# Denoising and Regularization via Exploiting the Structural Bias of Convolutional Generators

Reinhard Heckel\* and Mahdi Soltanolkotabi<sup>†</sup>

\*Dept. of Electrical and Computer Engineering, Technical University of Munich †Dept. of Electrical and Computer Engineering, University of Southern California

October 31, 2019

#### Abstract

Convolutional Neural Networks (CNNs) have emerged as highly successful tools for image generation, recovery, and restoration. This success is often attributed to large amounts of training data. However, recent experimental findings challenge this view and instead suggest that a major contributing factor to this success is that convolutional networks impose strong prior assumptions about natural images. A surprising experiment that highlights this architectural bias towards natural images is that one can remove noise and corruptions from a natural image without using any training data, by simply fitting (via gradient descent) a randomly initialized, over-parameterized convolutional generator to the single corrupted image. While this over-parameterized network can fit the corrupted image perfectly, surprisingly after a few iterations of gradient descent one obtains the uncorrupted image. This intriguing phenomena enables state-of-the-art CNN-based denoising and regularization of linear inverse problems such as compressive sensing. In this paper we take a step towards demystifying this experimental phenomena by attributing this effect to particular architectural choices of convolutional networks, namely convolutions with fixed interpolating filters. We then formally characterize the dynamics of fitting a two layer convolutional generator to a noisy signal and prove that earlystopped gradient descent denoises/regularizes. This results relies on showing that convolutional generators fit the structured part of an image significantly faster than the corrupted portion.

## 1 Introduction

Convolutional neural networks are extremely popular for image generation. The majority of image generating networks is convolutional, ranging from Deep Convolutional Generative Adversarial Networks (DC-GANs) [Rad+15] to the U-Net [Ron+15]. It is well known that convolutional neural networks incorporate implicit assumption about the signals they generate, such as pixels that are close being related. This makes them particularly well suited for representing sets of images or modeling distributions of images. It is less known, however, that those prior assumptions build into the architecture are so strong that convolutional neural networks are useful even without ever being exposed to training data.

The latter was first shown in the Deep Image Prior (DIP) paper [Uly+18]. Ulyanov et al. [Uly+18] observed that when 'training' an standard convolutional auto-encoder such as the popular Unet [Ron+15] on a single noisy image and regularizing by early stopping, the network performs image restoration such as denoising with state-of-the-art performance. This is based on the empirical observation that un-trained convolutional auto-decoders fit a single natural image faster when optimized with gradient descent than pure noise. A more recent paper [HH19] proposed a much simpler image generating network, termed the deep decoder. This network can be seen as the relevant part of a convolutional generator architecture to function as an image prior, and can

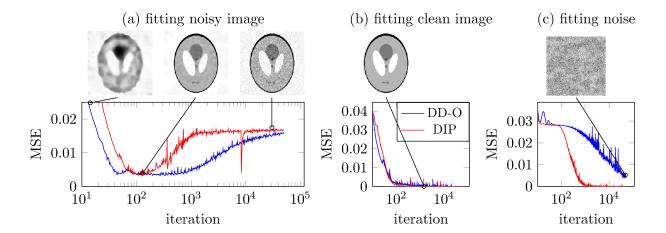


Figure 1: Fitting an over-parameterized Deep Decoder (DD-O) and the deep image prior (DIP) to a (a) noisy image, (b) clean image, and (c) pure noise. Here, MSE denotes Mean Square Error of the network output with respect to the clean image in (a) and fitted images in (b) and (c). While the network can fit the noise due to over-parameterization, it fits natural images in significantly fewer iterations than noise. Hence, when fitting a noisy image, the image component is fitted faster than the noise component which enables denoising via early stopping.

be obtained from a standard convolutional autoencoder by removing the encoder, the skip connections, and the trainable convolutional filters of spacial extent larger than one. The deep decoder does not use learned or trainable convolutional filters like conventional convolutional networks do, and instead only uses convolutions with fixed convolutional kernels to generate an image.

In this paper, we study a simple untrained convolutional network that only consists of convolutionallike operations, such as the deep image prior and the deep decoder. We consider the overparameterized regime where the network has sufficiently many parameters to represent an arbitrary image (including noise) perfectly and show that:

Fitting convolutional generators via early stopped gradient descent provably denoises "natural" images.

To prove this denoising capability we characterize how the network architecture governs the dynamics of fitting over-parameterized networks to a single (noisy) image. In particular we prove:

Convolutional generators optimized with gradient descent fit natural images faster than noise.

Here, by complex images, we mean unstructured images that consist of a large number of edges or variations such as noise.

We depict this phenomena in Figure 1 where we fit a randomly initialized over-parameterized convolutional generator to a signal via running gradient descent on the objective  $\mathcal{L}(\mathbf{C}) = \|G(\mathbf{C}) - \mathbf{y}\|_2^2$ . Here,  $G(\mathbf{C})$  is the convolutional generator with weight parameters  $\mathbf{C}$ , and  $\mathbf{y}$  is either a noisy image, a clean image, or noise. This experiment demonstrates that over-parameterized convolutional networks fit a natural image (Figure 1b) much faster than noise (Figure 1c). Thus, when fitting the noisy image (Figure 1a), early stopping the optimization enables image denoising. Interestingly, this effect is so strong that it gives state-of-the-art denoising performance, outperforming



Figure 2: Denoising with BM3D and various convolutional generators. The relative ranking of algorithms in this picture is representative and maintained on a large test set of images. DIP and DD-O are over-parameterized convolutional generators and with early stopping outperform the BM3D algorithm, the next best method that does not require training data.



Figure 3: The 1st, 2nt, 6th, and 21th trigonometric basis functions in dimension n = 300.

the BM3D algorithm [Dab+07], which is the next-best method that requires no training data (see Figure 2). Beyond denoising, this effect also enables significant improvements in regularizing a variety of inverse problems such as compressive sensing [Vee+18; Hec19].

#### 1.1 Contributions and overview of results

In this paper we take a step towards demystifying why fitting a convolutional generator via early stopped gradient descent leads to such surprising denoising and regularization capabilities.

- We first show experimentally that this denoising phenomena is primarily attributed to convolutions with fixed interpolation kernels, typically implicitly implemented by bi-linear upsampling operations in the convolutional network.
- We then study over-parameterized convolutional networks (with fixed convolutional filters) theoretically. Specifically, we show that fitting such a network via early stopped gradient descent to a signal provably denoises it. Specifically, let  $\mathbf{x} \in \mathbb{R}^n$  be a smooth signal that can be represented exactly as a linear combination of the p orthogonal trigonometric functions of lowest frequency (defined in equation (5), see Figure 3 for a depiction), and suppose our goal is to obtain an estimate of this signal from a noisy signal  $\mathbf{y} = \mathbf{x} + \mathbf{z}$  where  $\mathbf{z} \sim \mathcal{N}(0, \frac{\varsigma^2}{n}\mathbf{I})$ . Note that natural images are often smooth. Let  $\hat{\mathbf{y}} = G(\mathbf{C}^t)$  be the estimate obtained from early stopping the fitting of a two-layer convolutional generator to  $\mathbf{y}$ . We prove this estimate achieves the denoising rate

$$\|\hat{\mathbf{y}} - \mathbf{x}\|_2^2 \le c \frac{p}{n} \varsigma^2,$$

with c a constant. We note that this rate is optimal up to a constant factor.

• Our denoising result follows from a detailed analysis of gradient descent applied to fitting a two-layer convolutional generator  $G(\mathbf{C})$  with fixed convolutional filters to a noisy signal  $\mathbf{y} = \mathbf{x} + \mathbf{z}$  with  $\mathbf{x}$  representing the signal and  $\mathbf{z}$  the noise. Specifically, let  $\tilde{\mathbf{x}} \in \mathbb{R}^n$  and  $\tilde{\mathbf{z}} \in \mathbb{R}^n$  be the coefficients of the signal and noise in terms of trigonometric basis functions  $\mathbf{w}_1, \ldots, \mathbf{w}_n \in \mathbb{R}^n$  (precisely defined later, see Figure 3 for a depiction). We show that there is a dual filter  $\sigma \in \mathbb{R}^n$ , depending only on the convolutional filter used, whose entries typically obey  $\sigma_1 \geq \sigma_2 \geq \ldots \geq \sigma_n > 0$  and can be thought of as weights associated with each of those basis function. These weights in turn determine the speed at which the associated components of the signal are fitted. Specifically, we show that the dynamics of gradient descent are approximately given by

$$G(\mathbf{C}_{\tau}) - \mathbf{x} \approx \underbrace{\sum_{i=1}^{n} \mathbf{w}_{i} \widetilde{x}_{i} (1 - \eta \sigma_{i}^{2})^{\tau}}_{\text{error in fitting signal}} + \underbrace{\sum_{i=1}^{n} \mathbf{w}_{i} \widetilde{z}_{i} ((1 - \eta \sigma_{i}^{2})^{\tau} - 1)}_{\text{fit of noise}}.$$

The filters commonly used are typically such that the dual filter decays rather quickly, implying that low-frequency components in the trigonometric expansion are fitted significantly faster than high frequency ones. So when the signal mostly consists of low-frequency components, we can choose an early stopping time such that the error in fitting the signal is very low, and thus the signal part is well described, whereas at the same time only a small part of the noise, specifically the part aligned with the low-frequency components has been fitted.

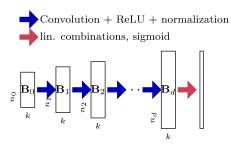
# 2 Convolutional generators

A convolutional generator maps an input tensor  $\mathbf{B}_0$  to an image only using upsampling and convolutional operations, followed by channel normalization (a special case of batch normalization) and applications of non-linearities, see Figure 4. All previously mentioned convolutional generator networks [Rad+15; Ron+15] including the networks considered in the DIP paper [Uly+18] primarily consist of those operations.

For motivating the architecture of the convolutional generators studied in this paper, we first demonstrate in Section 2.1 that convolutions with fixed interpolation filters are critical to the denoising performance with early stopping. Specifically, we empirically show that convolutions with *fixed* convolutional kernels are critical for convolutional generators to fit natural images faster than noise. Finally, in Section 2.2 we formally introduce the class of convolutional generators studied in this paper, and in Section 2.3 we introduce a minimal convolutional architecture which is the focus of our theoretical results.

#### 2.1 The importance of fixed convolutional filters

Convolutions with fixed convolutional kernels are critical for denoising with early stopping, because they are critical for the phenomena that natural images are fitted significantly faster than noise. To see this, consider the experiment in Figure 4 in which we fit an image and noise by minimizing the least-squares loss via gradient descent with i) a convolutional generator with only fixed convolutional filters (see Section 2.2 below for a precise description) and ii) a conventional convolutional generator with trainable convolutional filters (essentially the architecture from the popular DC-GAN generators, see Appendix A for details and additional numerical evidence). Figure 4 shows



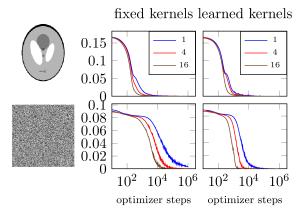


Figure 4: **Left panel:** Convolutional generators. The output is generated through repeated convolutional layers, channel normalization, and applying ReLU non-linearities. **Right panel:** Fitting the phantom MRI and noise with different architectures of depth d=5, for different number of over-parameterization factors (1,4, and 16). Gradient descent on convolutional generators involving fixed convolutional matrixes fit an image significantly faster than noise.

that the convolutional network with *fixed filters* fits the natural image much faster than noise, whereas the network with learned convolutional filters, only fits it slightly faster, and this effect vanishes as the network becomes highly overparameterized. Thus, fixed convolutional filters enable un-trained convolutional networks to function as highly effective image priors. We note that the upsampling operation present in most architectures implicitly incorporates a convolution with a *fixed* convolutional (interpolation) filter.

#### 2.2 Architecture of convolutional generator with fixed convolutions

In this section, we describe the architecture of a convolutional network with fixed convolutional operators only (i.e., the deep decoder [HH19]). In this architecture, the channels in the (i + 1)-th layer are given by

$$\mathbf{B}_{i+1} = \operatorname{cn}(\operatorname{ReLU}(\mathbf{U}_i \mathbf{B}_i \mathbf{C}_i)), \quad i = 0, \dots, d-1,$$

and finally, the output of the d-layer network is formed as

$$\mathbf{x} = \mathbf{B}_d \mathbf{C}_{d+1}$$
.

Here, the coefficient matrices  $\mathbf{C}_i \in \mathbb{R}^{k \times k}$  and  $\mathbf{C}_{d+1} \in \mathbb{R}^{k \times k_{\text{out}}}$  contain the weights of the network. The number of channels, k, determines the number of weight parameters of the network, given by  $dk^2 + k_{\text{out}}k$ . Each column of the tensor  $\mathbf{B}_i\mathbf{C}_i \in \mathbb{R}^{n_i \times k}$  is formed by taking linear combinations of the channels of the tensor  $\mathbf{B}_i$  in a way that is consistent across all pixels, and the ReLU activation function is given by  $\text{ReLU}(t) = \max(0, t)$ . Then,  $\text{cn}(\cdot)$  performs the channel normalization operation which normalizes each channel individually and can be viewed as a special case of the popular batch normalization operation [IS15].

The operator  $\mathbf{U}_i \in \mathbb{R}^{n_{i+1} \times n_i}$  is a tensor implementing an upsampling and most importantly a convolution operation with a fixed kernel. This fixed kernel was chosen in all experiments above as a triangular kernel so that  $\mathbf{U}$  performs bi-linear 2x upsampling (this is the standard implementation

in the popular packages pytorch and tensorflow). As mentioned earlier this convolution with a *fixed* kernel is critical for fitting natural images faster than complex ones.

## 2.3 Two layer convolutional generator studied theoretically in this paper

The simplest model to study the denoising capability of convolutional generators and the phenomena that a natural image is fitted faster than a complex one theoretically is a network with only one hidden layers and one output channel i.e.,  $G(\mathbf{C}) \in \mathbb{R}^n$ . Then, the generator becomes

$$G(\mathbf{C}) = \text{ReLU}(\mathbf{UB}_1\mathbf{C}_1)\mathbf{c}_2,$$

where  $\mathbf{U} \in \mathbb{R}^{n \times n}$  is a circulant matrix that implements a convolution with a filter  $\mathbf{u}$ . In this paper we consider the over-parameterized regime where  $k \geq n$ . In this regime, using that the input  $\mathbf{B}_1$  is random, with probability one, the matrix  $\mathbf{B}_1$  has full column rank and thus spans  $\mathbb{R}^n$ . It follows that optimizing over  $\mathbf{C}_1$  and  $\mathbf{c}_2$  is equivalent to optimizing over the parameter  $\mathbf{C} \in \mathbb{R}^{n \times k}$  in

$$G(\mathbf{C}) = \text{ReLU}(\mathbf{UC})\mathbf{v},$$
 (1)

where  $\mathbf{v} = [1, \dots, 1, -1, \dots, -1]/\sqrt{k}$  is fixed and  $\mathbf{C} \in \mathbb{R}^{n \times k}$  is the new coefficient matrix we optimize over. Figure 11 in the appendix shows that even this simple two-layer convolutional network fits a simple image faster than noise. This is the simplest model in which the phenomena that a convolutional networks fits structure faster than noise can reliably observed. As a consequence, the dynamics of training the model (1) are the focus of the remainder of this paper.

# 3 Warmup: Dynamics of gradient descent on least squares

As a prelude for studying the dynamics of fitting convolutional generators via a *non-linear* least square problem, we study the dynamics of gradient descent applied to a *linear* least squares problem. We demonstrate how early stopping can lead to denoising capabilities even with a simple linear model. We consider a least-squares problem of the form

$$\mathcal{L}(\mathbf{c}) = \frac{1}{2} \|\mathbf{y} - \mathbf{J}\mathbf{c}\|_2^2,$$

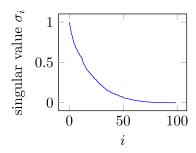
and study gradient descent with a constant step size  $\eta$  starting at  $\mathbf{c}_0 = \mathbf{0}$ . The updates are given by

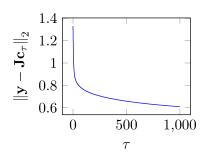
$$\mathbf{c}_{\tau+1} = \mathbf{c}_{\tau} - \eta \nabla \mathcal{L}(\mathbf{c}_{\tau}), \quad \nabla \mathcal{L}(\mathbf{c}) = \mathbf{J}^T (\mathbf{J}\mathbf{c} - \mathbf{y}).$$

The following simple proposition characterizes the trajectory of gradient descent.

