

# Learning Asymmetric and Local Features in Multi-Dimensional Data through Wavelets with Recursive Partitioning

Meng Li and Li Ma

**Abstract**—Effective learning of asymmetric and local features in images and other data observed on multi-dimensional grids is a challenging objective critical for a wide range of image processing applications involving biomedical and natural images. It requires methods that are sensitive to local details while fast enough to handle massive numbers of images of ever increasing sizes. We introduce a probabilistic model-based framework that achieves these objectives by incorporating adaptivity into discrete wavelet transforms (DWT) through Bayesian hierarchical modeling, thereby allowing wavelet bases to adapt to the geometric structure of the data while maintaining the high computational scalability of wavelet methods—linear in the sample size (e.g., the resolution of an image). We derive a recursive representation of the Bayesian posterior model which leads to an exact message passing algorithm to complete learning and inference. While our framework is applicable to a range of problems including multi-dimensional signal processing, compression, and structural learning, we illustrate its work and evaluate its performance in the context of image reconstruction using real images from the ImageNet database, two widely used benchmark datasets, and a dataset from retinal optical coherence tomography and compare its performance to state-of-the-art methods based on basis transforms and deep learning.



## 1 INTRODUCTION

EFFECTIVE learning of asymmetric and local features in images and other data observed on multi-dimensional grids plays a critical role in a wide range of applications. One such application is optical coherence tomography (OCT). OCT is a non-invasive imaging modality widely used in ophthalmology to visualize cross-sections of tissue layers. These tissue layers—such as the inner nuclear layer and outer nuclear layer—are often mostly homogeneous horizontally while involving large vertical contrasts. These contrasts across layers are key for ophthalmologists to make a diagnosis based on the (algorithmically reconstructed) image. Furthermore, local structures in such images can indicate ocular diseases, and their proper quantitative assessment is an important reference for monitoring the progression of the disease in clinical practice [1], [2], [3], [4], [5], [6]. Many other applications of 2D and 3D image analyses in biomedicine and beyond also involve asymmetric and local features to various extents. The effective analysis of such multi-dimensional observations can be greatly enhanced by incorporating adaptivity into the algorithm or method to take into account such features.

A further challenge in modern applications involving multi-dimensional observations is the ever increasing size of the datasets. For example, both the number of images analyzed as well as the resolution—i.e., the total number of pixels—of each image have been expanding rapidly. Many traditional methods and models become computational impractical for modern data as they scale polynomially with the resolution. Effective methods for analyzing such data

must scale well with both the resolution of each image and the number of images.

The primary aim of this work is to present a general-purpose generative probabilistic model for data on multi-dimensional grids that can be used to address these challenges in inference and learning—being able to effectively adapt to the asymmetric and local nature of interesting features, while achieving a highly efficient linear computational budget.

Our starting point is a well-known strategy for representing functions—a multi-resolution representation through the discrete wavelet transform (DWT). Wavelet analysis is hardly a new topic [7], [8], [9] and it has played an important role in the context of signal processing and image analysis. Its linear computational scalability is well-suited for analyzing massive data. However, traditional statistical wavelet analyses have mostly been focusing on effective modeling and inference on the wavelet coefficients *given a fixed* wavelet transform of the original data [10], [11], [12], [13], [14], [15]. A predetermined fixed wavelet transform, however, cannot adapt to the structure of the data and consequently suffers in its ability to effectively maintain the local structures in the original observation. Also, classical wavelet transforms on multi-dimensional grids are generally symmetric with respect to the dimensions, rendering them ineffective for preserving asymmetric features. No downstream statistical analyses can recover what has already been lost at the upstream wavelet transform stage.

In this work, we show that it is possible to incorporate the desired adaptivity into the wavelet transform stage while maintaining the computational scalability of the statistical analysis through a very simple hierarchical modeling strategy—starting the model “one level up”, that is, by incorporating the wavelet transform itself as an unknown

- *M. Li is with the Department of Statistics, Rice University, Houston, TX, 77025. E-mail: meng@rice.edu.*
- *L. Ma is with the Department of Statistical Science, Duke University, Durham, NC, 27708. E-mail: li.ma@duke.edu.*

quantity of interest into the probabilistic model, and learn it based on the data. Specifically, we consider latent (1D) wavelet transforms that can “twist and turn” (or “warp”) over the multi-dimensional grid, or the *index space*, and adopt a Bayesian prior on the path of its twisting and turning. In other words, we place a prior on the local directionality of the 1D transform to allow the “warping” to adapt to the geometric structure of the underlying function, e.g., the true image, through the Bayesian machinery.

