

Riesz-Quincunx-Unet Variational Auto-Encoder for Satellite Image Denoising

Duy H. Thai¹, Xiqi Fei¹, Minh Tri Le¹, Andreas Züfle², Konrad Wessels¹

¹George Mason University, Department of Geography and Geoinformation Science, USA

²Emory University, Department of Computer Science, USA

Abstract—Multiresolution deep learning approaches, such as the U-Net architecture, have achieved high performance in classifying and segmenting images. However, these approaches do not provide a latent image representation and cannot be used to decompose, denoise, and reconstruct image data. The U-Net and other convolutional neural network (CNNs) architectures commonly use pooling to enlarge the receptive field, which usually results in irreversible information loss. This study proposes to include a Riesz-Quincunx (RQ) wavelet transform, which combines 1) higher-order Riesz wavelet transform and 2) orthogonal Quincunx wavelets (which have both been used to reduce blur in medical images) inside the U-net architecture, to reduce noise in satellite images and their time-series. In the transformed feature space, we propose a variational approach to understand how random perturbations of the features affect the image to further reduce noise. Combining both approaches, we introduce a hybrid RQUNet-VAE scheme for image and time series decomposition used to reduce noise in satellite imagery. We present qualitative and quantitative experimental results that demonstrate that our proposed RQUNet-VAE was more effective at reducing noise in satellite imagery compared to other state-of-the-art methods. We also apply our scheme to several applications for multi-band satellite images, including: image denoising, image and time-series decomposition by diffusion and image segmentation.

Keywords: Quincunx wavelet, high order Riesz transform, image time series decomposition, variational auto-encoder, deep neural networks, Sentinel-2, Unet

I. INTRODUCTION

The temporal frequency of medium resolution, optical satellite imagery, such as Landsat 8 and 9 and Sentinel2A&B, has increased significantly in the past four years from one observation every 16 days with Landsat 8 to an observation every 2.9 days on average [16], [25]. Such time series of multi-spectral satellite data enables many novel applications such as global agriculture monitoring and land cover change at the appropriate resolution of land management impacts that has not been possible with low resolution (250-500m GSD) time-series such as MODIS [14], [28], [37]. Moreover, the harmonized Landsat Sentinel2 products open up new avenues for real-time monitoring of a wider variety of change phenomenon [5]. There is a wide variety of time series analysis methods for change detection, which focus primarily on the temporal domain while the spatial domain has been largely

neglected (for review see [46]). The availability of hyper-temporal medium resolution imagery allows new application in the spatial domain, such as semantic segmentation with UNET [24] to track objects of change through time.

Change detection methods are hampered by significant noise in the time series that remains despite various processing efforts to reduce the noise [46]. This noise is firstly caused by clouds, and cloud shadows that result in data gaps even when they are correctly detected and serious change artifacts are caused when some clouds or their edges remain undetected [26]. Second, atmospheric variability, notably water vapor and aerosols that cause variability in top of atmosphere reflectance, despite best efforts at atmospheric correction [39], [45]. Third, BRDF variation due to variation in sun-sensor geometry and non-Lambertian reflectance properties of the target. Fourth, variations in spectral reflectance due to seasonal vegetation phenology, which become hard to model when clouds result in missing data that makes the time-series irregular and unpredictable. These issues are often addressed by multi-date composites or temporal interpolation to fill in clouds, shadows and other missing data. Despite efforts to reduce the impact of spatio-temporal noise in the satellite time series, it often limits the accuracy of timely land cover change detection over very large areas using either conventional time-series analysis [42], as well as the new generation of machine learning methods [46].

Various cutting-edge techniques have been developed for image decomposition that can be used to remove noise from conventional images or videos, for example, wavelet smoothing techniques [36], [8], [2] and regularization methods [27], but these methods require parameter selection that vary greatly between datasets and cannot optimally represent varying signals over space and time commonly found in earth observation images. On the other hand, a neural network (UNet) can automatically learn optimal local representation of signal and noise [24], [32], but incurs high computational cost and requires prohibitively large sets of domain-specific training data. The approach proposed in this paper combines the strengths of convolutional neural networks and conventional smoothing techniques. For this purpose we introduce a novel non-subsampled high order Riesz-Quincunx wavelet with variational auto-encoder UNet (RQUNet-VAE) as a hybrid combination of deterministic wavelet expansion (two-dimensional Riesz transformation [36], Quincunx wavelet) and a variational version of a convolutional neural network [22]. The proposed

RQUNet-VAE approach uses framelet decomposition [43] to map an image into a sparse feature space and leverages UNet [24] to enable learning of the frames of the feature space. The rationale is that such a hybrid method should provide better feature representation and artifact reduction compared to conventional approaches. Therefore, while previous studies used UNet-VAE for image segmentation [15], we use it to mimic the properties of latent factor model in our RQUNet-VAE decomposition.

To implement this concept, we first introduce a non-subsampled version of high-order Riesz wavelet expansion [36], [34], [35] with Quincunx sampling [8]. Quincunx sampling is used to reduce redundancy of an expansion while high order Riesz transform is used to increase the directional property of wavelet expansion. The hybrid model reduces computational cost with deterministic bases as predefined parameters instead of letting all bases be learnable parameters.

Next, this Riesz-Quincunx wavelet is integrated into the skip-connections of the deep neural network UNet-VAE [13], [15], [7] for learning new bases from the training dataset. Our rationale of using wavelet expansion is that signals extracted from the UNet-VAE encoder, at the skip-connection level, contains both the main signal and details (as well as noise) of an input image which are separated into scaling and wavelet coefficients. Truncation of these coefficients eliminates small wavelet coefficients which contain noise and detailed texture. By decoding the remaining coefficients back into image space, we obtain a denoised version of the original image.

The theoretical framework of RQUNet-VAE is based on Hankel matrix algebra [20], framelet decomposition [43], [44] and proximal operators [21]. Framelet decomposition uses isotropic family-matrix convolution to combine all channels of a multi-band image in a convolutional operation with learnable frames. Furthermore, there is a connection between framelet decomposition and sampling with a finite rate of innovation [40] via Hankel matrix theory and annihilating filter. We prove that RQUNet-VAE also relates to latent factor model whose loss function is defined by Kullback-Leibler divergence [38].

Finally we demonstrate how to apply our proposed RQUNet-VAE to satellite image time series, that is, sequences of images of the same area. Reducing noise satellite image time series is challenging due to the severe background noise resulting from spectral variability caused by changing environmental conditions due to atmospheric and seasonal variability, remaining small clouds, as well as variable sun-sensor geometry [25], [45]. The level of noise is very different from that of conventional videos which contains slow motion changes between subsequent images that can be removed with time delay embedding, for example. In contrast, the satellite time series are constituted by discrete frames that are independently capture several days apart, causing large variability in reflectance properties of the background and objects.

To test the effectiveness of our proposed concept, the RQUNet-VAE was applied to image decomposition, image

denoising and segmentation of satellite images and their time-series. Our experimental results show that our hybrid method provides better feature representation and artifact reduction than traditional approaches. The objectives of this paper are:

- 1) introduce RQUNet-VAE as a generalized wavelet expansion approach;
- 2) extend RQUNet-VAE to a diffusion process [23] and enable spectral decomposition [9];
- 3) apply it to image denoising and segmentation in noisy environment for multi-band satellite images and their time-series.

Organization of the paper is: Section II provides the RQUNet-VAE expansion, including mathematical properties and image or time-series decomposition; Section III gives numerical examples and comparisons of image denoising and segmentation in noisy environments for multi-band satellites images. Section IV gives the conclusions, Mathematical background, proofs, and additional experiments are provided in our supplemental material (SM).

II. RIESZ-QUINCUNX-UNET VARIATIONAL AUTO-ENCODER (RQUNET-VAE)

A. Notations and Definitions:

The following list provides an overview of notations used throughout this work.

- continuous coordinate in the spatial domain: $x = (x_1, x_2) \in \mathbb{R}^2$,
- discrete coordinate in the spatial domain: $k = (k_1, k_2), m = (m_1, m_2) \in \mathbb{Z}^2$,
- Fourier coordinate: $\omega = (\omega_1, \omega_2) \in [-\pi, \pi]^2$,
- Complex number: $j = \sqrt{-1}$,
- Image domain: $\Omega = \{1, \dots, n_1\} \times \{1, \dots, n_2\}, |\Omega| = n_1 \times n_2$,
- Images and stacks of images:
 - a gray-scale image $\underline{f} = (f_1 \dots f_{n_2}) \in \mathbb{R}^{|\Omega|}, f_i \in \mathbb{R}^{n_1}$,
 - a multi-channel image $\underline{\underline{f}} = \{\underline{f}_1, \dots, \underline{f}_P\} \in \mathbb{R}^{|\Omega| \times P}, \underline{f}_i \in \mathbb{R}^{|\Omega|}$,
 - a set of observed images $\mathfrak{F} := \underline{\underline{f}} = \{\underline{\underline{f}}_i\}_{i=1}^T \in \mathbb{R}^{|\Omega| \times P \times T}$,
- Distribution and Lebesgue density: $\underline{f} \stackrel{\text{i.i.d.}}{\sim} \mathbb{P}(\underline{f}) := \mathbb{P}(\underline{\underline{df}}) = \mathbb{p}(\underline{f}) \underline{\underline{df}}$,
- Continuous Fourier transform of a continuous function $a \in L_2(\mathbb{R}^2)$ is:

$$a(x) \xleftrightarrow{\mathcal{F}} \hat{a}(\omega) = \int_{\mathbb{R}^2} a(x) e^{-j\langle x, \omega \rangle_{\ell_2}} dx,$$

and its discrete version is computed via Poisson summation formulae:

$$a[k] := a(x) |_{x=k \in \mathbb{Z}^2}$$

$$\xleftrightarrow{\mathcal{F}} \hat{A}(e^{j\omega}) = \sum_{k \in \mathbb{Z}^2} a[k] e^{-j\langle k, \omega \rangle_{\ell_2}} = \sum_{k \in \mathbb{Z}^2} \hat{a}(2\pi k + \omega).$$

Additional definitions can be found in the supplemental material.

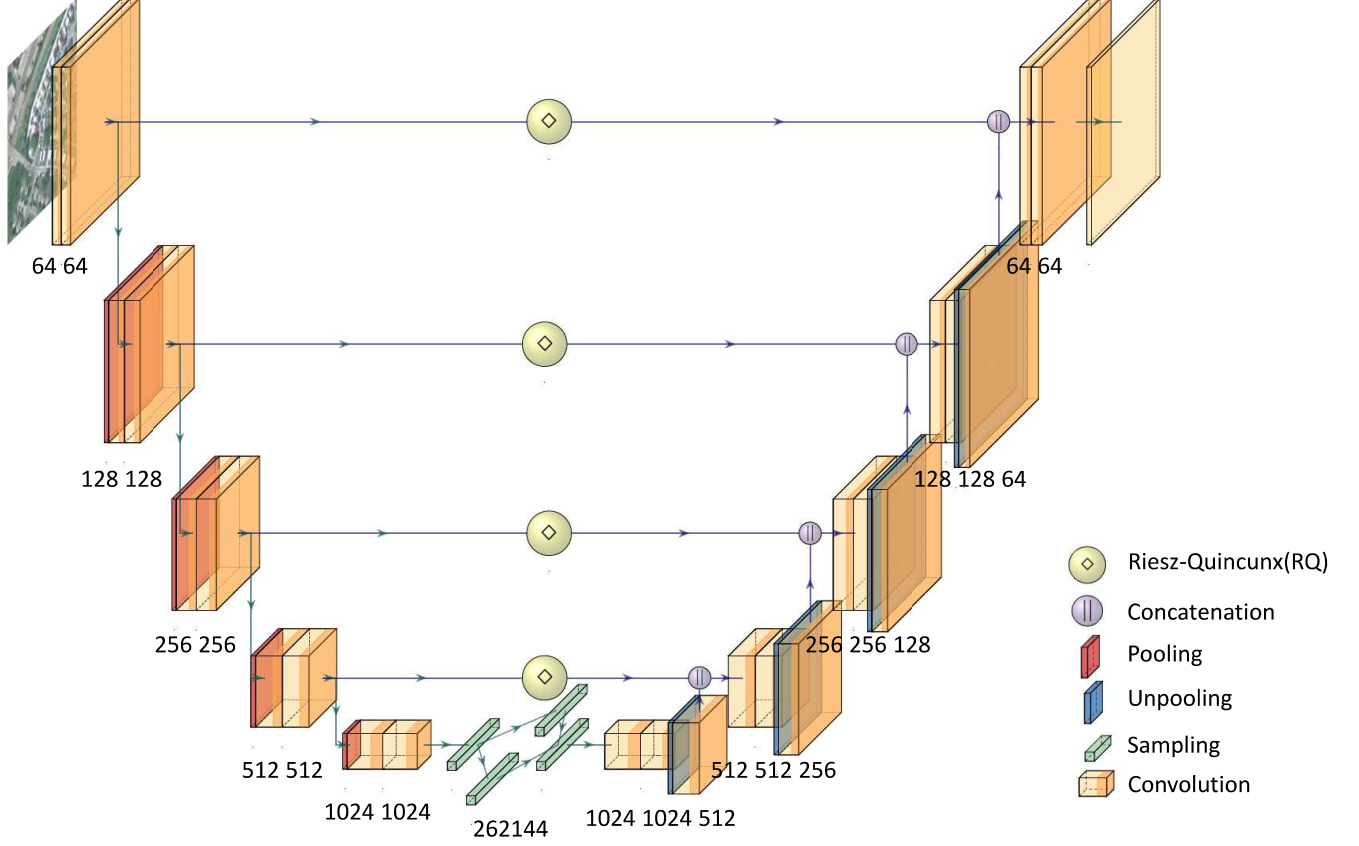


Fig. 1. RQUNet-VAE architecture.

B. RQUNet-VAE Architecture Overview

The network architecture of our proposed RQUNet-VAE is illustrated in Figure 1. The network utilizes a UNet [24] as our primary backbone architecture with the modified skip-connection signals, between the encoder and decoder paths, and bottom layer. The network introduces the Variational Autoencoder (VAE) in the bottom layer to learn the signal distribution in the latent space during training on large datasets. The VAE layer combined with the UNet backbone allows the model to learn and retain much of the input information, similar to generative models, to perform image reconstruction. The reconstructed image then can go through a convolution layer, or a classifier, to produce the final classification or segmentation output. The primary objective of our proposed network is to remove certain levels of noise in the input images in order to produce improved, quality reconstructed images using the Riesz-Quincunx (RQ) scheme in the skip-connections, before performing image classification. Since the RQ scheme is computational heavy, we only applied the RQ computation on a pre-trained network to produce predictions on input images.

C. N -th order Riesz Quincunx non-subsampled wavelet

Before proposing our RQUNet-VAE expansion, we firstly introduce framelet decomposition [43] in the following proposition.

Proposition II.1. Given an image $\underline{f} \in \mathbb{R}^{n_1 \times n_2}$, all wavelet filter banks $\underline{\Phi}, \tilde{\underline{\Phi}} \in \mathbb{R}^{n_2 d_1 \times d_2}$ (see definitions of wavelet filter banks in Section 4 of the supplemental material) and local basis $\underline{\Xi}, \tilde{\underline{\Xi}} \in \mathbb{R}^{n_1 \times d}$ (see Equation 5 in the supplemental material for details) satisfy the unity conditions

$$\tilde{\underline{\Xi}} \underline{\Xi}^T = \sum_{i=1}^d \tilde{\xi}_i \xi_i^T = Id_{n_1 \times n_1}, \quad \underline{\Phi} \tilde{\underline{\Phi}}^T = \sum_{i=1}^{d_2} \tilde{\phi}_i \phi_i^T = Id_{n_2 d_1 \times n_2 d_1} \quad (1)$$

then, a framelet decomposition is:

$$\begin{aligned} \underline{f} &= \mathcal{H}_{d_1|n_2}^\dagger \left(\tilde{\underline{\Xi}} \underline{c}_f \tilde{\underline{\Phi}}^T \right) \\ &= \frac{1}{d_1} \sum_{s=1}^{d_2} \sum_{l=1}^d \sum_{i=1}^{n_2} \begin{pmatrix} \langle f_k, \mathbf{e}_{\phi_s^i}(\xi_l) \rangle_{\ell_2} & \mathbf{e}_{\tilde{\phi}_s^1}(\tilde{\xi}_l) \\ \vdots & \vdots \\ \langle f_k, \mathbf{e}_{\phi_s^i}(\xi_l) \rangle_{\ell_2} & \mathbf{e}_{\tilde{\phi}_s^{n_2}}(\tilde{\xi}_l) \end{pmatrix} \\ &= \frac{1}{d_1} \sum_{s=1}^{d_2} \left(\mathbf{e}_{\tilde{\phi}_s^1}(\tilde{\underline{\Xi}} \underline{c}_{f,s}) \quad \dots \quad \mathbf{e}_{\tilde{\phi}_s^{n_2}}(\tilde{\underline{\Xi}} \underline{c}_{f,s}) \right) \end{aligned} \quad (2)$$

where $\mathcal{H}_{d_1|n_2}^\dagger$ is the extended Hankel matrix of image \underline{f} as described in the SM. Framelet coefficients are:

$$\underline{c}_f := (c_{f,1} \quad \dots \quad c_{f,d_2}) = \underline{\Xi}^T \mathcal{H}_{d_1|n_2} \left(\underline{f} \right) \underline{\Phi}. \quad (3)$$

Proof. For self-containment, we provide a proof of Proposition II.1 in Section 5.1. in the SM. \square

Inspired by Proposition II.1, we introduce the proposed RQUNet-VAE expansion in the following.

a) *Riesz-Quincunx wavelet expansion*:: To have wavelet expansion, firstly we provide a definition of frame by isotropic polyharmonic Bspline and N -th order Riesz transform, see [36]. In particular, an isotropic polyharmonic B-splines basis is defined in the Fourier domain as: $x \in \mathbb{R}^2$,

$$\beta_\gamma(x) \xleftrightarrow{\mathcal{F}} \hat{\beta}_\gamma(\omega) = \frac{\hat{V}^{\text{iso}}(e^{j\omega})^{\frac{\gamma}{2}}}{\|\omega\|_{\ell_2}^\gamma} \quad (4)$$

with a 2D localization operator:

$$\begin{aligned} \hat{V}^{\text{iso}}(e^{j\omega}) &= \frac{10}{3} \\ &- \frac{1}{3} [4 \cos \omega_1 + 4 \cos \omega_2 + \cos(\omega_1 + \omega_2) + \cos(\omega_1 - \omega_2)] . \end{aligned}$$

Its dual function is defined via an auto-correlation function: $k \in \mathbb{Z}^2$,

$$\begin{aligned} \tilde{\beta}_\gamma(x) &\xleftrightarrow{\mathcal{F}} \hat{\tilde{\beta}}_\gamma(\omega) = \frac{\hat{\beta}_\gamma(\omega)}{\hat{A}(e^{j\omega})}, \\ a(k) &\xleftrightarrow{\mathcal{F}} \hat{A}(e^{j\omega}) = \sum_{m \in \mathbb{Z}^2} \hat{\beta}_{2\gamma}(2\pi m + \omega). \end{aligned} \quad (5)$$

An impulse response of the L -th order Riesz transform is:

$$\mathcal{R}^l\{\delta\}(x) \xleftrightarrow{\mathcal{F}} \hat{\mathcal{R}}^l(\omega) = (-j)^L \sqrt{\frac{L!}{n!(L-l)!}} \frac{\omega_1^l \omega_2^{L-l}}{(\omega_1^2 + \omega_2^2)^{\frac{L}{2}}},$$

for $l = 0, \dots, L$ and where $\delta(\cdot)$ is the Dirac delta function. Secondly, wavelet expansion is defined as follows: Given a 2D dyadic sampling matrix $\underline{D} = \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix}$, we have $\hat{\rho}_0 = [0, 0]^T$, $\hat{\rho}_1 = [1, 0]^T$ and

$$\underline{D}^i = \begin{cases} 2^{\frac{i}{2}} \text{Id}_2, & i \text{ is even} \\ 2^{\frac{i-1}{2}} \underline{D}, & i \text{ is odd} \end{cases}, \quad |\det \underline{D}^i| = 2^i.$$

Non-subsampled scaling and wavelet spaces $\mathcal{V}_i^{\text{ns}} = \text{span}_{m \in \mathbb{Z}^2} \{\varphi_i(\cdot - m)\}$ and $\mathcal{W}_{il}^{\text{ns}} = \text{span}_{m \in \mathbb{Z}^2} \{\psi_{il}(\cdot - m)\}$ (for $i = 0, \dots, I; l = 0, \dots, L$) satisfy a multiscale decomposition $\mathcal{V}_{i-1}^{\text{ns}} = \mathcal{V}_i^{\text{ns}} \oplus \sum_{l=0}^L \mathcal{W}_{il}^{\text{ns}}$, $i = 1, \dots, I$ where primal/dual scaling and wavelet functions are:

$$\begin{aligned} \varphi_i(x) &= 2^{-\frac{i}{2}} \beta_\gamma(\underline{D}^{-i} x), \quad \tilde{\varphi}_i(x) = 2^{-\frac{i}{2}} \tilde{\beta}_\gamma(\underline{D}^{-i} x), \\ \psi_{il}(x) &= 2^{-\frac{i}{2}} \mathcal{R}^l\{\psi\}(\underline{D}^{-i} x), \quad \tilde{\psi}_{il}(x) = 2^{-\frac{i}{2}} \mathcal{R}^l\{\tilde{\psi}\}(\underline{D}^{-i} x). \end{aligned}$$

Their Fourier transforms are

$$\hat{\varphi}_i(\omega) = 2^{\frac{i}{2}} \hat{\beta}_\gamma(\underline{D}^{iT} \omega), \quad \hat{\tilde{\varphi}}_i(\omega) = 2^{\frac{i}{2}} \hat{\tilde{\beta}}_\gamma(\underline{D}^{iT} \omega),$$

$$\hat{\psi}_{il}(\omega) = 2^{\frac{i}{2}} \hat{\mathcal{R}}^l(\underline{D}^{iT} \omega) \hat{\psi}(\underline{D}^{iT} \omega),$$

$$\hat{\tilde{\psi}}_{il}(\omega) = 2^{\frac{i}{2}} \hat{\mathcal{R}}^l(\underline{D}^{iT} \omega) \hat{\tilde{\psi}}(\underline{D}^{iT} \omega).$$

Due to a discrete Fourier transform of continuous functions, by Poisson summation we have the following proposition:

Proposition II.2. *The scaling and wavelet functions satisfy the unity condition in the Fourier domain:*

$$\hat{\varphi}_I^*(\omega) \hat{\varphi}_I(\omega) + \sum_{i=0}^I \sum_{l=0}^L \hat{\psi}_{il}^*(\omega) \hat{\psi}_{il}(\omega) + \hat{e}(\omega) = 1 \quad (6)$$

up to a discretization error:

$$\begin{aligned} \hat{e}(\omega) &= \hat{\varphi}_I^*(\omega) \sum_{k \in \mathbb{Z}^2 \setminus \{0\}} \hat{\varphi}_I(2\pi k + \omega) \\ &+ \sum_{m \in \mathbb{Z}^2 \setminus \{0\}} \hat{\varphi}_I^*(2\pi m + \omega) \left(\hat{\varphi}_I(\omega) + \sum_{k \in \mathbb{Z}^2 \setminus \{0\}} \hat{\varphi}_I(2\pi k + \omega) \right) \\ &+ \sum_{i=0}^I \sum_{l=0}^L \left[\hat{\psi}_{il}^*(\omega) \sum_{k \in \mathbb{Z}^2 \setminus \{0\}} \hat{\psi}_{il}(2\pi k + \omega) \right. \\ &\left. + \sum_{m \in \mathbb{Z}^2 \setminus \{0\}} \hat{\psi}_{il}^*(2\pi m + \omega) \left(\hat{\psi}_{il}(\omega) + \sum_{k \in \mathbb{Z}^2 \setminus \{0\}} \hat{\psi}_{il}(2\pi k + \omega) \right) \right]. \end{aligned} \quad (7)$$

Proof. We provide a proof of Proposition II.2 in Section 5.2 in SM. \square

To compensate the error in Equation 7 in the unity condition (Equation 6), wavelet function at scale $i = 0$ is defined as:

$$\hat{\psi}_0(\omega) = \frac{1}{\hat{\psi}_0^*(\omega)} \left(1 - \hat{\varphi}_I^*(\omega) \hat{\varphi}_I(\omega) - \sum_{i=1}^I \hat{\psi}_i^*(\omega) \hat{\psi}_i(\omega) \right);$$

then, we have a wavelet expansion for an image $f \in \ell_2(\mathbb{Z}^2)$:

$$\begin{aligned} f[k] &= \sum_{m \in \mathbb{Z}^2} \langle f, \tilde{\varphi}_I(\cdot - m) \rangle_{\ell_2} \varphi_I(k - m) \\ &+ \sum_{i=0}^I \sum_{l=0}^L \sum_{m \in \mathbb{Z}^2} \langle f, \tilde{\psi}_{il}(\cdot - m) \rangle_{\ell_2} \psi_{il}(k - m). \end{aligned} \quad (8)$$

