# A signal processing interpretation of noise-reduction convolutional neural networks

Luis A. Zavala-Mondragón *Student member IEEE*, Peter H.N. de With, *Fellow IEEE*, and Fons van der Sommen, *Member IEEE*

*Abstract*—Encoding-decoding CNNs play a central role in data-driven noise reduction and can be found within numerous deep-learning algorithms. However, the development of these CNN architectures is often done in *ad-hoc* fashion and theoretical underpinnings for important design choices is generally lacking. Up to this moment there are different existing relevant works that strive to explain the internal operation of these CNNs. Still, these ideas are either scattered and/or may require significant expertise to be accessible for a bigger audience. In order to open up this exciting field, this article builds intuition on the theory of deep convolutional framelets and explains diverse ED CNN architectures in a unified theoretical framework. By connecting basic principles from signal processing to the field of deep learning, this self-contained material offers significant guidance for designing robust and efficient novel CNN architectures.

*Index Terms*—Denoising, convolutional neural network, encoding-decoding.

## I. INTRODUCTION

A well-known image-processing application is noise/artifact reduction of images, which is consists in estimating a noise/artifact-free signal out of a noisy observation. In order to achieve this, conventional signal processing algorithms often employ explicit assumptions on the signal and noise characteristics, which has resulted in well-known algorithms such as wavelet shrinkage [1], sparse dictionaries [2], total-variation minimization [3] and low-rank approximation [4]. With the advent of deep learning techniques, signal processing algorithms applied to image denoising, have been regularly outperformed and increasingly replaced by encoding-decoding convolutional neural networks (CNNs).

In this article, rather than conventional signal processing algorithms, we focus on the so-called encoding-decoding CNNs. These models contain an *encoder* which maps the input to a multi-channel/redundant representations and a *decoder*, which maps the encoded signal back to the original domain. In both, the encoder and decoder, sparsifying non-linearities which suppress parts of the signal are applied. In contrast with conventional signal processing algorithms, encoding-decoding CNNs are often presented as a solution which does not make explicit assumptions on the signal and noise. For example, in supervised algorithms, an encoding-decoding CNN *learns* the optimal parameters to filter the signal from a set of paired

Luis A. Zavala-Mondragón, Fons van der Sommen and P.H.N de With are with the Electrical Engineering department of the Eindhoven University of Technology, VCA lab, Flux 5, Eindhoven (corresponding e-mail: lzavala905@gmail.com).

examples of noise/artifact-free images and images contaminated with noise/artifacts [5], [6], [7], which highly simplifies the solution of the noise-reduction problems, since this circumvents the use of explicit modeling of the signal and noise. Furthermore, the good performance and simple use of encoder-decoder CNNs/autoencoders have enabled additional data-driven noise-reduction algorithms, where CNNs are embedded as part of a larger system. Examples of such approaches are unsupervised noise-reduction [8], denoising based on generative adversarial networks [9]. Besides this, smoothness in signals can be obtained also by advanced regularization using CNNs, e.g. by exploiting data-driven model-based iterative reconstruction [10].

Despite of the impressive noise-reduction performance and flexibility of encoding-decoding convolutional neural networks, these models have also downsides that should be considered. First, the complexity and heuristic nature of such designs often offers restricted understanding of the internal operation of such architectures [11]. Second, the training and deployment of CNNs requires specialized hardware and the use of significant computational resources. Third and final, the restricted understanding of the signal modeling in encoding-decoding CNNs does not clearly reveal the limitations of such models and, consequently, it is not obvious how to overcome these problems.

In order to overcome the limitations of encoding-decoding CNNs, new research has tackled the lack of explainability of these models by acknowledging the similarity of the building blocks of encoding-decoding CNNs applied to image noise reduction and the elements of well-known signal processing algorithms, such as wavelet decomposition, low-rank approximation [12], [13], [14], variational methods [15], lower-dimensional manifolds [8] and convolutional sparse coding [16]. Furthermore, practical works based on shrinkage-based CNNs inspired in well-established wavelet shrinkage algorithms has further deepened the connections between signal processing and CNNs [17], [18]. This unified treatment of signal processing-inspired CNNs has resulted in more explainable [6], [8], better performing [6] and more memory-efficient designs [19].

This article has three main objectives. First, to summarize the diverse explanations of the components of encoding-decoding convolutional neural networks applied to image noise reduction based on the concept of deep convolutional framelets [12] and on elementary signal processing concepts.

Both aspects are considered with the aim of achieving an in-depth understanding of the internal operation of encoding-decoding CNNs and to show that the design choices have *implicit* assumptions about the signal behavior inside the CNN. A second objective is to offer practitioners tools for optimizing their CNN designs with signal processing concepts. Third and final, the aim is to show practical use cases, where existing CNNs are analyzed in a unified framework, thereby enabling a better comparison of different designs by making their internal operation explicitly visible. Our analysis are based on existing works [12], [6], [20], who analyzed CNNs where the non-linearities are ignored. In this article, we overcome this limitation and present a complete analysis including the non-linear activations, which reveals important assumptions implicit in the analyzed models.

The structure of this article is as follows. Section II introduces the notation that is used in this text. Section III describes the signal model and the architecture of encoding-decoding networks. Afterwards, Section IV addresses fundamental aspects of signal processing, such as singular value decomposition, low-rank approximation, framelets, as well as the estimation of signals in the framelet domain. All the concepts of Sections III and IV converge in Section V, where the encoding-decoding CNNs are interpreted in terms of a data-driven low-rank approximation and of wavelet shrinkage. Afterwards, based on the learnings from Section V, Section VI shows the analysis of diverse architectures from a signal processing perspective and under a set of explicit assumptions. Afterwards, Section VII explores if some of the theoretical properties exposed here are related to trained models. Based on the diverse described models and theoretical operation of CNNs, Section VIII addresses a design criterion which can be used to design or choose new models and briefly describes the state-of-the art for noise reduction with CNNs. Finally, Section IX elaborates concluding remarks and discusses diverse elements that have not yet been (widely) explored by current CNN designs.

## II. NOTATION

Convolutional neural networks (CNNs) are composed by basic elements, such as convolution, activation and down/up-sampling layers. In order to achieve better clarity in the explanations given in this paper, we define the mathematical notation to represent the basic operations of CNNs. Part of the definitions presented here are based on the work of Zavala *et al.* [19].

In the following, a scalar is represented by a lower-case letter (e.g. $a$), while a vector is represented by an underlined lower-case letter (e.g. $\underline{a}$). Furthermore, a matrix, such as an image or convolution mask, is represented by a boldface lowercase letter (e.g. variables $\mathbf{x}$ and $\mathbf{y}$). Finally, a tensor is defined by a boldface uppercase letter. For example, the two arbitrary tensors $\mathbf{A}$ and $\mathbf{Q}$ are defined by

$$\mathbf{A} = \begin{pmatrix} \mathbf{a}_0^0 & \cdots & \mathbf{a}_{N_C-1}^0 \\ \vdots & \ddots & \vdots \\ \mathbf{a}_0^{N_R-1} & \cdots & \mathbf{a}_{N_C-1}^{N_R-1} \end{pmatrix}, \mathbf{Q} = \begin{pmatrix} \mathbf{q}^0 \\ \vdots \\ \mathbf{q}^{N_R-1} \end{pmatrix}. \quad (1)$$

Here, entries $\mathbf{a}_c^r$ and $\mathbf{q}^r$ represent two-dimensional arrays (matrices). Since the defined tensors are used in the context of CNNs, matrices $\mathbf{a}_c^r$ and $\mathbf{q}^r$ are learned filters, which have dimensions $(N_V \times N_H)$, where $N_V$ and $N_H$ denote the filter dimensions in the vertical and horizontal directions, respectively. Finally, we define the total tensor dimension of $\mathbf{A}$ and $\mathbf{Q}$ by $(N_C \times N_R \times N_V \times N_H)$ and $(N_R \times 1 \times N_V \times N_H)$, where $N_R$ and $N_C$ are the number of row and column entries, respectively. If the tensor $\mathbf{A}$ contains the convolution weights in a CNN, the row-entry dimensions represent the input number of channels to a layer, while the number of column elements denotes the number of output channels.

Having defined the notation for the variables, we focus on a few relevant operators. First, the transpose of a tensor $(\cdot)^\mathsf{T}$, expressed by

$$\mathbf{Q}^\mathsf{T} = \begin{pmatrix} \mathbf{q}^0 & \cdots & \mathbf{q}^{N_R-1} \end{pmatrix}. \quad (2)$$

Furthermore, the convolution of two tensors is written as $\mathbf{AQ}$ and specified by

$$\mathbf{AQ} = \begin{pmatrix} \sum_{r=0}^{N_R-1} \mathbf{a}_r^0 * \mathbf{q}^r \\ \vdots \\ \sum_{r=0}^{N_R-1} \mathbf{a}_r^{N_R-1} * \mathbf{q}^r \end{pmatrix}. \quad (3)$$

Here, the symbol $*$ defines the convolution between two matrices (images).

In this paper, images which are 2D arrays (matrices), are often convolved with 4D tensors. When this operation is performed, images are considered to have dimensions $(1 \times 1 \times N_V \times N_H)$. In addition, in this paper matrix $\mathbf{I}$ is the identity signal for the convolution operator, which for a 2D image is the Kronecker delta/discrete impulse (an image with a single non-zero pixel with unity amplitude at the center of the image). Furthermore, we indicate that variables in the *decoding path* of a CNN are distinguished with a tilde (e.g. $\tilde{\mathbf{K}}$, $\underline{\tilde{b}}$).

Additional symbols that will be used throughout the article are the down- and up-sampling operations by a factor $s$, which are denoted by $f_{(s\downarrow)}(\cdot)$ for down-sampling and for up-sampling $f_{(s\uparrow)}(\cdot)$. In this paper, both operations are defined in the same way as in multi-rate filter banks. For example, consider the signal

$$\underline{x} = \begin{pmatrix} 1, & 2, & 3, & 4, & 5, & 6, & 7, & 8, & 9, & 10 \end{pmatrix}. \quad (4)$$

If we apply the down-sampling operator to $\underline{x}$ by a factor 2, this results into

$$\underline{z} = f_{(2\downarrow)}(\underline{x}) = \begin{pmatrix} 1, & 3, & 5, & 7, & 9 \end{pmatrix}, \quad (5)$$

where $\underline{z}$ is the down-sampled version of $\underline{x}$. Conversely, the result of applying the up-sample operator $f_{(2\uparrow)}(\cdot)$ gives as result

$$f_{(2\uparrow)}(\underline{z}) = \begin{pmatrix} 1, & 0, & 3, & 0, & 5, & 0, & 7, & 0, & 9, & 0 \end{pmatrix}. \quad (6)$$

Additional operators used in the article are the *rectified linear unit* (ReLU), the *shrinkage/thresholding* and the *clipping*, which are represented by $(\cdot)_+$, $\tau_{(\cdot)}(\cdot)$ and $\mathcal{C}_{(\cdot)}(\cdot)$, respectively.

For better clarity, the most important symbols used in this article are summarized in Table I. In addition, the graphical
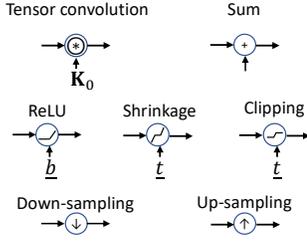
Fig. 1. Symbols used for the schematic representations of the CNNs addressed in this article.

representations of some of the symbols that will be used to graphically describe CNNs are shown in Fig. 1.

TABLE I
RELEVANT SYMBOLS USED IN THIS PAPER

| Symbol | Meaning |
|---|---|
| $f_{(2\downarrow)}(\cdot)$ | Down-sampling operation. |
| $f_{(2\uparrow)}(\cdot)$ | Up-sampling operation. |
| $\mathbf{I}$ | Convolution identity. |
| $\mathbf{K}$ | Encoding convolution kernel. |
| $\tilde{\mathbf{K}}$ | Decoding convolution kernel. |
| $\mathbf{W}$ | Filters for the forward discrete wavelet transform. |
| $\tilde{\mathbf{W}}$ | Filters for the inverse discrete wavelet transform. |
| $\mathbf{W}_\mathrm{H}$ | High-pass filters of the forward discrete wavelet transform. |
| $\mathbf{W}_\mathrm{L}$ | Low-pass filter of the forward discrete wavelet transform. |
| $\mathbf{x}$ | Noiseless image. |
| $\mathbf{y}$ | Noisy image. |
| $\boldsymbol{\eta}$ | Additive noise. |
| $\underline{b}$ | Bias vector. |
| $t$ | Threshold level. |
| $*$ | Image convolution. |
| $\mathbf{Kx}$ | Tensor convolution between tensor $\mathbf{K}$ and signal $\mathbf{x}$. |
| $(\cdot)^\mathsf{T}$ | Transpose of a tensor. |
| $(\cdot)_+$ | ReLU activation. |
| $\tau_{(\cdot)}(\cdot)$ | Generic thresholding/shrinkage operation. |
| $\mathcal{C}_{(\cdot)}(\cdot)$ | Generic clipping operation. |

## III. ENCODING-DECODING CNNS

### A. Signal model and noise reduction configurations

In noise-reduction applications, the common additive signal model is defined by

$$\mathbf{y} = \mathbf{x} + \boldsymbol{\eta}, \tag{7}$$

where the observed signal $\mathbf{y}$ is the result of contaminating a noiseless image $\mathbf{x}$ with additive noise $\boldsymbol{\eta}$. Assume that the noiseless signal $\mathbf{x}$ is to be estimated from the noisy observation $\mathbf{y}$. In deep learning applications, this is often achieved by models with the form

$$\hat{\mathbf{x}} = G(\mathbf{y}). \tag{8}$$

Here, $G(\cdot)$ is a generic encoding-decoding CNN. We refer to this form of noise reduction as *non-residual*. Alternatively, it is possible to find $\hat{\mathbf{x}}$ by training $G(\cdot)$ to estimate the noise component $\hat{\boldsymbol{\eta}}$, and subtract it from the noisy image $\mathbf{y}$ to estimate the noiseless image $\hat{\mathbf{x}}$, or equivalently

$$\hat{\mathbf{x}} = \mathbf{y} - G(\mathbf{y}). \tag{9}$$

This model is referred to as *residual* [5], [21], [7], because the output of the network is subtracted from its input. For reference, Fig. 2 portrays the difference of the placement of the encoding-decoding structure in residual and non-residual configurations.
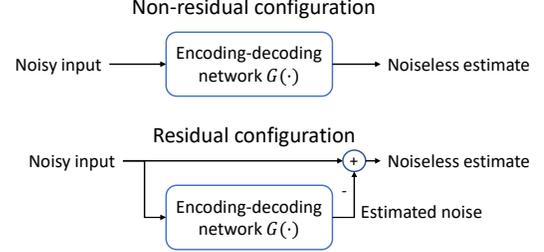


Fig. 2. Residual and non-residual network configurations. Note that the main difference between both designs is the global skip connection occurring in the residual structure. Still, it can be observed that the network $G(\cdot)$ may contain skip connections *internally*.

### B. Encoding-decoding CNNs

Encoding-decoding (convolutional) neural networks are rooted in the techniques for data-dimensionality reduction and unsupervised feature extraction, where a given signal is mapped to an alternative space via a non-linear transformation. This space should have properties which are somehow attractive for the considered task. For example, for dimensionality reduction, the alternative space should be lower-dimensional than the original input. In this article, we are interested in models that are useful for noise-reduction applications. Specifically, this manuscript addresses models that are referred to as encoding-decoding *convolutional* neural networks such as the model by Ranzato *et al.* [22], in which the encoder uses convolution filters to produce multi-channel/redundant representations, in which sparsifying non-linearities are applied. The sparsified signal is later mapped back to the original representation. It should be noted that despite the fact that the origins of the encoding-decoding CNNs are linked to feature extraction, this type of architecture quickly showed to be useful for other applications such as noise reduction, which is the topic of this article. For the rest of this manuscript whenever we mention an encoding-decoding CNN, we are referring to a design which follows the same basic principles of Ranzato's design.

It can be observed that encoding-decoding CNNs are constituted of three main parts. (1) The **encoder**, which maps the incoming image to a representation with more image channels with a convolution layer. Every channel of the resulting redundant representation contains a fraction of the content of the original signal. It should be noted that the encoder often (but not necessarily) decreases the resolution of the higher dimensional representation, to enable multi-resolution processing and to decrease the memory requirements of the design. (2) The **decoder**, which maps the multi-channel representation back to the original space. (3) The **non-linearities**, which suppress specific parts of the signal. In summary, the most basic encoding-decoding step in a CNN $G(\cdot)$ is expressed by

$$\boxed{G(\mathbf{y}) = G_\mathrm{dec}(G_\mathrm{enc}(\mathbf{y}))} , \tag{10}$$

where $G_{\text{enc}}(\cdot)$ is the encoder, which is generally defined by

$$\begin{aligned}
\mathbf{C}_0 &= E_0(\mathbf{y}), \\
\mathbf{C}_1 &= E_1(\mathbf{C}_0), \\
\mathbf{C}_2 &= E_2(\mathbf{C}_1), \\
&\vdots \\
\mathbf{C}_{N-1} &= E_{N-1}(\mathbf{C}_{N-2}), \\
G_{\text{enc}}(\mathbf{y}) &= \mathbf{C}_{N-1}.
\end{aligned} \tag{11}$$

Here, $\mathbf{C}_n$ represents the code generated by the $n$-th encoding $E_n(\cdot)$, which can be expressed by the equation

$$\boxed{\mathbf{C}_n = E_n(\mathbf{C}_{n-1}) = f_{(s\downarrow)}\big(A_{(\underline{b}_{n-1})}(\mathbf{K}_{n-1}\mathbf{C}_{n-1})\big)} \quad . \tag{12}$$

Here, the function $A(\cdot)_{(\cdot)}$ is a generic activation used in the encoder and $f_{(s\downarrow)}(\cdot)$ is a down-sampling function by factor $s$. Complementary to the encoder, the decoder network maps the multi-channel sparse signal back to the original domain. Here, we define the decoder by

$$\begin{aligned}
\tilde{\mathbf{C}}_{N-2} &= D_{N-1}(\mathbf{C}_{N-1}), \\
&\vdots \\
\tilde{\mathbf{C}}_1 &= D_2(\tilde{\mathbf{C}}_2), \\
\tilde{\mathbf{C}}_0 &= D_1(\tilde{\mathbf{C}}_1), \\
G(\mathbf{y}) &= D_0(\tilde{\mathbf{C}}_0),
\end{aligned} \tag{13}$$

where $\hat{\mathbf{C}}_n$ is the $n-$th decoded signal, which is produced by the $n$-th decoder layer, yielding the general expression:

$$\boxed{\tilde{\mathbf{C}}_{n-1} = D_n(\tilde{\mathbf{C}}_n) = \tilde{A}_{(\underline{\tilde{b}})}\big(\tilde{\mathbf{K}}_n^{\mathsf{T}} f_{(s\uparrow)}(\tilde{\mathbf{C}}_n)\big)} \quad . \tag{14}$$

In the above, $\tilde{A}(\cdot)_{(\cdot)}$ is the activation function used in the decoder and $f_{(s\uparrow)}(\cdot)$ is an up-sampling function of factor $s$.