**Proposition 1.** Let  $\mathbf{J} \in \mathbb{R}^{n \times m}$  be a matrix with left singular vectors  $\mathbf{w}_1, \dots, \mathbf{w}_n \in \mathbb{R}^n$  and corresponding singular values  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n$ . Then the residual after  $\tau$  steps,  $\mathbf{r}_{\tau} = \mathbf{y} - \mathbf{J}\mathbf{c}_{\tau}$ , of gradient descent starting at  $\mathbf{c}_0 = 0$  is

$$\mathbf{r}_{\tau} = \sum_{i=1}^{n} \mathbf{w}_{i} \left\langle \mathbf{w}_{i}, \mathbf{y} \right\rangle (1 - \eta \sigma_{i}^{2})^{\tau}.$$





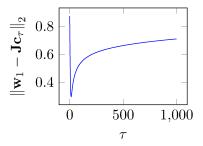


Figure 5: Gradient descent on the least squares problem of minimizing  $\|\mathbf{y} - \mathbf{J}\mathbf{c}_{\tau}\|^2$ , where  $\mathbf{J} \in \mathbb{R}^{100 \times 100}$  has decaying singular values (left panel) and the observation is the sum of a signal component, equal to the leading singular vector  $\mathbf{w}_1$  of  $\mathbf{J}$ , and a noisy component  $\mathbf{z} \sim \mathcal{N}(0, (1/n)\mathbf{I})$ , i.e.,  $\mathbf{y} = \mathbf{w}_1 + \mathbf{z}$ . The signal component  $\mathbf{w}_1$  is fitted significantly faster than the other components (right panel), thus early stopping enables denoising.

Suppose that the signal  $\mathbf{y}$  lies in the column span of  $\mathbf{J}$ , and that the stepsize is chosen sufficiently small (i.e.,  $\eta \leq 1/\|\mathbf{J}\|^2$ ). Then, by Proposition 1, gradient descent converges to a zero-loss solution and thus fits the signal perfectly. More importantly, gradient descent fits the components of  $\mathbf{y}$  corresponding to large singular values faster than it fits the components corresponding to small singular values.

To explicitly show how this observation enables regularization via early stopped gradient descent, suppose our goal is to find a good estimate of a signal x from a noisy observation

$$\mathbf{v} = \mathbf{x} + \mathbf{z}$$
.

where the signal  $\mathbf{x}$  lies in a signal subspace that is spanned by the p leading left-singular vectors of  $\mathbf{J}$ . Then, by Proposition 1, the signal estimate after  $\tau$  iterations,  $\mathbf{J}\mathbf{c}_{\tau}$ , obeys

$$\|\mathbf{J}\mathbf{c}_{\tau} - \mathbf{x}\|_{2} \le (1 - \eta\sigma_{p}^{2})^{\tau} \|\mathbf{x}\|_{2} + E(\mathbf{z}), \quad E(\mathbf{z}) \coloneqq \sqrt{\sum_{i=1}^{n} ((1 - \eta\sigma_{i}^{2})^{\tau} - 1)^{2} \langle \mathbf{w}_{i}, \mathbf{z} \rangle^{2}}.$$
 (2)

Thus, after a few iterations most of the signal has been fitted (i.e.,  $(1-\eta\sigma_p)^{\tau}$  is small). Furthermore, if we assume that the ratio  $\sigma_p/\sigma_{p+1}$  is sufficiently large so that the spread between the two singular values separating the signal subspace from the rest is sufficiently large, most of the noise outside the signal subspace has not been fitted (i.e.,  $((1-\eta\sigma_i^2)^{\tau}-1)^2\approx 0$  for  $i=p+1,\ldots,n$ ).

In particular, suppose the noise vector has a Gaussian distribution given by  $\mathbf{z} \sim \mathcal{N}(0, \frac{\varsigma^2}{n}\mathbf{I})$ . Then  $E(\mathbf{z}) \approx \varsigma \sqrt{\frac{p}{n}}$  so that after order  $\tau = \log(\epsilon)/\log(1 - \eta \sigma_p^2)$  iterations, with high probability,

$$\|\mathbf{J}\mathbf{c}_{\tau} - \mathbf{x}\|_{2} \le \epsilon \|\mathbf{x}\|_{2} + c\varsigma\sqrt{\frac{p}{n}}.$$

This demonstrates, that provided the signal lies in a subspace spanned by the leading singular vectors, early stoped gradient descent reaches the optimal denoising rate of  $\varsigma \sqrt{p/n}$  after a few iterations. See Figure 5 for a numerical example demonstrating this phenomena.

## 4 Dynamics of gradient descent on convolutional generators

We are now ready to study the implicit bias of gradient descent towards natural/structured images theoretically. Consider a two-layer decoder network (introduced in Section 2.3) of the form

$$G(\mathbf{C}) = \text{ReLU}(\mathbf{UC})\mathbf{v},$$

where  $\mathbf{v} = [1, \dots, 1, -1, \dots, -1]/\sqrt{k}$ , and with weight parameter  $\mathbf{C} \in \mathbb{R}^{n \times k}$ , and recall that  $\mathbf{U}$  is a circulant matrix implementing a convolution with a kernel  $\mathbf{u} \in \mathbb{R}^n$ . We consider minimizing the non-linear least squares objective

$$\mathcal{L}(\mathbf{C}) = \frac{1}{2} \|\mathbf{y} - G(\mathbf{C})\|_2^2$$
(3)

with (early-stopped gradient) descent with a constant stepsize  $\eta$  starting at a random initialization  $\mathbf{C}_0$  of the weights. The iterates are given by

$$\mathbf{C}_{\tau+1} = \mathbf{C}_{\tau} - \eta \nabla \mathcal{L}(\mathbf{C}_{\tau}). \tag{4}$$

In our warmup section on linear least squares we saw that the singular vectors and values of the matrix  $\mathbf{J}$  determine the speed at which different components of the noisy signal  $\mathbf{y}$  are fitted by gradient descent. The main insight that enables us to extend this intuition to the nonlinear case is that the role of the feature matrix  $\mathbf{J}$  can be replaced with the Jacobian of the generator, defined as  $\mathcal{J}(\mathbf{C}) := \frac{\partial}{\partial \mathbf{C}} G(\mathbf{C})$ . Contrary to the linear least squares problem, however, in the nonlinear case, the Jacobian is not constant and changes across iterations. Nevertheless, we show that the eigen-values and vectors of the Jacobian at the random initialization govern the dynamics of fitting the network throughout the iterative updates.

For the two-layer convolutional generator that we consider, the left eigenvectors of the Jacobian mapping can be well approximated by the trigonometric basis functions, defined below, throughout the updates. Interestingly, the form of these eigenvectors only depends on the network architecture and not the convolutional kernel used.

**Definition 1.** The trigonometric basis functions  $\mathbf{w}_1, \dots, \mathbf{w}_n$  are defined as

$$[\mathbf{w}_{i}]_{j} = \frac{1}{\sqrt{n}} \begin{cases} 1 & k = 0\\ \sqrt{2}\cos(2\pi ji/n) & k = 1,\dots,n/2 - 1\\ (-1)^{j} & k = n/2\\ \sqrt{2}\sin(2\pi ji/n) & k = n/2 + 1,\dots,n - 1 \end{cases}$$
 (5)

Figure 3 depicts some of these eigenvectors.

In addition to the left eigenvectors we can also approximate the spectrum of the Jacobian throughout the updates by an associated filter/kernel that only depends on the original filter/kernel used in the network architecture.

**Definition 2** (Dual kernel). Associated with a kernel  $\mathbf{u} \in \mathbb{R}^n$  we define the dual kernel  $\sigma \in \mathbb{R}^n$  as

$$\sigma = \|\mathbf{u}\|_{2} \sqrt{\left|\mathbf{F}g\left(\frac{\mathbf{u} \otimes \mathbf{u}}{\|\mathbf{u}\|_{2}^{2}}\right)\right|} \quad with \quad g(t) = \frac{1}{2} \left(1 - \frac{\cos^{-1}(t)}{\pi}\right) t.$$

Here, for two vectors  $\mathbf{u}, \mathbf{v} \in \mathbb{R}^n$ ,  $\mathbf{u} \circledast \mathbf{v}$  denotes their circular convolution, the scalar non-linearity q is applied entrywise, and  $\mathbf{F}$  is the discrete Fourier transform matrix.

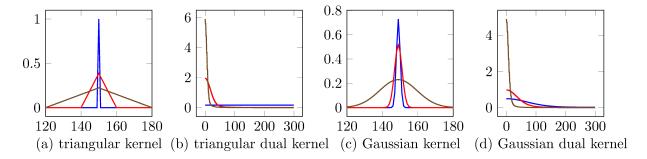


Figure 6: Triangular and Gaussian kernels and the weights associated to low-frequency trigonometric functions they induce, for a generator network of output dimension n = 300. The wider the kernels are, the more the weights are concentrated towards the low-frequency components of the signal.

In Figure 6, we depict two commonly used interpolation kernels  $\mathbf{u}$ , namely a triangular and a Gaussian kernel (recall that the standard upsampling operator is a convolution with a triangle), along with the induced dual kernel  $\sigma$ . The figure shows that the dual filter  $\sigma$  induced by these kernels has a few large values associated with the low frequency trigonometric functions, and the other values associated with high frequencies are small. This in turn implies that the Jacobian spectrum throughout training decreases rapidly.

With these definitions in place we are ready to state our main denoising result. A denoising result requires a signal model—we assume a low-frequency signal  $\mathbf{x}$  that can be represented as a linear combination of the first p-trigonometric basis functions.

**Theorem 1** (Denoising with early stopping). Let  $\mathbf{x} \in \mathbb{R}^n$  be a signal in the span of the first p trigonometric basis functions, and consider a noisy observation

$$y = x + z$$

where  $\mathbf{z}$  is Gaussian noise with distribution  $\mathcal{N}(\mathbf{0}, \frac{\varsigma^2}{n}\mathbf{I})$ . To denoise this signal, we fit a two layer generator network  $G(\mathbf{C}) = \text{ReLU}(\mathbf{UC})\mathbf{v}$  with, for some  $\epsilon > 0$ ,

$$k \ge C_{\mathbf{u}} n/\epsilon^4,$$
 (6)

channels and with convolutional kernel  $\mathbf{u}$  of the convolutional operator  $\mathbf{U}$  and associated kernel  $\sigma$ , to the noisy signal  $\mathbf{y}$ . Here,  $C_{\mathbf{u}}$  a constant only depending on the convolutional kernel  $\mathbf{u}$ . Then, with probability at least  $1 - e^{-k^2} - \frac{1}{n^2}$ , the reconstruction error obtained after  $\tau = \log(1 - \sqrt{p}/n)/\log(1 - \eta\sigma_{p+1}^2)$  iterations of gradient descent (4) with step size  $\eta \leq \frac{1}{\|\mathbf{F}\mathbf{u}\|_{\infty}^2}$  (Fu is the Fourier transform of

**u**) starting from  $\mathbf{C}_0$  with i.i.d.  $\mathcal{N}(0,\omega^2)$ , entries,  $\omega \propto \frac{\|\mathbf{y}\|_2}{\sqrt{n}}$ , is bounded by

$$\|G(\mathbf{C}_{\tau}) - \mathbf{x}\|_{2} \le (1 - \eta \sigma_{p}^{2})^{\tau} \|\mathbf{x}\|_{2} + \varsigma \sqrt{\frac{2p}{n}} + \epsilon, \tag{7}$$

for C a fixed numerical constant.

Note that for this choice of stopping time, provided that dual kernel decays sharply around the p-th singular value, the first term in the bound (7) (i.e.,  $(1 - \eta \sigma_p^2)^{\tau} \approx 0$ ) essentially vanishes and the error bound becomes  $O(\varsigma \sqrt{\frac{p}{n}})$ . The dual kernel decays sharply around the leading eigenvalues provided the kernel is for example a sufficiently wide triangular or Gaussian kernel (see Figure 6).

This result demonstrates that when the noiseless signal  $\mathbf{x}$  is sufficiently structured (e.g. contains only the p lowest frequency components in the trigonometric basis) and the convolutional generator has sufficiently many channels, then early stopped gradient descent achieves a near optimal denoising performance proportional to  $\varsigma\sqrt{\frac{p}{n}}$ . This theorem is obtained from a more general result stated in the appendix which characterizes the evolution of the reconstruction error obtained by the convolutional generator.

**Theorem 2** (Reconstruction dynamics of convolutional generators). Consider the setting and assumptions of Theorem 1 but now with a fixed noise vector  $\mathbf{z}$ , and without an explicit stopping time. Then, for all iterates  $\tau$  obeying  $\tau \leq \frac{1}{\eta \sigma_p^2}$  and provided that  $k \geq C_{\mathbf{u}} n/\epsilon^4$ , for some  $\epsilon > 0$ , with probability at least  $1 - e^{-k^2} - \frac{1}{n^2}$ , the reconstruction error obeys

$$\|G(\mathbf{C}_{\tau}) - \mathbf{x}\|_{2} \le (1 - \eta \sigma_{p}^{2})^{\tau} \|\mathbf{x}\|_{2} + \sqrt{\sum_{i=1}^{n} ((1 - \eta \sigma_{i}^{2})^{\tau} - 1)^{2} \langle \mathbf{w}_{i}, \mathbf{z} \rangle^{2}} + \epsilon.$$

This theorem characterizes the reconstruction dynamics of convolutional generators throughout the updates. In particular, it helps explains why convolutional generators fit a natural signal significantly faster than noise, and thus early stopping enables denoising and regularization. To see this note that as mentioned previously each of the basis functions  $\mathbf{w}_i$  have a (positive) weight, singular value, or dual kernel element  $\sigma_i > 0$  associated with them that only depend on the convolutional kernel used in the architecture (through the definition of the dual kernel). These weights determine how fast the different components of the noisy signal is fitted by gradient descent. As we demonstrated earlier in Figure 6 for typical convolutional filters those weights decay very quickly from low to high frequency basis functions. As a result when the signal  $\mathbf{x}$  is sufficiently structured (i.e. lies in the range of the p trigonometric functions with lowest frequencies), after a few iterations most of the signal is fitted (i.e.,  $(1 - \eta \sigma_p^2)^{\tau}$  is small), while most of the noise has not been fitted (i.e.,  $((1 - \eta \sigma_i^2)^{\tau} - 1)^2 \approx 0$  for  $i = p+1, \ldots, n$ ). Thus, early stopping achieves denoising.

### 4.1 The spectrum of the Jacobian for multilayer networks

Our theoretical results show that for single hidden-layer networks, the leading singular vectors of the Jacobian are the trigonometric functions throughout all iterations, and that the associated weights, singular values, or dual kernel values are concentrated towards the low frequency components. In this section, we show that for a multilayer network, the spectrum of the Jacobian is concentrated towards singular vectors/functions that are similar to the low-frequency components. We also show that throughout training those functions do vary, albeit the low frequency components do not change significantly and the spectrum remains concentrated towards the low frequency components. This shows that the implications of our theory continue to apply to muli-layer networks.

In more detail, we take a standard one dimensional deep decoder with d=4 layers with output in  $\mathbb{R}^{512}$  and with k=64 channels in each layer. Recall that the standard one dimensional decoder

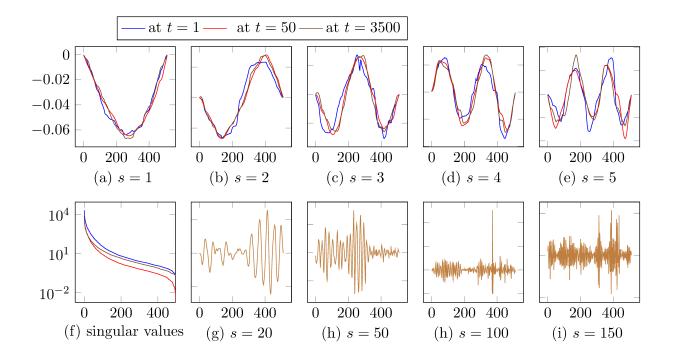


Figure 7: The Singular value distribution of the Jacobian of a four-layer deep decoder after t = 50 and t = 3500 iterations of gradient descent (panel (f)), along with the corresponding singular vectors/function. The singular functions corresponding to the large singular vectors are close to the low-frequency Fourier modes and do not change significantly through training.

obtains layer i+1 from layer i by linearly combining the channels of layer i with learnable coefficients followed by linear upsampling (which involves convolution with the triangular kernel [1/2, 1, 1/2]). The number of parameters is  $d \times k^2 = 32 \cdot 512$ , so the network is over-parameterized by a factor of 32. In Figure 7, we display the singular values as well as the leading singular vectors/function of the Jacobian at initialization (t=1) and after t=50 and t=3500 iterations of gradient descent. As can be seen the leading singular vectors (s=1-6) are close to the trigonometric basis functions and do not change dramatically throughout training. The singular vectors corresponding to increasingly smaller singular values (s=20,50,100,150) contain increasingly higher frequency components but are far from the high-frequency trigonometric basis functions.

