# Deep Convolutional Framelet Denosing for Low-Dose CT via Wavelet Residual Network

Eunhee Kang, Won Chang, Jaejun Yoo, and Jong Chul Ye*, *Senior Member, IEEE*

*Abstract*—Model based iterative reconstruction (MBIR) algorithms for low-dose X-ray CT are computationally expensive. To address this problem, we recently proposed a deep convolutional neural network (CNN) for low-dose X-ray CT and won the second place in 2016 AAPM Low-Dose CT Grand Challenge. However, some of the texture were not fully recovered. To address this problem, here we propose a novel framelet-based denoising algorithm using wavelet residual network which synergistically combines the expressive power of deep learning and the performance guarantee from the framelet-based denoising algorithms. The new algorithms were inspired by the recent interpretation of the deep convolutional neural network (CNN) as a cascaded convolution framelet signal representation. Extensive experimental results confirm that the proposed networks have significantly improved performance and preserves the detail texture of the original images.

*Index Terms*—Deep learning, low-dose CT, framelet denoising, convolutional neural network (CNN), convolution framelets

## I. INTRODUCTION

**X**-RAY computed tomography (CT) is one of the most valuable imaging techniques in clinics. It is used in various ways, including whole-body diagnostic CT, C-arm CT for interventional imaging, dental CT, etc. However, X-ray CT causes potential cancer risks due to radiation exposure. To ensure patient safety, X-ray dose reduction techniques have been extensively studied, and the reduction in the number of X-ray photons using tube current modulation is considered one of the solutions. A drawback of this approach is, however, the low signal-to-noise ratio (SNR) of projections, which induces noise in the reconstructed image. Various model based iterative reconstruction (MBIR) methods [1], [2], [3] have been investigated to obtain a clear reconstructed image. However, these approaches are usually computationally expensive due to the iterative applications of forward and backward projections.

Recently, deep learning approaches have been actively explored for various computer vision applications through the use of extensive data and powerful graphical processing units (GPUs). Deep networks have achieved great successes in computer vision applications such as classification [4], denoising [5], [6], [7], segmentation [8], and super-resolution [9], etc.

In MR image reconstruction, Wang et al [10] was the first to apply deep learning to compressed sensing MRI (CS-MRI). Deep network architecture using unfolded iterative compressed sensing (CS) algorithm was also proposed [11], [12]. In CT restoration problems, our group introduced the deep learning approach for low-dose X-ray CT [13], whose performance has been rigorously confirmed by winning the second place award in 2016 AAPM Low-Dose CT Grand Challenge. Since then, several pioneering deep learning approaches for low-dose CT have been proposed by many researchers [14], [15], [16], [17], [18], [19], [20], [21], [22], [23]. Some algorithms uses generative adversarial network (GAN) loss [21], [22]. Recent proposal is to incorporate deep neural network within iterative steps [16], [24]. However, existing algorithms consider a deep network as a black-box, so it is difficult to understand the role of deep networks within iterative steps.

Therefore, one of the main contributions of this paper is to show that a feed-forward deep learning-based denoising is indeed the first iteration of a special instance of frame-based denoising algorithm using deep convolutional framelets [25]. Frame-based denoising approaches using wavelet frames have been an extensive research topics in applied mathematics community due to its proven convergence [26], [27]. On the other hand, the theory of deep convolutional framelet [25] was recently proposed to explain the mathematical origin of deep neural network as a multi-layer realization of the convolution framelets [28]. Accordingly, the main goal of this paper is to synergistically combine the expressive power of deep neural network and the performance guarantee from the framelet-based denoising algorithms. In particular, we show that the performance of the deep learning-based denoising algorithm can be improved with iterative steps similar to the classical framelet-based denoising approaches [29], [30]. Furthermore, we can provide the theoretical guarantee of the algorithm to converge.

Compared to the recent proposals of learning-based optimization approaches [16], [24], one of the important advantages of our work is that our deep network is no more a black box but can be optimized for specific restoration tasks by choosing optimal framelet representation. Thus, we can employ an improved wavelet residual network (WavResNet) structure [31] in our deep convolutional framelet denoising thanks to its effectiveness in recovering the directional components. We confirm our theoretical reasoning using extensive numerical experiments.

## II. Theory

For simplicity, we derive our theory for 1-D signals, but the extension to 2-D image is straightforward.

### A. Frame-based Denoising

Consider an analysis operator $\mathcal{W}$ given by $\mathcal{W}^\top = \begin{bmatrix} w_1 & \cdots & w_m \end{bmatrix}$, where the superscript $^\top$ denotes the Hermitian transpose and $\{w_k\}_{k=1}^m$ is a family of function in a Hilbert space $H$. Then, $\{w_k\}_{k=1}^m$ is called a frame if it satisfies the following inequality [32]:

$$\alpha\|f\|^2 \leq \|\mathcal{W}f\|^2 \leq \beta\|f\|^2, \quad \forall f \in H, \tag{1}$$

where $\alpha, \beta > 0$ are called the frame bounds. Then the recovery of the original signal can be done from the frame coefficient $c = \mathcal{W}f$ using the dual frame $\tilde{\mathcal{W}}$ satisfying the frame condition: $\tilde{\mathcal{W}}^\top \mathcal{W} = I$, since $f = \tilde{\mathcal{W}}^\top c = \tilde{\mathcal{W}}^\top \mathcal{W}f = f$. This condition is often called the perfect reconstruction (PR) condition. We often call $\tilde{\mathcal{W}}^\top$ as the synthesis operator. The frame is said to be tight, if $\alpha = \beta$ in (1). This is equivalent to $\tilde{\mathcal{W}} = \mathcal{W}$ or $\mathcal{W}^\top \mathcal{W} = I$.

Suppose that noisy measurement $g \in \mathbb{R}^n$ is given by

$$g = f^* + e$$

where $f^* \in \mathbb{R}^n$ is a unknown ground-truth image and $e \in \mathbb{R}^n$ denotes the noise. Then, the classical tight frame-based denosing approaches [26], [27] solve the following alternating minimization problem:

$$\min_{f,\alpha} \frac{\mu}{2}\|g - f\|^2 + \frac{1-\mu}{2}\left\{\|\mathcal{W}f - \alpha\|^2 + \lambda\|\alpha\|_1\right\} \tag{2}$$

where $\lambda, \mu > 0$ denote the regularization parameters. The corresponding proximal update equation is then given by [26], [27]:

$$f_{n+1} = \mu g + (1-\mu)\mathcal{W}^\top T_\lambda(\mathcal{W}f_n), \tag{3}$$

where $T_\lambda(\cdot)$ denotes the soft-thresholding operator with the threshold value of $\lambda$, and $f_n$ refers to the $n$-th update. Thus, the frame-based denoising algorithm in (3) is designed to remove insignificant parts of frame coefficients through shrinkage operation, by assuming that most of the meaningful signal has large frame coefficients and noises are distributed across all frame coefficients.

One of the most important advantages of this framelet-based denoising is its proven convergence [26], [27]. Our goal is thus to exploit the proven convergence of these approaches for our CNN based low-dose CT denoising. Toward this goal, in the next section we show that CNN is closely related to the frame bases.

### B. Deep Convolutional Framelets

Here, the theory of the deep convolutional framelet [25] is briefly reviewed to make this work self-contained. To avoid special treatment of boundary condition, our theory is mainly derived using circular convolution. Specifically, let $f = [f[1], \cdots, f[n]]^T \in \mathbb{R}^n$ be an input signal and $\overline{\psi} = [\psi[d], \cdots, \psi[1]]^T \in \mathbb{R}^d$ denotes filter represented as the flipped version of vector $\psi$. Then, the convolution operation in CNN can be represented using Hankel matrix operation [25]. Specifically, a single-input single-output (SISO) convolution with the filter $\overline{\psi}$ is given by a matrix vector multiplication:

$$y = f \circledast \overline{\psi} = \mathbb{H}_d(f)\psi, \tag{4}$$

where $\mathbb{H}_d(f)$ is a wrap-around Hankel matrix

$$\mathbb{H}_d(f) = \begin{bmatrix} f[1] & f[2] & \cdots & f[d] \\ f[2] & f[3] & \cdots & f[d+1] \\ \vdots & \vdots & \ddots & \vdots \\ f[n] & f[1] & \cdots & f[d-1] \end{bmatrix}.$$

Similarly, multi-input multi-output (MIMO) convolution with the matrix input $F := [f_1 \cdots f_p]$ and the multi-channel filter matrix $\overline{\Psi}$

$$\overline{\Psi} := \begin{bmatrix} \overline{\psi}_1^1 & \cdots & \overline{\psi}_q^1 \\ \vdots & \ddots & \vdots \\ \overline{\psi}_1^p & \cdots & \overline{\psi}_q^p \end{bmatrix} \in \mathbb{R}^{dp \times q} \tag{5}$$

can be represented as

$$Y = F \circledast \overline{\Psi} = \mathbb{H}_{d|p}(F)\Psi \tag{6}$$

where $\mathbb{H}_{d|p}(F)$ is an *extended Hankel matrix* by stacking $p$ Hankel matrices side by side:

$$\mathbb{H}_{d|p}(F) := \begin{bmatrix} \mathbb{H}_d(f_1) & \mathbb{H}_d(f_2) & \cdots & \mathbb{H}_d(f_p) \end{bmatrix} \tag{7}$$

In (5), $\overline{\psi}_i^j \in \mathbb{R}^d, i = 1, \cdots, q; j = 1, \cdots, p$ refer to the $j$-th input channel filters to generate the $i$-th output channel. Note that the convolutional representation using an extended Hankel matrix in (6) is equivalent to the multi-channel filtering operations commonly used in CNN [25].

