcal{A}^* denotes the adjoint of \mathcal{A} . Mathematically, $\mathcal{DC}_{\lambda_k, \mathbf{y}, \mathcal{A}}$ performs a gradient step on the loss $\mathbf{x} \mapsto \frac{\lambda_k}{2} \|\mathcal{A}\mathbf{x} - \mathbf{y}\|_2^2$, promoting *data consistent* solutions. Therefore, the alternating cascade of ItNet can be seen as a proximal gradient descent scheme, where the proximal operator is replaced by a trainable enhancement network. Here, the U-Net architecture is used again, due to its omnipresence in image-to-image processing tasks. Unrolled methods in the spirit of ItNet are frequently used to solve inverse problems, e.g., see [9], [10], [14], [38], [48]–[51].

3.2 Neural Network Training

The learnable parameters of the networks are trained from sample data pairs $\{(\mathbf{y}^i = \mathcal{A}\mathbf{x}_0^i + \mathbf{e}^i, \mathbf{x}_0^i)\}_{i=1}^M$ by minimizing an empirical loss function. Depending on the use case, the signals \mathbf{x}_0^i are either drawn from a fixed publicly available training dataset or according to a synthetic probability distribution. If $\text{Net}[\boldsymbol{\theta}]: \mathbb{R}^m \rightarrow \mathbb{R}^N$ denotes a reconstruction network with all learnable parameters collected in $\boldsymbol{\theta}$, then the training amounts to (approximately) solving

$$\min_{\boldsymbol{\theta}} \sum_{i=1}^M \ell(\text{Net}[\boldsymbol{\theta}](\mathbf{y}^i), \mathbf{x}_0^i) + \mu \cdot \|\boldsymbol{\theta}\|_2^2 \quad (3)$$

for some cost function $\ell: \mathbb{R}^N \times \mathbb{R}^N \rightarrow \mathbb{R}_{\geq 0}$, which is the squared distance unless stated otherwise. Overfitting is addressed by ℓ_2 -regularization (weight decay) with a hyperparameter $\mu \geq 0$. In order to solve (3), we utilize mini-batch stochastic gradient descent and the Adam optimizer [52]. We found that larger mini-batches were beneficial for the training performance during later epochs, which is achieved by gradient accumulation.

Due to the ubiquitous presence of noise in inverse problems, it is natural to account for it in the training data. In many applications, measurement noise is modeled as an independent random variable, for instance, following a Gaussian distribution. Therefore, the perturbation \mathbf{e}^i is treated as statistical noise during the training phase, i.e., a fresh realization is randomly drawn in each epoch. This technique is well known as *jittering* in machine learning research, where it is primarily used to avoid overfitting [53]–[55]; see also [56]. In Section 5.1, we relate jittering to the phenomenon of inverse crimes and demonstrate its importance for the robustness of learned reconstruction schemes. Due to varying noise levels in the evaluation of our models, we design \mathbf{e}^i as a centered Gaussian vector with random variance, such that its expected norm $\mathbb{E}\|\mathbf{e}^i\|_2$ is distributed uniformly in a range $[0, \tilde{\eta}]$ for some $\tilde{\eta} \geq 0$.

3.3 Total-Variation Minimization

Dating back to the seminal work of Rudin et al. [57], *total-variation (TV) minimization* has become a standard tool for solving signal and image reconstruction tasks [58], [59]. We apply it to the problem (1) in the following form:

$$\begin{aligned} \text{TV}[\eta]: \mathbb{R}^m &\rightarrow \mathbb{R}^N, \\ \mathbf{y} &\mapsto \underset{\mathbf{x} \in \mathbb{R}^N}{\text{argmin}} \|\nabla \mathbf{x}\|_1 \quad \text{s.t.} \quad \|\mathcal{A}\mathbf{x} - \mathbf{y}\|_2 \leq \eta, \end{aligned} \quad (4)$$

where ∇ denotes a discrete gradient operator. Crucial to the above optimization problem is the use of the ℓ_1 -norm, which is known to promote gradient-sparse solutions. Indeed, under suitable assumptions on \mathcal{A} , compressed sensing theory suggests an error bound of the form (2) for a gradient-sparse signal \mathbf{x}_0 and $\text{Rec} = \text{TV}[\eta]$, e.g., see [60]–[63]. In other words, TV minimization is provably robust with a near-optimal dependence on η . This particularly justifies its use as a reference method, allowing us to empirically characterize the robustness of learned reconstruction schemes.

In our numerical simulations, the problem of (4) is solved by the *alternating direction method of multipliers (ADMM)* [64], [65]. For 1D signals, $\nabla \in \mathbb{R}^{2N \times N}$ is chosen as a forward finite difference operator with Neumann boundary conditions, extended by a constant row vector to capture the mean of

the signal. For image signals, $\nabla \in \mathbb{R}^{2N \times N}$ corresponds to a forward finite difference operator with periodic boundary conditions. Finally, we emphasize that $\text{TV}[\eta]$ is explicitly adapted to the amount of perturbation of the measurements.

3.4 Adversarial Noise

In the setup of (1), adversarial noise for a given reconstruction method $\text{Rec}: \mathbb{R}^m \rightarrow \mathbb{R}^N$ can be computed by solving an optimization problem: for a fixed signal $\mathbf{x}_0 \in \mathbb{R}^N$ and noise level $\eta \geq 0$, find an additive perturbation $\mathbf{e}_{\text{adv}} \in \mathbb{R}^m$ of the noiseless measurements $\mathbf{y}_0 = \mathcal{A}\mathbf{x}_0$ that maximizes the reconstruction error, i.e.,

$$\mathbf{e}_{\text{adv}} = \underset{\mathbf{e} \in \mathbb{R}^m}{\text{argmax}} \|\text{Rec}(\mathbf{y}_0 + \mathbf{e}) - \mathbf{x}_0\|_2 \quad \text{s.t.} \quad \|\mathbf{e}\|_2 \leq \eta. \quad (5)$$

Such an attack strategy is a straightforward adaption of a common approach in adversarial machine learning [34]. In contrast to [24], we consider a constrained optimization problem that avoids shortcomings of an unconstrained formulation; in particular, this allows for precise control over the noise level. Moreover, (5) explores a natural perturbation model, operating directly in the measurement domain, cf. the discussion in [26].

In order to solve the problem (5), we use the projected gradient descent algorithm in conjunction with the Adam optimizer, which was found to be most effective (cf. [66]). The non-convexity of (5) is accounted for by choosing the worst perturbation out of multiple runs with random initialization. Assuming a whitebox model (i.e., Rec is fully accessible), we use PyTorch’s automatic differentiation [32] to compute gradients of the considered NN schemes.

A central aspect of our work is that the above perturbation strategy is also applied to $\text{TV}[\eta]$. This is non-trivial, since the gradient of the implicit map $\mathbf{y} \mapsto \text{TV}[\eta](\mathbf{y})$ has to be computed. The large-scale nature of imaging problems prevents us from using the recent concept of differentiable convex optimization layers [67]. Instead, we rely on unrolling the ADMM scheme for TV minimization, which again enables automatic differentiation. However, a large number of iterations might be required to ensure convergence of ADMM. This leads to numerical difficulties when calculating the gradient of the unrolled algorithm, e.g., memory & time constraints and error accumulation. We address this issue by decreasing the number of ADMM iterations in combination with a pre-initialization of the primal and dual variables.

4 MAIN RESULTS

This section studies the robustness of NN-based solution methods for three different instances of the inverse problem (1). The goal of our experiments is to assess the loss of reconstruction accuracy caused by noise. To that end, we rely on two types of visualization:

- *Noise-to-error curves* are generated by plotting the relative noise level $\eta/\|\mathcal{A}\mathbf{x}_0\|_2$ against the relative reconstruction error $\|\mathbf{x}_0 - \text{Rec}(\mathcal{A}\mathbf{x}_0 + \mathbf{e})\|_2/\|\mathbf{x}_0\|_2$.
- *Individual reconstruction results* are shown for different relative noise levels and a randomly selected signal from the test set.

In both cases, the perturbation vector e is either of *statistical* or *adversarial* type. The former means that e is a random vector such that $\mathbb{E}[\|e\|_2^2] = \eta^2$, whereas the latter is found by (5). While noise-to-error curves are of quantitative nature, individual reconstructions facilitate a qualitative judgment of robustness. Note that the sensitivity to noise is different in each considered scenario. Therefore, we have selected the maximal level of adversarial noise such that the benchmark of TV minimization does not yield a (subjectively) acceptable performance anymore. A specification of all empirically selected hyper-parameters can be found in the supplementary material (see Table S9–S11).

4.1 Case Study A: Compressed Sensing With Gaussian Measurements

Our first study is devoted to sparse recovery of 1D signals from Gaussian measurements. This means that the entries of the forward operator \mathcal{A} in (1) are independent Gaussian random variables with zero mean and variance $1/m$. Although a toy problem, this setup is a folkloric benchmark in the field of compressed sensing (CS) theory [5].

We consider two different scenarios based on (approximately) gradient-sparse signals; note that such a model is canonical for TV minimization and compatible with the local connectivity of our convolutional NN schemes.

Scenario A1: We draw x_0 from a synthetic distribution of *piecewise constant signals* with zero boundaries and well-controlled random jumps, see Fig. 3 for an example. In this

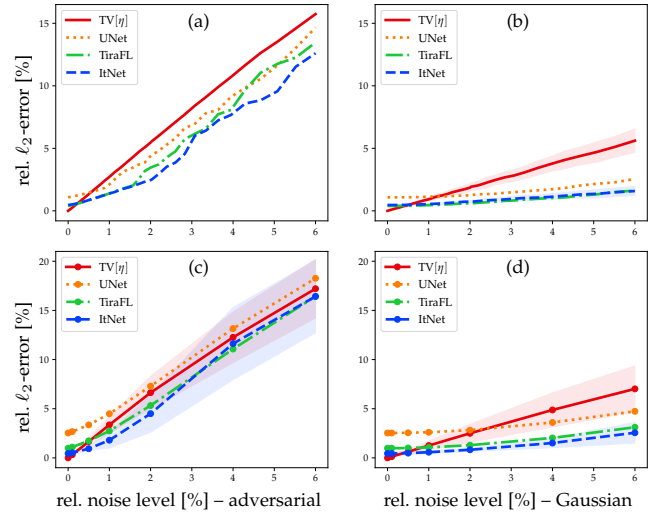


Fig. 2. **Scenario A1 – CS with 1D signals.** (a) shows the adversarial noise-to-error curve for the randomly selected signal of Fig. 3. (b) shows the corresponding Gaussian noise-to-error curve, where the mean and standard deviation are computed over 200 draws of e . (c) and (d) display the respective curves averaged over 50 signals from the test set. For the sake of clarity, we have omitted the standard deviations for UNet and TiraFL, which behave similarly.

scenario, we choose $m = 100$, $N = 256$, and use $M = 200k$ training samples.

Scenario A2: We sample $x_0 \in [0, 1]^{28 \times 28}$ from the widely used *MNIST database* [68] with $M = 60k$ training

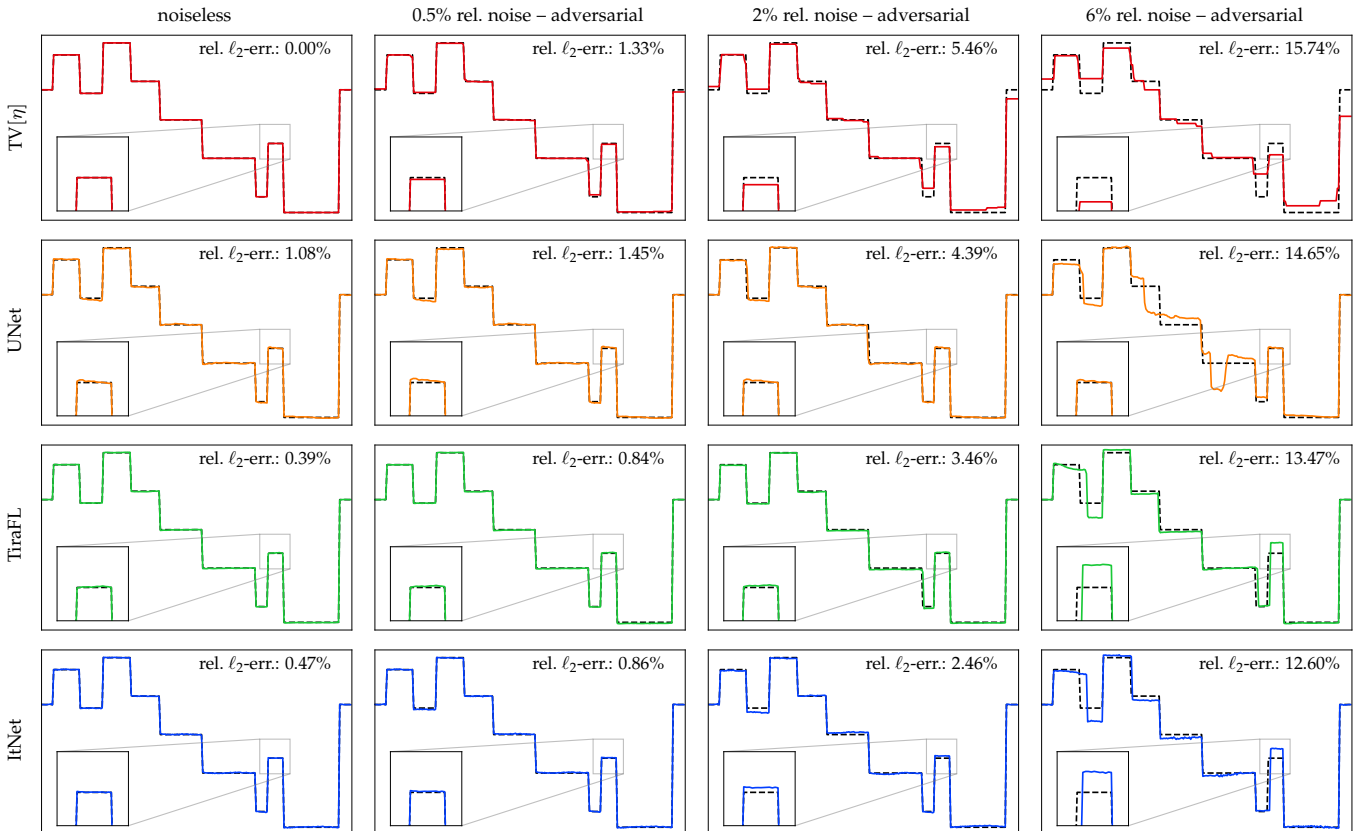


Fig. 3. **Scenario A1 – CS with 1D signals.** Individual reconstructions of a randomly selected signal from the test set for different levels of adversarial noise. The ground truth signal is visualized by a dashed line.

images of handwritten digits. In the context of (1), the images are treated as 1D signals of dimension $N = 28^2 = 784$. For visual purposes, all reconstructions are displayed as images, see Fig. 5. The number of Gaussian measurements is $m = 300$.

In both scenarios, we chose the model-based, linear inversion layer of the networks as a generalized Tikhonov matrix, i.e., $\mathcal{A}^\dagger = (\mathcal{A}^T \mathcal{A} + \alpha \cdot \nabla^T \nabla)^{-1} \mathcal{A}^T \in \mathbb{R}^{N \times m}$ with the empirically chosen regularization parameter $\alpha = 0.02$. We were not able to train the NNs to a comparable reconstruction accuracy with other natural choices, such as $\mathcal{A}^\dagger = \mathcal{A}^T$. The above matrix is also used to initialize the inversion layer $\mathcal{L} \in \mathbb{R}^{N \times m}$ of the fully-learned schemes.

Fig. 2 shows the noise-to-error curves for *Scenario A1 (CS with 1D signals)*; see also Table S1 and S2. The associated individual reconstructions for adversarial noise are displayed in Fig. 3; see Fig. S2 for the corresponding results with Gaussian noise. Fig. S1 supplements the simulation of Fig. 2(b) and (d) by two additional types of random noise, drawn from the uniform and Bernoulli distribution. Both exhibit results that are virtually indistinguishable from the Gaussian case. Fig. 4 shows the noise-to-error curves for *Scenario A2 (CS with MNIST)*; see also Table S3 and S4. The associated individual reconstructions for adversarial noise are displayed in Fig. 5; see Fig. S3 for two additional digits and Fig. S4 for the corresponding results with Gaussian noise.

Conclusions: The above results confirm that the considered NN-based schemes are as least as robust to adversarial perturbations as the benchmark of TV minimization. Although $\text{TV}[\eta]$ is perfectly tuned to each noise level η , it is clearly outperformed in the case of statistical noise. The gap between statistical and adversarial perturbations is comparable for all methods.

TV minimization is a perfect match for Scenario A1. In particular, exact recovery from noiseless measurements is guaranteed by CS theory [63], [69]. Although this cannot be expected for NN-based solvers, they still come with an overall superior robustness against noise. The situation is even more striking in Scenario A2. Here, TV minimization performs worse, since the signals are only approximately gradient-sparse. In contrast, the NN-based reconstruction schemes adapt well to the simple MNIST database, leading to significantly better outcomes in every regard. Hence, the increase in accuracy by learned methods does not necessarily imply a loss of robustness.

The performance ranking of the considered deep NNs is as one might expect: First, data consistency as encouraged by the ItNet-architecture is beneficial. Furthermore, Table S1–S4 reveal that the Tiramisu architecture is superior to a simple U-Net, and that a learnable inversion layer improves the recovery. The latter observation is not surprising, since Tikhonov regularization is known to work poorly in conjunction with subsampled Gaussian measurements.