An important remark is that the encoder-decoder CNN does not always contain down/up-sampling layers in which case, the decimation factor $s$ is unity, which causes $f_{(1\uparrow)}(\mathbf{x}) = f_{(1\downarrow)}(\mathbf{x}) = \mathbf{x}$ for any matrix $\mathbf{x}$. Furthermore, it should be noted also that we assume that the number of channels of the code $\mathbf{C}_N$ is always larger than the previous one $\mathbf{C}_{N-1}$. Furthermore, it should be noted that a single encoder layer $E_n(\cdot)$ and its corresponding decoder layer $D_n(\cdot)$ can be considered a single-layer encoder-decoder network/pair.

For this article, the encoding convolution filter for a given layer $\mathbf{K}$ has dimensions $(N_o \times N_i \times N_h \times N_v)$, where $N_i$ and $N_o$ are the number of input and output channels for a convolution layer, respectively. Similarly, $N_h$ and $N_v$ are the number of elements in the horizontal and vertical directions, respectively. Note that the encoder increases the number of channels of the signal (e.g. $N_o > N_i$), akin to Ranzatto's design [22]. Furthermore, it is assumed that the *decoder* is symmetric in the number of channels to the *encoder*, therefore, the dimensions of the decoding convolution kernel $\tilde{\mathbf{K}}^{\mathsf{T}}$ are $(N_i \times N_o \times N_h \times N_v)$. The motivation of this symmetry is to emphasize the similarity between the signal processing and the CNN elements.

## IV. SIGNAL PROCESSING FUNDAMENTALS

As shown by Ye *et al.* [12], within encoding-decoding CNNs, the signal is treated akin to well-known sparse representations, where the coefficients used for the transformation are directly learned from the training data. Prior to addressing this important concept in more detail, relevant supporting concepts such as sparsity, sparse transformations and non-linear signal estimation in the wavelet domain are explained.

### A. Sparsity

A sparse image is a signal where most coefficients are small and the relatively few large coefficients capture most of the information [23]. This characteristic allows to discard low-amplitude components with relatively small perceptual changes. Hereby, the use of sparse signals is attractive for applications such as image compression, denoising and suppression of artifacts.

Despite the convenient characteristics of sparse signals, natural images are often non-sparse. Still, there are numerous transformations, which allow to map the signal to a sparse domain and which are analogous to the internal operations of CNNs. For example, *singular value decomposition* factorizes the image in terms of two sets of orthogonal bases of which few basis pairs contain most of the energy of the image. An alternative transformation is based on *framelets*, where an image is decomposed in a multi-channel representation, whereby each resulting channel contains a fragment of the Fourier spectrum. In the remainder of this section we will address all of these representations in more detail.

### B. Sparse signal representations

*1) Singular value decomposition (SVD) and low-rank approximation:* Assume that an image (patch) is represented by a matrix $\mathbf{y}$ with dimensions $(N_r \times N_c)$, where $N_r$ and $N_c$ are the number of rows and columns, respectively. Then, the singular value decomposition factorizes $\mathbf{y}$ as

$$\mathbf{y} = \sum_{n=0}^{N_{\text{SV}}-1} (\underline{u}_n \underline{v}_n^{\mathsf{T}}) \cdot \underline{\sigma}[n], \tag{15}$$

in which $N_{\text{SV}}$ is the number of singular values, $n$ is a scalar index, while $\underline{u}_n$ and $\underline{v}_n$ are the $n^{\text{th}}$ left- and right-singular vectors, respectively. Furthermore, vector $\underline{\sigma}$ contains the singular values and each of its entries $\underline{\sigma}[n]$ is the weight assigned to every basis pair $\underline{u}_n$, $\underline{v}_n$. This means that the product $(\underline{u}_n \underline{v}_n^{\mathsf{T}})$ contributes more to the image content for higher values of $\underline{\sigma}[n]$. It is customary that the singular values are ranked in descending order and the amplitudes of the singular values $\underline{\sigma}$ are sparse, therefore $\underline{\sigma}[0] \gg \underline{\sigma}[N_{\text{SV}}-1]$. The reason for this sparsity is because image (patches) intrinsically have high correlation. For example, many images contain repetitive patterns (e.g. a wall with bricks, a fence, the tiles of a rooftop or the stripes of a zebra) or uniform regions (for example, the sky, the skin of a person). This means an image patch may contain only a few linearly independent vectors that describe most of the image content. Consequently, a higher weight is assigned to such image bases.

Given that the amplitudes of the singular values of $\mathbf{y}$ in SVD are sparse, it is possible approximate $\hat{\mathbf{y}}$ with only a few bases $(\underline{u}_n \underline{v}_n^{\mathsf{T}})$. Note that this procedure reduces the rank of signal $\mathbf{y}$ and hence it is known as *low-rank approximation*. This process is equivalent to

$$\hat{\mathbf{y}} = \sum_{n=0}^{N_{\mathrm{LR}}-1} (\underline{u}_n \underline{v}_n^{\mathsf{T}}) \cdot \sigma[n], \tag{16}$$

where $N_{\mathrm{SV}} > N_{\mathrm{LR}}$. Note that this effectively cancels the product $(\underline{u}_n \underline{v}_n^{\mathsf{T}})$ where the weight given by $\sigma[n]$ is low. Alternatively, it is possible to assign a weight of zero to the product $(\underline{u}_n \underline{v}_n^{\mathsf{T}})$ for $n \geq N_{\mathrm{LR}}$.

The low-rank representation of a matrix is desirable for diverse applications among which we can find image denoising. The motivation for using low-rank approximation for this application results from the fact that –as mentioned earlier– natural images are considered low-rank due to the strong spatial correlation between pixels, whereas noise is high-rank (it is spatially uncorrelated). In consequence, reducing the rank/number of singular values decreases the presence of noise, while still providing a good approximation of the noise-free signal, as exemplified in Fig. 3.
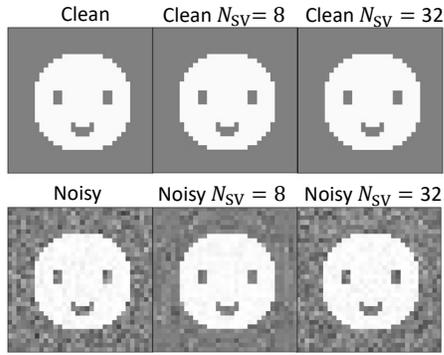


Fig. 3. SVD reconstruction of clean and corrupted images with a different number of singular values. Note that the reconstruction of the clean image with 8 or 32 singular values ($N_{\mathrm{SV}} = 8$ or $N_{\mathrm{SV}} = 32$, respectively) yields reconstructions indistinguishable from the original image. This contrasts with their noisy counterparts, where $N_{\mathrm{SV}} = 8$ reconstructs a smoother image in which the noise is attenuated, while $N_{\mathrm{SV}} = 32$ reconstructs the noise texture perfectly.

*2) Framelets:* Just as SVD, framelets are also commonly used for image processing. In a nutshell, a framelet transform is a signal representation that factorizes/decomposes an arbitrary signal into multiple bands/channels. Each of these channels contain a segment of the energy of the original signal. In image and signal processing, the framelet bands are the result of convolving the analyzed signal with a group of discrete filters that have finite length/support. In this article, the most important characteristic that the filters of the framelet transform should comply with, is that the bands they generate capture *all* the energy contained on the input to the decomposition. This is important to avoid the loss of information of the decomposed signal. In this text, we refer to framelets that comply with the previous characteristics as *tight* framelets and the following paragraphs will describe this property in more detail.

In its decimated version, the framelet decomposition for *tight* frames is represented by

$$\mathbf{Y}_{\mathrm{fram}} = f_{(2\downarrow)}(\mathbf{Fy}), \tag{17}$$

in which $\mathbf{Y}_{\mathrm{fram}}$ is the decomposed signal and $\mathbf{F}$ is the framelet basis (tensor). Note that the signal $\mathbf{Y}_{\mathrm{fram}}$ has more channels than $\mathbf{y}$. Furthermore, the original signal $\mathbf{y}$ is recovered from $\mathbf{Y}_{\mathrm{fram}}$ by

$$\mathbf{y} = \tilde{\mathbf{F}}^{\mathsf{T}} f_{(2\uparrow)}(\mathbf{Y}_{\mathrm{fram}}) \cdot c. \tag{18}$$

Here, $\tilde{\mathbf{F}}$ is the filter of the inverse framelet transform and $c$ denotes an arbitrary constant. If $c = 1$ the framelet is *normalized*. Finally, note that the framelet transform can also be *undecimated*. This means that in undecimated representations, the down-sampling and up-sampling layers $f_{(2\downarrow)}(\cdot)$ and $f_{(2\uparrow)}(\cdot)$ are not used. An important property of the undecimated representation is that it is less prone to aliasing than its decimated counterpart, but more computationally expensive. Therefore, for efficiency reasons, the decimated framelet decomposition is often preferred over the undecimated representation. In summary, the decomposition and synthesis of the decimated framelet decomposition is represented by

$$\boxed{\mathbf{y} = \tilde{\mathbf{F}}^{\mathsf{T}} f_{(2\uparrow)}\big(f_{(2\downarrow)}(\mathbf{Fy})\big) \cdot c} \;, \tag{19}$$

while for the undecimated framelet it holds that

$$\boxed{\mathbf{y} = \tilde{\mathbf{F}}^{\mathsf{T}}(\mathbf{Fy}) \cdot c} \;. \tag{20}$$

A notable normalized framelet is the *discrete wavelet transform* (DWT), where variables $\mathbf{F}$ and $\tilde{\mathbf{F}}$ are replaced by tensors $\mathbf{W} = (\mathbf{w}_{\mathrm{LL}}, \mathbf{w}_{\mathrm{LH}}, \mathbf{w}_{\mathrm{HL}}, \mathbf{w}_{\mathrm{HH}})$ and $\tilde{\mathbf{W}} = (\tilde{\mathbf{w}}_{\mathrm{LL}}, \tilde{\mathbf{w}}_{\mathrm{LH}}, \tilde{\mathbf{w}}_{\mathrm{HL}}, \tilde{\mathbf{w}}_{\mathrm{HH}})$, respectively. Here, $\mathbf{w}_{\mathrm{LL}}$ is the filter for the low-frequency band, while $\mathbf{w}_{\mathrm{LH}}, \mathbf{w}_{\mathrm{HL}}, \mathbf{w}_{\mathrm{HH}}$ are the filters used to extract the detail in the horizontal, vertical and diagonal directions, respectively. Finally, $\tilde{\mathbf{w}}_{\mathrm{LH}} \, \tilde{\mathbf{w}}_{\mathrm{LH}}, \tilde{\mathbf{w}}_{\mathrm{HL}}, \tilde{\mathbf{w}}_{\mathrm{HH}}$ are the filters of the inverse decimated DWT.

In order to understand the DWT more intuitively, Fig. 4 shows the decimated framelet decomposition using the filters of the discrete wavelet transform. Note that the convolution $\mathbf{Wy}$ results in a four-channel signal, where each channel contains only a fraction of the spectrum of image $\mathbf{y}$. This allows to down-sample each channel with minimal aliasing. Furthermore, to recover the original signal, each individual channel is up-sampled, thereby introducing aliasing, which is then removed by the filters of the inverse transform. Finally, all the channels are added and the original signal is recovered.

Analogous to the low-rank approximation, in framelets, the reduction of noise is achieved by setting noisy components to zero. These components are typically assumed to have low-amplitude compared with the amplitude of the sparse signal, as expressed by

$$\hat{\mathbf{y}} = \tilde{\mathbf{F}}^{\mathsf{T}} f_{(2\uparrow)}\big(\tau_{(\underline{t})}\big(f_{(2\downarrow)}(\mathbf{Fy})\big)\big) \cdot c, \tag{21}$$

where $\tau_{\underline{t}}(\cdot)$ is a generic thresholding/shrinkage function, which sets each of the pixels in $f_{(2\downarrow)}(\mathbf{Fy})$ to zero when values are lower than the threshold level $\underline{t}$.
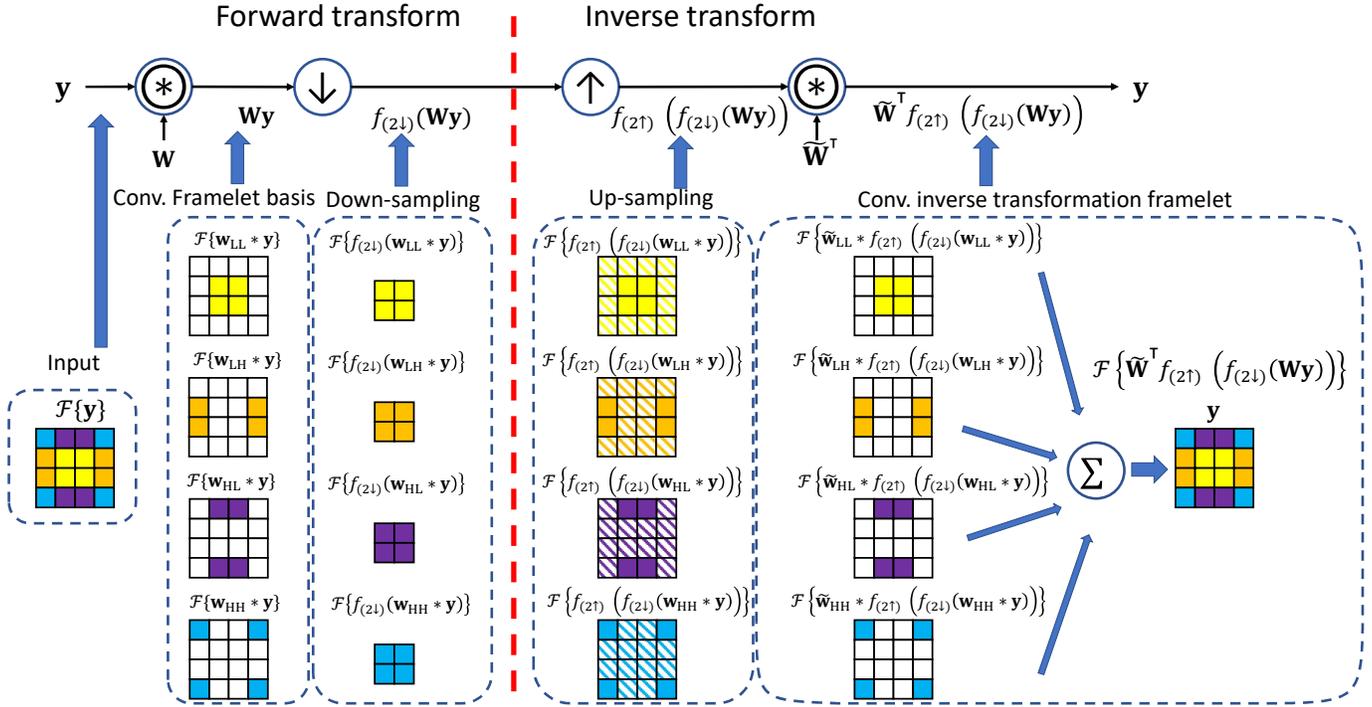
Fig. 4. 2D spectrum analysis of the decimated discrete framelet decomposition and reconstruction. In the figure, function $\mathcal{F}\{\cdot\}$ stands for the amplitude Fourier spectrum of the input argument. The yellow squares indicate a region in the low-frequency area of the Fourier spectrum, while the orange, purple and blue squares indicate the high-pass/detail bands. For these images, ideal orthogonal bases are assumed. Note that the forward transform is composed by two steps. First, the signal is convolved with the wavelet basis ($\mathbf{W}\mathbf{y}$). Afterwards, down-sampling is applied to the signal ($f_{(2\downarrow)}(\mathbf{W}\mathbf{y})$). During the inverse transformation, the signal is up-sampled by inserting zeros between each sample ($f_{(2\uparrow)}(f_{(2\downarrow)}(\mathbf{W}\mathbf{y}))$), which causes spatial aliasing (dashed blocks). Finally, the spatial aliasing is removed by the inverse transform filter $\tilde{\mathbf{W}}$ and all the channels are added ($\tilde{\mathbf{W}}^{\mathsf{T}} f_{(2\uparrow)}(f_{(2\downarrow)}(\mathbf{W}\mathbf{y}))$).

### C. Nonlinear signal estimation in the framelet domain

As mentioned in Section IV-B2, framelets decompose a given image $\mathbf{y}$ by convolving it with a tensor $\mathbf{F}$. Note that many of the filters that compose $\mathbf{F}$ have a high-pass nature. Images often contain approximately uniform regions in which the variation is low, therefore, convolving a signal $\mathbf{y}$ with a high-pass filter $\mathbf{f}_h$ –where $\mathbf{f}_h \in \mathbf{F}$– produces the sparse detail band $\mathbf{d} = \mathbf{f}_h * \mathbf{y}$ in which uniform regions have low amplitudes, while transitions i.e. edges contain most of the energy of the bands.

