

Unsupervised Denoising for Spectral CT Images using a U-Net with Block-Based Training

Raziye Kubra Kumrular^{a*} and Thomas Blumensath^a

^aInstitute of Sound and Vibration Research, The University of Southampton, SO17 1BJ, U.K.

ABSTRACT

Spectral Computed Tomography (CT) is a versatile imaging technique increasingly utilized in industry, medicine, and scientific research. This technique allows us to observe the energy-dependent X-ray attenuation throughout an object by using Photon Counting Detector (PCD) technology. However, a major drawback of Spectral CT is the increase in noise due to a lower photon count per channel, as increasing the number of energy channels without also increasing scan time reduces the photon count per channel. This challenge often complicates quantitative material identification, which is a major application of the technology. In this study, we investigate the use of unsupervised image denoising approaches and demonstrate the applicability of the Noise2Inverse method, an unsupervised denoising method for tomographic imaging. These approaches have the advantage over supervised machine learning methods in that they do not require any additional clean or noisy training data, which can be very difficult to collect in Spectral CT imaging. Our model uses a U-Net paired with a block-based training approach. In particular, we demonstrate that the block-based models can be efficiently trained using small image blocks, each block incorporating spectral information. This training process is performed on images that are reconstructed from subsets of measured Spectral tomography data. The experiments used two simulated Spectral CT phantoms, each with a unique shape and material decomposition. Upon evaluation using the peak signal-to-noise ratio (PSNR) and structural similarity index (SSIM) performance metrics, our approach exhibited improvements compared to two alternative approaches: the unsupervised Low2High method previously employed in sparse Spectral CT imaging and a traditional Iterative reconstruction method that imposes a Total Variation (TV) constraint.

Keywords: Spectral Computed Tomography, Unsupervised Denoising, Block-Based Training, Noise2Inverse

1. INTRODUCTION

Spectral Computed Tomography (CT) imaging has been a very active field of research with thousands of papers published over the past few decades, as it allows us to observe the energy dependence of the object being imaged using Photon Counting Detector (PCD) technology.¹ The applicability of PCD technology to obtain energy-resolved images has been shown in different fields. The first clinical PCD-CT system has demonstrated better resolution and noise characteristics in four different clinical applications than similarly configured energy-integrating CT (EID-CT).² Spectral imaging is also widely used in threat detection during airport luggage security screening,^{3,4} as well as in different applications of Non-Destructive Testing (NDT).^{5,6}

In Spectral imaging, the projection data is intrinsically noisy because there are fewer photons in each energy channel. The energy channel must be carefully chosen to minimize noise because wider energy channels integrate more photons and thus have a lower noise level. Consequently, there is a trade-off between the width of energy channels and noise level. To address these challenges, specialized noise-robust spectral reconstruction techniques have been developed. Some of these methods have focused on dictionary learning methods,⁷ prior-based methods⁸ and tensor-based nuclear norm regularization.⁹ However, these studies⁷⁻⁹ use detectors with a maximum of eight energy channels, each with a width of several keVs, which is far from the ideal assumption of monochromatic acquisitions. More recent detectors provide a significantly finer energy resolution, but increasing the number of energy channels leads to significant computational challenges when using the above iterative algorithms, which operate jointly across the channels.

Kubra Kumrular is thankful for the support from the Republic of Turkiye Ministry of National Education.

(* corresponding author, e-mail: r.k.kumrular@soton.ac.uk)

Investigating alternative approaches based on channel-wise reconstruction thus remains crucial, especially when working with large multi-channel datasets. Recently data-driven approaches have been applied for spectral imaging. A supervised deep learning-based Spectral CT method, which includes information in the spectral domain, was designed to improve reconstructions when the signal is affected by Poisson noise.¹⁰ The challenges of obtaining high-quality reconstructions from sparse measurements for a 64-channel PCD-CT were addressed using an unsupervised denoising method called Low2High¹¹ that can be applied after single-channel reconstruction. In this study, we instead utilise the Noise2Inverse framework.¹² We train a U-Net architecture¹³ using a block-based training approach. As we do not assume the availability of clean training data, we instead utilise pairs of noisy images, each reconstructed from mutually exclusive subsets of projections. Our approach aims to improve image quality and thus help accurate material identification in spectral imaging applications where clean training data sets are not available.

2. BACKGROUND

2.1 Spectral Imaging

The attenuation of an X-ray beam travelling through an object is often modelled using the Beer-Lambert law. For a poly-energetic X-ray spectrum used in Spectral Imaging, an adapted version of the Beer-Lambert Law is:

$$I(E) = I_0(E) e^{-\int_L \mu(E,r) dr} \quad (1)$$

where $I(E)$ and $I_0(E)$ are the transmitted (measured) X-ray intensity and the initial intensity emitting from the X-ray source, at energy level E respectively, both of which include the detector sensitivity. $\mu(E, r)$ is the linear attenuation coefficient (LAC) of the object at energy E and $\int_L \mu(E, r) dr$ represents the line integral of attenuation along one ray path from the source to one detector element at one rotation angle. This line integral sums attenuation along the path r . It is critical to note that both the energy dependency of the X-ray source spectrum $I_0(E)$ and the energy sensitivity of the detector significantly influence the system's overall spectral response. Consequently, Spectral CT images are generated by tomographic reconstruction of each energy channel individually using the measured sinograms $I(E)$, which enables the incorporation of detailed energy information into the images.^{1,6}