Let $\Phi = [\phi_1, \cdots, \phi_n]$ and $\tilde{\Phi} = [\tilde{\phi}_1, \cdots, \tilde{\phi}_n] \in \mathbb{R}^{n \times n}$ (resp. $\Psi = [\psi_1, \cdots, \psi_q]$ and $\tilde{\Psi} = [\tilde{\psi}_1, \cdots, \tilde{\psi}_q] \in \mathbb{R}^{d \times q}$) are frames and its duals satisfying the frame condition:

$$\tilde{\Phi}\Phi^\top = I_{n \times n} \quad, \quad \Psi\tilde{\Psi}^\top = I_{d \times d}. \tag{8}$$

Accordingly, we can obtain the following matrix identity:

$$\mathbb{H}_d(f) = \tilde{\Phi}\Phi^\top \mathbb{H}_d(f)\Psi\tilde{\Psi}^\top = \tilde{\Phi}C\tilde{\Psi}^\top \tag{9}$$

where $C := \Phi^\top \mathbb{H}_d(f)\Psi$ denotes the framelet coefficient. This results in the following encoder-decoder layer structure [25]:

$$f = \left(\tilde{\Phi}C\right) \circledast \nu(\tilde{\Psi}), \tag{10}$$

$$C := \Phi^\top\left(f \circledast \overline{\Psi}\right) \tag{11}$$

where $\overline{\Psi}$ is from (5) by setting $q = 1$, and

$$\nu(\tilde{\Psi}) := \frac{1}{d}\begin{bmatrix} \tilde{\psi}_1 \\ \vdots \\ \tilde{\psi}_r \end{bmatrix}. \tag{12}$$

Similarly, for a given matrix input $Z \in \mathbb{R}^{n \times p}$, we can also derive the paired encoder-decoder structure [25]:

$$C = \Phi^\top\left(Z \circledast \overline{\Psi}\right) \tag{13}$$

$$Z = (\Phi C) \circledast \nu(\tilde{\Psi}) \tag{14}$$

where the encoder filter is given by (5) and the decoder filters is defined by

$$\nu(\tilde{\Psi}) \quad := \quad \frac{1}{d} \begin{bmatrix} \tilde{\psi}_1^1 & \cdots & \tilde{\psi}_1^p \\ \vdots & \ddots & \vdots \\ \tilde{\psi}_q^1 & \cdots & \tilde{\psi}_q^p \end{bmatrix} \in \mathbb{R}^{dq \times p} \tag{15}$$

such that they satisfy the frame condition

$$\Psi \tilde{\Psi}^\top = I_{dp \times dp} \quad . \tag{16}$$

The simple convolutional framelet expansion using (11), (10), (13) and (14) is so powerful that the deep CNN architecture emerges from them. Specifically, by inserting the pair (13) and (14) between the pair (11) and (10), we can derive a deep network structure. For more detail, see [25].

### C. Deep Convolutional Framelet Denoising

Now, note that the computation of our deep convolutional framelet coefficients can be represented by analysis operator:

$$\mathcal{W}f := C = \Phi^\top (f \circledast \overline{\Psi})$$

whereas the synthesis operator is given by the decoder part of convolution:

$$\tilde{\mathcal{W}}^\top C := (\Phi C) \circledast \nu(\tilde{\Psi}).$$

If the frame conditions (8) or (16) are met at each layer, we can therefore use the classical update algorithm in (3) for denosing. Then, what is the shrinkage operator that corresponds to $T_\lambda(\cdot)$ in (3)? One of the unique aspects of deep convolutional framelets is that by changing the number of filter channels, we can achieve the shrinkage behaviour [25]. More specifically, low-rank shrinkage behaviour emerges when the number of output filter channels are not sufficient. Therefore, the explicit application of the shrinkage operator is no more necessary.

To understand this claim, consider the following regression problem under low-rank Hankel structured matrix constraint:

$$\min_{f \in \mathbb{R}^n} \quad \|f^* - f\|^2$$
$$\text{subject to} \quad \text{RANK}\mathbb{H}_d(f) \le r < d. \tag{17}$$

where $f^* \in \mathbb{R}^n$ denotes the ground-truth signal, $r$ is the upper bound of the rank, and $\mathbb{H}_d(f) \in \mathbb{R}^{n \times d}$. The low-rank Hankel structured matrix constraint in (17) is known for its excellent performance in image denoising [33], artifact removal [34] and deconvolution [35].

A classical approach to address (17) is using the explicit singular value shrinkage operation to impose the low-rankness [36], [37]. However, using deep convolutional framelets, we do not need such explicit shrinkage operation. More specifically, let $V \in \mathbb{R}^{d \times r}$ denote the basis for $R\left((\mathbb{H}_d(f))^\top\right)$ where $R(\cdot)$ denote the range space. Then, there always exist two matrices pairs $\Phi, \tilde{\Phi} \in \mathbb{R}^{n \times n}$ and $\Psi, \tilde{\Psi} \in \mathbb{R}^{d \times r}$ satisfying the conditions

$$\tilde{\Phi}\Phi^\top = I_{n \times n}, \qquad \Psi\tilde{\Psi}^\top = P_{R(V)} \tag{18}$$

where $R(V)$ denote the range space of $V$ and $P_{R(V)}$ represents a projection onto $R(V)$. Note that the bases matrix $\tilde{\Psi} \in \mathbb{R}^{d \times r}$ in (18) does not satisfy the frame condition (8) due to the insufficient number channels, i.e. $r < d$. However, we still

have the following matrix equality that is essential for deep convolutional framelet expansion [25]:

$$\mathbb{H}_d(f) = \tilde{\Phi}\Phi^\top \mathbb{H}_d(f)\Psi\tilde{\Psi}^\top.$$

Accordingly, we can define a space $\mathcal{H}_r$ by collecting signals that can be decomposed to the single layer deep convolutional framelet expansion:

$$\mathcal{H}_r = \left\{ f \in \mathbb{R}^n \mid f = \left(\tilde{\Phi}C\right) \circledast \nu(\tilde{\Psi}), C = \Phi^\top \left(f \circledast \overline{\Psi}\right) \right\}$$

Then, the regression problem in (17) can be equivalently represented by

$$\min_{f \in \mathcal{H}_r} \|f^* - f\|^2 \,, \tag{19}$$

which implies that the explicit rank condition is embedded as a single layer convolutional framelets.

However, (19) holds for any signals that can be represented by arbitary $(\Phi, \tilde{\Phi})$ and $(\Psi, \tilde{\Psi})$ satisfying (18), and we should find ones that are optimized for given data. In our deep convolutional framelets, $\Phi$ and $\tilde{\Phi}$ correspond to the generalized pooling and unpooling which are chosen based on the application-specific knowledges [25], so we are interested in only estimating the filters $\Psi$, $\tilde{\Psi}$. Then, the main goal of the neural network training is to learn $(\Psi, \tilde{\Psi})$ from training data $\{(f_{(i)}, f_{(i)}^*)\}_{i=1}^N$ assuming that $\{f_{(i)}^*\}$ are associated with rank-$r$ Hankel matrices. Thus, (19) can be modified for the training data as follows:

$$\min_{\{f_{(i)}\} \in \mathcal{H}_r} \sum_{i=1}^N \|f_{(i)}^* - f_{(i)}\|^2 \tag{20}$$

which can be converted to the neural network training problem:

$$\min_{(\Psi, \tilde{\Psi})} \sum_{i=1}^N \left\| f_{(i)}^* - \mathcal{Q}(f_{(i)}; \Psi, \tilde{\Psi}) \right\|^2 \tag{21}$$

where

$$\mathcal{Q}(f_{(i)}; \Psi, \tilde{\Psi}) = \left( \tilde{\Phi}C[f_{(i)}] \right) \circledast \nu(\tilde{\Psi}) \tag{22}$$
$$C[f_{(i)}] = \Phi^\top \left( f_{(i)} \circledast \overline{\Psi} \right). \tag{23}$$

The idea can be further extended to the multi-layer deep convolutional framelet expansion with nonlinearity. Then, (21) is equivalently to the standard neural network training. Once the network is fully trained, the inference for a given noisy input $f$ is simply done by $\mathcal{Q}(f; \Psi, \tilde{\Psi})$, which is equivalent to find a denoised solution. Therefore, using deep convolutional framelets with insufficient channels, we do not need an explicit shrinkage operation and the update equation (3) can be replaced by

$$f_{n+1} = \mu g + (1 - \mu)\mathcal{Q}(f_n; \Psi, \tilde{\Psi}) \,, \tag{24}$$

where $\mathcal{Q}(f_n)$ is the deep convolutional framelet output.

However, one of the main differences of (24) from (3) is that our deep convolutional framelet does not satisfy the tight frame condition that is required to guarantees the convergence of (3). Therefore, to guarantee the convergence, we need to relax the iteration using Krasnoselskii-Mann (KM) method [38] as

described in Algorithm 1. Then, using the standard tools of proximal optimization [38], we can show that the sequence generated by Algorithm 1 converges to a fixed point.