Following [34], primal and dual wavelet functions are defined as:

$$\begin{aligned} \hat{\psi}(\omega) &= 2^{-\frac{1}{2}} \hat{G}(e^{j\underline{D}^{-T}\omega}) \hat{\beta}_\gamma(\underline{D}^{-T}\omega), \\ \hat{\tilde{\psi}}(\omega) &= 2^{-\frac{1}{2}} \hat{\tilde{G}}(e^{j\underline{D}^{-T}\omega}) \hat{\tilde{\beta}}_\gamma(\underline{D}^{-T}\omega), \end{aligned} \quad (9)$$

where the refinement and highpass filters are

$$\begin{aligned} h[k] &\xleftrightarrow{\mathcal{F}} \widehat{H}(e^{j\omega}) = 2^{\frac{1}{2}} \frac{\widehat{\beta}_\gamma(\mathbf{D}^T \omega)}{\widehat{\beta}_\gamma(\omega)}, \\ \tilde{h}[k] &\xleftrightarrow{\mathcal{F}} \widehat{\tilde{H}}(e^{j\omega}) = \frac{\widehat{A}(e^{j\omega})}{\widehat{A}(e^{j\mathbf{D}^T \omega})} \widehat{L}(e^{j\omega}), \\ g[k] &\xleftrightarrow{\mathcal{F}} \widehat{G}(e^{j\omega}) = -e^{-j\omega_1} \widehat{H}(e^{-j(\omega+\pi)}) \widehat{A}(e^{j(\omega+\pi)}), \\ \tilde{g}[k] &\xleftrightarrow{\mathcal{F}} \widehat{\tilde{G}}(e^{j\omega}) = -e^{-j\omega_1} \frac{\widehat{H}(e^{-j(\omega+\pi)})}{\widehat{A}(e^{j\mathbf{D}^T \omega})}, \end{aligned}$$

and a scaled auto-correlation function is:

$$\widehat{A}(e^{j\mathbf{D}^T \omega}) = \frac{1}{2} \left| \widehat{H}(e^{j\omega}) \right|^2 \widehat{A}(e^{j\omega}) + \frac{1}{2} \left| \widehat{H}(e^{j(\omega+\pi)}) \right|^2 \widehat{A}(e^{j(\omega+\pi)}). \quad (10)$$

b) Non-subsampled Riesz-Quincunx wavelet smoothing and Hankel matrix: Denote $\underline{\mathfrak{C}}_\phi := \underline{\mathfrak{C}}_{\varphi_I}^* \underline{\mathfrak{C}}_{\tilde{\varphi}_I}$ where a matrix kernel $\phi \in \mathbb{R}^{n_1 \times n_2}$ and its matrix form $\underline{\Phi} := \underline{\tilde{\Phi}}_I \underline{\Phi}_I^T$ are defined in Equation 3 in the SM where $\underline{\tilde{\Phi}}_I$ and $\underline{\Phi}_I$ are matrix forms of $\tilde{\varphi}_I$ and φ_I , respectively. Similarly, we denote wavelet kernel tensors $\underline{\psi} = \{\psi_p\}_{p=1}^P := \{\psi_{il}\}_{i=0,\dots,L}^{l=0,\dots,L}$ and $\underline{\tilde{\psi}} = \{\tilde{\psi}_p\}_{p=1}^P := \{\tilde{\psi}_{il}\}_{i=0,\dots,L}^{l=0,\dots,L}$ with $\underline{\psi}_p, \underline{\tilde{\psi}}_p \in \mathbb{R}^{|\Omega|}$ ($p = 1, \dots, P$) and their matrix form are $\underline{\Psi} = \begin{pmatrix} \underline{\Psi}_1 & \dots & \underline{\Psi}_P \end{pmatrix}$ and $\underline{\tilde{\Psi}} = \begin{pmatrix} \underline{\tilde{\Psi}}_1 & \dots & \underline{\tilde{\Psi}}_P \end{pmatrix}$ where block element matrices $\underline{\Psi}_p, \underline{\tilde{\Psi}}_p \in \mathbb{R}^{n_2 n_1 \times n_2}$ are also defined in Equation 3 in the supplemental material. For an image $\underline{f} \in \mathbb{R}^{n_1 \times n_2}$, from proposition II.1 we have the following proposition:

Proposition II.3. *A non-subsampled Riesz Quincunx wavelet has a form of framelet decomposition (2):*

$$\begin{aligned} \underline{f} &= \underline{\mathfrak{C}}_\phi(\underline{f}) + \underline{\mathfrak{C}}_\psi^* \circ \text{prox}_{\mu \mathcal{D}} \circ \underline{\mathfrak{C}}_{\tilde{\psi}}(\underline{f}) \\ &= n_1 \mathcal{H}_{n_1|n_2}^\dagger \left(\mathcal{H}_{n_1|n_2}(\underline{f}) \left(\underline{\Phi} + \underline{\tilde{\Psi}} \underline{\tilde{\Psi}}^T \right) \right). \end{aligned} \quad (11)$$

Scaling and wavelet filter bank matrices satisfy the unity condition:

$$\underline{\Phi} + \underline{\tilde{\Psi}} \underline{\tilde{\Psi}}^T = \frac{1}{n_1} \text{Id}_{n_1 n_2 \times n_1 n_2}. \quad (12)$$

Proof. We provide a proof of Proposition II.3 in Section 5.3 in SM. \square

A smoothing version of a framelet decomposition (Equation 11) is defined with proximity operators $\text{prox}_{\mu \mathcal{D}}$ (for $\mu > 0$) (as defined in Equation 13 in the supplemental material):

$$\underline{\tilde{f}} = n_1 \mathcal{H}_{n_1|n_2}^\dagger \left(\mathcal{H}_{n_1|n_2}(\underline{f}) \underline{\Phi} + \text{prox}_{\mu \mathcal{D}} \left\{ \mathcal{H}_{n_1|n_2}(\underline{f}) \underline{\tilde{\Psi}} \right\} \underline{\tilde{\Psi}}^T \right).$$

Its iterative scheme, called generalized intersection algorithm with fixpoints is described in [23].

D. RQUNet-VAE expansion

Given a multi-channel image $\underline{f} = \{f_1, \dots, f_P\} \in \mathbb{R}^{|\Omega| \times P}$, we introduce our RQUNet-VAE expansion via mappings for skip-connecting signal, latent variables, and a reconstructed signal. Then, we propose RQUNet-VAE expansion and its functional space for regularization.

a) Mappings:: We introduce filter banks in an encoder $\underline{\theta}^{1(i)} \in \mathbb{R}^{d_1 \times d_2 \times P \times 2^i L}$ and $\underline{\theta}^{2(i)} \in \mathbb{R}^{d_1 \times n_2 \times P \times 2^i L}$ whose matrix forms are $\underline{\Theta}^{1(i)} \in \mathbb{R}^{d_1 d_2 P \times d_2 \times 2^i L}$, $\underline{\Theta}^{1(i)} \in \mathbb{R}^{d_1 n_2 P \times n_2 \times 2^i L}$ (as defined in Equation 6 in the supplemental material). Similar for filter banks $\left\{ \left(\underline{\tilde{\theta}}^{1(i)}, \underline{\tilde{\theta}}^{2(i)} \right) \right\}_{i=0}^{I-1}$ in a decoder. For batch-normalization and dropout layers, we refer the readers to [11], [29].

a.1. Skip-connecting signal (encoder): Given a local basis $\underline{\Xi}^{(i)}$ and an analysis operator at scale $i = 1, \dots, I$:

$$\begin{aligned} \mathcal{T}^{(i)} &:= \mathcal{R}_p \circ \mathcal{B} \circ \text{prox}_{\text{ReLU}} \circ \underline{\Xi}^{2^{(i)}}_{\underline{\theta}^{2(i)}} \circ \text{prox}_{\text{ReLU}} \circ \underline{\Xi}^{1^{(i)}}_{\underline{\theta}^{1(i)}} \\ &: \mathbb{R}^{n_1 \times n_2 \times P} \rightarrow \mathbb{R}^{2^{-i} n_1 \times 2^{-i} n_2 \times 2^i L}, \end{aligned} \quad (13)$$

we define an iterated mapping:

$$\begin{aligned} \mathcal{C}^{(i)} &= \mathcal{T}^{(i)} \circ \underline{\Xi}^{(i-1),T} \mathcal{T}^{(i-1)} \circ \dots \circ \underline{\Xi}^{(1),T} \mathcal{T}^{(1)} \\ &: \mathbb{R}^{n_1 \times n_2 \times P} \rightarrow \mathbb{R}^{2^{-(i-1)} n_1 \times 2^{-(i-1)} n_2 \times 2^{(i-1)} L}. \end{aligned}$$

Then, we have the following proposition describing a mapping for the skip-connecting signal:

Proposition II.4. *A mapping for skip-connecting signal in RQUNet-VAE is:*

$$\begin{aligned} \mathcal{C}_{\gamma_c} : \mathbb{R}^{n_1 \times n_2 \times P} &\rightarrow \left\{ \mathbb{R}^{2^{-(i-1)} n_1 \times 2^{-(i-1)} n_2 \times 2^{(i-1)} L} \right\}_{i=1}^I; \\ \underline{\underline{c}} &= \mathcal{C}_{\gamma_c}(\underline{f}) = \left\{ \underline{\underline{c}}^{(i)} = \mathcal{C}^{(i)}(\underline{f}) \right\}_{i=0}^{I-1} = \underline{\mathfrak{C}}_\phi(\underline{c}) + \underline{\mathfrak{C}}_{\tilde{\psi}}^* \underline{\mathfrak{C}}_{\tilde{\psi}}(\underline{\underline{c}}) \end{aligned} \quad (14)$$

which is equivalent to:

$$\underline{\underline{c}}^{(i)} = \left\{ \underline{c}_l^{(i)} \right\}_{l=0}^{2^{i-1} L - 1}, \quad \underline{c}_l^{(i)} = \underline{\mathfrak{C}}_\phi(\underline{c}_l^{(i)}) + \underline{\mathfrak{C}}_{\tilde{\psi}}^* \underline{\mathfrak{C}}_{\tilde{\psi}}(\underline{c}_l^{(i)})$$

for $i = 0, \dots, I-1$ where unknown filter banks are $\gamma_c := \left\{ \left(\underline{\theta}^{1(i)}, \underline{\theta}^{2(i)} \right) \right\}_{i=1}^I$.

Proof. We provide a proof of Proposition II.4 in Section 5.4 in SM. \square

a.2. Variational term (encoder): We firstly define linear mappings as a perceptron network for latent variable:

$$\begin{aligned} \mathcal{F}^\mu(\cdot) &= \underline{W}^\mu \text{vec}(\cdot) + b^\mu, \\ \mathcal{F}^\sigma(\cdot) &= \underline{W}^\sigma \text{vec}(\cdot) + b^\sigma : \mathbb{R}^{2^{-i-1} n_1 \times 2^{-i-1} n_2 \times 2^i L} \rightarrow \mathbb{R}^d; \end{aligned} \quad (15)$$

where $\underline{W}^\mu, \underline{W}^\sigma \in \mathbb{R}^{d \times 2^{-(I+1)} n_1 n_2 L}$ and $b^\mu, b^\sigma \in \mathbb{R}^d$ and d is latent dimension and $\text{vec}(\cdot)$ is a vectorize operation. Denote Hadamard product \odot and model's parameters as:

$$\gamma_c \cup \gamma_s \cup \gamma_m := \left\{ \underline{W}^\mu, b^\mu, \underline{W}^\sigma, b^\sigma, \left\{ \left(\underline{\theta}^{1(i)}, \underline{\theta}^{2(i)} \right) \right\}_{i=1}^I \right\};$$

then, we have the following proposition for a variational term:

Proposition II.5. *A latent variable in RQUNet-VAE is sampled from a distribution:*

$$z = \mathcal{M}_{\gamma_m} \left(\underline{f} \right) + \mathcal{S}_{\gamma_s}^{\frac{1}{2}} \left(\underline{f} \right) \odot \epsilon, \quad \epsilon \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}_d(\mathbf{0}_d, \text{Id}_d) \quad (16)$$

where maps for the mean and variance for the latent variable are:

$$\mathcal{M}_{\gamma_m} := \mathcal{F}^\mu \circ \mathcal{R}_p \circ \underline{\Xi}^{(I),T} \mathcal{T}^{(I)} \circ \dots \circ \underline{\Xi}^{(1),T} \mathcal{T}^{(1)}, \quad (17)$$

$$\mathcal{S}_{\gamma_s}^{\frac{1}{2}} := \mathcal{F}^\sigma \circ \mathcal{R}_p \circ \underline{\Xi}^{(I),T} \mathcal{T}^{(I)} \circ \dots \circ \underline{\Xi}^{(1),T} \mathcal{T}^{(1)}. \quad (18)$$

Proof. We provide a proof of Proposition II.5 in Section 5.5 in SM. \square

a.3. Decoder: Firstly, we define a mapping in the decoder:

$$\tilde{\mathcal{T}}^{(i)} = \left(\mathcal{Z}_1^{-1}, \text{prox}_{\text{ReLU}} \circ \underline{\mathcal{C}}_{\hat{\theta}^{1(i+1)}}^{\text{iso}} \circ \text{prox}_{\text{ReLU}} \circ \underline{\mathcal{C}}_{\hat{\theta}^{2(i+1)}}^{\text{iso}} \circ \mathcal{P}_{\hat{\underline{s}}_{\text{aug}}^{(i)}} \right),$$

$$\tilde{\mathcal{T}}^{(0)} = \text{prox}_{\text{ReLU}} \circ \underline{\mathcal{C}}_{\hat{\theta}^{1(1)}}^{\text{iso}} \circ \text{prox}_{\text{ReLU}} \circ \underline{\mathcal{C}}_{\hat{\theta}^{2(1)}}^{\text{iso}} \circ \mathcal{P}_{\hat{\underline{s}}_{\text{aug}}^{(0)}},$$

where \mathcal{Z}_1^{-1} is scale-delayed 1 step back of the 1st argument and a concatenate layer is defined as an operation $\mathcal{P}_{\hat{\underline{s}}_{\text{aug}}^{(i)}} \left(\underline{\hat{c}}_{\text{aug}}^{(i)} \right) = \left(\underline{\hat{c}}_{\text{aug}}^{(i)} \parallel \mathcal{B} \circ \underline{\hat{\Xi}}^{(i)} \underline{\hat{s}}^{(i+1)} \right)$, i.e. input signal at concatenation layer includes bypass signal $\underline{c}^{(i)}$ and unpooled-batchnormed lowpass signal $\mathcal{B} \circ \underline{\hat{\Xi}}^{(i)} \underline{\hat{s}}^{(i+1)}$.

Next, we define an iterated mapping for a skip-connecting signal \underline{c} and a lowpass signal $\underline{s}^{(I)}$ at scale I as:

$$\tilde{\mathcal{T}}_I \left(\underline{c}, \underline{s}^{(I)} \right) = \tilde{\mathcal{T}}^{(0)} \circ \tilde{\mathcal{T}}^{(1)} \circ \dots \circ \tilde{\mathcal{T}}^{(I-1)} \left(\underline{c}^{(I-1)}, \underline{s}^{(I)} \right).$$

Note that $\underline{s}^{(I)}$ is reconstructed from a latent variable z . Denote unknown parameters in a decoder as:

$$\alpha := \left\{ \underline{W}^s, b^s, \underline{\theta}^0, \left\{ \left(\underline{\tilde{\theta}}^{1(i)}, \underline{\tilde{\theta}}^{2(i)} \right) \right\}_{i=1}^I \right\}.$$

Then, we have the following proposition:

Proposition II.6. *Given a skip-connecting signal \underline{c} and a latent variable y from an encoder, a decoder mapping in RQUNet-VAE is defined as*

$$\mathcal{D}_\alpha \left(\underline{c}, z \right) := \text{prox}_{\text{ReLU}} \circ \underline{\mathcal{C}}_{\hat{\theta}^0}^{\text{iso}} \circ \tilde{\mathcal{T}}_I \left(\underline{c}, \text{uvec} \left(\underline{W}^s z + b^s \right) \right). \quad (19)$$

Proof. We provide a proof of Proposition II.6 in Section 5.6 in SM. \square

b) RQUNet-VAE expansion: An auto-encoder is recast as a latent factor model (a decoder) whose latent variable is computed from the observed data as an encoder:

$$z = \mathcal{M}_{\gamma_m} \left(\underline{f} \right) + \mathcal{S}_{\gamma_s}^{\frac{1}{2}} \left(\underline{f} \right) \odot \epsilon, \quad \epsilon \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}_d(\mathbf{0}_d, \text{Id}_d), \quad (20)$$

$$\underline{c} = \mathcal{C}_{\gamma_c} \left(\underline{f} \right), \quad (21)$$

$$\underline{\tilde{f}} = \mathcal{D}_\alpha \left(\underline{c}, z \right) + \sigma \underline{e}, \quad \underline{e} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}_{|\Omega| \times P}(0, \text{Id}), \quad (22)$$

$$\underline{\tilde{f}} = \underline{f} \in \mathfrak{F}, \quad (23)$$

with a known standard deviation $\sigma > 0$ and the data set $\mathfrak{F} = \left\{ \underline{f}_i \right\}_{i=1}^T \subset \mathbb{R}^{|\Omega| \times P}$. Combining Equation 22, Equation 20, Equation 21, and Equation 23, we obtain a variational auto-encoder:

$$\underline{\tilde{f}} = \mathcal{D}_\alpha \left(\mathcal{C}_{\gamma_c} \left(\underline{f} \right), \mathcal{M}_{\gamma_m} \left(\underline{f} \right) + \mathcal{S}_{\gamma_s}^{\frac{1}{2}} \left(\underline{f} \right) \odot \epsilon \right) + \sigma \underline{e}, \quad (24)$$

with standard normal random variables $\epsilon \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}_d(\mathbf{0}_d, \text{Id}_d)$ and $\underline{e} = [e_{l,c}]_{l \in \Omega}^{c=1, \dots, P}$, $e_{l,c} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1)$.

The auto-encoder system (Equations 22-23) is recast as Bayesian inference:

$$\begin{aligned} \underline{\tilde{f}} | z &\stackrel{\text{i.i.d.}}{\sim} \mathcal{N}_{n_1 \times n_2 \times P} \left(\mathcal{D}_\alpha \left(\underline{c}, z \right), \sigma^2 \text{Id} \right) \\ &= \mathbb{H}_\alpha \left(\underline{\tilde{f}} | z \right) = \mathbf{h}_\alpha \left(\underline{\tilde{f}} | z \right) d\underline{\tilde{f}}, \end{aligned} \quad (25)$$

$$\begin{aligned} z | \underline{\tilde{f}} &\stackrel{\text{i.i.d.}}{\sim} \mathbb{K}_\alpha \left(z | \underline{\tilde{f}} \right) = \mathbf{k}_\alpha \left(z | \underline{\tilde{f}} \right) dy \\ &\approx \mathcal{N}_d \left(\mathcal{M}_{\gamma_m} \left(\underline{\tilde{f}} \right), \text{diag} \left\{ \mathcal{S}_{\gamma_s} \left(\underline{\tilde{f}} \right) \right\} \right). \end{aligned} \quad (26)$$

An explanation for the above hierarchical model is: an observed signal $\underline{\tilde{f}}$ is assumed to be sampled from a normal distribution whose latent variable z is sampled from an unknown distribution \mathbb{K}_α which is approximated by a normal distribution (Equation 26) parameterized by (γ_m, γ_s) . The key idea for $z | \underline{\tilde{f}} \sim \mathbb{K}_\alpha$ to depend on a parameter α is because of the auto-encoder (Equations 20-23), i.e. $\underline{\tilde{f}} = \underline{f}$. To see this, Bayes' rule for a conditional density of random variable $z | \underline{\tilde{f}}$ is:

$$\begin{aligned} \mathbf{p} \left(z | \underline{\tilde{f}} \right) &= \frac{\mathbf{p} \left(\underline{\tilde{f}} | z \right) \mathbf{p}(z)}{\mathbf{p} \left(\underline{\tilde{f}} \right)} = \frac{\mathbf{p} \left(\underline{\tilde{f}} | z \right) \mathbf{p}(z)}{\mathbf{p} \left(\underline{\tilde{f}} \right)} \\ &\stackrel{(25)}{=} \frac{\mathbf{h}_\alpha \left(\underline{\tilde{f}} | z \right) \mathbf{p}(z)}{\mathbf{p} \left(\underline{\tilde{f}} \right)} := \mathbf{k}_\alpha \left(z | \underline{\tilde{f}} \right), \end{aligned}$$

which depends on a decoder's parameter α . Since the above density is intractable because of the incomputable integral

in a marginal distribution $p(\underline{f}) = \int_{\mathbb{R}^d} h_\alpha(\underline{f} | z) p(z) dz$, a distribution \mathbb{K}_α is approximated by a normal distribution:

$$\begin{aligned} k_\alpha(z | \underline{f}) &= \frac{h_\alpha(\underline{f} | z) p(z)}{p(\underline{f})} \\ &\approx \mathcal{N}\left(z; \mathcal{M}_{\gamma_m}(\underline{f}), \text{diag}\left\{\mathcal{S}_{\gamma_s}(\underline{f})\right\}\right). \end{aligned} \quad (27)$$

Denote variance vector as

$$\begin{aligned} \mathcal{S}_{\gamma_s}(\underline{f}) &= \mathcal{S}_{\gamma_s}^{\frac{1}{2}}(\underline{f}) \odot \mathcal{S}_{\gamma_s}^{\frac{1}{2}}(\underline{f}) \\ &= \left(\mathcal{S}_{\gamma_s}(\underline{f})_1 \quad \dots \quad \mathcal{S}_{\gamma_s}(\underline{f})_d\right)^T \in \mathbb{R}^d. \end{aligned}$$

Choose a standard normal prior distribution $\mathbb{P}(z) = p(z)dz = \mathcal{N}_d(\mathbf{0}_d, \text{Id}_d)$, we have the following proposition for finding model's parameters:

Proposition II.7. *Unknown parameters in RQUnet-VAE are obtained from the following minimization problem:*

$$(\gamma_c^\dagger, \gamma_m^\dagger, \gamma_s^\dagger, \alpha^\dagger) = \text{argmin} \mathcal{L}(\gamma_c, \gamma_m, \gamma_s, \alpha) \quad (28)$$

where:

$$\begin{aligned} \mathcal{L}(\cdot) &:= \frac{1}{2\sigma^2} \sum_{i=1}^T \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \text{Id})} \left[\left\| \underline{f}_i - \mathcal{D}_\alpha \left(\mathcal{E}_{\gamma_c}(\underline{f}_i), \mathcal{M}_{\gamma_m}(\underline{f}_i) + \mathcal{S}_{\gamma_s}^{\frac{1}{2}}(\underline{f}_i) \odot \epsilon \right) \right\|_{\ell_2}^2 \right] \\ &+ \sum_{i=1}^T KL \left[\mathcal{N}_d \left(\mathcal{M}_{\gamma_m}(\underline{f}), \text{diag} \left\{ \mathcal{S}_{\gamma_s}(\underline{f}) \right\} \right) \parallel \mathcal{N}_d(\mathbf{0}_d, \text{Id}_d) \right] \end{aligned}$$

and KL-divergence is defined as:

$$KL[\cdot \parallel \cdot] = \frac{1}{2} \left(\left\| \mathcal{M}_{\gamma_m}(\underline{f}) \right\|_{\ell_2}^2 - d + \sum_{i=1}^d \left(\mathcal{S}_{\gamma_s}(\underline{f})_i - \log \mathcal{S}_{\gamma_s}(\underline{f})_i \right) \right). \quad (29)$$

Proof. We provide a proof of Proposition II.7 in Section 5.7 in SM. \square

It is clear that the only one random variable in the above minimization is ϵ , so the gradient descent method can be applied for model parameters $(\gamma_c, \gamma_m, \gamma_s, \alpha)$. Since high order Riesz-Quincunx wavelet expansion is an identity operator, for training these model parameters, we remove this layer in the training procedure. But, later we use it in its smoothing version with the trained parameters $(\gamma_c^\dagger, \gamma_m^\dagger, \gamma_s^\dagger, \alpha^\dagger)$ to truncate small wavelet coefficients of signals in skip-connection. A summary of RQUnet-VAE with a training procedure (without Riesz-Quincunx wavelet) is shown in the SM for encoder and decoder, respectively.

c) *Generalized Besov space by RQUnet-VAE and proximal operators:* We note that Equation 28 is a non-convex minimization problem with potentially many local minimas. Assume an existence of minima $(\alpha^\dagger, \gamma_c^\dagger, \gamma_m^\dagger, \gamma_s^\dagger)$ such that a condition of perfect reconstruction occurs, i.e. $\underline{f} = \underline{f}$; then, the auto-encoder (Equations 22-23) plays as a generalized wavelet expansion with learnable parameters:

$$\begin{aligned} \lim_{\sigma \rightarrow 0} \left\| \underline{f} - \mathcal{D}_{\alpha^\dagger} \left(\mathcal{E}_{\gamma_c^\dagger}(\underline{f}), \mathcal{M}_{\gamma_m^\dagger}(\underline{f}) + \mathcal{S}_{\gamma_s^\dagger}^{\frac{1}{2}}(\underline{f}) \odot \epsilon \right) + \sigma \underline{e} \right\|_{\ell_2}^2 \\ = 0, \end{aligned}$$

having its deterministic version:

$$\underline{f} = \mathcal{D}_{\alpha^\dagger} \left(\mathcal{E}_{\gamma_c^\dagger}(\underline{f}), \mathcal{M}_{\gamma_m^\dagger}(\underline{f}) + \mathcal{S}_{\gamma_s^\dagger}^{\frac{1}{2}}(\underline{f}) \odot \epsilon \right),$$

where encoder and decoder are forward and backward wavelet transform. Then, it induces a generalized Besov space:

$$\mathfrak{B} := \left\{ \underline{f} \in \mathbb{R}^{|\Omega| \times P} : \left\| \underline{f} \right\|_{\mathfrak{B}} := \mathcal{P} \circ \mathcal{E}_{\gamma_c}(\underline{f}) < \infty \right\}, \quad (30)$$

where function $\mathcal{P}(\cdot)$ acts on the skip-connecting signal $\underline{c} = \mathcal{E}_{\gamma_c}(\underline{f})$ for proximal mapping as described in Equation 13 in the supplemental material.