4.2 Case Study B: Image Recovery of Phantom Ellipses

Our second set of experiments concerns the recovery of phantom ellipses from Fourier or Radon measurements. These tasks correspond to popular simulation studies for biomedical imaging, e.g., see [13], [18], [49], [70]. We sample

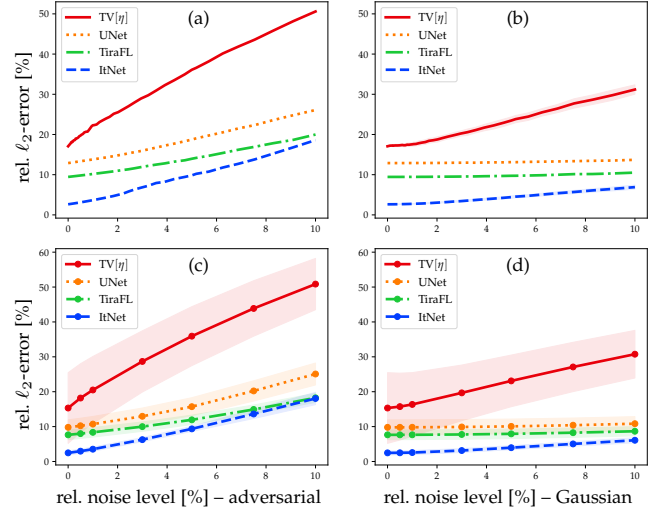


Fig. 4. **Scenario A2 – CS with MNIST.** (a) shows the adversarial noise-to-error curve for the randomly selected digit 3 of Fig. 5. (b) shows the corresponding Gaussian noise-to-error curve, where the mean and standard deviation are computed over 200 draws of e . (c) and (d) display the respective curves averaged over 50 signals from the test set.

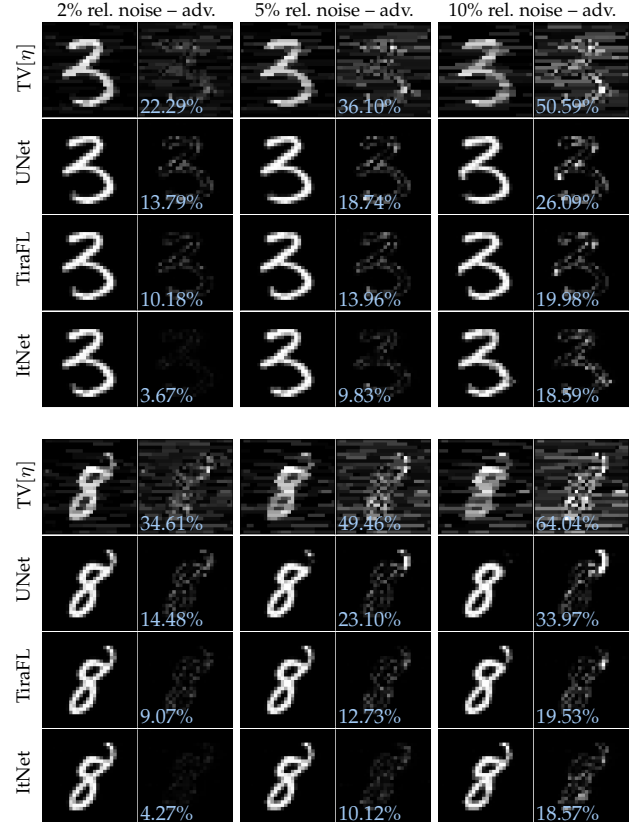


Fig. 5. **Scenario A2 – CS with MNIST.** Individual reconstructions of two randomly selected digits from the test set for different levels of adversarial noise. The reconstructed digits and their error plots (with relative ℓ_2 -error) are displayed in the windows $[0, 1]$ and $[0, 0.6]$, respectively. The horizontal line artifacts in the $\text{TV}[\eta]$ -solutions are due to the fact that the MNIST images are treated as vectorized 1D signals. Remarkably, although relying on 1D convolutional filters, the NN-based reconstructions do not suffer from these artifacts.

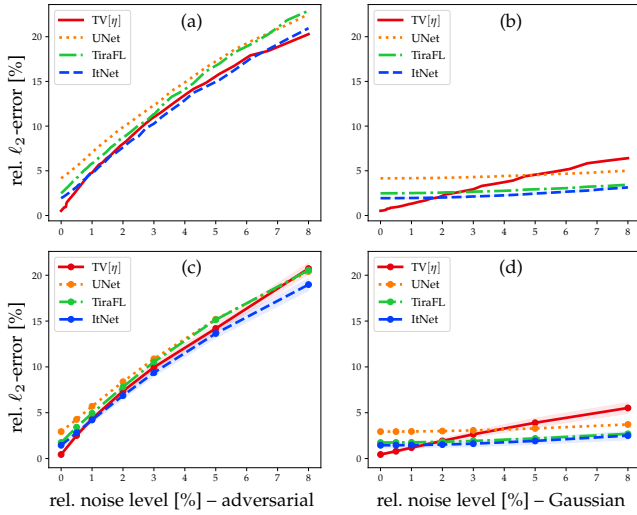


Fig. 6. **Scenario B1 – Fourier meas. with ellipses.** (a) shows the adversarial noise-to-error curve for the randomly selected image of Fig. 7. (b) shows the corresponding Gaussian noise-to-error curve, where the mean and (almost imperceptible) standard deviation are computed over 50 draws of e . (c) and (d) display the respective curves averaged over 50 images from the test set. For the sake of clarity, we have omitted the standard deviations for UNet and TiraFL, which behave similarly.

$x_0 \in [0, 1]^{256 \times 256}$ from a distribution of superimposed random ellipses with mild linear intensity gradients and well-controlled geometric properties, see Fig. 7 for an example. The training is performed on $M = 25k$ images. We consider the following two measurement scenarios for (1), associated with the problems of *compressed sensing MRI* [1] and *low-dose computed tomography (CT)* [13], [71], respectively:

Scenario B1: The forward operator takes the form $\mathcal{A} = P\mathcal{F} \in \mathbb{C}^{m \times N}$, where $\mathcal{F} \in \mathbb{C}^{N \times N}$ is the 2D discrete Fourier transform and $P \in \{0, 1\}^{m \times N}$ is a subsampling operator defined by a golden-angle radial mask with 40 lines ($m = 10941$ and $N = 256^2 = 65536$). Note that the entire data processing is complex-valued, while the actual reconstructions are computed as real-valued magnitude images, as common in MRI. We use the canonical inversion layer $\mathcal{A}^\dagger = \mathcal{A}^* = \mathcal{F}^{-1}P^* \in \mathbb{C}^{N \times m}$.

Scenario B2: The forward operator $\mathcal{A} \in \mathbb{R}^{m \times N}$ is given by a sparse-angle Radon transform with 60 views ($m = 21780$ and $N = 65536$). The non-linear inversion layer $\mathcal{A}^\dagger: \mathbb{R}^m \rightarrow \mathbb{R}^N$ is chosen as the filtered back-projection algorithm (FBP) with a Hann filter.

In contrast to Case Study A, the aforementioned problems are of significantly higher dimensionality, requiring several adaptations. First, fully-learned schemes are difficult to realize, since the size of the inversion layer scales multiplicatively in the image dimensions. In the Fourier case, the number of free parameters can be reduced by enforcing a Kronecker product structure on $\mathcal{L} \in \mathbb{C}^{N \times m}$; this exploits the fact that \mathcal{F} is a tensor product of two 1D Fourier transforms, cf. [47]. Furthermore, due to the non-separability of $\|\nabla(\cdot)\|_1$, the formulation of $\text{TV}[\eta]$ in (4) becomes computationally infeasible for finding adversarial noise. Hence, we solve the unconstrained version of $\text{TV}[\eta]$ instead, i.e., the objective function is changed to $x \mapsto \lambda \cdot \|\nabla x\|_1 + \|\mathcal{A}x - y\|_2^2$. Note that this strategy is theoretically equivalent [5, Appx. B],

but requires an appropriate choice of the regularization parameter $\lambda > 0$. A near-optimal selection with respect to the relative ℓ_2 -error is determined by grid searches over the test set and a densely sampled range of noise levels η .

Fig. 6 shows the noise-to-error curves for *Scenario B1 (Fourier meas. with ellipses)*; see also Table S5 and S6. The associated individual reconstructions for $\text{TV}[\eta]$ and ItNet with adversarial noise are displayed in Fig. 7; see Fig. S5 for the remaining networks and Fig. S6 for the corresponding results with Gaussian noise. In the tables and individual reconstructions, we have also reported the *peak signal-to-noise ratio (PSNR)* and *structural similarity index measure (SSIM)* [72]. In the case of *Scenario B2 (Radon meas. with ellipses)*, we only present individual reconstructions based on $\text{TV}[\eta]$ and UNet; see Fig. 8 for adversarial noise and Fig. S7 for the common Poisson noise model. This restriction is due to the more complicated nature of the Radon transform, and in particular, the need for automatic differentiation. The used implementation [73] requires significantly more computational effort, compared to the fast Fourier transform.

Conclusions: The main findings of Case Study A remain valid: (i) the adversarial robustness of NN-based methods and TV minimization is similar with respect to the ℓ_2 -error; (ii) NNs are more resilient against statistical perturbations in mid- to high-noise regimes (see also the individual reconstructions in Fig. S6 and S7); (iii) there is a clear gap between adversarial and statistical noise that is comparable for model-based and learned schemes.

The individual reconstruction results in Fig. 7 and 8 allow for further insights. First, the effect of adversarial noise for $\text{TV}[\eta]$ manifests itself in the well-known staircasing phenomenon, a considerable loss of resolution as well as point-like artifacts (see the zoomed region in Fig. 7). In contrast, NN-based methods always produce sharp images, with almost imperceptible visual errors up to 3% relative noise in the case of Fourier measurements (1% noise in the case of Radon measurements). For the highest noise level, on the other hand, they exhibit unnatural ellipsoidal artifacts.

At first sight, this observation might indicate a vulnerability to adversarial noise. However, a simple *transferability test* refutes this conclusion (cf. [74]): plugging the perturbed measurements for ItNet into $\text{TV}[\eta]$ leads to the same ellipsoidal artifacts; see Fig. 7 and Fig. S8. Furthermore, Fig. 8 reveals that the corresponding artifacts are already present in the FBP inversion and are not caused by the post-processing network. This shows that the learned solvers do not suffer from undesired instabilities, but the observed artifacts are due to actual features in the corrupted measurements. Interestingly, adversarial perturbations found for $\text{TV}[\eta]$ do not transfer to NN-based methods, see Fig. S8. Overall, the attack strategy of (5) has different qualitative effects on each reconstruction paradigm: while known flaws of TV minimization are amplified, the NNs are perturbed by adding “real” ellipsoidal features to the measurements.

On a final note, we confirm the ranking of architectures as pointed out in Case Study A. Nevertheless, there is no clear superiority of the fully learned schemes as in case of Gaussian measurements, since the inverse Fourier transform appears to be a near-optimal choice of model-based inversion layer.

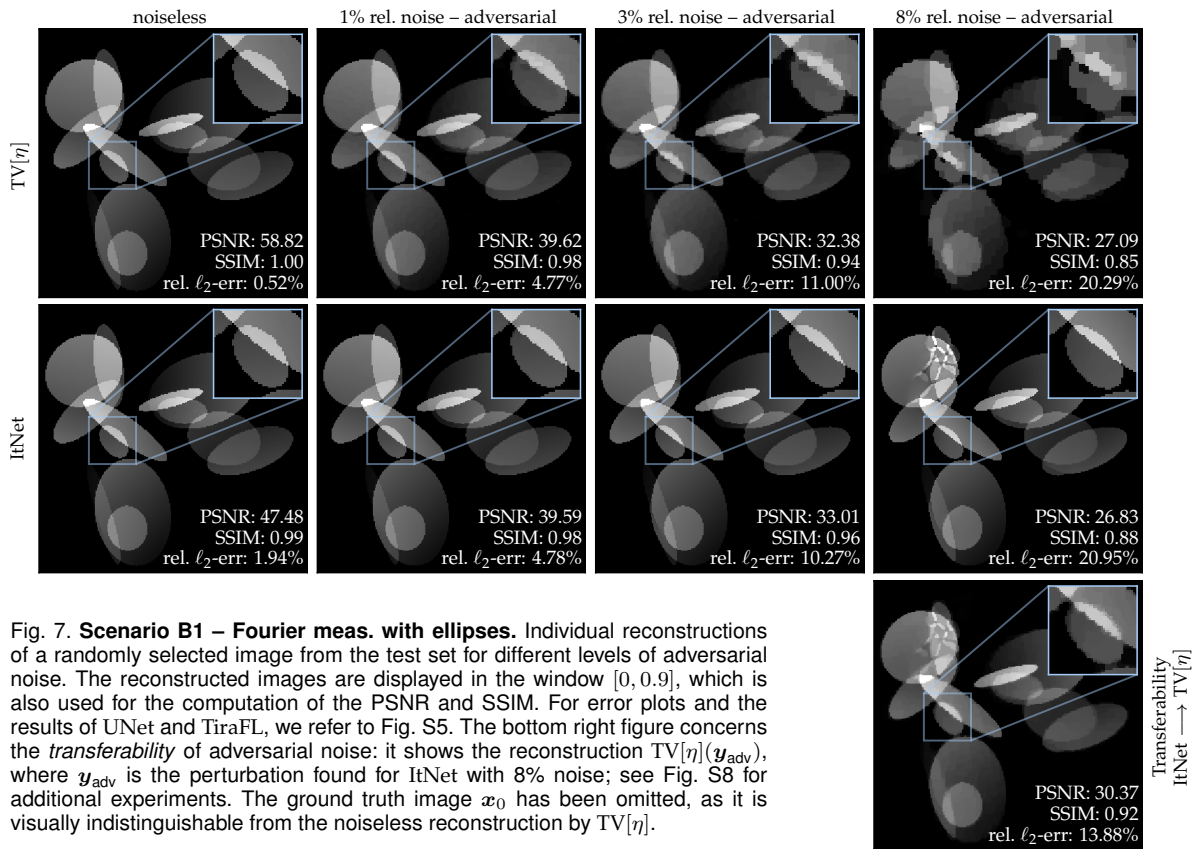


Fig. 7. **Scenario B1 – Fourier meas. with ellipses.** Individual reconstructions of a randomly selected image from the test set for different levels of adversarial noise. The reconstructed images are displayed in the window $[0, 0.9]$, which is also used for the computation of the PSNR and SSIM. For error plots and the results of UNet and TiraFL, we refer to Fig. S5. The bottom right figure concerns the *transferability* of adversarial noise: it shows the reconstruction $TV[\eta](y_{adv})$, where y_{adv} is the perturbation found for ItNet with 8% noise; see Fig. S8 for additional experiments. The ground truth image x_0 has been omitted, as it is visually indistinguishable from the noiseless reconstruction by $TV[\eta]$.

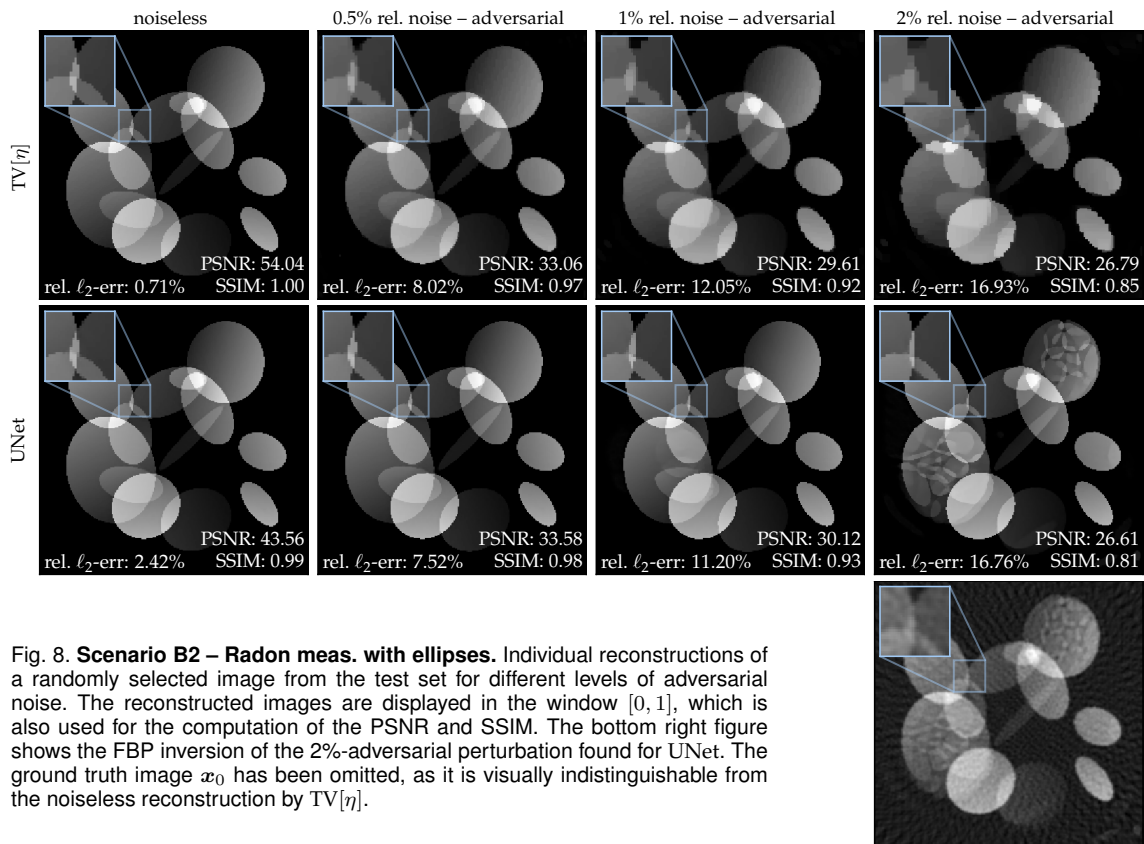


Fig. 8. **Scenario B2 – Radon meas. with ellipses.** Individual reconstructions of a randomly selected image from the test set for different levels of adversarial noise. The reconstructed images are displayed in the window $[0, 1]$, which is also used for the computation of the PSNR and SSIM. The bottom right figure shows the FBP inversion of the 2%-adversarial perturbation found for UNet. The ground truth image x_0 has been omitted, as it is visually indistinguishable from the noiseless reconstruction by $TV[\eta]$.

4.3 Case Study C: MRI on Real-World Data (fastMRI)

The third case study of this article is devoted to a real-world MRI scenario. To this end, we use the publicly available *fastMRI* knee dataset, which consists of 1594 multi-coil diagnostic knee MRI scans.² Our experiments are based on the subset of 796 coronal proton-density weighted scans without fat-suppression, resulting in $M \approx 17\text{k}$ training images. We draw magnitude images $\mathbf{x}_0 \in \mathbb{R}^{320 \times 320}$, obtained from fully-sampled multi-coil³ data, and consider subsampled Fourier measurements as in Scenario B1 with 50 radial lines ($m = 17178$ and $N = 320^2 = 102400$). As before, the data processing is complex-valued, while the actual reconstructions are computed as real-valued magnitude images. The model-based and learned inversion layers are realized as in Scenario B1. As common in the *fastMRI* challenge, we have trained all networks with a cost function based on a combination of the ℓ_1 - and SSIM-distance, see also [75]. TV minimization is solved in the unconstrained formulation, with the regularization parameter determined by a grid search over a subset of the validation set.