Assuming a model in which a single pixel $d \in \mathbf{d}$ is observed, which is contaminated with additive noise $\eta$. Then, the resulting observed pixel $z$ is defined by

$$z = d + \eta. \tag{22}$$

In order to recover the noiseless pixel $d$ from observation $z$, it is possible to use the point-*maximum a posteriori* (MAP) estimate [1], [24], defined by the maximization problem

$$\hat{d} = \underset{d}{\arg \max}\left[ \ln\big(P(d|z)\big)\right]. \tag{23}$$

Here, the log-posterior $\ln\big(P(d|z)\big)$ is defined by

$$\ln\big(P(d|z)\big) = \ln\big(P(z|d)\big) + \ln\big(P(d)\big), \tag{24}$$

where the conditional probability density function (PDF) $P(z|d)$ expresses the noise distribution, which is often as-

sumed Gaussian and is defined by

$$P(z|d) \propto \exp\left(-\frac{(z-d)^2}{2\sigma_\eta^2}\right). \tag{25}$$

Here, $\sigma_\eta^2$ is the noise variance. Furthermore, as prior probability, it is assumed that the distribution of $P(d)$ corresponds to a Laplacian distribution, which has been used in wavelet-based denoising [1]. Therefore, $P(d)$ is mathematically described by

$$P(d) \propto \exp\left(-\frac{|d|}{\sigma_d}\right), \tag{26}$$

where $\sigma_d$ is the dispersion measure of the Laplace distribution. For reference, Fig. 5 portrays an example of a Gaussian and a Laplacian PDF. Note that the Laplacian distribution has a higher probability of zero elements to occur than the Gaussian distribution for the same standard deviation. Finally,
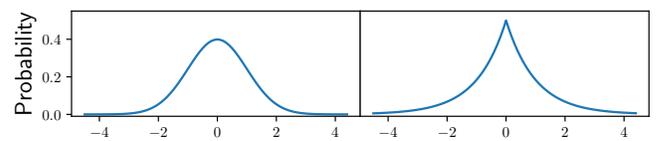


Fig. 5. Probability density function for Gaussian (left) and Laplacian (right) distributions.

substituting Eq. (25) and Eq. (26) in Eq (24) results in

$$\ln\big(P(d|z)\big) \propto -\frac{(z-d)^2}{2\sigma_\eta^2} - \frac{|d|}{\sigma_d}. \tag{27}$$

In the above, maximizing $d$ in $\ln\left(P(d|z)\right)$ with the first derivative criterion –in an (un) constrained way– leads to two common activations in noise-reduction CNNs: the ReLU and the soft-shrinkage function. Furthermore, the solution also can be used to derive the so-called *clipping* function, which is useful in residual networks.

For reference and further understanding, Fig. (6) portrays the elements composing the noise model of Eq. (22), the signal transfer characteristics of the ReLU, soft-shrinkage and clipping functions, and the effect that these functions have on the signal of the observed noisy detail band $z$.
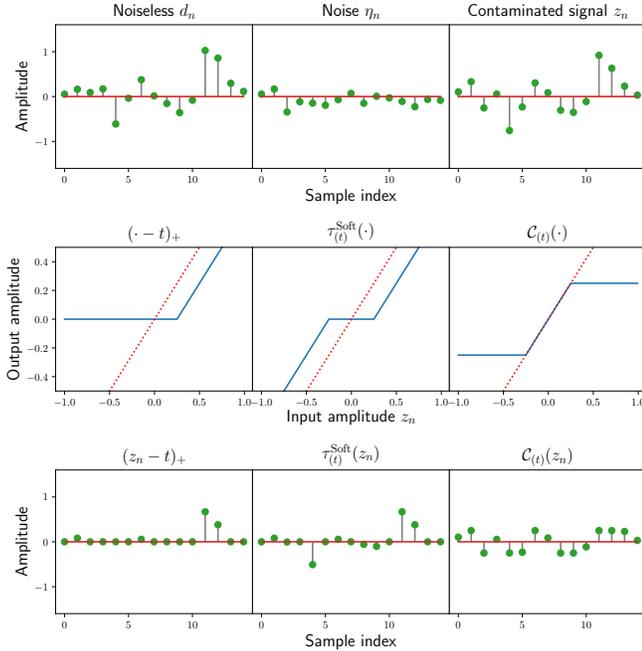


Fig. 6. Signals involved in the additive noise model, input/output transfer characteristics of activation layers and estimates produced by the activation layers when applied to the noise-contaminated signal. The first row shows the signals involved in the additive noise model. The second row depicts the output amplitude of activation functions with respect to the input amplitude. Finally, the last row depicts the application of the activation functions to the noisy observation $z$.

*1) Rectified linear unit (ReLU):* If Eq. (27) is solved for $d$ while constraining the estimator to be positive, the noiseless estimate $\hat{d}$ becomes

$$\boxed{\hat{d} = (z - t)_+} \quad , \tag{28}$$

which is also expressed by

$$(z - t)_+ = \begin{cases} z - t, & \text{if } z \geq t, \\ 0, & \text{if } t > z. \end{cases} \tag{29}$$

Here, the threshold level is defined by

$$t = \sigma_\eta^2 / \sigma_d. \tag{30}$$

Note that this estimator cancels the negative and low-amplitude elements of $d$ lower than the magnitude of the threshold level $t$. For example, if the signal content on the feature map is low, then $\sigma_d \to 0$. In such case, $t \to +\infty$ and consequently $\hat{d} \to 0$. This means that the channel is suppressed. Alternatively,

if the feature map has strong signal presence i.e. $\sigma_d \to \infty$, consequently $t \to 0$ and then $\hat{d} \to (z)_+$.

A final remark is made on the modeling of functions of a CNN. It should be noted that the estimator of Eq. (28) is analogous to the activation function of a CNN, known as *rectified linear unit* (ReLU). However, in a CNN the value of $t$ would be the bias $b$ learned from the training data.

*2) Soft-shrinkage/thresholding:* If Eq. (27) is maximized in an unconstrained way, the estimate $\hat{d}$ is

$$\boxed{\hat{d} = \tau_{(t)}^{\text{Soft}}(z) = (z - t)_+ - (-z - t)_+} \quad . \tag{31}$$

Here, $\tau_{(t)}^{\text{Soft}}(\cdot)$ denotes the soft-shrinkage/-thresholding function, which is often also written in the form

$$\tau_{(t)}^{\text{Soft}}(z) = \begin{cases} z + t, & \text{if } z \geq t, \\ 0, & \text{if } t > z \geq -t, \\ z - t, & \text{if } -t > z. \end{cases} \tag{32}$$

It can be observed that the soft-threshold enforces the low-amplitude components whose magnitude is lower than the magnitude threshold level $t$ to zero. In this case, $t$ is also defined by Eq. (30). It should be noted that the soft-shrinkage estimator can also be obtained from a variational perspective [25]. Finally, it can be observed that the soft-shrinkage is the superposition of two ReLU functions, which has been pointed out by Fan *et al.* [18].

*3) Soft clipping:* In Section IV-C1 and Section IV-C2, the estimate $\hat{d}$ is obtained directly from the noisy observation $z$. Alternatively, it is possible to estimate the noise $\eta$ and subtract it from $z$ akin to the residual CNNs represented by Eq. (9). This can be achieved by solving the model

$$\boxed{\hat{\eta} = z - \hat{d} = z - \tau_{(t)}^{\text{Soft}}(z)} \quad , \tag{33}$$

which is equivalent to

$$\hat{\eta} = \mathcal{C}_{(t)}^{\text{Soft}}(z) = z - ((z - t)_+ - (-z - t)_+), \tag{34}$$

where $\mathcal{C}_{(t)}^{\text{Soft}}(\cdot)$ is the *soft clipping* function. Note that this function also can be expressed by

$$\mathcal{C}_{(t)}^{\text{Soft}}(z) = \begin{cases} t, & \text{if } z \geq t, \\ z, & \text{if } t \geq z > -t, \\ -t, & \text{if } -t \geq z. \end{cases} \tag{35}$$

*4) Other thresholding layers:* One of the main drawbacks of the soft-threshold activation is that it is a biased estimator. This limitation has been addressed by the hard and semi-hard thresholds, which are (asymptotically) unbiased estimators for large input values. In this section, we focus solely on the semi-hard threshold and avoid the hard variant, because is discontinuous and, therefore, not suited for models that rely on gradient-based optimization, such as CNNs.

Among the semi-hard thresholds, two notable examples are the *garrote shrink* and the shrinkage functions generated by derivatives of Gaussians (DoG) [26], [19]. The garrote shrink function $\tau_{(\cdot)}^{\text{Gar}}(\cdot)$ is defined by

$$\tau_{(t)}^{\text{Gar}}(z) = \frac{(z^2 - t^2)_+}{z}. \tag{36}$$

Furthermore, an example of a shrinkage function based on the derivative of Gaussians is given by

$$\tau_{(t)}^{\text{DoG}}(z) = z - \mathcal{C}_{(t)}^{\text{DoG}}(z), \tag{37}$$

where the semi-hard clipping function with the derivative of Gaussians $\mathcal{C}_{(\cdot)}^{\text{DoG}}(\cdot)$ is given by

$$\mathcal{C}_{(t)}^{\text{DoG}}(z) = z \cdot \exp\left(-\frac{z^p}{t^p}\right), \tag{38}$$
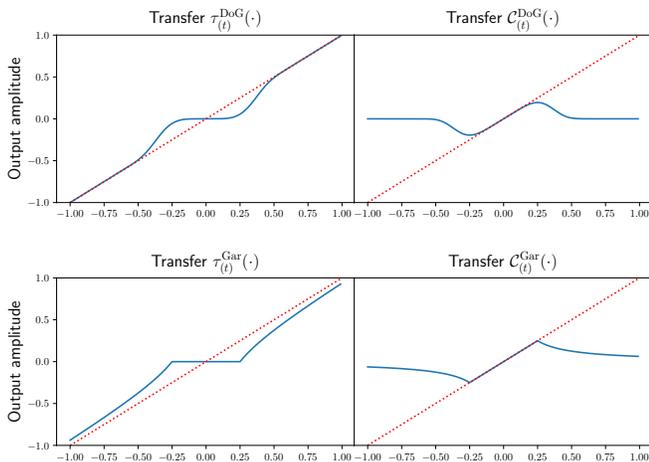
in which $p$ is an even number.



Fig. 7. Transfer characteristics of the semi-hard thresholds based on the difference of Gaussians and of the garrote shrink, as well as their clipping counterparts. Note that in contrast with the soft-shrinkage and clipping functions shown in Fig. 6, in the semi-hard thresholds tend to unity for large values, while the semi-hard clipping functions tend to zero for large signal intensities.

The garrote and semi-hard DoG shrinkage function are shown in Fig. 7, as well as their clipping counterparts. Note that the shrinkage functions approximate unity for $|z| \gg t$. Therefore, they are asymptotically unbiased for large signal values.

The final thresholding function addressed in this section is the linear expansion of thresholds proposed by Blu and Luisier [26]. This technique combines multiple thresholding functions to improve the performance. This approach is known as *linear expansion of thresholds* (LET) and it is defined by

$$\tau_{(\underline{t})}^{\text{LET}}(z) = \sum_{n=0}^{N_T-1} a_n \cdot \tau_{(t_n)}(z), \tag{39}$$

where $a_n$ is the weighting factor assigned to each threshold, where all weighting factor should add up to unity.

## V. BRIDGING THE GAP BETWEEN SIGNAL PROCESSING AND CNNS: DEEP CONVOLUTIONAL FRAMELETS AND SHRINKAGE-BASED CNNS

This section addresses the theoretical operation of noise-reduction convolutional neural networks based on ReLUs and shrinkage/thresholding functions. The first part of this section describes the theory of deep convolutional framelets [12], which is the most extensive study on the operation of

encoding-decoding ReLU-based CNNs up to this moment. Afterwards, the section concentrates on the operation of networks which use shrinkage functions instead of ReLUs [17], [18], [19], with the aim of mimicking well-established denoising algorithms [1]. Finally, the last part of this section addresses the connections between both methods and additional links between convolutional neural networks and signal processing.

### A. Theory of deep convolutional framelets

The theory of deep convolutional framelets [12] describes the operation of encoding-decoding ReLU-based CNNs. Its most relevant contributions are as follows. (1) To establish the equivalence of framelets and the convolutional layers of CNNs. (2) The theory of deep convolutional framelets provides the conditions to preserve the signal integrity within a ReLU CNN. (3) Explain how ReLU and convolution layers reduce noise within an encoding-decoding CNN.

The similarity between framelets and the encoding and decoding convolutional filters can be observed when comparing Eqs. (12), (14) with Eqs. (17), (18), where it becomes visible that the convolution structure of encoding-decoding CNNs is analogous to the forward and inverse framelet decomposition.

Regarding the signal reconstruction characteristics, the theory of deep convolutional framelets [12] states the following. First, in order to be able to recover an arbitrary signal $\mathbf{y} \in \mathbb{R}^N$, the number of output channels of a convolution layer with ReLU activation should *at least* duplicate the number of input channels. Second, the encoding convolution kernel $\mathbf{K}$ should be composed of pairs of filters with opposite phase. These two requirements ensure that any negative and positive values propagate through the network. Under these conditions, the encoding and decoding convolution filters $\mathbf{K}$ and $\tilde{\mathbf{K}}$ should comply with

$$\boxed{\mathbf{y} = \tilde{\mathbf{K}}^\mathsf{T}(\mathbf{K}\mathbf{y})_+ \cdot c} \ . \tag{40}$$

It can be noticed that Eq. (40) is an extension of Eq. (20), which describes the reconstruction characteristics of tight framelets. From this point, we refer to convolutional kernels compliant with Eq (40) as *phase-complementary tight framelets*. As a final remark, it should be noted that a common practice in CNN designs is also to use ReLU non-linearities in the decoder, in such case the phase-complementary tight-framelet condition can still be met as long as the pixels $y \in \mathbf{y}$ comply with $y \geq 0$, which is equivalent to

$$\boxed{\mathbf{y} = (\mathbf{y})_+ = \left(\tilde{\mathbf{K}}^\mathsf{T}(\mathbf{K}\mathbf{y})_+ \cdot c\right)_+} \ . \tag{41}$$

It can be observed that the relevance of the properties defined in Eqs. (40) and (41) is that they ensure that a CNN can propagate any arbitrary signal, which is important to avoid any distortions (such as image blur) in the processed images.

An additional element of the theory of deep convolutional framelets regarding the reconstruction of the signal, is to show that conventional pooling layers (e.g. average pooling) discard high-frequency information of the signal, which effectively blurs the processed signals. Furthermore, Ye *et al.* [12] have demonstrated that this can be fixed by replacing the conventional up/down-sampling layers by reversible operations, such

## Fitting low-rank approximation in ReLU CNNs

*From low-rank approximation to an encoding-decoding CNN.* In order to further understand the analogy between CNNs and low-rank approximation established by the theory of deep convolutional framelets, we can use as starting point the definition of singular value decomposition, which is expressed in Eq. (15), by

$$\mathbf{y} = \sum_{n=0}^{N_{\text{SV}}-1} (\underline{u}_n \underline{v}_n^{\mathsf{T}}) \cdot \underline{\sigma}[n].$$

Given that left and right singular vector pairs $\underline{u}_n \underline{v}_n^{\mathsf{T}}$ generate an image $\mathbf{D}[n]$, then Eq. (15) can be rewritten to

$$\mathbf{y} = \sum_{n=0}^{N_{\text{SV}}-1} \mathbf{D}[n] \cdot \underline{\sigma}[n], \qquad (42)$$

where tensor $\mathbf{D} = \left((\underline{u}_0 \underline{v}_0^{\mathsf{T}}) \dots (\underline{u}_{N_{\text{SV}}-1} \underline{v}_{N_{\text{SV}}-1}^{\mathsf{T}})\right)^{\mathsf{T}}$ contains the products of the left and right singular vectors and has dimensions $(N_{\text{SV}} \times 1 \times M \times N)$. Furthermore, the equation can be further reformulated to

$$\mathbf{y} = \tilde{\mathbf{K}}^{\mathsf{T}} \mathbf{D}, \qquad (43)$$

in which $\tilde{\mathbf{K}}^{\mathsf{T}} = \left((\underline{\sigma}[0]) \dots (\underline{\sigma}[N_{\text{SV}}-1])\right)$, where the brackets of the $(1 \times 1)$ filters have been excluded for simplicity. In addition, it is now assumed that it is desirable to perform low-rank approximation of signal $\mathbf{y}$ based on the reformulation of Eq. (43). If we assume that $\mathbf{D} \in \mathbb{R}_{\geq 0}^{N}$, then the low-rank approximation can be expressed by

$$\hat{\mathbf{y}} = \tilde{\mathbf{K}}^{\mathsf{T}} (\mathbf{D} + \underline{b})_+, \qquad (44)$$

in which the values $\underline{b}$ are set to zero for the channels of $\mathbf{D}$ that have high contributions to the image content. Conversely, the channels of $\mathbf{D}[n]$ with less perceptual relevance are then cancelled by assigning large negative values to the corresponding entries of $\underline{b}$. As final reformulation, we can assume that the basis images $\mathbf{D}$ are the result of decomposing the input image $\mathbf{y}$ with a set of convolution filters i.e. $\mathbf{D} = \mathbf{K}\mathbf{y}$, this transforms Eq. (44) into

$$\boxed{\hat{\mathbf{x}} = \tilde{\mathbf{K}}^{\mathsf{T}} (\mathbf{K}\mathbf{y} + \underline{b})_+} \quad . \qquad (45)$$

Here, it is visible that Eq. (45) is analogous to the encoding-decoding architecture defined in Eqs. (10) to (14) and the encoder and decoder filters are akin to the framelet

formulation of Section IV-B2. Note that Eq. (45) assumes that the entries $\mathbf{D} = \mathbf{K}\mathbf{y}$ are positive, which may be not always true. In this situation, tensor $\mathbf{D}$ requires redundant channels in which their respective phases are inverted to avoid the signal loss. Furthermore, it should also be noticed that in a CNN, the bias/threshold level is not inferred from the statistics of the feature maps, but learned from the data presented to the network during training.

*Multi-layer designs.* It should be noted that CNNs contain multiple layers, which recursively decompose/reconstruct the signal. This may pose an advantage with respect to conventional low-rank approximation algorithms for a few reasons. First, the data-driven nature of CNNs allows to learn the basis functions which optimally decompose and suppress noise in the signal. Second, since networks are deep, the incoming signal is recursively decomposed and sparsified. This multi-decomposition scheme is very similar to the designs used in noise-reduction algorithms based on framelets. It can be noted that the recursive sparsifying principles have been observed in the past in methods such as the (learned) iterative soft-thresholding algorithm [27], [28] as well as convolutional sparse coding. In fact, convolutional sparse-coding approach, which also has been used for interpreting the operation of CNNs [16]

*What about practical implementations?* When training a CNN, the parameters of the model (i.e. $\mathbf{K}$, $\tilde{\mathbf{K}}^{\mathsf{T}}$ and $\underline{b}$) are updated to reduce the loss between the processed noisy signal and the ground truth, which does not warranty that the numerical values of the convolution filters and biases of the trained model comply with the assumptions performed here. This is because CNNs do not have mechanisms to enforce that filters have properties such as sparsity or perfect reconstruction and negative values for the biases. Consequently, CNNs may not necessarily perform a low-rank approximation of the signal, although the mathematical formulation of the low-rank approximation and the single-layer encoding-decoding are similar. Hence, the analysis presented here should be treated as insight on the mathematical formulation and/or potential properties that can be enforced for specific applications and not as a literal description of what trained models do.

as the discrete wavelet transform. To exemplify this property, we refer to Fig. 4. If only an average pooling layer followed by an up-sampling stage would be applied, the treatment of the signal would be equivalent to the low-frequency branch of the DWT. Consequently, only the low-frequency spectrum of the signal would be recovered and images processed with that structure would become blurred. In contrast, if the full forward and inverse wavelet transform of Fig. 4 is used for up- and down-sampling, it is possible to reconstruct any signal, irrespective of its frequency content.