2.2 Unsupervised Learning Methods

Data-driven approaches to image denoising can be divided into three categories: unsupervised, supervised, and semi-supervised methods. Here, we focus on the unsupervised methods since there is often a lack of low-noise high-quality reference data in CT imaging applications^{12,14} that could be used for supervised training. Spectral CT applications exemplify this challenge. Thus, denoising methods that may be trained with noisy reference data of paired¹⁵ or single¹⁶ image become of interest. In the Noise2Noise¹⁵ training method, each pair of images contains independent noise of the same slice, which is also generally unavailable in CT. In contrast, the Noise2Self¹⁶ approach uses a single noisy image, which is based on the assumption that noise in one pixel is statistically independent of noise in another pixel, in the training process. When it comes to examining the Noise2Inverse approach, we can see that this method is a new framework that can be applied especially for linear reconstruction methods in tomography imaging. The main concept of Noise2Inverse^{12,14} is that the model is trained on reconstructed images using data that possesses element-wise independent and mean-zero noise in the measurement domain.

Converting the Beer-Lambert law from intensity to absorption and discretizing the integral, we have a linear system

$$\tilde{y} = Ax + \epsilon \quad (2)$$

where \tilde{y} contains the measurements corrupted by a noise that is element-wise independent and zero-mean conditional on the data. A represents the projection operator and x is the discretized absorption image. Noise2Inverse has been implemented in single-energy tomographic reconstruction by reconstructing image pairs

from mutually exclusive subsets of the measurement data. Reconstructed noise in the image pairs is then assumed to be uncorrelated. For denoising, we use a parameterised deep neural network Λ_θ , that is trained by optimising the parameters:

$$\theta^* = \arg \min_{\theta} \frac{1}{|\mathcal{J}|} \sum_{J \in \mathcal{J}} \|\Lambda_\theta(\tilde{x}_{J^C}) - (\tilde{x}_J)\|_2^2 \quad (3)$$

that provide the best prediction of the target reconstruction \tilde{x}_J (reconstructed from the projections in set J), from the input \tilde{x}_{J^C} (reconstructed from the projections in set J^C).

Based on similar reasoning, the Low2High¹¹ approach has been introduced for Sparse Multi-Spectral imaging. In the Low2High method, a different strategy is used to produce two pairs of reconstructed images from the same set of measurements. This is done using a filtered backprojection (FBP) algorithm with two different filters, the standard FBP Ram-Lak(s=1) filter (high) and a Hann(s=0.2) low-pass filter that removes the higher image frequencies (low). The reasoning here is that noise is predominately concentrated in higher frequencies so that the low-frequency image does not contain significant noise and the network is then trained to predict the coherent high frequencies from the low-frequency content whilst the noise averages out in the same way as in Noise2Inverse.