Fig. 9 shows the noise-to-error curves; see also Table S7 and S8. The associated individual reconstructions for TV[η] and TiraFL with adversarial noise are displayed in Fig. 10; see Fig. S9 for the remaining networks and Fig. S10 for the corresponding results with Gaussian noise.

Conclusions: Our experimental results show that the main findings of Case Study A and B carry over to real-world data. The noise-to-error curves in Fig. 9 reveal a superior robustness of the learned reconstruction schemes over TV minimization, even for noiseless measurements (cf. Scenario A2). Fig. 10 underpins this observation from a qualitative viewpoint: the model-based prior of TV[η] tends to blur fine details in the reconstructed images—this “oil painting” effect becomes stronger with larger perturbations. In contrast, the NN-based reconstructions always yield high resolution images. Despite adversarial noise, the central image region—which is of main medical interest—remains largely unaffected, whereas tiny vessel structures appear in the outside (fat) region. Such an amplification of existing patterns is comparable to the ellipsoidal artifacts in Case Study B. We emphasize that this phenomenon only occurs for large adversarial perturbations, where the benchmark of TV minimization already suffers from severe distortions. In particular, the performance of the learned

2. Data used in the preparation of this article were obtained from the NYU *fastMRI* Initiative database [30], [31] (<https://fastmri.med.nyu.edu>). As such, NYU *fastMRI* investigators provided data but did not participate in analysis or writing of this article. The primary goal of *fastMRI* is to test whether machine learning can aid in the reconstruction of medical images.

3. Note that our measurement model actually corresponds to the simpler modality of subsampled single-coil MRI. While the *fastMRI* challenge also provides single-coil data, it is based on retrospective masking of *emulated* Fourier measurements. The subsampling is done by omitting k -space lines in the phase-encoding direction, which we found less suitable for our robustness analysis; see Section 5.3 for an experiment with the original setup. Since emulating single-coil measurements is unavoidable, we have decided to sample from the multi-coil magnitude reconstructions in favor of higher image quality. This was found to be particularly important to ensure that TV minimization can serve as a competitive benchmark method, at least for noiseless measurements.

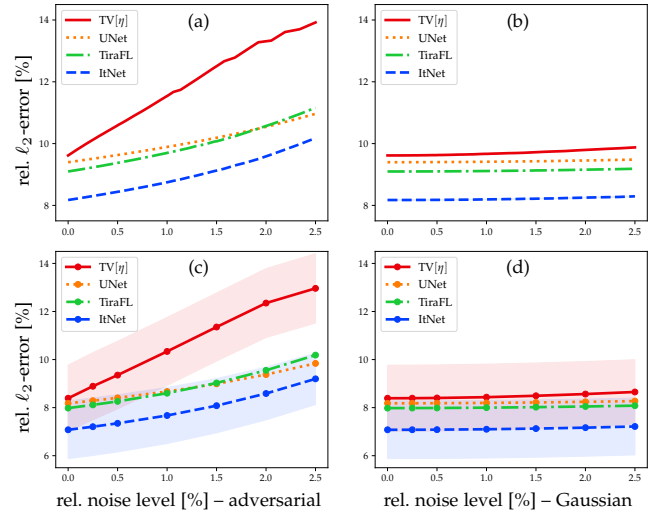


Fig. 9. **Case Study C – fastMRI.** (a) shows the adversarial noise-to-error curve for the randomly selected image of Fig. 10. (b) shows the corresponding Gaussian noise-to-error curve, where the mean and (almost imperceptible) standard deviation are computed over 50 draws of e . (c) and (d) display the respective curves averaged over 30 images from the validation set. For the sake of clarity, we have omitted the standard deviations for UNet and TiraFL, which behave similarly.

methods is not impaired by the same amount of Gaussian noise (see Fig. S10).

5 FURTHER ASPECTS OF ROBUSTNESS

This section presents several additional experiments that allow for further insights into the robustness of learned methods.

5.1 Training Without Noise – An Inverse Crime?

In this section, the importance of *jittering* for the stability of deep-learning-based reconstruction schemes is discussed (see Section 3.2). We have found that this technique can be beneficial for promoting adversarial robustness, in particular, for iterative architectures. The previous claim is verified by an ablation study, comparing two versions of ItNet for Scenario A2, one trained with jittering and the other without. The resulting noise-to-error curves in Fig. 11 reveal that noiseless training data can have drastic consequences. Indeed, the relative recovery error blows up at $\sim 15\%$ adversarial noise if jittering is not used. In a similar experiment, we analyze the adversarial robustness of image recovery from Radon measurements as in Scenario B2. The results of Fig. 12 show a clear superiority of the UNet that was subjected to noise during training (see also Fig. S7 for the effect of Poisson noise). Without jittering, almost imperceptible distortions in the FBP inversions are intensified by the post-processing network (see blue arrows).

The above observations can be related to the notion of *inverse crimes* in the literature on inverse problems, e.g., see [76], [77]. This term is commonly used to explain the phenomenon of exact, but highly unstable, recovery from noiseless, simulated measurements. In a similar way, networks seem to learn accurate, but unstable, reconstruction rules if they are trained with noiseless data. We note that this

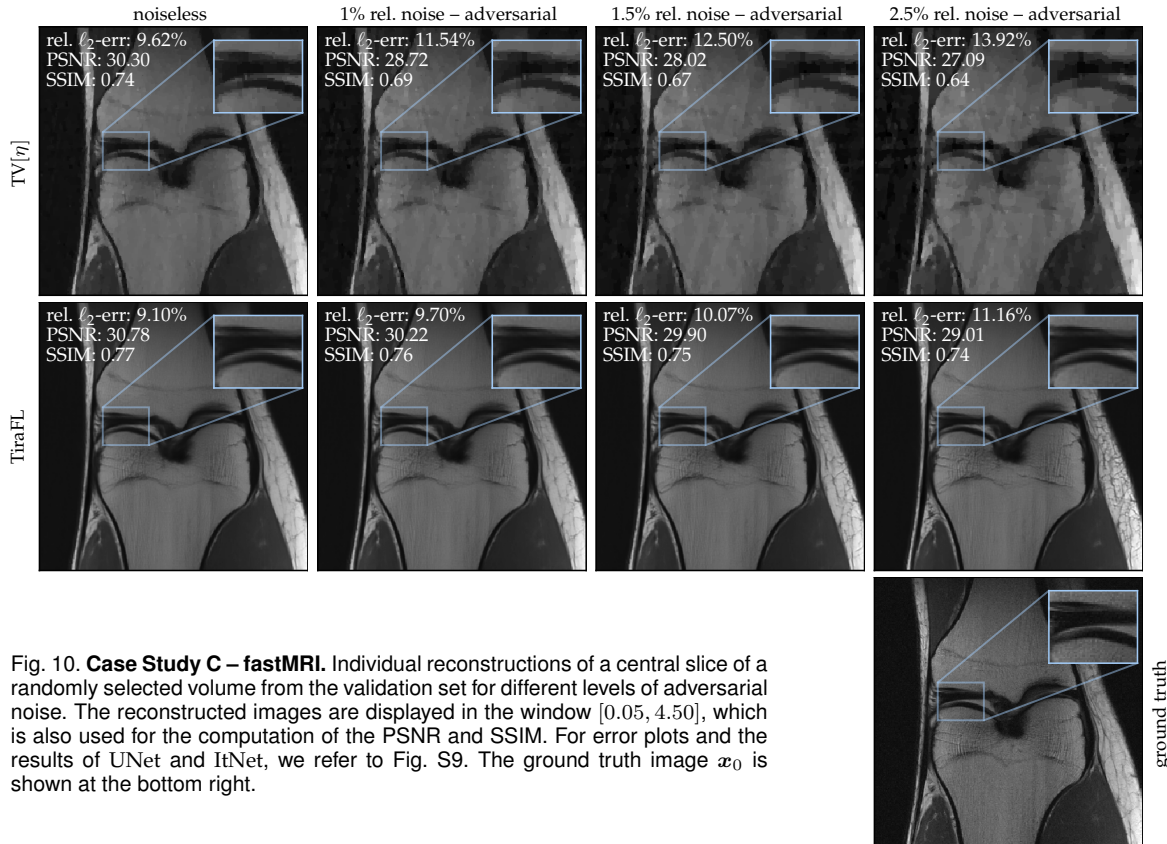


Fig. 10. **Case Study C – fastMRI.** Individual reconstructions of a central slice of a randomly selected volume from the validation set for different levels of adversarial noise. The reconstructed images are displayed in the window $[0.05, 4.50]$, which is also used for the computation of the PSNR and SSIM. For error plots and the results of UNet and ItNet, we refer to Fig. S9. The ground truth image x_0 is shown at the bottom right.

does not only concern simulated phantom data but also real-world scenarios. Indeed, in medical imaging applications, one often acquires fully sampled (noisy) reference scans $\{\tilde{y}^i\}_{i=1}^M$, which are used to generate the ground truth training images $x_0^i = \mathcal{A}_{\text{full}}^{-1} \tilde{y}^i$. The measurements are usually subsampled retrospectively by $y^i = P \tilde{y}^i$, where P denotes an appropriate selection operator. NN-based solution methods for the limited data problem (1) with $\mathcal{A} = P \mathcal{A}_{\text{full}}$ are then obtained by training on $\{(y^i, x_0^i)\}_{i=1}^M$. Importantly, such data pairs also “commit” an inverse crime, since they follow the noiseless forward model $\mathcal{A} x_0^i = P \mathcal{A}_{\text{full}} x_0^i = y^i$. Hence, we believe that simulating additional noise might be helpful in the situation of real-world measurements as well. Jittering is a simple and natural remedy in that regard that can additionally reduce overfitting [53]. The exploration of further regularization techniques or more sophisticated ways of injecting noise during training is left to future research.

5.2 Adversarial Examples for Classification From Compressed Measurements

In medical healthcare, image recovery is merely one component of the entire data-processing chain. Indeed, machine learning techniques are particularly suitable for automated diagnosis or personalized treatment recommendations. As argued in the introduction of this article, the study of adversarial examples for such classification tasks differs from the robustness analysis of reconstruction methods. In this section, we shed further light on this subject by analyzing classification from compressed measurements—think of detecting a tumor from a subsampled MRI scan.

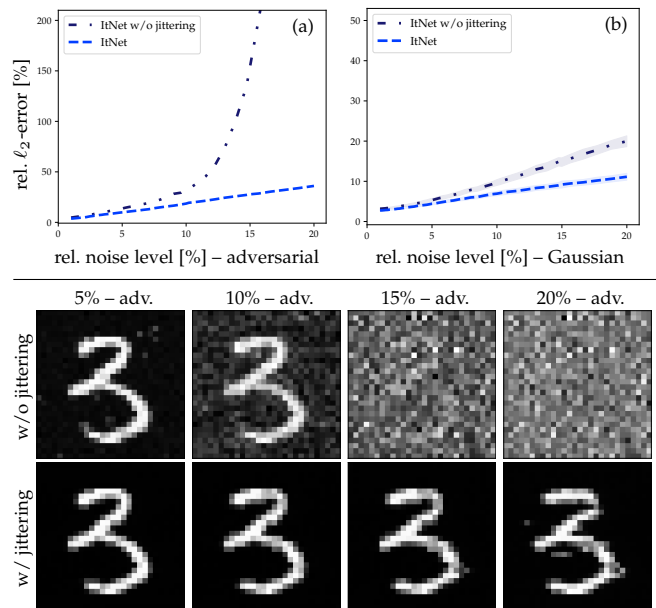


Fig. 11. **An inverse crime?** A comparison between ItNet trained with and without jittering. The above noise-to-error curves are generated for the MNIST-digit 3 from Fig. 5 with (a) adversarial and (b) Gaussian noise. Individual reconstructions for adversarial noise are shown below (the intermediate steps performed by ItNet are visualized in Fig. S11).

To this end, we revisit the toy model of Scenario A2, with the goal to predict MNIST digits from their Gaussian measurements. This is realized by training a basic convolutional NN classifier ConvNet: $\mathbb{R}^N \rightarrow [0, 1]^{10}$, mapping images to class probabilities for each of the 10 digits. The concatena-

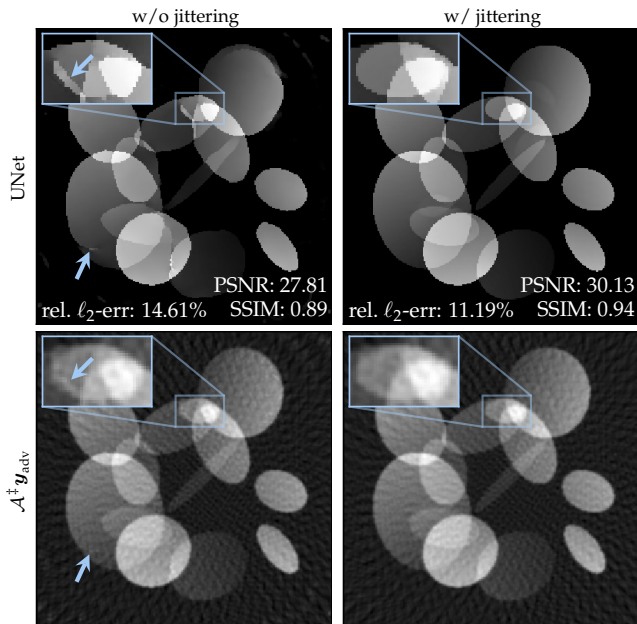


Fig. 12. **An inverse crime?** A comparison between UNet trained with and without jittering for image recovery from sparse-angle Radon measurements, see Fig. 8 in Scenario B2. The reconstructions are obtained for 1% adversarial noise. The bottom figures show the FBP inversions of the found perturbations, respectively. The blue arrows highlight tiny distortions that are amplified by the post-processing network.

tion with a reconstruction method $\text{Rec}: \mathbb{R}^m \rightarrow \mathbb{R}^N$ then yields the following classification map:

$$\text{CC}: \mathbb{R}^m \rightarrow [0, 1]^{10}, \mathbf{y} \mapsto [\text{ConvNet} \circ \text{Rec}](\mathbf{y}). \quad (6)$$

The approach of CC can be seen as a simplified model for the automated diagnosis from subsampled measurements; see also [78] and the references therein for the related problem of *compressed classification*.

Inspired by [66], we adapt the attack strategy (5) to the classification setting by (approximately) solving

$$\mathbf{e}_{\text{adv}} = \underset{\|\mathbf{e}\|_2 \leq \eta}{\text{argmax}} \max_{k \neq c} [\text{CC}(\mathbf{y}_0 + \mathbf{e})]_k - [\text{CC}(\mathbf{y}_0 + \mathbf{e})]_c$$

where $c \in \{0, 1, \dots, 9\}$ is the true class label of \mathbf{x}_0 . Fig. 13 shows a noise-to-accuracy curve visualizing the relative amount of correct classifications for different choices of Rec. The corresponding image reconstructions $\text{Rec}(\mathbf{y}_0 + \mathbf{e}_{\text{adv}})$ as well as the predicted classes $\text{argmax}_k [\text{CC}(\mathbf{y}_0 + \mathbf{e}_{\text{adv}})]_k$ for an example digit are presented below.

All classifiers exhibit a transition behavior: the success rate is almost perfect for small perturbations and then drops to zero at some point. The associated images show that we have found adversarial examples in the ordinary sense of machine learning. Indeed, every visualized reconstruction is still recognizable as the digit 9. In other words, although being stable, each of the recovery methods is capable of producing slightly perturbed images that fool the ConvNet-part. Remarkably, this phenomenon occurs independently of using a model-based or learned solver for (1). We conclude that deep-learning-based data-processing pipelines (as in medical healthcare) remain vulnerable to adversarial attacks, even if provably robust reconstruction schemes are employed.

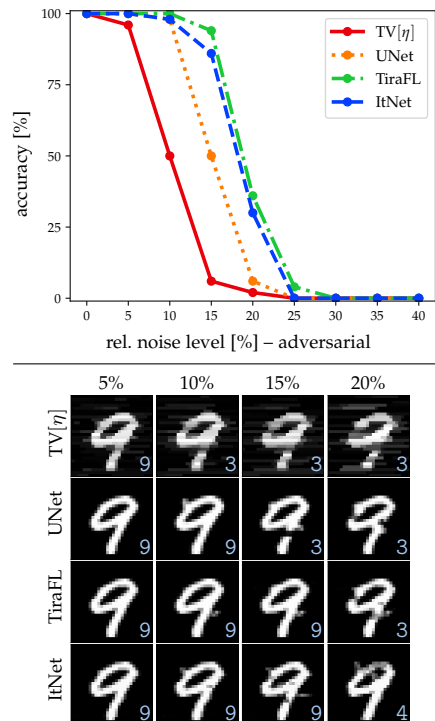


Fig. 13. **Classification from compressed measurements.** The above curve plots the relative adversarial noise level against the prediction accuracy of the classifier (6) for different recovery methods (averaged over 50 digits from the test set). The intermediate reconstructions of a randomly selected digit are shown below for different noise levels. Their predicted class labels are displayed in the bottom right corner.

5.3 The Original fastMRI Challenge Setup

This section demonstrates that the original fastMRI challenge data for single-coil MRI is more susceptible to adversarial noise. In contrast to Case Study C, the challenge measurement setup is based on omitting k-space lines in the phase-encoding direction (corresponding to 4-fold acceleration), i.e., the subsampling mask is defined by vertical lines. The resulting undersampling ratio of $\sim 23\%$ is higher than in Case Study C ($\sim 17\%$). Fig. 14 shows individual image reconstructions for $\text{TV}[\eta]$ and Tira.⁴ Compared to Fig. 10, the outcomes indicate a loss of adversarial robustness, as the reconstructed images exhibit undesired line-shaped artifacts (see blue box in Fig. 14). This phenomenon occurs regardless of using a model-based ($\text{TV}[\eta]$) or learned method (Tira). In fact, the observed artifacts are a consequence of the underlying measurement system: the anisotropic mask pattern implies that vertical image features become more “aligned” with the kernel of the forward operator. Hence, clearly visible distortions may be caused by relatively small perturbations of the measurements (cf. [25]). This confirms that the design of sampling patterns does not only influence the accuracy of a reconstruction method (e.g., see [79]), but also its adversarial robustness.