Then, given an image $\underline{f} \in \mathbb{R}^{|\Omega| \times P}$ with an expansion in a space \mathfrak{B} (30), we have a regularization in that space \mathfrak{B} with the optimally trained parameters as generalized wavelet smoothing via a proximal operator:

$$\begin{aligned} \tilde{\underline{c}} &= \underset{\underline{w}}{\text{argmin}} \left\{ \mathcal{P}(\underline{w}) + \frac{1}{2\mu} \left\| \mathcal{E}_{\gamma_c}(\underline{f}) - \underline{w} \right\|_F^2 \right\} \\ &= \text{prox}_{\mu \mathcal{P}} \circ \mathcal{E}_{\gamma_c}(\underline{f}) \\ &= \underline{c}_{\phi} \circ \mathcal{E}_{\gamma_c}(\underline{f}) + \underline{c}_{\psi}^* \circ \text{prox}_{\mu \mathcal{P}} \circ \underline{c}_{\psi} \circ \mathcal{E}_{\gamma_c}(\underline{f}), \end{aligned} \quad (31)$$

where its element form is:

$$\begin{aligned} \tilde{\underline{c}} &= \left\{ \tilde{c}_l^{(i)} \right\}_{i=0}^{I-1}, \quad \tilde{c}_l^{(i)} = \left\{ \tilde{c}_l^{(i)} \right\}_{l=0}^{2^{i-1}L-1}, \\ \tilde{c}_l^{(i)} &= \underline{c}_{\phi} \left(c_l^{(i)} \right) + \underline{c}_{\psi}^* \circ \text{prox}_{\mu \mathcal{P}} \circ \underline{c}_{\psi} \left(c_l^{(i)} \right) \end{aligned}$$

for $l = 0, \dots, 2^{i-1}L - 1$, $i = 0, \dots, I - 1$. Note that $\text{prox}_{\mu \mathcal{P}}(\cdot)$ in Equation 31 acts on wavelet coefficients only. A reason is scaling coefficients contain main energy of a signal whose values are large. These scaling coefficients should be preserved during a shrinking process while wavelet coefficients mainly contain oscillating signals in different scales, including noise which should be removed.

Then, we have a smoothed image as:

$$\tilde{\underline{f}} = \mathcal{D}_{\alpha^\dagger} \left(\tilde{\underline{c}}, \mathcal{M}_{\gamma_m}(\underline{f}) + \mathcal{S}_{\gamma_s}^{\frac{1}{2}}(\underline{f}) \odot \epsilon \right) \in \mathbb{R}^{|\Omega| \times P}.$$

In summary, we have the following proposition:

Proposition II.8. Given a learned parameter $(\alpha^\dagger, \gamma_c^\dagger, \gamma_m^\dagger, \gamma_s^\dagger)$ obtain by training RQUNet-VAE with a dataset, we have RQUNet-VAE smoothing (with a parameter $\mu > 0$) for an image $\underline{\underline{f}}$ as a solution of a regularization in a generalized Besov space \mathfrak{B} (30):

$$\begin{aligned} \underline{\underline{f}} = \mathcal{D}_{\alpha^\dagger} \left(\underline{\underline{\mathcal{C}}}_\phi \circ \mathcal{C}_{\gamma_c^\dagger} \left(\underline{\underline{f}} \right) + \underline{\underline{\mathcal{C}}}_\psi^* \circ \text{prox}_{\mu\mathcal{P}} \circ \underline{\underline{\mathcal{C}}}_{\tilde{\psi}} \circ \mathcal{C}_{\gamma_c^\dagger} \left(\underline{\underline{f}} \right) \right. \\ \left. , \mathcal{M}_{\gamma_m^\dagger} \left(\underline{\underline{f}} \right) + \mathcal{S}_{\gamma_s^\dagger}^{\frac{1}{2}} \left(\underline{\underline{f}} \right) \odot \epsilon \right). \end{aligned} \quad (32)$$

Proof. A proof of Proposition II.8 is directly obtained from the above paragraphs. \square

E. RQUNet-VAE iterative shrinkage Lagrangian system

a) *Decomposition for multi-band image:* RQUNet-VAE iterative shrinkage algorithm for multi-band image $\underline{\underline{f}} \in \mathbb{R}^{|\Omega| \times P}$ is described in Algorithm 1 in the SM. This is equivalent to a nonlinear mapping with unknown model's parameters Γ : $\underline{\underline{u}}^{(0)} = \underline{\underline{f}}, \underline{\underline{\lambda}}^{(0)} = \mathbf{0}$ for $\tau = 1, \dots, N$ and

$$\left(\underline{\underline{u}}^{(\tau)}, \underline{\underline{\lambda}}^{(\tau+1)} \right) = \mathcal{K}_\mu \left(\underline{\underline{f}}, \underline{\underline{u}}^{(\tau-1)}, \underline{\underline{\lambda}}^{(\tau)}; \Gamma \right). \quad (33)$$

b) *Diffusion process and spectral decomposition for multiband image:* For a multiband image $\underline{\underline{f}} \in \mathbb{R}^{|\Omega| \times P}$, a diffusion process by a Lagrangian system (34) is: $\underline{\underline{u}}^{(0)} = \underline{\underline{f}}, \tau = 1, \dots, N$,

$$\left(\underline{\underline{u}}^{(\tau)}, \underline{\underline{\lambda}}^{(\tau+1)} \right) = \mathcal{K}_\mu \left(\underline{\underline{u}}^{(\tau-1)}, \underline{\underline{u}}^{(\tau-1)}, \underline{\underline{\lambda}}^{(\tau)}; \Gamma \right). \quad (34)$$

Given $\left\{ \underline{\underline{u}}^{(\tau)} \right\}_{\tau=1}^N$ generated by diffusion process (34) in Algorithm 2, we have a discrete TV-like transform as in [9]:

$$\underline{\underline{\phi}}^{(\tau)} := \frac{\tau}{\beta} \left(\underline{\underline{u}}^{(\tau+1)} - 2\underline{\underline{u}}^{(\tau)} + \underline{\underline{u}}^{(\tau-1)} \right),$$

whose inverse transform is: $\underline{\underline{f}}^{(N)} := (1 + N) \underline{\underline{u}}^{(N)} - N \underline{\underline{u}}^{(N+1)}$,

$$\underline{\underline{f}} = \underline{\underline{f}}^{(N)} + \beta \sum_{\tau=1}^N \underline{\underline{\phi}}^{(\tau)}.$$

Its filtered version

$$\underline{\underline{f}}_{HN} = H^{(N)} \underline{\underline{f}}^{(N)} + \beta \sum_{\tau=1}^N H^{(\tau)} \underline{\underline{\phi}}^{(\tau)}$$

is defined, e.g. by the ideal spectral-filter $H^{(\tau)}$:

$$H_{\text{lowpass}}^{(\tau)} = \begin{cases} 0, & \tau \in \{0, \dots, \tau_1\}, \\ 1, & \tau \in \{\tau_1, \dots, N\}, \end{cases}$$

$$H_{\text{highpass}}^{(\tau)} = \begin{cases} 1, & \tau \in \{0, \dots, \tau_1\}, \\ 0, & \tau \in \{\tau_1, \dots, N\}, \end{cases}$$

$$H_{\text{bandpass}}^{(\tau)} = \begin{cases} 0, & \tau \in \{0, \dots, \tau_1\}, \\ 1, & \tau \in \{\tau_1, \dots, \tau_2\}, \\ 0, & \tau \in \{\tau_2, \dots, N\} \end{cases},$$

$$H_{\text{bandstop}}^{(\tau)} = \begin{cases} 1, & \tau \in \{0, \dots, \tau_1\}, \\ 0, & \tau \in \{\tau_1, \dots, \tau_2\}, \\ 1, & \tau \in \{\tau_2, \dots, N\}, \end{cases}$$

where the time threshold τ_i are selected by the TV-spectrum $S^{(\tau)} := \left\| \underline{\underline{\phi}}^{(\tau)} \right\|_{\ell_1}$.

c) *Decomposition for multiband time series:* Next, we show how our RQUNet-VAE can be applied to satellite image time series, that is, a sequence of images of the same area.

Given a multi-band video $\underline{\underline{f}} = \{ \underline{\underline{f}}_1, \dots, \underline{\underline{f}}_T \} \in \mathbb{R}^{|\Omega| \times P \times T}$, $\underline{\underline{f}}_t = \{ \underline{\underline{f}}_{t,1}, \dots, \underline{\underline{f}}_{t,P} \}$, $\underline{\underline{f}}_{t,p} \in \mathbb{R}^{|\Omega|}$ and scaling and wavelet bases, e.g. 1D Haar bases:

$$\begin{aligned} \phi_{I,t,m} &= \begin{cases} 1, & t \geq 2^I m \text{ \& } t < 2^I(1+m), \\ 0, & \text{else} \end{cases}, \\ \xi_{i,t,m} &= \begin{cases} 1, & t \geq 2^i m \text{ \& } t < 2^i(\frac{1}{2} + m) \\ -1, & t \geq 2^i(\frac{1}{2} + m) \text{ \& } t < 2^i(1+m), \\ 0, & \text{else} \end{cases} \end{aligned} \quad (35)$$

a time-wavelet smoothing expansion for a video $\underline{\underline{f}} \in \mathbb{R}^{|\Omega| \times P \times T}$ by proximal operator $\mu\mathcal{P}$ is defined as: $t = 1, \dots, T$,

$$\begin{aligned} \underline{\underline{f}}_t &= \frac{1}{2^{\frac{I}{2}}} \sum_{s=0}^{2^{-I}T-1} \frac{1}{2^{\frac{I}{2}}} \sum_{t'=1}^T \underline{\underline{f}}_{t'} \phi_{I,t',s} \phi_{I,t,s} \\ &+ \sum_{i=1}^I \frac{1}{2^{\frac{i}{2}}} \sum_{s=0}^{2^{-i}T-1} \text{prox}_{\mu\mathcal{P}} \left\{ \underbrace{\frac{1}{2^{\frac{i}{2}}} \sum_{t'=1}^T \underline{\underline{f}}_{t'} \xi_{i,t',s}}_{:= \underline{\underline{w}}_{i,m}} \right\} \xi_{i,t,s}, \end{aligned}$$

which is equivalent to:

$$\underline{\underline{f}} := \mathcal{G}_\phi \left(\underline{\underline{f}} \right) + \mathcal{W}_\xi^* \circ \text{prox}_{\mu\mathcal{P}} \circ \mathcal{W}_\xi \left(\underline{\underline{f}} \right) = \{ \underline{\underline{f}}_1, \dots, \underline{\underline{f}}_T \}. \quad (36)$$

Wavelet coefficient in (36) is $\underline{\underline{w}}_i = \{ \underline{\underline{w}}_{i,1}, \dots, \underline{\underline{w}}_{i,T} \} := \mathcal{W}_\xi \left\{ \underline{\underline{f}} \right\} \in \mathbb{R}^{|\Omega| \times P \times 2^{-i}T}$, $\underline{\underline{w}}_{i,m} = \{ \underline{\underline{w}}_{i,m,1}, \dots, \underline{\underline{w}}_{i,m,P} \}$

and $w_{i,m,p} \in \mathbb{R}^{|\Omega|}$. Combined with Equation (34), RQUNet-VAE's iterative shrinkage algorithm for multi-band video $f \in \mathbb{R}^{n_1 \times n_2 \times P \times T}$ is described in Algorithm 2 in SM, called RQUNet-VAE scheme 2.

F. RQUNet-VAE based segmentation

In this section, we propose a mathematical framework for a segmentation problem with our RQUNet-VAE which serves as a smoothing term. Given a dataset of original images (without artificially added noise) with ground-truth masks, we first train the UNet-VAE to obtain the optimal model parameters. Next we predict a segmented image from a noisy image using the RQUNet-VAE with a smoothing term as a truncation scheme of the N -th order Riesz-Quincunx wavelet expansion on the skip-connecting signals. If the training data is $(\mathcal{F}, \mathcal{F}^{\text{gt}}) := \left\{ \underline{\underline{f}}_i, \underline{\underline{f}}_i^{\text{gt}} \right\}_{i=1}^n \subset \mathbb{R}^{|\Omega| \times P} \times \mathbb{R}^{|\Omega| \times K}$ and K -dimensional hot key tensors of the ground-truth masks, i.e. pixel intensity and its allocation are:

$$\underline{\underline{f}}_i := [f_{i,l}]_{l \in \Omega}, f_{i,l} \in \mathbb{R}^P, \underline{\underline{f}}_i^{\text{gt}} := [f_{i,l}^{\text{gt}}]_{l \in \Omega},$$

where $\underline{\underline{f}}_i^{\text{gt}} = [f_{i,l,k}^{\text{gt}}]_{k=1}^K \in \mathbb{R}^K$ is a one-hot-key vector. Given the loss function: $\underline{\underline{x}}, \underline{\underline{y}} \in \mathbb{R}^{|\Omega| \times K}$,

$$\mathcal{H}(\underline{\underline{x}}, \underline{\underline{y}}) = \sum_{l \in \Omega} \sum_{k=1}^K x_{l,k} \log \frac{\exp(y_{l,k})}{\sum_{h=1}^K \exp(y_{l,h})};$$

similar to proposition II.7, a loss function of our RQUNet-VAE based segmentation problem is as follows:

Proposition II.9. *Unknown parameters in RQUNet-VAE based segmentation are obtained from the following minimization problem:*

$$(\underline{\underline{\theta}}^\dagger, \gamma_c^\dagger, \gamma_m^\dagger, \gamma_s^\dagger, \alpha^\dagger) = \operatorname{argmin}_{\theta, \gamma_c, \gamma_m, \gamma_s, \alpha} \sum_{i=1}^n \mathcal{L}(\underline{\underline{\theta}}, \gamma_c, \gamma_m, \gamma_s, \alpha) \quad (37)$$

where:

$$\begin{aligned} \mathcal{L}(\cdot) := & \text{KL} \left[\mathcal{N}_d \left(\mathcal{M}_{\gamma_m}(\underline{\underline{f}}_i), \text{diag} \left\{ \mathcal{S}_{\gamma_s}(\underline{\underline{f}}_i) \right\} \right) \parallel \mathcal{N}_d(\mathbf{0}_d, Id_d) \right] \\ & - \frac{1}{2n\sigma^2} \mathbb{E}_{\epsilon \sim \mathcal{N}(\mathbf{0}, Id)} \left[\mathcal{H} \left(\underline{\underline{f}}_i^{\text{gt}}, \underline{\underline{\theta}}^{\text{iso}} \circ \mathcal{D}_\alpha \left(\mathcal{C}_{\gamma_c}(\underline{\underline{f}}_i), \mathcal{M}_{\gamma_m}(\underline{\underline{f}}_i) \right. \right. \right. \\ & \left. \left. \left. + \mathcal{S}_{\gamma_s}(\underline{\underline{f}}_i)^{\frac{1}{2}} \odot \epsilon \right) \right) \right], \end{aligned}$$

and the K-L divergence is defined in (29).

Proof. We provide a proof of Proposition II.9 in Section 5.8 in SM. \square

After training the model parameters by minimizing the loss function (37), the segmentation of a new noisy image is predicted:

$$f^{\text{new}} = f_0^{\text{new}} + \sigma \epsilon,$$

with standard Gaussian noise $\epsilon = [\epsilon_{l,c}]_{l \in \Omega}^{c=1, \dots, P} \in \mathbb{R}^{|\Omega| \times P}$ and $\epsilon_{l,c} \sim \mathcal{N}(0, 1)$. From proposition II.8 by adding N -th order Riesz Quincunx wavelet truncation with proximal operator $\text{prox}_{\mu \mathcal{P}}(\cdot)$ parameterized by a smoothing parameter μ as:

$$\begin{aligned} u(\underline{\underline{f}}^{\text{new}}) &= \left\{ u_1(\underline{\underline{f}}^{\text{new}}), \dots, u_K(\underline{\underline{f}}^{\text{new}}) \right\} \in \mathbb{R}^{|\Omega| \times K} \\ &= \underline{\underline{\theta}}^{\text{iso}} \circ \mathcal{D}_{\alpha^\dagger} \left(\underline{\underline{\theta}}^{\text{iso}} \circ \mathcal{C}_{\gamma_c^\dagger}(\underline{\underline{f}}^{\text{new}}) + \underline{\underline{\theta}}^* \circ \text{prox}_{\mu \mathcal{P}} \circ \underline{\underline{\theta}}^{\text{iso}} \circ \mathcal{C}_{\gamma_c^\dagger}(\underline{\underline{f}}^{\text{new}}) \right), \\ &\quad \mathcal{M}_{\gamma_m^\dagger}(\underline{\underline{f}}^{\text{new}}) + \mathcal{S}_{\gamma_s^\dagger}^{\frac{1}{2}}(\underline{\underline{f}}^{\text{new}}) \odot \epsilon, \end{aligned}$$

$$\underline{\underline{y}}^{\text{new}} = \operatorname{argmax}_{k=1, \dots, K} u_k(\underline{\underline{f}}^{\text{new}}) \in \mathbb{R}^{|\Omega| \times K}, u_k(\underline{\underline{f}}^{\text{new}}) \in \mathbb{R}^{|\Omega|}.$$

III. EXPERIMENTAL RESULTS

A. Experimental Setup

This section describes the dataset used for our experiments, introduces competitor solutions, and defines evaluation metrics to measure the ability of an algorithm to reduce the noise of an image.

1) *Datasets:* We used Sentinel-2 satellite images (S30:MSI harmonized, V1.5) processed as part of the Harmonized Landsat Sentinel-2 (HLS) dataset obtained from USGS DAAC (<https://lpdaac.usgs.gov/data/get-started-data/collection-overview/missions/harmonized-landsat-sentinel-2-hls-overview/>) [5]. The S30:MSI harmonized surface reflectance product is resampled from the original 10m to 30m resolution and adjusted to Landsat8 spectral response function in order to ultimately create a harmonized time series with a 2-3 day revisit. Radiometric and geometric corrections are applied to convert data to surface reflectance, adjust for BRDF differences, and spectral bandpass differences. The study was focused on the HLS tile 18STH covering a large part of Northern Virginia in the US, for 2016-2020. From the main HLS tile's time series, 510 images were randomly created with a size of 256×256 pixels and including only the three visible bands (4R, 3G, 2B). Quality Assessment (QA) layers were used to exclude images with more than 30 percent cloud shadow, adjacent cloud, cloud and cirrus clouds. The 500 images were used for training of our RQUNet-VAE and competitor approaches, while ten images were used for testing in the denoising and decomposition experiments. The ten test images are visualized in SM.

For the time series decomposition experiments the above Sentinel-2 data were used to create 80, randomly located time series of length 99 images and image size 40×40 pixels. Images were normalized from 0 to 1 using image minimum and maximum. The RQUNet-VAE was trained on each image in the dataset to obtain all unknown model parameters. All the Sentinel-2 datasets are available on Github¹.

¹<https://github.com/trile83/RQUNetVAE>

For our image segmentation experiments in Section III-D we used National Agriculture Imagery Program (NAIP) images coinciding with high resolution ground truth land cover data of each pixel in an image [18]. We acquired NAIP Nature Color imagery of northern Virginia consisting of RGB bands which is similar to some commercial satellite imagery (e.g. BlackSky, Planet). For training and validation (ground truth) we used the 1m resolution land cover dataset for the Chesapeake Bay watershed [1]. The dataset contains six classes - Water, Tree and shrubs, Herbaceous, Barren, Impervious (roads) and Impervious (other). We combined Impervious (other) and Impervious (roads) together into an integrated Impervious class. This provided us with four classes Water, Tree and Shrubs, Grass, and Impervious. For the segmentation experiments these classes were grouped into three classes, Vegetated (tree, grass, shrubs), Water and Impervious.

2) *Evaluated Algorithms:* Since RQUNet-VAE Scheme 1 is based on harmonic analysis, we compared it to state-of-the-art approaches including wavelet CDF 9/7 [33], [30], [6], curvelet [3] and Riesz Dyadic wavelet kernel [23], while the iterative Scheme 2 method was compared to state-of-the-art iterative methods such as directional TV-L2 [27], [10], [31] and GIAF [23].

Note that since the competitor methods were designed for gray-scale images, the methods were applied independently to each band of every multi-band image. The RQUNet-VAE, on the other hand, was designed for multi-channels images.

3) Evaluation Metrics:

a) *Image reconstruction:* To evaluate the ability of an algorithm to reduce the noise of an image, this section defines two commonly used measures. Given a clean multi-channel image $\underline{f} \in \mathbb{R}^{|\Omega| \times P}$ as ground-truth and a denoised image $\underline{f}^\dagger \in \mathbb{R}^{|\Omega| \times P}$, we use

- peak-signal-to-noise-ratio (PSNR):

$$\text{PSNR} = 10 \log \frac{\max(\underline{f})}{\text{MSE}}, \text{MSE} = \frac{1}{n_1 n_2 P} \|\underline{f} - \underline{f}^\dagger\|_{\ell_2}^2$$

- Structural similarity index (SSIM) [41]:

$$\text{SSIM} = \frac{(2\mu_f \mu_{f^\dagger} + c_1)(2\sigma_{ff^\dagger})}{(\mu_f^2 + \mu_{f^\dagger}^2 + c_1)(\sigma_f^2 + \sigma_{f^\dagger}^2 + c_2)}$$

where (μ_f, σ_f^2) and $(\mu_{f^\dagger}, \sigma_{f^\dagger}^2)$ are mean and variance of ground-truth \underline{f} and its denoised image \underline{f}^\dagger and σ_{ff^\dagger} are their covariance. $c_1 = (0.01r)^2$, $c_2 = (0.03r)^2$ where r is the dynamic range of pixel-values.

b) *Image segmentation:* To evaluate our RQUNet-VAE for segmentation of noisy test data, given that our statistical method provides uncertainty quantification, we propose the following evaluation framework:

Given a noisy test image $\underline{f} \in \mathbb{R}^{|\Omega| \times P}$ as an input of the trained RQUNet-VAE segmentation and its ground-truth $\underline{f}^{\text{gt}} = [\underline{f}_l^{\text{gt}}]_{l \in \Omega} \in \mathbb{R}^{n_1 \times n_2 \times K}$, $\underline{f}_l \in \mathbb{R}^K$ and the trained RQUNet-VAE by minimizing the loss function (37) from the

clean dataset, we run prediction n time to obtain segmented masks $u_l^{(i)}(\underline{f}) = \left[u_l^{(i)}(\underline{f}) \right]_{l \in \Omega} \in \mathbb{R}^{|\Omega| \times K}$, $u_l^{(i)}(\underline{f}) \in \mathbb{R}^K$ for $i = 1, \dots, n$. For segmentation problem, class-balanced accuracy is defined for every pixel $l \in \Omega$:

$$\hat{p}_l^n = \frac{1}{n} \sum_{i=1}^n \delta_{\{u_l^{(i)}(\underline{f}), f_l^{\text{gt}}\}}, \delta_{\{u_l^{(i)}(\underline{f}), f_l^{\text{gt}}\}} = \begin{cases} 1, & u_l^{(i)}(\underline{f}) = f_l^{\text{gt}} \\ 0, & \text{else} \end{cases} \quad (38)$$

which is modeled by a Binomial random variable approximated by Normal distribution (for large n):

$$Y_i := \delta_{\{u_l^{(i)}(\underline{f}), f_l^{\text{gt}}\}} \sim \text{Bernoulli}(\hat{p}_l^n, 1), \quad (39)$$

$$p_l^n = \frac{1}{n} \sum_{i=1}^n Y_i \sim \frac{1}{n} \text{Binomial}(\hat{p}_l^n, n) \approx \mathcal{N}\left(\hat{p}_l^n, \frac{\hat{p}_l^n(1 - \hat{p}_l^n)}{n}\right). \quad (40)$$

4) *Training Procedure:* For comparison with other state-of-the-art-methods, all parameters have been optimized for minimal mean-square-error (MSE) via heuristic search for each individual training image.

Since RQUNet-VAE is a hybrid model of deterministic high-order Riesz-Quincunx wavelet, we apply a training procedure for UNet-VAE with 100 epochs and a batch size of 16 using the Adam optimization method [12]. The source code of our training in PyTorch can be found at <https://github.com/trile83/RQUNetVAE>.

Image Set		Sentinel-2, std = 0.04
Number of Images		10
RQUNet-VAE scheme 1	PSNR	38.693
	SSIM	0.969
Riesz Dyadic	PSNR	37.433
	SSIM	0.959
curvelet	PSNR	36.314
	SSIM	0.942
wavelet CDF 9/7	PSNR	36.061
	SSIM	0.945
RQUNet-VAE scheme 2	PSNR	39.087
	SSIM	0.971
GIAF-Riesz Dyadic	PSNR	38.987
	SSIM	0.969
TV-L2 ($L = 2$)	PSNR	38.522
	SSIM	0.968
TV-L2 ($L = 9$)	PSNR	38.53
	SSIM	0.968

TABLE I
PSNR AND SSIM: MEAN OVER THE THREE IMAGE SETS (1ST ND 2ND BEST IN BOLD).

B. Denoising Experiments

To impose artificial noise on images for evaluation, Gaussian noise was added to each band of the Sentinel-2 images with a standard deviation of $\text{std} = 0.04$. All images were normalized to interval $[0, 1]$ using the image minimum and maximum. To give some examples, original image shown in Figure 2(a) and the corresponding image with added noise in Figure 2(b). (All ten test images with added noise are visualized in SM).