The ultimate key contribution of the theory of deep convolutional framelets is the explanation of the operation of ReLU-based noise-reduction CNNs. For the non-residual configuration, ReLU CNNs perform the following operations. (1) The convolution filters decompose the incoming signal into a sparse multi-channel representation. (2) The feature maps which are uncorrelated to the signal, contain mainly noise. In this case, the bias and the ReLU activation cancel the noisy feature maps in a process analogous to the MAP estimate shown in Section IV-C1. (3) The decoder reconstructs the
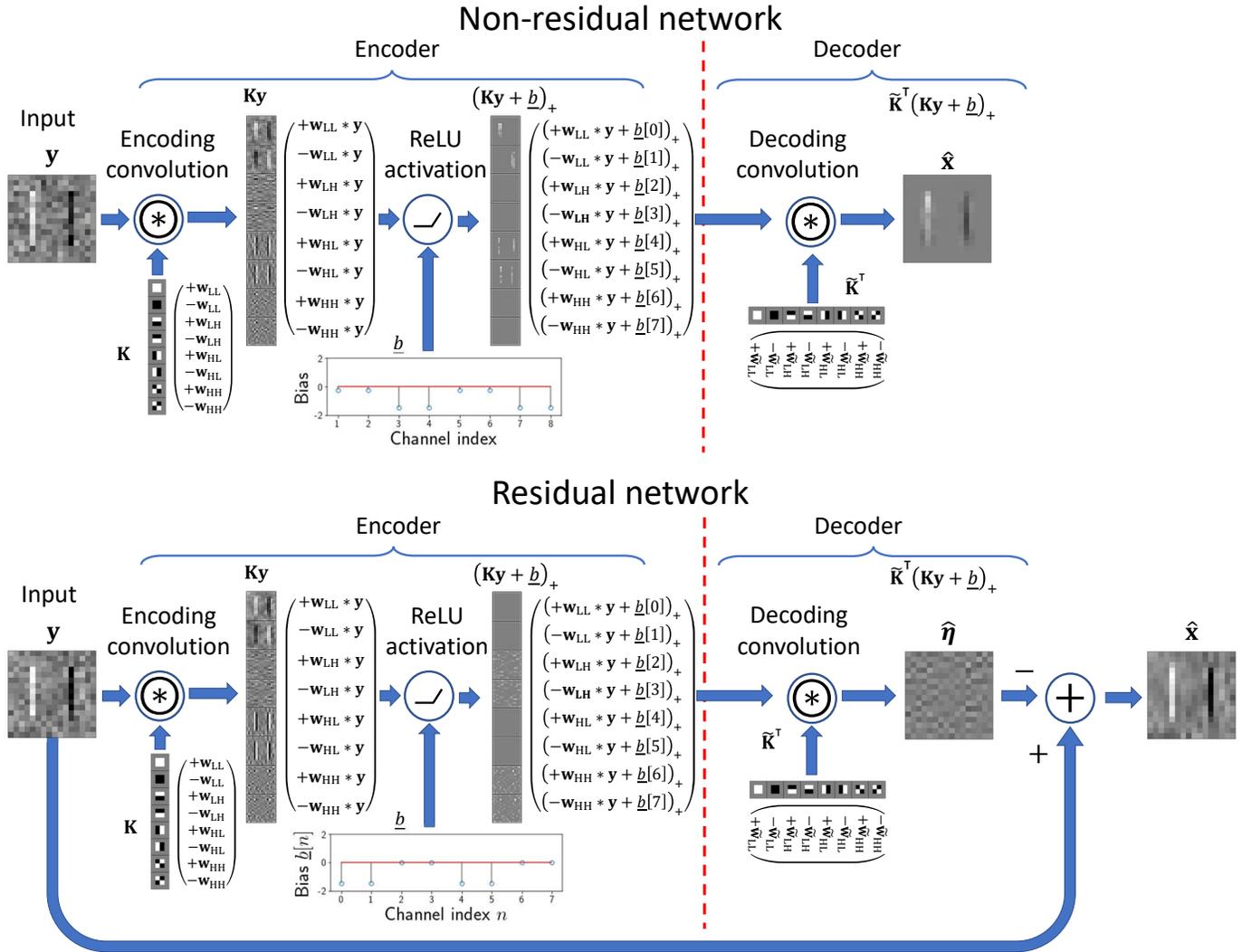
Fig. 8. Operation of a simplified denoising (non-) residual ReLU CNN according to the theory of deep convolutional framelets (TDCF). In the figure, the noisy observation $\underline{y}$ is composed by two vertical bars plus uncorrelated Gaussian noise. Furthermore, for this example, the encoding and decoding convolution filters ($\mathbf{K}$ and $\widetilde{\mathbf{K}}$, respectively) are the Haar basis of the 2D discrete wavelet transform and its phase-inverted counterparts. Given the content of the image, the image in the decomposed domain $\mathbf{K}\mathbf{y}$ produces only a weak activation for the vertical and diagonal filters ($\mathbf{w_{LH}}$ and $\mathbf{w_{HH}}$, respectively) and those feature maps contain mainly noise. In the case of the non-residual network, the ReLUs and biases suppress the channels with low activation (see column $(\mathbf{K}\mathbf{y} + \underline{b})_+$), which is akin to the low-rank approximation. In contrast, in the residual example, the channels with image content are suppressed, while preserving the uncorrelated noise. Finally, the decoding section reconstructs the noise-free estimate $\tilde{\mathbf{x}}$ for the non-residual network or the noise estimate $\hat{\boldsymbol{\eta}}$ for the residual example, where it is subtracted from $\mathbf{y}$ to compute the noiseless estimate $\hat{\mathbf{x}}$.

filtered image. Note that this process is analogous to the low-rank decomposition described in Section IV-B1. In the case of *residual* networks, the CNN learns to estimate the noise, which means that in that configuration the ReLU non-linearities suppress the channels with high activation.

A visual example of the low-rank approximation in ReLU CNNs is shown in Fig. 8, which illustrates the operation of an idealized single-layer encoding-decoding ReLU CNN operating both, in residual and non-residual way. It can be noted The ReLU activation suppresses specific channels in the sparse decomposition provided by the encoder, thereby preserving the low-rank structures in the non-residual network. Alternatively, in the residual example, the ReLUs eliminate the feature maps with high activation, which results in a noise estimate that is subtracted from the input to estimate the

noiseless signal.

### B. Shrinkage and clipping-based CNNs

Just as ReLU networks, the encoder of shrinkage networks [17], [18], [19] separates the input signal in a multi-channel representation. As a second processing stage, the shrinkage networks estimate the noiseless encoded signal by cancelling the low-amplitude pixels in the feature maps in a process akin to the MAP estimate of Section IV-C2. As final step, the encoder reconstructs the estimated noiseless image. Note that the use of shrinkage functions reduces the number of channels required by ReLU counterparts to achieve perfect signal reconstruction, because the shrinkage activation preserves positive and negative values, while ReLU only preserves the positive part of the signal.
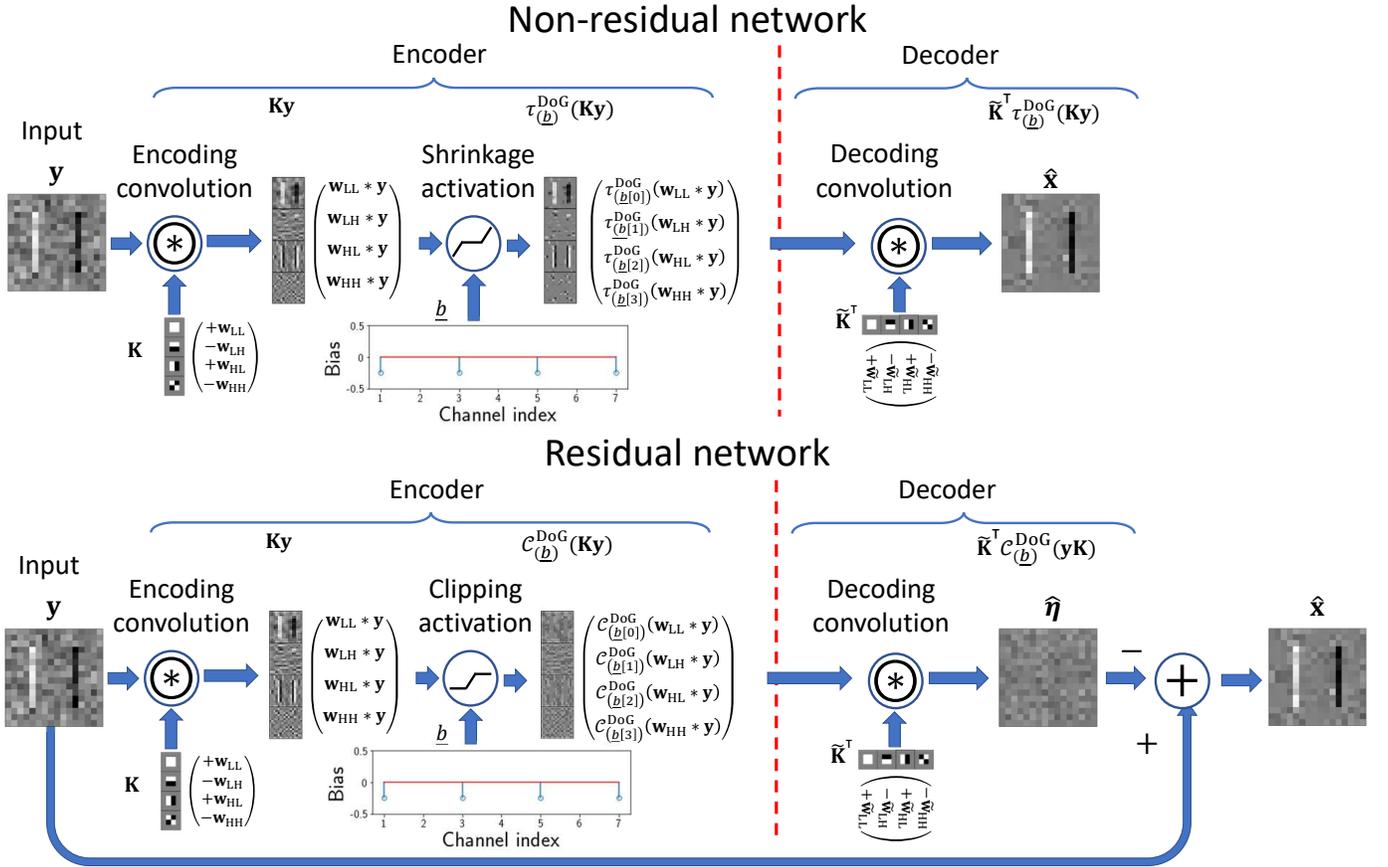
Fig. 9. Operation of denoising in shrinkage and clipping networks. In the non-residual configuration, the noisy signal $\mathbf{y}$ is decomposed by a set of convolution filters, which for this example are the 2D Haar basis functions of the discrete wavelet transform ($\mathbf{Ky}$). As a second step, the semi-hard shrinkage produces a MAP estimate of the noiseless detail bands/feature maps ($\tau_{(\underline{b})}^{\text{DoG}}(\mathbf{Ky})$). As third and final step, the decoder maps the estimated noiseless encoded signal to the original image domain. In the residual network, the behavior is similar, but the activation layer is a clipping function which performs a MAP estimate of the noise in the feature maps, which is reconstructed by the decoder to generate the noise estimate $\hat{\boldsymbol{\eta}}$. After reconstruction, the noise estimate is subtracted from the noisy observation $\mathbf{y}$ to generate the noise-free estimate $\tilde{\mathbf{x}}$.

As shown in Section III-A, in residual learning, a given encoding-decoding network estimates the noise signal $\boldsymbol{\eta}$, so that it can be subtracted from the noisy observation $\mathbf{y}$ to generate the noiseless estimate $\hat{\mathbf{x}}$. As shown in Section IV-C3, in the framelet domain this is achieved by preserving the low-amplitude values of the feature maps by clipping the signal. Therefore in residual networks, the shrinkage functions can be explicitly replaced by clipping activations.

Visual examples of the operation of a single-layer shrinkage and clipping networks are presented in Fig. 8, where it can be noted that the operation of shrinkage and clipping networks is analogous to their ReLU counterparts, with the main difference that shrinkage and clipping networks do not require phase-complements in the encoding and decoding layers as ReLU-based CNNs do.

### C. Shrinkage and clipping in ReLU networks

As addressed in Section IV-C, the soft-threshold function is the superposition of two ReLU activations. As a consequence, it is feasible that in ReLU CNNs shrinkage behavior could arise in addition to the low-rankness enforcement mentioned in Section V-A. It should be noted that this only can happen if the number of channels of the encoder and decoder complies with

the redundancy constraints of theory of deep convolutional framelets and if the decoder is linear. To prove this, Eq. (31) is reparametrized as

$$\hat{\mathbf{d}} = \tilde{\mathbf{K}}^{\mathsf{T}}(\mathbf{Ky} + \underline{b})_{+}, \tag{50}$$

where convolution filters $\mathbf{K}$ and $\tilde{\mathbf{K}}^{\mathsf{T}}$ are defined by $\mathbf{K} = \left(\begin{pmatrix} \mathbf{I} & -\mathbf{I} \end{pmatrix}\right)^{\mathsf{T}}$ and $\tilde{\mathbf{K}}^{\mathsf{T}} = \left(\begin{pmatrix} \mathbf{I} & -\mathbf{I} \end{pmatrix}\right)$, respectively, and $\underline{b} = \begin{pmatrix} -t & -t \end{pmatrix}^{\mathsf{T}}$ represents the threshold value.

In addition to the soft-shrinkage, note that the clipping function described by Eq. (34) also can be expressed by Eq. (50) if $\mathbf{K} = \left(\begin{pmatrix} \mathbf{I} & -\mathbf{I} & \mathbf{I} & -\mathbf{I} \end{pmatrix}\right)^{\mathsf{T}}$, $\tilde{\mathbf{K}}^{\mathsf{T}} = \left(\begin{pmatrix} \mathbf{I} & -\mathbf{I} & -\mathbf{I} & \mathbf{I} \end{pmatrix}\right)$ and $\underline{b} = \begin{pmatrix} 0 & 0 & -t & -t \end{pmatrix}^{\mathsf{T}}$. It can be noted that representing the clipping function in convolutional form requires *four* times more channels than the original input signal.

It should be noted that the ability of ReLUs to approximate other signals has also been noticed Daubechies *et al.* [31], who have proven that deep ReLU CNNs are universal function approximators. In addition, Ye and Sung [13] have demonstrated that the ReLU function is the main source of the high-approximation power of CNNs.

---

**Network depth**

*The relationship between network depth and low-rank approximation.* It should be noted that one of the key elements of CNNs is the network depth, which we address in this section. To illustrate the effect of network depth, assume an arbitrary $N$-layer encoding-decoding CNN, in which the encoding layers are defined by

$$\mathbf{E}_0 = (\mathbf{K}_0\mathbf{y} + \underline{b}_0)_+,$$
$$\mathbf{E}_1 = (\mathbf{K}_1\mathbf{E}_0 + \underline{b}_1)_+,$$
$$\mathbf{E}_2 = (\mathbf{K}_2\mathbf{E}_1 + \underline{b}_2)_+, \tag{46}$$
$$\vdots$$
$$\mathbf{E}_{N-1} = (\mathbf{K}_{N-1}\mathbf{E}_{N-2} + \underline{b}_{N-1})_+,$$

$$\boxed{\mathbf{E}_n = (\mathbf{K}_n\mathbf{E}_{n-1} + \underline{b}_n)_+} \quad. \tag{47}$$

Here, $\mathbf{E}_n$ represents the encoded signal at the $n$-th decomposition level, while $\mathbf{K}_n$, $\underline{b}_n$ are the convolution weights and biases for the $n$-th encoding layer, respectively. As addressed in Sections IV-C1 and V-A, the role of the ReLU activations is to enforce sparsity and non-negativity, which can be interpreted as the process of suppressing non-informative bases in the low-rank approximation algorithm. Consequently, every encoded signal $\mathbf{E}_n$, is an encoded *sparsified* version of the signal $\mathbf{E}_{n-1}$. In order to recover the signal, we apply the decoder part of the CNN, given by

$$\tilde{\mathbf{E}}_{N-1} = (\tilde{\mathbf{K}}_{N-1}^\intercal\mathbf{E}_{N-1} + \underline{\tilde{b}}_{N-1})_+,$$
$$\vdots$$
$$\tilde{\mathbf{E}}_1 = (\tilde{\mathbf{K}}_2^\intercal\tilde{\mathbf{E}}_2 + \underline{\tilde{b}}_2)_+, \tag{48}$$
$$\tilde{\mathbf{E}}_0 = (\tilde{\mathbf{K}}_1^\intercal\tilde{\mathbf{E}}_1 + \underline{\tilde{b}}_1)_+,$$
$$\hat{\mathbf{x}} = (\tilde{\mathbf{K}}_0^\intercal\tilde{\mathbf{E}}_0 + \underline{\tilde{b}}_0)_+,$$

$$\boxed{\tilde{\mathbf{E}}_{n-1} = (\tilde{\mathbf{K}}_n^\intercal\tilde{\mathbf{E}}_n + \underline{\tilde{b}}_n)_+} \quad. \tag{49}$$

Here, $\hat{\mathbf{x}}$ is the low-rank estimate/denoised version of the input signal $\mathbf{y}$, while $\tilde{\mathbf{E}}_n$, $\tilde{\mathbf{K}}_n^\intercal$, $\underline{\tilde{b}}_n$ are the decoded

signal components at the $n$-th composition level and the decoder convolution weights and biases for the $n$-th layer, respectively. In Eq. (49) every decoded signal $\tilde{\mathbf{E}}_n$ is the low-rank estimate of the encoded layer $\mathbf{E}_{(n-1)}$. It should be noted that the activation of each of the decoder layers $(\cdot + \underline{\tilde{b}}_n)_+$ can further enforce sparsity on the low-rank estimates $\tilde{\mathbf{E}}_{(n-1)}$.