4. Since the fastMRI challenge setup does not rely on a fixed subsampling mask, the fully-learned approach for Tiramisu is not available here. Our Tira-net performs competitively in the fastMRI public leaderboard: We have achieved an SSIM of 0.765, whereas the leading method has 0.783 (<https://fastmri.org/leaderboards/>, teamname AnItalianDessert, accessed on 2020-11-08).

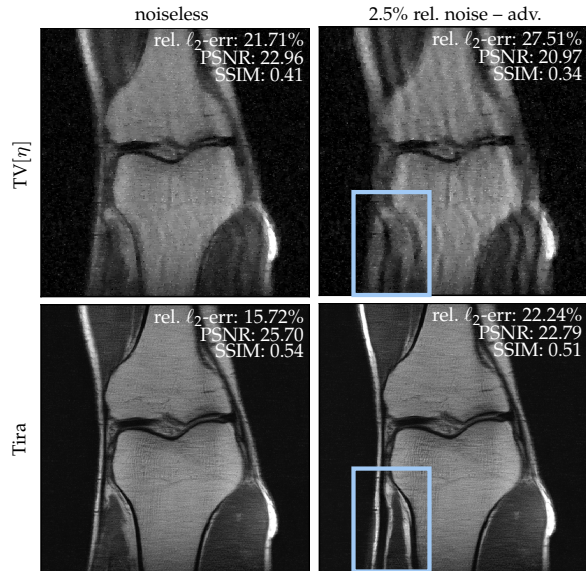


Fig. 14. **The original fastMRI challenge setup.** Reconstructions of a randomly selected image from the validation set. Compared to the analogous experiment in Fig. 10, the Fourier subsampling operator is based on vertical lines in the k-space instead of a radial mask. The reconstructed images are displayed in the window $[0.05, 4.50]$, which is also used for the computation of the PSNR and SSIM. Note that the data are given as emulated single-coil (ESC) measurements, whereas the reconstructions in Fig. 10 are based on multi-coil images. Hence, the signal-to-noise ratios are not directly comparable.

6 DISCUSSION

In an extensive series of experiments, this work has analyzed the robustness of deep-learning-based solution methods for inverse problems. Central to our approach was to study the effect of adversarial noise, i.e., worst-case perturbations of the measurements that maximize the reconstruction error. A systematic comparison with a model-based reference method has shown that standard deep NN schemes are remarkably resilient against statistical and adversarial distortions. On the other hand, we have demonstrated that instabilities might be caused by the “inverse crime” of training with noiseless data. A simple remedy in that regard is jittering—a standard regularization and robustification technique in deep learning [8]. However, it is well known that this does not cure the adversarial vulnerability of deep NN classifiers, which requires more sophisticated defense strategies [80]. While such defenses may also improve the robustness in the context of image recovery [26], our results allow for a surprising conclusion: Injecting Gaussian random noise in the training phase seems sufficient to obtain solution methods for inverse problems that are resistant to other types of noise, including adversarial perturbations.

Admittedly, there are several aspects that go beyond the scope of our study: (i) We are restricted to a selection of end-to-end NN architectures, excluding other approaches, such as generative models [16], the deep image prior [81], or learned regularizers [82]. However, since these algorithms typically involve more model-based components, we expect their robustness to be comparable to the schemes considered in the present work. (ii) Due to the non-convexity of (5), a theoretical optimality certificate for our attack strategy is lacking. Nevertheless, our results provide empirical evi-

dence that we have solved the problem adequately: The gap between worst-case and statistical perturbations appears consistent across all considered scenarios. More importantly, we have verified the ability to detect an error blowup caused by adversarial noise (see Fig. 11). (iii) Our analysis takes a mathematical perspective on robustness, thereby relying on standard similarity measures, in particular, the Euclidean norm. It is well known that such quantitative metrics are insensitive to several types of visual distortions. For instance, even the winning networks of the fastMRI challenge were unable to capture certain tiny pathological features that rarely appear in the data [83]. While some of these details are possibly lost in the subsampling process, this failure could also be due to instabilities of deep learning techniques; see [84] for recent progress in that direction.

The relevance of artificial intelligence for future health-care is undeniable. In this field, reliable reconstruction methods are indispensable, since errors caused by instabilities can be fatal. Although there is typically no “adversary” in medical imaging (i.e., an agent that intentionally manipulates the measurements), it is reassuring to know the limits of what could go wrong in principle. Of more practical interest is the robustness against random perturbations, which is a more realistic noise model for common imaging modalities. We believe that our work makes progress in both regards, by showing optimistic results on the use of deep NNs for inverse problems in imaging.

ACKNOWLEDGMENTS

M.G. and M.M. acknowledge support by the DFG Priority Programme DFG-SPP 1798. We express our gratitude to the Institute of Mathematics of the Technical University of Berlin for providing us hardware resources to realize the numerical experiments presented in this work.

REFERENCES

- [1] M. Lustig, D. L. Donoho, J. M. Santos, and J. M. Pauly, “Compressed sensing MRI,” *IEEE Signal Process. Mag.*, vol. 25, no. 2, pp. 72–82, 2008.
- [2] J. Haupt, W. U. Bajwa, M. Rabbat, and R. Nowak, “Compressed sensing for networked data,” *IEEE Signal Process. Mag.*, vol. 25, no. 2, pp. 92–101, 2008.
- [3] J. L. Starck, E. Pantin, and F. Murtagh, “Deconvolution in astronomy: A review,” *Publ. Astron. Soc. Pac.*, vol. 114, no. 800, pp. 1051–1069, 2002.
- [4] A. Tarantola and B. Valetta, “Inverse problems = quest for information,” *J. Geophys.*, vol. 50, no. 1, pp. 159–170, 1981.
- [5] S. Foucart and H. Rauhut, *A Mathematical Introduction to Compressive Sensing*, ser. Applied and Numerical Harmonic Analysis. Birkhäuser Basel, 2013.
- [6] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “ImageNet classification with deep convolutional neural networks,” in *Advances in Neural Information Processing Systems 25*, F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, Eds. Curran Associates, Inc., 2012, pp. 1097–1105.
- [7] Y. LeCun, Y. Bengio, and G. Hinton, “Deep learning,” *Nature*, vol. 521, no. 7553, pp. 436–444, 2015.
- [8] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*. MIT Press, 2016.
- [9] K. Gregor and Y. LeCun, “Learning fast approximations of sparse coding,” in *Proceedings of the 27th International Conference on International Conference on Machine Learning (ICML)*, J. Fürnkranz and T. Joachims, Eds., 2010, pp. 399–406.

- [10] Y. Yang, J. Sun, H. Li, and Z. Xu, "Deep ADMM-Net for compressive sensing MRI," in *Advances in Neural Information Processing Systems* 29, D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, Eds. Curran Associates, Inc., 2016, pp. 10–18.
- [11] C. Dong, C. C. Loy, K. He, and X. Tang, "Image super-resolution using deep convolutional networks," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 38, no. 2, pp. 295–307, 2016.
- [12] E. Kang, J. Min, and J. C. Ye, "A deep convolutional neural network using directional wavelets for low-dose X-ray CT reconstruction," *Med. Phys.*, vol. 44, no. 10, pp. e360–e375, 2017.
- [13] K. H. Jin, M. T. McCann, E. Froustey, and M. Unser, "Deep convolutional neural network for inverse problems in imaging," *IEEE Trans. Image Process.*, vol. 26, no. 9, pp. 4509–4522, 2017.
- [14] K. Hammernik, T. Klatzer, E. Kobler, M. P. Recht, D. K. Sodickson, T. Pock, and F. Knoll, "Learning a variational network for reconstruction of accelerated MRI data," *Magn. Reson. Med.*, vol. 79, no. 6, pp. 3055–3071, 2018.
- [15] H. Chen, Y. Zhang, W. Zhang, P. Liao, K. Li, J. Zhou, and G. Wang, "Low-dose CT via convolutional neural network," *Biomed. Opt. Express*, vol. 8, no. 2, pp. 679–694, 2017.
- [16] A. Bora, A. Jalal, E. Price, and A. G. Dimakis, "Compressed sensing using generative models," in *Proceedings of the 34th International Conference on Machine Learning (ICML)*, D. Precup and Y. W. Teh, Eds., vol. 70, 2017, pp. 537–546.
- [17] B. Zhu, J. Z. Liu, S. F. Cauley, B. R. Rosen, and M. S. Rosen, "Image reconstruction by domain-transform manifold learning," *Nature*, vol. 555, no. 7697, pp. 487–492, 2018.
- [18] T. A. Bubba, G. Kutyniok, M. Lassas, M. März, W. Samek, S. Siltanen, and V. Srinivasan, "Learning the invisible: A hybrid deep learning-shearlet framework for limited angle computed tomography," *Inverse Probl.*, vol. 35, no. 6, p. 064002, 2019.
- [19] S. Arridge, P. Maass, O. Öktem, and C.-B. Schönlieb, "Solving inverse problems using data-driven models," *Acta Numer.*, vol. 28, pp. 1–174, 2019.
- [20] M. Elad, "Deep, Deep Trouble: Deep Learning's Impact on Image Processing, Mathematics, and Humanity," SIAM News, available online: <https://sinews.siam.org/Details-Page/deep-deep-trouble>, 2017.
- [21] H. Chen, Y. Zhang, M. K. Kalra, F. Lin, Y. Chen, P. Liao, J. Zhou, and G. Wang, "Low-dose CT with a residual encoder-decoder convolutional neural network," *IEEE Trans. Med. Imag.*, vol. 36, no. 12, pp. 2524–2535, 2017.
- [22] A. Hauptmann, J. Adler, S. R. Arridge, and O. Öktem, "Multi-scale learned iterative reconstruction," *IEEE Trans. Comput. Imag.*, 2020, available online: <https://doi.org/10.1109/TCI.2020.2990299>.
- [23] Y. Huang, T. Würfl, K. Breininger, L. Liu, G. Lauritsch, and A. Maier, "Some investigations on robustness of deep learning in limited angle tomography," in *Medical Image Computing and Computer Assisted Intervention – MICCAI 2018*, A. F. Frangi, J. A. Schnabel, C. Davatzikos, C. Alberola-López, and G. Fichtinger, Eds. Springer Cham, 2018, pp. 145–153.
- [24] V. Antun, F. Renna, C. Poon, B. Adcock, and A. C. Hansen, "On instabilities of deep learning in image reconstruction and the potential costs of AI," *Proc. Natl. Acad. Sci.*, 2020, available online: <https://doi.org/10.1073/pnas.1907377117>.
- [25] N. M. Gottschling, V. Antun, B. Adcock, and A. C. Hansen, "The troublesome kernel: why deep learning for inverse problems is typically unstable," 2020, preprint arXiv:2001.01258.
- [26] A. Raj, Y. Bresler, and B. Li, "Improving robustness of deep-learning-based image reconstruction," in *Proceedings of the 37th International Conference on Machine Learning (ICML)*, H. Daumé and A. Singh, Eds., 2017.
- [27] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus, "Intriguing properties of neural networks," 2014, preprint arXiv:1312.6199.
- [28] A. Kurakin, I. Goodfellow, and S. Bengio, "Adversarial examples in the physical world," 2017, preprint arXiv:1607.02533.
- [29] K. Eykholt, I. Evtimov, E. Fernandes, B. Li, A. Rahmati, C. Xiao, A. Prakash, T. Kohno, and D. Song, "Robust physical-world attacks on deep learning visual classification," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 1625–1634.
- [30] J. Zbontar, F. Knoll, A. Sriram, T. Murrell, Z. Huang, M. J. Muckley, A. Defazio, R. Stern, P. Johnson, M. Bruno, M. Parente, K. J. Geras, J. Katsnelson, H. Chandarana, Z. Zhang, M. Drozdal, A. Romero, M. Rabbat, P. Vincent, N. Yakubova, J. Pinkerton, D. Wang, E. Owens, C. L. Zitnick, M. P. Recht, D. K. Sodickson, and Y. W. Lui, "fastMRI: an open dataset and benchmarks for accelerated MRI," 2018, preprint arXiv:1811.08839.
- [31] F. Knoll, J. Zbontar, A. Sriram, M. J. Muckley, M. Bruno, A. Defazio, M. Parente, K. J. Geras, J. Katsnelson, H. Chandarana, Z. Zhang, M. Drozdal, A. Romero, M. Rabbat, P. Vincent, J. Pinkerton, D. Wang, N. Yakubova, E. Owens, C. L. Zitnick, M. P. Recht, D. K. Sodickson, and Y. W. Lui, "fastMRI: a publicly available raw k-space and DICOM dataset of knee images for accelerated MR image reconstruction using machine learning," *Radiology Artif. Intell.*, vol. 2, no. 1, p. e190007, 2020.
- [32] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer, "Automatic differentiation in PyTorch," Contribution to the NIPS 2017 Autodiff Workshop, available online: <https://openreview.net/forum?id=BjJsrmlfCZ>, 2017.
- [33] N. Carlini, "A Complete List of All (arXiv) Adversarial Example Papers," available online: <https://nicholas.carlini.com/writing/2019/all-adversarial-example-papers.html>, accessed on 2020-11-02, 2020.
- [34] X. Yuan, P. He, Q. Zhu, and X. Li, "Adversarial examples: Attacks and defenses for deep learning," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 30, no. 9, pp. 2805–2824, 2019.
- [35] G. Ortiz-Jimenez, A. Modas, S.-M. Moosavi-Dezfooli, and P. Frossard, "Optimism in the face of adversity: Understanding and improving deep learning through adversarial robustness," 2020, preprint arXiv:2010.09624.
- [36] A. Arnab, O. Miksik, and P. H. Torr, "On the robustness of semantic segmentation models to adversarial attacks," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 42, no. 12, pp. 3040–3053, 2020.
- [37] E. T. Quinto, "Singularities of the X-Ray Transform and Limited Data Tomography in \mathbb{R}^2 and \mathbb{R}^3 ," *SIAM J. Math. Anal.*, vol. 24, no. 5, pp. 1215–1225, 1993.
- [38] J. Schlemper, J. Caballero, J. V. Hajnal, A. N. Price, and D. Rueckert, "A deep cascade of convolutional neural networks for dynamic mr image reconstruction," *IEEE Trans. Med. Imag.*, vol. 37, no. 2, pp. 491–503, 2017.
- [39] E. Kobler, A. Effland, K. Kunisch, and T. Pock, "Total deep variation: A stable regularizer for inverse problems," 2020, arXiv:2006.08789.
- [40] G. Ongie, A. Jalal, R. G. Baraniuk, C. A. Metzler, A. G. Dimakis, and R. Willett, "Deep learning techniques for inverse problems in imaging," *IEEE J. Sel. Areas Inf. Theory*, vol. 1, no. 1, pp. 39–56, 2020.
- [41] O. Ronneberger, P. Fischer, and T. Brox, "U-Net: convolutional networks for biomedical image segmentation," in *Medical Image Computing and Computer Assisted Intervention – MICCAI 2015*, N. Navab, J. Hornegger, W. M. Wells, and A. F. Frangi, Eds. Springer Cham, 2015, pp. 234–241.
- [42] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 770–778.
- [43] J. M. Wolterink, T. Leiner, M. A. Viergever, and I. Išgum, "Generative adversarial networks for noise reduction in low-dose CT," *IEEE Trans. Med. Imag.*, vol. 36, no. 12, pp. 2536–2545, 2017.
- [44] Q. Yang, P. Yan, Y. Zhang, H. Yu, Y. Shi, X. Mou, M. K. Kalra, Y. Zhang, L. Sun, and G. Wang, "Low-dose CT image denoising using a generative adversarial network with Wasserstein distance and perceptual loss," *IEEE Trans. Med. Imag.*, vol. 37, no. 6, pp. 1348–1357, 2018.
- [45] S. Jégou, M. Drozdal, D. Vazquez, A. Romero, and Y. Bengio, "The one hundred layers tiramisú: Fully convolutional densenets for semantic segmentation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 11–19.
- [46] G. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger, "Densely connected convolutional networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 4700–4708.
- [47] J. Schlemper, I. Oksuz, J. R. Clough, J. Duan, A. P. King, J. A. Schnabel, J. V. Hajnal, and D. Rueckert, "dAUTOMAP: decomposing AUTOMAP to achieve scalability and enhance performance," 2019, preprint arXiv:1909.10995.
- [48] H. K. Aggarwal, M. P. Mani, and M. Jacob, "MoDL: model-based deep learning architecture for inverse problems," *IEEE Trans. Med. Imag.*, vol. 38, no. 2, pp. 394–405, 2018.
- [49] J. Adler and O. Öktem, "Learned primal-dual reconstruction," *IEEE Trans. Med. Imag.*, vol. 37, no. 6, pp. 1322–1332, 2018.