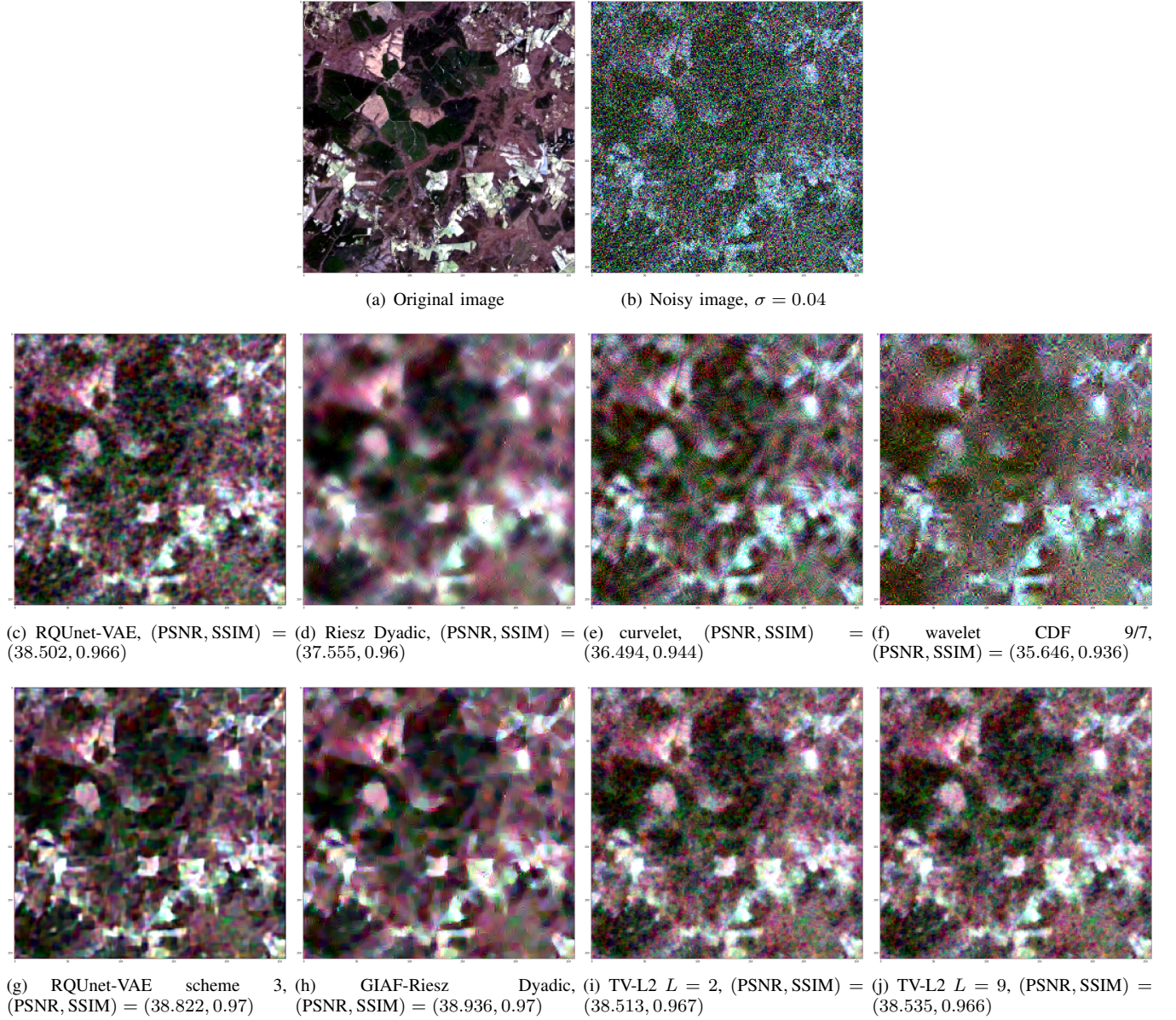


Fig. 2. Qualitative image denoising results for RQUNet-VAE and competitor approaches. Image denoising by the RQUNet-VAE with standard deviation $\sigma = 0.04$ and optimal α^* .

1) *Quantitative Results:* Table I provides the results of the comparison of the RQUNet-VAE Schemes 1 and 2 with the state-of-the-art methods described in Section III-A2 using the evaluation metrics described in Section III-A3. For the ten test images described in Section III-A1 the proposed RQUNet-VAE Scheme 2 yields the highest Peak-Signal-to-Noise-Ratio (PSNR) and the highest Structural Similarity Index (SSIM) of all competitors (Table 1). Among all the approaches based on harmonic analysis, RQUNet-VAE Scheme 1 yields the best results. Note GIAF also iteratively computes scaling and wavelet coefficients, similar to our RQUNet-VAE Scheme 2. The main difference, however, is that GIAF employs Riesz Dyadic wavelet kernel for scaling and wavelet functions whereas RQUNet-VAE Scheme 2 uses our adapted RQUNet-VAE. Using this scheme, allows more redundant parameters to better learn specific features from the data and thus, better adapt to specific images. (Denoising results for all other test

images are visualized in SM).

2) *Qualitative Results:* Figures 1, 2 and 3 in SM are original, noisy and denoised images with the 2nd scheme of our RQUNet-VAE. Figures 2(c)-2(j) provide qualitative results to assess how well the RQUNet-VAE schemes are able to reduce artificially added noise individual Sentinel2 images. Figure 2(c) shows the result of denoising the noisy image of Figure 2(b) using our RQUNet-VAE Scheme 1. Since Scheme 1 is based on harmonic analysis, we compare it to noise reduction using a Riesz Dyadic wavelet kernel [23] in Figure 2(d), using curvelets [3] in Figure 2(e), and using wavelet CDF 9/7 [33], [30], [6] in Figure 2(f). The RQUNet-VAE Scheme 1 best preserves the edges of objects of the original image. The Riesz Dyadic wavelet kernel in Figure 2(d) yields a blurred denoised image. In contrast, wavelets (Figure 2(f)) better retain edges, but much of the noise remains. Across all these figures, it is clearly discernible

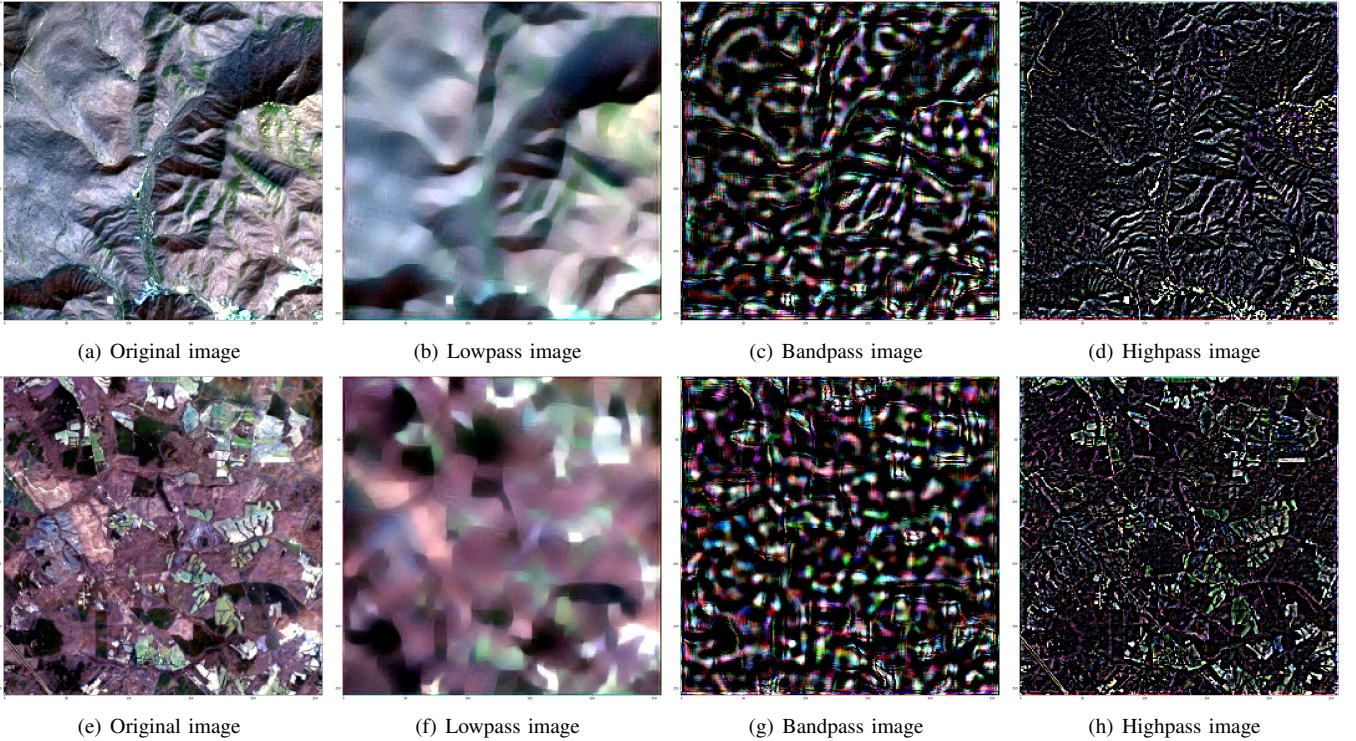


Fig. 3. Image decomposition by RQUNet-VAE by spectral analysis with a threshold $T = 3$ for spectral histogram. Parameters: number of iteration $\text{Iter} = 30$, a smoothing parameter $\alpha = 0.4$.

that the proposed RQUNet-VAE Scheme 1 provides the best combination of noise reduction and delineation of edges and objects.

Figure 2(g) shows the denoised image using the RQUNet-VAE Scheme 2, compared to other iterative methods including GIAF [23] in Figure 2(h) and directional TV-L2 [27], [10], [31] using a number of directions $L = 2$ in Figure 2(i) and $L = 9$ in Figure 2(j). Figure 2 (g) illustrates that RQUNet-VAE significantly reduces artefacts while preserving texture pattern, contrast and sharp edges of objects in the reconstructed image, such that it is most similar to the original image. Note that the GIAF-Riesz dyadic method in Figure 2 (h) also performs very well, but a reconstructed images also contained some artefacts. The RQUNet-VAE Scheme 2 again provides the best trade-off between reduction of noise and delineation of edges and objects, where GIAF (Figure 2(h)) oversmooths edges between objects while TV-L2 (Figure 2(i) and Figure 2(j)) still retain clearly visible noise. This demonstrates that the RQUNet-VAE with learned and deterministic frames increases sparsity in a generalized Besov space. Section III-D will demonstrate how this balance of noise reduction and discrimination of edges improves machine learning results when applying the RQUNet-VAE to segmentation of high resolution images for land cover mapping.

C. Image Decomposition Experiments

This section evaluates how RQUNet-VAE Scheme 2 decomposes and subsequently denoises images. RQUNet-VAE Scheme 2 uses an iterative scheme to decompose an image into (i) a lowpass image, (ii) a bandpass image, and (iii) a highpass image. Denoising can be achieved by truncating highpass features. The RQUNet-VAE was applied

to the Sentinel-2 image dataset described in Section III-A1. Figure 3 shows examples of the decomposition for four images. Figures 3(a,e) shows the original Sentinel2 images of a typical rural landscape with forest cover, cultivated fields and small, distributed structures. The corresponding lowpass images in Figures 3(b,f) captures the larger land cover parcels, without the fine-scale texture which has been removed. The corresponding bandpass in Figures 3(c,g) captures most of the signal of texture, whereas the highpass image in Figure Figures 3(d,h) captures very fine-scale texture, oscillating patterns, along with noise.

a) Time series decomposition: The RQUNet-VAE was applied to a Sentinel2 time series decomposition by using Haar wavelet smoothing in time as a diffusion process, (Section II-E). The proposed smoothing technique is therefore simultaneously applied in spatial domain (by RQUNet-VAE) and temporal domain (by Haar wavelet in time), following Algorithm 2 in SM. The smoothing parameter was set at $\alpha = 0.03$ for all decompositions.

The Sentinel2 time series was comprised as follows: $\left\{ \underline{\underline{f_i}} \right\}_{i=1}^{80}$, $\underline{\underline{f_i}} \in \mathbb{R}^{40 \times 40 \times 3 \times 99}$, where each of the 99 images have three channels of image size 40×40 . The time series is then padded to $\underline{\underline{f_i}} \in \mathbb{R}^{64 \times 64 \times 3 \times 120}$. The RQUNet-VAE is trained on each $\underline{\underline{f_i}}$ in the padded dataset to obtain all unknown model parameters. An Adam optimizer was used to train RQUNet-VAE with 200 epochs and batch-size of 16.

The trained parameters are then applied to Algorithm 2 with 10 iterations for spectral decomposition. Note that Algorithm 2 is an extension of generalized intersection algorithms with

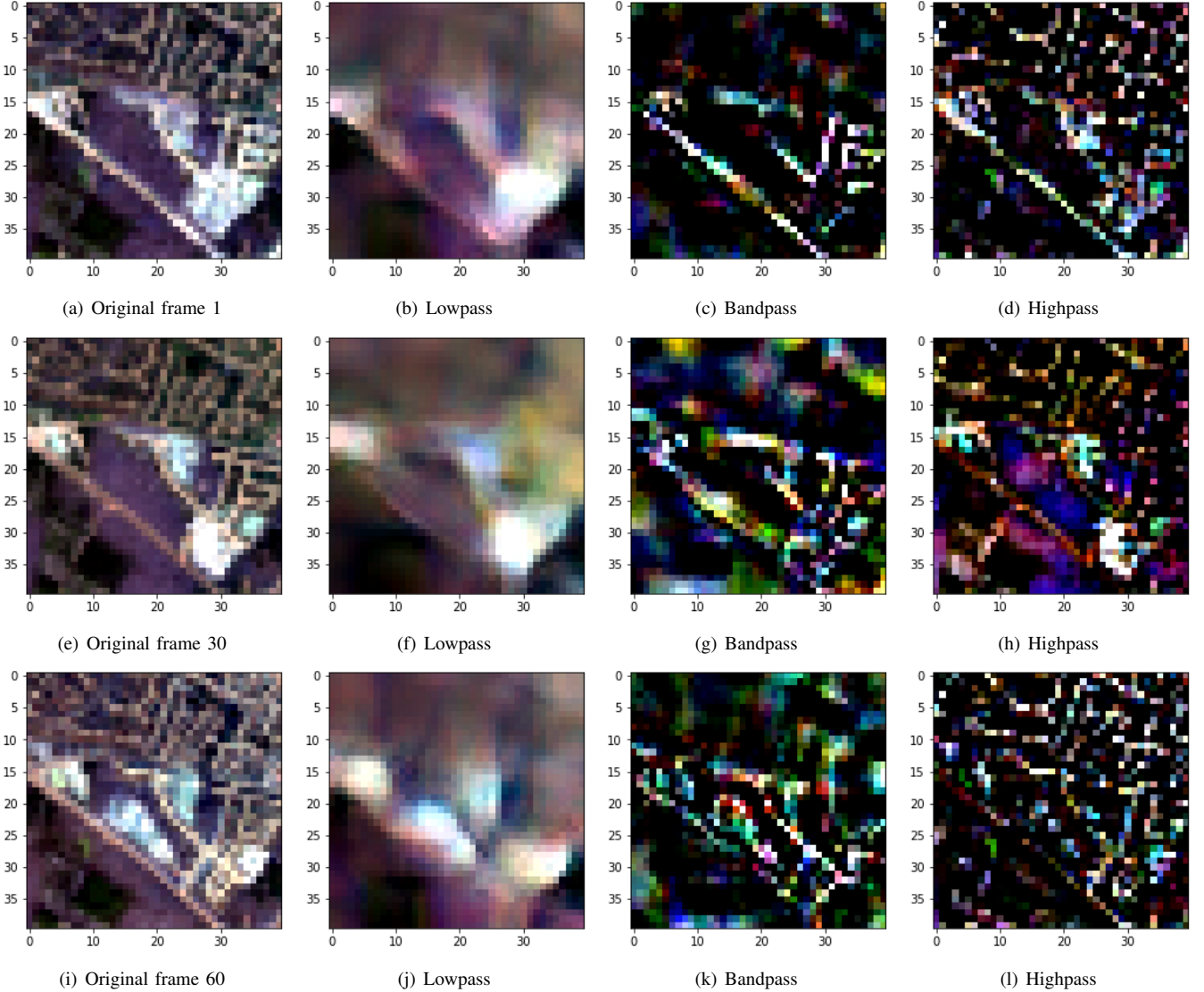


Fig. 4. Time series decomposition by RQUNet-VAE with GIAF based spectral decomposition for site 296.

fixpoints [23] for video processing via spectral decomposition into lowpass, bandpass and highpass videos.

Figure 4 illustrates the time series decomposition with the RQUNet-VAE and diffusion process. The smoothing parameter was set at $\alpha = 0.03$ for all decompositions. Similar to image decomposition in the previous sections, this time series decomposition extracts a homogeneous time series and a residual time series which contains small objects (e.g. roads), noise, and texture.

D. Image Segmentation by RQUNet-VAE

The RQUNet-VAE was applied to the task of high resolution aerial image segmentation. The goal is to demonstrate that the RQUNet-VAE makes the image segmentation more robust to noise compared to the existing U-Net architecture. We first provide a brief background to the problem of image segmentation in Section III-D1. Then, in Section III-D2 we apply 1) the traditional U-Net architecture and 2) the RQUNet-VAE to the problem of image segmentation and show that RQUNet-VAE yields better segmentation results where artificial noise is added to images. This shows that our proposed RQUNet-

VAE makes the existing U-Net architecture more robust to noise.

1) *Background on Image Segmentation:* The conventional image segmentation procedure has the following steps, (i) pre-processing to remove noise and unwanted small objects, (ii) segmentation, (iii) post-processing with morphological operators. This 3-step process requires a large number of parameters that need to be selected by an expert. Therefore, attempts have been made to incorporate a smoothing term into various models. The Mumford-Shah model [17] was proposed to introduce a smoothing term in a minimization approach for segmentation. However, computation of this smoothing term is an NP-hard problem and therefore not feasible [17]. Later, the Chan-Vese segmentation model [4] was proposed that was solvable by relying on the level set method [19]. By adding a smoothing term in a variational formulae, this model successfully segments objects in a noisy background. Inspired by Chan-Vese model [4], we use RQUNet-VAE for segmentation by introducing a smoothing term, i.e. Riesz-Quincunx wavelet truncation, directly into the UNet-VAE. The

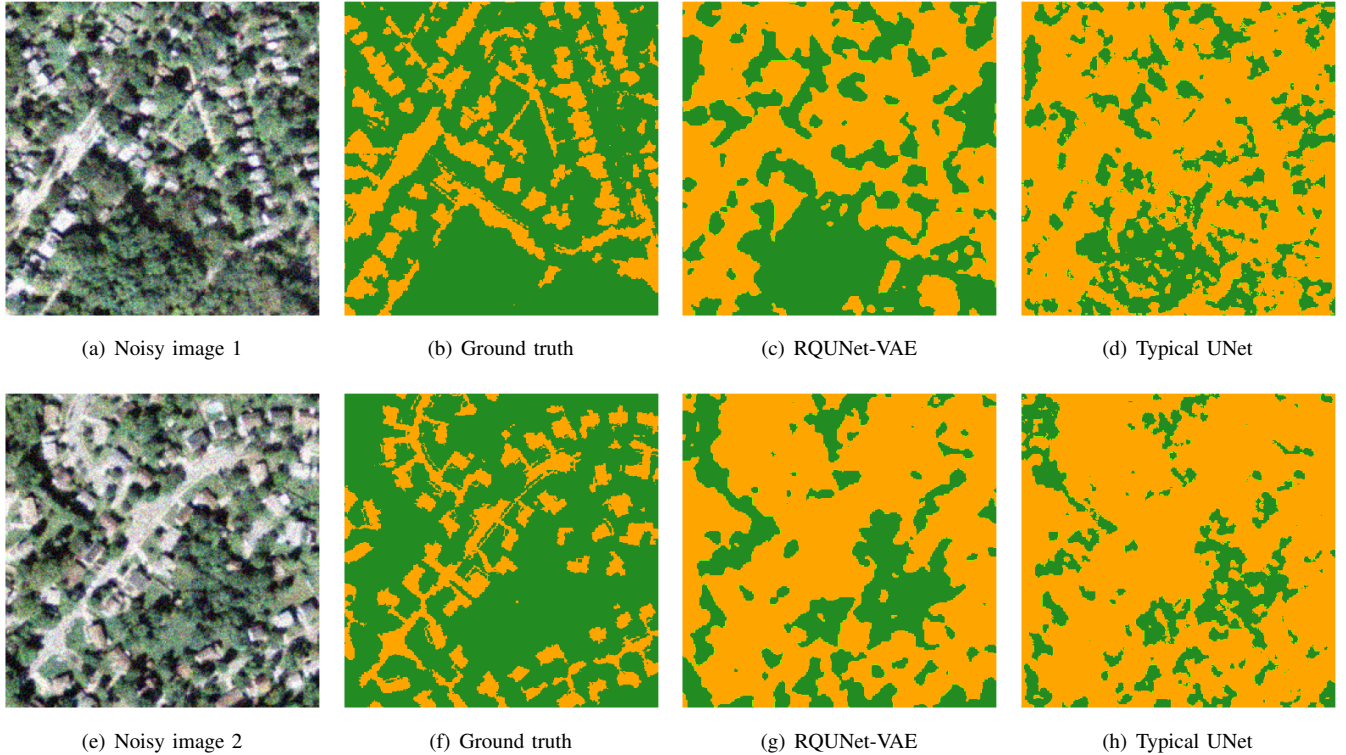


Fig. 5. Segmentation result with RQUNet-VAE as a smoothing term for noisy images. a,e) NAIP images with artificial noise added ($\text{std} = 0.08$) as input to segmentation; b,f) ground-truth segmentation masks; (c,g) segmentation masks returned by our RQUNet-VAE with a smoothing parameter $\alpha = 0.5$; (d,h) segmentation mask returned by the typical UNet architecture [24].

advantage is that the RQUNet-VAE includes the smoothing term inside the Unet and thus combines (1) the ability of a Unet to learn features from an image, with (2) denoising capabilities enabled by the Riesz-Quincunx wavelet truncation. This smoothing term should eliminate small scale objects in a segmented image, e.g. texture and noise. Then, there is no need for separate pre- or post-processing steps, as these are all performed by the RQUNet-VAE segmentation.

2) *Experimental Evaluation of RQUNet-VAE for Image Segmentation:* We apply both 1) the traditional U-Net [24] and 2) RQUNet-VAE to the problem of segmentation using the 20 NAIP images described in Section III-A1 using two classes of land cover, impervious and vegetated. Figure 5 shows the segmentation results for two sample images after noise has been added (Section III-B, standard deviation of 0.08). The corresponding ground truth land cover masks are shown in Figures 5(b) and 5(f). Table II provides quantitative segmentation results, including overall accuracy, precision, recall, and F1-score for both of the two classes. Due to the added noise, the segmentation accuracy of the U-Net is very low at 0.6640 and 0.6170 for the two images, respectively Table II. The precision and recall for the two classes (non-impervious and impervious), indicate that this low accuracy is attributed to a bias towards the impervious class, as suggested by the low recall for the non-impervious class (0.3780 for the traditional UNet). This implies that more than two out of five impervious pixels are incorrectly classified as vegetated.

The performance of the traditional UNet is much lower with the added noise than on the original, which indicates that the addition of the noise substantially confuses the U-

Net. In contrast, the RQUNet-VAE coped much better with the noise, yielding higher accuracy values of 0.7057 ± 0.0013 and 0.6634 ± 0.0012 for the two images, respectively. The negative impact of noise on the segmentation accuracy is substantially reduced when using the RQUNet-VAE which makes the U-Net more robust to noise. Furthermore, the RQUNet-VAE is better able to discriminate the vegetated class, shown by the higher F1-scores Table II.

A more comprehensive evaluation of RQUNet-VAE vs U-Net on 20 images with various levels of noise is given in Figure 6. This figure shows the mean accuracy (solid line) of the U-Net and the RQUNet-VAE across all 20 images. When using the original image with no noise ($\sigma = 0$), the RQUNet-VAE only provides a marginal improvement over U-Net. As the noise level σ is increased, both RQUNet-VAE and U-Net exhibit a drop in accuracy. However, that the drop in accuracy is substantially slower for RQUNet-VAE. In particular, for noise values around $\sigma = 0.1$, the RQUNet-VAE has up to 10% higher accuracy than U-Net. This improvement stems from the variational auto-encoder approach, which, through variational changes to the latent representation of an image, allows the identification of pixels with a high probability of belonging to one class while being assigned to another class by the deterministic U-Net. Figure 6 gives the results using two parameters setting of RQUNet-VAE, $\alpha = 0.5$ and $\alpha = 1.0$. In both cases the resulting accuracy is nearly identical (the two lines are almost perfectly on top of each other), showing that the the performance is not very sensitive to the choice of α and thus that RQUNet-VAE is robust to non-optimal

TABLE II

QUANTITATIVE COMPARISON BETWEEN RQUNET-VAE AND THE TRADITIONAL UNET ARCHITECTURE. SINCE RQUNET-VAE USES A VARIATIONAL (NON-DETERMINISTIC) APPROACH, WE PROVIDE MEAN AND STANDARD DEVIATION AGGREGATED OVER 20 RUNS. PRECISION, RECALL, AND F1-SCORE ARE PROVIDED FOR BOTH CLASSES [VEGETATED, IMPERVIOUS].

Algorithm	Accuracy (overall)	Precision (per class)	Recall (per class)	F1-Score (per class)
RQUNet-VAE Mean (Image 1)	0.7057	[0.9047 0.4498]	[0.5457 0.8660]	[0.6807 0.5921]
RQUNet-VAE Stdev (Image 1)	± 0.0013	$\pm [0.0008 \ 0.0014]$	$\pm [0.0030 \ 0.0016]$	$\pm [0.0024 \ 0.0012]$
Traditional UNet (Image 1)	0.6640	[0.9460 0.3960]	[0.3780 0.9500]	[0.5400 0.5590]
RQUNet-VAE Mean (Image 2)	0.6734	[0.9837 0.3258]	[0.3662 0.9806]	[0.5337 0.4890]
RQUNet-VAE Stdev (Image 2)	± 0.0012	$\pm [0.0006 \ 0.0009]$	$\pm [0.0027 \ 0.0007]$	$\pm [0.0028 \ 0.0010]$
Traditional UNet (Image 2)	0.6170	[0.9860 0.2900]	[0.2440 0.9890]	[0.3920 0.4490]

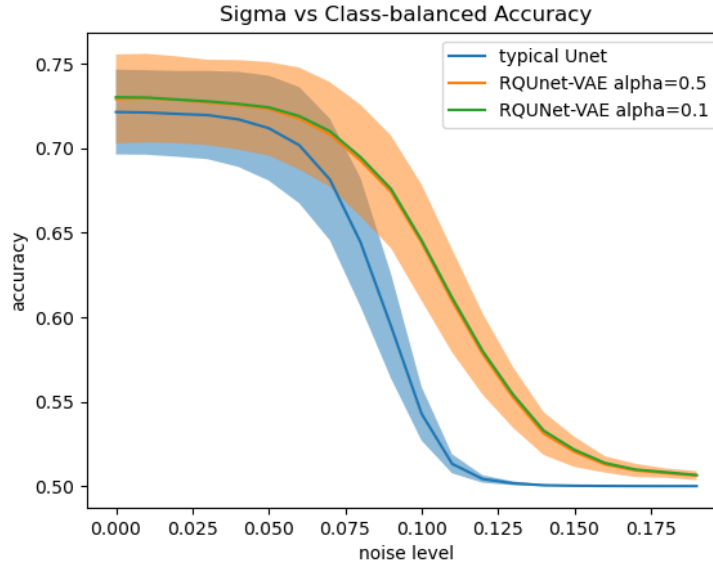


Fig. 6. Outcome of segmentation experiments using the RQUNet-VAE and traditional UNet applied to 20 NAIP images with increasing added noise (0-0.20). The confidence intervals (blue and orange bands) are calculated from the mean and standard deviation of accuracy of 50 prediction iterations for each noise level and alpha value (0.5, 0.1) in the RQ scheme. This illustrates the results of variational terms of RQUNet-VAE.

choices of α . In summary, Figure 6 shows that our proposed RQUNet-VAE is much more robust to Gaussian noise added to an image, as the reduction of segmentation accuracy with higher noise levels, is less pronounced.