In designing an appropriate prior for the local directionality, we note that “warping” a 1D wavelet transform through the grid points is equivalent to fixing the 1D wavelet transform while shuffling grid points in the multivariate index space of the observation—i.e., through applying a given 1D wavelet transform to a permuted version of the observation. This connection implies that probabilistic models on “warping” can be induced from distributions on the space of permutations of the index points or locations. Moreover, we draw a further connection between permutations and recursive dyadic partitioning on the index space to construct a prior on permutations induced by random recursive partitioning over the index space. This prior takes advantage of the fact that multi-dimensional images tend to be piecewise smooth to strike a balance between flexibility and computational tractability, allowing us to complete exact Bayesian inference through a recursive message passing algorithm with a computational budget linear in the resolution and sample size.

Due to the connection to recursive partitioning, we shall refer to our approach as WARP, or *WAvelets with Recursive Partitioning*. Through extensive numerical studies involving a large number of natural images from the ImageNet database, two additional benchmark data sets, and an OCT data set, we show that WARP often outperforms the existing state-of-the-art approaches by a substantial margin while maintaining the computational efficiency of classical wavelet analyses with fixed wavelet transforms. While we focus on 2D and 3D images in our motivation and numerical examples, our framework is readily applicable to observations of more than three dimensions without modification.

The rest of the paper is organized as follows. Section 2 introduces the WARP framework. In Section 2.1 we review the key components of Bayesian wavelet regression models, introduce permutation of the index space as a way to incorporate adaptivity into wavelet analysis, and construct a class of priors on permutations induced by recursive dyadic partitioning on the index space. We derive the corresponding posterior model and provide computational recipes for exact Bayesian inference under the WARP model with Haar wavelets in Section 2.2. In Section 3, we carry out an extensive numerical study and compare our method to existing state-of-the-art wavelet and non-wavelet methods including a deep learning method on a variety of real images. In Section 5 we carry out a case study by applying WARP to analyze an OCT data set, and compares its performance to a number of state-of-the-art approaches. Section 6 concludes with some brief remarks. The C++ source code along with a Matlab toolbox and R package to implement the proposed method is available online at <https://github.com/MaStatLab/WARP>.

## 2 METHOD

### 2.1 A Bayesian hierarchical wavelet regression model with recursive dyadic partitions

#### 2.1.1 Background and overview

We use  $\Omega$  to denote a space of indices or locations (e.g., pixels in images) where we obtain numerical measurements (e.g., intensities of pixels). Throughout this work, we assume  $\Omega$  to be an  $m$ -dimensional rectangular tube consisting of  $n_i = 2^{J_i}$  grid points in the  $i$ th dimension for  $i = 1, 2, \dots, m$ , that is, the function values are observed on a multi-dimensional equidistant grid. To simplify notation, we shall use  $[a, b]$  to represent the set  $\{a, a + 1, \dots, b\}$  for two integers  $a$  and  $b$  with  $a \leq b$ . Then the index space  $\Omega$  is of the form

$$\Omega = [0, 2^{J_1} - 1] \times [0, 2^{J_2} - 1] \times \dots \times [0, 2^{J_m} - 1].$$

The locations in  $\Omega$  can be placed into a vector of length  $n = 2^J$ . For example, we can map the location  $s = (s_1, s_2, \dots, s_m) \in \Omega$  to the  $t$ th element in the vector, where  $t = s_1 + \sum_{l=2}^m (\prod_{i=1}^{l-1} n_i) s_l$ . Correspondingly, any function  $f : \Omega \rightarrow \mathbb{R}$  can be represented as a vector  $\mathbf{f}$  of length  $n = 2^J$  whose  $t$ th element is  $f(s)$ .