- [50] K. Hammernik, J. Schlemper, C. Qin, J. Duan, R. M. Summers, and D. Rueckert, " σ -net: Systematic evaluation of iterative deep neural networks for fast parallel MR Image reconstruction," 2019, preprint arXiv:1912.09278.
- [51] I. Y. Chun, Z. Huang, H. Lim, and J. Fessler, "Momentum-Net: fast and convergent iterative neural network for inverse problems," *IEEE Trans. Pattern Anal. Mach. Intell.*, 2020, available online: <https://doi.org/10.1109/TPAMI.2020.3012955>.
- [52] D. P. Kingma and J. Ba, "Adam: a method for stochastic optimization," 2014, preprint arXiv:1412.6980.
- [53] J. Sietsma and R. J. Dow, "Creating artificial neural networks that generalize," *Neural Netw.*, vol. 4, no. 1, pp. 67–79, 1991.
- [54] L. Holmstrom and P. Koistinen, "Using additive noise in back-propagation training," *IEEE Trans. Neural Netw.*, vol. 3, no. 1, pp. 24–38, 1992.
- [55] C. M. Bishop, "Training with noise is equivalent to Tikhonov regularization," *Neural Comput.*, vol. 7, no. 1, pp. 108–116, 1995.
- [56] P. Vincent, H. Larochelle, I. Lajoie, Y. Bengio, and P.-A. Manzagol, "Stacked denoising autoencoders: Learning useful representations in a deep network with a local denoising criterion," *J. Mach. Learn. Res.*, vol. 11, pp. 3371–3408, 2010.
- [57] L. I. Rudin, S. Osher, and E. Fatemi, "Nonlinear total variation based noise removal algorithms," *Physica D: Nonlinear Phenomena*, vol. 60, no. 1–4, pp. 259–268, 1992.
- [58] A. Chambolle and P.-L. Lions, "Image recovery via total variation minimization and related problems," *Numer. Math.*, vol. 76, no. 2, pp. 167–188, 1997.
- [59] M. Benning and M. Burger, "Modern regularization methods for inverse problems," *Acta Numer.*, vol. 27, pp. 1–111, 2018.
- [60] E. J. Candès, J. K. Romberg, and T. Tao, "Robust uncertainty principles: exact signal reconstruction from highly incomplete frequency information," *IEEE Trans. Inf. Theory*, vol. 52, no. 2, pp. 489–509, 2006.
- [61] D. Needell and R. Ward, "Near-optimal compressed sensing guarantees for total variation minimization," *IEEE Trans. Imag. Proc.*, vol. 22, no. 10, pp. 3941–3949, 2013.
- [62] C. Poon, "On the role of total variation in compressed sensing," *SIAM J. Imag. Sci.*, vol. 8, no. 1, pp. 682–720, 2015.
- [63] M. Genzel, M. März, and R. Seidel, "Compressed sensing with 1D total variation: Breaking sample complexity barriers via non-uniform recovery," 2020, preprint arXiv:2001.09952.
- [64] R. Glowinski and A. Marroco, "Sur l'approximation, par éléments finis d'ordre un, et la résolution, par pénalisation-dualité d'une classe de problèmes de Dirichlet non linéaires," *RAIRO Anal. Numer.*, vol. 9, no. R2, pp. 41–76, 1975.
- [65] D. Gabay and B. Mercier, "A dual algorithm for the solution of nonlinear variational problems via finite element approximation," *Comput. Math. Appl.*, vol. 2, no. 1, pp. 17–40, 1976.
- [66] N. Carlini and D. Wagner, "Towards evaluating the robustness of neural networks," in *IEEE Symposium on Security and Privacy*, 2017, pp. 39–57.
- [67] A. Agrawal, B. Amos, S. Barratt, S. Boyd, S. Diamond, and J. Z. Kolter, "Differentiable convex optimization layers," in *Advances in Neural Information Processing Systems 32*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché Buc, E. Fox, and R. Garnett, Eds. Curran Associates, Inc., 2019, pp. 9562–9574.
- [68] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proc. IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
- [69] D. Amelunxen, M. Lotz, M. B. McCoy, and J. A. Tropp, "Living on the edge: phase transitions in convex programs with random data," *Inf. Inference*, vol. 3, no. 3, pp. 224–294, 2014.
- [70] J. Adler and O. Öktem, "Solving ill-posed inverse problems using iterative deep neural networks," *Inverse Probl.*, vol. 33, no. 12, p. 124007, 2017.
- [71] E. Y. Sidky and X. Pan, "Image reconstruction in circular cone-beam computed tomography by constrained, total-variation minimization," *Phys. Med. Biol.*, vol. 53, no. 17, pp. 4777–4807, 2008.
- [72] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli, "Image quality assessment: from error visibility to structural similarity," *IEEE Trans. Image Process.*, vol. 13, no. 4, pp. 600–612, 2004.
- [73] P. Ernst, "Pytorch implementation of scikit-image's radon function, version 0.1.4," available online: https://github.com/phernst/pytorch_radon, 2020.
- [74] N. Papernot, P. McDaniel, and I. Goodfellow, "Transferability in machine learning: from phenomena to black-box attacks using adversarial samples," 2016, preprint arXiv:1605.07277.
- [75] H. Zhao, O. Gallo, I. Frosio, and J. Kautz, "Loss functions for image restoration with neural networks," *IEEE Trans. Comput. Imag.*, vol. 3, no. 1, pp. 47–57, 2017.
- [76] J. Kaipio and E. Somersalo, *Statistical and computational inverse problems*, ser. Applied Mathematical Sciences. Springer New York, 2006, vol. 160.
- [77] J. L. Mueller and S. Siltanen, *Linear and nonlinear inverse problems with practical applications*. SIAM, 2012.
- [78] A. S. Bandeira, D. G. Mixon, and B. Recht, "Compressive classification and the rare eclipse problem," in *Compressed Sensing and its Applications: Second International MATHEON Conference 2015*, ser. Applied and Numerical Harmonic Analysis, H. Boche, G. Caire, R. Calderbank, M. März, G. Kutyniok, and R. Mathar, Eds. Springer Cham, 2017, pp. 197–220.
- [79] C. Boyer, N. Chauffert, P. Ciuciu, J. Kahn, and P. Weiss, "On the generation of sampling schemes for magnetic resonance imaging," *SIAM J. Imaging Sci.*, vol. 9, no. 4, pp. 2039–2072, 2016.
- [80] G. R. Machado, E. Silva, and R. R. Goldschmidt, "Adversarial machine learning in image classification: A survey towards the defender's perspective," 2020, preprint arXiv:2009.03728.
- [81] D. Ulyanov, A. Vedaldi, and V. Lempitsky, "Deep image prior," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 9446–9454.
- [82] H. Li, J. Schwab, S. Antholzer, and M. Haltmeier, "NETT: solving inverse problems with deep neural networks," *Inverse Probl.*, vol. 36, no. 6, p. 065005, 2020.
- [83] F. Knoll, T. Murrell, A. Sriram, N. Yakubova, J. Zbontar, M. Rabbat, A. Defazio, M. J. Muckley, D. K. Sodickson, C. L. Zitnick, and M. P. Recht, "Advancing machine learning for MR image reconstruction with an open competition: Overview of the 2019 fastMRI challenge," *Magn. Reson. Med.*, vol. 84, no. 6, pp. 3054–3070, 2020.
- [84] K. Cheng, F. Calivá, R. Shah, M. Han, S. Majumdar, and V. Pedoia, "Addressing the false negative problem of deep learning MRI reconstruction models by adversarial attacks and robust training," in *Proceedings of the 3rd Conference on Medical Imaging with Deep Learning (MIDL)*, T. Arbel, I. B. Ayed, M. de Bruijne, M. Descoteaux, H. Lombaert, and C. Pal, Eds., vol. 121, 2020, pp. 121–135.

Supplementary Material

The supplementary material is organized as follows:

- Section S1–S3 contain supplementary results for Case Study A–C, respectively.
- Section S4 contains supplementary results for Section 5.
- Section S5 provides an overview of all empirically selected hyper-parameters for the considered network architectures, training processes, and adversarial attacks.

S1 SUPPLEMENTARY RESULTS FOR CASE STUDY A (CS WITH GAUSSIAN MEASUREMENTS)

rel. noise – adversarial		0.0%	0.1%	0.5%	1.0%	2.0%	4.0%	6.0%
TV[η]	rel. ℓ_2 -err. [%]	0.00±0.00	0.32±0.08	1.66±0.42	3.36±0.86	6.63±1.67	12.26±2.57	17.21±2.98
UNet	rel. ℓ_2 -err. [%]	2.53±1.97	2.67±2.01	3.36±2.11	4.49±2.31	7.29±2.52	13.15±3.24	18.27±3.58
UNetFL	rel. ℓ_2 -err. [%]	2.01±1.70	2.14±1.73	2.82±1.84	3.95±2.01	6.46±2.21	11.91±2.54	16.98±2.72
Tira	rel. ℓ_2 -err. [%]	1.22±1.15	1.33±1.18	1.95±1.33	3.05±1.64	5.90±2.23	11.97±3.13	17.18±3.27
TiraFL	rel. ℓ_2 -err. [%]	0.98±0.88	1.10±0.92	1.73±1.15	2.74±1.42	5.32±1.89	11.07±2.82	16.43±3.42
ItNet	rel. ℓ_2 -err. [%]	0.45±0.18	0.52±0.18	0.93±0.17	1.80±0.80	4.50±1.90	11.62±3.67	16.42±3.70

TABLE S1

Scenario A1 – CS with 1D signals. A numerical representation of the results of Fig. 2(c), including the additional methods UNetFL and Tira. The smallest relative error per noise level is highlighted in bold.

rel. noise – Gaussian		0.0%	0.1%	0.5%	1.0%	2.0%	4.0%	6.0%
TV[η]	rel. ℓ_2 -err. [%]	0.00±0.00	0.18±0.06	0.90±0.30	1.79±0.60	3.61±1.18	6.84±2.15	9.74±2.64
UNet	rel. ℓ_2 -err. [%]	2.53±1.97	2.58±1.99	2.79±2.03	3.09±2.07	3.82±2.11	5.67±2.25	7.70±2.49
UNetFL	rel. ℓ_2 -err. [%]	2.01±1.70	2.05±1.71	2.24±1.75	2.53±1.79	3.23±1.91	4.99±2.05	7.03±2.21
Tira	rel. ℓ_2 -err. [%]	1.22±1.15	1.26±1.16	1.43±1.21	1.71±1.25	2.35±1.43	4.15±1.79	6.42±2.23
TiraFL	rel. ℓ_2 -err. [%]	0.98±0.88	1.02±0.90	1.20±0.95	1.48±1.05	2.10±1.24	3.86±1.63	5.98±2.19
ItNet	rel. ℓ_2 -err. [%]	0.45±0.18	0.47±0.18	0.59±0.17	0.80±0.17	1.38±0.50	2.91±0.99	5.33±2.15

TABLE S2

Scenario A1 – CS with 1D signals. A numerical representation of the results of Fig. 2(d), including the additional methods UNetFL and Tira. The smallest relative error per noise level is highlighted in bold.

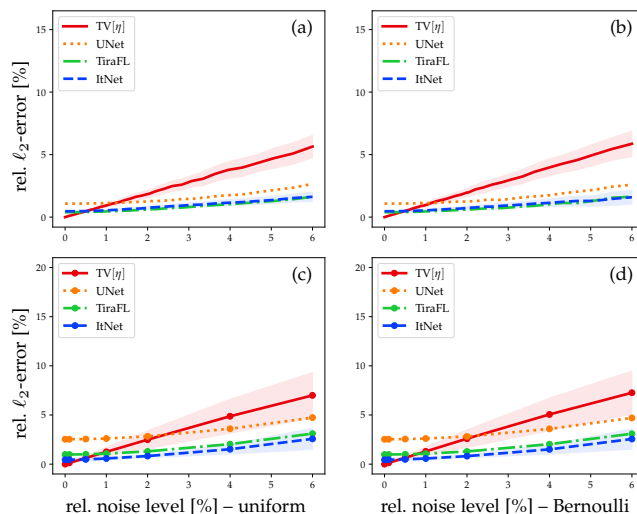


Fig. S1. **Scenario A1 – CS with 1D signals.** (a) and (b) show uniform and Bernoulli noise-to-error curves, respectively, for the signal of Fig. 3. In the latter case, we have generated symmetrized Bernoulli noise with $p = 0.025$. The mean and standard deviation are computed over 200 draws of e . (c) and (d) display the respective curves averaged over 50 signals from the test set. For the sake of clarity, we have omitted the standard deviations for UNet and TiraFL, which behave similarly.

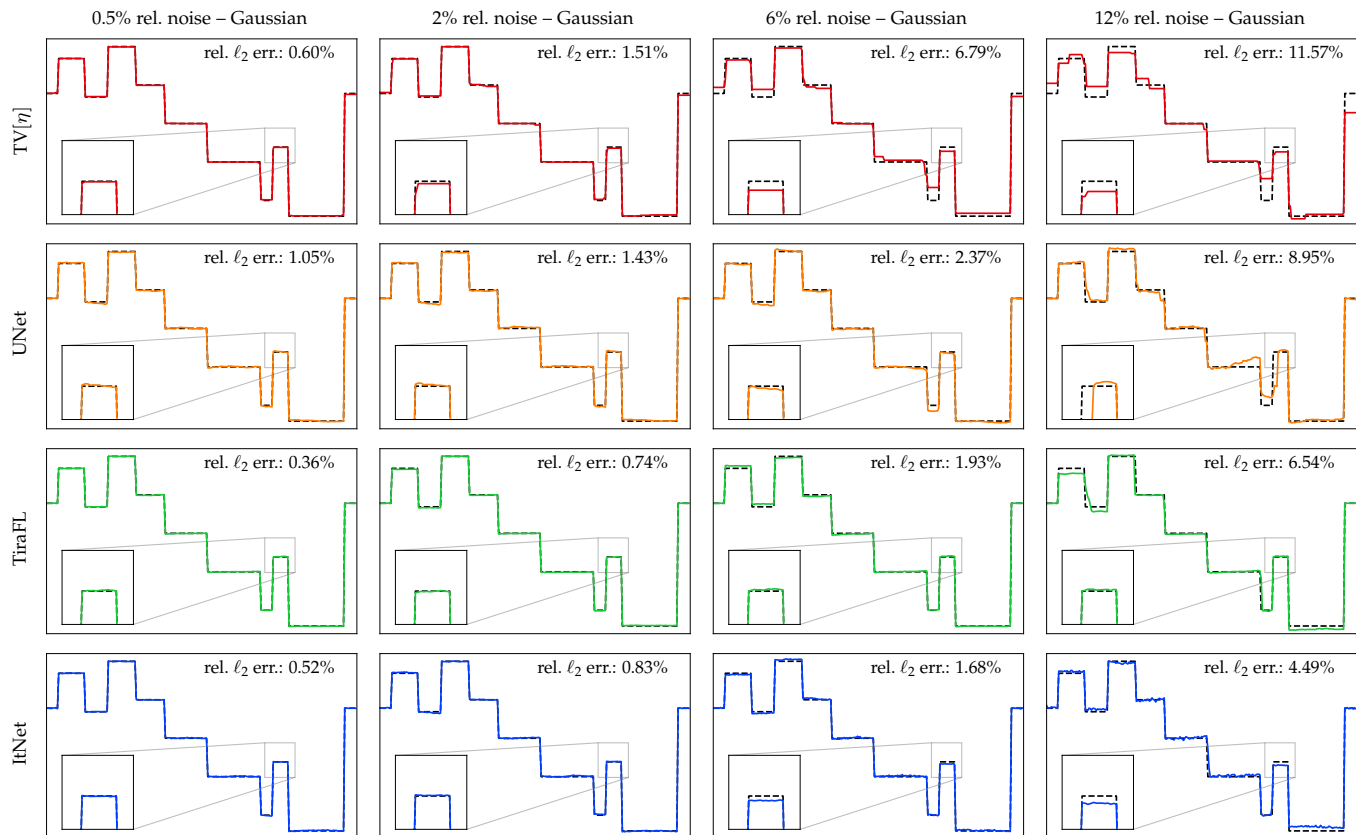


Fig. S2. **Scenario A1 – CS with 1D signals.** Individual reconstructions of the signal from Fig. 3 under Gaussian noise. The ground truth signal is visualized by a dashed line. In favor of the more insightful noise level 12%, we have omitted the noiseless case.

rel. noise – adversarial		0.0%	0.5%	1.0%	3.0%	5.0%	7.5%	10.0%
TV[η]	rel. ℓ_2 -err. [%]	15.32 \pm 10.13	18.18 \pm 9.83	20.50 \pm 9.60	28.68 \pm 8.77	35.92 \pm 8.45	43.87 \pm 7.95	50.85 \pm 7.35
UNet	rel. ℓ_2 -err. [%]	9.79 \pm 2.14	10.24 \pm 2.17	10.71 \pm 2.19	12.96 \pm 2.37	15.71 \pm 2.58	20.23 \pm 2.91	25.08 \pm 3.15
UNetFL	rel. ℓ_2 -err. [%]	7.88 \pm 1.42	8.23 \pm 1.42	8.60 \pm 1.42	10.23 \pm 1.42	12.13 \pm 1.45	14.97 \pm 1.47	18.28 \pm 1.51
Tira	rel. ℓ_2 -err. [%]	8.56 \pm 1.77	8.95 \pm 1.78	9.37 \pm 1.79	11.21 \pm 1.81	13.50 \pm 1.80	16.87 \pm 1.95	20.66 \pm 2.11
TiraFL	rel. ℓ_2 -err. [%]	7.64 \pm 1.38	7.99 \pm 1.37	8.36 \pm 1.36	9.99 \pm 1.34	11.94 \pm 1.34	14.91 \pm 1.28	18.25 \pm 1.30
ItNet	rel. ℓ_2 -err. [%]	2.47\pm0.58	2.96\pm0.60	3.53\pm0.60	6.26\pm0.59	9.35\pm0.72	13.62\pm1.31	18.06\pm1.77

TABLE S3

Scenario A2 – CS with MNIST. A numerical representation of the results of Fig. 4(c), including the additional methods UNetFL and Tira. The smallest relative error per noise level is highlighted in bold.

rel. noise – Gaussian		0.0%	0.5%	1.0%	3.0%	5.0%	7.5%	10.0%
TV[η]	rel. ℓ_2 -err. [%]	15.32 \pm 10.13	16.55 \pm 9.71	17.48 \pm 9.31	21.52 \pm 8.19	25.53 \pm 7.75	30.00 \pm 7.48	34.20 \pm 7.18
UNet	rel. ℓ_2 -err. [%]	9.79 \pm 2.14	9.87 \pm 2.15	9.96 \pm 2.14	10.36 \pm 2.13	10.86 \pm 2.18	11.54 \pm 2.16	12.37 \pm 2.11
UNetFL	rel. ℓ_2 -err. [%]	7.88 \pm 1.42	7.88 \pm 1.42	7.89 \pm 1.42	7.99 \pm 1.41	8.16 \pm 1.40	8.49 \pm 1.38	8.92 \pm 1.36
Tira	rel. ℓ_2 -err. [%]	8.56 \pm 1.77	8.56 \pm 1.77	8.57 \pm 1.76	8.67 \pm 1.75	8.85 \pm 1.72	9.17 \pm 1.69	9.59 \pm 1.65
TiraFL	rel. ℓ_2 -err. [%]	7.64 \pm 1.38	7.70 \pm 1.38	7.77 \pm 1.37	8.12 \pm 1.35	8.52 \pm 1.36	9.18 \pm 1.35	9.88 \pm 1.35
ItNet	rel. ℓ_2 -err. [%]	2.47\pm0.58	2.58\pm0.59	2.72\pm0.58	3.60\pm0.58	4.65\pm0.66	6.00\pm0.73	7.32\pm0.80

TABLE S4

Scenario A2 – CS with MNIST. A numerical representation of the results of Fig. 4(d), including the additional methods UNetFL and Tira. The smallest relative error per noise level is highlighted in bold.