*Summary.* In conclusion, the mathematical formulation of deep networks is analogous to a recursive data-driven low-rank approximation, where the input to the successive encoding-decoding pairs is the low-rank approximated encoded signal generated by the encoder of the previous level. Still, just as mentioned in the text box *Fitting low-rank approximation in ReLU CNNs*, low-rank approximation algorithms and CNNs are similar in terms of mathematical formulation, but we cannot ensure that the values obtained during training for the encoding, decoding filters and their biases have the properties needed to ensure that a CNN is an exact recursive data-driven low-rank approximation. For example, it is possible that the filters of the encoder and decoder do not reconstruct the signal perfectly, because this may not be necessary to reduce the loss function used to optimize the network.

*Is it possible to impose a tighter relation between low-rank approximation and CNNs?* In specific applications where signal preservation and interpretability is required (e.g. medical imaging) it is desirable that the operation of CNNs is closer to the low-rank approximation description. In order to achieve this, the CNNs embedded in frameworks such as the convolutional analysis operator [29] and FISTA-Net [30] explicitly train the filters $\mathbf{K}_n$ and $\tilde{\mathbf{K}}_n$ to have properties such as perfect signal reconstruction and sparsity. By enforcing these characteristics, the mathematical descriptions of the low-rank behavior and of CNNs are more similar and the models become inherently more interpretable and predictable on their operation.

---

### D. Additional links between encoding-decoding CNNs and existing signal processing techniques

Up to this moment, it has been assumed that the operation of the encoding and decoding convolution filters is limited to map the input image to a multi-channel representation and to reconstruct it (i.e. $\mathbf{K}$ and $\tilde{\mathbf{K}}^\intercal$ comply with $\tilde{\mathbf{K}}^\intercal(\mathbf{K})_+ = \mathbf{I} \cdot c$). Still, it is possible that –besides decomposition and synthesis tasks– the encoding-decoding structure also filters/colors the signal in a way that improves the image estimates. It should be noted that this implies that the perfect reconstruction encoding-decoding structure is no longer preserved. For example, considering the following linear encoding-decoding structure

$$\hat{\mathbf{x}} = \tilde{\mathbf{K}}^\intercal(\mathbf{K}\mathbf{y}), \tag{51}$$

which can be reduced to

$$\hat{\mathbf{x}} = \mathbf{k} * \mathbf{y}. \tag{52}$$

Here, $\mathbf{k} = \tilde{\mathbf{K}}^\intercal\mathbf{K}$, is optimized to reduce the distance between $\mathbf{y}$ and the ground-truth $\mathbf{x}$. Consequently, the equivalent filter $\mathbf{k}$ can be considered a *Wiener filter*. It should be noted that this text is not the first in addressing the potential Wiener-like behavior of a CNN. For example, Mohan *et al.* [14] suggested that by eliminating the bias of the convolution layers, the CNN could behave more akin to the Wiener filter and to be able to generalize better to unseen noise levels. It should be noted that by doing that the CNN can also behave akin to the switching behavior described by the theory of deep convolutional framelets, which can be described by the

equation

$$(z)_+ = \begin{cases} z, & \text{if } z \geq 0, \\ 0 & \text{if } z < 0, \end{cases} \tag{53}$$

where $z$ is a pixel which belongs to the signal $\mathbf{z} = \mathbf{k} * \mathbf{x}$. It can be observed that in contrast with the low-rank behavior described in Section V-A, in this case the switching behavior is only dependent on the correlation between signal $\mathbf{x}$ and the filter $\mathbf{k}$. Consequently, if the value of $z$ is positive, its value is preserved. On the contrary, if the correlation between $\mathbf{x}$ and $\mathbf{k}$ is negative, then the value of $z$ is cancelled. Consequently, the noise reduction becomes independent/invariant of the noise level. It can be observed that this effect can be considered a non-linear extension of the so-called signal annihilation filters [32].

It should be noticed that besides the low-rank approximation interpretation of ReLU-based CNNs, additional links to other techniques can be derived. For example, the decomposition and synthesis provided by the encoding-decoding structure is also akin to the non-negative matrix factorization (NMF) [33], in which a signal is factorized as a weighted sum of positive bases. In this conception, the feature maps are the bases, which are constrained to be positive by the ReLU function. Furthermore, an additional interpretation of encoding-decoding CNNs can be obtained by analyzing them from a low-dimensional manifold representation perspective [8]. Here, the convolution layers of CNNs are interpreted as two operations. On one hand, they can provide a Hankel representation. On the other hand, they provide a bottleneck which reduces the dimensionality of the manifold of image patches. It should be noticed that the Hankel-like structure that is attributed to the convolution layers of CNNs, has also been noticed by the theory of the deep convolutional framelets [12]. Two final connection with signal-processing and CNNs is the variational formulation combined with kernel-based methods [15] and the convolutional sparse coding interpretation of CNNs by Papyan *et al.* [16].

## VI. MATHEMATICAL ANALYSIS OF RELEVANT DESIGNS

In order to demonstrate the application of the principles summarized in Sections IV and V, this section analyzes relevant designs of ReLU and shrinkage CNNs. The analyses focus on three main aspects, which are: (1) the overall descriptions of the network architecture, (2) the signal reconstruction characteristics provided by the convolutional layers of the encoder and decoder sub-networks, and (3) the number operations $\mathcal{O}(\cdot)$ executed by the trainable parts of the network, since this will give insight on the computational requirements to execute each network and its overall complexity.

The signal reconstruction analysis provides a theoretical indication that a given CNN design can propagate *any* arbitrary signal when considering the use of ideal filters (i.e. they provide perfect reconstruction and are maximally sparse). In other words, for a fixed network architecture, there exists a selection of parameters (weights and biases) that make the neural network equal to the identity function. This result is important, because a design that cannot propagate arbitrary signals under ideal conditions, will potentially distort the signals that propagate through it *by design*. Consequently, this
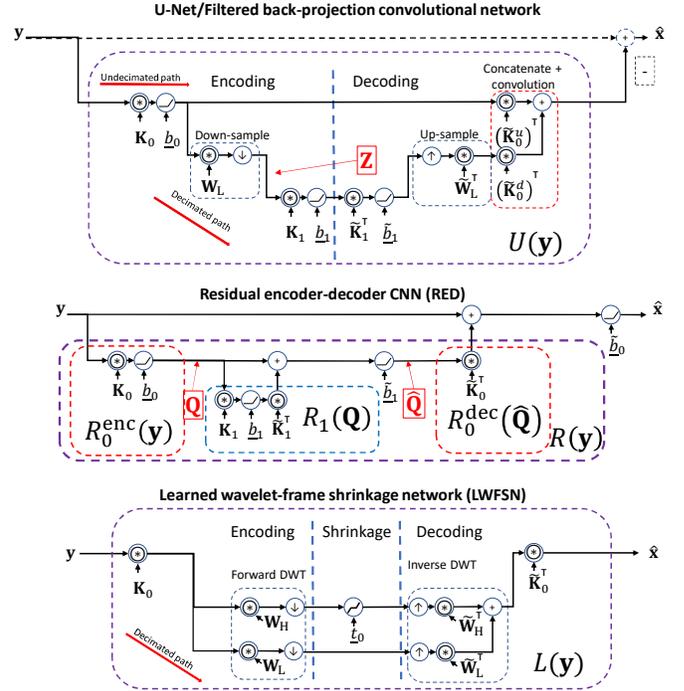


Fig. 10. Simplified structure of encoding-decoding ReLU CNNs. The displayed networks are the U-Net/filtered back-projection network the encoder-decoder residual CNN (RED) and finally, the learned wavelet-frame shrinkage network. Note that for all the designs, the encoding-decoding structures are indicated by dashed blocks. It should be beared in mind that the drawings are simplified, they do not contain normalization layers, are shallow, commonly appearing dual convolutions are drawn as one layer.

cannot be fixed by training with large datasets and/or with the application of any special loss term. In order to understand better the signal reconstruction analysis, we provide a brief example, where it is a non-residual CNN $G(\cdot)$, where we propagate a noiseless signal $\mathbf{x}$ contaminated with noise $\boldsymbol{\eta}$, so that

$$\mathbf{x} \approx G(\mathbf{x} + \boldsymbol{\eta}). \tag{54}$$

Here, an ideal CNN allows to propagate *any* $\mathbf{x}$, while cancelling the noise component $\boldsymbol{\eta}$, irrespective of the content of $\mathbf{x}$. If we switch our focus to an ideal *residual CNN* $R(\cdot)$, it is possible to observe that

$$\hat{\mathbf{x}} \approx R(\mathbf{y}) = \mathbf{y} - G(\mathbf{y}). \tag{55}$$

Here, $G(\cdot)$ is the encoding-decoding section of the residual network $R(\cdot)$. Consequently, it is desirable that the network $G(\cdot)$ is able to propagate the noise $\boldsymbol{\eta}$, while suppressing the noiseless signal $\mathbf{x}$, which is equivalent to

$$\boldsymbol{\eta} \approx G(\mathbf{x} + \boldsymbol{\eta}). \tag{56}$$

It should be noted that in both residual and non-residual cases, there are two behaviors. On one hand, there is a signal which the network decomposes and reconstructs (almost) perfectly, and on the other hand a signal is suppressed. The signal reconstruction analysis focuses on the signals that the network can propagate or reconstruct, rather than the signal cancellation behavior. In consequence, we focus on the linear part of $G(\cdot)$ (i.e. its convolution structure), of which, according to

Section V-A, we assume that it handles the decomposition and reconstruction of the signal within the CNN. It should be noted that the idealized model assumed here, is only considered for analysis purposes, since practical implementations do no guarantee that this exact behavior is factually obtained. The reader is referred to Section V-D and the text boxes *Fitting low-rank approximation in ReLU CNNs* and *Network depth*.

In order to test the perfect reconstruction in non-residual CNNs, we propose the following procedure. (1) We assume an *idealized* model $G(\cdot)$, where its convolution filters $\mathbf{K}_n$ and $\tilde{\mathbf{K}}_n$ comply with the phase-complementary tight framelet condition and where the biases and non-linearities suppress low-amplitude (and negative for ReLU activations) samples from the feature maps. (2) The biases/thresholds of ReLU/shrinkage CNNs are set to zero (or to infinity for clipping activations). It can be observed that this condition prevents the low-rank (or high-rank for residual models) approximation behavior of the idealized CNN. Under this circumstance, it should be possible to prove that the analyzed CNN can perfectly reconstruct *any* signal. (3) The last step involves simplifying the mathematical description of the resulting model of the previous point. The mathematical simplification of the model should lead to the identity function if the model complies with the perfect reconstruction property.

To conclude the explanation on the perfect reconstruction analysis, we provide two relevant considerations. First, it can be claimed that residual networks, such as the model $R(\mathbf{y}) = \mathbf{y} - G(\mathbf{y})$ discussed in Eq. (55), is able to reconstruct any signal when $G(\mathbf{y}) = 0$ for any $\mathbf{y} = \mathbf{x} + \boldsymbol{\eta}$. Still, this does not convey information about the behavior of the encoding-decoding network $G(\cdot)$, which should be able to perform perfect decomposition and reconstruction of the noise signal $\boldsymbol{\eta}$, as discussed in Eq. (55). To avoid this trivial solution, instead of analyzing the network $R(\cdot)$, the analysis described for non-residual models is applied to the encoding-decoding structure $G(\cdot)$, which means that the residual connection is excluded from the analysis.

The second concluding remark is that in order to distinguish the equations of the perfect signal reconstruction analysis from other models, we specify the analyzed designs of the perfect reconstruction models, in which the low-rank approximation behavior is avoided by setting the bias values to zero, with a special operator $\mathscr{P}\{\cdot\}$.

For the analyses regarding the total number of operations of the trainable parameters, it is assumed that the tensors $\mathbf{K}_0$, $\tilde{\mathbf{K}}_0^{\mathsf{T}}$, $(\tilde{\mathbf{K}}_0^{\mathrm{u}})^{\mathsf{T}}$, $(\tilde{\mathbf{K}}_0^{\mathrm{d}})^{\mathsf{T}}$, $\mathbf{K}_1$ and $\tilde{\mathbf{K}}_1^{\mathsf{T}}$ shown in Fig. 10 have dimensions $(C_0 \times 1 \times N_{\mathrm{f}} \times N_{\mathrm{f}})$, $(1 \times C_0 \times N_{\mathrm{f}} \times N_{\mathrm{f}})$, $(1 \times C_0 \times N_{\mathrm{f}} \times N_{\mathrm{f}})$, $(1 \times C_0 \times N_{\mathrm{f}} \times N_{\mathrm{f}})$, $(C_1 \times C_0 \times N_{\mathrm{f}} \times N_{\mathrm{f}})$, $(C_0 \times C_1 \times N_{\mathrm{f}} \times N_{\mathrm{f}})$, respectively. Here, $C_0$ and $C_1$ represent the number of channels after the first and second convolution layer, all the convolution filters are assumed to be square with $N_{\mathrm{f}} \times N_{\mathrm{f}}$ pixels. Furthermore, the input signal $\mathbf{x}$ has dimensions ($1 \times 1 \times N_{\mathrm{r}} \times N_{\mathrm{c}}$), where $N_{\mathrm{r}}$ and $N_{\mathrm{c}}$ denote the numbers of rows and columns, respectively.

The analyses shown for the different networks in this article have the following limitations. (1) The analyzed networks have only enough decomposition levels and convolution layers to understand their basic operation. The motivation for

this simplification is to keep the analyses short and clear. Moreover, the same principles can be extended to deeper networks. Since the same single-decomposition CNNs would be recursively embedded within the given architectures. (2) The normalization layers are not considered, because they are linear operators which provide mean shifts and amplitude scaling. Consequently, for analysis purposes it can be assumed that they are embedded in the convolution weights. (3) For every encoder convolution kernel it is assumed that there is at least one decoder filter. (4) No co-adaptation between the filters of the encoder and decoder layers is considered.

The remainder of this section shows analyses of a selection of a few representative designs. Specifically, the chosen designs are the *U-Net* [34] and its residual counterpart, the *filtered back-projection network* [21] [1]. Additional designs analyzed here are the *the residual encoder-decoder CNN* [5] [2], as well as the *learned wavelet-shrinkage network* (LWFSN) [3]. For reference, all the designs are portrayed in Fig. 10.

### A. U-Net/filtered back-projection network

*1) U-Net – overview of the design:* The first networks analyzed are the U-Net and filtered back-projection networks both of which share the encoding-decoding structure $U(\cdot)$. However, they differ in the fact that the U-Net is non-residual, while the filtered back-projection network operates in residual configuration. Therefore, the estimate of the noiseless signal $\hat{\mathbf{x}}$ from the noisy observation $\mathbf{y}$ in the conventional U-Net is achieved by

$$\hat{\mathbf{x}} = U(\mathbf{y}), \tag{57}$$

whereas in the filtered back-projection network, $U(\cdot)$ is used in residual configuration, which is equivalent to

$$\hat{\mathbf{x}} = \mathbf{y} - U(\mathbf{y}). \tag{58}$$

If we now switch focus to the encoding-decoding structure of the U-Net $U(\mathbf{y})$, it can be shown that it is described by

$$U(\mathbf{y}) = U^{\mathrm{u}}(\mathbf{y}) + U^{\mathrm{d}}(\mathbf{y}), \tag{59}$$

where $U^{\mathrm{u}}(\mathbf{y})$ corresponds to the *undecimated path* and is defined by

$$U^{\mathrm{u}}(\mathbf{y}) = \left(\tilde{\mathbf{K}}_0^{\mathrm{u}}\right)^{\mathsf{T}} \left(\mathbf{K}_0 \mathbf{y} + \underline{b}_0\right)_+, \tag{60}$$

while the decimated path is

$$U^{\mathrm{d}}(\mathbf{y}) = \left(\tilde{\mathbf{K}}_0^{d}\right)^{\mathsf{T}} \tilde{\mathbf{W}}_{\mathrm{L}}^{\mathsf{T}} f_{(2\uparrow)} \left( \left(\tilde{\mathbf{K}}_1^{\mathsf{T}} (\mathbf{K}_1 \mathbf{Z} + \underline{b}_1)_+ + \underline{\tilde{b}}_1\right)_+ \right). \tag{61}$$

Here, signal $\mathbf{Z}$ is defined by

$$\mathbf{Z} = f_{(2\downarrow)} \left(\mathbf{W}_{\mathrm{L}} (\mathbf{K}_0 \mathbf{y} + \underline{b}_0)_+\right). \tag{62}$$

Note the decimated path contains two nested encoding-decoding architectures, as observed by Jin *et al.* [21], who

[1]Matlab implementation by their authors available at https://github.com/panakino/FBPConvNet

[2]Pytorch implementation by their authors available at https://github.com/SSinyu/RED-CNN

[3]Pytorch implementation available at https://github.com/LuisAlbertZM/demo_LWFSN_TMI and interactive demo available at IEEE's code ocean https://codeocean.com/capsule/9027829/tree/v1. The demo also includes as reference pytorch implementations of FBPConvNet and the tight-frame U-Net.

has acknowledged that the nested filtering structure is akin to the (learned) iterative soft thresholding algorithm [27], [28].