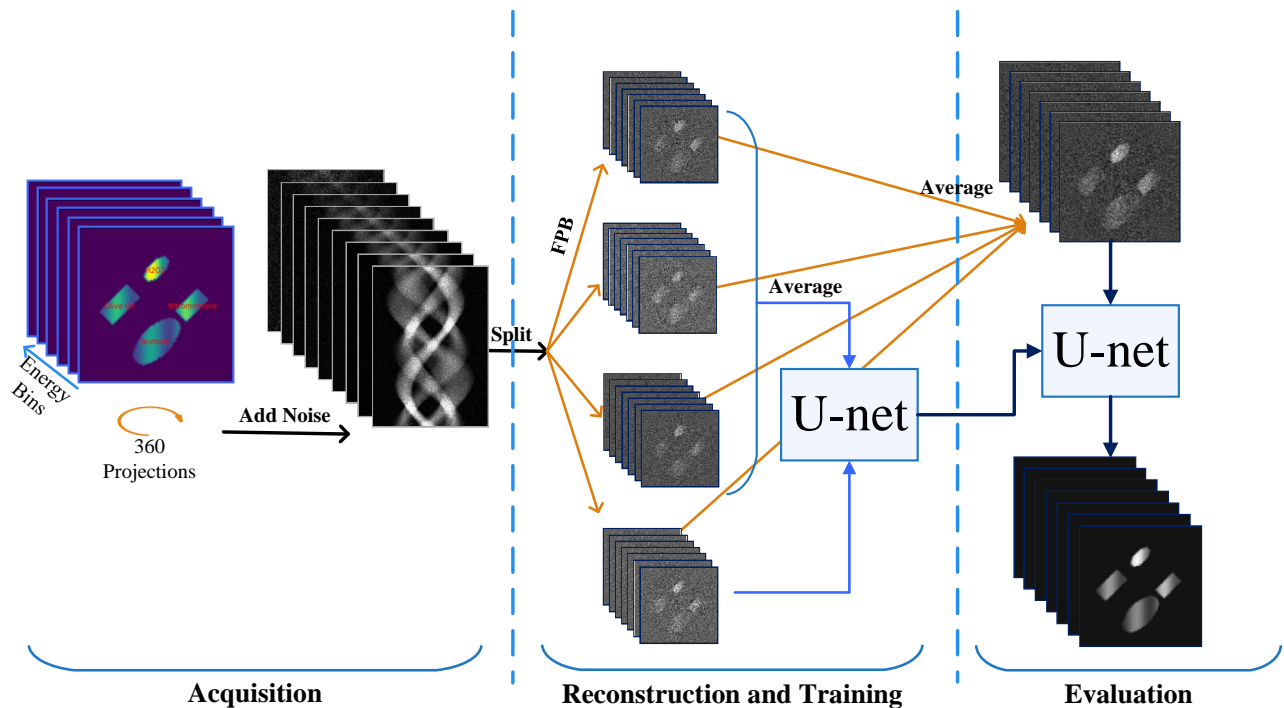


Figure 1: Our approach: The spectral sinogram is obtained over 360 degrees and split into 4 mutually exclusive sets, which are reconstructed independently for each energy channel using FBP. The network is trained using images generated by averaging all possible combinations of 3 reconstructions out of the 4 images as network input to predict the 4th spectral image that was not used to generate the current network input. Once trained, all 4 images are averaged and denoised by the model.

2.3 Our Approach

The Noise2Inverse approach is well-suited for denoising in tomographic imaging and our approach proposes an unsupervised learning strategy for Spectral Imaging based on the Noise2Inverse method. The approach requires us to generate input and target images with independent noise. To achieve this, for a given set of projections acquired over an angular range of 360°, we split the sinogram into K different subsets, $\tilde{y}_{E_1,1}, \dots, \tilde{y}_{E_N,K}$ where each

split contains mutually exclusive projections at equally spaced angles for the same energy channels. After splitting the sinogram, each subset of the sinogram is reconstructed using energy channel-wise FBP, $\tilde{x}_{E_1,1}, \dots, \tilde{x}_{E_N,K}$. For the training step, we generate K different network input images by averaging over all K-1 different subgroups of reconstructions where each subgroup contains K-1 images, whilst the target image is the reconstruction from the set that has not been included in the network input. With this strategy, the input is less noisy than the target. To generate the final denoised image $x_{E_1}^*, \dots, x_{E_N}^*$, all inputs used in training are averaged and used as input of the trained network. The schematic diagram of our approach for $K = 4$ is detailed in Figure 1.

3. EXPERIMENTS

3.1 X-ray source spectrum

To simulate an X-ray source spectrum, we used the SpeckPy software (v2.0).¹⁷ The tube voltage was set to 150 kVp using a tungsten reflection target at an angle of 12 degrees with filtering of 4 mm Aluminum, 1 mm Beryllium and 1000 mm Air. The width of the energy bin [keV] was selected as 0.5 keV and the exposure setting was selected as 1 mAs. To simulate a spectral resolution of 1 keV, we created 131 spectral energy bins between 20 and 150 keV. Specifically, we interpolated the source spectrum between 19550 and 150450 eV with an initial resolution of 0.1 keV before averaging the X-ray flux over 10 neighbouring energy bands. To normalise the X-ray fluence of the source we assumed an X-ray exposure time that would guarantee the detection of 60000 photons for each pixel when summed over all energy channels.

3.2 Synthetic Spectral Data

We created two 2D phantoms (of spatial size 100×100) containing 4 different objects each, with different objects having different materials. We utilized the X-ray DB Python library, which provides attenuation profiles of materials for various elements and compounds, to simulate our phantoms. For the materials we selected in our simulation, we sourced their densities from the PubChem database.¹⁸ Densities were spatially modulated using a sinusoidal function to simulate relative density variations throughout the object. We assigned the X-ray attenuation coefficients of the material to objects (shown in Fig 2) and the background to zero. The choice of the 8 materials used was inspired by the study of the Multi-Spectral dataset.¹⁹ Figure 2 illustrates the visualization of the two phantoms, each showing a distinct energy level (45 keV and 70 keV), highlighting the differentiation in the appearance of the phantoms and their respective objects when observed in different energy channels. Water, olive oil, nitromethane and acetone are selected for the first phantom, and methanol, ethylenediamine, aluminium and nitrobenzene are chosen for the second phantom.

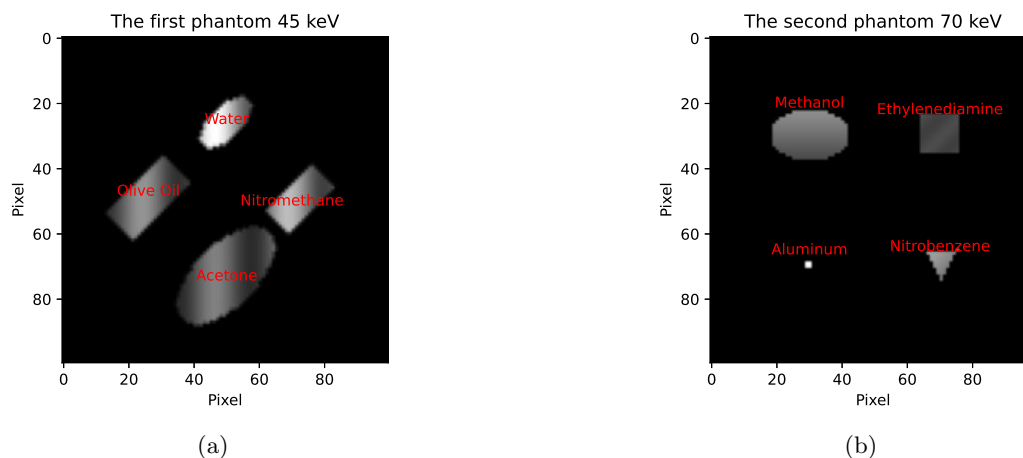


Figure 2: Examples of phantoms at energy levels of 45 keV and 70 keV, featuring different materials and shapes: (a) contains water, olive oil, nitromethane, and acetone; (b) contains methanol, ethylenediamine, aluminium, and nitrobenzene