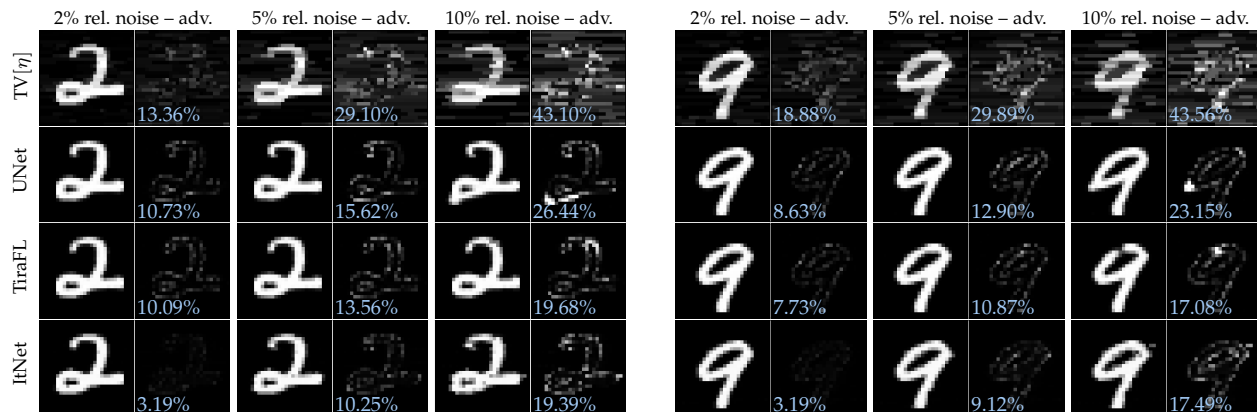


Fig. S3. **Scenario A2 – CS with MNIST.** Individual reconstructions of two additional digits from the test set for different levels of adversarial noise (see Fig. 5). The reconstructed digits and their error plots (with relative ℓ_2 -error) are displayed in the windows $[0, 1]$ and $[0, 0.6]$, respectively.



Fig. S4. **Scenario A2 – CS with MNIST.** Individual reconstructions of the digits from Fig. 5 and S3 under Gaussian noise. The reconstructed digits and their error plots (with relative ℓ_2 -error) are displayed in the windows $[0, 1]$ and $[0, 0.6]$, respectively. In favor of the more insightful noise level 25%, we have omitted 2%.

S2 SUPPLEMENTARY RESULTS FOR CASE STUDY B (IMAGE RECOVERY OF PHANTOM ELLIPSES)

rel. noise – adversarial		0.0%	0.5%	1.0%	2.0%	3.0%	5.0%	8.0%
TV $_{[\eta]}$	rel. ℓ_2 -err. [%]	0.44±0.11	2.51±0.24	4.34±0.35	7.35±0.45	9.96±0.46	14.19±0.44	20.72±0.63
	PSNR	60.00±3.26	44.73±2.20	39.98±2.16	35.38±2.00	32.74±1.91	29.66±1.82	26.37±1.93
	SSIM	1.00±0.00	0.99±0.00	0.98±0.01	0.96±0.02	0.94±0.03	0.93±0.03	0.93±0.02
UNet	rel. ℓ_2 -err. [%]	2.94±0.63	4.27±0.55	5.70±0.53	8.38±0.52	10.88±0.53	15.20±0.70	20.41±0.96
	PSNR	43.52±2.06	40.15±1.79	37.61±1.72	34.24±1.67	31.97±1.66	29.07±1.69	26.51±1.78
	SSIM	0.99±0.01	0.98±0.01	0.97±0.01	0.96±0.01	0.94±0.02	0.91±0.02	0.85±0.03
UNetFL	rel. ℓ_2 -err. [%]	2.72±0.50	4.12±0.44	5.57±0.43	8.35±0.43	10.97±0.43	15.58±0.67	21.03±1.02
	PSNR	44.13±2.29	40.46±1.96	37.80±1.87	34.28±1.76	31.90±1.73	28.85±1.76	26.25±1.82
	SSIM	0.99±0.00	0.99±0.00	0.98±0.01	0.97±0.01	0.95±0.02	0.91±0.03	0.85±0.04
Tira	rel. ℓ_2 -err. [%]	1.74±0.37	3.33±0.33	4.85±0.37	7.73±0.42	10.42±0.53	15.01±0.74	20.39±0.99
	PSNR	48.05±2.50	42.27±1.93	39.01±1.85	34.94±1.80	32.35±1.76	29.18±1.76	26.52±1.83
	SSIM	1.00±0.00	0.99±0.00	0.99±0.00	0.97±0.01	0.95±0.02	0.91±0.03	0.87±0.04
TiraFL	rel. ℓ_2 -err. [%]	1.75±0.39	3.42±0.34	4.94±0.41	7.82±0.44	10.54±0.51	15.14±0.69	20.55±0.95
	PSNR	48.05±2.58	42.05±1.93	38.85±1.86	34.85±1.83	32.24±1.80	29.10±1.75	26.45±1.81
	SSIM	1.00±0.00	0.99±0.00	0.99±0.01	0.97±0.01	0.95±0.02	0.91±0.03	0.87±0.04
ItNet	rel. ℓ_2 -err. [%]	1.45±0.29	2.81±0.28	4.21±0.32	6.87±0.37	9.37±0.40	13.65±0.49	18.98±0.65
	PSNR	49.63±1.80	43.76±1.62	40.23±1.64	35.97±1.63	33.27±1.67	30.00±1.64	27.13±1.65
	SSIM	0.99±0.00	0.99±0.00	0.98±0.00	0.97±0.01	0.96±0.01	0.92±0.02	0.88±0.03

TABLE S5

Scenario B1 – Fourier meas. with ellipses. A numerical representation of the results of Fig. 6(c), including the additional methods UNetFL and Tira. The best relative error/PSNR/SSIM per noise level is highlighted in bold. Note that the high SSIM values for TV $_{[\eta]}$ for 5% and 8% can be explained by the fact that adversarial perturbations for TV $_{[\eta]}$ cause point-like artifacts, see the zoomed region in Fig. 7. In contrast to the PSNR, the SSIM seems to be less sensitive to such types of errors.

rel. noise – Gaussian		0.0%	0.5%	1.0%	2.0%	3.0%	5.0%	8.0%
TV $_{[\eta]}$	rel. ℓ_2 -err. [%]	0.44±0.11	0.80±0.12	1.18±0.15	1.93±0.25	2.64±0.35	3.90±0.50	5.52±0.58
	PSNR	60.00±3.26	54.73±2.39	51.29±2.44	47.08±2.46	44.33±2.47	40.94±2.25	37.91±2.17
	SSIM	1.00±0.00	1.00±0.00	0.99±0.00	0.98±0.01	0.96±0.02	0.96±0.02	0.92±0.04
UNet	rel. ℓ_2 -err. [%]	2.94±0.63	2.94±0.63	2.95±0.63	3.00±0.63	3.07±0.62	3.28±0.61	3.72±0.61
	PSNR	43.52±2.06	43.50±2.06	43.47±2.06	43.33±2.07	43.12±2.07	42.51±2.09	41.41±2.12
	SSIM	0.99±0.01	0.99±0.01	0.99±0.01	0.99±0.01	0.99±0.01	0.99±0.00	0.98±0.00
UNetFL	rel. ℓ_2 -err. [%]	2.72±0.50	2.73±0.50	2.74±0.50	2.80±0.51	2.88±0.51	3.13±0.53	3.61±0.55
	PSNR	44.13±2.29	44.12±2.29	44.07±2.29	43.90±2.28	43.64±2.28	42.92±2.26	41.65±2.20
	SSIM	0.99±0.00	0.99±0.00	0.99±0.00	0.99±0.00	0.99±0.00	0.99±0.00	0.99±0.00
Tira	rel. ℓ_2 -err. [%]	1.74±0.37	1.75±0.37	1.77±0.37	1.83±0.38	1.92±0.38	2.19±0.40	2.70±0.45
	PSNR	48.05±2.50	48.02±2.50	47.95±2.49	47.65±2.48	47.20±2.46	46.03±2.41	44.19±2.36
	SSIM	1.00±0.00	1.00±0.00	1.00±0.00	1.00±0.00	1.00±0.00	0.99±0.00	0.99±0.00
TiraFL	rel. ℓ_2 -err. [%]	1.75±0.39	1.76±0.39	1.77±0.39	1.83±0.40	1.92±0.40	2.19±0.42	2.70±0.46
	PSNR	48.05±2.58	48.02±2.57	47.94±2.57	47.66±2.55	47.20±2.54	46.04±2.47	44.21±2.39
	SSIM	1.00±0.00	1.00±0.00	1.00±0.00	1.00±0.00	1.00±0.00	1.00±0.00	0.99±0.00
ItNet	rel. ℓ_2 -err. [%]	1.45±0.29	1.46±0.29	1.47±0.29	1.53±0.30	1.62±0.32	1.92±0.37	2.50±0.45
	PSNR	49.63±1.80	49.60±1.80	49.52±1.82	49.19±1.89	48.65±1.98	47.18±2.18	44.89±2.27
	SSIM	0.99±0.00	0.99±0.00	0.99±0.00	0.99±0.00	0.99±0.00	0.99±0.00	0.99±0.00

TABLE S6

Scenario B1 – Fourier meas. with ellipses. A numerical representation of the results of Fig. 6(d), including the additional methods UNetFL and Tira. The best relative error/PSNR/SSIM per noise level is highlighted in bold.

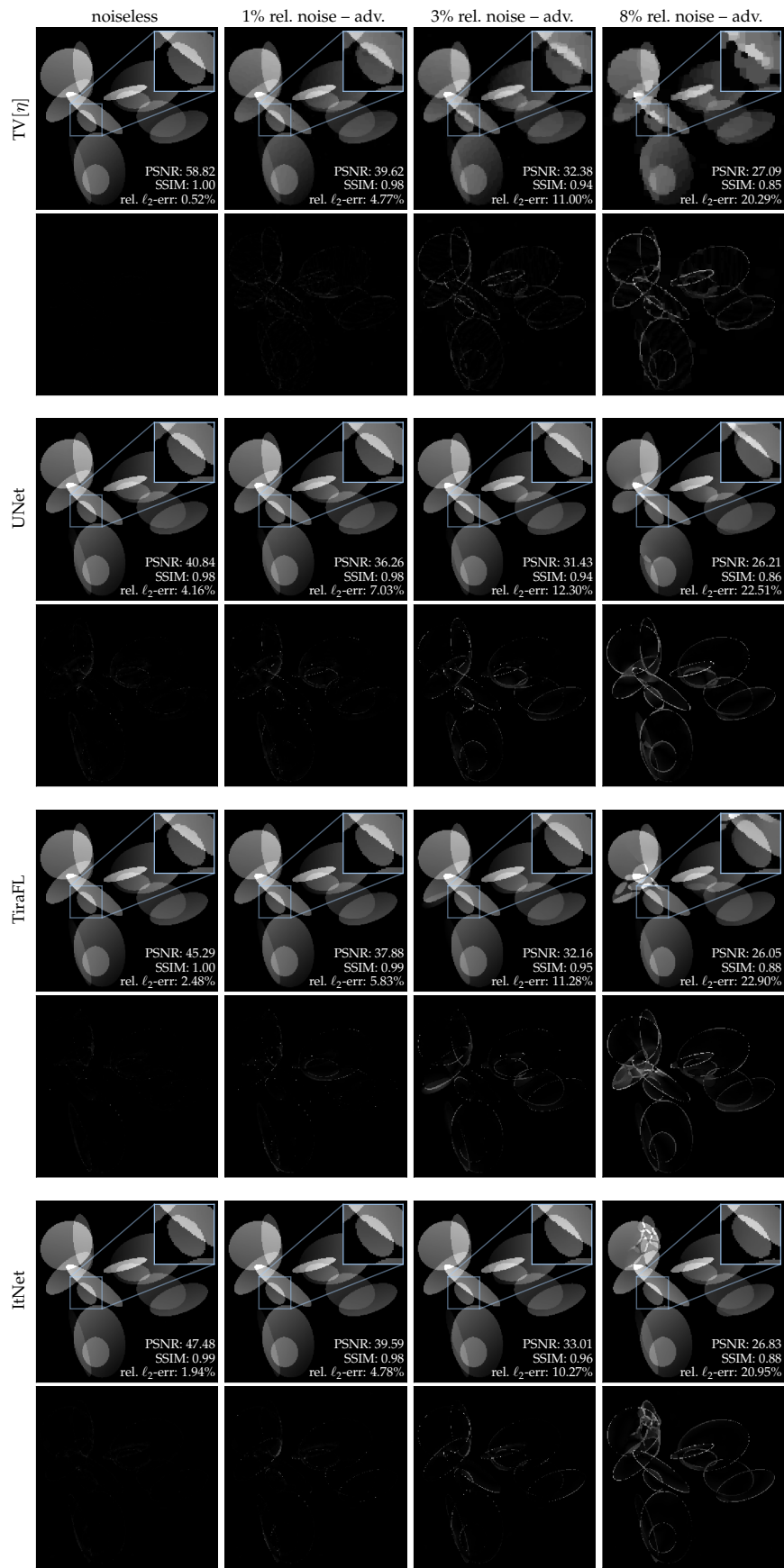


Fig. S5. **Scenario B1 – Fourier meas. with ellipses.** Individual reconstructions of the image from Fig. 7 for different levels of adversarial noise. The reconstructed images are displayed in the window $[0, 0.9]$, which is also used for the computation of the PSNR and SSIM. The error plots shown below each reconstruction are displayed in the window $[0, 0.6]$.

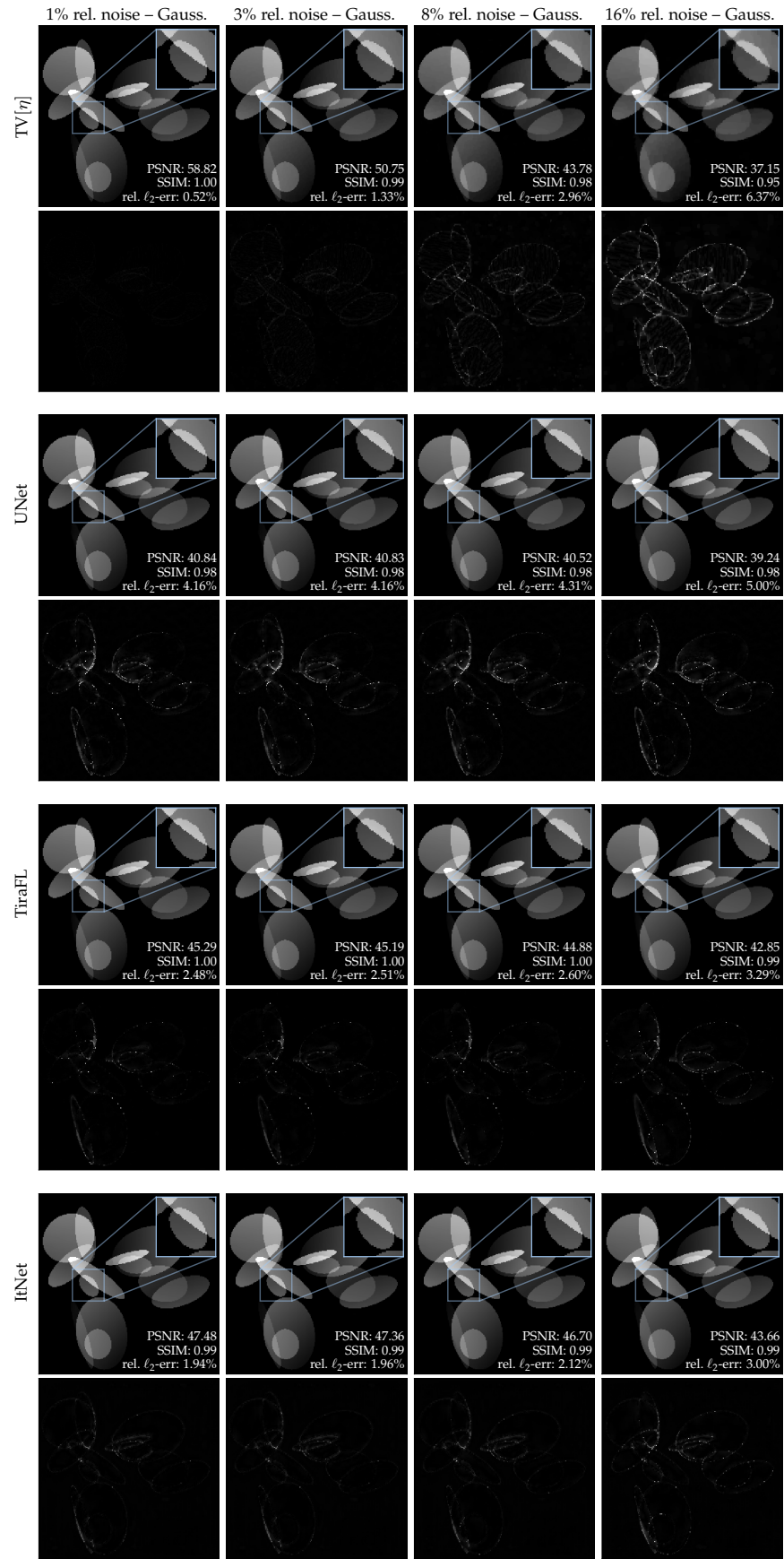


Fig. S6. **Scenario B1 – Fourier meas. with ellipses.** Individual reconstructions of the image from Fig. 7 under Gaussian noise. The reconstructed images are displayed in the window $[0, 0.9]$, which is also used for the computation of the PSNR and SSIM. The error plots shown below each reconstruction are displayed in the window $[0, 0.15]$. In favor of the more insightful noise level 16%, we have omitted the noiseless case.