## 5 Related literature

As mentioned before, the DIP paper [Uly+18] was the first to show that over-parameterized convolutional networks enable solving denoising, inpainting, and super-resolution problems well even without any training data. Subsequently, the paper [HH19] proposed a much simpler image generating network, termed the deep decoder. The papers [Vee+18; Hec19; JH19] have shown that the DIP and the deep decoder also enable solving or regularizing compressive sensing problems and other inverse problems.

Since the convolutional generators considered here are image-generating deep networks, our work

is also related to methods that rely on trained deep image models. Deep learning based methods are either trained end-to-end for tasks ranging from compression [Tod+16; Agu+17; The+17; Bur+12; Zha+17] to denoising [Bur+12; Zha+17], or are based on learning a generative image model (by training an autoencoder or GAN [HS06; Goo+14]) and then using the resulting model to solve inverse problems such as compressed sensing [Bor+17; HV18], denoising [Hec+18; Hua+18], or phase retrieval [Han+18; SA18], by minimizing an associated loss. In contrast to the method studied here, where the optimization is over the weights of the network, in all the aforementioned methods, the weights are adjusted only during training and then are fixed upon solving the inverse problem.

A large body of work focuses on understanding the optimization landscape of the simple non-linearities or neural networks [Can+15; Sol17; ABG17; Zho+17; Oym18; Fu+18; Tu+15] when the labels are created according to a planted model. Our proofs rely on showing that the dynamics of gradient descent on an over-parameterized network can be related to that of a linear network or a kernel problem. This proof technique has been utilized in a variety of recent publication [Sol+18; Ven+18; Du+18; OS18; OS19; Aro+19; Oym+19]. Two recent publication have used this proof technique to show that functions are learned at different rates: Basri et al. [Bas+19] have shown that functions of different frequencies are learned at different speeds, and Arora et al. [Aro+19] has provided a theoretical explanation of the empirical observation that a simple 2-layer network fits random labels slower than actual labels in the context of classification. A recent publication [Li+19] focuses on demonstrating how early stopping leads to robust classification in the presence of label corruption under a cluster model for the input data. Neither of the aforementioned publication however, does address denoising in a regression setting or fitting convolutional generators of the form studied in this paper.

## Code

Code to reproduce the experiments is available at https://github.com/MLI-lab/overparameterized\_convolutional\_generators.

# Acknowledgements

R. Heckel is supported by NSF award IIS-1816986, and acknowledges support of the NVIDIA Corporation in form of a GPU. M. Soltanolkotabi is supported by the Packard Fellowship in Science and Engineering, a Sloan Research Fellowship in Mathematics, an NSF-CAREER under award #1846369, the Air Force Office of Scientific Research Young Investigator Program (AFOSR-YIP) under award #FA9550-18-1-0078, an NSF-CIF award #1813877, and a Google faculty research award.

#### References

[Agu+17] E. Agustsson, F. Mentzer, M. Tschannen, L. Cavigelli, R. Timofte, L. Benini, and L. V. Gool. "Soft-to-hard vector quantization for end-to-end learning compressible representations". In: Advances in Neural Information Processing Systems. 2017, pp. 1141–1151.

- [ABG17] A. Alon Brutzkus and A. Globerson. "Globally optimal gradient descent for a convnet with Gaussian inputs". In: *arXiv:1702.07966* (2017).
- [Aro+19] S. Arora, S. S. Du, W. Hu, Z. Li, and R. Wang. "Fine-grained analysis of optimization and generalization for overparameterized two-layer neural networks". In: *International Conference on Machine Learning*. 2019.
- [Bas+19] R. Basri, D. Jacobs, Y. Kasten, and S. Kritchman. "The convergence rate of neural networks for learned functions of different frequencies". In: arXiv:1906.00425 (2019).
- [Bor+17] A. Bora, A. Jalal, E. Price, and A. G. Dimakis. "Compressed sensing using generative models". In: arXiv:1703.03208 (2017).
- [Bur+12] H. C. Burger, C. J. Schuler, and S. Harmeling. "Image denoising: Can plain neural networks compete with BM3D?" In: *IEEE Conference on Computer Vision and Pattern Recognition*. 2012, pp. 2392–2399.
- [Can+15] E. J. Candes, X. Li, and M. Soltanolkotabi. "Phase retrieval via Wirtinger flow: Theory and algorithms". In: *IEEE Transactions on Information Theory* 64 (4 2015).
- [Dab+07] K. Dabov, A. Foi, V. Katkovnik, and K. Egiazarian. "Image denoising by sparse 3-D transform-domain collaborative filtering". In: *IEEE Transactions on Image Processing* 16.8 (2007), pp. 2080–2095.
- [Dan+16] A. Daniely, R. Frostig, and Y. Singer. "Toward Deeper Understanding of Neural Networks: The Power of Initialization and a Dual View on Expressivity". In: Advances in Neural Information Processing Systems. 2016, pp. 2253–2261.
- [Du+18] S. S. Du, X. Zhai, B. Poczos, and A. Singh. "Gradient Descent Provably Optimizes Over-parameterized Neural Networks". In: *International Conference on Learning Representations*. 2018.
- [Fu+18] H. Fu, Y. Chi, and Y. Liang. "Local geometry of one-hidden-layer neural networks for logistic regression". In: arXiv:1802.06463 (2018).
- [Goo+14] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A Courville, and Y. Bengio. "Generative adversarial nets". In: Advances in Neural Information Processing Systems. 2014, pp. 2672–2680.
- [HV18] P. Hand and V. Voroninski. "Global guarantees for enforcing deep generative priors by empirical risk". In: *Conference on Learning Theory*. arXiv:1705.07576. 2018.
- [Han+18] P. Hand, O. Leong, and V. Voroninski. "Phase Retrieval Under a Generative Prior". In: arXiv preprint arXiv:1807.04261 (2018).
- [Hec19] R. Heckel. "Regularizing linear inverse problems with convolutional neural networks". In: arXiv:1907.03100 (2019).
- [HH19] R. Heckel and P. Hand. "Deep Decoder: Concise Image Representations from Untrained Non-convolutional Networks". In: *ICLR*. 2019.
- [Hec+18] R. Heckel, W. Huang, P. Hand, and V. Voroninski. "Deep denoising: Rate-optimal recovery of structured signals with a deep prior". In: arXiv:1805.08855 (2018).
- [HS06] G. E. Hinton and R. R. Salakhutdinov. "Reducing the dimensionality of data with neural networks". In: *Science* 313.5786 (2006), pp. 504–507.

- [Hua+18] W. Huang, P. Hand, R. Heckel, and V. Voroninski. "A Provably Convergent Scheme for Compressive Sensing under Random Generative Priors". In: arXiv:1812.04176 (2018).
- [IS15] S. Ioffe and C. Szegedy. "Batch normalization: Accelerating deep network training by reducing internal covariate shift". In: *International Conference on Machine Learning*. 2015, pp. 448–456.
- [JH19] G. Jagatap and C. Hegde. "Algorithmic guarantees for inverse imaging with untrained network priors". In: arXiv:1906.08763 (2019).
- [Li+19] M. Li, M. Soltanolkotabi, and S. Oymak. "Gradient descent with early stopping is provably robust to label noise for overparameterized neural networks". In: arXiv:1903.11680 (2019).
- [Oym18] S. Oymak. "Stochastic gradient descent learns state equations with nonlinear activations". In: arXiv:1809.03019 (2018).
- [OS18] S. Oymak and M. Soltanolkotabi. "Overparameterized nonlinear learning: Gradient descent takes the shortest path?" In: arXiv:1812.10004 (2018).
- [OS19] S. Oymak and M. Soltanolkotabi. "Towards moderate overparameterization: Global convergence guarantees for training shallow neural networks". In: arXiv:1902.04674 (2019).
- [Oym+19] S. Oymak, Z. Fabian, M. Li, and M. Soltanolkotabi. "Generalization guarantees for neural networks via harnessing the low-rank structure of the Jacobian". In: arXiv:1906.05392 (2019).
- [Rad+15] A. Radford, L. Metz, and S. Chintala. "Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks". In: arXiv:1511.06434 [cs] (2015).
- [Ron+15] O. Ronneberger, P. Fischer, and T. Brox. "U-Net: Convolutional Networks for Biomedical Image Segmentation". In: Lecture Notes in Computer Science. 2015.
- [SA18] F. Shamshad and A. Ahmed. "Robust compressive phase retrieval via deep generative priors". In: arXiv preprint arXiv:1808.05854 (2018).
- [Sol17] M. Soltanolkotabi. "Learning ReLUs via Gradient Descent". In: Advances in Neural Information Processing Systems. 2017.
- [Sol+18] M. Soltanolkotabi, A. Javanmard, and J. D. Lee. "Theoretical insights into the optimization landscape of over-parameterized shallow neural networks". In: *IEEE Transactions on Information Theory* (2018).
- [The+17] L. Theis, W. Shi, A. Cunningham, and F. Huszár. "Lossy image compression with compressive autoencoders". In: arXiv:1703.00395 (2017).
- [Tod+16] G. Toderici, S. M. OMalley, S. J. Hwang, D. Vincent, D. Minnen, S. Baluja, M. Covell, and R. Sukthankar. "Variable rate image compression with recurrent neural networks". In: *International Conference on Learning Representations*. 2016.
- [Tro11] J. A. Tropp. "User-Friendly Tail Bounds for Sums of Random Matrices". In: Foundations of Computational Mathematics 12.4 (2011), pp. 389–434.
- [Tu+15] S. Tu, R. Boczar, M. Simchowitz, M. Soltanolkotabi, and B. Recht. "Low-rank solutions of linear matrix equations via procrustes flow". In: arXiv:1507.03566 (2015).

- [Uly+18] D. Ulyanov, A. Vedaldi, and V. Lempitsky. "Deep image prior". In: Conference on Computer Vision and Pattern Recognition. 2018.
- [Vee+18] D. V. Veen, A. Jalal, M. Soltanolkotabi, E. Price, S. Vishwanath, and A. G. Dimakis. "Compressed sensing with Deep Image Prior and learned regularization". In: arXiv:1806.06438 (2018).
- [Ven+18] L. Venturi, A. Bandeira, and J. Bruna. "Spurious valleys in two-layer neural network optimization landscapes". In: arXiv:1802.06384 (2018).
- [Zha+17] K. Zhang, W. Zuo, Y. Chen, D. Meng, and L. Zhang. "Beyond a Gaussian denoiser: Residual learning of deep CNN for image denoising". In: *IEEE Transactions on Image Processing* 26.7 (2017), pp. 3142–3155.
- [Zho+17] K. Zhong, Z. Song, P. Jain, P. L. Bartlett, and I. S. Dhillon. "Recovery guarantees for one-hidden-layer neural networks". In: arXiv:1706.03175 (2017).

# A numerical study of the implicit bias of convolutional networks

In this section, we empirically demonstrate that convolutions with fixed convolutional kernels are critical for convolutional generators to fit natural images faster than noise. Towards this goal, we study numerically the following four closely related architectural choices, which differ in the upsampling/no-upsampling and convolutional operations which generate the activations in the (i+1)-st layer,  $\mathbf{B}_{i+1}$ , from the activations in the i-th layer,  $\mathbf{B}_i$ :

- i) Bilinear upsampling and linear combinations. Layer i+1 is obtained by linearly combining the channels of layer i with learnable coefficients (i.e., performing one-times-one convolutions), followed by bi-linear upsampling. This is the deep decoder architecture from [Hec+18].
- ii) Fixed interpolation kernels and linear combinations. Layer i + 1 is obtained by linearly combining the channels of layer i with learnable coefficients followed by convolving each channel with the same 4x4 interpolation kernel that is used in the linear upsampling operator.
- iii) **Parameterized convolutions:** Layer i + 1 is obtained from layer i though a convolutional layer.
- iv) **Deconvolutional network:** Layer i + 1 is obtained from layer i though a deconvolution layer. This is essentially the DC-GAN [Rad+15] generator architecture.

To emphasize that architectures i)-iv) are structurally very similar operations, we recall that each operation consists only of upsampling and convolutional operations. Let  $\mathbf{T}(\mathbf{c}) \colon \mathbb{R}^n \to \mathbb{R}^n$  be the convolutional operator with kernel  $\mathbf{c}$ , let  $\mathbf{u}$  the linear upsampling kernel (equal to  $\mathbf{u} = [0.5, 1, 0.5]$  in the one-dimensional case), and let  $\mathbf{U} \colon \mathbb{R}^n \to \mathbb{R}^{2n}$  be an upsampling operator, that in the one dimensional case transforms  $[x_1, x_2, \dots, x_n]$  to  $[x_1, 0, x_2, 0, \dots, x_n, 0]$ . In each of the architectures i)-iv), the  $\ell$ -th channel of layer i+1 is obtained from the channels in the i-th layer as:  $\mathbf{b}_{i+1,\ell} = \text{ReLU}\left(\sum_{j=1}^k \mathbf{M}(\mathbf{c}_{ij\ell})\mathbf{b}_i\right)$ , where the linear operator  $\mathbf{M}$  is defined as follows for the four architectures

i) 
$$\mathbf{M}(c) = c\mathbf{T}(\mathbf{u})\mathbf{U}$$
, ii)  $\mathbf{M}(c) = c\mathbf{T}(\mathbf{u})$ , iii)  $\mathbf{M}(\mathbf{c}) = \mathbf{T}(\mathbf{c})$ , iv)  $\mathbf{M}(\mathbf{c}) = \mathbf{T}(\mathbf{c})\mathbf{U}$ .

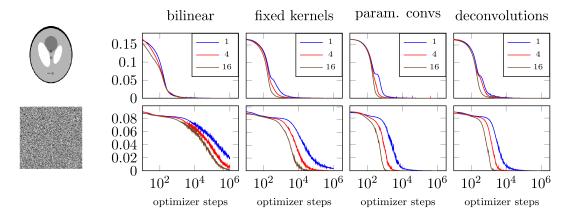


Figure 8: Fitting the phantom MRI and noise with different architectures of depth d = 5, for different number of over-parameterization factors (1,4, and 16). Gradient descent on convolutional generators involving fixed convolutional matrixes fit an image significantly faster than noise.

The coefficients associated with the *i*-th layer are given by  $\mathbf{C}_i = \{\mathbf{c}_{ij\ell}\}$ , and all coefficients of the networks are  $\mathbf{C} = \{\mathbf{c}_{ij\ell}\}$ . Note that here, the coefficients or parameters of the networks are the weights and not the input to the network.

## A.1 Demonstrating implicit bias of convolutional generators

We next show that convolutional generators with fixed convolutional operations fit natural or simple images significantly faster than complex images or noise. Throughout this section, for each image or signal  $\mathbf{x}^*$  we fit weights by minimizing the loss

$$\mathcal{L}(\mathbf{C}) = \|G(\mathbf{C}) - \mathbf{x}^*\|_2^2$$

with respect to the network parameters C using plain gradient descent with a fixed stepsize.

In order to exemplify the effect, we fit the phantom MRI image as well as noise for each of the architectures above for a 5-layer network. We choose the number of channels, k, such that the over-parameterization factor (i.e., the ratio of number of parameters of the network over the output dimensionality) is 1, 4, and 16, respectively. The results in Figure 8 show that for architectures i) and ii) involving fixed convolutional operations, gradient descent requires more than one order of magnitude fewer iterations to obtain a good fit of the phantom MRI image relative to noise. For architectures ii) and iii), with trainable convolutional filters, we see a smaller effect, but the effect essentially vanishes when the network is highly over-parameterized.

This effect continues to exist for natural images in general, as demonstrated by Figure 10 which depicts the average and standard deviation of the loss curves of 100 randomly chosen images from the imagenet dataset.

We also note that the effect continues to exist in the sense that highly structured images with a large number of discontinuities are difficult to fit. An example is the checkerboard image in which each pixel alternates between 1 and 0; this image leads to the same loss curves as noise.

In our next experiment, we highlight that the distance between final and initial network weights is a key feature that determines the difference of fitting a natural image and noise. Towards this goal,

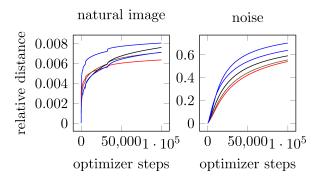


Figure 9: The relative distances of the weights in each layer from its random initialization. The weights need to change significantly more to fit the noise, compared to an image, thus a natural image lies closer to a random initialization than noise.