*2) U-Net – signal reconstruction analysis:* To prove if the U-Net can perfectly reconstruct *any* signal, we assume that the biases are equal to zero, on this condition the network $\mathscr{P}\{U\}(\mathbf{y})$ is defined by

$$\mathscr{P}\{U\}(\mathbf{y}) = \mathscr{P}\{U^{\mathrm{u}}\}(\mathbf{y}) + \mathscr{P}\{U^{\mathrm{d}}\}(\mathbf{y}), \tag{63}$$

where sub-network $\mathscr{P}\{U^{\mathrm{u}}\}(\cdot)$ is defined by

$$\mathscr{P}\{U^{\mathrm{u}}\}(\mathbf{y}) = \big(\tilde{\mathbf{K}}_0^{\mathrm{u}}\big)^{\mathsf{T}}(\mathbf{K}_0\mathbf{y})_+. \tag{64}$$

Assuming that $\big(\mathbf{K}_0, \tilde{\mathbf{K}}_0^{\mathrm{u}}\big)$ is a complementary-phase tight-framelet pair then $\mathscr{P}\{U^{\mathrm{u}}\}(\mathbf{y})$ is simplified to

$$\mathscr{P}\{U^{\mathrm{u}}\}(\mathbf{y}) = \mathbf{y} \cdot c_0. \tag{65}$$

Furthermore, the low-frequency path is

$$U^{\mathrm{d}}(\mathbf{y}) = \big(\tilde{\mathbf{K}}_0^d\big)^{\mathsf{T}}\tilde{\mathbf{W}}_{\mathrm{L}}^{\mathsf{T}} f_{(2\uparrow)}\Big(\big(\tilde{\mathbf{K}}_1^{\mathsf{T}}(\mathbf{K}_1\mathscr{P}\{\mathbf{Z}\})_+\big)_+\Big), \tag{66}$$

Where $\mathscr{P}\{\mathbf{Z}\}$ is defined by

$$\mathscr{P}\{\mathbf{Z}\} = f_{(2\downarrow)}\big(\mathbf{W}_{\mathrm{L}}(\mathbf{K}_0\mathbf{y})_+\big). \tag{67}$$

If $\mathbf{K}_1$ is a phase-complementary tight frame, we know that $\tilde{\mathbf{K}}_1^{\mathsf{T}}(\mathbf{K}_1\mathbf{Z})_+ = \mathbf{Z} \cdot c_1$. Consequently, Eq. (66) becomes

$$\mathscr{P}\{U^{\mathrm{d}}\}(\mathbf{y}) = \big(\tilde{\mathbf{K}}_0^d\big)^{\mathsf{T}}\tilde{\mathbf{W}}_{\mathrm{L}}^{\mathsf{T}} f_{(2\uparrow)}(f_{(2\downarrow)}\big((\mathbf{W}_{\mathrm{L}}\mathbf{K}_0\mathbf{y})_+\big))\cdot c_1. \tag{68}$$

Here, it can be noticed that if $\mathbf{K}_0$ is a phase-complementary tight framelet, then $\mathscr{P}\{U^{\mathrm{d}}\}(\mathbf{y})$ approximates a low-pass version of $\mathbf{y}$, or equivalently

$$\mathscr{P}\{U^{\mathrm{d}}\}(\mathbf{y}) \approx \mathcal{W}_{\mathrm{L}}\mathbf{y} \cdot c_1, \tag{69}$$

where $\mathcal{W}_{\mathrm{L}}$ is a low-pass filter. Finally, substituting Eq. (65) and Eq. (69) in Eq. (63) results in

$$\mathscr{P}\{U\}(\mathbf{y}) \approx (\mathbf{I} \cdot c_0 + \mathcal{W}_{\mathrm{L}} \cdot c_1)\mathbf{y}. \tag{70}$$

This result proves that the design of the U-Net cannot evenly reconstruct all the frequency of $\mathbf{y}$ unless $c_1 = 0$, in which case, the whole low-frequency branch of the network is ignored. Note that this limitation is inherent to its design and cannot be circumvented by training with large datasets and/or with any loss function.

*3) U-Net – number of operations:* It can be noted that encoding filter $\mathbf{K}_0$, convolves $\mathbf{x}$ at its original resolution and maps it to a tensor with $C_0$ channels. Therefore, the number of operations $\mathcal{O}(\cdot)$ for kernel $\mathbf{K}_0$ is $\mathcal{O}(\mathbf{K}_0) = C_0 \cdot N_{\mathrm{r}} \cdot N_{\mathrm{c}} \cdot N_{\mathrm{f}}^2$ [FLOPS] (floating-point operations). Conversely, due to the symmetry between encoder and decoder filters, $\mathcal{O}(\tilde{\mathbf{K}}_0^{\mathrm{u}}) = \mathcal{O}(\tilde{\mathbf{K}}_0^{\mathrm{d}}) = \mathcal{O}(\mathbf{K}_0)$. Furthermore, for this design, filter $\mathbf{K}_1$ processes the signal encoded by $\mathbf{K}_0$, which is down-sampled by a factor of one half and maps it from $C_0$ to $C_1$ channels, this results in the estimated operation cost $\mathcal{O}(\mathbf{K}_1) = \mathcal{O}(\tilde{\mathbf{K}}_1) = C_0 \cdot C_1 \cdot N_{\mathrm{r}} \cdot N_{\mathrm{c}} \cdot N_{\mathrm{f}}^2 \cdot (2)^{-2}$ [FLOPs]. Finally, adding the contributions of filters $\mathbf{K}_0$, $\tilde{\mathbf{K}}_0^{\mathrm{u}}$, $\tilde{\mathbf{K}}_0^{\mathrm{d}}$, $\mathbf{K}_1$ and $\tilde{\mathbf{K}}_1$ results in

$$\mathcal{O}(U) = (3 + 2^{-1} \cdot C_1) \cdot C_0 \cdot N_{\mathrm{r}} \cdot N_{\mathrm{c}} \cdot N_{\mathrm{f}}^2 \text{ [FLOPS]} \tag{71}$$

*4) U-Net – Concluding remarks:* The U-Net/FBPConvNet is a flexible multi-resolution architecture. Still, as it has been shown, the pooling structure of this CNN may be sub-optimal for noise reduction applications because this configuration does not allow to recover the frequency information of the signal evenly. This has been noted and fixed by Han and Ye [6], who introduced the so-called tight-frame U-Net. In which the down/up-sampling structure is replaced by the discrete wavelet transform and its inverse. This simple modification overcomes the limitations of the U-Net and improved its performance for artifact removal in compressed sensing imaging.

### B. Residual encoder-decoder CNN

*1) Residual encoder-decoder CNN – overview of the design:* The residual encoder-decoder CNN shown in Fig. 10 consists of nested single-layer residual encoding-decoding networks. For example, in the network showcased in Fig. 10 it is visible that network $R_1(\cdot)$ is nested into $R_0(\cdot)$. Furthermore, for this case the image estimate is given by

$$\hat{\mathbf{x}} = (\mathbf{y} + R_0(\mathbf{y}) + \tilde{\underline{b}}_0)_+, \tag{72}$$

in which $R_0(\cdot)$ is the outer residual network and $\tilde{\underline{b}}_0$ is the bias for the output layer. Note that the ReLU placed at the output layer intrinsically assumes that the estimated signal $\hat{\mathbf{x}}$ is positive.

From Eq. (72), the output of the sub-network $R_0(\cdot)$ is defined by

$$\mathbf{Z} = R_0^{\mathrm{dec}}(\hat{\mathbf{Q}}). \tag{73}$$

Here, the decoder $R_0^{\mathrm{dec}}(\cdot)$ is defined by

$$R_0^{\mathrm{dec}}(\hat{\mathbf{Q}}) = \tilde{\mathbf{K}}_0^{\mathsf{T}}\hat{\mathbf{Q}}. \tag{74}$$

In the above, $\hat{\mathbf{Q}}$ is the noiseless estimate of the intermediate signal $\mathbf{Q}$ and it is defined by

$$\hat{\mathbf{Q}} = (\mathbf{Q} + R_1(\mathbf{Q}) + \tilde{\underline{b}}_1)_+, \tag{75}$$

where the network $R_1(\cdot)$ is

$$R_1(\mathbf{Q}) = \tilde{\mathbf{K}}_1^{\mathsf{T}}\big(\mathbf{K}_1\mathbf{Q} + \underline{b}_1\big)_+. \tag{76}$$

Furthermore, $\mathbf{Q}$ represents the signal encoded by $R_0(\cdot)$, or equivalently

$$\mathbf{Q} = R_0^{\mathrm{enc}}(\mathbf{y}), \tag{77}$$

where $R_0^{\mathrm{enc}}(\cdot)$ is defined by

$$R_0^{\mathrm{enc}}(\mathbf{y}) = \mathbf{K}_0\mathbf{y}. \tag{78}$$

*2) Residual encoder-decoder CNN – signal reconstruction analysis:* As mentioned earlier, the residual encoder-decoder CNN is composed by nested residual blocks, which are independently analyzed to study the reconstruction characteristics of this network. First, block $R_1(\cdot)$, is given by

$$\mathscr{P}\{R_1\}(\mathbf{Q}) = \tilde{\mathbf{K}}_1^{\mathsf{T}}(\mathbf{K}_1\mathbf{Q})_+. \tag{79}$$

Under complementary-phase tight-frame assumptions for the pair $(\mathbf{K}_1, \tilde{\mathbf{K}}_1)$, Eq. (79) reduces to

$$\mathscr{P}\{R_1\}(\mathbf{Q}) = \mathbf{Q}, \tag{80}$$

which shows that the encoder and decoder $R_1(\cdot)$ can approximately reconstruct any signal. Now switching to $R_0$ it can be observed that the linear part is

$$\mathscr{P}\{R_0\}(\mathbf{y}) = \tilde{\mathbf{K}}_0^\intercal (\mathbf{K}_0 \mathbf{y})_+. \tag{81}$$

Just as with $R_1(\cdot)$, it is assumed that the convolution kernels are tight-framelets. Therefore Eq. (81) becomes

$$\mathscr{P}\{R_0\}(\mathbf{y}) = \mathbf{y}. \tag{82}$$

Consequently, $R_0(\cdot)$ and $R_1(\cdot)$ can reconstruct any arbitrary signal under complementary-phase tight-frame assumptions.

*3) Residual encoder-decoder CNN – number of operations:* In this case, all the convolution layers operate at the original resolution of image $\mathbf{x}$. Therefore, the number of operations $O(\cdot)$ for kernel $\mathbf{K}_0$ and $\tilde{\mathbf{K}}_0$ is $O(\mathbf{K}_0) = O(\tilde{\mathbf{K}}_0) = C_0 \cdot N_\mathrm{r} \cdot N_\mathrm{c} \cdot N_\mathrm{f}^2$ [FLOPs], while $\mathbf{K}_1$ and $\tilde{\mathbf{K}}_1$ are requiring $O(\mathbf{K}_1) = O(\tilde{\mathbf{K}}_1) = C_0 \cdot C_1 \cdot N_\mathrm{r} \cdot N_\mathrm{c} \cdot N_\mathrm{f}^2$ [FLOPs]. By adding the contributions of both encoding-decoding pairs, the total operations for the residual encoder-decoder becomes

$$\mathcal{O}(R) = 2 \cdot (1 + C_1) \cdot C_0 \cdot N_\mathrm{r} \cdot N_\mathrm{c} \cdot N_\mathrm{f}^2 \text{ [FLOPS].} \tag{83}$$

*4) Residual encoder-decoder CNN – concluding remarks:* The residual encoder-decoder network consists on a set of nested single-resolution residual encoding-decoding CNNs. The single-resolution design increases its computation cost with respect to multi-resolution designs such as the U-Net. In addition, it should be noted that the use of a ReLU as output layer of the encoder-decoder residual network forces the signal estimates to be positive, but this is not always convenient. For example in computed tomography imaging, it is common that images contain positive and negative values.

### C. Learned wavelet-frame shrinkage network

*1) Learned wavelet-frame shrinkage network – description of the architecture:* The learned wavelet-frame shrinkage network is a multi-resolution architecture in which the discrete wavelet transform is used for down/up-sampling and also as part of the decomposition where shrinkage is applied. In this CNN, the noiseless estimates are produced by

$$\hat{\mathbf{x}} = L(\mathbf{y}), \tag{84}$$

where $L(\cdot)$ represents the encoding-decoding structure of the learned wavelet-frame shrinkage network and the encoding-decoding network $L(\cdot)$ is

$$L(\mathbf{y}) = L_\mathrm{L}(\mathbf{y}) + L_\mathrm{H}(\mathbf{y}). \tag{85}$$

Here, the high-frequency path is given by

$$L_\mathrm{H}(\mathbf{y}) = \tilde{\mathbf{K}}_0^\intercal \tilde{\mathbf{W}}_\mathrm{H}^\intercal f_{(2\uparrow)}\big(\tau_{(\underline{t}_0)}^\mathrm{LET}\big(f_{(2\downarrow)}\big(\mathbf{W}_\mathrm{H}\mathbf{K}_0\mathbf{y}\big)\big)\big). \tag{86}$$

Note that in this design the encoder leverages the filter $\mathbf{W}_\mathrm{H}$ for generating a sparse signal prior to the shrinkage stage, i.e. $\tau_{(\underline{t}_0)}^\mathrm{LET}\big(f_{(2\downarrow)}\big(\mathbf{W}_\mathrm{H}\mathbf{K}_0\mathbf{y}\big)\big)$. Meanwhile, the low-frequency path $L_\mathrm{L}(\cdot)$ is

$$L_\mathrm{L}(\mathbf{y}) = \tilde{\mathbf{K}}_0^\intercal \tilde{\mathbf{W}}_\mathrm{L}^\intercal f_{(2\uparrow)}\big(f_{(2\downarrow)}\big(\mathbf{W}_\mathrm{L}\mathbf{K}_0\mathbf{y}\big)\big). \tag{87}$$
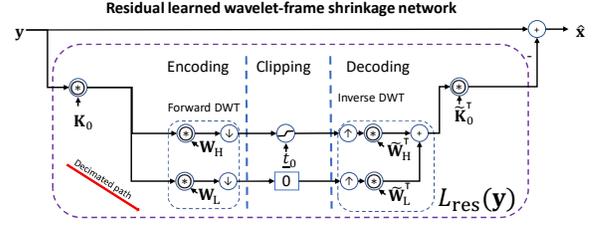


Fig. 11. Residual version of the learned wavelet-frame shrinkage network. It can be noticed that the low-frequency branch of the network is nulled. In deeper networks it would further decomposed and the nulling would be activated at the deepest level (lowest resolution).

*2) Learned wavelet-frame shrinkage network – signal reconstruction analysis:* When analyzing the signal propagation of the learned wavelet-frame shrinkage network, we set the threshold level $\underline{t}_0 = 0$. This turns Eq. (85) turns into

$$\mathscr{P}\{L\}(\mathbf{y}) = \mathscr{P}\{L_\mathrm{L}\}(\mathbf{y}) + \mathscr{P}\{L_\mathrm{H}\}(\mathbf{y}). \tag{88}$$

Here, $\mathscr{P}\{L_\mathrm{H}\}(\cdot)$ is defined by

$$\mathscr{P}\{L_\mathrm{H}\}(\mathbf{y}) = \tilde{\mathbf{K}}_0^\intercal \tilde{\mathbf{W}}_\mathrm{H}^\intercal f_{(2\uparrow)}\big(f_{(2\downarrow)}\big(\mathbf{W}_\mathrm{H}\mathbf{K}_0\mathbf{y}\big)\big), \tag{89}$$

while the low-frequency path $\mathscr{P}\{L_\mathrm{L}\}(\cdot)$ is mathematically described by

$$\mathscr{P}\{L_\mathrm{L}\}(\mathbf{y}) = \tilde{\mathbf{K}}_0^\intercal \tilde{\mathbf{W}}_\mathrm{L}^\intercal f_{(2\uparrow)}\big(f_{(2\downarrow)}\big(\mathbf{W}_\mathrm{L}\mathbf{K}_0\mathbf{y}\big)\big). \tag{90}$$

Substituting Eq. (89) and Eq. (90) in Eq. (88) results in the equation

$$\mathscr{P}\{L\}(\mathbf{y}) = \tilde{\mathbf{K}}_0^\intercal \tilde{\mathbf{W}}^\intercal f_{(2\uparrow)}\big(f_{(2\downarrow)}\big(\mathbf{W}\mathbf{K}_0\mathbf{y}\big)\big). \tag{91}$$

For the discrete wavelet transform, it holds that $\mathbf{Q} = \tilde{\mathbf{W}}^\intercal f_{(2\uparrow)}\big(f_{(2\downarrow)}\big(\mathbf{W}\mathbf{Q}\big)\big)$. Consequently, Eq. (91) is simplified to

$$\mathscr{P}\{L\}(\mathbf{y}) = \tilde{\mathbf{K}}_0^\intercal \mathbf{K}_0\mathbf{y}. \tag{92}$$

Assuming that $\mathbf{K}_0$ is a tight framelet i.e. $\tilde{\mathbf{K}}_0^\intercal \mathbf{K}_0 = \mathbf{I} \cdot c$, with $c = 1$, then

$$\mathscr{P}\{L\}(\mathbf{y}) = \mathbf{y}. \tag{93}$$

This proves that the encoding-decoding section of the learned wavelet-frame shrinkage network allows for perfect signal reconstruction.

*3) Learned wavelet-frame shrinkage network – number of operations:* The learned wavelet-frame shrinkage network contains a simpler convolution structure than the networks reviewed up to this moment. Therefore, for a single-level decomposition architecture, the total number of operations is

$$\mathcal{O}(L) = 2 \cdot C_0 \cdot N_\mathrm{r} \cdot N_\mathrm{c} \cdot N_\mathrm{f}^2 \text{ [FLOPS].} \tag{94}$$

*4) Learned wavelet-frame shrinkage network – residual variant:* To illustrate the use of clipping activations in residual noise reduction, the residual version of the learned wavelet-frame shrinkage network is included. Note that there are two main differences with the conventional learned wavelet-frame shrinkage network. First, the shrinkage functions are replaced by clipping activations. Second, the low-frequency signal is suppressed. This is performed because the original design of the learned wavelet-frame shrinkage network does
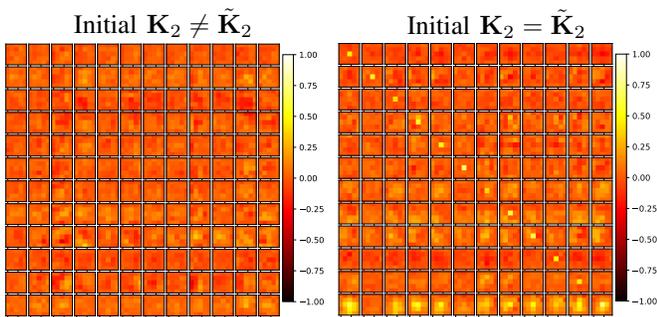
Fig. 12. Phase-complementary tight-framelet test for the trained toy network, initialized with random weights. The left figure shows the product $\tilde{\mathbf{K}}_2^{\mathsf{T}}(\mathbf{K}_2)_+$, where the initialization of $\mathbf{K}_2$ and $\tilde{\mathbf{K}}_2$ is different. It can be seen that the pair $(\mathbf{K}_2, \tilde{\mathbf{K}}_2)$ does not comply with the complementary-phase framelet criterion of Eq. (95). This contrasts with the right result, which displays the result of the product $\tilde{\mathbf{K}}_2^{\mathsf{T}}(\mathbf{K}_2)_+$, for the same CNN, but where the initial values of $\tilde{\mathbf{K}}_2$ and $\mathbf{K}_2$ are identical. For this initialization, the filters approximate the complementary-phase tight-framelet criterion.
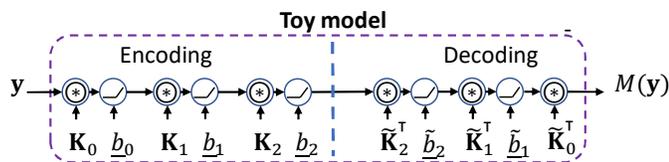


Fig. 13. Toy model used for experiment on the properties of the filters of a trained CNN. The dimensions for tensors $\mathbf{K}_0$, $\mathbf{K}_1$ and $\mathbf{K}_2$ are $(6 \times 1 \times 3 \times 3)$, $(12 \times 6 \times 3 \times 3)$ and $(24 \times 12 \times 3 \times 3)$. The network is symmetric and the filter dimensions for the decoder convolution kernels $\tilde{\mathbf{K}}_n$ are the same as there corresponding encoding kernel $\mathbf{K}_n$.

not have any non-linearities in that section. This is akin to the low-frequency nulling proposed by Kwon and Ye [35]. The modified learned wavelet-frame shrinkage network is shown in Fig. 11. It can be observed that by setting to zero the low-frequency branch of the design the model is inherently assuming that the noise is high-pass.