To generate simulated test data, we use the 2D spectral phantoms and generate 1D sinograms $y_{E_1}, y_{E_2}, \dots, y_{E_N}$ using the geometry described below over the full angular range of 360° with 1° increments and corrupt these with Poisson noise using the source spectrum $I_0(E)$ discussed above. If $p(E)$ is the simulated clean X-ray attenuation value for one pixel, then the noisy pixel $\tilde{p}(E)$ for that energy is distributed as:

$$I_0(E)e^{-\tilde{p}(E)} \sim \text{Poisson} \left(I_0(E)e^{-p(E)} \right). \quad (4)$$

All noisy projections were split into four sets and each of them was reconstructed with the FBP algorithm for our training strategy.

3.3 X-ray Imaging Setup

A linear array detector was simulated with 0.8 mm wide pixels in a 256-pixel array. The scanning geometry used a 57.50 cm distance between the X-ray source and the object and a 58.05 cm distance between the object and the detector, again simulating the setup in.¹⁹

3.4 Comparative reconstruction approaches

For comparison of our method, we also employed a traditional iterative reconstruction method that imposes a Total Variation (TV) constraint as well as the Low2High approach,¹¹ both of which used all 360 noisy projections as inputs. The total variations approach minimises the following cost function

$$x_{\text{reco}} = \arg \min_x \left\{ \frac{1}{2} \|Ax - \tilde{y}\|_2^2 + \alpha \text{TV}(x) \right\}, \quad (5)$$

where the parameter α controls the regularization strength and is chosen empirically for each phantom to optimise denoising performance, which requires knowledge of the clean image, which is not available in real applications. All experiments were performed using the Core Imaging Library (CIL).²⁰

3.5 Network Implementation and Training

We utilized the U-Net architecture from,¹³ implemented using PyTorch, which remains state of the art in many biomedical image-denoising applications. We cropped our original images before starting the training process because having dimensions that are powers of two significantly simplifies various computational processes, especially in deep learning architectures that involve sub-sampling. Each image has 128 energy channels and 96×96 pixels in the spatial domain. Specifically, this training process has been conducted using a block-based approach. Inputs and targets have been divided into blocks of size $4 \times 16 \times 16$, where 4 is the energy dimension. Selected blocks had a 75 % overlap. We trained the model using the Adam optimizer with a learning rate of 10^{-4} using 100 epochs.

3.6 Image Quality Assessment

The quality of the denoised images was assessed against the ground truth phantoms using the structural similarity index (SSIM) and peak signal noise ratio (PSNR) metrics applied channel-wise. We further analyzed the overall image quality by computing the mean and standard deviation of SSIM and PSNR metrics. Additionally, by examining the LACs of various materials across the energy channel, the accuracy of recovering the linear attenuation coefficients (LACs) profiles, which can be used to identify different materials, was assessed.

4. RESULTS

We applied our new Noise2Inverse-based training method to the two phantoms and compared the approach to the Low2High method as well as to a traditional Iterative reconstruction method. The TV constraint inverse problem was solved using the FISTA solver with the optimal parameters for the TV minimization found to be 1 for the first phantom and 0.5 for the second phantom. The iterative method was run for 100 iterations.

Figures 3 and 4 show the denoised images of each phantom at three different energies (55, 85 and 125 keV) for all different methods. Interestingly, in the channel-wise SSIM and PSNR metrics shown in figs 5 and 6, our

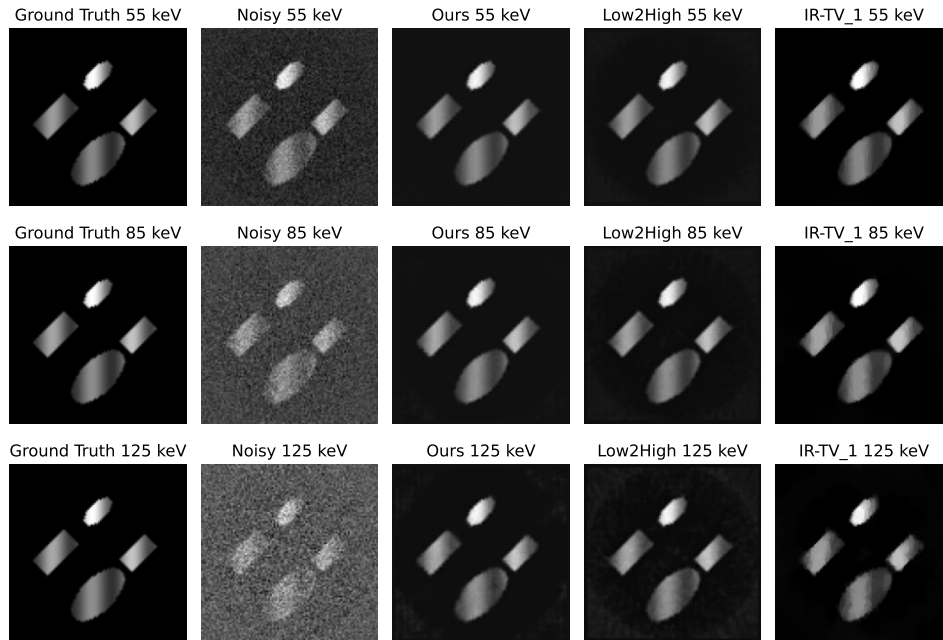


Figure 3: Channel-wise reconstruction for the first phantom. The first column is the ground truth and the second column reconstruction of full noisy projection with FBP. The third column shows our method and the fourth one shows the unsupervised Low2High method. The last column represents the iterative reconstruction method, and 1 here indicates the alpha value that is selected for TV minimization.