IV. CONCLUSION

In this work, we introduced the RQUNet-VAE which augments the existing UNet architecture with a generalized wavelet expansion approach that we extended to a diffusion process to enable spectral decomposition. To the best of our knowledge this is the first approach that enables image decomposition and image denoising using a variational variant of the UNet architecture. An important application of this decomposition is denoising, achieved by truncating highpass features, that is, discarding information of decomposed images having the highest variance. We apply our proposed RQUNet-VAE to image denoising and segmentation of multi-band satellite images and their time-series which often contain noise due to multiple causes. During quantitative comparisons of noise reduction the RQUNet-VAE yields the highest PSNR and SSIM of all competitor methods. Among all the approaches based on harmonic analysis, RQUNet-VAE Scheme 2 yields the best results. For the application of satellite image denoising our proposed RQUNet-VAE provides superior qualitative performance compared to other competitors. The denoising

by RQUNet-VAE Scheme 1 was visually compared to the noise reduction using a Riesz Dyadic wavelet and curvelets and using wavelet CDF and it provided the best combination of noise reduction and delineation of edges and objects. Furthermore the propose RQUNet-VAE Scheme 2, was compared to other iterative methods including GIAF and directional TV-L2 and it resulted in the best trade-off between reduction of noise and delineation of edges and objects, whereas the GIAF oversmoothed edges between objects while TV-L2 clearly retained visible noise.

To quantitatively measure the improvement of the RQUNet-VAE against the traditional UNet, we applied it to high resolution aerial image segmentation. Our experiments show only a slight improvement over the traditional UNet for segmentation of the original images with little noise. However, as artificial noise is added to images, we observe that the segmentation quality of the UNet decreases more rapidly than when using the RQUNet-VAE. The superior performance of RQUNet-VAE is due to a neural network in UNet-VAE with a learned dictionary obtained from the training procedure to increase the level of sparsity for an input image in a generalized Besov space. This property is due to a fundamental concept in signal processing, that the signal is sparse in some transformed domain. This demonstrates that the RQUNet-VAE architecture is

able to substantially increase the robustness to noise compared to the traditional UNet for the application of satellite image segmentation.

V. ACKNOWLEDGMENTS

This research is based upon work supported in part by the Office of the Director of National Intelligence (ODNI), Intelligence Advanced Research Projects Activity (IARPA), Space-based Machine Automated Recognition Technique (SMART) program, via contract 2021-20111000003. The views and conclusions contained herein are those of the authors and should not be interpreted as necessarily representing the official policies, either expressed or implied, of ODNI, IARPA, or the U.S. Government. The U.S. Government is authorized to reproduce and distribute reprints for governmental purposes notwithstanding any copyright annotation therein. The research was conducted in collaboration with BlackSky. We thank Diego Torrejon and Isaac Corley of BlackSky for valuable comments on the manuscript.

REFERENCES

- [1] J. Allenby and C. Phelan. Implementing technology and precision conservation in the Chesapeake bay. *Chesapeake Conservancy*, 2013.
- [2] J.-F. Cai, R. H. Chan, and Z. Shen. A framelet-based image inpainting algorithm. *Applied and Computational Harmonic Analysis*, 24(2):131–149, 2008.
- [3] E. J. Candès and D. L. Donoho. New tight frames of curvelets and optimal representations of objects with piecewise c_2 singularities. *Communications on Pure and Applied Mathematics: A Journal Issued by the Courant Institute of Mathematical Sciences*, 57(2):219–266, 2004.
- [4] T. F. Chan and L. A. Vese. Active contours without edges. *IEEE Transactions on image processing*, 10(2):266–277, 2001.
- [5] M. Claverie, J. Ju, J. G. Masek, J. L. Dungan, E. F. Vermote, J.-C. Roger, S. V. Skakun, and C. Justice. The harmonized Landsat and Sentinel-2 surface reflectance data set. *Remote sensing of environment*, 219:145–161, 2018.
- [6] I. Daubechies and W. Sweldens. Factoring wavelet transforms into lifting steps. *Journal of Fourier analysis and applications*, 4(3):247–269, 1998.
- [7] P. Esser, E. Sutter, and B. Ommer. A variational U-Net for conditional appearance and shape generation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8857–8866, 2018.
- [8] M. Feilner, D. Van De Ville, and M. Unser. An orthogonal family of Quincunx wavelets with continuously adjustable order. *IEEE Transactions on Image Processing*, 14(4):499–510, 2005.
- [9] G. Gilboa. A total variation spectral framework for scale and texture analysis. *SIAM journal on Imaging Sciences*, 7(4):1937–1961, 2014.
- [10] T. Goldstein and S. Osher. The split Bregman method for L_1 -regularized problems. *SIAM journal on imaging sciences*, 2(2):323–343, 2009.
- [11] S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International conference on machine learning*, pages 448–456. PMLR, 2015.
- [12] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [13] D. P. Kingma and M. Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- [14] W. Kleynhans, B. P. Salmon, J. C. Olivier, F. Van den Bergh, K. J. Wessels, T. L. Grobler, and K. C. Steenkamp. Land cover change detection using autocorrelation analysis on MODIS time-series data: Detection of new human settlements in the gauteng province of south africa. *IEEE Journal of selected topics in applied earth observations and remote sensing*, 5(3):777–783, 2012.
- [15] S. Kohl, B. Romera-Paredes, C. Meyer, J. De Fauw, J. R. Ledsam, K. Maier-Hein, S. Eslami, D. Jimenez Rezende, and O. Ronneberger. A probabilistic U-Net for segmentation of ambiguous images. *Advances in neural information processing systems*, 31, 2018.
- [16] J. Li and D. P. Roy. A global analysis of Sentinel-2A, Sentinel-2B and Landsat-8 data revisit intervals and implications for terrestrial monitoring. *Remote Sensing*, 9(9):902, 2017.
- [17] D. B. Mumford and J. Shah. Optimal approximations by piecewise smooth functions and associated variational problems. *Communications on pure and applied mathematics*, 1989.
- [18] NAIP. National agriculture imagery program (NAIP) information sheet, 2015. https://www.fsa.usda.gov/Internet/FSA_File/naip_info_sheet_2015.pdf.
- [19] S. Osher and J. A. Sethian. Fronts propagating with curvature-dependent speed: Algorithms based on hamilton-jacobi formulations. *Journal of computational physics*, 79(1):12–49, 1988.
- [20] J. R. Partington, J. R. Partington, et al. *An introduction to Hankel operators*. Cambridge University Press, 1988.
- [21] N. G. Polson, J. G. Scott, and B. T. Willard. Proximal algorithms in statistics and machine learning. *Statistical Science*, 30(4):559–581, 2015.
- [22] Y. Pu, Z. Gan, R. Henao, X. Yuan, C. Li, A. Stevens, and L. Carin. Variational autoencoder for deep learning of images, labels and captions. *Advances in neural information processing systems*, 29, 2016.
- [23] R. Richter, D. H. Thai, and S. F. Huckemann. Generalized intersection algorithms with fixpoints for image decomposition learning. *CoRR*, abs/2010.08661, 2020.
- [24] 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*, pages 234–241. Springer, 2015.
- [25] D. P. Roy, H. Huang, L. Boschetti, L. Giglio, L. Yan, H. H. Zhang, and Z. Li. Landsat-8 and Sentinel-2 burned area mapping-A combined sensor multi-temporal change detection approach. *Remote Sensing of Environment*, 231:111254, 2019.
- [26] D. P. Roy, Y. Jin, P. Lewis, and C. Justice. Prototyping a global algorithm for systematic fire-affected area mapping using MODIS time series data. *Remote sensing of environment*, 97(2):137–162, 2005.
- [27] L. I. Rudin, S. Osher, and E. Fatemi. Nonlinear total variation based noise removal algorithms. *Physica D: nonlinear phenomena*, 60(1-4):259–268, 1992.
- [28] B. P. Salmon, W. Kleynhans, F. Van den Bergh, J. C. Olivier, T. L. Grobler, and K. J. Wessels. Land cover change detection using the internal covariance matrix of the extended Kalman filter over multiple spectral bands. *IEEE Journal of selected topics in applied earth observations and remote sensing*, 6(3):1079–1085, 2013.
- [29] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *The journal of machine learning research*, 15(1):1929–1958, 2014.
- [30] W. Sweldens. The lifting scheme: A construction of second generation wavelets. *SIAM journal on mathematical analysis*, 29(2):511–546, 1998.
- [31] D. H. Thai and C. Gottschlich. Directional global three-part image decomposition. *EURASIP Journal on Image and Video Processing*, 2016(1):1–20, 2016.
- [32] C. Tian, L. Fei, W. Zheng, Y. Xu, W. Zuo, and C.-W. Lin. Deep learning on image denoising: An overview. *Neural Networks*, 131:251–275, 2020.
- [33] M. Unser and T. Blu. Mathematical properties of the jpeg2000 wavelet filters. *IEEE transactions on image processing*, 12(9):1080–1090, 2003.
- [34] M. Unser, N. Chenouard, and D. Van De Ville. Steerable pyramids and tight wavelet frames in $L_2(\mathbb{R}^d)$. *IEEE Transactions on Image Processing*, 20(10):2705–2721, Oct. 2011.
- [35] M. Unser, D. Sage, and D. Van De Ville. Multiresolution monogenic signal analysis using the Riesz-Laplace wavelet transform. *IEEE Transactions on Image Processing*, 18(11):2402–2418, Nov. 2009.
- [36] M. Unser and D. Van De Ville. Wavelet steerability and the higher-order Riesz transform. *IEEE Transactions on Image Processing*, 19(3):636–652, Mar. 2010.
- [37] F. Van Den Bergh, K. J. Wessels, S. Miteff, T. L. Van Zyl, A. D. Gazendam, and A. K. Bachoo. Hitempo: A platform for time-series analysis of remote-sensing satellite data in a high-performance computing environment. *International Journal of Remote Sensing*, 33(15):4720–4740, 2012.
- [38] T. Van Erven and P. Harremos. Rényi divergence and Kullback-Leibler divergence. *IEEE Transactions on Information Theory*, 60(7):3797–3820, 2014.
- [39] E. Vermote, C. Justice, M. Claverie, and B. Franch. Preliminary analysis of the performance of the Landsat 8/OLI land surface reflectance product. *Remote Sensing of Environment*, 185:46–56, 2016.
- [40] M. Vetterli, P. Marziliano, and T. Blu. Sampling signals with finite rate of innovation. *IEEE transactions on Signal Processing*, 50(6):1417–1428, 2002.
- [41] Z. Wang and A. C. Bovik. Mean squared error: Love it or leave it? a new look at signal fidelity measures. *IEEE signal processing magazine*, 26(1):98–117, 2009.
- [42] C. E. Woodcock, T. R. Loveland, M. Herold, and M. E. Bauer. Transitioning from change detection to monitoring with remote sensing: A paradigm shift. *Remote Sensing of Environment*, 238:111558, 2020.

- [43] J. C. Ye, Y. Han, and E. Cha. Deep convolutional framelets: A general deep learning framework for inverse problems. *SIAM Journal on Imaging Sciences*, 11(2):991–1048, 2018.
- [44] 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*, 10(2):711–750, 2017.
- [45] H. K. Zhang, D. P. Roy, L. Yan, Z. Li, H. Huang, E. Vermote, S. Skakun, and J.-C. Roger. Characterization of Sentinel-2A and Landsat-8 top of atmosphere, surface, and nadir BRDF adjusted reflectance and NDVI differences. *Remote sensing of environment*, 215:482–494, 2018.
- [46] Z. Zhu. Change detection using Landsat time series: A review of frequencies, preprocessing, algorithms, and applications. *ISPRS Journal of Photogrammetry and Remote Sensing*, 130:370–384, 2017.

Riesz-Quincunx-Unet Variational Auto-Encoder for Satellite Image Denoising: Supplemental Material

Duy H. Thai¹, Xiqi Fei¹, Tri M. Le¹, Andreas Züfle², Konrad Wessels¹

¹George Mason University, Department of Geography and Geoinformation Science, USA

²Emory University, Department of Computer Science, USA

1 Additional Images

This section provides details on the ten satellite images used for testing of our proposed RQUNet-VAE approach. Figure 1 shows the original images, Figure 2 shows the same images with artificial noise added, and Figure 3 shows the denoised images using our proposed RQUNet-VAE.

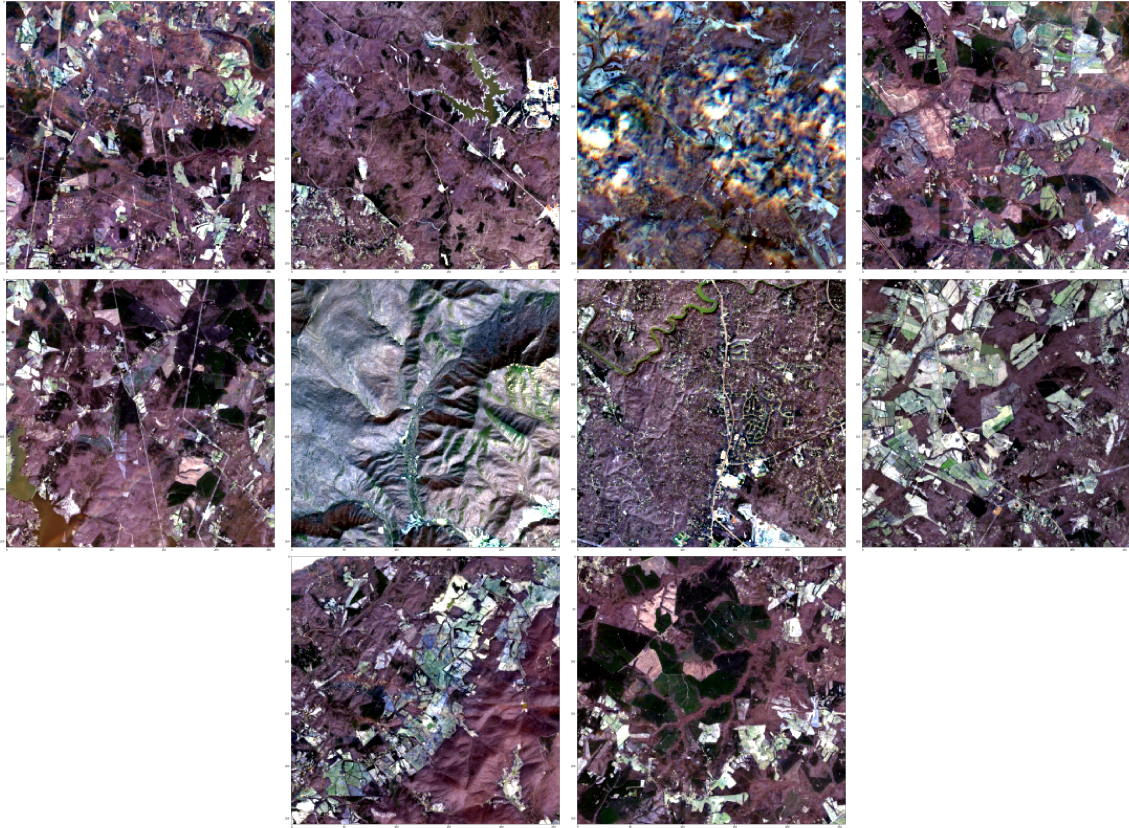


Figure 1: Original images.

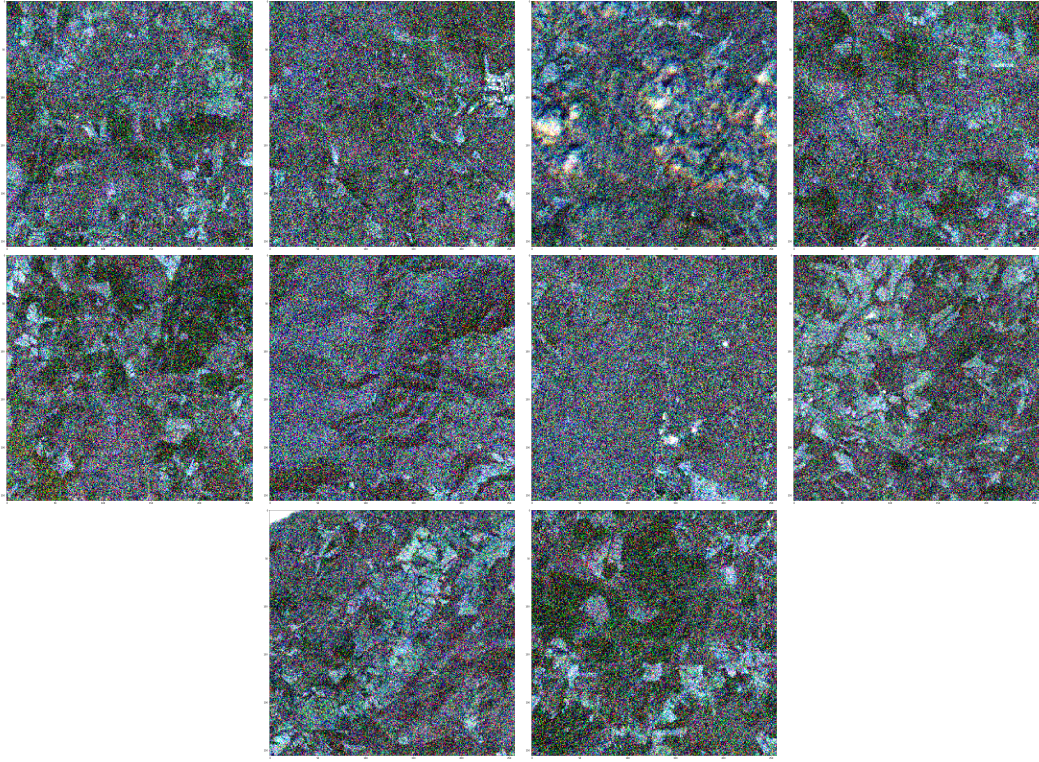


Figure 2: Noisy images with variance $\sigma = 0.04$.

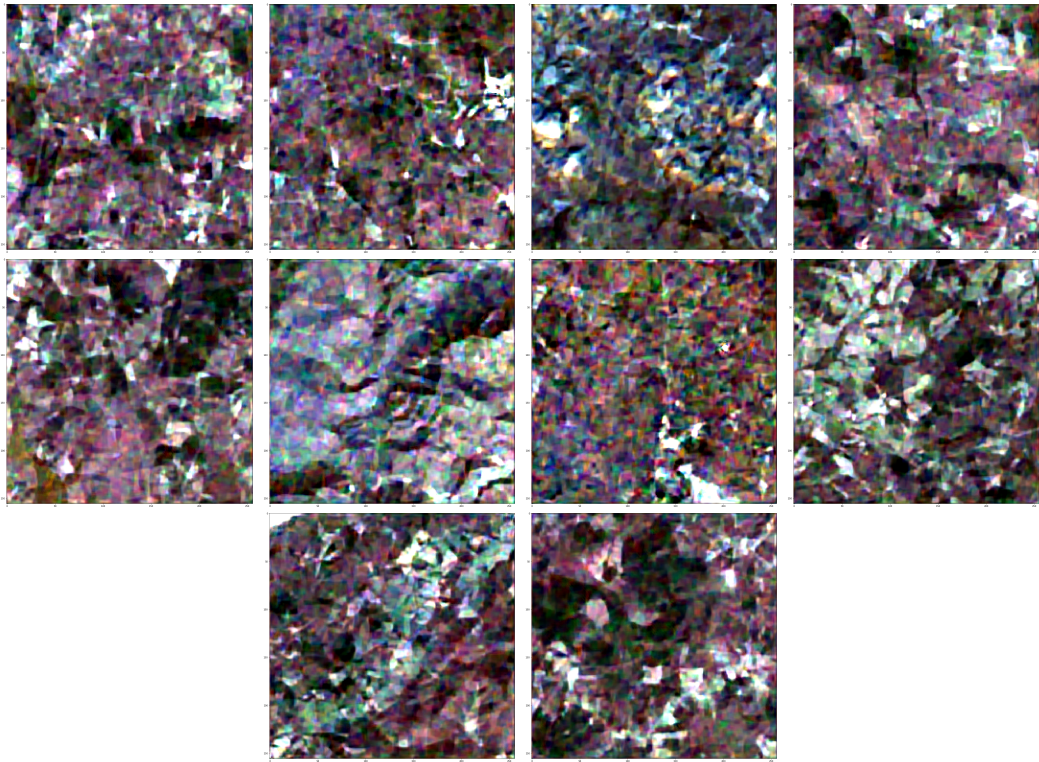


Figure 3: Denoised images for $\sigma = 0.04$ by RQUnet-VAE smoothing.

2 Riesz Quincunx Filter Banks

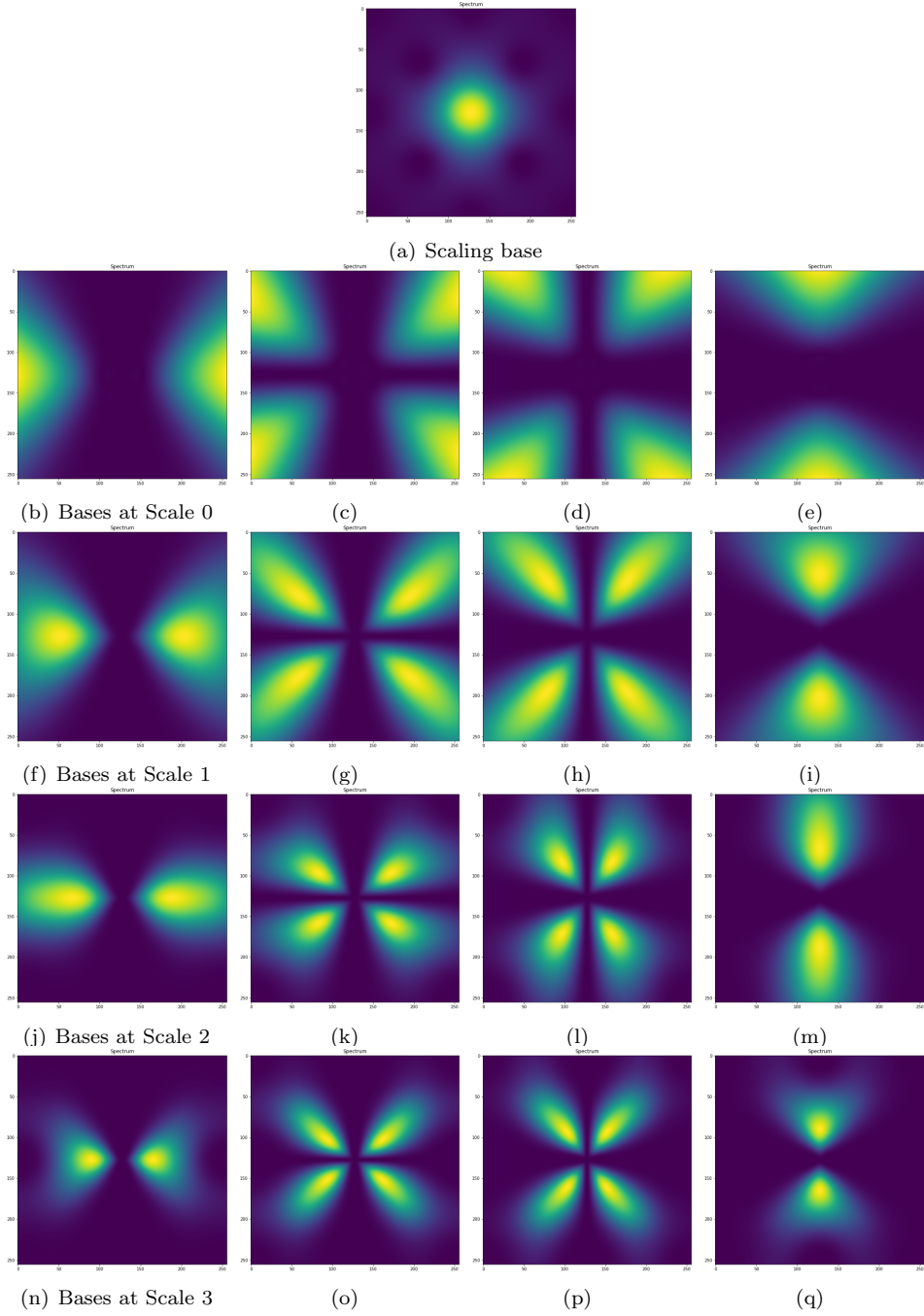


Figure 4: This figure illustrates N -th order Riesz Quincunx Filter banks in the Fourier domain with 3 scales, a fractional order of Bspline $\gamma = 1.2$ and order of Riesz transform $N = 3$.