*5) (Residual) Learned wavelet-frame shrinkage network – concluding remarks:* The (residual) learned wavelet-frame shrinkage network is a design which explicitly mimics wavelet-shrinkage algorithms. It can be observed that the (r)LWFSN inherently assume that noise is high-frequency and explicitly avoid non-linear processing in the low-frequency band. Follow-up experiments also included non-linearities in the low-frequency band of the learned wavelet-frame shrinkage network [36] and obtained similar results to the original design.

## VII. WHAT HAPPENS IN TRAINED MODELS?

### A. Properties of convolution kernels and low-rank approximation

The assumption that the convolution filters of a CNN behave as (complementary-phase) tight framelets is useful for analyzing the theoretical ability of a CNN to propagate signals. Albeit, it is difficult to prove that trained models comply with this assumption, because there are diverse elements affecting the optimization of the model, e.g. the initialization of the network, the data presented to the model, the optimization algorithm as well as its parameters. In addition, in real CNNs,

there may be co-adaptation between the diverse CNN layers, which may prevent that the individual filters of the CNN behave as tight framelets, since the decomposition and filtering performed by one layer is not independent from the rest [37].

To test the behavior if the filters of trained CNN can converge to complementary-phase tight framelets, at least, on a simplified environment, we propose to train a toy model, as displayed in Fig. 13. If the trained filters of an encoder-decoder pair of the toy model ($\mathbf{K}_l$, $\tilde{\mathbf{K}}_l$), (where $l$ denotes one of the decomposition levels) behave as a complementary-phase tight framelet, then the pair ($\mathbf{K}_l$, $\tilde{\mathbf{K}}_l$) approximately complies with the condition presented in Eq. (40), which for identity input $\mathbf{I}$ simplifies to

$$\tilde{\mathbf{K}}_n^{\mathsf{T}}(\mathbf{K}_n)_+ = \mathbf{I} \cdot c_n , \tag{95}$$

in which $c_n$ is an arbitrary constant.

The toy model is trained on images that contain multiple randomly-generated overlapping triangles. All the images were scaled to the range [0,1]. For this experiment, the input to the images is the noise-contaminated image and the objective/desired output is the noiseless image. For training the CNNs, normally-distributed noise with standard deviation of 0.1 was added to the ground truth. For every epoch, a batch of 192 training images is generated. As validation and test images we use the "astronaut" and "cameraman" images included in the software package *Scipy*. The model is optimized with Adam for 25 epochs with a linearly decreasing learning rate. The initial learning rate for the optimizer is set to $10^{-3}$ and the batch size is set to 1 sample. The convolution kernels were initialized with Xavier initialization using a uniform distribution (see Glorot and Bengio [38]). The code is available at IEEE's code ocean https://codeocean.com/capsule/7845737/tree.

With the described settings, we have trained the toy model, and have tested if the phase-complementary tight-framelet property holds for the filters of the deepest level $l$=2. The results for the operation $\tilde{\mathbf{K}}_2^{\mathsf{T}}(\mathbf{K}_2)_+$ are displayed in Fig. 12 (left), which shows that when the weights of the encoder and decoder have different initial values, the kernel pair ($\mathbf{K}_2$, $\tilde{\mathbf{K}}_2$) are not complementary-phase tight framelets. We have observed that the forward and inverse filters of wavelets/framelets are often the same or at least very similar. Based on this reasoning, we have initialized the toy model with the same initial values of the kernel pairs ($\mathbf{K}_n$, $\tilde{\mathbf{K}}_n$). As shown by Fig. 12 (right), with the proposed initialization, the filters of the CNN converge to tensors with properties reminiscent of complementary-phase tight-framelets. This suggests that the initialization of the CNN has an important influence on the convergence of the model to a specific solution.

Fig. 14 displays a test image processed with two toy models, one trained trained with different and one trained with the same initial values for the encoding-decoding pairs. It can be observed that there are no significant differences between the images produced by both models. In the same figure (lower row), we have set the bias of both networks to zero. In this case, it is expected that the networks reconstruct the noisy input, as confirmed by the figure, where both CNNs partly reconstruct the original noisy signal. This result suggests that the ReLU plus bias pairs operate akin to the low-rank approx-
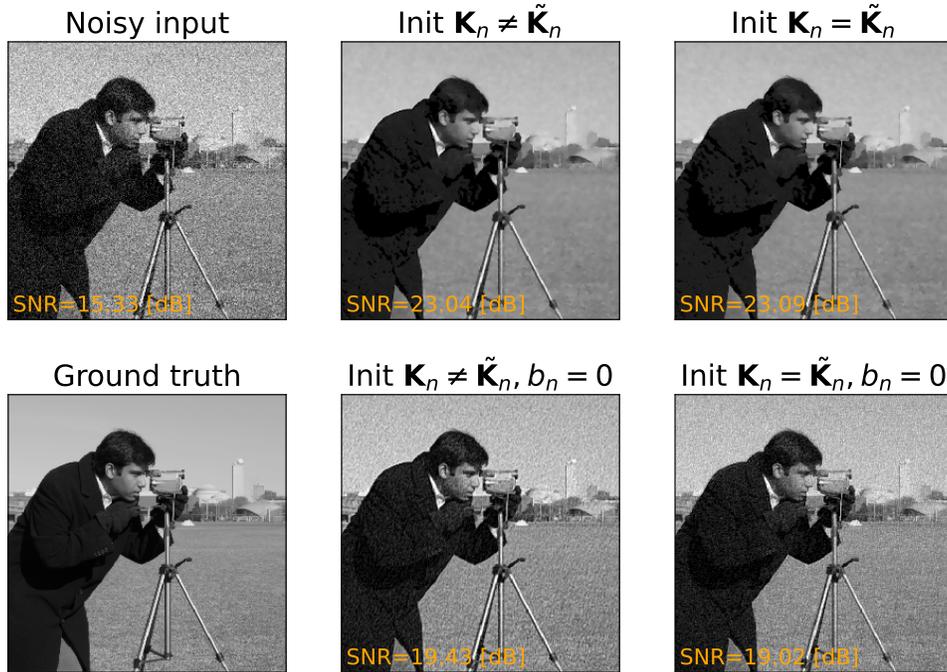
Fig. 14. Processed "cameraman" image for (in)dependently sampled initialization for the encoding and decoding filters. The top-left picture represents the noise-contaminated input ($\sigma_\eta = 0.1$) and the bottom-left, the noiseless reference. The middle-column images are the processed noisy image with the toy model trained with different initialization for its convolution filters, while the right-column images are processed with the model where the same initial values are used for the encoding and decoding filters. The top-middle and top-right images are nearly identical in terms of quality and SNR, so that initialization has no effect. The middle and bottom-right images are the same model presented that processed the middle and top-right figures, but where its bias is set to zero. As expected, the noise is partly reconstructed.

imation mechanism proposed the theory of deep convolutional framelets.

The following conclusions can be drawn from this experiment. First, the filters of the CNN may not necessarily converge to complementary-phase tight framelets. This is possibly caused by the initialization of the network and/or the interaction/co-adaptation between the multiple encoder/decoder layers. Second, we confirm that for our experimental setting, the low-rank approximation behavior in the CNN can be observed. For example, when setting the biases and thresholds to zero, part of the noise texture (high-rank information) is recovered. Third, it is possible that linear filtering happens in the network as well, which may explain why noise texture is not fully recovered when setting the biases to zero. Fourth and final, we have observed that the behavior of the trained models change drastically depending on factors such as the learning rate and the initialization values of the model. For this reason, we consider this experiment and its outcome more as a proof of concept, where further investigation is needed.

### B. Generalization

From the explanations in Section IV-C, it can be noted that the bias/threshold used in CNNs can modulate how much of the signal is suppressed by the nonlinearities. In addition, Section V-D established that there are additional mechanisms for noise reduction within the CNN, such as the Wiener-like behavior observed by Mohan *et al.* [14]. This raises the question how robust conventional CNNs are to noise-level changes different from the level that the model has been trained with. To perform such experiment, we have trained two variants of the toy model. The first variant is inspired by the multi-scale sparse coding network by Mentl *et al.* [17], where the biases of each of the nonlinearities (ReLU in this case) are multiplied by an estimate of the standard deviation of the noise. In the design of this example, the noise estimate $\hat{\sigma}_\eta$, which in accordance to Chang *et al.* [1] is defined by

$$\hat{\sigma}_\eta = 1.4826 \cdot \text{Median}(|\mathbf{f}_{\text{HH}} * \mathbf{x}|). \tag{96}$$

Here, variable $\mathbf{f}_{\text{HH}}$ is the diagonal convolution filter of the discrete wavelet transform with Haar basis. For comparison purposes, we will refer to this model as *adaptive* toy model. The second variant of the toy model being tested, examines the case where the convolution layers of the model do not add bias to the signal. This model is based in the so-called bias-free CNNs proposed by Mohan *et al.*, in which the bias of every convolution filter is set to zero during training. This setting has the purpose of achieving better generalization on the model, since it is claimed that this modification causes the model to behave independently of the noise level.

We have trained the described variants of the toy models with the same settings of the experiment in Section VII-A. The three models are evaluated on the test image with varying noise levels $\sigma_n \in [0.100, 0.150, 0.175, 0.200, 0.225]$, the result for this evaluation is displayed in Fig. 15. These results confirm that the performance of the original toy model degrades for higher noise levels. In contrast, the adaptive and bias-free toy
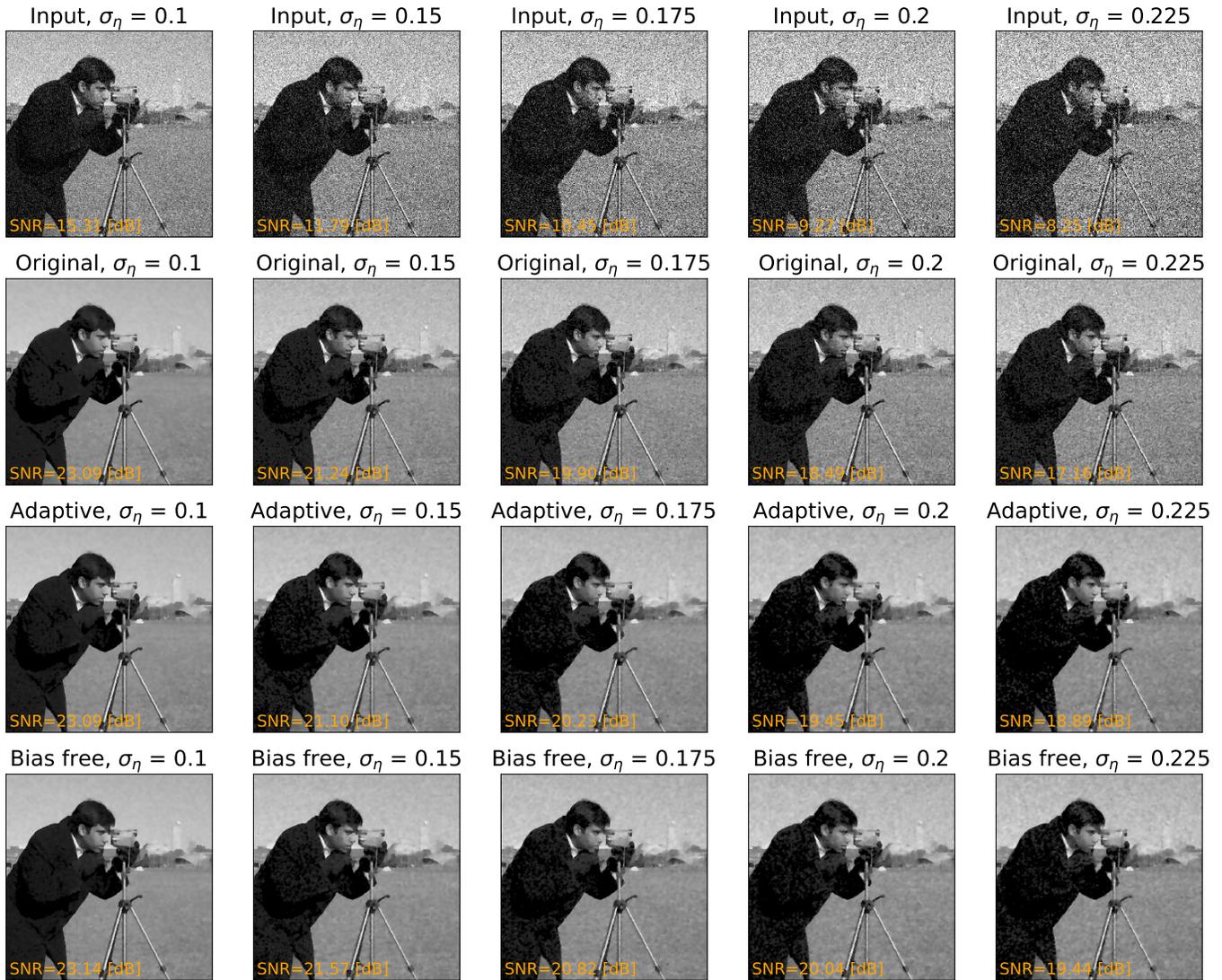
Fig. 15. Comparison of the baseline (original) toy model against its adaptive and bias-free variants. The models are evaluated in the cameraman picture with increasing noise levels. The top row displays the noisy input. The second top-row represents the images processed with the original toy model. Meanwhile, the third row are the results of the adaptive toy model. Finally, the bottom-row are the results corresponding to the bias-free model. It can be observed that the performance original toy model degrades as the noise level increases, while the performance adaptive and bias-free model degrade less with increased noise levels, resulting in pictures with lower noise levels.

model perform better than the original toy model for most noise levels.

The results of this experiment confirm the diverse noise-reduction mechanisms within a CNN, as well as showing that the CNNs have certain modeling limitations. For example, noise invariance, which can be addressed by further incorporating prior knowledge to the model, such as the case of the adaptive model, or by forcing the model to have a more Wiener-like behavior such as the case of the bias-free model. In the case of the bias-free model, note that, theoretically, it should be possible to obtain exactly the same behavior with the original toy model if the biases of the model would have converged to zero. This reasoning suggests that the large amount of free parameters and non-linear behavior of the model can potentially prevent to find the optimal/robust solution, in which case the incorporation of prior knowledge

can help to improve the model.

## VIII. WHICH NETWORK FITS TO MY PROBLEM?

### A. Design elements

When choosing or designing a CNN for a specific noise-reduction application, multiple choices and design elements should be considered. For example, the required performance, the memory required to train/deploy models, if certain signal preservation characteristics are required, the target execution time for the model, the characteristics of the images being processed, etc. Based on these requirements diverse design elements of CNNs can be more or less desirable, for example, the activation functions, the use of single/multi-resolution models, the need for skip connections, and so forth. This section briefly discusses such elements by focusing on the impact that such elements have in terms of performance

TABLE II
DESIGN ELEMENTS AND THEIR IMPACT IN PERFORMANCE AND COMPUTATION COST.

| Design elements | | Expressivity | Performance | No. parameters | Receptive field per layer |
|---|---|---|---|---|---|
| Activation | ReLU | High | Best | High | N/A |
| | Shrinkage | Low | Good | Medium | N/A |
| | Clipping | Low | Good | Medium | N/A |
| Scale | Single-scale | High | Good | High | Big |
| | Multi-scale | High | Good | Medium/high | Small |
| Topology | Non-residual | High | Good | Higher | N/A |
| | Residual | High | Best | Lower | N/A |

and potential computational cost. A summary of the main conclusions of these elements is included in Table II.

*1) Nonlinearity:* In literature, the most common activation function in CNNs is the ReLU. There are two main advantages of the ReLU with respect to other activations. First, ReLUs potentially enforce more sparsity in the feature maps than –for example– the soft shrinkage, because ReLUs cancel not only small values of the feature maps like the shrinkage functions do, but also all the negative values. The second advantage of the ReLU is its capacity to approximate other functions (see Section V-C). Note that the high-capacity of the ReLU to represent other functions [31], [13] (often referred to as *expressivity*) may also be one of the reasons why these models are prone to overfitting.

The better expressivity of the ReLU-CNNs may be the reason why –at the time of writing this manuscript– ReLU-based CNNs perform marginally better than the shrinkage-based models in terms of metrics such as signal-to-noise ratio or the structural similarity index metric [39], [40], [19]. Despite this small benefit, the visual characteristics of estimates produced by ReLU and shrinkage-based networks are very similar. Furthermore, the computational cost of ReLU-based designs is potentially higher those with shrinkage functions, because ReLUs require more feature maps to preserve the signal integrity. For example, the LWFSN shown in Section VI-C achieves a performance very close to the FBPConvNet and the tight-frame U-Net for noise reduction in computed tomography, but only with a small fraction of the total trainable parameters, which allows for a faster and less computation-expensive model [19].