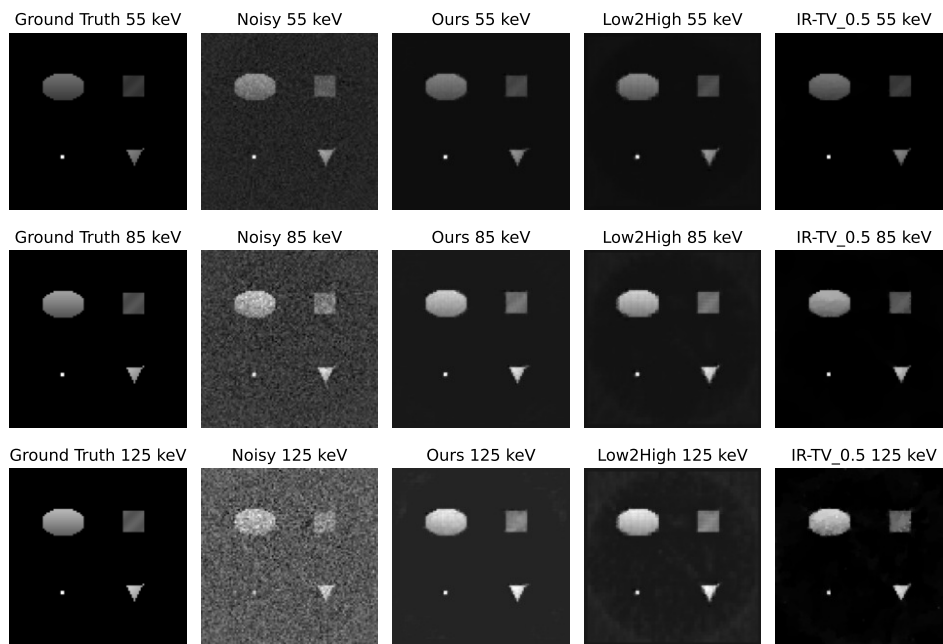


Figure 4: Channel-wise reconstruction for the second phantom. The first column is the ground truth and the second column reconstruction of the full noisy projection with FBP. The third column shows our method and fourth one shows the unsupervised Low2High method. The last column represents the iterative reconstruction method, and 0.5 here indicates the alpha value that is selected for TV minimization.

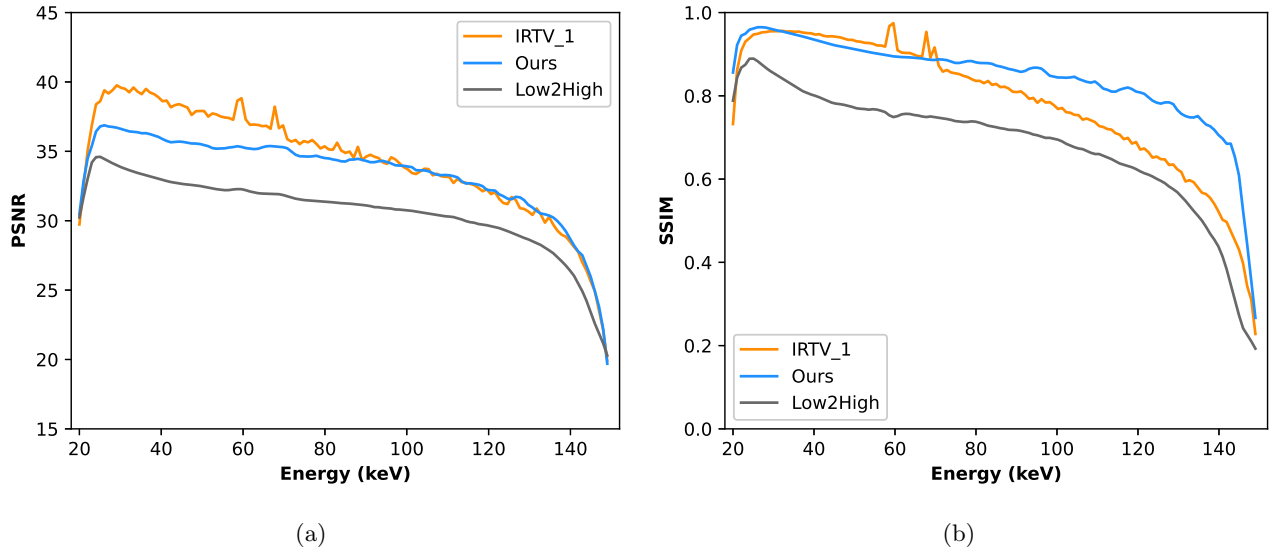


Figure 5: Comparative analysis of a) Channel-wise PSNR for the first phantom, and b) Channel-wise SSIM for the first phantom.