Now, we consider a regression model

$$\mathbf{y} = \mathbf{f} + \boldsymbol{\epsilon} \quad \text{with } \boldsymbol{\epsilon} \sim N(\mathbf{0}, \Sigma_{\boldsymbol{\epsilon}}), \quad (1)$$

where  $\mathbf{y} = (y_0, y_1, \dots, y_{2^J-1})'$  are the observations on  $\Omega$ ,  $\mathbf{f} = (f_0, f_1, \dots, f_{2^J-1})'$  the underlying unknown function mean (or the signal), and  $\boldsymbol{\epsilon} = (\epsilon_0, \epsilon_1, \dots, \epsilon_{2^J-1})'$  the noise. For ease of illustration, we assume homogeneous white noise, i.e.,  $\Sigma_{\boldsymbol{\epsilon}} = \sigma^2 I_n$ , though our model and inference algorithms do not rely on this assumption at all and can be readily apply to models with heterogeneous variance; see Section 6 for further discussion.

One can apply a 1D discrete wavelet transform (DWT) to the observation vector  $\mathbf{y}$  through multiplying the corresponding orthonormal matrix  $W$  to both sides of Eq. (1), obtaining  $\mathbf{w} = \mathbf{z} + \mathbf{u}$  where  $\mathbf{w} = W\mathbf{y}$  is the vector of empirical wavelet coefficients,  $\mathbf{z} = W\mathbf{f}$  the mean vector for wavelet coefficients and  $\mathbf{u} = W\boldsymbol{\epsilon}$  the noise vector in the wavelet domain. This model can be rewritten in a location-scale form:  $w_{j,k} = z_{j,k} + u_{j,k}$  for  $j = 0, 1, \dots, J - 1$  and  $k = 0, 1, \dots, 2^j - 1$ , where  $w_{j,k}$ ,  $z_{j,k}$ ,  $u_{j,k}$  are the  $k$ th wavelet coefficient, signal, and noise at the  $j$ th scale in the wavelet (i.e., location-scale) domain, respectively.

It will generally be unreasonable to treat multi-dimensional observations simply as a vector with an arbitrary ordering of the locations; see [16], [17], [18]. Such a vectorization ignores the structure of the underlying function, and thus will result in less effective “energy concentration”, i.e., producing a wavelet decomposition of  $\mathbf{f}$  that is not very sparse—with many non-zero  $z_{j,k}$ ’s of small to moderate sizes, reducing the signal-to-noise ratio at those  $(j, k)$  combinations.

For each specific data set at hand, however, there typically exists some orderings of the locations that effectively reorganize the data so that the corresponding vectorization of the data provides an efficient representation of the underlying function; see Figure 1 for an illustration. Adopting a Bayesian modeling perspective, one can think of the underlying “good” vectorizations as latent structures of interest.

Also, one can view the wavelet regression model under each given index permutation as a competing generative model for the observed data given the latent structure. This perspective inspires us to incorporate a prior on the permutations, thereby allowing us to compute a posterior on the space of competing wavelet regression models, and then based on the goal of the analysis proceed with the common devices for Bayesian inference. Two particular useful tools are (i) Bayesian model selection [19]—learning a good permutation for representing the image based on its posterior probability; and (ii) Bayesian model averaging—estimating the underlying function based on averaging over the different permutations using their posterior probabilities [20].

This Bayesian approach does incur a practical challenge commonly arising in high-dimensional problems—the space of all permutations is so massive that brute-force enumeration over the space is computationally prohibitive. In the current context, effective exploration of the model (i.e., permutation) space becomes possible, however, once we realize that the vast majority of the permutations will lead to wavelet regression models that ignore the spatial smoothness of the underlying function—i.e., close locations in  $\Omega$  often correspond to similar values in  $\mathbf{f}$ . In particular, we can focus attention on a subclass of permutations that to various extents preserve spatial smoothness, and design a model space prior supported on this manageable subclass. To this end, we appeal to a relationship between recursive dyadic partitioning (RDP) and permutations, and shall consider the collection of permutations induced by RDPs on  $\Omega$ .

Next we introduce some basic notions regarding RDPs on  $\Omega$ , which are then used to construct a prior on permutations. In reading the next two subsections, the reader may refer to Figure 1 for an illustration of the key notions and notations.