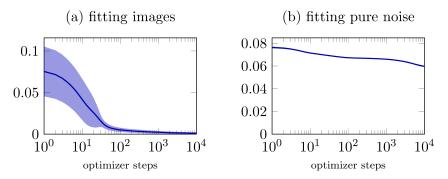


Figure 10: The loss curves for architecture i), a convolutional generator with linear upsampling operations, averaged over  $100~3\times512\times512$  (color) images from the Imagenet dataset. The error band is one standard deviation. Convolutional generators fit natural images significantly faster than noise.

we again fit the phantom MRI image and noise for the architecture i) and an over-parameterization factor of 4 and record, for each layer i the relative distance  $\|\mathbf{C}_i^{(t)} - \mathbf{C}_i^{(0)}\| / \|\mathbf{C}_i^{(0)}\|$ , where  $\mathbf{C}_i^{(0)}$  are the weights at initialization (we initialize randomly), and  $\mathbf{C}_i^{(t)}$  are the weights at the optimizer step t. The results, plotted in Figure 9, show that to fit the noise, the weights have to change significantly, while for fitting a natural image they only change slightly.

# B The spectrum of the Jacobian of the deep decoder and deep image prior

Our theoretical results predict that for over-parameterized networks, the parts of the signal that is aligned with the leading singular vectors of the Jacobian at initialization is fitted fastest. In this section we briefly show that natural images are much more aligned with the leading singular vectors than with Gaussian noise, which is equally aligned with each of the singular vectors.

Towards this goal, we compute the norm of the product of the Jacobian at a random initalization,  $C_0$ , with a signal y as this measures the extent to which the signal is aligned with the leading

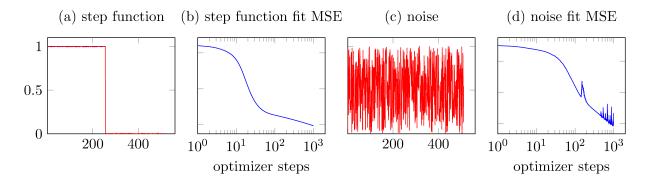


Figure 11: Fitting a step function and noise with a two-layer deep decoder: Even for a two-layer network, the simple image (step function) is fitted significantly faster than the noise.

singular vectors, due to

$$\left\| \mathcal{J}^T(\mathbf{C}_0) \mathbf{y} \right\|^2 = \left\| \mathbf{V} \Sigma \mathbf{W}^T \mathbf{y} \right\|^2 = \sum_i \sigma_i^2 \left\langle \mathbf{w}_i, \mathbf{y} \right\rangle^2,$$

where  $\mathcal{J}(\mathbf{C}_0) = \mathbf{W} \Sigma \mathbf{V}^T$  is the singular value decomposition of the Jacobian.

Figure 12 depicts the distribution of the norm of the product of the Jacobian at initialization,  $\mathcal{J}(\mathbf{C}_0)$ , with an image  $\mathbf{y}^*$  or noise  $\mathbf{z}$  of equal norm ( $\|\mathbf{y}^*\| = \|\mathbf{z}\|$ ). Since for both the deep decoder and the deep image prior, the norm of product of the Jacobian and the noise (i.e.,  $\|\mathcal{J}^T(\mathbf{C}_0)\mathbf{z}\|$ ) is significantly smaller than that with a natural image (i.e.,  $\|\mathcal{J}^T(\mathbf{C}_0)\mathbf{y}^*\|$ ), it follows that a structured image is much better aligned with the leading singular vectors of the Jacobian than the Gaussian noise, which is approximately equally aligned with any of the singular vectors. Thus, the Jacobian at random initialization has an approximate low-rank structure, with natural images lying in the space spanned by the leading singular vectors.

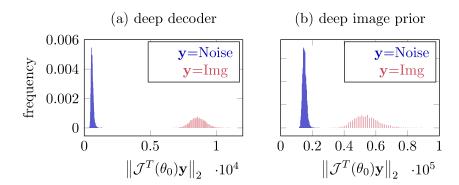


Figure 12: The distribution of the  $\ell_2$ -norm of the inner product of the Jacobian of a deep decoder (a) and a deep image prior (b) at a random initialization, with an image  $\mathbf{y} = \mathbf{x}^*$  and noise  $\mathbf{y} = \mathbf{z}$ , both of equal norm. For both deep decoder and deep image prior, this quantity is significantly smaller for noise than for a natural image. Thus, a natural image is better aligned with the leading singular vectors of the Jacobian than noise. This demonstrates that the Jacobian of the networks is approximately low-rank, with natural images lying in the space spanned by the leading singular vectors.

## C Proofs and formal statement of results

The results stated in the main text are obtained from a slightly more general result which applies beyond convolutional networks. Specifically, we consider neural network generators of the form

$$G(\mathbf{C}) = \text{ReLU}(\mathbf{UC})\mathbf{v},$$
 (8)

with  $\mathbf{C} \in \mathbb{R}^{n \times k}$ , and  $\mathbf{U} \in \mathbb{R}^{n \times n}$  an arbitrary fixed matrix, and  $\mathbf{v} \in \mathbb{R}^k$ , with half of the entries of  $\mathbf{v}$  equal to  $+1/\sqrt{k}$  and the other half equal to  $-1/\sqrt{k}$ .

The (transposed) Jacobian of ReLU(Uc) is  $\mathbf{U}^T \operatorname{diag}(\operatorname{ReLU}'(\mathbf{Uc}))$ . Thus the Jacobian of  $G(\mathbf{C})$  is given by

$$\mathcal{J}^{T}(\mathbf{C}) = \begin{bmatrix} v_{1}\mathbf{U}^{T}\operatorname{diag}(\operatorname{ReLU}'(\mathbf{U}\mathbf{c}_{1})) \\ \vdots \\ v_{k}\mathbf{U}^{T}\operatorname{diag}(\operatorname{ReLU}'(\mathbf{U}\mathbf{c}_{k})) \end{bmatrix} \in \mathbb{R}^{nk \times n},$$
(9)

where ReLU' is the derivative of the activation function. Next we define the neural tangent kernel associated with this generator.

**Definition 3** (Neural Tangent Kernel (NTK)). Associated with a generator  $G(\mathbf{C}) = \text{ReLU}(\mathbf{UC})\mathbf{v}$ , we define the neural tangent kernel

$$\Sigma(\mathbf{U}) \coloneqq \mathbb{E}\left[\mathcal{J}(\mathbf{C})\mathcal{J}^T(\mathbf{C})\right],$$

where expectation is over  $\mathbf{C}$  with iid  $\mathcal{N}(0,1)$  entries.

Consider the eigenvalue decomposition of the NTK given by

$$\mathbf{\Sigma}(\mathbf{U}) = \sum_{i=1}^{n} \sigma_i^2 \mathbf{w}_i \mathbf{w}_i^T.$$

Our results depend on the largest and smallest eigenvalue of the NTK, defined throughout as

$$\alpha = \sigma_n^2, \quad \beta = \sigma_1^2 = \|\mathbf{U}\|.$$

With these definitions in place we are now ready to state our result about neural generators.

**Theorem 3.** Consider a noisy signal  $\mathbf{y} \in \mathbb{R}^n$  given by

$$y = x + z$$

where  $\mathbf{x} \in \mathbb{R}^n$  is assumed to lie in the signal subspace spanned by the p leading singular vectors  $\mathbf{w}_1, \dots, \mathbf{w}_n$  of  $\mathbf{\Sigma}(\mathbf{U})$ , and  $\mathbf{z} \in \mathbb{R}^n$  is an arbitrary noise vector. Suppose that the number of channels obeys

$$k \ge 2^{34} n(\beta/\alpha)^{28} \xi^{-4} (\beta/\alpha + \xi T)^2$$
 (10)

for an error tolerance parameter  $0 < \xi \le \sqrt{8 \log\left(\frac{2n}{\delta}\right)}$  and for a constant T obeying  $1 \le T \le 16 \frac{\beta^8}{\xi^2 \alpha^8}$ . We fit the neural generator  $G(\mathbf{C})$  to the noisy signal  $\mathbf{y} \in \mathbb{R}^n$  by minimizing a loss of the form

$$\mathcal{L}(\mathbf{C}) = \frac{1}{2} \|G(\mathbf{C}) - \mathbf{y}\|_2^2 \tag{11}$$

via running gradient descent with iterations  $\mathbf{C}_{\tau+1} = \mathbf{C}_{\tau} - \eta \nabla \mathcal{L}(\mathbf{C}_{\tau})$ , starting from  $\mathbf{C}_0$  with i.i.d.  $\mathcal{N}(0, \omega^2)$  entries,  $\omega = \frac{\|\mathbf{y}\|_2}{\sqrt{n\beta}} \xi \frac{\alpha^2}{\beta^2}$ , and step size obeying  $\eta \leq 1/\beta^2$ . Then, with probability at least  $1 - e^{-k^2} - \delta$ , for all iterations  $\tau \leq T$ ,

$$\|\mathbf{x} - G(\mathbf{C}_{\tau})\|_{2} \leq (1 - \eta \sigma_{p}^{2})^{\tau} \|\mathbf{x}\|_{2} + \sqrt{\sum_{i=1}^{n} ((1 - \eta \sigma_{i}^{2})^{\tau} - 1)^{2} \langle \mathbf{w}_{i}, \mathbf{z} \rangle^{2}} + 2\xi \frac{\alpha^{2}}{\beta^{2}} \sqrt{8 \log(2n/\delta)} \|\mathbf{y}\|_{2}.$$
(12)

## C.1 Proof of Theorem 2

Theorem 2 stated in the main text follow directly from Theorem 3 above as follows. We first note that for a convolutional generator (where **U** implements a convolution and thus is circulant) the eigenvectors of the NTK are given by the trigonometric basis functions per Definition 1 and the eigenvalues are the square of the entries of the dual kernel (Definition 2). To see this, we note that as detailed in Section **G**, the neural tangent kernel is given by

$$\left[\mathbf{\Sigma}(\mathbf{U})\right]_{ij} = \frac{1}{2} \left(1 - \cos^{-1} \left(\frac{\langle \mathbf{u}_i, \mathbf{u}_j \rangle}{\|\mathbf{u}_i\| \|\mathbf{u}_i\|}\right) / \pi\right) \langle \mathbf{u}_i, \mathbf{u}_j \rangle. \tag{13}$$

Because the matrix  $\mathbf{U}$  implements a convolution with a kernel  $\mathbf{u}$  that is equal to its first column, the neural tangent kernel  $\Sigma(\mathbf{U})$  is again a circulant matrix and is also Hermitian. Thus, its spectrum is given by the Fourier transform of the first column of the circulant matrix, and its left-singular vectors are given by the trigonometric basis functions defined in equation (5) and depicted in Figure 3.

Furthermore, using the fact that the eigenvalues of a circulant matrix are given by its discrete Fourier transform we can substitute  $\beta = \|\mathbf{U}\| = \|\mathbf{F}\mathbf{u}\|_{\infty}$  and  $\alpha = \sigma_n$  to conclude that  $C_{\mathbf{u}} \propto \left(\frac{\|\mathbf{F}\mathbf{u}\|_{\infty}}{\sigma_n}\right)^{30}$ , where  $\mathbf{F}\mathbf{u}$  is the discrete Fourier transform of  $\mathbf{u}$ . This yields Theorem 2.

#### C.2 Proof of Theorem 1

Finally, we note that to obtain the simplified final expression in Theorem 1 from Theorem 2 we also used the fact that for a Gaussian vector  $\mathbf{z}$ , the vector  $\mathbf{W}^T\mathbf{z}$  is also Gaussian. Furthermore, by the concentration of Lipschitz functions of Gaussians with high probability we have

$$\sum_{i=1}^{n} ((1 - \eta \sigma_i^2)^{\tau} - 1)^2 \langle \mathbf{w}_i, \mathbf{z} \rangle^2 \approx \mathbb{E} \left[ \sum_{i=1}^{n} ((1 - \eta \sigma_i^2)^{\tau} - 1)^2 \langle \mathbf{w}_i, \mathbf{z} \rangle^2 \right]$$

$$\stackrel{\text{(i)}}{=} \frac{\varsigma^2}{n} \sum_{i=1}^{n} ((1 - \eta \sigma_i^2)^{\tau} - 1)^2$$

$$\stackrel{\text{(ii)}}{\leq} \varsigma^2 \frac{2p}{n}.$$

Here, equation (i) follows from  $\langle \mathbf{w}_i, \mathbf{z} \rangle$  being zero mean Gaussian with variance  $\varsigma^2/n$  (since  $\mathbf{z} \sim \mathcal{N}(0, (\varsigma^2/n)\mathbf{I})$ , and  $\|\mathbf{w}_i\|_2 = 1$ ). Finally, (ii) follows by choosing the early stop time so that it obeys  $\tau \leq \log(1 - \sqrt{p}/n)/\log(1 - \eta\sigma_{p+1}^2)$  which in turn implies that  $(1 - \eta\sigma_i^2)^{\tau} \geq 1 - \sqrt{p}/n$ , for all i > p, yielding that  $((1 - \eta\sigma_i^2)^{\tau} - 1)^2 \leq p/n^2$ , for all i > p.

# D The dynamics of linear and nonlinear least-squares

Theorem 3 builds on a more general result on the dynamics of a non-linear least squares problem which applies beyond convolutional networks that is stated and discussed in this section. Consider a nonlinear least-squares fitting problem of the form

$$\mathcal{L}(\theta) = \frac{1}{2} \| f(\theta) - \mathbf{y} \|_2^2.$$

Here,  $f: \mathbb{R}^p \to \mathbb{R}^n$  is a non-linear model with parameters  $\theta \in \mathbb{R}^p$ . To solve this problem, we run gradient descent with a fixed stepsize  $\eta$ , starting from an initial point  $\theta_0$ , with updates of the form

$$\theta_{\tau+1} = \theta_{\tau} - \eta \nabla \mathcal{L}(\theta_{\tau}) \quad \text{where} \quad \nabla \mathcal{L}(\theta) = \mathcal{J}^{T}(\theta)(f(\theta) - \mathbf{y}).$$
 (14)

Here,  $\mathcal{J}(\theta) \in \mathbb{R}^{n \times p}$  is the Jacobian associated with the nonlinear map f with entries given by  $[\mathcal{J}(\theta)]_{i,j} = \frac{\partial f_i(\theta)}{\partial \theta_j}$ . In order to study the properties of the gradient descent iterates (14), we relate the non-linear least squares problem to a linearized one in a ball around the initialization  $\theta_0$ . We note that this general strategy has been utilized in a variety of recent publications [Du+18; OS19; Oym+19]. The associated linearized least-squares problem is defined as

$$\mathcal{L}_{\text{lin}}(\theta) = \frac{1}{2} \| f(\theta_0) + \mathbf{J}(\theta - \theta_0) - \mathbf{y} \|_2^2.$$

$$\tag{15}$$

Here,  $\mathbf{J} \in \mathbb{R}^{n \times p}$ , referred to as the reference Jacobian, is a fixed matrix independent of  $\theta$  that approximates the Jacobian mapping at initialization,  $\mathcal{J}(\theta_0)$ , and is formally defined later. Starting from the same initial point  $\theta_0$ , the gradient descent updates of the linearized problem are

$$\widetilde{\theta}_{\tau+1} = \widetilde{\theta}_{\tau} - \eta \mathbf{J}^{T} \left( f(\theta_{0}) + \mathbf{J} (\widetilde{\theta}_{\tau} - \theta_{0}) - \mathbf{y} \right)$$

$$= \widetilde{\theta}_{\tau} - \eta \mathbf{J}^{T} \mathbf{J} \left( \widetilde{\theta}_{\tau} - \theta_{0} \right) - \eta \mathbf{J}^{T} \left( f(\theta_{0}) - \mathbf{y} \right).$$
(16)

To show that the non-linear updates (14) are close to the linearized iterates (16), we make the following assumptions:

**Assumption 1** (Bounded spectrum). We assume the singular values of the reference Jacobian obey

$$\alpha \le \sigma_n \le \sigma_1 \le \beta. \tag{17}$$

Furthermore, we assume that the Jacobian mapping associated with the nonlinear model f obeys

$$\|\mathcal{J}(\theta)\| < \beta \quad \text{for all} \quad \theta \in \mathbb{R}^p.$$
 (18)

**Assumption 2** (Closeness of the reference and initialization Jacobians). We assume the reference Jacobian and the Jacobian of the nonlinearity at initialization  $\mathcal{J}(\theta_0)$  are  $\epsilon_0$ -close in the sense that

$$\|\mathcal{J}(\theta_0)\mathcal{J}^T(\theta_0) - \mathbf{J}\mathbf{J}^T\| \le \epsilon_0^2, \quad and \quad \|\mathcal{J}(\theta_0) - \mathbf{J}\| \le \epsilon_0.$$
 (19)

**Assumption 3** (Bounded variation of Jacobian around initialization). We assume that within a radius R around the initialization, the Jacobian varies by no more than  $\epsilon$  in the sense that

$$\|\mathcal{J}(\theta) - \mathcal{J}(\theta_0)\| \le \frac{\epsilon}{2}, \quad \text{for all} \quad \theta \in \mathcal{B}_R(\theta_0),$$
 (20)

where  $\mathcal{B}_R(\theta_0) := \{\theta : \|\theta - \theta_0\| \le R\}$  is the ball with radius R around  $\theta_0$ .