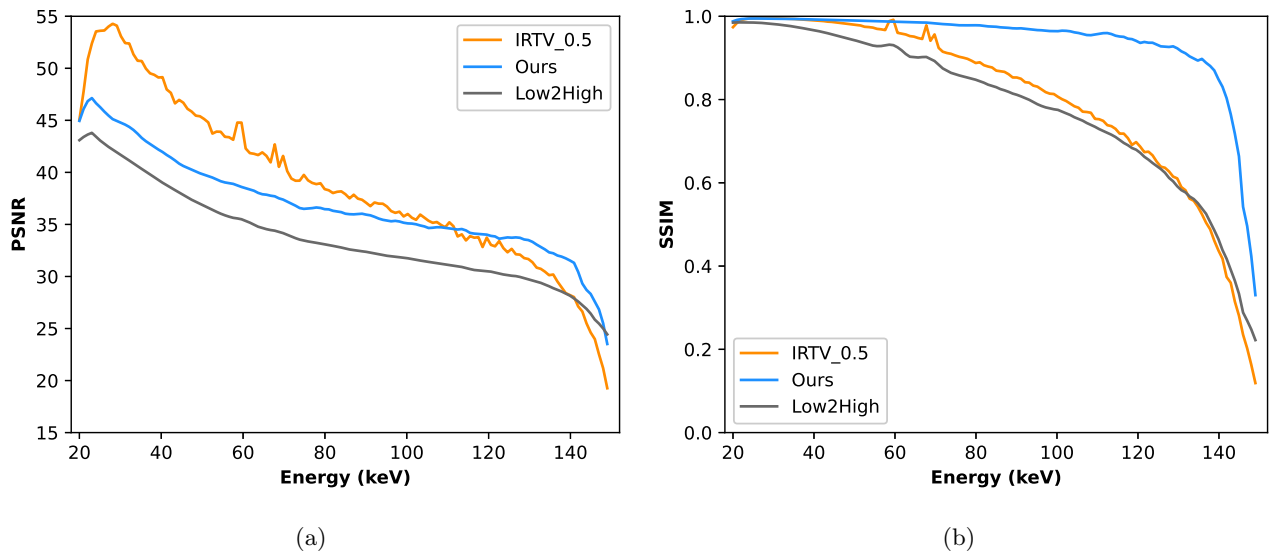


Figure 6: Comparative analysis of a) Channel-wise PSNR for the second phantom, and b) Channel-wise SSIM for the second phantom.

Table 1: SSIM and PSNR of different methods for the first phantom (mean \pm SD) .

Method	SSIM	PSNR (dB)
IRTV_1	0.79 ± 0.16	34.4 ± 3.9
Ours	0.84 ± 0.11	33.5 ± 3.0
Low2High	0.69 ± 0.15	30.7 ± 2.7

Table 2: SSIM and PSNR of different methods for the second phantom (mean \pm SD).

Method	SSIM	PSNR (dB)
IRTV_0.5	0.81 ± 0.21	39.0 ± 7.9
Ours	0.94 ± 0.10	36.9 ± 4.6
Low2High	0.79 ± 0.18	33.8 ± 4.5

method had a better performance, especially for high noise (i.e. low photon count) energy channels (the low and high energy channels, where the source spectrum has limited flux), though the average PSNR performance was found to be still better for the iterative method.

The results (tables 1 and 2) indicate that our method exceeds the performance of the Low2High and traditional iterative methods for each phantom, as measured in average SSIM values across the energy channel. In addition, our method demonstrated comparable performance in terms of average PSNR for the first phantom; however, for the second phantom, the average PSNR was inferior to that achieved by the traditional iterative reconstruction method. To evaluate the denoising performance across the energy channel, pixels were randomly selected within the objects of interest, and the attenuation profiles for two materials, one from each phantom, were compared with the ground truth as illustrated in Figures 7 and 8. The LAC values were better preserved over the energy channels with unsupervised methods (Low2High and our method) compared to the IR-TV method in terms of noise. The LAC profile of the Low2High method is slightly lower than our method compared with the ground truth, this shrinking comes from the use of a filter in the reconstruction part. Notably, our method exhibits a closer alignment with the ground truth over the energy channels. Despite the presence of noise within the LAC profile, the IR-TV method for each phantom yielded superior average PSNR outcomes. This is because wherein a singular pixel was analyzed across the energy channel for the LAC profile, but PSNR averaged over the energy channel.