### 2.1.2 Recursive dyadic partitioning on the location space

A *partition* of  $\Omega$  is a collection of nonempty sets  $\{A_1, A_2, \dots, A_H\}$  such that  $\Omega = \cup_{h=1}^H A_h$  and  $A_{h_1} \cap A_{h_2} = \emptyset$  for any  $h_1 \neq h_2$ . Now let  $\mathcal{T}^0, \mathcal{T}^1, \mathcal{T}^2, \dots, \mathcal{T}^j, \dots$  be a *sequence of partitions* of  $\Omega$ . We say that this sequence is a *recursive dyadic partition* (RDP) if it satisfies the following two conditions: (i)  $\mathcal{T}^j$  consists of  $2^j$  blocks:  $\mathcal{T}^j = \{A_{j,k} : k = 0, 1, \dots, 2^j - 1\}$ ; (ii)  $\mathcal{T}^{j+1}$  is obtained by dividing each set in  $\mathcal{T}^j$  into two pieces, i.e.,  $A_{j,k} = A_{j+1,2k} \cup A_{j+1,2k+1}$  for all  $j \geq 0$  and  $k = 0, 1, \dots, 2^j - 1$ .

We call an RDP *canonical* if the sequence of partitions satisfy two additional conditions: (iii) if the partition blocks  $A_{j,k}$  are rectangles of the form

$$A_{j,k} = [a_{j,k}^{(1)}, b_{j,k}^{(1)}] \times [a_{j,k}^{(2)}, b_{j,k}^{(2)}] \times \dots \times [a_{j,k}^{(m)}, b_{j,k}^{(m)}].$$

and (iv)  $A_{j+1,2k}$  and  $A_{j+1,2k+1}$  are produced by dividing  $A_{j,k}$  into two halves at the middle of one of  $A_{j,k}$ 's *divisible* dimensions.

A rectangular partition block  $A_{j,k}$  is *divisible* in dimension  $d$  if  $A_{j,k}$  is supported on at least two values in that dimension, i.e.,  $a_{j,k}^{(d)} < b_{j,k}^{(d)}$ . In this case, if  $A_{j,k}$  is divided in dimension  $d$ , then its children  $A_{j+1,2k}$  and  $A_{j+1,2k+1}$  are given by

$$[a_{j+1,2k}^{(d)}, b_{j+1,2k}^{(d)}] = [a_{j,k}^{(d)}, (a_{j,k}^{(d)} + b_{j,k}^{(d)})/2]$$

and

$$[a_{j+1,2k+1}^{(d)}, b_{j+1,2k+1}^{(d)}] = [(a_{j,k}^{(d)} + b_{j,k}^{(d)})/2 + 1, b_{j,k}^{(d)}],$$

while

$$[a_{j+1,2k}^{(d')}, b_{j+1,2k}^{(d')}] = [a_{j+1,2k+1}^{(d')}, b_{j+1,2k+1}^{(d')}] = [a_{j,k}^{(d')}, b_{j,k}^{(d')}]$$

for all  $d' \neq d$ .

Any canonical RDP on  $\Omega$  will have exactly  $J + 1$  levels, i.e.,  $\mathcal{T}^0, \mathcal{T}^1, \dots, \mathcal{T}^J$ . The  $j$ th level partition  $\mathcal{T}^j$  consists of  $2^j$  rectangular pieces of equal size, each covering  $n/2^j$  locations in  $\Omega$ . From now on, we simply use RDP to refer to canonical ones when this causes no confusion.

### 2.1.3 RDPs and permutations

Each RDP can be represented by a  $J$  level bifurcating tree with the partition blocks in  $\mathcal{T}^j$  forming the  $2^j$  nodes in the  $j$ th level of the tree. As such, we can use  $\mathcal{T} = \cup_{j=0}^J \mathcal{T}^j$  to represent the RDP. Each node in the  $J$ th level corresponds to a unique location in  $\Omega$ , and is called “atomic” as it contains a single element. We shall interchangeably refer to an RDP as a “tree”, and to the partition blocks as “nodes”.