Figure 4 illustrates the N -th order Riesz Quincunx filter banks in the Fourier domain: rows are for 4 wavelet scales and columns are 4 directions for each scale corresponding to $N = 0, \dots, 3$ orders in Riesz transform.

3 Detailed Encoder and Decoder of RQUNet-VAE

Figure 5 shows the detailed architecture of the RQUNet-VAE encoder (left) and decoder (right).

4 Additional Background and definitions

4.1 Hankel matrices and convolutional operations

The following list provides definition of Hankel matrix and its inverse and a set of convolution operations used in this work.

- Extended Hankel matrix of an image $\underline{f} \in \mathbb{R}^{|\Omega|}$:

$$\mathcal{H}_{d_1|n_2}(\underline{f}) = \begin{pmatrix} \mathcal{H}_{d_1}(f_1) & \dots & \mathcal{H}_{d_1}(f_{n_2}) \end{pmatrix} \in \mathbb{R}^{n_1 \times n_2 d_1},$$

with a Hankel matrix:

$$\mathcal{H}_{d_1}(f_i) = \begin{pmatrix} f_i[1] & \dots & f_i[d_1] \\ \vdots & \ddots & \vdots \\ f_i[n_1] & \dots & f_i[n_1 + d_1 - 1] \end{pmatrix} \in \mathbb{R}^{n_1 \times d_1}.$$

- The extended Hankel matrix of a multi-channel image $\underline{\underline{f}} \in \mathbb{R}^{|\Omega| \times P}$:

$$\mathcal{H}_{n_1|d_2|P}(\underline{\underline{f}}) = \begin{pmatrix} \mathcal{H}_{n_1|d_2}(\underline{f}_1) & \dots & \mathcal{H}_{n_1|d_2}(\underline{f}_P) \end{pmatrix}.$$

Given $\{e_k\}_{k=1}^{n_1}$ as an orthonormal basis of \mathbb{R}^{n_1} , $\{\tilde{e}_k = \frac{1}{\sqrt{d_1}} \mathcal{H}_{d_1}(e_k)\}_{k=1}^{n_1} \subset \mathbb{R}^{n_1 \times d_1}$ is the orthonormal basis of $\mathcal{H}(n_1, d_1) = \{h : h \in \mathbb{R}^{n_1 \times d_1}\}$, i.e. $\langle \tilde{e}_k, \tilde{e}_l \rangle_F = \delta_{kl}$ and an orthonormal expansion:

$$\mathcal{H}_{d_1}(f_i) = \sum_{k=1}^{n_1} \underbrace{\langle \tilde{e}_k, \mathcal{H}_{d_1}(f_i) \rangle_F}_{=\sqrt{d_1} f_i[k]} \tilde{e}_k.$$

- The inverse of an extended Hankel matrices for a matrix $\underline{g} = \begin{pmatrix} \underline{g}_1 & \dots & \underline{g}_{n_2} \end{pmatrix} \in \mathbb{R}^{n_1 \times d_1 n_2}$, $\underline{g}_i \in \mathbb{R}^{n_1 \times d_1}$ as:

$$\mathcal{H}_{d_1|n_2}^\dagger(\underline{g}) = \begin{pmatrix} \mathcal{H}_{d_1}^\dagger(\underline{g}_1) & \dots & \mathcal{H}_{d_1}^\dagger(\underline{g}_{n_2}) \end{pmatrix} \in \mathbb{R}^{|\Omega|}, \quad (1)$$

having an inverse Hankel matrix:

$$\mathcal{H}_{d_1}^\dagger(\underline{g}_i) = \frac{1}{\sqrt{d_1}} \begin{pmatrix} \langle \tilde{e}_1, \underline{g}_i \rangle_F \\ \vdots \\ \langle \tilde{e}_{n_1}, \underline{g}_i \rangle_F \end{pmatrix} \in \mathbb{R}^{n_1}. \quad (2)$$

- Filter banks and their family of matrices for the Riesz-Quincunx wavelet:

- scaling filter bank $\underline{\phi} \in \mathbb{R}^{d_1 \times d_2}$,
- primal wavelet filter bank $\underline{\underline{\psi}} = \{\underline{\psi}_1, \dots, \underline{\psi}_P\} \in \mathbb{R}^{d_1 \times d_2 \times P}$, $\underline{\psi}_i \in \mathbb{R}^{d_1 \times d_2}$,

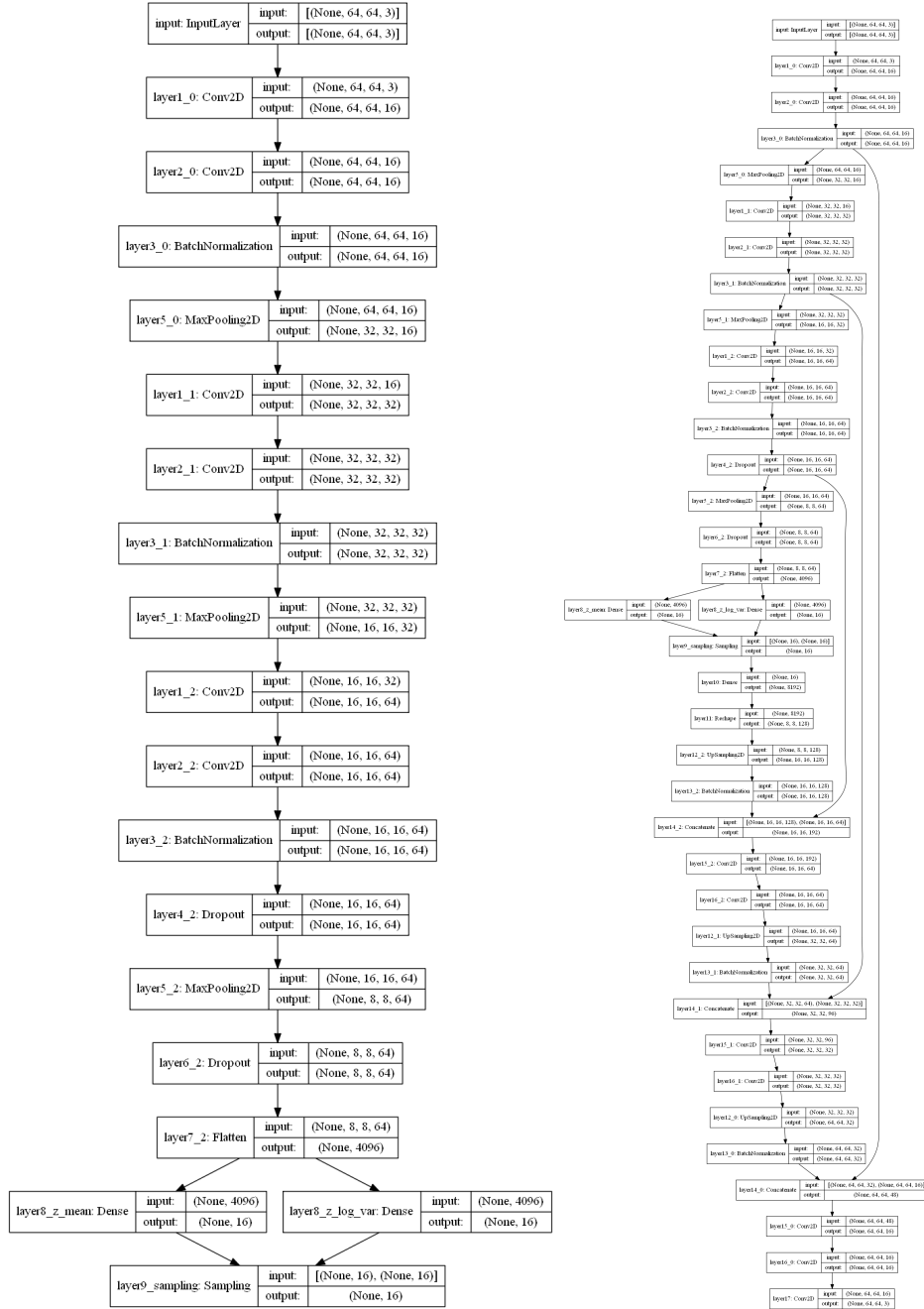


Figure 5: RQUnet-VAE Encoder (left) and Decoder (right)

- dual wavelet filter bank $\underline{\tilde{\psi}} = \{\underline{\tilde{\psi}}_1, \dots, \underline{\tilde{\psi}}_P\} \in \mathbb{R}^{d_1 \times d_2 \times P}$, $\underline{\tilde{\psi}}_i \in \mathbb{R}^{d_1 \times d_2}$,
- having matrix forms
 $\underline{\underline{\Psi}} = \{\underline{\Psi}_1, \dots, \underline{\Psi}_P\}$, $\underline{\underline{\tilde{\Psi}}} = \{\underline{\tilde{\Psi}}_1, \dots, \underline{\tilde{\Psi}}_P\}$ such that for $p \in 1, \dots, P$:

$$\underline{\Psi}_p = \begin{pmatrix} \underline{\Psi}_p^1 \\ \vdots \\ \underline{\Psi}_p^{d_2} \end{pmatrix} = \begin{pmatrix} \psi_{p,1}^1 & \dots & \psi_{p,d_2}^1 \\ \vdots & \ddots & \vdots \\ \psi_{p,1}^{d_2} & \dots & \psi_{p,d_2}^{d_2} \end{pmatrix}, \underline{\tilde{\Psi}}_p = \begin{pmatrix} \underline{\tilde{\Psi}}_p^1 \\ \vdots \\ \underline{\tilde{\Psi}}_p^{d_2} \end{pmatrix} = \begin{pmatrix} \tilde{\psi}_{p,1}^1 & \dots & \tilde{\psi}_{p,d_2}^1 \\ \vdots & \ddots & \vdots \\ \tilde{\psi}_{p,1}^{d_2} & \dots & \tilde{\psi}_{p,d_2}^{d_2} \end{pmatrix} \in \mathbb{R}^{d_2 d_1 \times d_2} \quad (3)$$

and

$$\underline{\Phi} = \begin{pmatrix} \underline{\Phi}^1 \\ \vdots \\ \underline{\Phi}^{d_2} \end{pmatrix} = \begin{pmatrix} \phi_1^1 & \dots & \phi_{d_2}^1 \\ \vdots & \ddots & \vdots \\ \phi_1^{d_2} & \dots & \phi_{d_2}^{d_2} \end{pmatrix}, \underline{\tilde{\Phi}} = \begin{pmatrix} \underline{\tilde{\Phi}}^1 \\ \vdots \\ \underline{\tilde{\Phi}}^{d_2} \end{pmatrix} = \begin{pmatrix} \tilde{\phi}_1^1 & \dots & \tilde{\phi}_{d_2}^1 \\ \vdots & \ddots & \vdots \\ \tilde{\phi}_1^{d_2} & \dots & \tilde{\phi}_{d_2}^{d_2} \end{pmatrix} \in \mathbb{R}^{d_2 d_1 \times d_2} \quad (4)$$

with $\phi_{i_2}^{i_1}[i_3] := \phi[i_3, i_2 - i_1]$, $\psi_{p,i_2}^{i_1}[l] := \psi_p[i_3, i_2 - i_1]$, $\tilde{\psi}_{p,i_2}^{i_1}[i_3] := \tilde{\psi}_p[i_3, i_2 - i_1]$ and $\phi_{i_2}^{i_1}, \psi_{p,i_2}^{i_1}, \tilde{\psi}_{p,i_2}^{i_1} \in \mathbb{R}^{d_1}$, $i_1 = 1, \dots, d_2$.

- Local bases are defined as: $\xi_i, \tilde{\xi}_i \in \mathbb{R}^{n_1}$,

$$\underline{\Xi} = (\xi_1 \quad \dots \quad \xi_d), \quad \underline{\tilde{\Xi}} = (\tilde{\xi}_1 \quad \dots \quad \tilde{\xi}_d) \in \mathbb{R}^{n_1 \times d}, \quad (5)$$

- Filter banks and their family of matrices for Unet-VAE:

$$\underline{\underline{\theta}} = \{\underline{\underline{\theta}}_1, \dots, \underline{\underline{\theta}}_Q\} \in \mathbb{R}^{d_1 \times d_2 \times P \times Q}, \underline{\underline{\theta}}_q = \{\underline{\theta}_{q,1}, \dots, \underline{\theta}_{q,P}\} \in \mathbb{R}^{d_1 \times d_2 \times P}, \underline{\theta}_{q,p} \in \mathbb{R}^{d_1 \times d_2},$$

whose matrix form is:

$$\underline{\underline{\Theta}} = \{\underline{\Theta}_1, \dots, \underline{\Theta}_Q\} \in \mathbb{R}^{d_1 d_2 P \times d_2 \times Q}, \underline{\Theta}_q = \begin{pmatrix} \underline{\Theta}_{q,1} \\ \vdots \\ \underline{\Theta}_{q,P} \end{pmatrix} \in \mathbb{R}^{d_1 d_2 P \times d_2}, \quad (6)$$

$$\underline{\Theta}_{q,p} = \begin{pmatrix} \underline{\Theta}_{q,p}^1 \\ \vdots \\ \underline{\Theta}_{q,p}^{d_2} \end{pmatrix} = \begin{pmatrix} \theta_{q,p,1}^1 & \dots & \theta_{q,p,d_2}^1 \\ \vdots & \ddots & \vdots \\ \theta_{q,p,1}^{d_2} & \dots & \theta_{q,p,d_2}^{d_2} \end{pmatrix} = (\theta_{q,p,1} \quad \dots \quad \theta_{q,p,d_2}) \in \mathbb{R}^{d_1 d_2 \times d_2},$$

$$\theta_{q,p,i_2}^{i_1}[i_3] := \theta_{q,p}[i_3, i_2 - i_1], \quad i_1 = 1, \dots, d_2, \quad \theta_{q,p,i_2}^{i_1} \in \mathbb{R}^{d_1}.$$

- Convolution Operations: Given filter bank matrices $(\underline{\phi}, \underline{\psi}, \underline{\theta})$ with their matrix forms $(\underline{\Phi}, \underline{\Psi}, \underline{\Theta})$ in (Equations 3-6) respectively, we have convolution operations for a matrix $\underline{f} \in \mathbb{R}^{n_1 \times n_2}$ and a multi-band image $\underline{\underline{f}} = \{\underline{f}_1, \dots, \underline{f}_P\} \in \mathbb{R}^{|\Omega| \times P}$:

– **1D convolution:**

$$\mathfrak{C}_{\phi_{i_1}^{i_2}} : \mathbb{R}^{n_1} \rightarrow \mathbb{R}^{n_1}; \quad \mathfrak{C}_{\phi_{i_1}^{i_2}}(f_{i_2}) = \left(\sum_{k_1=1}^{d_1} f_{i_2}[k_1] \check{\phi}_{i_1}^{i_2}[k_2 - k_1] \right)_{k_2=1}^{n_1} = \mathcal{H}_{d_1}(f_{i_2}) \phi_{i_1}^{i_2}, \quad (7)$$

– **2D convolution:**

$$\underline{\mathfrak{C}}_{\underline{\phi}} : \mathbb{R}^{n_1 \times n_2} \rightarrow \mathbb{R}^{n_1 \times d_2}; \quad \underline{\mathfrak{C}}_{\underline{\phi}}(\underline{f}) = \left(\sum_{k_1=1}^{d_1} \sum_{k_2=1}^{d_2} f[k_1, k_2] \check{\phi}[r_1 - k_1, r_2 - k_2] \right)_{\substack{r_1=1, \dots, n_1 \\ r_2=1, \dots, d_2}}, \quad (8)$$

– **Matrix-family convolution:**

$$\underline{\mathfrak{C}}_{\underline{\psi}} : \mathbb{R}^{n_1 \times n_2} \rightarrow \mathbb{R}^{n_1 \times d_2 \times P}; \quad \underline{\mathfrak{C}}_{\underline{\psi}}(\underline{f}) = \left\{ \underline{\mathfrak{C}}_{\underline{\psi}_1}(\underline{f}), \dots, \underline{\mathfrak{C}}_{\underline{\psi}_P}(\underline{f}) \right\} = \mathcal{H}_{n_1|d_2}(\underline{f}) \underline{\Psi}, \quad (9)$$

– **Anisotropic matrix-family convolution:**

$$\underline{\mathfrak{C}}_{\underline{\psi}}^{\text{ani}} : \mathbb{R}^{n_1 \times n_2 \times P} \rightarrow \mathbb{R}^{n_1 \times d_2 \times P};$$

$$\underline{\mathfrak{C}}_{\underline{\psi}}^{\text{ani}}(\underline{f}) = \left\{ \underline{\mathfrak{C}}_{\underline{\psi}_1}(\underline{f}_1), \dots, \underline{\mathfrak{C}}_{\underline{\psi}_P}(\underline{f}_P) \right\} = \left\{ \mathcal{H}_{n_1|d_2}(\underline{f}_1) \underline{\Psi}_1, \dots, \mathcal{H}_{n_1|d_2}(\underline{f}_P) \underline{\Psi}_P \right\}.$$

– **Isotropic matrix-family convolution:**

$$\underline{\mathfrak{C}}_{\underline{\Theta}}^{\text{iso}} : \mathbb{R}^{n_1 \times n_2 \times P} \rightarrow \mathbb{R}^{n_1 \times d_2 \times Q}, \quad \underline{\tilde{f}} = \left\{ \underline{\tilde{f}}_1, \dots, \underline{\tilde{f}}_Q \right\} = \underline{\mathfrak{C}}_{\underline{\Theta}}^{\text{iso}}(\underline{f}) = \mathcal{H}_{n_1|d_2|P}(\underline{f}) \underline{\Theta}$$

with:

$$\underline{\tilde{f}}_q = \sum_{p=1}^P \underline{\mathfrak{C}}_{\underline{\Theta}_{q,p}}(\underline{f}_p) = \sum_{p=1}^P \mathcal{H}_{n_1|d_2}(\underline{f}_p) \underline{\Theta}_{q,p} = \mathcal{H}_{n_1|d_2|P}(\underline{f}) \underline{\Theta}_q.$$

Then, we have the following relation for 2D convolution operation:

Proposition 4.1. *A 2D convolution is defined by Hankel matrix by describing it via a 1D convolution:*

$$\underline{\mathfrak{C}}_{\underline{\phi}}(\underline{f})[k_1, k_2] = \sum_{i=1}^{d_2} \mathfrak{C}_{\phi_{k_2}^i}(f_i)[k_1] = \sum_{i=1}^{d_2} \left(\mathcal{H}_{d_1}(f_i) \phi_{k_2}^i \right)[k_1], \quad k_1 = 1, \dots, n_1, \quad k_2 = 1, \dots, d_2$$

which is written in a matrix form as:

$$\underline{\mathfrak{C}}_{\underline{\phi}}(\underline{f}) = \mathcal{H}_{n_1|d_2}(\underline{f}) \underline{\Phi}. \quad (10)$$

Proof. We provide a proof of Proposition 4.1 in Section 5.9. \square

Moreover, note that we have an adjoint operator $\underline{\mathfrak{C}}_{\underline{\psi}}^* = \sum_{p=1}^P \mathfrak{C}_{\psi_p}^* : \mathbb{R}^{n_1 \times d_2 \times P} \rightarrow \mathbb{R}^{n_1 \times n_2}$. And, for $d_1 = n_1, d_2 = n_2$, we have $\underline{\mathfrak{C}}_{\underline{\phi}}(\underline{f}) = \mathcal{F}^{-1}(\hat{\underline{P}} \odot \hat{\underline{F}})$, $\hat{\underline{P}} = [\hat{P}(e^{j\omega})]_{\omega \in [-\pi, \pi]^d}$ and $\hat{\underline{F}} = [\hat{F}(e^{j\omega})]_{\omega \in [-\pi, \pi]^d}$. Its adjoint operator $\underline{\mathfrak{C}}_{\underline{\phi}}^* = \underline{\mathfrak{C}}_{\check{\phi}} : \mathbb{R}^{n \times m} \rightarrow \mathbb{R}^{n \times m}$ is defined with a discrete time-reversed kernel

$$\check{\phi}(x_1, x_2) = \phi(-x_1, -x_2) \xleftrightarrow{\mathcal{F}} \hat{\phi}(\omega_1, \omega_2) = \hat{\phi}^*(\omega_1, \omega_2)$$

whose discretized version is:

$$\check{\phi}[k_1, k_2] = \phi[-k_1, -k_2] \xleftrightarrow{\mathcal{F}} \hat{P}^*(e^{j\omega_1}, e^{j\omega_2}) = \hat{P}(e^{-j\omega_1}, e^{-j\omega_2}).$$

4.2 Proximal Operators and Moreau-Yosida envelope

Given a non-smooth function $\mathcal{P}(\cdot)$ which can be convex or non-convex, its Moreau-Yosida envelope ($\frac{1}{\mu} > 0$ -Lipschitz differentiable) is

$$\mathcal{P}_\mu(\cdot) := \inf_{\underline{u} \in \mathbb{R}^{n_1 \times n_2}} \left\{ \mathcal{P}(\underline{u}) + \frac{1}{2\mu} \|\underline{u} - \cdot\|_F^2 \right\}, \quad (11)$$

with $\lim_{\mu \rightarrow 0} \mathcal{P}_\mu(\theta) = \mathcal{P}(\theta)$ and $\|\cdot\|_F$ is Frobenius norm. Its gradient

$$\nabla \mathcal{P}_\mu(\cdot) = \frac{1}{\mu} \left(\cdot - \text{prox}_\mu \mathcal{P}(\cdot) \right) : \mathbb{R}^{n_1 \times n_2} \rightarrow \mathbb{R}^{n_1 \times n_2} \quad (12)$$

is defined by a proximity operator

$$\text{prox}_\mu \mathcal{P}(\cdot) := \underset{\underline{u} \in \mathbb{R}^{n_1 \times n_2}}{\text{argmin}} \left\{ \mathcal{P}(\underline{u}) + \frac{1}{2\mu} \|\underline{u} - \cdot\|_F^2 \right\}. \quad (13)$$

We also note that some activation functions in neural network can be well written as a proximal operator associated with its abstract function, e.g. Rectified Linear Unit (ReLU) activation function:

$$\begin{aligned} \text{prox}_{\text{ReLU}}(\underline{f}) &= \underset{\underline{u} \in \mathbb{R}^{n_1 \times n_2}}{\text{argmin}} \left\{ \mathcal{P}_{\text{ReLU}}(\underline{u}) + \frac{1}{2} \|\underline{u} - \underline{f}\|_F^2 \right\} \\ &= \left[\max(0, f_{i_1, i_2}) \right]_{i_1, i_2=0}^{n_1-1, n_2-1}, \quad \underline{f} = [f_{i_1, i_2}]_{i_1, i_2=0}^{n_1-1, n_2-1}, \end{aligned}$$

where $\mathcal{P}_{\text{ReLU}}(\cdot)$ is a non-defined function.

4.3 Kullback Leibler divergence:

Given distributions $\mathbf{F}(\text{d}z) = \mathbf{f}(z)\text{d}z$ and $\mathbf{G}(\text{d}z) = \mathbf{g}(z)\text{d}z$ on a domain \mathbb{R}^d , their KL-divergence is:

$$\begin{aligned} \text{KL}(\mathbf{F} \parallel \mathbf{G}) &= \mathbb{E}_{Z \sim \mathbf{F}} \left[\log \frac{\mathbf{f}(Z)}{\mathbf{g}(Z)} \right] = \int_{\mathbb{R}^d} \mathbf{f}(z) \log \frac{\mathbf{f}(z)}{\mathbf{g}(z)} \text{d}z \\ &\geq 0. \end{aligned} \quad (14)$$

5 Proofs

5.1 Proof of Proposition II.1

Due to the unity condition (Equation 1) and $\mathcal{H}_{d_1|n_2}^\dagger \circ \mathcal{H}_{d_1|n_2} = \text{Id}$, convolutional framelet decomposition is:

$$\underline{f} = \mathcal{H}_{d_1|n_2}^\dagger \left(\tilde{\Xi} \Xi^\text{T} \mathcal{H}_{d_1|n_2} \left(\underline{f} \right) \Phi \tilde{\Phi}^\text{T} \right)$$

where:

$$\begin{aligned} c_f &= \Xi^\text{T} \mathcal{H}_{d_1|n_2} \left(\underline{f} \right) \Phi = \begin{pmatrix} \xi_1^\text{T} \\ \vdots \\ \xi_d^\text{T} \end{pmatrix} \mathcal{H}_{d_1|n_2} \left(\underline{f} \right) (\phi_1 \quad \dots \quad \phi_{d_2}) = \begin{pmatrix} \xi_1^\text{T} \mathcal{H}_{d_1|n_2} \left(\underline{f} \right) \phi_1 & \dots & \xi_1^\text{T} \mathcal{H}_{d_1|n_2} \left(\underline{f} \right) \phi_{d_2} \\ \vdots & \ddots & \vdots \\ \xi_d^\text{T} \mathcal{H}_{d_1|n_2} \left(\underline{f} \right) \phi_1 & \dots & \xi_d^\text{T} \mathcal{H}_{d_1|n_2} \left(\underline{f} \right) \phi_{d_2} \end{pmatrix} \\ &= (c_{f,1} \quad \dots \quad c_{f,d_2}) = [c_f^{l,s}]_{l=1,\dots,d}^{s=1,\dots,d_2} \end{aligned}$$

where:

$$\begin{aligned} c_f^{l,s} &= \xi_l^\text{T} \mathcal{H}_{d_1|n_2} \left(\underline{f} \right) \phi_s = \xi_l^\text{T} (\mathcal{H}_{d_1}(f_1) \quad \dots \quad \mathcal{H}_{d_1}(f_{n_2})) \begin{pmatrix} \phi_s^1 \\ \vdots \\ \phi_s^{n_2} \end{pmatrix} = \sum_{i=1}^{n_2} \xi_l^\text{T} \mathcal{H}_{d_1}(f_i) \phi_s^i \\ &= \sum_{i=1}^{n_2} \left\langle f_i, \mathfrak{C}_{\phi_s^i}(\xi_l) \right\rangle_{\ell_2}. \end{aligned}$$

The last equality is due to $u^\text{T} \mathcal{H}_{d_1}(a)v = u^\text{T} \mathfrak{C}_v(a) = \langle a, \mathfrak{C}_v(u) \rangle_{\ell_2}$ for $u, v, a \in \mathbb{R}^d$.