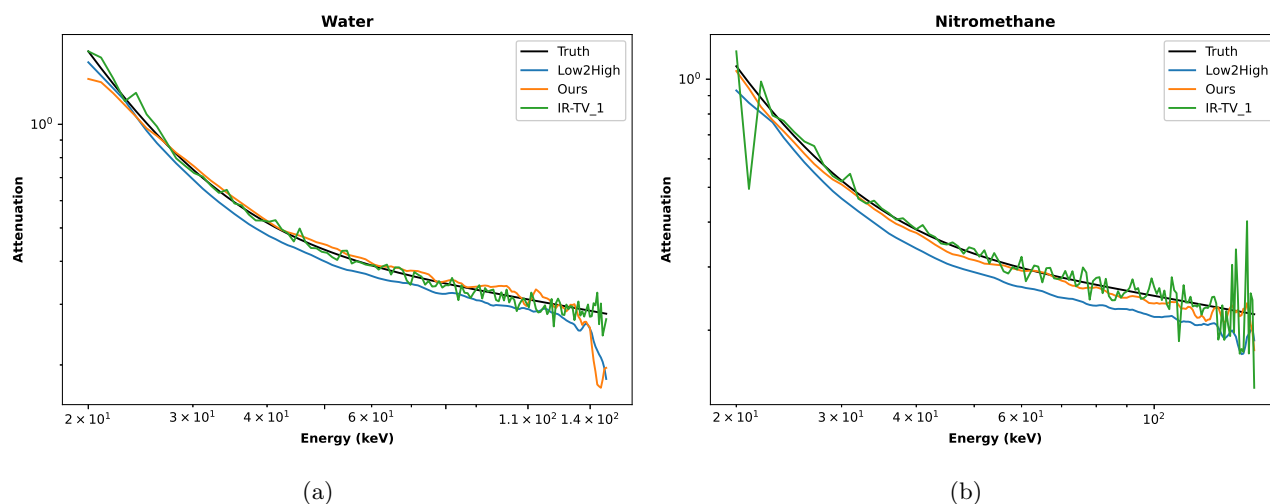


Figure 7: Examples of the linear attenuation coefficient of different materials over the energy channels in the first phantom. Axes are shown in log scale. a) Water, b) Nitromethane.

5. DISCUSSION AND CONCLUSION

The employment of energy information in Spectral imaging is significant as it enables the decomposition of materials with similar attenuation properties and enhances the accuracy of material decomposition. However, the projection data is intrinsically noisy because there aren't many photons within each energy channel. Here, to address the difficulty of collecting clean data, we studied the feasibility of a learning-based denoising approach that does not require additional clean and noisy training data. We were able to demonstrate that the image quality across spectral channels is preserved in two different spectral phantoms without the necessity for parameter adjustment. The U-Net denoises the FBP reconstruction using a noisy sinogram splitting strategy. The self-supervised U-Net was robust to the varying noise levels of the different energy channels, especially the first and last energy channels.

Whilst the traditional total variation-constrained reconstruction was also found to perform similarly well or even slightly better than our approach, that was only achieved with significant and time-consuming parameter tuning. As the optimal parameter is highly sensitive to the image structure, this approach would not be feasible

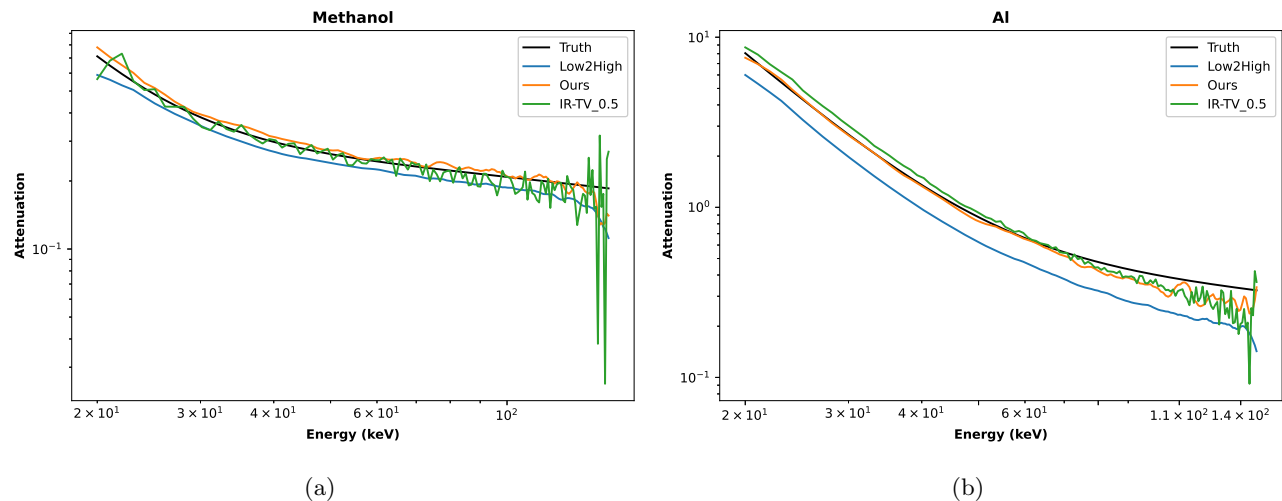


Figure 8: Examples of the linear attenuation coefficient of different materials over the energy channels in the second phantom. Axes are shown in log scale. a) Methanol, b) Aluminium

without the knowledge of the ground truth image and is thus not easily applicable in real applications. Furthermore, it is well known that TV regularisation leads to biased results, which can introduce further errors in the estimated spectra, which, for quantitative applications, can lead to unacceptable errors.

In this study, we have demonstrated the application of Noise2Inverse to spectral imaging using a block-based training approach with a U-Net, showing that Noise2Inverse in spectral imaging does offer a significant improvement in image quality. This was achieved without needing to fine-tune regularisation parameters which is a drawback of traditional iterative approaches. Future work will focus on applying the learned spectral denoising techniques to data with three spatial dimensions and validating them with real experimental data.