Given the RDP  $\mathcal{T}$ , each location  $\mathbf{s} \in \Omega$  falls into a unique branch of  $\mathcal{T}$ , that is,  $\Omega = A_{0,0} \supset A_{1,k_1(\mathbf{s})} \supset A_{2,k_2(\mathbf{s})} \supset \dots \supset A_{J,k_J(\mathbf{s})} = \{\mathbf{s}\}$ , with  $A_{j,k_j(\mathbf{s})}$  being the node in the  $j$ th level to which  $\mathbf{s}$  belongs. Accordingly, the RDP  $\mathcal{T}$  induces a unique vectorization of the locations in  $\Omega$  such that  $\mathbf{s}$  corresponds to the  $t(\mathbf{s})$ th element of the vector where  $t(\mathbf{s}) = \sum_{l=1}^J 2^{J-l} \cdot e_l(\mathbf{s})$  with  $e_l = k_l(\mathbf{s}) \bmod 2$ , indicating the branch of the tree  $\mathbf{s}$  falls into at level  $l$ . As such,  $\mathcal{T}$  induces a permutation of the  $n$  locations, and we let  $\pi_{\mathcal{T}}$  denote this permutation.

As an illustration, Figure 1 presents an RDP and the induced permutation using a toy  $4 \times 4$  image (so  $m = 2$  and  $J_1 = J_2 = 2$ ). We index pixels in the true image from 0 to 15. In addition, we assume that the underlying function takes only two values—1 and 2—on the 16 locations, represented by the white and the red colors, respectively. The demonstrated RDP corresponds well to the structure of the underlying signal, which would result in an effective 1D wavelet analysis on the vectorized observation.

We shall now utilize the relationship between RDPs and permutations to construct a prior on the latter. Before that, we shall simplify our notations a little. Note that while what the  $(j, k)$ th node  $A_{j,k}$  is depends on the RDP  $\mathcal{T}$ , different RDPs can share common nodes—the  $(j, k)$ th node in one  $\mathcal{T}$  may be the same as the  $(j, k')$ th node in another. (Note that the level of the node must be the same in either RDP.) In the following, we will need to specify quantities that only depend on the node regardless of the RDP tree  $\mathcal{T}$  it arises from. A succinct way for expressing such quantities is to write them as a mapping from  $\mathcal{A}$  to  $\mathbb{R}$ , where  $\mathcal{A}$  denotes the collection of all sets that *could* be nodes in *some* RDP, or equivalently,  $\mathcal{A}$  is the totality of nodes in all RDPs. (This is to be distinguished from the collection of nodes in any particular RDP, which is denoted by  $\mathcal{T}$ .) It is worth noting  $\mathcal{A}$  is a finite set.

Now we may define  $\rho_{j,k}$  in such a way that its value only depends on what the set  $A_{j,k}$  is, regardless of the RDP  $\mathcal{T}$  to which it belongs. In this case we can let  $\rho_{j,k} = \rho(A_{j,k})$ , where  $\rho(\cdot)$  is a mapping from  $\mathcal{A}$  to  $[0, 1]$ . While a set  $A \in \mathcal{A}$