Now, we expand \underline{f} :

$$\begin{aligned} \underline{f} &= \mathcal{H}_{d_1|n_2}^\dagger \left(\tilde{\Xi} \underline{c_f} \tilde{\Phi}^\text{T} \right) = \left(\mathcal{H}_{d_1}^\dagger \left(\tilde{\Xi} \underline{c_f} \tilde{\Phi}^{1,\text{T}} \right) \quad \dots \quad \mathcal{H}_{d_1}^\dagger \left(\tilde{\Xi} \underline{c_f} \tilde{\Phi}^{n_2,\text{T}} \right) \right) \\ &= \left(\mathcal{H}_{d_1}^\dagger \left(\sum_{s=1}^{d_2} \sum_{l=1}^d c_f^{l,s} \tilde{\xi}_l \tilde{\phi}_s^{1,\text{T}} \right) \quad \dots \quad \mathcal{H}_{d_1}^\dagger \left(\sum_{s=1}^{d_2} \sum_{l=1}^d c_f^{l,s} \tilde{\xi}_l \tilde{\phi}_s^{n_2,\text{T}} \right) \right) \\ &= \left(\frac{1}{d_1} \sum_{s=1}^{d_2} \sum_{l=1}^d c_f^{l,s} \mathfrak{C}_{\tilde{\phi}_s^1}(\tilde{\xi}_l) \quad \dots \quad \frac{1}{d_1} \sum_{s=1}^{d_2} \sum_{l=1}^d c_f^{l,s} \mathfrak{C}_{\tilde{\phi}_s^{n_2}}(\tilde{\xi}_l) \right) \\ &= \frac{1}{d_1} \sum_{s=1}^{d_2} \left(\mathfrak{C}_{\tilde{\phi}_s^1}(\tilde{\Xi} c_{f,s}) \quad \dots \quad \mathfrak{C}_{\tilde{\phi}_s^{n_2}}(\tilde{\Xi} c_{f,s}) \right). \end{aligned}$$

This concludes the proof. Note that the 3rd equality is due to:

$$\tilde{\Xi} \underline{c_f} \tilde{\Phi}^{i,\text{T}} = \begin{pmatrix} \tilde{\xi}_1 & \dots & \tilde{\xi}_m \end{pmatrix} \begin{pmatrix} c_f^{1,1} & \dots & c_f^{1,d_2} \\ \vdots & \ddots & \vdots \\ c_f^{d,1} & \dots & c_f^{d,d_2} \end{pmatrix} \begin{pmatrix} \tilde{\phi}_{1,\text{T}}^i \\ \vdots \\ \tilde{\phi}_{d_2,\text{T}}^i \end{pmatrix} = \sum_{s=1}^{d_2} \sum_{l=1}^d \tilde{\xi}_l c_f^{l,s} \tilde{\phi}_s^{i,\text{T}};$$

and, the 4th equality is due to:

$$\begin{aligned}
\mathcal{H}_{d_1}^\dagger \left(\sum_{s=1}^{d_2} \sum_{l=1}^d c_f^{l,s} \tilde{\xi}_l \tilde{\phi}_s^{n_2, \text{T}} \right) &= \frac{1}{\sqrt{d_1}} \sum_{s=1}^{d_2} \sum_{l=1}^d c_f^{l,s} \begin{pmatrix} \left\langle \frac{1}{\sqrt{d_1}} \mathcal{H}_{d_1}(\underline{e}_1), \tilde{\xi}_l \tilde{\phi}_s^{n_2, \text{T}} \right\rangle_{\text{F}} \\ \vdots \\ \left\langle \frac{1}{\sqrt{d_1}} \mathcal{H}_{d_1}(\underline{e}_{n_1}), \tilde{\xi}_l \tilde{\phi}_s^{n_2, \text{T}} \right\rangle_{\text{F}} \end{pmatrix} \\
&= \frac{1}{d_1} \sum_{s=1}^{d_2} \sum_{l=1}^d c_f^{l,s} \begin{pmatrix} \left\langle \underline{e}_1, \mathfrak{C}_{\tilde{\phi}_s^{n_2}}(\tilde{\xi}_l) \right\rangle_{\text{F}} \\ \vdots \\ \left\langle \underline{e}_{n_1}, \mathfrak{C}_{\tilde{\phi}_s^{n_2}}(\tilde{\xi}_l) \right\rangle_{\text{F}} \end{pmatrix} = \frac{1}{d_1} \sum_{s=1}^{d_2} \sum_{l=1}^d c_f^{l,s} \mathfrak{C}_{\tilde{\phi}_s^{n_2}}(\tilde{\xi}_l).
\end{aligned}$$

5.2 Proof of Proposition II.2

Given a wavelet expansion acting on a discrete function $f \in \ell_2(\mathbb{Z}^2)$:

$$f[k] = \sum_{m \in \mathbb{Z}^2} \langle f, \tilde{\varphi}_I(\cdot - m) \rangle_{\ell_2} \varphi_I(k - m) + \sum_{i=0}^I \sum_{l=0}^L \sum_{m \in \mathbb{Z}^2} \langle f, \tilde{\psi}_{il}(\cdot - m) \rangle_{\ell_2} \psi_{il}(k - m),$$

we compute its scaling and wavelet coefficients in the Fourier domain:

$$c_I[m] = \langle f, \tilde{\varphi}_I(\cdot - m) \rangle_{\ell_2} \xleftrightarrow{\mathcal{F}} \hat{C}_I(e^{j\omega}) = \sum_{m \in \mathbb{Z}^2} c_I[m] e^{-j\langle m, \omega \rangle_{\ell_2}} = \sum_{k \in \mathbb{Z}^2} f[k] \sum_{m \in \mathbb{Z}^2} \tilde{\varphi}_I(m - k) e^{-j\langle m, \omega \rangle_{\ell_2}}. \quad (15)$$

To compute a Fourier transform of $\tilde{\varphi}_I(m - k)$, we find a Fourier transform of its continuous version for $x \in \mathbb{R}^d$:

$$\tilde{\varphi}_I(x - k) \xleftrightarrow{\mathcal{F}} e^{-j\langle k, \omega \rangle_{\ell_2}} \int_{\mathbb{R}^2} \tilde{\varphi}_I(x) e^{-j\langle x, \omega \rangle_{\ell_2}} dx = \hat{\varphi}_I^*(\omega) e^{-j\langle k, \omega \rangle_{\ell_2}}.$$

By Poisson summation formulae, we have the following identity:

$$\sum_{m \in \mathbb{Z}^d} \tilde{\varphi}_I(m - k) e^{-j\langle m, \omega \rangle_{\ell_2}} = \sum_{m \in \mathbb{Z}^d} \hat{\varphi}_I^*(2\pi m + \omega) e^{-j\langle k, 2\pi m + \omega \rangle_{\ell_2}}$$

Note $e^{j2\pi m} = 1, \forall m \in \mathbb{Z}^d$; then, scaling coefficient $c_I[m] \xleftrightarrow{\mathcal{F}} \hat{C}_I(e^{j\omega})$ in (15) is

$$\begin{aligned}
\hat{C}_I(e^{j\omega}) &= \sum_{k \in \mathbb{Z}^2} f[k] \sum_{m \in \mathbb{Z}^2} \hat{\varphi}_I^*(2\pi m + \omega) e^{-j\langle k, 2\pi m + \omega \rangle_{\ell_2}} = \sum_{m \in \mathbb{Z}^2} \hat{\varphi}_I^*(2\pi m + \omega) \hat{F}(e^{j(2\pi m + \omega)}) \\
&= \hat{F}(e^{j\omega}) \left[\hat{\varphi}_I^*(\omega) + \sum_{m \in \mathbb{Z}^2 \setminus \{0\}} \hat{\varphi}_I^*(2\pi m + \omega) \right], \quad (16)
\end{aligned}$$

where $f[k] \xleftrightarrow{\mathcal{F}} \hat{F}(e^{j\omega})$. Similarly, wavelet coefficient $d_{il}[m] \xleftrightarrow{\mathcal{F}} \hat{D}_{il}(e^{j\omega})$ is:

$$\hat{D}_{il}(e^{j\omega}) = \hat{F}(e^{j\omega}) \left[\hat{\psi}_{il}^*(\omega) + \sum_{m \in \mathbb{Z}^2 \setminus \{0\}} \hat{\psi}_{il}^*(2\pi m + \omega) \right]. \quad (17)$$

Now, we compute wavelet expansion in the Fourier domain as:

$$\begin{aligned}
\widehat{F}(e^{j\omega}) &= \sum_{k \in \mathbb{Z}^2} f[k] e^{-j\langle k, \omega \rangle_{\ell_2}} \\
&= \sum_{m \in \mathbb{Z}^2} c_I[m] \sum_{k \in \mathbb{Z}^2} \varphi_I(k-m) e^{-j\langle k, \omega \rangle_{\ell_2}} + \sum_{i=0}^I \sum_{l=0}^L \sum_{m \in \mathbb{Z}^2} d_{il}[m] \sum_{k \in \mathbb{Z}^2} \psi_{il}(k-m) e^{-j\langle k, \omega \rangle_{\ell_2}} \\
&= \sum_{k \in \mathbb{Z}^2} \widehat{\varphi}_I(2\pi k + \omega) \sum_{m \in \mathbb{Z}^2} c_I[m] e^{-j\langle m, 2\pi k + \omega \rangle_{\ell_2}} + \sum_{i=0}^I \sum_{l=0}^L \sum_{k \in \mathbb{Z}^2} \widehat{\psi}_{il}(2\pi k + \omega) \sum_{m \in \mathbb{Z}^2} d_{il}[m] e^{-j\langle m, 2\pi k + \omega \rangle_{\ell_2}} \\
&= \widehat{C}_I(e^{j\omega}) \left(\widehat{\varphi}_I(\omega) + \sum_{k \in \mathbb{Z}^2 \setminus \{0\}} \widehat{\varphi}_I(2\pi k + \omega) \right) + \sum_{i=0}^I \sum_{l=0}^L \widehat{D}_{il}(e^{j\omega}) \left(\widehat{\psi}_{il}(\omega) + \sum_{k \in \mathbb{Z}^2 \setminus \{0\}} \widehat{\psi}_{il}(2\pi k + \omega) \right).
\end{aligned}$$

The 3rd equality is from Poisson summation formulae. We obtain the unity condition (Equation 6) from Equation 16 and Equation 17. To compensate the error in Equation 6, we modify a primal wavelet frame at scale 0 as:

$$\begin{aligned}
1 &= \widehat{\varphi}_I^*(\omega) \widehat{\varphi}_I(\omega) + \sum_{i=1}^I \sum_{l=0}^L \widehat{\psi}_{il}^*(\omega) \widehat{\psi}_{il}(\omega) + \sum_{l=0}^L \widehat{\psi}_{0l}^*(\omega) \widehat{\psi}_{0l}(\omega) + \widehat{e}(\omega) \\
&\Leftrightarrow \widehat{\psi}_0(\omega) := \widehat{\psi}_0(\omega) + \frac{\widehat{e}(\omega)}{\widehat{\psi}_0^*(\omega)}.
\end{aligned}$$

This is due to a unity property of L -th Riesz transform $\sum_{l=0}^L |\widehat{\mathcal{R}}^l(\omega)|^2 = 1$. Finally, the unity condition (Equation 6) becomes:

$$\widehat{\varphi}_I^*(\omega) \widehat{\varphi}_I(\omega) + \widehat{\psi}_0^*(\omega) \widehat{\psi}_0(\omega) + \sum_{i=1}^I \sum_{l=0}^L \widehat{\psi}_{il}^*(\omega) \widehat{\psi}_{il}(\omega) = 1.$$

Then, multiply 2 sides with $F(e^{j\omega})$ and take an inverse Fourier transform, we obtain a wavelet expansion:

$$\begin{aligned}
\widehat{F}(e^{j\omega}) &= \widehat{F}(e^{j\omega}) \widehat{\varphi}_I^*(\omega) \widehat{\varphi}_I(\omega) + \sum_{i=0}^I \sum_{l=0}^L \widehat{F}(e^{j\omega}) \widehat{\psi}_{il}^*(\omega) \widehat{\psi}_{il}(\omega) \\
\stackrel{\mathcal{F}}{\longleftrightarrow} f[k] &= \sum_{m \in \mathbb{Z}^2} \langle f, \tilde{\varphi}_I(\cdot - m) \rangle_{\ell_2} \varphi_I(k-m) + \sum_{i=0}^I \sum_{l=0}^L \sum_{m \in \mathbb{Z}^2} \langle f, \tilde{\psi}_{il}(\cdot - m) \rangle_{\ell_2} \psi_{il}(k-m)
\end{aligned}$$

with $\widehat{\psi}_{0n}(\omega) = \widehat{\psi}_{0n}(\omega)$.

5.3 Proof of Proposition II.3

Denote $\underline{\psi}_0 := \underline{\varphi}_I$ and $\tilde{\underline{\psi}}_0 := \tilde{\underline{\varphi}}_I$ whose matrix forms are $\underline{\Psi}_0 = \underline{\Phi}_I$, $\tilde{\underline{\Psi}}_0 = \tilde{\underline{\Phi}}_I \in \mathbb{R}^{n_1 n_2 \times n_2}$ as in (3), respectively. Note that $\underline{\mathfrak{C}}_\phi := \underline{\mathfrak{C}}_{\underline{\varphi}_I}^* \underline{\mathfrak{C}}_{\tilde{\underline{\varphi}}_I}$. A convolutional form of an expansion (??) is recast with

Hankel matrix as:

$$\begin{aligned}
\underline{f} &= \underline{\mathfrak{C}}_{\underline{\varphi}_I}^* \underline{\mathfrak{C}}_{\underline{\tilde{\varphi}}_I}(\underline{f}) + \sum_{p=1}^P \underline{\mathfrak{C}}_{\underline{\psi}_p}^* \underline{\mathfrak{C}}_{\underline{\tilde{\psi}}_p}(\underline{f}) = \sum_{p=0}^P \underline{\mathfrak{C}}_{\underline{\psi}_p}^* \underline{\mathfrak{C}}_{\underline{\tilde{\psi}}_p}(\underline{f}) = \sum_{p=0}^P \mathcal{H}_{n_1|n_2}(\underline{w}_{f,p}) \underline{\tilde{\Psi}}_p \\
&= \sum_{p=0}^P \left(\mathcal{H}_{n_1}(\underline{w}_{f,p,1}) \quad \dots \quad \mathcal{H}_{n_1}(\underline{w}_{f,p,n_2}) \right) \begin{pmatrix} \underline{\tilde{\Psi}}_p^1 \\ \vdots \\ \underline{\tilde{\Psi}}_p^{n_2} \end{pmatrix} = \sum_{p=0}^P \sum_{i=1}^{n_2} \mathcal{H}_{n_1}(\underline{w}_{f,p,i}) \underbrace{\left(\underline{\tilde{\psi}}_{p,1}^i \quad \dots \quad \underline{\tilde{\psi}}_{p,d_2}^i \right)}_{=\underline{\tilde{\Psi}}_p^i} \\
&= n_1 \sum_{p=0}^P \frac{1}{n_1} \sum_{i=1}^{n_2} \left(\underline{\mathfrak{C}}_{\underline{\tilde{\psi}}_{p,1}^i}(\underline{\tilde{\Xi}} \underline{w}_{f,p,i}) \quad \dots \quad \underline{\mathfrak{C}}_{\underline{\tilde{\psi}}_{p,n_2}^i}(\underline{\tilde{\Xi}} \underline{w}_{f,p,i}) \right) \\
&= n_1 \sum_{p=0}^P \mathcal{H}_{n_1|n_2}^\dagger(\underline{w}_{f,p} \underline{\tilde{\Psi}}_p^T) \\
&= n_1 \mathcal{H}_{n_1|n_2}^\dagger \left(\sum_{p=0}^P \underline{w}_{f,p} \underline{\tilde{\Psi}}_p^T \right) = n_1 \mathcal{H}_{n_1|n_2}^\dagger(\underline{w}_f \underline{\tilde{\Psi}}^T),
\end{aligned}$$

where wavelet coefficient is

$$\begin{aligned}
\underline{w}_{f,p} &:= \underline{\mathfrak{C}}_{\underline{\tilde{\psi}}_p}(\underline{f}) = \mathcal{H}_{n_1|n_2}(\underline{f}) \underline{\tilde{\Psi}}_p \\
&= (\underline{w}_{f,p,1} \quad \dots \quad \underline{w}_{f,p,n_2}) \in \mathbb{R}^{n_1 \times n_2}, \quad p = 0, \dots, P \\
&\Leftrightarrow \underline{w}_f := (\underline{w}_{f,0} \quad \dots \quad \underline{w}_{f,P}) = \mathcal{H}_{n_1|n_2}(\underline{f}) \underline{\tilde{\Psi}}.
\end{aligned}$$

Note that the last 3rd equality is due to proposition ?? with an identity local basis $\underline{\tilde{\Xi}} = \text{Id}$. The 2nd last equality is because $\mathcal{H}_{n_1|n_2}^\dagger$ is a linear operator: Given a matrix $\underline{g}_p = (\underline{g}_{p,1} \quad \dots \quad \underline{g}_{p,n_2})$ and constant $\{a_p\}_{p=0}^P \subset \mathbb{R}$, we have:

$$\sum_{p=0}^P a_p \mathcal{H}_{n_1|n_2}^\dagger(\underline{g}_p) = \left(\sum_{p=0}^P a_p \mathcal{H}_{n_1}^\dagger(\underline{g}_{p,1}) \quad \dots \quad \sum_{p=0}^P a_p \mathcal{H}_{n_1}^\dagger(\underline{g}_{p,n_2}) \right).$$

And, by a definition of an inverse Hankel matrix (2), we have:

$$\sum_{p=0}^P a_p \mathcal{H}_{n_1}^\dagger(\underline{g}_{p,i}) = \frac{1}{\sqrt{n_1}} \begin{pmatrix} \left\langle \underline{\tilde{e}}_1, \sum_{p=0}^P a_p \underline{g}_{p,i} \right\rangle_F \\ \vdots \\ \left\langle \underline{\tilde{e}}_{n_1}, \sum_{p=0}^P a_p \underline{g}_{p,i} \right\rangle_F \end{pmatrix} = \mathcal{H}_{n_1}^\dagger \left(\sum_{p=0}^P a_p \underline{g}_{p,i} \right).$$

Thus, we have a linear property of an extended Hankel matrix:

$$\begin{aligned} \sum_{p=0}^P a_p \mathcal{H}_{n_1|n_2}^\dagger(\underline{g}_p) &= \begin{pmatrix} \mathcal{H}_{n_1}^\dagger\left(\sum_{p=0}^P a_p \underline{g}_{p,1}\right) & \cdots & \mathcal{H}_{n_1}^\dagger\left(\sum_{p=0}^P a_p \underline{g}_{p,n_2}\right) \end{pmatrix} \\ &\stackrel{(1)}{=} \mathcal{H}_{n_1|n_2}^\dagger\left(\begin{pmatrix} \sum_{p=0}^P a_p \underline{g}_{p,1} & \cdots & \sum_{p=0}^P a_p \underline{g}_{p,n_2} \end{pmatrix}\right) \\ &= \mathcal{H}_{n_1|n_2}^\dagger\left(\sum_{p=0}^P a_p \underline{g}_p\right). \end{aligned}$$

In the end, we have the proposed non-subsampled Riesz-Quincunx wavelet is a framelet decomposition with an identity local basis:

$$\underline{f} = n_1 \mathcal{H}_{n_1|n_2}^\dagger \left(\mathcal{H}_{n_1|n_2}(\underline{f}) \underline{\tilde{\Psi}} \underline{\tilde{\Psi}}^T \right).$$

Since $\mathcal{H}_{n_1|n_2}^\dagger \circ \mathcal{H}_{n_1|n_2} = \text{Id}$, this implies the unity condition:

$$\frac{1}{n_1} \text{Id}_{n_1 n_2 \times n_1 n_2} = \underline{\tilde{\Psi}} \underline{\tilde{\Psi}}^T = \underline{\tilde{\Phi}}_I \underline{\tilde{\Phi}}_I^T + \sum_{p=1}^P \underline{\tilde{\Psi}}_p \underline{\tilde{\Psi}}_p^T.$$

5.4 Proof of Proposition II.4

The mapping $\mathcal{T}^{(i)}$ in Equation 13 is defined with $\underline{\mathfrak{C}}_{\underline{\theta}^{1(i)}}^{\text{iso}} : \mathbb{R}^{n_1 \times n_2 \times P} \rightarrow \mathbb{R}^{n_1 \times d_2 \times 2^i L}$ and $\underline{\mathfrak{C}}_{\underline{\theta}^{2(i)}}^{\text{iso}} : \mathbb{R}^{n_1 \times d_2 \times 2^i L} \rightarrow \mathbb{R}^{n_1 \times n_2 \times 2^i L}$. Given a local basis $\underline{\Xi}_{\text{aug}}^{(i)} = \begin{pmatrix} \text{Id} & \underline{\Xi}^{(i)} \end{pmatrix}$, an encoder is defined with a lowpass signal and skip-connecting signal as outputs:

$$\underline{c}_{\text{aug}}^{(i)} = \underline{\Xi}_{\text{aug}}^{(i)T} \mathcal{T}^{(i)}(\underline{s}^{(i-1)}) = \begin{pmatrix} \mathcal{T}^{(i)}(\underline{s}^{(i-1)}) \\ \underline{\Xi}^{(i),T} \mathcal{T}^{(i)}(\underline{s}^{(i-1)}) \end{pmatrix} = \begin{pmatrix} \underline{c}^{(i)} \\ \underline{s}^{(i)} \end{pmatrix}, \quad i = 1, \dots, I, \quad \underline{s}^{(0)} = \underline{f}, \quad \underline{s}^{(I)} = \mathcal{R}_p(\underline{s}^{(I)}) \quad (18)$$

$$\Leftrightarrow \begin{cases} \underline{s}^{(I)} = \mathcal{R}_p \circ \underline{\Xi}^{(I),T} \mathcal{T}^{(I)} \circ \dots \circ \underline{\Xi}^{(1),T} \mathcal{T}^{(1)}(\underline{f}), \\ \underline{c}^{(i)} = \mathcal{T}^{(i)} \circ \underline{\Xi}^{(i-1),T} \mathcal{T}^{(i-1)} \circ \dots \circ \underline{\Xi}^{(1),T} \mathcal{T}^{(1)}(\underline{f}), \quad i = 1, \dots, I, \end{cases} \quad (19)$$

where

$$\underline{c}^{(i)} \in \mathbb{R}^{2^{-(i-1)} n_1 \times 2^{-(i-1)} n_2 \times 2^{(i-1)} L}, \quad \underline{s}^{(i)} \in \mathbb{R}^{2^{-(i-1)-1} n_1 \times 2^{-(i-1)-1} n_2 \times 2^{(i-1)} L}.$$

Each subband of the skip connection $\underline{c}^{(i)} = \{c_0^{(i)}, \dots, c_{2^i L - 1}^{(i)}\}$, $c_l^{(i)} \in \mathbb{R}^{2^{-(i-1)} n_1 \times 2^{-(i-1)} n_2}$, is passed to the N -th order Riesz Quincunx wavelet expansion (Equation 11) at scale I' : $i = 0, \dots, I-1, l = 0, \dots, 2^i L - 1, k \in \mathbb{Z}^2$:

$$\begin{aligned} \underline{c}_l^{(i)} &= \underline{\mathfrak{C}}_\phi(\underline{c}_l^{(i)}) + \sum_{p=1}^P \underline{\mathfrak{C}}_{\underline{\psi}_p}^* \underline{\mathfrak{C}}_{\underline{\tilde{\psi}}_p}(\underline{c}_l^{(i)}) = \underline{\mathfrak{C}}_\phi(\underline{c}_l^{(i)}) + \underline{\mathfrak{C}}_{\underline{\tilde{\psi}}}^* \underline{\mathfrak{C}}_{\underline{\tilde{\psi}}}(\underline{c}_l^{(i)}) \\ &= \left(\mathcal{H}_{2^{-(i-1)} n_1 | 2^{-(i-1)} n_2}(\underline{c}_l^{(i)}) \underline{\Phi} \right) + \sum_{p=1}^P \left(\mathcal{H}_{2^{-(i-1)} n_1 | 2^{-(i-1)} n_2} \left(\mathcal{H}_{2^{-(i-1)} n_1 | 2^{-(i-1)} n_2}(\underline{c}_l^{(i)}) \underline{\tilde{\Psi}}_p \right) \underline{\tilde{\Psi}}_p \right) \end{aligned}$$

$$\Leftrightarrow \underline{c}_l^{(i)} = \mathfrak{W} \mathfrak{W}^{-1} \left\{ \underline{c}_l^{(i)} \right\},$$

which satisfies the condition of perfect reconstruction. Scaling and wavelet functions are defined in the previous section.

5.5 Proof of Proposition II.5

The proof is straight-forward by noting that a lowpass signal as an encoder's output in (19)

$$\underline{\underline{s}}^{(I)} = \mathcal{R}_p \circ \underline{\underline{\Xi}}^{(I),T} \mathcal{T}^{(I)} \circ \dots \circ \underline{\underline{\Xi}}^{(1),T} \mathfrak{T}^{(1)} \left(\underline{\underline{f}} \right) .$$

And, unknown parameters are: $\gamma_m := \{\gamma_c, \underline{W}^\mu, b^\mu\}$ and $\gamma_s := \{\gamma_c, \underline{W}^\sigma, b^\sigma\}$.