REFERENCES

- [1] Garnett, R., “A comprehensive review of dual-energy and multi-spectral computed tomography,” *Clinical Imaging* **67**, 160–169 (2020).
- [2] Rajendran, K., Petersilka, M., Henning, A., Shanblatt, E. R., Schmidt, B., Flohr, T. G., Ferrero, A., Baffour, F., Diehn, F. E., Yu, L., et al., “First clinical photon-counting detector ct system: technical evaluation,” *Radiology* **303**(1), 130–138 (2022).
- [3] Kehl, C., Mustafa, W., Kehres, J., Dahl, A. B., Olsen, U. L., and DTU, D. T. U., “Distinguishing malicious fluids in luggage via multi-spectral ct reconstructions,” *3D-NordOst, GFAI-Gesellschaft zur Förderung angewandter Informatik eV, Berlin, Germany* (2018).
- [4] Martin, L., Tuysuzoglu, A., Karl, W. C., and Ishwar, P., “Learning-based object identification and segmentation using dual-energy ct images for security,” *IEEE Transactions on Image Processing* **24**(11), 4069–4081 (2015).
- [5] Richtsmeier, D., Guliyev, E., Iniewski, K., and Bazalova-Carter, M., “Contaminant detection in non-destructive testing using a czt photon-counting detector,” *JINST16 P* **1011** (2021).
- [6] Schumacher, D., Zscherpel, U., and Ewert, U., “Photon counting and energy discriminating x-ray detectors-benefits and applications,” in *[19th World Conference on Non-Destructive Testing, 2016, Proceedings]*, **2016**, Tu-2 (2016).
- [7] Wu, W., Chen, P., Wang, S., Vardhanabhuti, V., Liu, F., and Yu, H., “Image-domain material decomposition for spectral ct using a generalized dictionary learning,” *IEEE transactions on radiation and plasma medical sciences* **5**(4), 537–547 (2020).
- [8] Xi, Y., Chen, Y., Tang, R., Sun, J., and Zhao, J., “United iterative reconstruction for spectral computed tomography,” *IEEE transactions on medical imaging* **34**(3), 769–778 (2014).

- [9] Rigie, D. S. and La Riviere, P. J., “Joint reconstruction of multi-channel, spectral ct data via constrained total nuclear variation minimization,” *Physics in Medicine & Biology* **60**(5), 1741 (2015).
- [10] Wu, W., Hu, D., Niu, C., Broeke, L. V., Butler, A. P., Cao, P., Atlas, J., Chernoglazov, A., Vardhanabhuti, V., and Wang, G., “Deep learning based spectral ct imaging,” *Neural Networks* **144**, 342–358 (2021).
- [11] Inkinen, S. I., Brix, M. A., Nieminen, M. T., Arridge, S., and Hauptmann, A., “Unsupervised denoising for sparse multi-spectral computed tomography,” *arXiv preprint arXiv:2211.01159* (2022).
- [12] Hendriksen, A. A., Pelt, D. M., and Batenburg, K. J., “Noise2inverse: Self-supervised deep convolutional denoising for tomography,” *IEEE Transactions on Computational Imaging* **6**, 1320–1335 (2020).
- [13] Ronneberger, O., Fischer, P., and Brox, T., “U-net: Convolutional networks for biomedical image segmentation,” in [*Medical Image Computing and Computer-Assisted Intervention–MICCAI 2015: 18th International Conference, Munich, Germany, October 5–9, 2015, Proceedings, Part III 18*], 234–241, Springer (2015).
- [14] Hendriksen, A. A., Bühner, M., Leone, L., Merlini, M., Vigano, N., Pelt, D. M., Marone, F., Di Michiel, M., and Batenburg, K. J., “Deep denoising for multi-dimensional synchrotron x-ray tomography without high-quality reference data,” *Scientific reports* **11**(1), 11895 (2021).
- [15] Lehtinen, J., Munkberg, J., Hasselgren, J., Laine, S., Karras, T., Aittala, M., and Aila, T., “Noise2noise: Learning image restoration without clean data,” *arXiv preprint arXiv:1803.04189* (2018).
- [16] Batson, J. and Royer, L., “Noise2self: Blind denoising by self-supervision,” in [*International Conference on Machine Learning*], 524–533, PMLR (2019).
- [17] Bujila, R., Omar, A., and Poludniowski, G., “A validation of spekpy: A software toolkit for modelling x-ray tube spectra,” *Physica Medica* **75**, 44–54 (2020).
- [18] Kim, S., Chen, J., Cheng, T., Gindulyte, A., He, J., He, S., Li, Q., Shoemaker, B. A., Thiessen, P. A., Yu, B., et al., “Pubchem 2023 update,” *Nucleic acids research* **51**(D1), D1373–D1380 (2023).
- [19] Kehl, C., Mustafa, W., Kehres, J., Dahl, A. B., and Olsen, U. L., “Multi-spectral imaging via computed tomography (music)-comparing unsupervised spectral segmentations for material differentiation,” *arXiv preprint arXiv:1810.11823* (2018).
- [20] Jørgensen, J. S., Ametova, E., Burca, G., Fardell, G., Papoutsellis, E., Pasca, E., Thielemans, K., Turner, M., Warr, R., Lionheart, W. R., et al., “Core imaging library-part i: a versatile python framework for tomographic imaging,” *Philosophical Transactions of the Royal Society A* **379**(2204), 20200192 (2021).