Our first result shows that under these assumptions the nonlinear iterative updates (14) are intimately related to the linear iterative updates (16). Specifically, we show that the residuals associated with these two problems defined below

nonlinear residual: 
$$\mathbf{r}_{\tau} := f(\theta_{\tau}) - \mathbf{y}$$
 (21)

linear residual: 
$$\widetilde{\mathbf{r}}_{\tau} := (\mathbf{I} - \eta \mathbf{J} \mathbf{J}^T)^{\tau} \mathbf{r}_0.$$
 (22)

are close in the proximity of the initialization.

**Theorem 4** (Closeness of linear and nonlinear least-squares problems). Assume the Jacobian mapping  $\mathcal{J}(\theta) \in \mathbb{R}^{n \times p}$  associated with the function  $f(\theta)$  obeys Assumptions 1, 2, and 3 around an initial point  $\theta_0 \in \mathbb{R}^p$  with respect to a reference Jacobian  $\mathbf{J} \in \mathbb{R}^{n \times p}$  and with parameters  $\alpha, \beta, \epsilon_0, \epsilon$ , and R. Furthermore, assume the radius R is given by

$$\frac{R}{2} := \left\| \mathbf{J}^{\dagger} \mathbf{r}_0 \right\|_2 + \frac{\kappa}{2} \left( \frac{1}{\alpha^2} + \eta T \right) \left\| \mathbf{r}_0 \right\|_2 \tag{23}$$

with  $\kappa$  and T constants obeying  $\kappa > 0$  and  $1 \leq T \leq \frac{1}{2\eta\epsilon^2}$ , and  $\mathbf{J}^{\dagger}$  the pseudo-inverse of  $\mathbf{J}$ . Also assume that  $\alpha, \beta, \epsilon_0, \epsilon$ , and  $\kappa$  obey

$$\epsilon \le \frac{1}{8} \frac{\kappa \alpha^2}{\beta^2} \quad and \quad \epsilon_0 \le \min\left(\frac{1}{4}\kappa, \sqrt{\frac{1}{8} \frac{\kappa \alpha^2}{\beta}}\right).$$
(24)

We run gradient descent with stepsize  $\eta \leq \frac{1}{\beta^2}$  on the linear and non-linear least squares problem, starting from the same initialization  $\theta_0$ . Then, for all  $\tau \leq T$  the iterates of the original and the linearized problems and their corresponding residuals obey

$$\|\mathbf{r}_{\tau} - \widetilde{\mathbf{r}}_{\tau}\| \le \frac{1}{2} \frac{\kappa}{\beta} \|\mathbf{r}_0\|_2, \tag{25}$$

$$\left\| \theta_{\tau} - \widetilde{\theta}_{\tau} \right\| \le \frac{\kappa}{2} \left( \frac{1}{\alpha^2} + \eta \tau \right) \left\| \mathbf{r}_0 \right\|_2. \tag{26}$$

Moreover, for all iterates  $\tau \leq T$ ,

$$\|\theta_{\tau} - \theta_0\| \le \frac{R}{2}.\tag{27}$$

The above theorem formalizes that in a (small) radius around the initialization, the non-linear problem behaves similarly as its linearization. Thus to characterize the dynamics of the nonlinear problem, it suffices to characterize the dynamics of the linearized problem. This is the subject of our next theorem.

**Theorem 5.** Consider a linear least squares problem (15) and let  $\mathbf{J} = \mathbf{W} \mathbf{\Sigma} \mathbf{V}^T \in \mathbb{R}^{n \times p} = \sum_{i=1}^n \sigma_i \mathbf{w}_i \mathbf{v}_i^T$  be the singular value decomposition of the matrix  $\mathbf{J}$ . Then the residual  $\tilde{\mathbf{r}}_{\tau}$  after  $\tau$  iterations of gradient descent with updates (16) is

$$\widetilde{\mathbf{r}}_{\tau} = \sum_{i=1}^{n} \left( 1 - \eta \sigma_i^2 \right)^{\tau} \mathbf{w}_i \left\langle \mathbf{w}_i, \mathbf{r}_0 \right\rangle. \tag{28}$$

Moreover, using a step size satisfying  $\eta \leq \frac{1}{\sigma_1^2}$ , the linearized iterates (16) obey

$$\left\|\widetilde{\theta}_{\tau} - \theta_0\right\|_2^2 \le \sum_{i=1}^n \left( \langle \mathbf{w}_i, \mathbf{r}_0 \rangle \frac{1 - (1 - \eta \sigma_i^2)^{\tau}}{\sigma_i} \right)^2.$$
 (29)

In the next section we will show we can combine these two general theorems to provide guarantees for denoising using general neural networks.

### D.1 Proof of Theorem 4 (closeness of linear and non-linear least-squares)

The proof is by induction. We suppose the statement, in particular the bounds (25), (26), and (27) hold for iterations  $t \leq \tau - 1$ . We then show that they continue to hold for iteration  $\tau$  in four steps. In Step I, we show that a weaker version of (27) holds, specifically that  $\|\theta_{\tau} - \theta_{0}\|_{2} \leq R$ . In Steps II and III we show that the bounds (25) and (26) hold, respectively. Finally, in Step IV we utilize Steps I-III to complete the proof of equation (27).

Step I: Next iterate obeys  $\theta_{\tau} \in \mathcal{B}_{R}(\theta_{0})$ . To prove  $\theta_{\tau} \in \mathcal{B}_{R}(\theta_{0})$ , first note that by the triangle inequality and the induction assumption (27) we have

$$\|\theta_{\tau} - \theta_{0}\|_{2} \leq \|\theta_{\tau} - \theta_{\tau-1}\|_{2} + \|\theta_{\tau-1} - \theta_{0}\|_{2},$$
  
$$\leq \|\theta_{\tau} - \theta_{\tau-1}\|_{2} + \frac{R}{2}.$$

So to prove  $\|\theta_{\tau} - \theta_0\|_2 \le R$  it suffices to show that  $\|\theta_{\tau} - \theta_{\tau-1}\|_2 \le R/2$ . To this aim note that

$$\|\theta_{\tau} - \theta_{\tau-1}\|_{2} = \eta \|\nabla \mathcal{L}(\theta_{\tau-1})\|_{2}$$

$$= \eta \|\mathcal{J}^{T}(\theta_{\tau-1})\mathbf{r}_{\tau-1}\|_{2}$$

$$\leq \eta \left(\|\mathcal{J}^{T}(\theta_{\tau-1})\widetilde{\mathbf{r}}_{\tau-1}\|_{2} + \|\mathcal{J}(\theta_{\tau-1})\|\|\mathbf{r}_{\tau-1} - \widetilde{\mathbf{r}}_{\tau-1}\|_{2}\right)$$

$$\leq \eta \left(\|\mathbf{J}^{T}\widetilde{\mathbf{r}}_{\tau-1}\|_{2} + \|\mathcal{J}(\theta_{\tau-1}) - \mathbf{J}\|\|\widetilde{\mathbf{r}}_{\tau-1}\|_{2} + \|\mathcal{J}(\theta_{\tau-1})\|\|\mathbf{r}_{\tau-1} - \widetilde{\mathbf{r}}_{\tau-1}\|_{2}\right)$$

$$\stackrel{(i)}{\leq} \eta \|\mathbf{J}^{T}\widetilde{\mathbf{r}}_{\tau-1}\|_{2} + \eta(\epsilon + \epsilon_{0})\|\widetilde{\mathbf{r}}_{\tau-1}\|_{2} + \eta\beta\|\mathbf{r}_{\tau-1} - \widetilde{\mathbf{r}}_{\tau-1}\|_{2}$$

$$\stackrel{(ii)}{\leq} \eta \|\mathbf{J}^{T}\widetilde{\mathbf{r}}_{\tau-1}\|_{2} + \eta\frac{1}{2}\kappa\|\widetilde{\mathbf{r}}_{\tau-1}\|_{2} + \eta\frac{1}{2}\kappa\|\mathbf{r}_{0}\|_{2}$$

$$\stackrel{(iii)}{\leq} \eta \|\mathbf{J}^{\dagger}\mathbf{r}_{0}\|_{2} + \eta\kappa\|\mathbf{r}_{0}\|_{2}$$

$$\leq \frac{R}{2}.$$

$$(30)$$

In the above (i) follows from Assumptions 1, 2, and 3, (ii) from  $\epsilon_0 + \epsilon \leq \frac{1}{2}\kappa$  (which is a consequence of (24)) combined with with the induction hypothesis (25), and (iii) follows from  $\|\tilde{\mathbf{r}}_{\tau-1}\| \leq \|\mathbf{r}_0\|$  as well as from the bound

$$\begin{split} \left\| \mathbf{J}^T \widetilde{\mathbf{r}}_{\tau-1} \right\|_2 &= \left\| \mathbf{J}^T (\mathbf{I} - \eta \mathbf{J} \mathbf{J}^T)^{\tau-1} \mathbf{r}_0 \right\|_2 \\ &= \left\| \mathbf{\Sigma} (\mathbf{I} - \eta \mathbf{\Sigma}^2)^{\tau-1} \mathbf{W}^T \mathbf{r}_0 \right\|_2 \\ &\leq \sqrt{\sum_{j=1}^n \sigma_j^2 \langle \mathbf{w}_j, \mathbf{r}_0 \rangle^2} \\ &= \left\| \mathbf{J}^\dagger \mathbf{r}_0 \right\|_2. \end{split}$$

Finally, the last inequality follows by definition of R in (23) together with the fact that  $T \geq 1$ .

**Step II: Original and linearized residuals are close:** In this step, we bound the deviation of the residuals of the original and linearized problem defined as

$$\mathbf{e}_{\tau} \coloneqq \mathbf{r}_{\tau} - \widetilde{\mathbf{r}}_{\tau}.$$

This step relies on the following lemma from [Oym+19].

**Lemma 1** (Bound on growth of perturbations, [Oym+19, Lem. 6.7]). Suppose that Assumptions 1, 2, and 3 hold and that  $\theta_{\tau}, \theta_{\tau+1} \in \mathcal{B}_R(\theta_0)$ . Then, provided the stepsize obeys  $\eta \leq 1/\beta^2$ , the deviation of the residuals obeys

$$\|\mathbf{e}_{\tau+1}\|_{2} \le \eta \left(\epsilon_{0}^{2} + \epsilon \beta\right) \|\widetilde{\mathbf{r}}_{\tau}\|_{2} + \left(1 + \eta \epsilon^{2}\right) \|\mathbf{e}_{\tau}\|_{2}.$$
 (31)

By the previous step,  $\theta_{\tau}$ ,  $\theta_{\tau+1} \in \mathcal{B}_R(\theta_0)$ . We next bound the two terms on the right hand side of (31). Regarding the first term, we note that an immediate consequence of Theorem 5 is the following bound on the residual of the linearized iterates:

$$\|\widetilde{\mathbf{r}}_{\tau}\|_{2} \le \left(1 - \eta \alpha^{2}\right)^{\tau} \|\mathbf{r}_{0}\|_{2}.\tag{32}$$

In order to bound the second term in (31), namely,  $\|\mathbf{e}_{\tau}\|_{2}$ , we used the following lemma, proven later in Section D.1.1.

**Lemma 2.** Suppose that for positive scalars  $\alpha, \eta, \rho, \xi > 0$ ,  $\eta \leq 1/\alpha^2$ , the sequences  $\tilde{r}_{\tau}$  and  $e_{\tau}$  obey

$$\widetilde{r}_{\tau} \le \left(1 - \eta \alpha^2\right)^{\tau} \rho \tag{33}$$

$$e_{\tau} \le \left(1 + \eta \epsilon^2\right) e_{\tau - 1} + \eta \xi \widetilde{r}_{\tau - 1} \tag{34}$$

Then, for all  $\tau \leq \frac{1}{2\eta\epsilon^2}$ , we have that

$$e_{\tau} \le 2\xi \frac{\rho}{\alpha^2}.\tag{35}$$

With these lemmas in place we now have all the tools to prove that the original and linear residuals are close. In particular, from Step I, we know that  $\theta_{\tau} \in \mathcal{B}_{R}(\theta_{0})$  so that the assumptions of Lemma 1 are satisfied. Lemma 1 implies that the assumption (34) of Lemma 2 is satisfied with

 $\xi = \epsilon_0^2 + \epsilon \beta$  and the bound (32) implies that the assumption (33) of Lemma 2 is satisfied with  $\rho = \|\mathbf{r}_0\|_2$ . Thus, Lemma 2 implies that for all  $\tau \leq \frac{1}{2\eta\epsilon^2}$ 

$$\|\mathbf{e}_{\tau}\|_{2} \leq 2 \frac{\left(\epsilon_{0}^{2} + \epsilon \beta\right)}{\alpha^{2}} \|\mathbf{r}_{0}\|_{2} \leq \frac{\kappa}{2\beta} \|\mathbf{r}_{0}\|_{2}, \tag{36}$$

where in the last inequality we used the assumption that  $\frac{2(\epsilon_0^2 + \epsilon \beta)}{\alpha^2} \leq \frac{\kappa}{2\beta}$ . This concludes the proof of (25).

**Step III: Original and linearized parameters are close:** First note that by the triangle inequality and Assumption 3 we have

$$\|\mathcal{J}(\theta_{\tau}) - \mathbf{J}\| \le \|\mathcal{J}(\theta_{\tau}) - \mathcal{J}(\theta_{0})\| + \|\mathcal{J}(\theta_{0}) - \mathbf{J}\| \le \epsilon_{0} + \epsilon. \tag{37}$$

The difference between the parameter of the original iterate  $\theta$  and the linearized iterate  $\widetilde{\theta}$  obey

$$\frac{1}{\eta} \left\| \boldsymbol{\theta}_{\tau} - \widetilde{\boldsymbol{\theta}}_{\tau} \right\|_{2} \leq \left\| \sum_{t=0}^{\tau-1} \nabla \mathcal{L}(\boldsymbol{\theta}_{t}) - \nabla \mathcal{L}_{\text{lin}}(\widetilde{\boldsymbol{\theta}}_{t}) \right\|_{2}$$

$$= \left\| \sum_{t=0}^{\tau-1} \mathcal{J}^{T}(\boldsymbol{\theta}_{t}) \mathbf{r}_{t} - \mathbf{J}^{T} \widetilde{\mathbf{r}}_{t} \right\|_{2}$$

$$\leq \sum_{t=0}^{\tau-1} \left\| (\mathcal{J}^{T}(\boldsymbol{\theta}_{t}) - \mathbf{J}^{T}) \widetilde{\mathbf{r}}_{t} \right\|_{2} + \left\| \mathcal{J}^{T}(\boldsymbol{\theta}_{t}) (\mathbf{r}_{t} - \widetilde{\mathbf{r}}_{t}) \right\|_{2}$$

$$\stackrel{(i)}{\leq} \sum_{t=0}^{\tau-1} (\epsilon_{0} + \epsilon) \|\widetilde{\mathbf{r}}_{t}\|_{2} + \beta \|\mathbf{e}_{t}\|_{2}$$

$$\stackrel{(ii)}{\leq} (\epsilon_{0} + \epsilon) \sum_{t=0}^{\tau-1} (1 - \eta \alpha^{2})^{\tau-1} \|\mathbf{r}_{0}\|_{2} + \beta \sum_{t=0}^{\tau-1} \|\mathbf{e}_{t}\|_{2}$$

$$= (\epsilon_{0} + \epsilon) \frac{1 - (1 - \eta \alpha^{2})^{\tau}}{\eta \alpha^{2}} \|\mathbf{r}_{0}\|_{2} + \beta \sum_{t=0}^{\tau-1} \|\mathbf{e}_{t}\|_{2}$$

$$\stackrel{(iii)}{\leq} \frac{(\epsilon_{0} + \epsilon)}{\eta \alpha^{2}} \|\mathbf{r}_{0}\|_{2} + \frac{\kappa \tau}{2} \|\mathbf{r}_{0}\|_{2}$$

$$\stackrel{(iv)}{\leq} \frac{\kappa}{2\eta \alpha^{2}} \|\mathbf{r}_{0}\|_{2} + \frac{\kappa \tau}{2} \|\mathbf{r}_{0}\|_{2},$$

where (i) follows from (37) combined with Assumption 1, (ii) from (32), (iii) from  $\eta \leq 1/\beta^2$  which implies  $(1 - \eta \alpha^2) \geq 0$  and (36). Finally, (iv) follows from the fact that  $\epsilon \leq \frac{1}{8}\kappa$  and  $\epsilon_0 \leq \frac{1}{4}\kappa$  per (24). This concludes the proof of (26).