5.6 Proof of Proposition II.6

Now, we define an “inverse” mapping of (Equation 15) as a linear perceptron network and a reshape operation mapping:

$$\mathcal{F}^s(\cdot) = \text{uvec}(\underline{W}^s \cdot + b^s) : \mathbb{R}^d \rightarrow \mathbb{R}^{2^{-i-1}n_1 \times 2^{-i-1}n_2 \times 2^i L} \quad (20)$$

with $\underline{W}^s \in \mathbb{R}^{2^{-(I+1)}n_1 n_2 L \times d}$ and $b^s \in \mathbb{R}^{2^{-(I+1)}n_1 n_2 L}$.

Then, combining Equations 16 and 20, a reconstructed lowpass image generated from the latent variable y in Equation 16:

$$\begin{aligned} \underline{\underline{\hat{s}}}^{(I)} = \mathcal{F}^s(y) &= \text{uvec} \left(\underline{W}^s \mathcal{M}_{\gamma_m} \left(\underline{\underline{f}} \right) + b^s \right) + \text{uvec} \left(\left(\underline{W}^s \mathcal{S}_{\gamma_s}^{\frac{1}{2}} \left(\underline{\underline{f}} \right) \right) \odot \epsilon \right) \\ &\in \mathbb{R}^{2^{-i-1}n_1 \times 2^{-i-1}n_2 \times 2^i L}, \quad \epsilon \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}_d(\mathbf{0}_d, \text{Id}_d) . \end{aligned}$$

The last equality is because $\text{uvec}(\cdot)$ is a linear operator.

Given an augmented local basis $\underline{\underline{\tilde{\Xi}}}_{\text{aug}}^{(i)} = \left(\text{Id}_N \quad \mathcal{B} \circ \underline{\underline{\tilde{\Xi}}}^{(i)} \right)$, a decoder at scale $i \in \{I-1, \dots, 0\}$ is defined by an output encoder and a skip connection $\underline{\underline{\hat{c}}}_{\text{aug}}^{(i)} = \left(\underline{\underline{c}}^{(i)} \quad \underline{\underline{\hat{s}}}^{(i+1)} \right)$ as: $\underline{\underline{\hat{s}}}^{(I)} = \mathfrak{T}^s(y)$,

$$\begin{cases} \underline{\underline{\hat{c}}}_{\text{aug}}^{(i-1)} = \tilde{\mathcal{T}}^{(i)} \left(\underline{\underline{\hat{c}}}_{\text{aug}}^{(i)} \right) = \left(\underline{\underline{c}}^{(i-1)} \quad \underline{\underline{\hat{s}}}^{(i)} \right), \quad i = I, \dots, 1, \\ \underline{\underline{\hat{s}}}^{(0)} = \tilde{\mathcal{T}}^{(0)} \left(\underline{\underline{\hat{c}}}_{\text{aug}}^{(0)} \right), \quad i = 0, \end{cases} \quad (21)$$

$$\Leftrightarrow \underline{\underline{\hat{s}}}^{(0)} = \tilde{\mathcal{T}}^{(0)} \circ \tilde{\mathcal{T}}^{(1)} \circ \dots \circ \tilde{\mathcal{T}}^{(I-1)} \left(\underline{\underline{c}}^{(I-1)} \quad \underline{\underline{\hat{s}}}^{(I)} \right) := \tilde{\mathcal{T}}_I \left(\underline{\underline{c}}, \mathfrak{T}^s(y) \right) . \quad (22)$$

Then, we have Proposition II.6.

5.7 Proof of Proposition II.7

Since $\mathbf{p}(y | \underline{\underline{f}}) = \mathbf{k}_\alpha(y | \underline{\underline{f}})$, a likelihood of the latent variable in an encoder (Equation 26) is computed by a marginal likelihood:

$$\mathbf{p}(z) = \int_{\mathcal{F}} \mathbf{p}(z | \underline{\underline{f}}) \mathbf{p}(\underline{\underline{f}}) d\underline{\underline{f}} = \mathbb{E}_{\underline{\underline{f}} \sim \mathfrak{F}} \left[\mathbf{k}_\alpha(z | \underline{\underline{f}}) \right] .$$

Assume we have a bunch of realization of latent variable z as $\mathfrak{Z} = \{z_i\}_{i=1}^{n_z} \subset \mathbb{R}^d$; then, we have a maximum of an expected log-marginal-likelihood $\mathbf{p}(z)$ as:

$$\alpha^\dagger = \operatorname{argmax}_{\alpha} \mathbb{E}_{\underline{f} \sim \mathfrak{F}} \left[\log \mathbf{k}_{\alpha} \left(\mathfrak{Z} \mid \underline{f} \right) \right] = \operatorname{argmax}_{\alpha} \mathbb{E}_{\underline{f} \sim \mathfrak{F}} \left[\frac{1}{n_y} \sum_{i=1}^{n_y} \log \mathbf{k}_{\alpha} \left(z_i \mid \underline{f} \right) \right], \quad z_i \stackrel{\text{i.i.d.}}{\sim} \mathbb{K}_{\alpha} \left(z \mid \underline{f} \right)$$

which is

$$\begin{aligned} \left(\gamma_m^\dagger, \gamma_s^\dagger, \alpha^\dagger \right) &= \operatorname{argmax}_{\gamma_m, \gamma_s, \alpha} \mathbb{E}_{\underline{f} \sim \mathfrak{F}} \left[\frac{1}{n_z} \sum_{i=1}^{n_z} \log \mathbf{k}_{\alpha} (z_i \mid \underline{f}) \right], \quad z_i \stackrel{\text{i.i.d.}}{\sim} \mathbb{G}_{\gamma_m, \gamma_s} (z \mid \underline{f}) = \mathbf{g}_{\gamma_m, \gamma_s} (z \mid \underline{f}) \, dz \\ &\approx \operatorname{argmax}_{\gamma_m, \gamma_s, \alpha} \mathbb{E}_{\underline{f} \sim \mathfrak{F}} \left[\mathbb{E}_{z \sim \mathbb{G}_{\gamma_m, \gamma_s} (\cdot \mid \underline{f})} \left[\log \mathbf{k}_{\alpha} (z \mid \underline{f}) \right] \right] \\ &= \operatorname{argmax}_{\gamma_m, \gamma_s, \alpha} \left\{ \mathbb{E}_{\underline{f} \sim \mathfrak{F}} \left[\mathbb{E}_{z \sim \mathbb{G}_{\gamma_m, \gamma_s} (\cdot \mid \underline{f})} \left[\log \mathbf{k}_{\alpha} (z \mid \underline{f}) \right] - \mathbb{E}_{z \sim \mathbb{G}_{\gamma_m, \gamma_s} (\cdot \mid \underline{f})} \left[\log \mathbf{g}_{\gamma_m, \gamma_s} (z \mid \underline{f}) \right] \right] \right\} \\ &= \operatorname{argmax}_{\gamma_m, \gamma_s, \alpha} \left\{ -\mathbb{E}_{\underline{f} \sim \mathfrak{F}} \left[\mathbb{E}_{z \sim \mathbb{G}_{\gamma_m, \gamma_s} (\cdot \mid \underline{f})} \left[\log \frac{\mathbf{g}_{\gamma_m, \gamma_s} (z \mid \underline{f})}{\mathbf{k}_{\alpha} (z \mid \underline{f})} \right] \right] \right\} \\ &= \operatorname{argmin}_{\gamma_m, \gamma_s, \alpha} \mathbb{E}_{\underline{f} \sim \mathfrak{F}} \left[\operatorname{KL} \left(\mathbb{G}_{\gamma_m, \gamma_s} (z \mid \underline{f}) \parallel \mathbb{K}_{\alpha} (z \mid \underline{f}) \right) \right]. \end{aligned} \quad (23)$$

In the end, we need to find $(\gamma_m, \gamma_s, \theta)$ that minimize the KL-distance between the true distribution and its approximated version in Equation 27 over all data \mathfrak{F} as:

$$\min_{\gamma_m, \gamma_s, \alpha} \mathbb{E}_{\underline{f} \sim \mathfrak{F}} \left[\operatorname{KL} \left(\mathbb{G}_{\gamma_m, \gamma_s} (z \mid \underline{f}) \parallel \mathbb{K}_{\alpha} (z \mid \underline{f}) \right) \right]. \quad (24)$$

From Bayes' rule (Equation 27), we derive an evidence lower bound of the KL-divergence in (24):

$$\begin{aligned} &\operatorname{KL} \left(\mathbb{G}_{\gamma_m, \gamma_s} (z \mid \underline{f}) \parallel \mathbb{K}_{\alpha} (z \mid \underline{f}) \right) \\ &= \mathbb{E}_{z \sim \mathbb{G}_{\gamma_m, \gamma_s} (\cdot \mid \underline{f})} \left[\log \frac{\mathbf{g}_{\gamma_m, \gamma_s} (z \mid \underline{f})}{\mathbf{k}_{\alpha} (z \mid \underline{f})} \right] \stackrel{(27)}{=} \mathbb{E}_{z \sim \mathbb{G}_{\gamma_m, \gamma_s} (\cdot \mid \underline{f})} \left[\log \frac{\mathbf{p}(\underline{f}) \mathbf{g}_{\gamma_m, \gamma_s} (z \mid \underline{f})}{\mathbf{h}_{\alpha} (\underline{f} \mid z) \mathbf{p}(z)} \right] \\ &= -\mathbb{E}_{z \sim \mathbb{G}_{\gamma_m, \gamma_s} (\cdot \mid \underline{f})} \left[\log \mathbf{h}_{\alpha} (\underline{f} \mid z) \right] + \mathbb{E}_{z \sim \mathbb{G}_{\gamma_m, \gamma_s} (\cdot \mid \underline{f})} \left[\log \frac{\mathbf{g}_{\gamma_m, \gamma_s} (z \mid \underline{f})}{\mathbf{p}(z)} \right] + \log \mathbf{p} (\underline{f}) \end{aligned}$$

which is equivalent to:

$$\begin{aligned} \mathcal{L} \left(\mathbb{G}_{\gamma_m, \gamma_s}, \alpha; \underline{f} \right) &:= \mathbb{E}_{z \sim \mathbb{G}_{\gamma_m, \gamma_s} (\cdot \mid \underline{f})} \left[\log \mathbf{h}_{\alpha} (\underline{f} \mid z) \right] - \operatorname{KL} \left[\mathbb{G}_{\gamma_m, \gamma_s} (z \mid \underline{f}) \parallel \mathbb{P}(z) \right] \\ &= \log \mathbf{p} (\underline{f}) - \operatorname{KL} \left(\mathbb{G}_{\gamma_m, \gamma_s} (z \mid \underline{f}) \parallel \mathbb{K}_{\alpha} (z \mid \underline{f}) \right) \\ &\leq \log \mathbf{p} (\underline{f}). \end{aligned}$$

The last inequality is due to $\text{KL}(\cdot \parallel \cdot) \geq 0$. Note that, via a reparamization trick from Equation 20

$$\begin{aligned} z &\stackrel{\text{i.i.d.}}{\sim} \mathcal{N}_d \left(\mathcal{M}_{\gamma_m}(\underline{f}), \text{diag} \left\{ \mathcal{S}_{\gamma_s}(\underline{f}) \right\} \right) = \mathbb{G}_{\gamma_m, \gamma_s} \left(z \mid \underline{f} \right), \\ \Leftrightarrow z &= \mathcal{M}_{\gamma_m}(\underline{f}) + \mathcal{S}_{\gamma_s}^{\frac{1}{2}}(\underline{f}) \odot \epsilon, \quad \epsilon \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \text{Id}) \end{aligned}$$

an evidence bound is recast as:

$$\mathcal{L} \left(\mathbb{G}_{\gamma_m, \gamma_s}, \alpha; \underline{f} \right) = \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \text{Id})} \left[\log \mathbf{h}_\alpha \left(\underline{f} \mid z = \mathcal{M}_{\gamma_m}(\underline{f}) + \mathcal{S}_{\gamma_s}^{\frac{1}{2}}(\underline{f}) \odot \epsilon \right) \right] - \text{KL} \left[\mathbb{G}_{\gamma_m, \gamma_s} \left(z \mid \underline{f} \right) \parallel \mathbb{P}(z) \right].$$

Then, a minimization (24) is equivalent to maximize an expected evidence bound over all observation \mathfrak{F} :

$$\begin{aligned} \left(\gamma_m^\dagger, \gamma_s^\dagger, \alpha^\dagger \right) &= \underset{\gamma_m, \gamma_s, \alpha}{\text{argmax}} \mathbb{E}_{\underline{f} \sim \mathfrak{F}} \left[\log \mathbf{p}(\underline{f}) - \text{KL} \left(\mathbb{G}_{\gamma_m, \gamma_s} \left(z \mid \underline{f} \right) \parallel \mathbb{K}_\alpha \left(z \mid \underline{f} \right) \right) \right] \\ &= \underset{\gamma_m, \gamma_s, \alpha}{\text{argmax}} \mathbb{E}_{\underline{f} \sim \mathfrak{F}} \left[\mathcal{L} \left(\mathbb{G}_{\gamma_m, \gamma_s}, \alpha; \underline{f} \right) \right] \\ &= \underset{\gamma_m, \gamma_s, \alpha}{\text{argmax}} \frac{1}{T} \sum_{i=1}^T \mathcal{L} \left(\mathbb{G}_{\gamma_m, \gamma_s}, \alpha; \underline{f}_i \right). \end{aligned} \quad (25)$$

This is because $\mathbf{p}(\underline{f})$ is independent to $(\gamma_m, \gamma_s, \alpha)$.

For choosing $\mathbb{P}(z) = \mathcal{N}_d(\mathbf{0}_d, \text{Id}_d)$ and $\mathbb{G}_{\gamma_m, \gamma_s} \left(z \mid \underline{f} \right) = \mathcal{N}_d \left(\mathcal{M}_{\gamma_m}(\underline{f}), \text{diag} \left\{ \mathcal{S}_{\gamma_s}(\underline{f}) \right\} \right)$, the KL-divergence term is defined above. From Equation 25, we have a likelihood:

$$\mathbf{h}_\alpha \left(\underline{f} \mid z \right) = \mathcal{N}_D \left(\underline{f}; \mathcal{D}_\alpha(\underline{c}, z), \sigma^2 \text{Id} \right) \propto \exp \left(-\frac{1}{2\sigma^2} \left\| \underline{f} - \mathcal{D}_\alpha(\underline{c}, z) \right\|_{\ell_2}^2 \right). \quad (26)$$

Then, a minimization (25) becomes the following minimization problem:

$$\left(\gamma_m^\dagger, \gamma_s^\dagger, \alpha^\dagger \right) = \underset{\gamma_m, \gamma_s, \alpha}{\text{argmax}} \left\{ \mathcal{L}(\gamma_m, \gamma_s, \alpha) = \sum_{i=1}^T \mathcal{L} \left(\mathbb{G}_{\gamma_m, \gamma_s}, \alpha; \underline{f}_i \right) \right\}. \quad (27)$$

5.8 Proof of Proposition II.9

We firstly define a vector-valued function for K -clusters:

$$u : \mathbb{R}^P \rightarrow \mathbb{R}^K; \quad u(f_{i,l}) = \begin{pmatrix} u_1(f_{i,l}) \\ \vdots \\ u_K(f_{i,l}) \end{pmatrix}, \quad u_k(f_{i,l}) \in \mathbb{R},$$

whose definition is defined in a matrix form by our proposed RQUnet-VAE as:

$$\begin{aligned}
u(\underline{\underline{f_i}}) &:= \left[u(f_{i,l}) \right]_{l \in \Omega} \\
&= \left\{ u_1(\underline{\underline{f_i}}), \dots, u_K(\underline{\underline{f_i}}) \right\} \in \mathbb{R}^{n_1 \times n_2 \times K}, \quad u_k(\underline{\underline{f_i}}) \in \mathbb{R}^{n_1 \times n_2} \\
&= \underset{\equiv}{\underset{\equiv}{\underset{\equiv}{\mathfrak{C}}}}^{\text{iso}} \circ \mathcal{D}_\alpha \left(\mathcal{C}_{\gamma_c}(\underline{\underline{f_i}}), \mathcal{M}_{\gamma_m}(\underline{\underline{f_i}}) + \mathcal{S}_{\gamma_s}^{\frac{1}{2}}(\underline{\underline{f_i}}) \odot \epsilon \right), \tag{28}
\end{aligned}$$

where an isotropic matrix-family convolution is for the K -clusters:

$$\underset{\equiv}{\underset{\equiv}{\underset{\equiv}{\mathfrak{C}}}}^{\text{iso}} : \mathbb{R}^{|\Omega| \times P} \rightarrow \mathbb{R}^{|\Omega| \times K}.$$

We define distributions for the set $(\mathfrak{F}, \mathfrak{F}^{\text{gt}})$ as $\mathbb{Q}(f_{i,l}) = q(f_{i,l}) \, \text{d}f_{i,l}$, $\mathbb{P}(f_{i,l}) = p(f_{i,l}) \, \text{d}f_{i,l}$ and their densities are:

$$\begin{aligned}
q(f_{i,l}) &= \left[q_k(f_{i,l}) \right]_{k=1}^K \in \mathbb{R}^K, \quad q_k(f_{i,l}) = \text{softmax} \left(u(f_{i,l}) \right)_k = \frac{\exp \left(u_k(f_{i,l}) \right)}{\sum_{h=1}^K \exp \left(u_h(f_{i,l}) \right)}, \\
p(f_{i,l}) &= \left[p_k(f_{i,l}) \right]_{k=1}^K \in \mathbb{R}^K, \quad p_k(f_{i,l}) = f_{i,l,k}^{\text{gt}} \in \{0, 1\}, \quad \sum_{k=1}^K p_k(f_{i,l}) = 1;
\end{aligned}$$

then, cross-entropy loss between distributions \mathbb{P} and \mathbb{Q} is:

$$\begin{aligned}
\mathcal{H}(\mathbb{P}, \mathbb{Q}) &= -\mathbb{E}_{\mathbb{P}} [\log \mathbb{Q}] := -\sum_{k=1}^K \mathbb{E}_{f \sim \mathfrak{F}} [p_k(f) \log q_k(f)] \\
&= \lim_{n \rightarrow \infty} -\frac{1}{n} \sum_{i=1}^n \mathcal{H} \left(\underline{\underline{f_i}}^{\text{gt}}, u(\underline{\underline{f_i}}) \right), \\
\mathcal{H} \left(\underline{\underline{f_i}}^{\text{gt}}, u(\underline{\underline{f_i}}) \right) &= \sum_{l \in \Omega} \sum_{k=1}^K f_{i,l,k}^{\text{gt}} \log \text{softmax} \left(u(f_{i,l}) \right)_k.
\end{aligned}$$

The last equality is due to Monte Carlo approximation method.

5.9 Proof of proposition 4.1

We start with a 2D convolution (8)

$$\tilde{\underline{\underline{f}}} = \underset{\phi}{\mathfrak{C}}(\underline{\underline{f}}) = \begin{pmatrix} \tilde{f}_1 & \dots & \tilde{f}_{d_2} \end{pmatrix} \in \mathbb{R}^{n_1 \times d_2}, \quad \tilde{f}_i \in \mathbb{R}^{n_1}.$$

Then, we rewrite an element of the above 2D convolution as a 1D convolution for a 2D image $\underline{f} = (f_1 \ \dots \ f_{n_2}) \in \mathbb{R}^{|\Omega|}$, $f_i \in \mathbb{R}^{n_1}$: $r_1 = 1, \dots, n_1$, $r_2 = 1, \dots, d_2$,

$$\begin{aligned} \tilde{f}[r_1, r_2] &= \underline{\mathfrak{C}}_{\underline{\phi}}(\underline{f})[r_1, r_2] = \sum_{k_1=1}^{d_1} \sum_{k_2=1}^{d_2} f[k_1, k_2] \check{\phi}[r_1 - k_1, r_2 - k_2] \\ &= \sum_{k_2=1}^{d_2} \left(\sum_{k_1=1}^{d_1} f_{k_2}[k_1] \check{\phi}_{r_2}^{k_2}[r_1 - k_1] \right) \\ &= \sum_{k_2=1}^{d_2} \mathfrak{C}_{\phi_{r_2}^{k_2}}(f_{k_2})[r_1] := \tilde{f}_{r_2}[r_1], \end{aligned}$$

where column vector notations are

$$f_{k_2}[k_1] := f[k_1, k_2], \ \phi_{r_2}^{k_2}[r_1 - k_1] := \phi_{k_2}[r_1 - k_1, r_2] := \phi[r_1 - k_1, r_2 - k_2].$$

In conclusion, we have a useful relation between 2D and 1D convolution operations: $k_1 = 1, \dots, n_1$, $k_2 = 1, \dots, d_2$,

$$\tilde{f}_{k_2}[k_1] = \underline{\mathfrak{C}}_{\underline{\phi}}(\underline{f})[k_1, k_2] = \sum_{i=1}^{d_2} \mathfrak{C}_{\phi_{k_2}^i}(f_i)[k_1] = \sum_{i=1}^{d_2} \left(\mathcal{H}_{d_1}(f_i) \phi_{k_2}^i \right) [k_1], \quad (29)$$

where $f_{k_2}[k_1] := f[k_1, k_2]$.

Then, a multi-input-numliti-output convolution is:

$$\begin{aligned} \underline{\tilde{f}} &= \left(\tilde{f}_1, \ \dots, \ \tilde{f}_{d_2} \right) \in \mathbb{R}^{n_1 \times d_2}, \ \tilde{f}_i \in \mathbb{R}^{n_1}, \\ &= \sum_{i=1}^{d_2} \left(\mathfrak{C}_{\phi_1^i}(f_i), \ \dots, \ \mathfrak{C}_{\phi_{d_2}^i}(f_i) \right) \stackrel{(7)}{=} \sum_{i=1}^{d_2} \mathcal{H}_{n_1}(f_i) \underbrace{\left(\phi_1^i, \ \dots, \ \phi_{d_2}^i \right)}_{=\underline{\Phi}^i} \\ &= \mathcal{H}_{n_1|d_2}(\underline{f}) \underline{\Phi}. \end{aligned}$$

This concludes a relation (8).

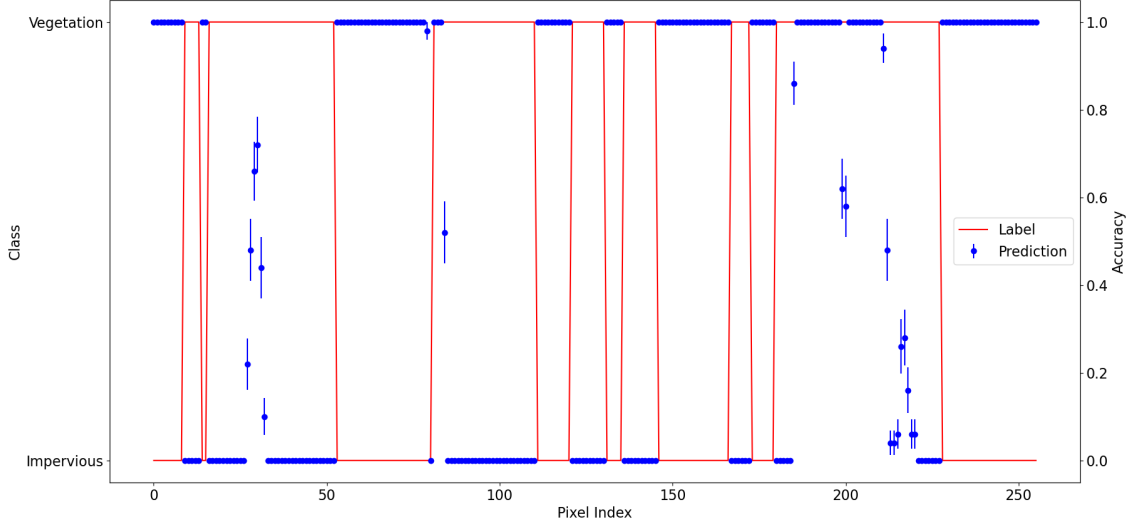


Figure 6: The probability that individual pixels are assigned to the impervious class along single horizontal line of pixels in an image during 50 RQUNet-VAE predictions compared to ground truth. Blue dots correspond to predicted accuracy (right y-axis) whereas the red line corresponds to the ground truth land cover class (left y-axis). This illustrates the results of variational terms of RQUNet-VAE.

6 Additional Experiments

While the main paper describes precision, recall, and F1-scores for entire images, this experiment looks at the accuracy of individual pixels an image to understand where our RQUNet-VAE incurs errors and where the variational approach is able to identify pixels having low confidence. For this purpose, we look at one horizontal line of pixels taken from the vertical center of the first noisy image shown in Figure 2. For each pixel, from the left to the right of the image, Figure 5.9 shows the ground truth label and the prediction accuracy. We see that for this line of pixels, the ground truth changes frequently between the impervious and the vegetation (non-impervious) class. We observe that there are many pixels where RQUNet-VAE provides a prediction accuracy of 1, implying that the class was predicted correctly in each of 50 runs of RQUNet-VAE (which is not deterministic due to the variational auto-encoder module). For example, on the very left of Figure 5.9 we observe that numerous vegetation pixels were classified correctly in all cases. Then however, the ground truth has a few pixels classified as impervious, which our RQUNet-VAE is not able to capture. For these pixels, we have an accuracy of 0, meaning that even the repeated iterations of our variational approach are not able to correctly classify these pixels. Across this horizontal line of pixels we also have numerous pixel where the variational approach provides a non-binary classification. For these cases, we provide the accuracy (as the fraction of correct classifications among the 50 runs) as well as the standard deviation of this accuracy denoted by the whiskers around the point estimate in the figure. We observe that in these cases, the variational approach allows our RQUNet-VAE to possibly make better decisions by allowing to base a decision on multiple runs and take the consensus of all runs.

While these results appear weak due to many miss-classifications, we again reiterate that the underlying image has been heavily obfuscated with noise (compare the original images in Figure 1 with the noised images in Figure 2) and that the experiments in the main paper show that the classic U-Net yields worse results, thus showing that our proposed RQUNet-VAE augments the traditional U-Net architecture by making it more robust to noise.