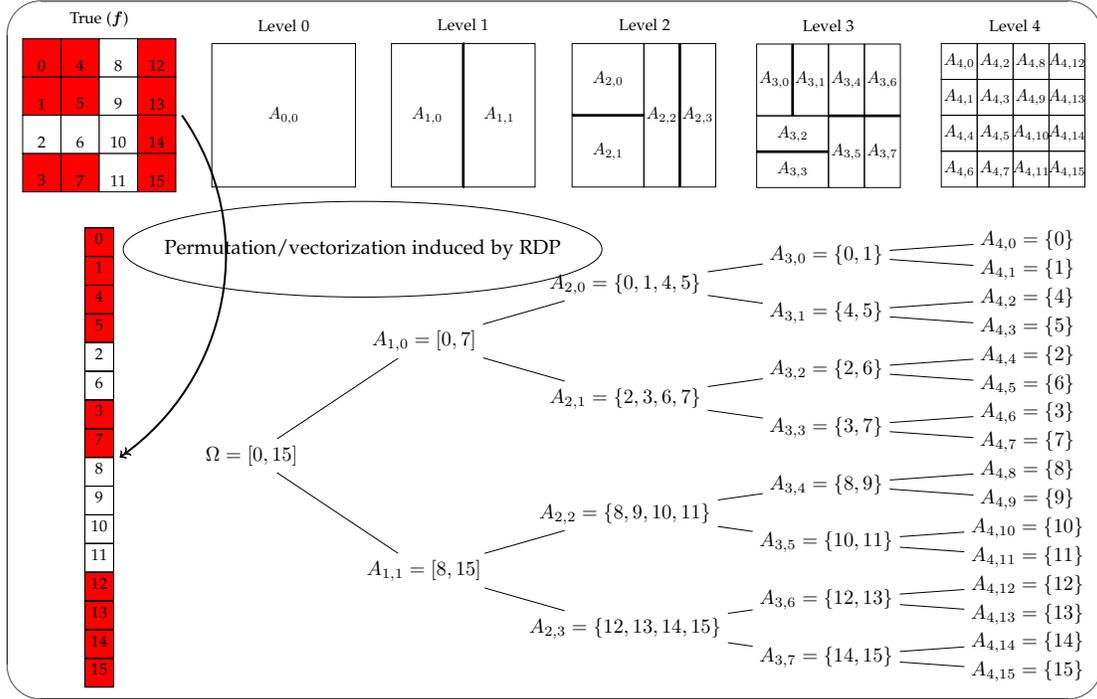


Fig. 1. Illustration of the correspondence between RDPs and permutations. In the tree representation,  $A_{2,0} = \{0, 1, 4, 5\}$  means the node  $A_{2,0}$  contains the (0, 1, 4, 5)th elements of  $\Omega$ . The coloring code for the observations is red for 2 and white for 1. From level 0 to level 3, edges that are thicker than others are the partitions of the current level; nodes at the last level are all atomic.

might be the  $A_{j,k}$  in one RDP and  $A_{j,k'}$  in another, the corresponding  $\rho(A)$  value will then be the same under this mapping based specification. The mapping-based notation such as  $\rho(\cdot)$  allows various parameters to be specified in a node-specific (rather than RDP-specific) manner. This observation has extremely important computational implications—as we will show later, the space of nodes  $\mathcal{A}$  for all canonical RDPs is of a cardinality linear in the size of  $\Omega$ , while that of canonical RDPs is exponential in  $n$ . (See Proposition 1 in the Supplementary Materials.) Therefore it is exactly the ability to carrying out the computation for the posterior in a node-specific manner that allows us to achieve linear complexity in our inference algorithm. Moreover, this notation will also help elucidate derivations on the posterior.

### 2.1.4 Priors on RDPs: random RDP

Our strategy of representing multi-dimensional functions using vectors will only pay off if the vectorization of  $\Omega$  can result in an efficient characterization of the data, thereby leading to stronger energy concentration under wavelet transforms. For example, the RDP illustrated in Figure 1 will lead to particularly efficient inference of the corresponding function. In general, the true optimal vectorization—or the corresponding RDP—is unknown, and one shall rely on the data to learn the RDPs that induce “good” vectorizations.

We aim to achieve this in a hierarchical Bayesian approach by treating the RDP as a latent structure and placing a prior on the RDP. We consider the following prior on the RDP originally proposed in the context of density estimation [21], [22], which is specified in a node-specific fashion and leads to very efficient node-based posterior inference algorithms that scale linearly in  $n$ , the size of  $\Omega$ .

We describe the prior as a simple generative procedure for an RDP in an inductive manner. First,  $\mathcal{T}^0 = \{\Omega\}$  by definition. Now suppose we have generated  $\mathcal{T}^0, \mathcal{T}^1, \dots, \mathcal{T}^j$  for some  $0 \leq j \leq J-1$ , then  $\mathcal{T}^{j+1}$  is generated as follows. For each  $A_{j,k} \in \mathcal{T}^j$ , let  $\mathcal{D}(A_{j,k}) \subset \{1, 2, \dots, m\}$  be the collection of its divisible dimensions. We randomly draw a dimension in  $\mathcal{D}(A_{j,k})$ , and divide  $A_{j,k}$  in that dimension to get  $A_{j+1,2k}$  and  $A_{j+1,2k+1}$ . In particular, we let  $\lambda_d(A_{j,k})$  be the probability for drawing the  $d$ th dimension, where  $\sum_{d=1}^m \lambda_d(A_{j,k}) = 1$  and  $\lambda_d(A_{j,k}) = 0$  for  $d \notin \mathcal{D}(A_{j,k})$ . In many problems, *a priori* one has no reason to favor dividing any particular dimension over another, and a default specification is to set

$$\lambda_d(A_{j,k}) = 1/|\mathcal{D}(A_{j,k})| \cdot \mathbf{1}_{\{d \in \mathcal{D}(A_{j,k})\}},$$

where  $\mathbf{1}_E$  is the indicator function of whether  $E$  holds or not. This completes the inductive generation of  $\mathcal{T}^{j+1}$ . The procedure will terminate after  $\mathcal{T}^J$  is generated as all nodes in  $\mathcal{T}^J$  are atomic with no divisible dimensions.