Step IV: Completing the proof of (27): By the triangle inequality

$$\|\boldsymbol{\theta}_{\tau} - \boldsymbol{\theta}_{0}\|_{2} \leq \|\widetilde{\boldsymbol{\theta}}_{\tau} - \boldsymbol{\theta}_{0}\|_{2} + \|\boldsymbol{\theta}_{\tau} - \widetilde{\boldsymbol{\theta}}_{\tau}\|_{2}$$

$$\stackrel{\text{(i)}}{\leq} \|\mathbf{J}^{\dagger}\mathbf{r}_{0}\|_{2} + \frac{\kappa}{2} \left(\frac{1}{\alpha^{2}} + \eta\tau\right) \|\mathbf{r}_{0}\|_{2}$$

$$\stackrel{\text{(ii)}}{\leq} R/2,$$

where inequality (i) follows from the bound (26), which we just proved, and the fact that, from equation (38) in Theorem 5,

$$\begin{aligned} \left\| \widetilde{\theta}_{\tau} - \theta_{0} \right\|_{2}^{2} &= \sum_{i=1}^{n} \langle \mathbf{w}_{i}, \mathbf{r}_{0} \rangle^{2} \frac{(1 - (1 - \eta \sigma_{i}^{2})^{\tau})^{2}}{\sigma_{i}^{2}} \\ &\leq \sum_{i=1}^{n} \langle \mathbf{w}_{i}, \mathbf{r}_{0} \rangle^{2} / \sigma_{i}^{2} \\ &= \left\| \mathbf{J}^{\dagger} \mathbf{r}_{0} \right\|^{2}. \end{aligned}$$

Moreover, inequality (ii) follows from  $\tau \leq T$ , by assumption as well as from the definition of R in equation (23).

### D.1.1 Proof of Lemma 2

We prove the result by induction. Assume equation (35) holds true for some  $\tau$ . We prove that then it also holds true for  $\tau + 1$ . By the two assumptions in the lemma,

$$\begin{aligned} e_{\tau+1} - e_{\tau} &\leq \eta \epsilon^2 e_{\tau} + \eta \xi \widetilde{r}_{\tau} \\ &\leq \eta \epsilon^2 e_{\tau} + \eta \xi (1 - \eta \alpha^2)^{\tau} \rho \\ &\stackrel{\text{(i)}}{\leq} \eta \xi \left( \epsilon^2 \frac{2\rho}{\alpha^2} + (1 - \eta \alpha^2)^{\tau} \rho \right), \end{aligned}$$

where (i) follows from the induction assumption (35). Summing up the difference of the errors gives

$$\begin{split} \frac{e_{\tau}}{\xi} &= \sum_{t=0}^{\tau-1} \frac{(e_{t+1} - e_t)}{\xi} \\ &\leq \tau \eta \epsilon^2 \frac{2\rho}{\alpha^2} + \eta \rho \sum_{t=0}^{\tau-1} (1 - \eta \alpha^2)^t \\ &= \tau \eta \epsilon^2 \frac{2\rho}{\alpha^2} + \eta \rho \frac{1 - (1 - \eta \alpha^2)^\tau}{\eta \alpha^2} \\ &\leq \frac{\rho}{\alpha^2} (\eta 2\tau \epsilon^2 + 1) \\ &\leq 2 \frac{\rho}{\alpha^2}, \end{split}$$

where the last inequality follows from the assumption that  $\tau \leq \frac{1}{2n\epsilon^2}$ .

#### D.2 Proof of Theorem 5

The proof of identity (28) is equivalent to the proof of Proposition (1). Regarding inequality (29), note that

$$\widetilde{\theta}_{\tau} - \widetilde{\theta}_{0} = -\eta \sum_{t=0}^{\tau-1} \nabla \mathcal{L}_{\text{lin}}(\widetilde{\theta}_{t}) = -\eta \sum_{t=0}^{\tau-1} \mathbf{J}^{T} \widetilde{\mathbf{r}}_{t} = -\eta \mathbf{V} \left( \sum_{t=0}^{\tau-1} \mathbf{\Sigma} \left( \mathbf{I} - \eta \mathbf{\Sigma}^{2} \right)^{t} \right) \begin{bmatrix} \langle \mathbf{u}_{1}, \mathbf{r}_{0} \rangle \\ \vdots \\ \langle \mathbf{u}_{n}, \mathbf{r}_{0} \rangle \end{bmatrix}.$$

Thus

$$\left\langle \mathbf{v}_{i}, \widetilde{\theta}_{\tau} - \widetilde{\theta}_{0} \right\rangle = -\eta \sigma_{i} \left\langle \mathbf{u}_{i}, \mathbf{r}_{0} \right\rangle \left( \sum_{t=0}^{\tau-1} (1 - \eta \sigma_{i}^{2})^{t} \right) = -\eta \sigma_{i} \left\langle \mathbf{u}_{i}, \mathbf{r}_{0} \right\rangle \frac{1 - (1 - \eta \sigma_{i}^{2})^{\tau}}{\eta \sigma_{i}^{2}}.$$
(38)

This in turn implies that

$$\left\|\widetilde{\theta}_{\tau} - \widetilde{\theta}_{0}\right\|_{2}^{2} \leq \sum_{i=1}^{n} \left(\left\langle \mathbf{u}_{i}, \mathbf{r}_{0} \right\rangle \frac{1 - (1 - \eta \sigma_{i}^{2})^{\tau}}{\sigma_{i}}\right)^{2}.$$

# E Proofs for neural network generators (proof of Theorem 3)

The proof of Theorem 3 relies on the fact that, in the overparameterized regime, the non-linear least squares problem is well approximated by an associated linearized least squares problem. Studying the associated linear problem enables us to prove the result.

We apply Theorem 4, which ensures that the associated linear problem is a good approximation of the non-linear least squarest problem, with the nonlinearity  $G(\mathbf{C}) = \text{ReLU}(\mathbf{UC})\mathbf{v}$  and with the parameter given by  $\theta = \mathbf{C}$ . Recall that  $\mathbf{v}$  is a fixed vector with half of the entries  $1/\sqrt{k}$ , and the other half  $-1/\sqrt{k}$ . As the reference Jacobian in the associated linear problem, we choose

$$\mathbf{J} \coloneqq \left( \mathbb{E} \left[ \mathcal{J}(\mathbf{C}) \mathcal{J}^T(\mathbf{C}) \right] \right)^{\frac{1}{2}},$$

so that it obeys  $\mathbf{JJ}^T = \mathbb{E}\left[\mathcal{J}(\mathbf{C})\mathcal{J}^T(\mathbf{C})\right]$ . We apply the theorem with parameters

$$\alpha = \sigma_n\left(\mathbf{\Sigma}(\mathbf{U})\right), \quad \beta = \|\mathbf{U}\|, \quad \kappa = \xi \frac{\alpha^2}{\beta}, \quad \epsilon_0 = 2\beta \left(\frac{\log(\frac{2n}{\delta})}{k}\right)^{1/4}, \quad \epsilon = \frac{\xi}{8} \|\mathbf{U}\| \frac{\alpha^4}{\beta^4}, \quad \omega = \frac{\|\mathbf{y}\|_2}{\sqrt{n\beta}} \xi \frac{\alpha^2}{\beta^2}.$$

We next verify by applying a series of Lemmas proven in Appendix H that the conditions of Theorem 4 are satisfied (specifically, Assumptions 1, 2, 3 and (24) on the Jacobians of the non-linear map and the associated linearized map, to be bounded and sufficiently close to each other).

**Bound on initial residual:** We start with bounding the initial residual by applying the following lemma.

**Lemma 3** (Initial residual). Consider  $G(\mathbf{C}) = \text{ReLU}(\mathbf{UC})\mathbf{v}$ , and let  $\mathbf{C} \in \mathbb{R}^{n \times k}$  be generated at random with i.i.d.  $\mathcal{N}(0,\omega^2)$  entries. Suppose half of the entries of  $\mathbf{v}$  are  $\nu/\sqrt{k}$  and the other half are  $-\nu/\sqrt{k}$ , for some constant  $\nu > 0$ . Then, with probability at least  $1 - \delta$ ,

$$\left\|G(\mathbf{C})\right\|_2 \leq \nu \omega \sqrt{8\log(2n/\delta)} \left\|\mathbf{U}\right\|_F.$$

Now, the initial residual can be upper bounded as

$$\|\mathbf{r}_0\|_2 \le \|\mathbf{y}\|_2 + \|G(\mathbf{C}_0)\|_2 \le 2\|\mathbf{y}\|_2,$$
 (39)

where we used that, by Lemma 3,

$$\|G(\mathbf{C}_0)\|_2 \le \omega \sqrt{8\log(2n/\delta)} \|\mathbf{U}\|_F \le \|\mathbf{y}\|_2 \xi \frac{\alpha^2}{\beta^2} \sqrt{8\log(2n/\delta)} \le \|\mathbf{y}\|_2.$$
 (40)

Here, the last inequality follows from  $\xi \leq 1/\sqrt{8\log(2n/\delta)}$  and  $\alpha/\beta \leq 1$ .

Verifying Assumption 1: To verify Assumption 1, note that by definition  $\mathbf{J}\mathbf{J}^T = \mathbf{\Sigma}(\mathbf{U})$  thus trivially  $\sigma_n(\mathbf{\Sigma}(\mathbf{U})) \geq \alpha$  holds. Furthermore, Lemma 4 below combined with the fact that  $\|\mathbf{v}\|_2 = 1$  implies that  $\|\mathbf{J}\| \leq \beta$  and  $\|\mathcal{J}(\mathbf{C})\| \leq \beta$  for all  $\mathbf{C}$ . This completes the verification of Assumption 1. In the next lemma we show that the Jacobian has bounded spectrum.

**Lemma 4** (Spectral norm of Jacobian). Consider  $G(\mathbf{C}) = \text{ReLU}(\mathbf{UC})\mathbf{v}$  with  $\mathbf{v} \in \mathbb{R}^k$  and  $\mathbf{U} \in \mathbb{R}^{n \times k}$  and associated Jacobian  $\mathcal{J}(\mathbf{C})$  (9), and let  $\mathbf{J}$  be any matrix obeying  $\mathbf{JJ}^T = \mathbb{E}\left[\mathcal{J}(\mathbf{C})\mathcal{J}^T(\mathbf{C})\right]$ , where the expectation is over a matrix  $\mathbf{C}$  with iid  $\mathcal{N}(0,\omega^2)$  entries. Then

$$\left\|\mathcal{J}(\mathbf{C})\right\| \leq \left\|\mathbf{v}\right\|_{2} \left\|\mathbf{U}\right\| \quad and \quad \left\|\mathbf{J}\right\| \leq \left\|\mathbf{v}\right\|_{2} \left\|\mathbf{U}\right\|.$$

**Verifying Assumption 2:** Verification of the two condition of the assumption requires the following two lemmas.

**Lemma 5** (Concentration lemma). Consider  $G(\mathbf{C}) = \text{ReLU}(\mathbf{UC})\mathbf{v}$  with  $\mathbf{v} \in \mathbb{R}^k$  and  $\mathbf{U} \in \mathbb{R}^{n \times k}$  and associated Jacobian  $\mathcal{J}(\mathbf{C})$  (9). Let  $\mathbf{C} \in \mathbb{R}^{n \times k}$  be generated at random with i.i.d.  $\mathcal{N}(0, \omega^2)$  entries. Then,

$$\left\| \mathcal{J}(\mathbf{C}) \mathcal{J}^T(\mathbf{C}) - \mathbf{\Sigma} \left( \mathbf{U} \right) \right\| \leq \left\| \mathbf{U} \right\|^2 \sqrt{\log \left( \frac{2n}{\delta} \right) \sum_{\ell=1}^k v_\ell^4},$$

holds with probability at least  $1 - \delta$ .

Note that Lemma 5 only ensures the first condition in (19). To see that it also implies the second condition in (19), we use the following lemma which establishes that the first condition implies the second.

**Lemma 6** ([Oym+19, Lem. 6.4]). Let  $\mathbf{X} \in \mathbb{R}^{n \times p}$ ,  $p \geq n$  and let  $\mathbf{B}$  be  $n \times n$  psd matrix obeying  $\|\mathbf{X}\mathbf{X}^T - \mathbf{B}\| \leq \epsilon^2$ , for a scalar  $\epsilon \geq 0$ . Then there exists a matrix  $\mathbf{Y} \in \mathbb{R}^{n \times p}$  obeying  $\mathbf{B} = \mathbf{Y}\mathbf{Y}^T$  such that

$$\|\mathbf{Y} - \mathbf{X}\| < 2\epsilon$$
.

We begin by verifying the first condition of the assumption, inequality (19). To this aim note that using the fact that  $\sum_{\ell}^{k} v_{\ell}^{4} = \frac{\nu^{4}}{k}$  by Lemma 5 we have

$$\|\mathcal{J}(\mathbf{C}_0)\mathcal{J}^T(\mathbf{C}_0) - \mathbf{\Sigma}(\mathbf{U})\| \le \|\mathbf{U}\|^2 \sqrt{\frac{\log\left(\frac{2n}{\delta}\right)}{k}} \le \epsilon_0^2,$$
 (41)

holds with probability at least  $1 - \delta$ . Thus, the first condition in (19) is satisfied. Furthermore, combining the first inequality in (41) with Lemma 6 we conclude that the second condition in (19) also holds for the chosen value of  $\epsilon_0$ , concluding the proof of Assumption 2 being satisfied.

**Verifying Assumption 3:** Verification of the assumption requires us to control the perturbation of the Jacobian matrix around a random initialization.

**Lemma 7** (Jacobian perturbation around initialization). Let  $C_0$  be a matrix with i.i.d.  $N(0, \omega^2)$  entries. Then, for all C obeying

$$\|\mathbf{C} - \mathbf{C}_0\| \le \omega \widetilde{R} \quad with \quad \widetilde{R} \le \frac{1}{2} \sqrt{k},$$

the Jacobian mapping (9) associated with the generator  $G(\mathbf{C}) = \text{ReLU}(\mathbf{UC})\mathbf{v}$  obeys

$$\|\mathcal{J}(\mathbf{C}) - \mathcal{J}(\mathbf{C}_0)\| \le \|\mathbf{v}\|_{\infty} 2(k\widetilde{R})^{1/3} \|\mathbf{U}\|,$$

with probability at least  $1 - ne^{-\frac{1}{2}\widetilde{R}^{4/3}k^{7/3}}$ .

In order to verify Assumption 3, first note that the radius in the theorem, defined in equation (23), obeys

$$R = 2 \left\| \mathbf{J}^{\dagger} \mathbf{r}_{0} \right\|_{2} + \kappa \left( \frac{1}{\alpha^{2}} + \eta T \right) \left\| \mathbf{r}_{0} \right\|_{2}$$

$$\stackrel{(i)}{\leq} \left( \frac{2}{\alpha} + \frac{\kappa}{\alpha^{2}} + \eta \kappa T \right) \left\| \mathbf{r}_{0} \right\|_{2}$$

$$\stackrel{(ii)}{\leq} 4 \left( \frac{1}{\alpha} + \frac{\xi}{\beta} T \right) \left\| \mathbf{y} \right\|_{2}$$

$$\stackrel{(iii)}{\leq} \omega \frac{\xi^{3}}{2^{15}} \frac{\alpha^{12}}{\beta^{12}} \sqrt{k}$$

$$:= \omega \widetilde{R}$$

Here, (i) follows from the fact that  $\left\|\mathbf{J}_{\mathcal{S}}^{\dagger}\mathbf{r}_{0}\right\|_{2} \leq \frac{1}{\alpha}\|\mathbf{r}_{0}\|_{2}$ , (ii) from the bound on the initial residual (39) and  $\kappa = \xi \frac{\alpha^{2}}{\beta}$ , and finally (iii) follows from the assumption (10).