**Theorem 2.1.** *There exists a parameter $\mu \in (0, 1)$ such that the deep convolutional framelet denoising algorithm in Algorithm 1 converges to a fixed point.*

*Proof.* See Appendix B. □

---

**Algorithm 1** Pseudocode implementation.
1: Train a deep network $\mathcal{Q}$ using training data set.
2: Set $0 \leq \mu \leq 1$ and $0 < \lambda_n < 1, \forall n$.
3: Set initial guess of $f_0$ and $f_1$.
4: **for** $n = 1, 2, \ldots$, until convergence **do**
5:     $q_n := \mathcal{Q}(f_n)$
6:     $\bar{f}_{n+1} := \mu g + (1 - \mu) q_n$
7:     $f_{n+1} := f_n + \lambda_n(\bar{f}_{n+1} - f_n)$
8: **end for**

---

Algorithm 1 corresponds to a recursive neural network (RNN) as shown in Fig. 1, which is related to an iterative network in [23]. If we use $\mu = 0, \lambda_n = 1$, the first iteration of Algorithm 1 corresponds to a feed-forward deep convolutional framelet denosing algorithm, which is also important by itself. In Experimental Results, we show the improvement using RNN. However, our feed-forward network is much faster with compatible image quality. Thus, we believe that both algorithms are useful in practice. Both algorithms have the same neural network backbone $\mathcal{Q}(\cdot)$, which will be described in detail in the following section.
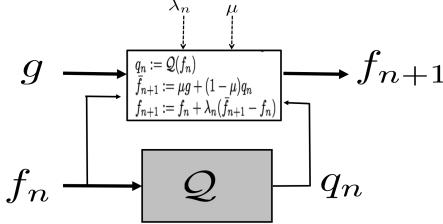


Fig. 1. Proposed RNN structure for deep convolutional framelet denoising.

### D. Optimizing the Network Architecture

In order to have the best denoising performance in frame-based denosing, the frame bases should have good energy compaction properties. For example, due to the vanishing moments of wavelets, wavelet transforms can annihilate the smoothly varying signals while maintaining the image edges, thus resulting in good energy compaction. Thus, wavelet frames such as contourlets [39] are often used for denoising. Furthermore, low-dose X-ray CT images exhibit streaking noise, so the contourlet transform [39] is good for detecting the streaking noise patterns by representing the directional edge information of X-ray CT images better. Thus, we are interested in using WavResNet [31] that employs the contourlet transform [39]. The proposed WavResNet architecture is illustrated in Fig. 2. WavResNet has three unique components: contourlet

transform, concatenation, and skipped connection. WavResNet is an extension of our prior work [13] that has similar network architecture except that residuals at each subband are estimated by the neural network [31]. In this paper, we provide a new interpretation of WavResNet using the theory of deep convolutional framelets [25].
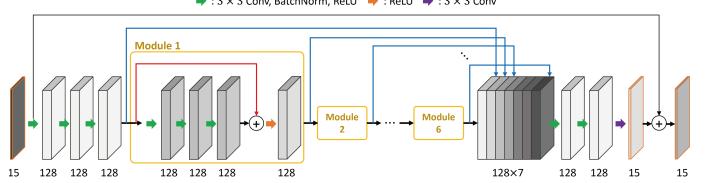


Fig. 2. The proposed WavResNet backbone (i.e. $\mathcal{Q}(f)$ in Algorithm 1 and Fig. 1) for low-dose X-ray CT restoration.

Specifically, for a given signal $f \in \mathbb{R}^n$, the directional sub-band transform $\{T_k\}_{k=1}^p, T_k \in \mathbb{R}^{n \times n}$ in contourlet transform satisfies the resolution of identity:

$$\sum_{k=1}^p \tilde{T}^{(k)\top} T^{(k)} = I_{n \times n}, \qquad (25)$$

which implies that there exists inverse transform $\{\tilde{T}_k\}_{k=1}^p$ to utilize all subband signals to recompose the original signal. Thus, the signal $f$ can be decomposed into directional components:

$$f = \sum_{k=1}^p \tilde{T}^{(k)\top} T^{(k)} f = \sum_{k=1}^p \tilde{T}^{(k)\top} f_k$$

where $f_k = T^{(k)} f \in \mathbb{R}^n$ corresponds to the $k$-th subband signals. Then, our goal is to obtain the deep convolutional framelet representation of the input matrix

$$T f := \begin{bmatrix} T^{(1)} f & \cdots & T^{(p)} f \end{bmatrix}.$$

Because this multi-channel input size is big, we further decompose the signal using patch extraction operator

$$Z_l = P_l T f = \begin{bmatrix} P_l T^{(1)} f & \cdots & P_l T^{(p)} f \end{bmatrix}.$$

where $\{P_l\}_l$ denotes the (overlapping) patch that has the same spatial location across all subbands. The WavResNet does not use any pooling, i.e. $\Phi = \tilde{\Phi} = I_{n \times n}$, since the global correlations have been removed using the contourlet transform. Thus, by inserting each $Z_l$ in (14) and (13), we have

$$Z_l = C \circledast \nu(\tilde{\Psi}) \qquad (26)$$
$$C = (Z_l) \circledast \overline{\Psi} . \qquad (27)$$

The successive layers are similarly implemented using the standard multi-channel convolution in CNN. The resulting patch-by-patch CNN processing are performed on all parts of images, and the final results are obtained by averaging.

Another important component of WavResNet is the *boosting* using the concatenation layer. This is closely related to the boosting scheme in classification that combines multiple weak classifiers to obtain a stronger classifier [40]. Specifically, suppose that perfect recovery (PR) condition satisfies for all cascade of encoder-decoder network. Then, the recovery

condition for deep convolutional framelets up to $L$-layer can be written by

$$Z_l = C^{(1)} \circledast \nu\left(\tilde{\Psi}^{(1)}\right)$$

$$\vdots$$

$$Z_l = C^{(L)} \circledast \nu\left(\tilde{\Psi}^{(L)}\right)\cdots \circledast \nu\left(\tilde{\Psi}^{(1)}\right)$$

where

$$C^{(i)} = \begin{cases} \left(C^{(i-1)} \circledast \overline{\Psi}^{(i)}\right), & 1 \le i \le L \\ Z_l, & i = 0 \end{cases} \qquad (28)$$

and the superscript $^{(i)}$ denotes the $i$-th layer. Thus, for a given intermediate encoder output $\{C^{(l)}\}_{l=1}^{L}$, by denoting $h^{(l)} := \nu\left(\tilde{\Psi}^{(l)}\right)\cdots \circledast \nu\left(\tilde{\Psi}^{(1)}\right)$, a *boosted* decoder can be constructed by combining multiple decoder representation:

$$Z_l = \sum_{i=1}^{L} w_i\left(C^{(i)} \circledast h^{(i)}\right), \qquad (29)$$

where $\sum_{i=1}^{L} w_i = 1$. This procedure can be performed using a single multi-channel convolution after concatenating encoder outputs, as shown in Fig. 3(a). In Experimental Results, we will show that this provides improved denoising performance thanks to the boosting effect.

Finally, WavResNet has the skipped connection [41] as shown in Fig. 3(b). In order to understand the role of the skipped connection, note that the low-dose input $f_{(i)}$ is contaminated with noise so that it can be written by

$$f_{(i)} = f_{(i)}^{*} + h_{(i)},$$

where $h_{(i)}$ denotes the noise components and $f_{(i)}^{*}$ refers to the noise-free ground-truth. Then, the network training (21) using the skipped connection can be equivalently represented as the network training to estimate the artifacts:

$$\min_{(\Psi, \tilde{\Psi})} \sum_{i=1}^{N} \left\| h_{(i)} - \tilde{\mathcal{Q}}(f_{(i)}; \Psi, \tilde{\Psi}) \right\|^2 \qquad (30)$$

where

$$\tilde{\mathcal{Q}}(f_{(i)}; \Psi, \tilde{\Psi}) = \left(\tilde{\Phi} C[f_{(i)}^{*} + h_{(i)}]\right) \circledast \nu(\tilde{\Psi}), \qquad (31)$$

$$C[f_{(i)}^{*} + h_{(i)}] = \Phi^{\top}\left((f_{(i)}^{*} + h_{(i)}) \circledast \overline{\Psi}\right).$$

Therefore, if we can find a convolution filter $\overline{\Psi}$ such that it approximately annihilates the true signal $f_{(i)}^{*}$ [42]:

$$f_{(i)}^{*} \circledast \overline{\Psi} \simeq 0 \implies C[f_{(i)}^{*} + h_{(i)}] \simeq C[h_{(i)}] \qquad (32)$$

then we can find the decoder filter $\tilde{\Psi}$ such that

$$\begin{aligned} \left(\tilde{\Phi} C[h_{(i)}]\right) \circledast \nu(\tilde{\Psi}) &= \left(\tilde{\Phi}\Phi^{\top}\left((h_{(i)}) \circledast \overline{\Psi}\right)\right) \circledast \nu(\tilde{\Psi}) \\ &= h_{(i)} \circledast \overline{\Psi} \circledast \nu(\tilde{\Psi}) \\ &\simeq h_{(i)}. \end{aligned}$$

Thus, our deep convolutional framelet with a skipped connection can estimate the artifact $h_{(i)}$ and remove it from $f_{(i)}$. On
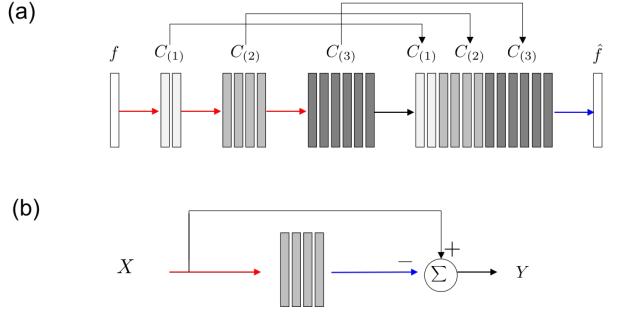


Fig. 3. (a) Concatenation layer, and (b) skipped connection.

the other hand, using the similar argument, we can see that if

$$h_{(i)} \circledast \overline{\Psi} \simeq 0 \implies C[f_{(i)}^{*} + h_{(i)}] \simeq C[f_{(i)}^{*}] \qquad (33)$$

then a deep convolutional framelet *without* the skipped connection can directly recover the ground-truth signal $f_{(i)}^{*}$, i.e. $\mathcal{Q}(f_{(i)}; \Psi, \tilde{\Psi}) \simeq f_{(i)}^{*}$. Then, which one is better ? In our case, the true underlying signal has lower dimensional structure compared to the random CT noises, so the annihlating filter relationship in (32) is more easier to achieve [42]. Therefore, we use the skipped connection as shown in Fig. 3(b).