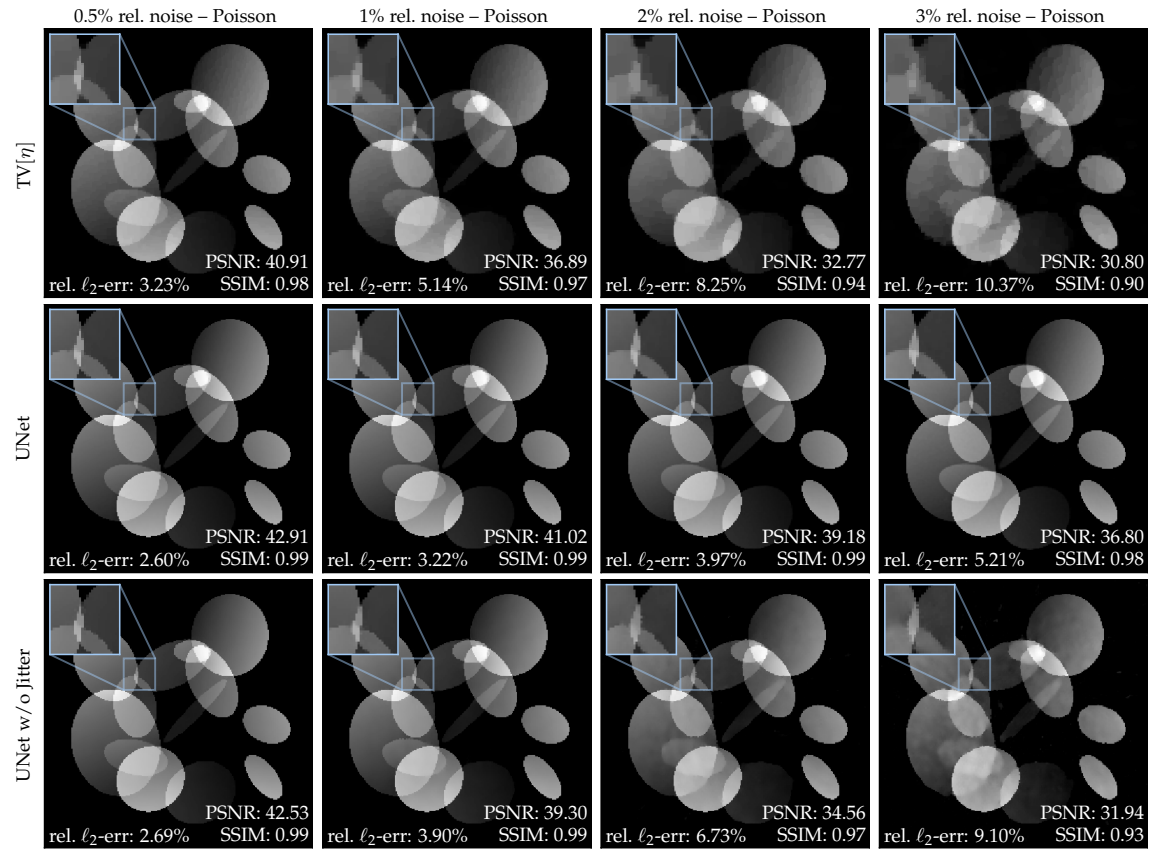


Fig. S7. **Scenario B2 – Radon meas. with ellipses.** Individual reconstructions of the image from Fig. 8 under Poisson noise. The reconstructed images are displayed in the window $[0, 1]$, which is also used for the computation of the PSNR and SSIM. In favor of the more insightful noise level 3%, we have omitted the noiseless case. The bottom row shows the corresponding reconstructions for a UNet that is trained without jittering; see also Section 5.1 on the inverse crime.

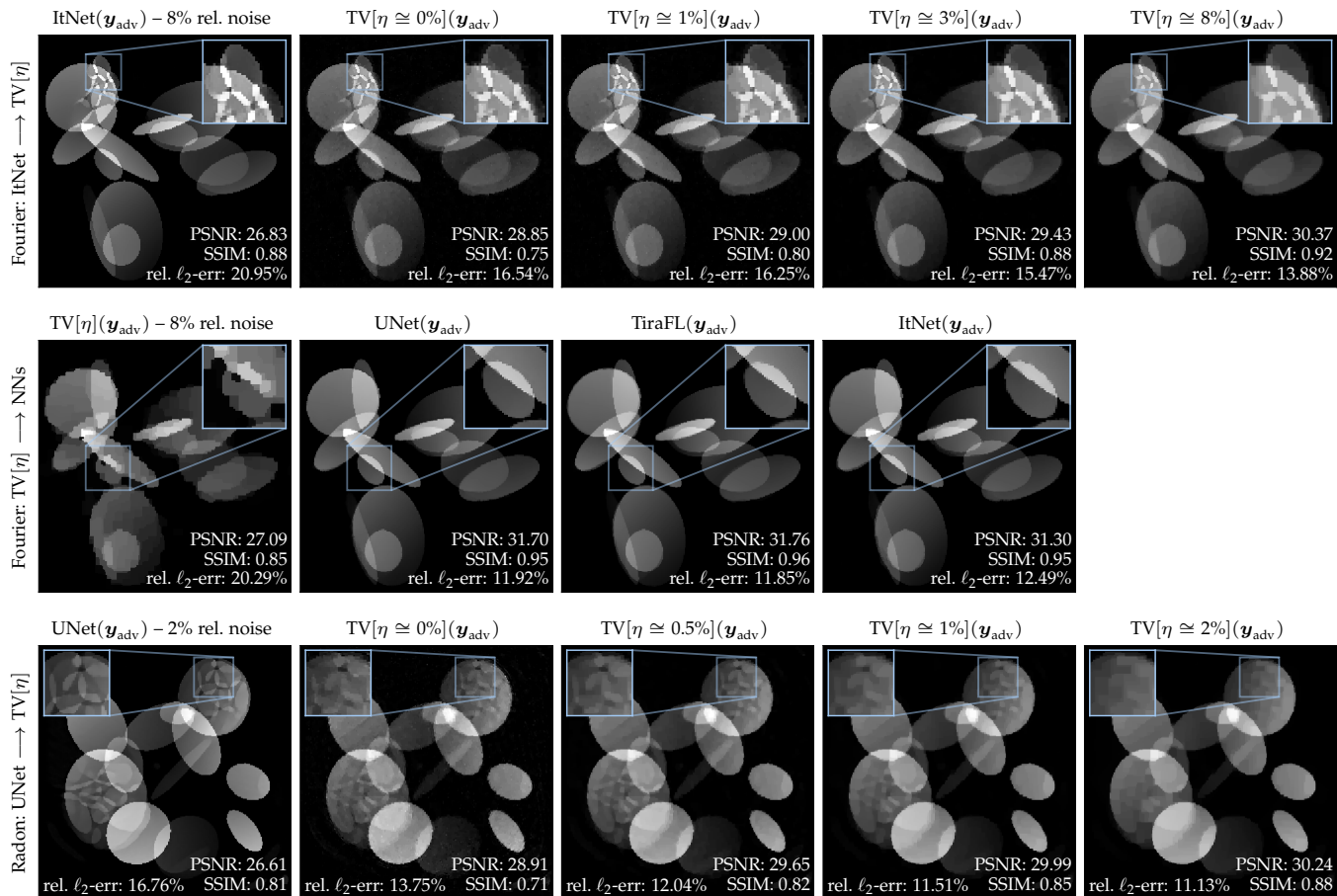


Fig. S8. **Case Study B – Transferability of perturbations.** This figure analyzes how adversarial noise transfers between TV minimization and NN-based solvers. The top row shows the recovery behavior of TV $[\eta]$ in the case of Fourier measurements when an adversarial perturbation \mathbf{y}_{adv} found for ItNet is used as input (cf. Fig. 7). Here, we also demonstrate the impact of the noise tuning parameter η , which controls the degree of regularization for TV minimization. The middle row presents the reverse experiment: an adversarial perturbation \mathbf{y}_{adv} found for TV $[\eta]$ is plugged into each considered NN. The bottom row is the analog of the top row in the case of Radon measurements (cf. Fig. 8).

S3 SUPPLEMENTARY RESULTS FOR CASE STUDY C (MRI ON REAL-WORLD DATA)

rel. noise – adversarial		0.0%	0.2%	0.5%	1.0%	1.5%	2.0%	2.5%
TV	rel. ℓ_2 -err. [%]	8.39±1.38	8.89±1.40	9.35±1.41	10.34±1.41	11.35±1.42	12.35±1.43	12.96±1.44
	PSNR	31.70±1.47	31.18±1.40	30.73±1.34	29.85±1.22	29.02±1.10	28.28±1.02	27.85±0.98
	SSIM	0.78±0.04	0.77±0.04	0.76±0.04	0.74±0.04	0.72±0.04	0.70±0.04	0.68±0.04
UNet	rel. ℓ_2 -err. [%]	8.18±1.27	8.29±1.27	8.41±1.27	8.67±1.26	8.99±1.25	9.38±1.22	9.84±1.18
	PSNR	31.90±1.38	31.79±1.37	31.66±1.35	31.38±1.30	31.06±1.23	30.69±1.15	30.26±1.06
	SSIM	0.80±0.04	0.80±0.03	0.79±0.03	0.79±0.03	0.78±0.03	0.78±0.03	0.77±0.03
UNetFL	rel. ℓ_2 -err. [%]	8.23±1.28	8.35±1.28	8.47±1.28	8.75±1.27	9.10±1.25	9.51±1.21	10.01±1.17
	PSNR	31.85±1.39	31.72±1.37	31.59±1.34	31.30±1.29	30.96±1.22	30.56±1.13	30.10±1.04
	SSIM	0.79±0.04	0.79±0.04	0.79±0.04	0.78±0.03	0.78±0.03	0.77±0.03	0.77±0.03
Tira	rel. ℓ_2 -err. [%]	7.97±1.26	8.10±1.26	8.24±1.26	8.58±1.25	9.00±1.21	9.52±1.17	10.16±1.12
	PSNR	32.13±1.41	31.99±1.39	31.84±1.36	31.48±1.30	31.05±1.19	30.54±1.09	29.97±0.97
	SSIM	0.80±0.03	0.80±0.03	0.80±0.03	0.79±0.03	0.79±0.03	0.78±0.03	0.77±0.03
TiraFL	rel. ℓ_2 -err. [%]	7.98±1.27	8.11±1.27	8.26±1.27	8.60±1.27	9.03±1.24	9.55±1.20	10.19±1.15
	PSNR	32.12±1.42	31.98±1.40	31.82±1.37	31.46±1.31	31.03±1.22	30.52±1.11	29.95±1.00
	SSIM	0.80±0.03	0.80±0.03	0.80±0.03	0.79±0.03	0.78±0.03	0.78±0.03	0.77±0.03
ItNet	rel. ℓ_2 -err. [%]	7.08±1.20	7.21±1.20	7.35±1.19	7.67±1.17	8.08±1.13	8.59±1.10	9.20±1.07
	PSNR	33.18±1.52	33.02±1.49	32.85±1.45	32.45±1.35	31.99±1.23	31.45±1.12	30.84±1.02
	SSIM	0.82±0.04	0.82±0.04	0.81±0.04	0.81±0.03	0.80±0.03	0.79±0.03	0.78±0.03

TABLE S7

Case Study C – fastMRI. A numerical representation of the results of Fig. 9(c), including the additional methods UNetFL and Tira. The best relative error/PSNR/SSIM per noise level is highlighted in bold.

rel. noise – Gaussian		0.0%	0.2%	0.5%	1.0%	1.5%	2.0%	2.5%
TV	rel. ℓ_2 -err. [%]	8.39±1.38	8.39±1.38	8.40±1.38	8.44±1.37	8.49±1.36	8.57±1.35	8.65±1.34
	PSNR	31.70±1.47	31.69±1.47	31.68±1.47	31.65±1.46	31.58±1.44	31.51±1.41	31.42±1.38
	SSIM	0.78±0.04	0.78±0.04	0.78±0.04	0.78±0.04	0.78±0.04	0.77±0.04	0.77±0.04
UNet	rel. ℓ_2 -err. [%]	8.18±1.27	8.18±1.27	8.18±1.27	8.20±1.26	8.22±1.26	8.24±1.26	8.27±1.26
	PSNR	31.90±1.38	31.90±1.38	31.90±1.38	31.89±1.38	31.86±1.37	31.84±1.37	31.80±1.36
	SSIM	0.80±0.04	0.80±0.04	0.80±0.04	0.80±0.04	0.80±0.04	0.80±0.04	0.79±0.04
UNetFL	rel. ℓ_2 -err. [%]	8.23±1.28	8.24±1.28	8.24±1.28	8.25±1.28	8.26±1.28	8.29±1.27	8.31±1.27
	PSNR	31.85±1.39	31.84±1.39	31.84±1.38	31.83±1.38	31.81±1.38	31.79±1.37	31.76±1.37
	SSIM	0.79±0.04	0.79±0.04	0.79±0.04	0.79±0.04	0.79±0.04	0.79±0.04	0.79±0.04
Tira	rel. ℓ_2 -err. [%]	7.97±1.26	7.97±1.26	7.98±1.26	7.99±1.26	8.01±1.26	8.04±1.25	8.07±1.25
	PSNR	32.13±1.41	32.13±1.41	32.13±1.41	32.11±1.40	32.09±1.40	32.05±1.39	32.02±1.38
	SSIM	0.80±0.03	0.80±0.03	0.80±0.03	0.80±0.03	0.80±0.03	0.80±0.03	0.80±0.03
TiraFL	rel. ℓ_2 -err. [%]	7.98±1.27	7.98±1.27	7.99±1.27	8.00±1.27	8.02±1.27	8.05±1.26	8.08±1.26
	PSNR	32.12±1.42	32.12±1.42	32.12±1.42	32.10±1.41	32.08±1.41	32.05±1.40	32.01±1.40
	SSIM	0.80±0.03	0.80±0.03	0.80±0.03	0.80±0.03	0.80±0.03	0.80±0.03	0.80±0.04
ItNet	rel. ℓ_2 -err. [%]	7.08±1.20	7.08±1.20	7.08±1.20	7.10±1.20	7.13±1.19	7.17±1.19	7.22±1.18
	PSNR	33.18±1.52	33.18±1.52	33.17±1.52	33.15±1.51	33.12±1.50	33.07±1.48	33.01±1.47
	SSIM	0.82±0.04	0.82±0.04	0.82±0.04	0.82±0.04	0.82±0.04	0.81±0.04	0.81±0.04

TABLE S8

Case Study C – fastMRI. A numerical representation of the results of Fig. 9(d), including the additional methods UNetFL and Tira. The best relative error/PSNR/SSIM per noise level is highlighted in bold.

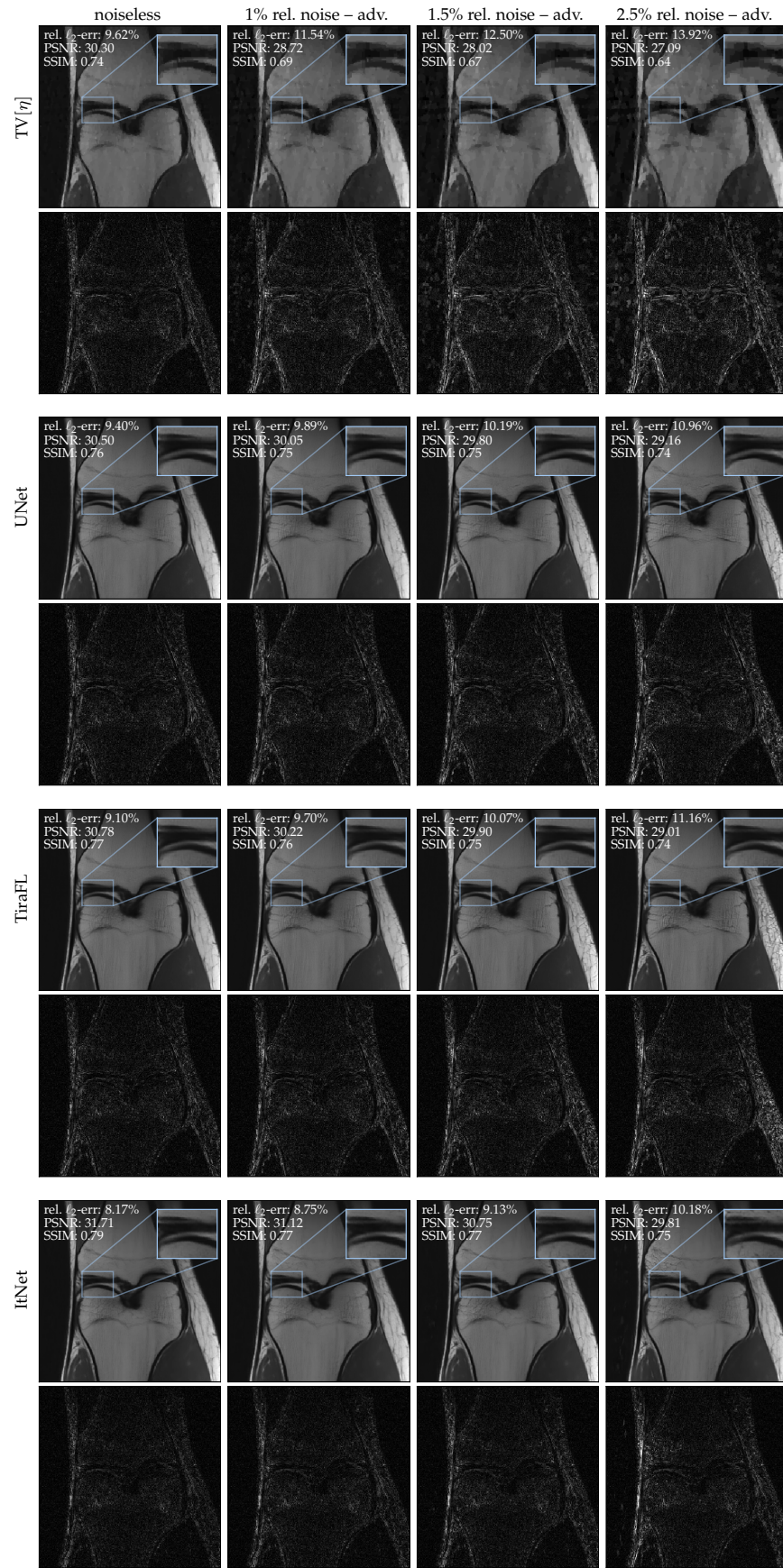


Fig. S9. **Case Study C – fastMRI**. Individual reconstructions of the image from Fig. 10 for different levels of adversarial noise. The reconstructed images are displayed in the window $[0.05, 4.50]$, which is also used for the computation of the PSNR and SSIM. The error plots shown below each reconstruction are displayed in the window $[0, 1.25]$.