Note that since  $\Delta = \frac{\xi^3}{2^{15}} \frac{\alpha^{12}}{\beta^{12}} \le \frac{1}{2}$  for this choice of radius by Lemma 7 we have

$$\|\mathcal{J}(\mathbf{C}) - \mathcal{J}(\mathbf{C}_0)\| \le \|\mathbf{v}\|_{\infty} 2(k\widetilde{R})^{1/3} \|\mathbf{U}\| = \frac{\xi}{16} \|\mathbf{U}\| \frac{\alpha^4}{\beta^4},$$

holds with probability at least

$$1 - ne^{-\frac{1}{2^{21}}\xi^4 \frac{\alpha^{16}}{\beta^{16}}k^{25/6}} \stackrel{\text{(i)}}{\ge} 1 - \delta,$$

where in (i) we used (10) together with  $\xi \leq \sqrt{8 \log\left(\frac{2n}{\delta}\right)}$ . Therefore, Assumption 3 holds with high probability by our choice of  $\epsilon = \frac{\xi}{8} \|\mathbf{U}\| \frac{\alpha^4}{\beta^4}$ .

Verifying the conditions of (24): We now turn our attention to verifying that the bounds on  $\epsilon$  and  $\epsilon_0$  in (24) are satisfied. We begin by showing the bound on  $\epsilon$  holds. To this aim note that our choice of  $\epsilon$  can alternatively be written as  $\epsilon = \frac{\xi}{8} \frac{\alpha^4}{\beta^3}$  so that the condition  $\epsilon \leq \frac{1}{8} \frac{\kappa \alpha^2}{\beta^2}$  on  $\epsilon$  is satisfied (recall that  $\kappa = \xi \frac{\alpha^2}{\beta}$ ). To verify the second condition in (24) note that

$$\epsilon_0 = 2\nu \|\mathbf{U}\| \sqrt[4]{\frac{\log(2n/\delta)}{k}} = 2\beta \sqrt[4]{\frac{\log(2n/\delta)}{k}} \stackrel{\text{(i)}}{\leq} \frac{1}{4}\kappa \stackrel{\text{(ii)}}{=} \min\left(\frac{1}{4}\kappa, \sqrt{\frac{1}{8}\frac{\kappa\alpha^2}{\beta}}\right),$$

where inequality (i) holds by (10) which implies  $8\sqrt[4]{\log(2n/\delta)/k} \le \xi \frac{\alpha^2}{\beta^2} = \frac{\kappa}{\beta}$  (ii) holds by the assumption  $\xi \le 2$ .

Verifying the bound on number of iterations: Finally we note that the constraint on the number of iterations,  $T, T \leq \frac{1}{2\eta\epsilon^2}$ , from Theorem 4 is satisfied under the number of constraints of Theorem 3 by  $\frac{1}{2\eta\epsilon^2} = 16 \frac{\beta^6}{\eta\xi^2\alpha^8} \geq 16 \frac{\beta^8}{\eta\xi^2\alpha^8}$ .

Concluding the proof of Theorem 3: Now that we have verified the conditions of Theorem 4 we can apply this theorem. This allows us to conclude that

$$\begin{split} \|G(\mathbf{C}_{\tau}) - \mathbf{x}\|_{2} &= \|G(\mathbf{C}_{\tau}) + \mathbf{z} - \mathbf{y}\|_{2} \\ &= \|\tilde{\mathbf{r}}_{\tau} + \mathbf{z} + \mathbf{r}_{\tau} - \tilde{\mathbf{r}}_{\tau}\|_{2} \\ &\stackrel{(i)}{\leq} \|\tilde{\mathbf{r}}_{\tau} + \mathbf{z} + \mathbf{r}_{\tau} - \tilde{\mathbf{r}}_{\tau}\|_{2} \\ &\stackrel{(ii)}{\leq} \|\tilde{\mathbf{r}}_{\tau} + \mathbf{z}\|_{2} + \|\mathbf{r}_{\tau} - \tilde{\mathbf{r}}_{\tau}\|_{2} \\ &\stackrel{(iii)}{\leq} \|\mathbf{W} (\mathbf{I} - \eta \boldsymbol{\Sigma}^{2})^{T} \mathbf{W}^{T} \mathbf{r}_{0} + \mathbf{z}\|_{2} + \frac{1}{2} \xi \frac{\alpha^{2}}{\beta^{2}} \|\mathbf{r}_{0}\|_{2} \\ &\stackrel{(iii)}{=} \|\mathbf{W} (\mathbf{I} - \eta \boldsymbol{\Sigma}^{2})^{T} \mathbf{W}^{T} (G(\mathbf{C}_{0}) - \mathbf{x}) - (\mathbf{W} (\mathbf{I} - \eta \boldsymbol{\Sigma}^{2})^{T} \mathbf{W}^{T} - \mathbf{I}) \mathbf{z}\|_{2} + \frac{1}{2} \xi \frac{\alpha^{2}}{\beta^{2}} \|\mathbf{r}_{0}\|_{2} \\ &\stackrel{(v)}{\leq} \|(\mathbf{I} - \eta \boldsymbol{\Sigma}^{2})^{T} \mathbf{W}^{T} \mathbf{x}\|_{2} + \|((\mathbf{I} - \eta \boldsymbol{\Sigma}^{2})^{T} - \mathbf{I}) \mathbf{W}^{T} \mathbf{z}\|_{2} + \|G(\mathbf{C}_{0})\|_{2} + \frac{1}{2} \xi \frac{\alpha^{2}}{\beta^{2}} \|\mathbf{r}_{0}\|_{2} \\ &\stackrel{(v)}{\leq} (1 - \eta \sigma_{p}^{2})^{\tau} \|\mathbf{x}\|_{2} + \sqrt{\sum_{i=1}^{n} ((1 - \eta \sigma_{i}^{2})^{\tau} - 1)^{2} \langle \mathbf{w}_{i}, \mathbf{z} \rangle^{2} + 2 \|\mathbf{y}\|_{2} \xi \frac{\alpha^{2}}{\beta^{2}} \sqrt{8 \log(2n/\delta)}. \end{split}$$

Here, (i) follows from the triangular inequality, (ii) from Theorem 4 equation (25), (iii) from Theorem 5, (iv) from  $\mathbf{r}_0 = G(\mathbf{C}_0) - \mathbf{y} = G(\mathbf{C}_0) - \mathbf{z} - \mathbf{z}$ , (v) from the triangular inequality, (vi) from the fact that  $\mathbf{x} \in \text{span}\{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_p\}$  combined with the fact that  $\|\mathbf{W} (\mathbf{I} - \eta \mathbf{\Sigma}^2)^T \mathbf{W}^T\| \leq 1$ . Finally, (vi) follows by bounding the last two terms with the bounds (39) and (40), respectively. This proves the final bound (12), as desired.

# F Proof of Proposition 1 and equation (2)

The residual of gradient descent at iteration  $\tau$  is

$$\mathbf{r}_{\tau} = \mathbf{y} - \mathbf{J}\mathbf{c}_{\tau}$$

$$= \mathbf{y} - \mathbf{J}(\mathbf{x}_{\tau-1} - \eta \mathbf{J}^{T}(\mathbf{J}\mathbf{c}_{\tau-1} - \mathbf{y}))$$

$$= (\mathbf{I} - \eta \mathbf{J}\mathbf{J}^{T})(\mathbf{y} - \mathbf{J}\mathbf{c}_{\tau-1})$$

$$= (\mathbf{I} - \eta \mathbf{J}\mathbf{J}^{T})^{\tau}(\mathbf{y} - \mathbf{J}\mathbf{c}_{0})$$

$$= (\mathbf{I} - \eta \mathbf{J}\mathbf{J}^{T})^{\tau}\mathbf{y}$$

$$= (\mathbf{I} - \eta \mathbf{W}\Sigma^{2}\mathbf{W}^{T})^{\tau}\mathbf{y}$$

where we used that  $\mathbf{c}_0 = 0$  and the SVD  $\mathbf{A} = \mathbf{W} \Sigma \mathbf{V}^T$ . Expanding  $\mathbf{y}$  in terms of the singular vectors  $\mathbf{w}_i$  (i.e., the columns of  $\mathbf{W}$ ), as  $\mathbf{y} = \sum_i \mathbf{w}_i \langle \mathbf{w}_i, \mathbf{y} \rangle$ , and noting that  $(\mathbf{I} - \eta \mathbf{W} \Sigma^2 \mathbf{W}^T)^{\tau} = \sum_i (1 - \eta \sigma_i^2)^{\tau} \mathbf{w}_i \mathbf{w}_i^T$  we get

$$\mathbf{r}_{\tau} = \sum_{i} (1 - \eta \sigma_{i}^{2})^{\tau} \mathbf{w}_{i} \left\langle \mathbf{w}_{i}, \mathbf{y} \right\rangle,$$

as desired.

**Proof of equation** (2): By proposition 1 and using that x and lies in the signal subspace

$$\mathbf{x} - \mathbf{J}\mathbf{c}_{\tau} = \sum_{i=1}^{p} \mathbf{w}_{i} (1 - \eta \sigma_{i}^{2})^{\tau} \langle \mathbf{w}_{i}, \mathbf{x} \rangle + \sum_{i=1}^{n} \mathbf{w}_{i} ((1 - \eta \sigma_{i}^{2})^{\tau} - 1) \langle \mathbf{w}_{i}, \mathbf{z} \rangle.$$

By the triangle inequality,

$$\|\mathbf{x} - \mathbf{J}\mathbf{c}_{\tau}\|_{2} \leq \left\| \sum_{i=1}^{p} \mathbf{w}_{i} (1 - \eta \sigma_{i}^{2})^{\tau} \left\langle \mathbf{w}_{i}, \mathbf{x} \right\rangle \right\|_{2} + \left\| \sum_{i=1}^{n} \mathbf{w}_{i} ((1 - \eta \sigma_{i}^{2})^{\tau} - 1) \left\langle \mathbf{w}_{i}, \mathbf{z} \right\rangle \right\|_{2}$$

$$\leq (1 - \eta \sigma_{p}^{2})^{\tau} \|\mathbf{x}\|_{2} + \sqrt{\sum_{i=1}^{n} ((1 - \eta \sigma_{i}^{2})^{\tau} - 1)^{2} \left\langle \mathbf{w}_{i}, \mathbf{z} \right\rangle^{2}},$$

where the second inequality follows by using orthogonality of the  $\mathbf{w}_i$  and by using  $(1 - \eta \sigma_i^2) \le 1$ , from  $\eta \le 1/\sigma_{\max}^2$ . This concludes the proof of equation (2).

# G The neural tangent kernel for convolutional generators

We first prove the closed form expression of the neural tangent kernel, i.e., equation (13). By the expression for the Jacobian in equation (9), we have that

$$\mathcal{J}(\mathbf{C})\mathcal{J}^{T}(\mathbf{C}) = \sum_{\ell=1}^{k} v_{\ell}^{2} \operatorname{diag}(\sigma'(\mathbf{U}\mathbf{c}_{\ell})) \mathbf{U}\mathbf{U}^{T} \operatorname{diag}(\sigma'(\mathbf{U}\mathbf{c}_{\ell}))$$

$$= \sum_{\ell=1}^{k} v_{\ell}^{2} \sigma'(\mathbf{U}\mathbf{c}_{\ell}) \sigma'(\mathbf{U}\mathbf{c}_{\ell})^{T} \odot \mathbf{U}\mathbf{U}^{T}$$

$$= \sigma'(\mathbf{U}\mathbf{C}) \operatorname{diag}(v_{1}^{2}, \dots, v_{k}^{2}) \sigma'(\mathbf{U}\mathbf{C})^{T} \odot \mathbf{U}\mathbf{U}^{T},$$

where  $\odot$  denotes the entrywise product of the two matrices. Then,

$$\mathbb{E}\left[\mathcal{J}(\mathbf{C})\mathcal{J}^{T}(\mathbf{C})\right] = \sum_{\ell=1}^{k} v_{\ell}^{2} \mathbb{E}\left[\sigma'(\mathbf{U}\mathbf{c}_{\ell})\sigma'(\mathbf{U}\mathbf{c}_{\ell})^{T}\right] \odot \mathbf{U}\mathbf{U}^{T}.$$
(42)

Next, we have with [Dan+16, Sec. 4.2] and using that the derivative of the ReLU function is the step function,

$$\left[ \mathbb{E} \left[ \sigma'(\mathbf{U}\mathbf{c}_{\ell})\sigma'(\mathbf{U}\mathbf{c}_{\ell})^T \right] \right]_{ij} = \frac{1}{2} \left( 1 - \cos^{-1} \left( \frac{\langle \mathbf{u}_i, \mathbf{u}_j \rangle}{\|\mathbf{u}_i\|_2 \|\mathbf{u}_j\|_2} \right) / \pi \right).$$

Using that  $\|\mathbf{v}\|_2 = 1$ , we get

$$\left[\mathbb{E}\left[\mathcal{J}(\mathbf{C})\mathcal{J}^{T}(\mathbf{C})\right]\right]_{ij} = \frac{1}{2}\left(1 - \cos^{-1}\left(\frac{\langle \mathbf{u}_{i}, \mathbf{u}_{j} \rangle}{\|\mathbf{u}_{i}\|_{2}\|\mathbf{u}_{j}\|_{2}}\right) / \pi\right) \langle \mathbf{u}_{i}, \mathbf{u}_{j} \rangle,$$

where  $\mathbf{u}_i$  are the rows of **U**. This concludes the proof of equation (13).

We next briefly comment on the singular value decomposition of a circulant matrix and explain that the singular vectors are given by Definition 1. Recall, that  $\mathbf{U} \in \mathbb{R}^{n \times n}$  is a circulant matrix, implementing the convolution with a filter. Assume for simplicity that n is even. It is well known that the discrete Fourier transform diagonalizes  $\mathbf{U}$ , i.e.,

$$\mathbf{U} = \mathbf{F}^{-1} \hat{\mathbf{U}} \mathbf{F},$$

where  $\mathbf{F} \in \mathbb{C}^{n \times n}$  is the DFT matrix with entries

$$[\mathbf{F}]_{jk} = e^{i2\pi jk/n}, \quad j, k = 0, \dots, n-1,$$

and  $\hat{\mathbf{U}}$  is a diagonal matrix with diagonal  $\mathbf{F}\mathbf{u}$ , where  $\mathbf{u}$  is the first column of the circulant matrix  $\mathbf{U}$ . From this, we can compute the singular value decomposition of  $\mathbf{U}$  by using that  $\hat{\mathbf{u}} = \mathbf{F}\mathbf{u}$  is conjugate symmetric (since  $\mathbf{u}$  is real) so that

$$[\hat{\mathbf{u}}]_{n-k+2} = \hat{\mathbf{u}}^*, \quad k = 2, \dots, n/2.$$

Let  $\mathbf{U} = \mathbf{W} \Sigma \mathbf{V}^T$  be the singular value decomposition of  $\mathbf{U}$ . The entries of the left singular vectors are given by the trigonometric basis functions defined in (5), and the singular values are given by the absolute values of  $\hat{\mathbf{u}}$ .

# H Proofs of Lemmas for neural network denoisers (Proofs of auxiliary lemmas in Section E)

#### H.1 Proof of Lemma 7: Jacobian perturbation around initialization

The proof follows that of [OS19, Lem. 6.9].

1. We start by relating the perturbation of the Jacobian to a perturbation of the activation patterns. For any C, C', we have that

$$\|\mathcal{J}(\mathbf{C}) - \mathcal{J}(\mathbf{C}')\| \le \|\mathbf{v}\|_{\infty} \|\mathbf{U}\| \max_{j} \|\sigma'(\mathbf{u}_{j}^{T}\mathbf{C}) - \sigma'(\mathbf{u}_{j}^{T}\mathbf{C}')\|^{2}.$$
(43)

To see this, first note that, by (9),

$$\mathcal{J}(\mathbf{C}) - \mathcal{J}(\mathbf{C}') = [\dots v_j(\operatorname{diag}(\sigma'(\mathbf{U}\mathbf{c}_j)) - \operatorname{diag}(\sigma'(\mathbf{U}\mathbf{c}_j))')\mathbf{U}\dots].$$

This in turn implies that

$$\begin{split} \left\| \mathcal{J}(\mathbf{C}) - \mathcal{J}(\mathbf{C}') \right\|^2 &= \left\| (\mathcal{J}(\mathbf{C}) - \mathcal{J}(\mathbf{C}'))(\mathcal{J}(\mathbf{C}) - \mathcal{J}(\mathbf{C}'))^T \right\| \\ &= \left\| (\sigma'(\mathbf{U}\mathbf{C}) - \sigma'(\mathbf{U}\mathbf{C}')) \mathrm{diag}(v_1^2, \dots, v_k^2)(\sigma'(\mathbf{U}\mathbf{C}) - \sigma'(\mathbf{U}\mathbf{C}'))^T \odot \mathbf{U}\mathbf{U}^T \right\|^2 \\ &\stackrel{(\mathrm{i})}{\leq} \left\| \mathbf{U} \right\|^2 \max_j \left\| (\sigma'(\mathbf{u}_j^T \mathbf{C}) - \sigma'(\mathbf{u}_j^T \mathbf{C}')) \mathrm{diag}(\mathbf{v}) \right\|^2 \\ &\leq \left\| \mathbf{v} \right\|_{\infty}^2 \left\| \mathbf{U} \right\|^2 \max_j \left\| \sigma'(\mathbf{u}_j^T \mathbf{C}) - \sigma'(\mathbf{u}_j^T \mathbf{C}') \right\|^2, \end{split}$$

where for (i) we used that for two positive semidefinite matrices  $\mathbf{A}, \mathbf{B}, \lambda_{\max}(\mathbf{A} \odot \mathbf{B}) \leq \lambda_{\max}(\mathbf{A}) \max_{i} \mathbf{B}_{ii}$ . This concludes the proof of equation (43).