By combining the contourlet transform, boosting and skipped connection, we conjecture that WavResNet can represent the signal much more effectively which makes the deep convolutional framelet denoising effective.
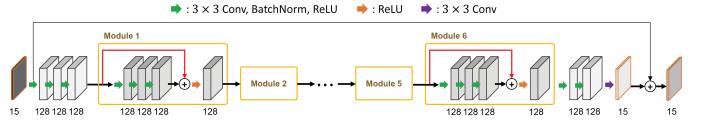


Fig. 4. A symmetric network architecture to investigate the importance of boosting layers in WavResNet.

## III. METHOD

### A. Proposed network architecture

In Fig. 2, we first apply non-subsampled contourlet transform to generate 15 channels inputs [39]. There is no downsampling or up-sampling in the filter banks; thus, it is a shift invariant. We used 4 level decomposition and 8, 4, 2, 1 directional separations for each level, which produces the total 15 bands. Thus, we have 15 subband channels.

The first convolution layer uses 128 set of $3 \times 3 \times 15$ convolution kernels to produce 128 channel feature maps. The shift invariant contourlet transform allows the patch processing and we used $55 \times 55 \times 15$ patches for the training and inference. Then, the final contourlet coefficients are obtained by taking patch averaging. Based on the calculation in [25], a sufficient condition to meet the PR is that the number of output channel should be 270, which is bigger than 128 channels. Thus, the first layer performs a low-rank approximation of the first layer Hankel matrix. Then, the following convolution layers use two $3 \times 3 \times 128$ convolution kernels, which is again believed to perform low rank approximation of the extended Hankel matrix approximation. Later, we will provide an empirical

result showing that the singular value spectrum of the extended Hankel matrix indeed becomes compressed as we go through the layers.

We have 6 set of main module composed of 3 sets of convolution, batch normalization, and ReLU layers, and 1 bypass connection with a convolution and ReLU layer. Finally, as shown in Fig. 2, our network has the end-to-end bypass connection [25] so that we can directly estimate the noise-free contourlet coefficients while exploiting the advantages of skipped connection [41]. Another uniqueness of the proposed network is that it has the concatenation layer as shown in Fig. 3(a). Specifically, our network concatenates the outputs of the individual modules, which is followed by the convolution layer with 128 set of $3 \times 3 \times 896$ convolution kernels. As discussed before, this corresponds to the signal boosting using multiple signal representation. In optimization aspect, this also provides various paths for gradient back-propagation. Finally, the last convolution layer uses 15 sets of $3 \times 3 \times 128$ convolution kernels. This may correspond to the pair-wise decoder layers with respect to the first two convolutional layers.

### B. Network training

We trained two networks: a feed-forward network and an RNN. We applied stochastic gradient descent (SGD) optimization method to train the proposed network. The size of mini-batch was 10. The convolution kernels were initialized by random Gaussian distribution. The learning rate was initially set to 0.01 and decreased continuously down to $10^{-7}$. The gradient clipping was employed in the range $[-10^{-3}, 10^{-3}]$ to use a high learning rate in the initial steps for fast convergence. For data augmentation, the training data were randomly flipped horizontally and vertically. Our network was implemented using MatConvNet [43] in MATLAB 2015a environment (Mathworks, Natick).

The training processes are composed by three stages. In stage 1, we trained the network using original database $DB_0$ which consists of a pair of quarter-dose and routine-dose CT images. After the network converged initially, stage 2 is proceeded sequentially. In stage 2, we add databases $DB_i$ gradually which consists of quarter-dose input, inference results from $\mathcal{Q}_k(f_i)$, and routine-dose CT images. Here, $\mathcal{Q}_k$ dentoes the trained network until $k$-th epochs and $f_i$ is the $i$-th inference results. Finally, in stage 3, we added a database whose both input and target images are routine-dose CT images. The theoretical background of such training is from the framelet nature of deep convolutional neural network [25]. Specifically, the neural network training is to learn the framelet bases from the training data that has the best representation of the signals. Thus, the learned bases should be robust enough to have near optimal representation for the given input data set. In our KM iteration, each iterative steps provided the improved images, which needs to be fed into the same neural network. Thus, the framelet bases should be trained to be optimal not only for the strongly aliased input but also for the near artifact-free images. The resulting network was used for both our RNN and feed-forward networks.

### C. Training dataset

We used projection data obtained from "2016 Low-Dose CT Grand Challenge". The raw projection data were measured by a 2D cylindrical detector that moves along a helical trajectory using a z-flying focal spot [44]. These projections were approximated into fanbeam projection data by a single slice rebinning technique [45]. We reconstructed X-ray CT images using conventional filtered backprojection algorithm. The number of pixels in X-ray CT images is $512 \times 512$ and the slice thickness is 3mm. We have 9 patient data sets of routine dose and quarter dose data for the training. Eight patient data were used for the training and validation, and the remaining one patient data was used for testing. Among the eight patient data, we used 3236 slices for the training and the remaining 350 slices for the validation.

For phantom studies, we used CT image data from Seoul National University Bundang Hospital, Korea. The number of pixels is $512 \times 512$ and the slice thickness is 4mm. The network was trained using 50 patient data sets of routine-dose, 13% dose, 25% dose, and 50% dose images. This study received technical support from Siemens Healthcare (Erlangen, Germany) to simulate CT images at various low dose levels. Low dose images were simulated by inserting Poisson noise into the projection data of routine dose and reconstructing images from those projection data using the filtered back projection method. The number of individual dose images is 7617 slices and we arbitrarily selected 540 slices and 60 slices for training and validation, respectively, for every 50 epoch. To quantify the resolution and contrast at test phase, we used the Catphan 500 (The Phantom Laboratory, Salem, NY, USA) composed of various modules. The spatial resolution was evaluated using a high-resolution module (CTP528), and contrast-to-noise ratio (CNR) were evaluated using a low contrast module (CTP515). The contrast was calculated as the difference of mean CT numbers between a supra-slice target (contrast of 1.0% and diameter of 15mm) and the adjacent background area. The noise was defined as the standard deviation of the CT numbers of the adjacent background area. Mean CT numbers and standard deviations were calculated from circular region-of-interest (ROI) having a diameter of 1 cm in the target and the adjacent background area. ROIs were placed at the exactly same locations on the images produced by different algorithms.

### D. Baseline algorithms

We compared the proposed method with the other denoising algorithms such as BM3D [46], MBIR regularized by total variation (TV), ALOHA [47], the image domain deep learning approach (RED-CNN) [14], and CNN with GAN loss [21]. MBIR regularized by TV was solved using an alternating direction method of multiplier (ADMM) [2] and Chambolle's proximal TV [48]. The details of RED-CNN and GAN were obtained from the original paper and we have implemented them accordingly [14], [21].