As concluding remark, it can be noted that regardless of the expressivity of the ReLU activation, it is not entirely clear if this means that ReLU activations outperform other functions such as the soft threshold *in general*. Because we could not find articles, which are specifically focused on comparing the performance of ReLU/shrinkage-based models. In spite of this, there are some works that compare shrinkage-based CNNs with other (architecturally different) models based on ReLUs that indicate that the compared ReLU-based designs slightly outperform the shrinkage-based ones. For example, Herbeteau and Kevrann [40] proposed the so-called DCT2-Net, which is a shrinkage-based CNN, which despite of its good performance, it is still outperformed by the ReLU-based DnCNN [7] CNN. A similar behavior was observed by Zavala *et al.* [19], where their shrinkage-based LWFSN could not outperform the ReLU-based FBPConvNet [21] and

the tight-frame U-Net [6]. Another similar case is the deep K-SVD network [39], which achieves performance close (but slightly less good) than the ReLU-based DnCNN. Among the few examples were we found that a ReLU CNN performed better than shrinkage-based models, we find the work of Fan *et al.* [18], where they compared variants of the so-called soft autoencoder and found that the shrinkage-based model outperformed the ReLU variant.

*2) Single/multi-scale designs:* Single-scale models have the advantage that they avoid aliasing because no down/up-sampling layers are used. Still this comes at the cost of more computations and memory. Furthermore, this approach may lead to models with larger filters and/or deeper networks to achieve the same receptive-field than multi-scale models, which may further increase the computation costs of single-scale models.

In the case of multi-scale models, the main consideration should be that the down/up-sampling structure should allow perfect signal reconstruction to avoid introducing aliasing and/or distortion to the image estimates (e.g. the discrete wavelet transform in the tight-frame U-Net and in the learned wavelet-frame shrinkage network).

*3) (Non-) residual models:* Residual noise-reduction CNNs often perform better than their non-residual counterparts (e.g. the U-Net vs FBPConvNet and the LWFSN vs the rLWFSN). This may be caused because the trained models have more freedom to learn the filters, because the design does not need to learn to reconstruct the noiseless signal, but only to estimate the noise [12]. Also, it can be observed that non-residual models potentially need more parameters than residual networks, because the propagation/reconstruction of the noiseless signal is also dependent on the number of channels of the network.

### B. State-of-the art

Defining the state-of-the art in image denoising with CNNs is challenging for diverse reasons. First, there is a wide variety of available CNNs, which often are not compared to each other. Second, the suitability of a CNN for a given task may depend on the image and noise characteristics, such as noise distribution and (non-) stationarity. Third, the large amount of variables, in terms of e.g. optimization, data and data augmentation, adds reproducibility issues which further complicate making a fair comparison between all available models [11]. In addition, it should be noted that for many of

the existing models, the performance gap between state-of-the art models and other CNNs is often small.

Despite of the mentioned challenges, we have found some models that could be regarded as the state-of-the art. The first model to be addressed is the so-called DRU-Net [41], which is a bias-free model [14] that incorporates a U-Net architecture with residual blocks. In addition, the DRU-Net uses an additional input to indicate to the network the noise intensity, which increases its generalization to different noise levels. An additional state-of-the art model is DnCNN [7]. This network is residual, single-scale, while also using ReLU activations. Another state-of-the art model is the multi-level-wavelet CNN [42], which has a design very similar to the tight-frame U-Net [6]. Both of these models are based on the original U-Net design [34], but are deployed in residual configuration and the down/up-sampling structure is based on the discrete wavelet transform. Furthermore, in addition to use encoding-decoding CNNs stand-alone, have used CNNs as proximal operator within model-based methods [43], [41], which further improves the denoising power of non-model-based encoding-decoding CNNs.

## IX. CONCLUSIONS AND FUTURE OUTLOOK

In this paper, the widely used encoding-decoding CNN architecture has been analyzed from several signal processing principles. This analysis has revealed the following conclusions. (1) Multiple signal processing concepts converge in the mathematical formulation encoding-decoding CNNs models. For example, the convolution and down/up-sampling structure of the encoder-decoder structure is akin to the framelet decomposition, the activation functions are rooted on classical signal estimators. In addition, linear filtering also may happen within the model. (2) The activations implicitly assume noise and signal characteristics of the feature maps. (3) There are still many signal processing developments that can be integrated with current CNNs, further improving their performance in terms of accuracy, efficiency or robustness.

Despite of the signal processing nature of the encoding-decoding CNNs, at the moment of this publication, the integration of CNNs and existing signal processing algorithms is at an early stage. A clear example of the signal modeling limitations of current CNN denoisers are the activation functions, where the estimators provided by current activation layers neglect spatial correlation of the feature maps. Possible alternatives to solve this limitation could be to perform an activation function inspired on denoisers working on principles such as Markov random fields [44], the locally-spatial indicators [45] and multi-scale shrinkage [24]. Further ideas are provided by the extensive survey on denoising algorithms by Pižurika and Philips [46]. Additional approaches that can be further explored are non-local [47] and collaborative filtering [48]. Both techniques exploit the redundancy in natural images and only a few models are exploring these properties [49], [50].

To finalize this manuscript we encourage the reader to actively consider the properties of the signals processed, the design requirements and the existing signal processing algorithms when designing new CNNs. By doing so, we expect that next the generation of CNNs denoisers not only will be better performing, but also more interpretable, and reliable.

## REFERENCES

[1] S. G. Chang, B. Yu, and M. Vetterli, "Adaptive wavelet thresholding for image denoising and compression," *IEEE Trans. Image Process.*, vol. 9, no. 9, pp. 1532–1546, 2000.

[2] M. Elad and M. Aharon, "Image denoising via learned dictionaries and sparse representation," in *2006 IEEE Comput. Society Conference on Comput. Vision and Pattern Recognition (CVPR'06)*, vol. 1. IEEE, 2006, pp. 895–900.

[3] L. I. Rudin, S. Osher, and E. Fatemi, "Nonlinear total variation based noise removal algorithms," *Physica D: nonlinear phenomena*, vol. 60, no. 1-4, pp. 259–268, 1992.

[4] K. H. Jin and J. C. Ye, "Annihilating filter-based low-rank hankel matrix approach for image inpainting," *IEEE Trans. Image Process.*, vol. 24, no. 11, pp. 3498–3511, 2015.

[5] H. Chen, Y. Zhang, M. K. Kalra, F. Lin, Y. Chen, P. Liao, J. Zhou, and G. Wang, "Low-dose ct with a residual encoder-decoder convolutional neural network," *IEEE Trans. Med. Imag.*, vol. 36, no. 12, pp. 2524–2535, 2017.

[6] Y. Han and J. C. Ye, "Framing u-net via deep convolutional framelets: Application to sparse-view ct," *IEEE Trans. Med. Imag.*, vol. 37, no. 6, pp. 1418–1429, 2018.

[7] K. Zhang, W. Zuo, Y. Chen, D. Meng, and L. Zhang, "Beyond a gaussian denoiser: Residual learning of deep cnn for image denoising," *IEEE Trans. Image Process.*, vol. 26, no. 7, pp. 3142–3155, 2017.

[8] T. Yokota, H. Hontani, Q. Zhao, and A. Cichocki, "Manifold modeling in embedded space: an interpretable alternative to deep image prior," *IEEE Transactions on Neural Networks and Learning Systems*, 2020.

[9] K. C. Kusters, L. A. Zavala-Mondragón, J. O. Bescós, P. Rongen, P. H. de With, and F. van der Sommen, "Conditional generative adversarial networks for low-dose ct image denoising aiming at preservation of critical image content," in *2021 43rd Annual International Conference of the IEEE Engineering in Medicine & Biology Society (EMBC)*. IEEE, 2021, pp. 2682–2687.

[10] H. Gupta, K. H. Jin, H. Q. Nguyen, M. T. McCann, and M. Unser, "Cnn-based projected gradient descent for consistent ct image reconstruction," *IEEE Trans. Med. Imag.*, vol. 37, no. 6, pp. 1440–1453, 2018.

[11] M. T. McCann, K. H. Jin, and M. Unser, "Convolutional neural networks for inverse problems in imaging: A review," *IEEE Signal Processing Magazine*, vol. 34, no. 6, pp. 85–95, 2017.

[12] J. C. Ye, Y. Han, and E. Cha, "Deep convolutional framelets: A general deep learning framework for inverse problems," *SIAM Journal on Imaging Sciences*, vol. 11, no. 2, pp. 991–1048, 2018.

[13] J. C. Ye and W. K. Sung, "Understanding geometry of encoder-decoder cnns," in *International Conference on Machine Learning*. PMLR, 2019, pp. 7064–7073.

[14] S. Mohan, Z. Kadkhodaie, E. P. Simoncelli, and C. Fernandez-Granda, "Robust and interpretable blind image denoising via bias-free convolutional neural networks," in *International Conference on Learning Representations*, 2020. [Online]. Available: https://openreview.net/forum?id=HJlSmC4FPS

[15] M. Unser, "From kernel methods to neural networks: A unifying variational formulation," *arXiv preprint arXiv:2206.14625*, 2022.

[16] V. Papyan, Y. Romano, and M. Elad, "Convolutional neural networks analyzed via convolutional sparse coding," *The Journal of Machine Learning Research*, vol. 18, no. 1, pp. 2887–2938, 2017.

[17] K. Mentl, B. Mailhé, F. C. Ghesu, F. Schebesch, T. Haderlein, A. Maier, and M. S. Nadar, "Noise reduction in low-dose ct using a 3d multiscale sparse denoising autoencoder," in *2017 IEEE 27th International Workshop on Machine Learning for Signal Processing (MLSP)*. IEEE, 2017, pp. 1–6.

[18] F. Fan, M. Li, Y. Teng, and G. Wang, "Soft autoencoder and its wavelet adaptation interpretation," *IEEE Transactions on Computational Imaging*, vol. 6, pp. 1245–1257, 2020.

[19] L. A. Zavala-Mondragón, P. Rongen, J. O. Bescos, P. H. De With, and F. Van der Sommen, "Noise reduction in ct using learned wavelet-frame shrinkage networks," *IEEE Trans. Med. Imag.*, 2022.

[20] L. A. Zavala-Mondragón, P. H. de With, and F. van der Sommen, "Image noise reduction based on a fixed wavelet frame and cnns applied to ct," *IEEE Trans. Image Process.*, vol. 30, pp. 9386–9401, 2021.

[21] K. H. Jin, M. T. McCann, E. Froustey, and M. Unser, "Deep convolutional neural network for inverse problems in imaging," *IEEE Trans. Image Process.*, vol. 26, no. 9, pp. 4509–4522, 2017.

[22] M. Ranzato, F. J. Huang, Y.-L. Boureau, and Y. LeCun, "Unsupervised learning of invariant feature hierarchies with applications to object recognition," in *2007 IEEE conference on computer vision and pattern recognition*. IEEE, 2007, pp. 1–8.

[23] E. J. Candes and M. B. Wakin, "An introduction to compressive sampling," *IEEE Signal Processing Magazine*, vol. 25, no. 2, pp. 21–30, 2008.

[24] L. Sendur and I. W. Selesnick, "Bivariate shrinkage functions for wavelet-based denoising exploiting interscale dependency," *IEEE Trans. Signal Process.*, vol. 50, no. 11, pp. 2744–2756, 2002.

[25] G. Steidl and J. Weickert, "Relations between soft wavelet shrinkage and total variation denoising," in *Joint pattern recognition symposium*. Springer, 2002, pp. 198–205.

[26] T. Blu and F. Luisier, "The sure-let approach to image denoising," *IEEE Trans. Image Process.*, vol. 16, no. 11, pp. 2778–2786, 2007.

[27] I. Daubechies, M. Defrise, and C. De Mol, "An iterative thresholding algorithm for linear inverse problems with a sparsity constraint," *Communications on Pure and Applied Mathematics: A Journal Issued by the Courant Institute of Mathematical Sciences*, vol. 57, no. 11, pp. 1413–1457, 2004.

[28] K. Gregor and Y. LeCun, "Learning fast approximations of sparse coding," in *Proceedings of the 27th international conference on international conference on machine learning*, 2010, pp. 399–406.

[29] I. Y. Chun, Z. Huang, H. Lim, and J. Fessler, "Momentum-net: Fast and convergent iterative neural network for inverse problems," *IEEE Trans. Pattern Anal. Mach. Intell.*, pp. 1–1, 2020.

[30] J. Xiang, Y. Dong, and Y. Yang, "Fista-net: Learning a fast iterative shrinkage thresholding network for inverse problems in imaging," *IEEE Transactions on Medical Imaging*, vol. 40, no. 5, pp. 1329–1339, 2021.

[31] I. Daubechies, R. DeVore, S. Foucart, B. Hanin, and G. Petrova, "Nonlinear approximation and (deep) relu networks," *Constructive Approximation*, vol. 55, no. 1, pp. 127–172, 2022.

[32] J. C. Ye, J. M. Kim, K. H. Jin, and K. Lee, "Compressive sampling using annihilating filter-based low-rank interpolation," *IEEE Transactions on Information Theory*, vol. 63, no. 2, pp. 777–801, 2016.

[33] A. Cichocki, R. Zdunek, and S.-i. Amari, "Nonnegative matrix and tensor factorization [lecture notes]," *IEEE Signal Processing Magazine*, vol. 25, no. 1, pp. 142–145, 2008.

[34] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *Medical Image Computing and Computer-Assisted Intervention – MICCAI 2015*, N. Navab, J. Hornegger, W. M. Wells, and A. F. Frangi, Eds. Cham: Springer International Publishing, 2015, pp. 234–241.

[35] T. Kwon and J. C. Ye, "Cycle-free cyclegan using invertible generator for unsupervised low-dose ct denoising," *IEEE Transactions on Computational Imaging*, vol. 7, pp. 1354–1368, 2021.

[36] L. A. Zavala-Mondragón, F. van der Sommen *et al.*, "On the performance of learned and fixed-framelet shrinkage networks for low-dose ct denoising," in *Medical Imaging with Deep Learning*, 2022.

[37] G. E. Hinton, N. Srivastava, A. Krizhevsky, I. Sutskever, and R. R. Salakhutdinov, "Improving neural networks by preventing co-adaptation of feature detectors," *arXiv preprint arXiv:1207.0580*, 2012.

[38] X. Glorot and Y. Bengio, "Understanding the difficulty of training deep feedforward neural networks," in *Proceedings of the thirteenth international conference on artificial intelligence and statistics*. JMLR Workshop and Conference Proceedings, 2010, pp. 249–256.

[39] M. Scetbon, M. Elad, and P. Milanfar, "Deep k-svd denoising," *IEEE Trans. Image Process.*, vol. 30, pp. 5944–5955, 2021.

[40] S. Herbreteau and C. Kervrann, "Dct2net: an interpretable shallow cnn for image denoising," *IEEE Trans. Image Process.*, vol. 31, pp. 4292–4305, 2022.

[41] K. Zhang, Y. Li, W. Zuo, L. Zhang, L. Van Gool, and R. Timofte, "Plug-and-play image restoration with deep denoiser prior," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 44, no. 10, pp. 6360–6376, 2021.

[42] P. Liu, H. Zhang, K. Zhang, L. Lin, and W. Zuo, "Multi-level wavelet-cnn for image restoration," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2018, pp. 773–782.

[43] V. Monga, Y. Li, and Y. C. Eldar, "Algorithm unrolling: Interpretable, efficient deep learning for signal and image processing," *IEEE Signal Processing Magazine*, vol. 38, no. 2, pp. 18–44, 2021.

[44] M. Malfait and D. Roose, "Wavelet-based image denoising using a markov random field a priori model," *IEEE Trans. Image Process.*, vol. 6, no. 4, pp. 549–565, 1997.

[45] A. Pižurica and W. Philips, "Estimating the probability of the presence of a signal of interest in multiresolution single- and multiband image denoising," *IEEE Trans. Image Process.*, vol. 15, no. 3, pp. 654–665, 2006.

[46] A. Pižurica, "Image denoising algorithms: From wavelet shrinkage to nonlocal collaborative filtering," *Wiley Encyclopedia of Electrical and Electronics Engineering*, pp. 1–17, 1999.

[47] K. Dabov, A. Foi, V. Katkovnik, and K. Egiazarian, "Image restoration by sparse 3d transform-domain collaborative filtering," in *Image Processing: Algorithms and Systems VI*, vol. 6812. International Society for Optics and Photonics, 2008, p. 681207.

[48] A. Buades, B. Coll, and J.-M. Morel, "A non-local algorithm for image denoising," in *2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'05)*, vol. 2. IEEE, 2005, pp. 60–65.

[49] D. Yang and J. Sun, "Bm3d-net: A convolutional neural network for transform-domain collaborative filtering," *IEEE Signal Processing Letters*, vol. 25, no. 1, pp. 55–59, 2018.

[50] H. Lee, H. Choi, K. Sohn, and D. Min, "Knn local attention for image restoration," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 2139–2149.

**Luis Albert Zavala-Mondragón** (lzavala905@gmail.com) is a PhD candidate at the Eindhoven University of Technology (TU/e). He holds a BSc and MSc in electrical engineering by the Universidad Nacional Autónoma de México (UNAM) and the Eindhoven University of Technology, respectively. Furthermore, he has experience in the field of hardware emulation (Intel, Mexico) and computer vision (Thirona, Netherlands). His current research interests are the development of efficient and explainable computer vision pipelines for healthcare applications. He is a student member of IEEE.

**Peter H. N. de With** (p.h.n.de.with@tue.nl) holds a MSc and PhD degree in electrical engineering from the Eindhoven University of Technology (TU/e) the Delft University of Technology, respectively. In 1984, he joined Philips Research Labs Eindhoven. In 1997-2000, he was a Full professor at University of Mannheim. From 2008 to 2010, he was VP of video technology at CycloMedia Technology.

In 2000 he became part-time professor at TU/e and since 2011, he is a full professor. He leads the Video Coding and Architectures (VCA) group. Peter has coauthored more than 70 refereed international book chapters and journal articles, over 500 conference publications and 40 international patents. He served as Technical Committee Member of the IEEE CES, ICIP, SPIE and is and co-recipient of multiple paper awards like the IEEE CES Paper Awards (several), VCIP and ICCE Best Paper Awards. He is member of the Royal Holland Society of Sciences and Humanities. He is a Fellow of IEEE.

**Fons van der Sommen** (fvdsommen@tue.nl) is an associate professor at Eindhoven University of Technology, holding a BSc and MSc in electrical engineering and a PhD in computer vision. Heading the healthcare cluster of the VCA research group, Fons has worked on a variety of image processing and computer vision applications, mainly in the medical domain. He has a strong interest in signal processing and information theory and strives to exploit methods from these fields to improve the robustness, efficiency and interpretability of modern-day AI architectures, such as Convolutional Neural Networks. He is a member of IEEE.