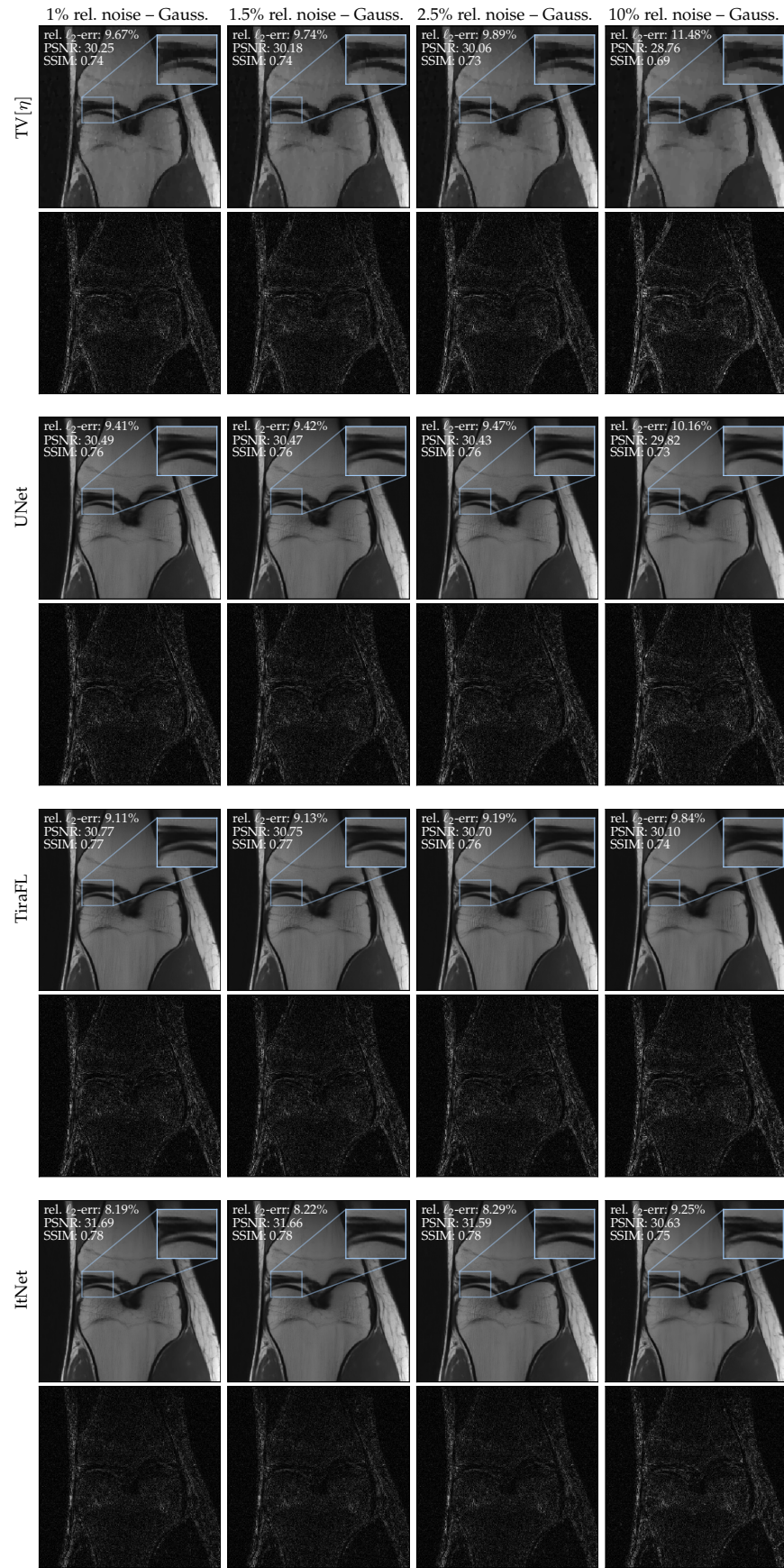


Fig. S10. **Case Study C – fastMRI**. Individual reconstructions of the image from Fig. 10 under Gaussian noise. The reconstructed images are displayed in the window $[0.05, 4.50]$, which is also used for the computation of the PSNR and SSIM. The error plots shown below each reconstruction are displayed in the window $[0, 1.25]$. In favor of the more insightful noise level 10%, we have omitted the noiseless case.

S4 SUPPLEMENTARY RESULTS FOR SECTION 5

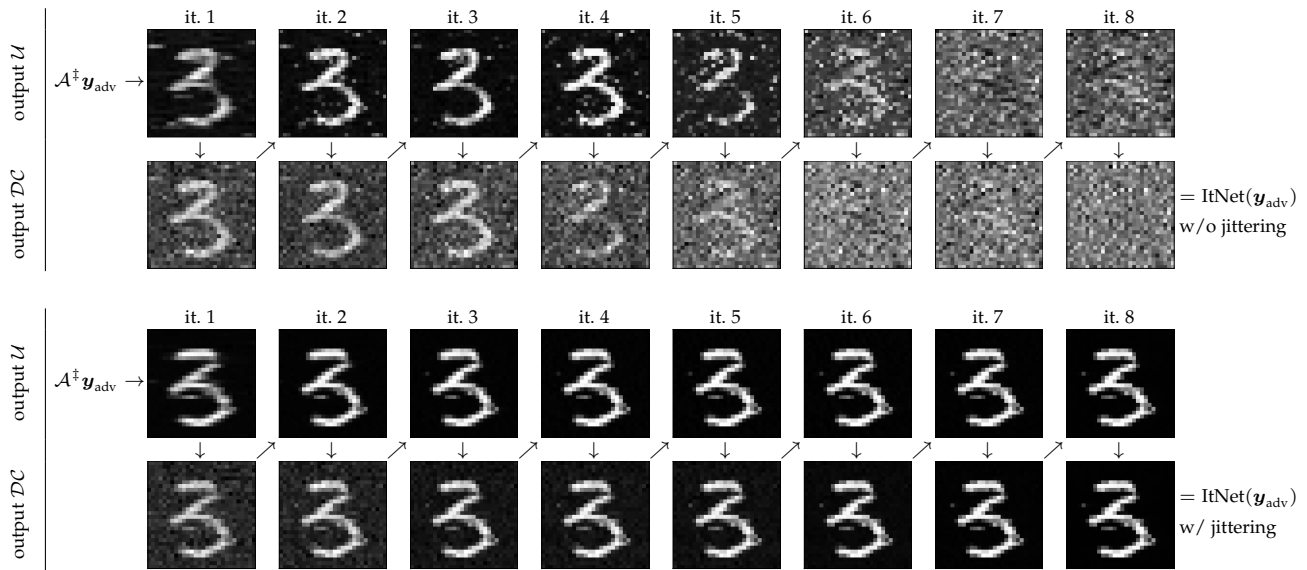


Fig. S11. **An inverse crime?** Intermediate steps performed by ItNet with and without jittering. The 20%-adversarial perturbations correspond to the individual reconstructions shown in Fig. 11.

S5 CHOICE OF HYPER-PARAMETERS

Table S9 summarizes all hyper-parameters concerning the considered network architectures. Table S10 shows the hyper-parameters selected for NN training. Table S11 contains relevant hyper-parameters for our adversarial attacks.

	Piecewise Constant	MNIST	Ellipses (Fourier)	Ellipses (Radon)	fastMRI (radial)	fastMRI (challenge)
UNet	inversion	Tikhonov (0.02)	\mathcal{A}^*	FBP (Hann-Filter)	\mathcal{A}^*	-
	U-Net levels channels per level	5 (64, 128, 256, 512, 1024)	Tikhonov (0.02) 5 (64, 128, 256, 512, 1024)	5 (36, 72, 144, 288, 576)	5 (24, 48, 96, 192, 384)	-
UNetFL	inversion	learned, init Tikh. (0.02)	learned	-	learned	-
	U-Net levels channels per level	5 (64, 128, 256, 512, 1024)	learned, init Tikh. (0.02) 5 (64, 128, 256, 512, 1024)	-	5 (24, 48, 96, 192, 384)	-
Tira	inversion	Tikhonov (0.02)	\mathcal{A}^*	-	\mathcal{A}^*	\mathcal{A}^*
	Tiramisu levels	5	5	-	5	5
	dense blocks per level	(5, 7, 9, 12, 15)	(5, 7, 9, 12, 15)	(5, 7, 9, 12, 15)	(5, 7, 9, 12, 15)	(6, 8, 10, 12, 14)
	initial channels	16	16	16	12	16
	channel growth rate	16	16	16	12	12
	bottleneck layers	25	25	20	18	20
TiraFL	inversion	learned, init Tikh. (0.02)	learned	-	learned	-
	Tiramisu levels	5	5	-	5	-
	dense blocks per level	(5, 7, 9, 12, 15)	(5, 7, 9, 12, 15)	(5, 7, 9, 12, 15)	(5, 7, 9, 12, 15)	-
	initial channels	16	16	16	12	-
	channel growth rate	16	16	16	12	-
	bottleneck layers	25	25	20	18	-
ItNet	inversion	Tikhonov (0.02)	\mathcal{A}^*	-	\mathcal{A}^*	-
	U-Net levels	5	5	-	5	-
	channels per level iterations	(64, 128, 256, 512, 1024) 8	(64, 128, 256, 512, 1024) 8	(32, 64, 128, 256, 512) 8	(24, 96, 192, 384) 8	-
ConvNet	convolutional layers channels per layer	-	4 (32, 32, 64, 64)	-	-	-
	fully connected layers features per layer	-	3 (200, 200, 10)	-	-	-
	remarks	-	dropout ($p = 0.5$) between fc layers	-	-	-

TABLE S9

Detailed description of hyper-parameters for all considered NN architectures. The convolution kernel sizes are 3 or 3×3 , the max-pooling sizes are 2 or 2×2 , and the activation functions are rectified linear units (ReLUs) for all networks.

	Piecewise Constant	MNIST	Ellipses (Fourier)	Ellipses (Radon)	fastMRI (radial)	fastMRI (challenge)
UNet	(200, 75) (40, 40) ($8 \cdot 10^{-5}, 5 \cdot 10^{-5}$) ($5 \cdot 10^{-3}, 5 \cdot 10^{-3}$) ($1 \cdot 10^{-5}, 1 \cdot 10^{-5}$) (1, 200) 2	(200, 20) (40, 40) ($8 \cdot 10^{-5}, 5 \cdot 10^{-5}$) ($5 \cdot 10^{-3}, 5 \cdot 10^{-3}$) ($1 \cdot 10^{-5}, 1 \cdot 10^{-5}$) (1, 200) 4	(100, 10) (40, 40) ($2 \cdot 10^{-4}, 5 \cdot 10^{-5}$) ($1 \cdot 10^{-5}, 1 \cdot 10^{-5}$) ($1 \cdot 10^{-5}, 1 \cdot 10^{-5}$) (1, 200) 10	(155, 5) (40, 40) ($2 \cdot 10^{-4}, 5 \cdot 10^{-5}$) ($5 \cdot 10^{-4}, 1 \cdot 10^{-4}$) ($1 \cdot 10^{-4}, 1 \cdot 10^{-4}$) (1, 200) 500	(100) (40) ($2 \cdot 10^{-4}$) ($1 \cdot 10^{-5}$) ($1 \cdot 10^{-4}$) (1) 150	-
UNetFL	(250, 50) (40, 40) ($2 \cdot 10^{-4}, 5 \cdot 10^{-5}$) ($5 \cdot 10^{-4}, 5 \cdot 10^{-4}$) ($1 \cdot 10^{-5}, 1 \cdot 10^{-5}$) (1, 200) 2	(200, 75) (40, 40) ($2 \cdot 10^{-4}, 5 \cdot 10^{-5}$) ($5 \cdot 10^{-3}, 5 \cdot 10^{-3}$) ($1 \cdot 10^{-5}, 1 \cdot 10^{-5}$) (1, 200) 4	(100, 10) (40, 40) ($2 \cdot 10^{-4}, 5 \cdot 10^{-5}$) ($1 \cdot 10^{-3}, 5 \cdot 10^{-4}$) ($1 \cdot 10^{-5}, 1 \cdot 10^{-5}$) (1, 200) 10	-	(45) (40) ($8 \cdot 10^{-5}$) ($1 \cdot 10^{-5}$) ($1 \cdot 10^{-4}$) (1) 150 init \mathcal{U} from UNet	-
Tira	(200, 50) (40, 40) ($2 \cdot 10^{-4}, 5 \cdot 10^{-5}$) ($5 \cdot 10^{-4}, 5 \cdot 10^{-4}$) ($1 \cdot 10^{-5}, 1 \cdot 10^{-5}$) (1, 200) 2	(200, 50) (40, 40) ($2 \cdot 10^{-4}, 5 \cdot 10^{-5}$) ($1 \cdot 10^{-5}, 1 \cdot 10^{-5}$) ($1 \cdot 10^{-5}, 1 \cdot 10^{-5}$) (1, 100) 4	(50, 10) (8, 5) ($2 \cdot 10^{-4}, 5 \cdot 10^{-5}$) ($1 \cdot 10^{-5}, 1 \cdot 10^{-5}$) ($1 \cdot 10^{-5}, 1 \cdot 10^{-5}$) (1, 200) 10	-	(50, 16) (6, 6) ($8 \cdot 10^{-5}, 5 \cdot 10^{-5}$) ($1 \cdot 10^{-6}, 1 \cdot 10^{-6}$) ($2 \cdot 10^{-4}, 2 \cdot 10^{-4}$) (1, 1) 150	(40) (4) ($1 \cdot 10^{-4}$) ($1 \cdot 10^{-5}$) ($1 \cdot 10^{-4}$) (1) 10
TiraFL	(200, 50) (40, 40) ($2 \cdot 10^{-4}, 5 \cdot 10^{-5}$) ($5 \cdot 10^{-4}, 5 \cdot 10^{-4}$) ($1 \cdot 10^{-5}, 1 \cdot 10^{-5}$) (1, 200) 2	(200, 50) (40, 40) ($2 \cdot 10^{-4}, 5 \cdot 10^{-5}$) ($5 \cdot 10^{-4}, 5 \cdot 10^{-4}$) ($1 \cdot 10^{-5}, 1 \cdot 10^{-5}$) (1, 200) 4	(30, 7) (10, 5) ($2 \cdot 10^{-4}, 5 \cdot 10^{-5}$) ($1 \cdot 10^{-4}, 1 \cdot 10^{-5}$) ($2 \cdot 10^{-4}, 1 \cdot 10^{-5}$) (1, 200) 10	-	(50, 20) (6, 6) ($8 \cdot 10^{-5}, 5 \cdot 10^{-5}$) ($1 \cdot 10^{-6}, 1 \cdot 10^{-6}$) ($2 \cdot 10^{-4}, 2 \cdot 10^{-4}$) (1, 1) 150	-
ItNet	(100, 5) (40, 40) ($5 \cdot 10^{-5}, 2 \cdot 10^{-5}$) ($5 \cdot 10^{-4}, 5 \cdot 10^{-4}$) ($1 \cdot 10^{-5}, 1 \cdot 10^{-5}$) (1, 200) 2	(100, 10) (40, 40) ($8 \cdot 10^{-5}, 5 \cdot 10^{-5}$) ($1 \cdot 10^{-3}, 1 \cdot 10^{-3}$) ($1 \cdot 10^{-5}, 1 \cdot 10^{-5}$) (1, 200) 4	(35, 6) (15, 15) ($5 \cdot 10^{-5}, 5 \cdot 10^{-5}$) ($1 \cdot 10^{-4}, 1 \cdot 10^{-4}$) ($2 \cdot 10^{-4}, 2 \cdot 10^{-4}$) (1, 200) 10	-	(15, 8) (10, 10) ($5 \cdot 10^{-5}, 5 \cdot 10^{-5}$) ($1 \cdot 10^{-6}, 1 \cdot 10^{-6}$) ($1 \cdot 10^{-4}, 1 \cdot 10^{-4}$) (1, 1) 150	-
ConvNet	-	(20, 10) (40, 40) ($2 \cdot 10^{-4}, 5 \cdot 10^{-5}$) ($1 \cdot 10^{-5}, 1 \cdot 10^{-5}$) ($1 \cdot 10^{-5}, 1 \cdot 10^{-5}$) (1, 1) no	-	-	-	-

TABLE S10

Detailed description of hyper-parameters for all NN trainings. All networks are trained in 1 or 2 phases and respective parameters are shown per training phase. Default parameters for the Adam optimizer are used except for the ϵ parameter, which is reported for all networks and training phases.

	Piecewise Constant	MNIST	Ellipses (Fourier)	Ellipses (Radon)	fastMRI (radial)	fastMRI (challenge)
TV[η]	ADMM iterations (rec.)	50000	5000	200	5000	5000
	ADMM iterations (init.)	50000	5000	200	5000	5000
	ADMM iterations (gradient)	2000	200	20	200	150
	Adam iterations	30	250	15	250	50
	step size	$5 \cdot 10^0$	$5 \cdot 10^0$	$5 \cdot 10^0$	$5 \cdot 10^0$	$5 \cdot 10^0$
	random initializations	200	100	6	6	6
Net	Adam iterations	100	1000	500	250	200
	step size	$5 \cdot 10^0$	$5 \cdot 10^0$	$5 \cdot 10^0$	$5 \cdot 10^0$	$5 \cdot 10^0$
	random initializations	200	6	6	6	6
ConvNet \circ TV[η]	ADMM iterations (rec.)	-	-	-	-	-
	ADMM iterations (init.)	50000	50000	-	-	-
	ADMM iterations (gradient)	2000	2000	-	-	-
	Adam iterations	100	100	-	-	-
	step size	$5 \cdot 10^{-1}$	$5 \cdot 10^{-1}$	-	-	-
	random initializations	100	100	-	-	-
ConvNet \circ Net	Adam iterations	-	-	-	-	-
	step size	$5 \cdot 10^{-1}$	$5 \cdot 10^{-1}$	-	-	-
	random initializations	100	100	-	-	-

TABLE S11

Detailed description of hyper-parameters for finding adversarial perturbations. The parameters reported for Net apply equally to all network types.