To verify the improvement of the new algorithm, we perform comparative study with our previous deep network in wavelet domain for "2016 Low-Dose CT Grand Challenge"
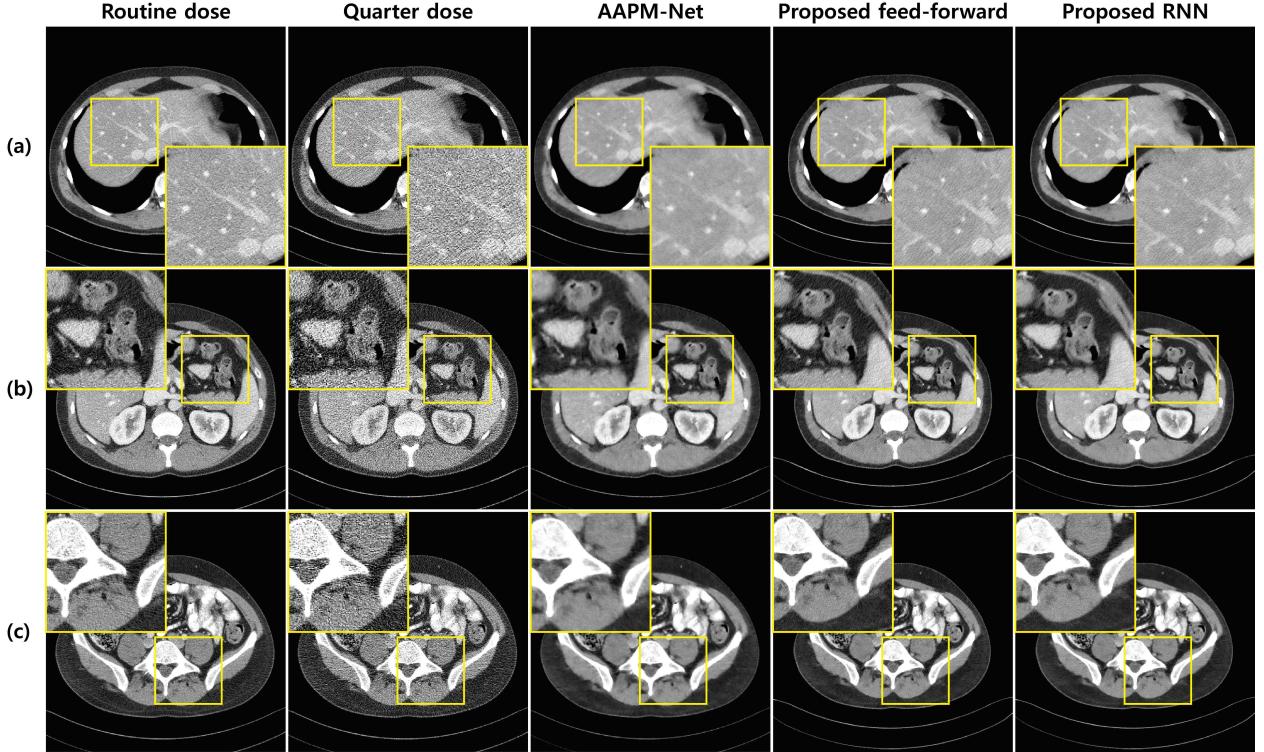
Fig. 5. Transverse view restoration results with routine-dose and quarter-dose images. AAPM-Net is the algorithm which we applied to the "2016 Low-Dose CT Grand Challenge". Intensity range is (-160,240) [HU] (Hounsfield Unit). (a) Example of liver, (b) example of intestine and (c) example of pelvic bone.

[13]. We call this network as AAPM-Net. The main difference between the proposed one and the AAPM-Net comes from the definition of the target images. In AAPM-Net, the target images was the original wavelet coefficients except the lowest frequency band. More specifically, in AAPM-Net, the lowest frequency band target is the residual, whereas the higher frequency band signals are the wavelet coefficients themselves. Therefore, this is not a ResNet from the perspective of deep convolutional framelets. On the other hand, in WavResNet, the residual wavelet coefficients between the routine-dose and low-dose inputs are estimated for every subband. In the current implementation, the final network output is the artifact-corrected images by subtracting the estimated artifacts using the end-to-end skipped connection as shown in Fig. 2. In order to demonstrate the importance of signal boosting, we also implemented a symmetric network as illustrated in Fig. 4. Except for the concatenation layers, the symmetric network also has identical 6 modules structures with symmetric encoder and encoder structures.

## IV. EXPERIMENTAL RESULTS

### A. Comparison with AAPM-net

To confirm the improvement over the AAPM-Net, we present the restoration images of one patient data in the test data set (see Fig. 5). This data have routine-dose images that can be used for subjective evaluation and objective evaluation using RMSE, PSNR, and SSIM.

In Fig. 5, various kinds of slices such as liver and pelvic bones are described and the magnified images are expressed in the yellow boxes. In AAPM-Net results, noise level was
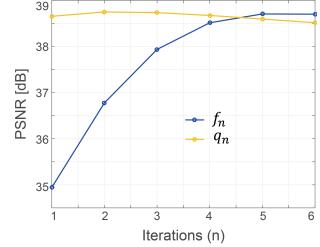


Fig. 6. PSNR values of restoration results according to iterations are plotted.

significantly reduced, but the results are blurry and loses some details. On the other hand, the result of the proposed networks (feed-forward and RNN) results clearly shows that the improved noise reduction while maintaining the edge details and the textures which is helpful for diagnostic purpose. More specifically, for the case of an liver image in Fig. 5 (a), the proposed network results retain the fine details such as vessels in the liver and it has better sharpness than the AAPM-Net. In Fig. 5 (b), the detail of internal structure of intestine was not observed in quarter-dose images and AAPM-Net results, while they are well-recovered in proposed network results. To examine the streaking noise reduction ability, we presented the slice which has the pelvic bone in Fig. 5(c). Proposed networks were again good at preserving the edge details such as inside region of the bones and the texture of the organ which located between the bones, while the streaking artifacts were completely removed. Among the feed-forward and RNN network structure, the results by RNN suppresses more streaking artifacts compared to the feedforward neural network. In Fig. 6, the PSNR plots for $\mathcal{Q}(f_n)$ and $f_n$ in
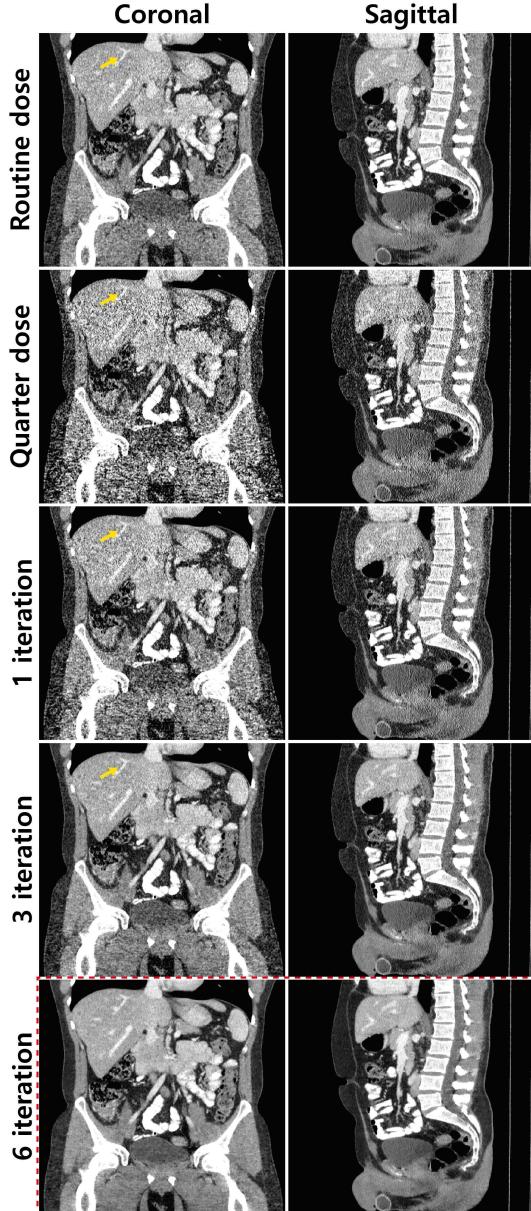
Fig. 7. Coronal and sagittal view restoration results along RNN iterations. Intensity range is (-160,240) [HU]. The red dashline boxed images come from the last iteration.
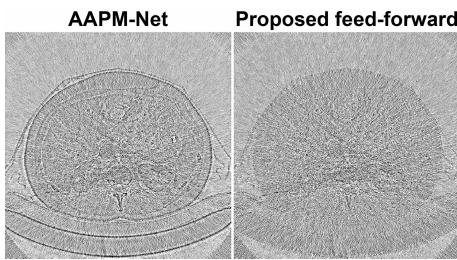


Fig. 8. Difference images between restoration images and routine dose image. Images are same slice in the second row of Fig. 5. Intensity range is (-1100,-950) [HU].

Algorithm 1 are illustrated. The result shows that averaged PSNR values for RNN ($f_n$) are increased according to the iterations and it converged after 5 iterations. On the other hand,

PSNR values for the direct network output $\mathcal{Q}(f_n)$ increase initially but started to decrease with iteration after the initial peak. Eventually, the PSNR values for RNN surpass the feed-forward network. This again confirms the convergence of the proposed RNN approach thanks to the KM iteration. However, our feed-forward network at the 1st iteration is also useful thanks to the computational advantage, so we provide the both results. In our KM iteration, $\mu$ is a parameter for incorporating the effect of the original low-dose image. This parameter also control the convergence behavior of KM iteration. The experimental results with various values of $\mu$ showed that the lower value tends to provide better results. However, if $\mu$ is less than 0.1, it did not converge and the image quality decreased in terms of iterations. Thus, we set $\mu = 0.1$ to get the best results while the algorithm retains the robustness.