2. Step one implies that we only need to control  $\sigma'(\mathbf{U}\mathbf{c}_j)$  around a neighborhood of  $\mathbf{c}'_j$ . Since  $\sigma'$  is the step function, we need to count the number of sign flips between the matrices  $\mathbf{U}\mathbf{C}$  and  $\mathbf{U}\mathbf{C}'$ . Let  $|\mathbf{v}|_{\pi(q)}$  be the q-th smallest entry of  $\mathbf{v}$  in absolute value.

**Lemma 8.** Suppose that, for all i, and  $q \leq k$ ,

$$\|\mathbf{C} - \mathbf{C}'\| \le \sqrt{q} \frac{|\mathbf{u}_i^T \mathbf{C}'|_{\pi(q)}}{\|\mathbf{u}_i\|}.$$

Then

$$\max_{i} \left\| \sigma'(\mathbf{u}_{i}^{T}\mathbf{C}) - \sigma'(\mathbf{u}_{i}^{T}\mathbf{C}') \right\| \leq \sqrt{2q}$$

*Proof.* Suppose that  $\sigma'(\mathbf{u}_i^T\mathbf{C})$  and  $\sigma'(\mathbf{u}_i^T\mathbf{C}')$  have 2q many different entries, then the conclusion of the statement would be violated. We show that this implies the assumption is violated as well, proving the statement by contraction. By the contradiction hypothesis,

$$\|\mathbf{C} - \mathbf{C}'\|^2 \ge \|\mathbf{u}_i^T (\mathbf{C} - \mathbf{C}')\|^2 / \|\mathbf{u}_i\|^2$$
$$\ge q \frac{|\mathbf{u}_i^T \mathbf{C}'|_{\pi(q)}^2}{\|\mathbf{u}_i\|^2},$$

where the last inequality follows by noting that at least 2q many entries have different signs, thus their difference is larger than their individual magnitudes, and at least q many individual magnitudes are lower bounded by the q-th smallest one.

3. Next, we note that, with probability at least  $1-ne^{-kq^2/2}$ , the q-th smallest entry of  $\mathbf{u}_i^T \mathbf{C}' \in \mathbb{R}^k$  obeys

$$\frac{|\mathbf{u}_i^T \mathbf{C}'|_{\pi(q)}}{\|\mathbf{u}_i\|} \ge \frac{q}{2k} \nu \quad \text{for all } i = 1, \dots, n.$$
(44)

We note that it is sufficient to prove this result for  $\nu = 1$ . This follows from anti-concentration of Gaussian random variables. Specifically, with the entries of  $\mathbf{C}'$  being iid  $\mathcal{N}(0,1)$  distributed, the entries of  $\mathbf{g} = \mathbf{u}_i^T \mathbf{C}' / \|\mathbf{u}_i\| \in \mathbb{R}^k$  are iid standard Gaussian random variables as well. We show that with probability at least  $1 - e^{-kq^2/2}$ , at most q entries are larger than  $\frac{q}{2k}$ . Let  $\gamma_{\delta}$  be the number for which  $P[|g_{\ell}| \leq \gamma_{\delta}] \leq \delta$ , where q is a standard Gaussian random variable. Note that  $\gamma_{\delta} \geq \sqrt{\pi/2} \delta \geq \delta$ . Define the random variable

$$\delta_{\ell} = \begin{cases} 1 & \text{if } |g_{\ell}| \leq \gamma_{\delta}, \\ 0 & \text{otherwise.} \end{cases}$$

with  $\delta = \frac{q}{2k}$ . With  $\mathbb{E}[\delta_{\ell}] = \delta$ , by Hoeffding's inequality,

$$P\left[\sum_{\ell=1}^{k} \delta_{\ell} \ge m\right] = P\left[\sum_{\ell=1}^{k} \delta_{\ell} - \mathbb{E}\left[\delta_{\ell}\right] \ge m/2\right] \le e^{-2k(m/2)^{2}} = e^{-km^{2}/2}.$$
 (45)

Thus, with probability at least  $1 - ke^{-kq^2/2}$  no more than m entries are smaller than  $\gamma_{\delta} \ge \delta = \frac{q}{2k}$ . The results now follows by taking the union bound over all  $i = 1, \ldots, n$ .

We are now ready to conclude the proof of the lemma. By equation (43),

$$\|\mathcal{J}(\mathbf{C}) - \mathcal{J}(\mathbf{C}')\| \le \|\mathbf{v}\|_{\infty} \|\mathbf{U}\| \max_{j} \|\sigma'(\mathbf{u}_{j}^{T}\mathbf{C}) - \sigma'(\mathbf{u}_{j}^{T}\mathbf{C}')\|$$
$$\le \|\mathbf{v}\|_{\infty} \|\mathbf{U}\| \sqrt{2q}$$

provided that

$$\|\mathbf{C} - \mathbf{C}'\| \le \sqrt{q} \frac{q}{2k} \nu,$$

with probability at least  $1 - ne^{-kq^2/2}$ . Setting  $q = (2kR)^{2/3}$  concludes the proof (note that the assumption  $R \leq \frac{1}{2}\sqrt{k}$  ensures  $q \leq k$ ).

## H.2 Proof of Lemma 4: Bounded Jacobian

By the expression of the Jacobian in equation (9),

$$\begin{split} \|\mathcal{J}(\mathbf{C})\|^2 &= \left\| \mathcal{J}(\mathbf{C})\mathcal{J}(\mathbf{C})^T \right\| \\ &= \left\| \sigma'(\mathbf{U}\mathbf{C}) \operatorname{diag}(v_1^2, \dots, v_k^2) \sigma'(\mathbf{U}\mathbf{C})^T \odot \mathbf{U}\mathbf{U}^T \right\|^2 \\ &\stackrel{\text{(i)}}{\leq} \|\mathbf{U}\|^2 \max_j \left\| \sigma'(\mathbf{u}_j^T \mathbf{C}) \operatorname{diag}(\mathbf{v}) \right\|_2^2 \\ &\leq \|\mathbf{v}\|_2^2 \|\mathbf{U}\|^2 \,, \end{split}$$

where for (i) we used that for two positive semidefinite matrices  $\mathbf{A}, \mathbf{B}, \lambda_{\max}(\mathbf{A} \odot \mathbf{B}) \leq \lambda_{\max}(\mathbf{A}) \max_i \mathbf{B}_{ii}$ . To prove the second inequality note that

$$\begin{aligned} \left\| \mathbf{J} \mathbf{J}^{T} \right\| &= \left\| \mathbf{\Sigma}(\mathbf{U}) \right\| \\ &= \left\| \mathbf{v} \right\|_{2}^{2} \left\| \mathbb{E}_{\mathbf{c} \sim \mathcal{N}(0,1)} \left[ \operatorname{ReLU}'(\mathbf{U}\mathbf{c}) \operatorname{ReLU}'(\mathbf{U}\mathbf{c})^{T} \right] \odot \left( \mathbf{U}\mathbf{U}^{T} \right) \right\| \\ &\stackrel{(i)}{\leq} \left\| \mathbf{v} \right\|_{2}^{2} \left\| \mathbf{U}\mathbf{U}^{T} \right\| \\ &= \left\| \mathbf{v} \right\|_{2}^{2} \left\| \mathbf{U}\mathbf{U}^{T} \right\|. \end{aligned}$$

Here, (i) follows from the fact that for two positive semidefinite matrices  $\mathbf{A}, \mathbf{B}, \lambda_{\max}(\mathbf{A} \odot \mathbf{B}) \leq \lambda_{\max}(\mathbf{A}) \max_i \mathbf{B}_{ii}$ .

## H.3 Proof of Lemma 5: Concentration lemma

We begin by defining the zero-mean random matrices

$$\mathbf{S}_{\ell} = v_{\ell}^{2} \left( \sigma' (\widetilde{\mathbf{U}} \mathbf{c}_{\ell}) \sigma' (\widetilde{\mathbf{U}} \mathbf{c}_{\ell})^{T} - \mathbb{E} \left[ \sigma' (\widetilde{\mathbf{U}} \mathbf{c}_{\ell}) \sigma' (\widetilde{\mathbf{U}} \mathbf{c}_{\ell})^{T} \right] \right) \odot \left( \widetilde{\mathbf{U}} \widetilde{\mathbf{U}}^{T} \right).$$

With this notation,

$$\begin{split} \mathcal{J}(\mathbf{C})\mathcal{J}^T(\mathbf{C}) - \mathbf{\Sigma}(\widetilde{\mathbf{U}}) &= \sum_{\ell=1}^k v_\ell^2 \left( \sigma'(\widetilde{\mathbf{U}} \mathbf{c}_\ell) \sigma'(\widetilde{\mathbf{U}} \mathbf{c}_\ell)^T - \mathbb{E} \left[ \sigma'(\widetilde{\mathbf{U}} \mathbf{c}_\ell) \sigma'(\widetilde{\mathbf{U}} \mathbf{c}_\ell)^T \right] \right) \odot \left( \widetilde{\mathbf{U}} \widetilde{\mathbf{U}}^T \right) \\ &= \sum_{\ell=1}^k \mathbf{S}_\ell. \end{split}$$

To show concentration we use the matrix Hoeffding inequality. To this aim note that the summands are centered in the sense that  $\mathbb{E}\left[\mathbf{S}_{\ell}\right] = \mathbf{0}$ . Next note that

$$\left(\sigma'(\widetilde{\mathbf{U}}\mathbf{c}_{\ell})\sigma'(\widetilde{\mathbf{U}}\mathbf{c}_{\ell})^{T} - \mathbb{E}\left[\sigma'(\widetilde{\mathbf{U}}\mathbf{c}_{\ell})\sigma'(\widetilde{\mathbf{U}}\mathbf{c}_{\ell})^{T}\right]\right) \odot \left(\widetilde{\mathbf{U}}\widetilde{\mathbf{U}}^{T}\right) \leq \left(\sigma'(\widetilde{\mathbf{U}}\mathbf{c}_{\ell})\sigma'(\widetilde{\mathbf{U}}\mathbf{c}_{\ell})^{T}\right) \odot \left(\widetilde{\mathbf{U}}\widetilde{\mathbf{U}}^{T}\right) \\
= \operatorname{diag}\left(\sigma'(\widetilde{\mathbf{U}}\mathbf{c}_{\ell})^{T}\right) \widetilde{\mathbf{U}}\widetilde{\mathbf{U}}^{T} \operatorname{diag}\left(\sigma'(\widetilde{\mathbf{U}}\mathbf{c}_{\ell})\right) \\
\leq B^{2}\widetilde{\mathbf{U}}\widetilde{\mathbf{U}}^{T}$$

Similarly,

$$\left(\sigma'(\widetilde{\mathbf{U}}\mathbf{c}_{\ell})\sigma'(\widetilde{\mathbf{U}}\mathbf{c}_{\ell})^{T} - \mathbb{E}\left[\sigma'(\widetilde{\mathbf{U}}\mathbf{c}_{\ell})\sigma'(\widetilde{\mathbf{U}}\mathbf{c}_{\ell})^{T}\right]\right) \odot \left(\widetilde{\mathbf{U}}\widetilde{\mathbf{U}}^{T}\right) \succeq -\mathbb{E}\left[\sigma'(\widetilde{\mathbf{U}}\mathbf{c}_{\ell})\sigma'(\widetilde{\mathbf{U}}\mathbf{c}_{\ell})^{T}\right] \odot \left(\widetilde{\mathbf{U}}\widetilde{\mathbf{U}}^{T}\right) \\
\succeq -B^{2}\widetilde{\mathbf{U}}\widetilde{\mathbf{U}}^{T}.$$

Thus,

$$-v_{\ell}^2 B^2 \widetilde{\mathbf{U}} \widetilde{\mathbf{U}}^T \preceq \mathbf{S}_{\ell} \preceq v_{\ell}^2 B^2 \widetilde{\mathbf{U}} \widetilde{\mathbf{U}}^T.$$

Therefore, using  $\mathbf{A}_{\ell} := v_{\ell}^2 B^2 \widetilde{\mathbf{U}} \widetilde{\mathbf{U}}^T$  we have

$$\mathbf{S}_{\ell}^2 \preceq \mathbf{A}_{\ell}^2$$
,

and

$$\sigma^2 := \left\| \sum_{\ell=1}^k \mathbf{A}_\ell^2 \right\| \le B^4 \left( \sum_{\ell=1}^k v_\ell^4 \right) \left\| \widetilde{\mathbf{U}} \right\|^4$$

To continue we will apply matrix Hoeffding inequality stated below.

**Theorem 6** (Matrix Hoeffding inequality, Theorem 1.3 [Tro11]). Consider a finite sequence  $\mathbf{S}_{\ell}$  of independent, random, self-adjoint matrices with dimension n, and let  $\{\mathbf{A}_{\ell}\}$  be a sequence of fixed self-adjoint matrices. Assume that each random matrix satisfies

$$\mathbb{E}[\mathbf{S}_{\ell}] = \mathbf{0}$$
 and  $\mathbf{S}_{\ell}^2 \preceq \mathbf{A}_{\ell}^2$  almost surely.

Then, for all  $t \geq 0$ ,

$$P\left[\left\|\sum_{\ell=1}^{k} \mathbf{S}_{\ell}\right\| \ge t\right] \le 2ne^{-\frac{t^2}{8\sigma^2}} \quad where \quad \sigma^2 := \left\|\sum_{\ell=1}^{k} \mathbf{A}_{\ell}^2\right\|.$$

Therefore, applying matrix Hoeffding inequality we get that

$$P\left[\left\|\sum_{\ell=1}^{k} \mathbf{S}_{\ell}\right\| \ge t\right] \le 2ne^{-\frac{t^2}{B^4 \|\mathbf{v}\|_4^4 \|\tilde{\mathbf{U}}\|^4}},$$

which concludes the proof.

## H.4 Proof of Lemma 3 (bound on initial residual)

Without loss of generality we prove the result for  $\nu = 1$ . First note that by the triangle inequality

$$\|\mathbf{r}_0\|_2 = \|\sigma(\mathbf{UC})\mathbf{v} - \mathbf{y}\|_2 \le \|\sigma(\mathbf{UC})\mathbf{v}\|_2 + \|\mathbf{y}\|_2.$$

We next bound  $\|\sigma(\mathbf{UC})\mathbf{v}\|_2$ . Consider the *i*-th entry of the vector  $\sigma(\mathbf{UC})\mathbf{v} \in \mathbb{R}^n$ , given by  $\sigma(\mathbf{u}_i^T\mathbf{C})\mathbf{v}$ , and note that  $q_j = (\sigma(\mathbf{u}_i^T\mathbf{c}_j) - \sigma(\mathbf{u}_i^T\mathbf{c}_{n-j}))/\|\mathbf{u}_i\|_2$  is sub-Gaussian with parameter 2, i.e.,  $P[|q_j| \geq t] \leq 2e^{-2t^2}$ . It follows that

$$P\left[\left|\sum_{j=1}^{k/2} q_j\right| \ge \beta\sqrt{k}\right] \le 2e^{-\frac{\beta^2}{8}}.$$

Thus,

$$P\left[\left|\sigma(\mathbf{u}_i^T\mathbf{C})\mathbf{v}\right| \ge \|\mathbf{u}_i\|_2 \xi \beta\right] \le 2e^{-\frac{\beta^2}{8}},$$

where we used that  $|v_j| = \xi/\sqrt{k}$ . Taking a union bound over all n entries,

$$\mathbf{P}\left[\left\|\sigma(\mathbf{U}\mathbf{C})\mathbf{v}\right\|_{2}^{2} \geq \left\|\mathbf{U}\right\|_{F}^{2} \xi^{2} \beta^{2}\right] \leq 2ne^{-\frac{\beta^{2}}{8}}.$$

Choosing  $\beta = \sqrt{8\log(2n/\delta)}$  concludes the proof.