The coronal and sagittal view of the restoration results by our RNN are described in Fig. 7. The quarter-dose images show that the noise levels are different depending on the slices. The lower part of the images exhibit a high noise level because the pelvic bones are included. The results shows that the noise level was reduced gradually according to the iterations. The last iteration reconstructed result maintains the edge details and textures which is helpful for diagnostic purpose. The yellow arrow indicates the vessel in the liver and the last iteration result has better sharpness. The proposed method can remove a wide range of noise levels and maintain the texture and edge information. The difference images between the result images and routine-dose images in Fig. 8 confirm the superiority of our method over AAPM-net. The difference images of proposed network only contains the noise of low-dose X-ray CT images, while the difference images of AAPM-Net also contains the edge information.
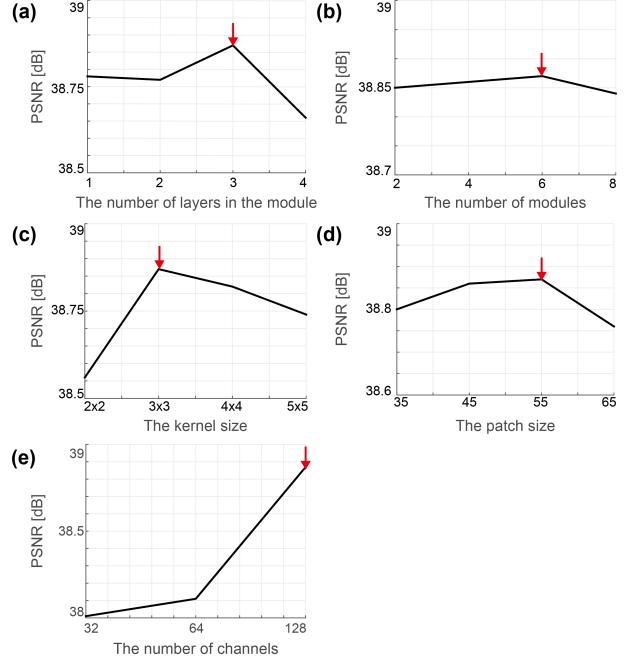


Fig. 9. Performance dependency on various hyper-parameters of the network.

## B. Ablation study

To demonstrate the advantages of the proposed method, we preformed the ablation study by excluding some structures from the network and applying the same training procedures. Table I presents the averaged RMSE, PSNR and SSIM index values of the results from 486 slices. The qualitative results shows that proposed feed-forward and RNN network have the best results, and among them the RNN was better. The PSNR and SSIM values of the symmetric network in Fig. 2 are lower than those of the proposed methods, which confirms the signal boosting effect from the concatenation.

In addition, we have investigated the effects of network hyper-parameters such as the number of channels, the number of layers in module, the number of modules, the kernel size, and the patch size as shown in Fig. 9. Here, network performance improves with more layers in each module until it reaches 3. With more than three layers we have found that the network is difficult to train due to many parameters to be optimized. As the number of modules increases, network performance improves slightly at the expense of increased reconstruction time. Given the compromise between performance and reconstruction time, we used six modules for our network. We have observed that the filter size $3 \times 3$ gave the best result with reasonable processing time for real applications. In addition, we found that the patch size is not critical. However, the reconstruction time and its receptive field lead us to choose the patch size of $55 \times 55$ for our network. Finally, with more channels, the performance improved. But due to the memory requirement as well as to prevent overfitting, we chose 128 channels.

### TABLE I
### ANALYSIS OF NETWORK STRUCTURE

| | RMSE | PSNR [dB] | SSIM index |
|---|---|---|---|
| Exclude external bypass connection | 48.09 | 33.63 | 0.828 |
| Exclude concatenation layer (symmetric) | 28.43 | 38.20 | 0.893 |
| Proposed feed-forward (128 channels) | 27.32 | 38.54 | 0.899 |
| Proposed RNN (128 channels) | 26.90 | 38.70 | 0.893 |

## C. Low-rank approximation property

To verify our theory that CNN is closely related to the Hankel matrix decomposition [25], we performed experiments to verify whether the trained network imposes low-rank approximation of the Hankel matrix. Specifically, we constructed extended Hankel matrices using the output channel images from each module in Fig. 2. Then, we plotted the singular value spectrum in Fig. 10. Blue dashed line is the result of the first module's extended Hankel matrix constructed from output feasture maps and the red solid line is the sixth module's extended Hankel matrix constructed from output features maps. We observed that the singular value spectrum becomes compressed as we go through layers, which indicates each layer of CNN performs the low-rank approximation of the Hankel matrix. These results clearly suggests the link between trained CNN and the low rank Hankel matrix decomposition.

## D. Comparison with existing algorithms

Fig. 11 shows the results by the comparative algorithms such as BM3D [46], MBIR regularized by TV, ALOHA [47],
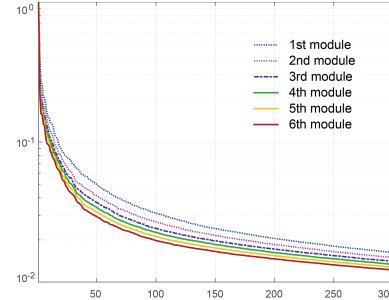


Fig. 10. The singular value spectrum of the extended Hankel matrix along layers.

RED-CNN [14], and GAN loss [21]. BM3D is a state-of-art of image denoising algorithm using nonlocal patch processing, MBIR is currently a standard algorithm of low-dose X-ray CT images and the RED-CNN is recently proposed deep network for low-dose X-ray CT and ALOHA is the latest low rank Hankel matrix method.

The intensity of the transverse view in Fig. 11 is adjusted to see inside structures of the lung. The result of BM3D loses the details in the lung such as vessels and exhibited some cartoon artifact. The result of MBIR appears a little blurred and textures are reconstructed incorrectly. On the other hand, deep learning based denoising algorithms have better performances than the other algorithms. However, RED-CNN results are somewhat blurry and exhibits remaining noises in the coronal view, while the proposed method provides clear restoration results.

### TABLE II
### EXECUTION TIME (MINI-BATCH: $55 \times 55 \times 10$, SLICE: $512 \times 512$)

| Training | Time [mini-batch/sec] | Implementation environment |
|---|---|---|
| RED-CNN | 0.19 | MatConvNet, GTX 1080 Ti |
| Proposed | 0.44 | MatConvNet, GTX 1080 Ti |
| **Restoration** | **Time [slice/sec]** | **Implementation environment** |
| BM3D | 2.73 | MATLAB, i7-4770 |
| MBIR TV | 9.45 | MATLAB, GTX 1080 |
| ALOHA | 1405 | MATLAB, GTX 1080 |
| RED-CNN | 0.38 | MatConvNet, GTX 1080 Ti |
| Proposed feed-forward | 2.05 | MatConvNet, GTX 1080 Ti |

With regard to the computation time, the CNN frameworks need learning to train the networks. Our method took 20 hours to train the network through 3 stages as described in Section III, and RED-CNN took 6 hours to train. For the restoration step, the CNN framework is advantageous compared to the other classical algorithms such as BM3D or MBIR TV. Our method takes approximately 2.05 seconds per slice for restoration which have $512 \times 512$ pixels with MATLAB implementation using a graphical processing unit (NVidia GeForce GTX 1080 Ti).

## E. Contrast and resolution loss study

To evaluate the contrast and spatial resolution loss, we compared the various algorithms using Catphan phantom at various
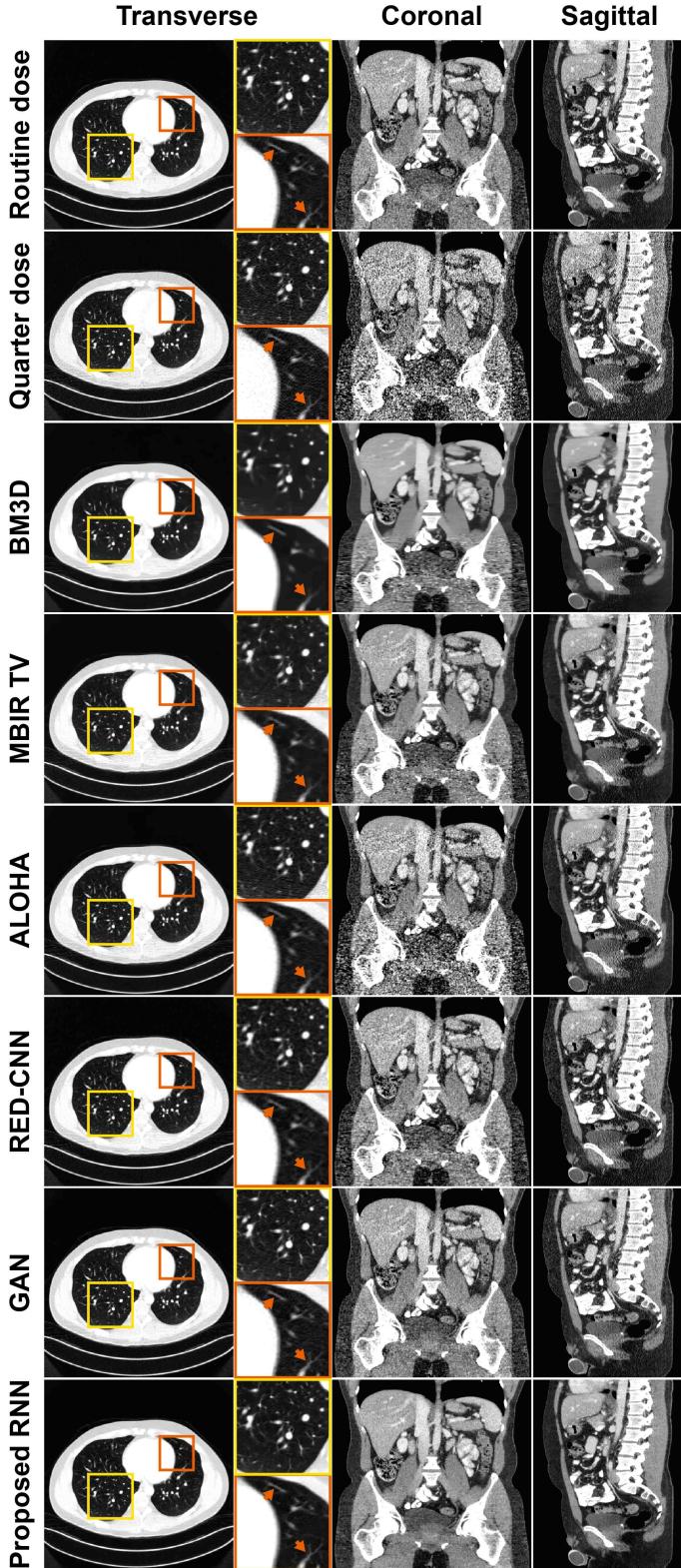
Fig. 11. Restoration results with routine-dose and quarter-dose image. Transverse view restoration images' intensity range is adjusted to see the details in the lung. Intensity range is (-1000,100) [HU].

CNN and the proposed method preserve the resolution lines better than other methods, and among them our method was better at all resolution grids.
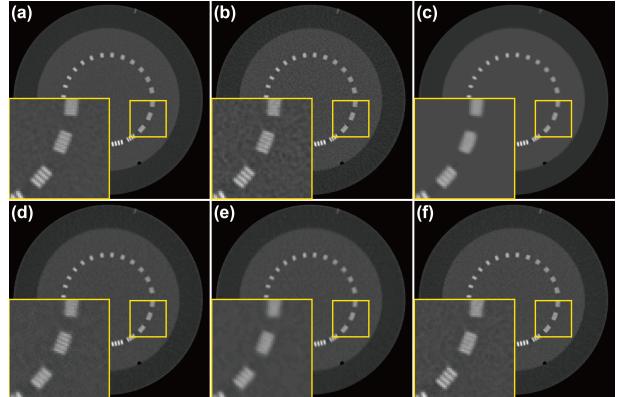


Fig. 12. Reconstruction results of Catphan data from %25 dose. Images by (a) routine dose, (b) quarter dose input, (c) BM3D, (d) RED-CNN, (e) GAN, and (f) the proposed feed-forward network.

To investigate the spatial resolution loss at lose dose level, Fig. 13 illustrates the intensity profile along the two resolution grids in Catphan phantom at various dose level. Down to quarter dose level, the proposed feed-forward network does not exhibit significant resolution loss. At 13% dose, we started to observe resolution loss especially at area (a). In addition, Table III shows the contrast to noise (CNR) variations at various radiation dose levels. The CNR values of our method gradually decreases from 50% to 13%, but they outperformed the CNR values of FBP results at the same dose levels. These results clearly confirm the robustness of the proposed method at various dose levels.
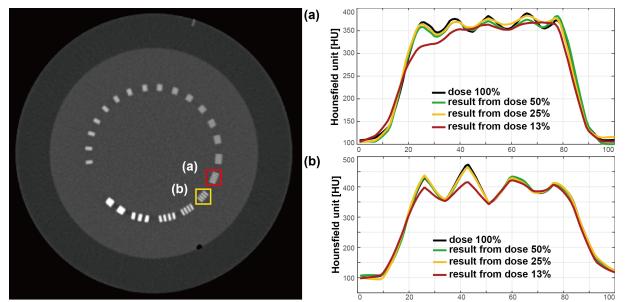


Fig. 13. Resolution profile at various level by the proposed feed-forward network.

TABLE III
CNR EVALUATION AT VARIOUS DOSE LEVELS

| Radiation dose level | 100% | 50% | | 25% | | 13% | |
|---|---|---|---|---|---|---|---|
| Alogrithm | FBP | FBP | Proposed | FBP | Proposed | FBP | Proposed |
| Contrast | 8.49 | 8.44 | 8.44 | 8.96 | 9.46 | 10.85 | 11.08 |
| Noise | 6.38 | 8.94 | 6.24 | 13.69 | 9.46 | 17.26 | 11.45 |
| CNR | 1.33 | 0.94 | 1.35 | 0.65 | 1.00 | 0.63 | 0.97 |

dose levels such as 13%, 25%, 50% of the original dose. In Fig. 12, the representative restoration results from 25% dose are illustrated, where the magnified areas are indicated by yellow boxes. By visual inspection, we can see that RED-

*F. Evaluation of lesion detection*

To verify the proposed method, we have performed task-driven experiment to evaluate the lesion detection performance

TABLE IV
EVALUATION OF LESION DETECTION

|  | Ground-truth | Quarter-dose | MBIR TV | Proposed |
|---|---|---|---|---|
| The number of lesions | 37 | 21 | 23 | 27 |
| Lesion detection rate | - | 57% | 62% | 73% |

by quarter-dose FBP, MBIR TV and a proposed method. We used 20 test data sets from the 2016 Low dose CT grand challenge. A board-certified radiologist (Won Chang) with seven years of experience in liver CT interpretation assessed the data set and recorded the exact locations of the lesions with blind to the reconstruction methods. The detection rates for solid focal hepatic lesions were compared using the McNemar test with Bonferroni correction and a difference with a $p > 0.017$ was considered significant. Since small number of radiologist involved in this study, it was not statistically significant, but the proposed method showed a significant higher lesion detection rate than FBP (73% vs. 57%, p-value=0.0412) and MBIR(73% vs. 62%, p-value=0.1336). With more radiologists involved, we are currently studying large scale statistical evaluation of the method, which will be reported later in a clinical journal.

## V. CONCLUSION

In this paper, we proposed a deep convolutional framelet-based denoising algorithm for low-dose X-ray CT restoration by synergistically combining the proven convergence of the classical framelet-based algorithm and the expressive power of deep learning. To provide the theoretical background for performance improvement, we employed the recent proposal of deep convolutional framelets that interprets a deep learning as a multilayer implementation of convolutional framelets with ReLU nonlinearity. Our theory resulted in two network structures: a feed-forward and RNN architectures. Moreover, by combining the redundant global transform, residual network (ResNet) and signal boosting from concatenation layers, the proposed feed-forward and RNN network provided significant improvement compared to the prior work by retaining the detailed texture. Using extensive experimental results, we showed that the proposed network is good at streaking noise reduction and preserving the texture details of the organs while the lesion information is maintained.

## APPENDIX A
### MATHEMATICAL PRELIMINARIES

For a given mapping $T : D \to \mathcal{H}$, the set of the *fixed points* of an operator $T : D \to D$ is denoted by $\mathrm{Fix}T = \{x \in D \mid Tx = x\}$. Then, $T$ is called *non-expansive* if

$$\|Tx - Ty\| \leq \|x - y\|, \quad \forall x, y \in D, \tag{34}$$

Then, we have the following convergence theorem for the non-expansive operator:

**Theorem 1.2** (Krasnoselski-Mann algorithm)**.** *[38] Let D be a nonempty closed convex subset of H, let $T : D \mapsto D$ be a nonexpansive operator such that $\mathrm{Fix}T \neq \emptyset$, let $(\lambda_n)$ be a sequence in $[0, 1]$ such that $\sum_{n=1}^{\infty} \lambda_n(1 - \lambda_n) = +\infty$ and let $f_0 \in D$. Consider the following seqeunce:*

$$f_{n+1} = f_n + \lambda_n(Tf_n - f_n). \tag{35}$$

*Then, the sequence $f_n$ converges to a point in $\mathrm{Fix}T$.*

## APPENDIX B
### PROOF OF THEOREM 2.1

Let the mapping $T$ be defined by

$$T(f) \quad := \quad \mu g + (1 - \mu)\mathcal{Q}(f)$$

where $\mathcal{Q}$ is the deep convolutional framelet network output. Our goal is to show that the operator $T$ is non-expansive. Note that

$$\begin{aligned} \|Tx - Ty\| &= \|(1 - \mu)\mathcal{Q}(x) - (1 - \mu)\mathcal{Q}(y)\| \\ &= (1 - \mu)\|\mathcal{Q}(x) - \mathcal{Q}(y)\| \\ &\leq (1 - \mu)\|\mathcal{Q}'(z)\|\|x - y\| \end{aligned}$$

where $\mathcal{Q}'(z)$ denotes the Jacobian of the network at $z$, and we use the mean value theorem for the last inequality. Therefore, if the Jacobian of the deep convolutional framelet is finite, we can choose $\mu = 1 - 1/\max_{z \in D} \|\mathcal{Q}'(z)\|$ such that

$$\|Tx - Ty\| \leq \|x - y\|,$$

i.e. $T$ is non-expansive. Thus, the remaining step is to show that $\max_{z \in D} \|\mathcal{Q}'(z)\| < \infty$. This is true because the convolutional framelets consists of convolutional filters with finite coefficients and ReLU, so the Jacobian can be represented as the product of the filter norms [49]. Thus, by denoting $\bar{f}_{n+1} = Tf_n$, Theorem 1.2 informs that our KM iteration for deep convolutional framelet inpainting converges. This concludes the proof.

## REFERENCES

[1] M. Beister, D. Kolditz, and W. A. Kalender, "Iterative reconstruction methods in x-ray CT," *Physica medica*, vol. 28, no. 2, pp. 94–108, 2012.
[2] S. Ramani and J. A. Fessler, "A splitting-based iterative algorithm for accelerated statistical x-ray CT reconstruction," *IEEE transactions on medical imaging*, vol. 31, no. 3, pp. 677–688, 2012.
[3] E. Y. Sidky and X. Pan, "Image reconstruction in circular cone-beam computed tomography by constrained, total-variation minimization," *Physics in medicine and biology*, vol. 53, no. 17, p. 4777, 2008.
[4] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in neural information processing systems*, 2012, pp. 1097–1105.

[5] H. C. Burger, C. J. Schuler, and S. Harmeling, "Image denoising: Can plain neural networks compete with BM3D?" in *2012 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2012, pp. 2392–2399.

[6] X. Mao, C. Shen, and Y.-B. Yang, "Image restoration using very deep convolutional encoder-decoder networks with symmetric skip connections," in *Advances in Neural Information Processing Systems*, 2016, pp. 2802–2810.

[7] K. Zhang, W. Zuo, Y. Chen, D. Meng, and L. Zhang, "Beyond a gaussian denoiser: Residual learning of deep CNN for image denoising," *arXiv preprint arXiv:1608.03981*, 2016.

[8] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2015, pp. 234–241.

[9] J. Kim, J. Kwon Lee, and K. Mu Lee, "Accurate image super-resolution using very deep convolutional networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 1646–1654.

[10] S. Wang, Z. Su, L. Ying, X. Peng, S. Zhu, F. Liang, D. Feng, and D. Liang, "Accelerating magnetic resonance imaging via deep learning," in *2016 IEEE 13th International Symposium on Biomedical Imaging (ISBI)*. IEEE, 2016, pp. 514–517.

[11] K. Hammernik, F. Knoll, D. Sodickson, and T. Pock, "Learning a variational model for compressed sensing MRI reconstruction," in *Proceedings of the International Society of Magnetic Resonance in Medicine (ISMRM)*, 2016.

[12] J. Sun, H. Li, Z. Xu *et al.*, "Deep admm-net for compressive sensing mri," in *Advances in Neural Information Processing Systems*, 2016, pp. 10–18.

[13] E. Kang, J. Min, and J. C. Ye, "A deep convolutional neural network using directional wavelets for low-dose x-ray ct reconstruction," *Medical Physics*, vol. 44, no. 10, 2017.

[14] 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 transactions on medical imaging*, vol. 36, no. 12, pp. 2524–2535, 2017.

[15] H. Chen, Y. Zhang, W. Zhang, P. Liao, K. Li, J. Zhou, and G. Wang, "Low-dose CT via convolutional neural network," *Biomedical optics express*, vol. 8, no. 2, pp. 679–694, 2017.

[16] J. Adler and O. Öktem, "Learned primal-dual reconstruction," *arXiv preprint arXiv:1707.06474*, 2017.

[17] H. Chen, Y. Zhang, W. Zhang, H. Sun, P. Liao, K. He, J. Zhou, and G. Wang, "Learned experts' assessment-based reconstruction network (" learn") for sparse-data ct," *arXiv preprint arXiv:1707.09636*, 2017.

[18] T. Würfl, F. C. Ghesu, V. Christlein, and A. Maier, "Deep learning computed tomography," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2016, pp. 432–440.

[19] Q. Yang, P. Yan, M. K. Kalra, and G. Wang, "Ct image denoising with perceptive deep neural networks," *arXiv preprint arXiv:1702.07019*, 2017.

[20] G. Wang, "A perspective on deep imaging," *IEEE Access*, vol. 4, pp. 8914–8924, 2016.

[21] Q. Yang, P. Yan, Y. Zhang, H. Yu, Y. Shi, X. Mou, M. K. Kalra, and G. Wang, "Low dose CT image denoising using a generative adversarial network with Wasserstein distance and perceptual loss," *arXiv preprint arXiv:1708.00961*, 2017.

[22] J. M. Wolterink, T. Leiner, M. A. Viergever, and I. Isgum, "Generative adversarial networks for noise reduction in low-dose ct," *IEEE Transactions on Medical Imaging*, 2017.

[23] D. Wu, K. Kim, G. El Fakhri, and Q. Li, "Iterative low-dose ct reconstruction with priors trained by artificial neural network," *IEEE transactions on medical imaging*, vol. 36, no. 12, pp. 2479–2486, 2017.

[24] H. Gupta, K. H. Jin, H. Q. Nguyen, M. T. McCann, and M. Unser, "CNN-based projected gradient descent for consistent image reconstruction," *arXiv preprint arXiv:1709.01809*, 2017.

[25] J. C. Ye, Y. S. Han, and E. Cha, "Deep convolutional framelets: A general deep learning framework for inverse problems," *SIAM Journal on Imaging Sciences (in press), also available as arXiv preprint arXiv:1707.00372*, 2018.

[26] M. Li, Z. Fan, H. Ji, and Z. Shen, "Wavelet frame based algorithm for 3d reconstruction in electron microscopy," *SIAM Journal on Scientific Computing*, vol. 36, no. 1, pp. B45–B69, 2014.

[27] B. Dong, Q. Jiang, and Z. Shen, "Image restoration: wavelet frame shrinkage, nonlinear evolution pdes, and beyond," *Multiscale Modeling & Simulation*, vol. 15, no. 1, pp. 606–660, 2017.

[28] R. Yin, T. Gao, Y. M. Lu, and I. Daubechies, "A tale of two bases: Local-nonlocal regularization on image patches with convolution framelets," *SIAM Journal on Imaging Sciences*, vol. 10, no. 2, pp. 711–750, 2017.

[29] J.-F. Cai, R. H. Chan, and Z. Shen, "A framelet-based image inpainting algorithm," *Applied and Computational Harmonic Analysis*, vol. 24, no. 2, pp. 131–149, 2008.

[30] J.-F. Cai, R. H. Chan, L. Shen, and Z. Shen, "Convergence analysis of tight framelet approach for missing data recovery," *Advances in Computational Mathematics*, vol. 31, no. 1-3, pp. 87–113, 2009.

[31] E. Kang, J. Yoo, and J. C. Ye, "Wavelet residual network for low-dose CT via deep convolutional framelets," *arXiv preprint arXiv:1707.09938*, 2017.

[32] R. J. Duffin and A. C. Schaeffer, "A class of nonharmonic Fourier series," *Transactions of the American Mathematical Society*, vol. 72, no. 2, pp. 341–366, 1952.

[33] K. H. Jin and J. C. Ye, "Sparse+ low rank decomposition of annihilating filter-based Hankel matrix for impulse noise removal," *arXiv preprint arXiv:1510.05559*, 2015.

[34] K. H. Jin, J.-Y. Um, D. Lee, J. Lee, S.-H. Park, and J. C. Ye, "Mri artifact correction using sparse+ low-rank decomposition of annihilating filter-based hankel matrix," *Magnetic resonance in medicine*, vol. 78, no. 1, pp. 327–340, 2017.

[35] J. Min, L. Carlini, M. Unser, S. Manley, and J. C. Ye, "Fast live cell imaging at nanometer scale using annihilating filter-based low-rank Hankel matrix approach," in *SPIE Optical Engineering+ Applications*. International Society for Optics and Photonics, 2015, pp. 95 970V–95 970V.

[36] E. J. Candès and B. Recht, "Exact matrix completion via convex optimization," *Foundations of Computational mathematics*, vol. 9, no. 6, p. 717, 2009.

[37] J.-F. Cai, E. J. Candès, and Z. Shen, "A singular value thresholding algorithm for matrix completion," *SIAM Journal on Optimization*, vol. 20, no. 4, pp. 1956–1982, 2010.

[38] H. H. Bauschke and P. L. Combettes, *Convex analysis and monotone operator theory in Hilbert spaces*. Springer, 2011.

[39] J. Zhou, A. L. Cunha, and M. N. Do, "Nonsubsampled contourlet transform: construction and application in enhancement," in *IEEE International Conference on Image Processing 2005*, vol. 1. IEEE, 2005, pp. I–469.

[40] R. E. Schapire, Y. Freund, P. Bartlett, W. S. Lee *et al.*, "Boosting the margin: A new explanation for the effectiveness of voting methods," *The annals of statistics*, vol. 26, no. 5, pp. 1651–1686, 1998.

[41] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 770–778.

[42] M. Vetterli, P. Marziliano, and T. Blu, "Sampling signals with finite rate of innovation," *IEEE transactions on Signal Processing*, vol. 50, no. 6, pp. 1417–1428, 2002.

[43] A. Vedaldi and K. Lenc, "MatConvNet: Convolutional neural networks for Matlab," in *Proceedings of the 23rd ACM international conference on Multimedia*. ACM, 2015, pp. 689–692.

[44] T. Flohr, K. Stierstorfer, S. Ulzheimer, H. Bruder, A. Primak, and C. H. McCollough, "Image reconstruction and image quality evaluation for a 64-slice CT scanner with z-flying focal spot," *Medical physics*, vol. 32, no. 8, pp. 2536–2547, 2005.

[45] F. Noo, M. Defrise, and R. Clackdoyle, "Single-slice rebinning method for helical cone-beam CT," *Physics in medicine and biology*, vol. 44, no. 2, p. 561, 1999.

[46] K. Dabov, A. Foi, V. Katkovnik, and K. Egiazarian, "Image denoising by sparse 3-d transform-domain collaborative filtering," *IEEE Transactions on image processing*, vol. 16, no. 8, pp. 2080–2095, 2007.

[47] K. H. Jin and J. C. Ye, "Random impulse noise removal using sparse and low rank decomposition of annihilating filter-based hankel matrix," in *Image Processing (ICIP), 2016 IEEE International Conference on*. IEEE, 2016, pp. 3877–3881.

[48] A. Chambolle, "An algorithm for total variation minimization and applications," *Journal of Mathematical imaging and vision*, vol. 20, no. 1, pp. 89–97, 2004.

[49] J. Sokolic, R. Giryes, G. Sapiro, and M. R. Rodrigues, "Robust large margin deep neural networks," *IEEE Transactions on Signal Processing*, 2017.