The above generative mechanism forms a probability distribution on the space of RDPs, which is called the *random recursive dyadic partition* (RRDP) distribution, and it is specified by the collection of selection probabilities  $\lambda_d(\cdot)$  defined on all *potential* nodes. We write

$$\mathcal{T} \sim \text{RRDP}(\boldsymbol{\lambda}),$$

where  $\{\boldsymbol{\lambda}(A) : A \in \mathcal{A}\}$ , and  $\boldsymbol{\lambda}(A) = (\lambda_1(A), \lambda_2(A), \dots, \lambda_m(A))'$ , that is,  $\boldsymbol{\lambda}$  is a mapping from  $\mathcal{A}$  to the  $(m-1)$ -dimensional simplex.

It is worth noting that the RRDP is a restricted version of the more general Bayesian classification and regression tree (CART) prior [23], [24]. The main constraint in RRDP compared to the general Bayesian CART is that the former

is supported on canonical RDPs only—that is, each dyadic partition must be an even split, occurring at the middle of the range in one of the divisible dimensions. This additional restriction ensures that the cardinality of  $\mathcal{A}$  is linear in  $n$ , thereby reducing the computational complexity required for inference to  $O(n)$ .

## 2.2 Recipes for Bayesian inference

In this section, we present recipes for deriving and sampling from the posterior of our Bayesian model, and for evaluating posterior summaries such as the posterior mean of  $\mathbf{f}$ . We note that the marginal posterior of the RDP  $\mathcal{T}$  is the key component for posterior inference, because once conditional on  $\mathcal{T}$ , our model reduces to a standard Bayesian wavelet regression for which closed-form conditional posteriors are readily available under common prior specifications.

Interestingly, when a Haar basis is adopted in the wavelet regression model, the marginal posterior of  $\mathcal{T}$  can be calculated analytically in closed form through a recursive algorithm that is operationally similar to Mallat’s pyramid algorithm, achieving a linear computational complexity  $O(n)$ .

### 2.2.1 Exact Bayesian inference under Haar basis

The Haar wavelet basis is unique in that the  $(j, k)$ th wavelet coefficient under the vectorization induced by any RDP  $\mathcal{T}$  is determined by only the locations inside the node  $A_{j,k}$ . We call this property of the Haar basis *node-autonomy* and say that inference under the Haar basis is *node-autonomous*. Specifically, for all RDPs in which  $A$  is a node and is divided in the  $d$ th direction, the corresponding Haar wavelet coefficient associated with the node  $A$  is given by

$$w_d(A) = 1/\sqrt{|A|} \cdot \left( \sum_{\mathbf{x} \in A_l^{(d)}} y(\mathbf{x}) - \sum_{\mathbf{x} \in A_r^{(d)}} y(\mathbf{x}) \right)$$

where  $A_l^{(d)}$  and  $A_r^{(d)}$  represent the two children nodes if  $A$  is divided in the  $d$ th dimension and  $|A| = 2^{J-j}$  is the total number of locations in  $A$ . In contrast, wavelet coefficients from wavelet bases with longer support than Haar are not node-autonomous—not only does the coefficient associated with  $A$  depend on the observations within  $A$  but on those in other (often but not always adjacent) nodes in  $\mathcal{T}$  as well.

Node-autonomy enables the posterior to be computed in a node-specific fashion, avoiding integration in the much larger space of RDPs. Consequently, exact inference can be completed in a computational complexity of the same scale as the total number of all potential nodes in RDPs, which is equal to  $\prod_{i=1}^m (2n_i - 1) = O(2^m n)$ .

Next we lay out the general strategy for inference. We show through two theorems that the marginal posterior of the RDP  $\mathcal{T}$  is computable in analytically through a recursive algorithm that resembles Mallat’s pyramid algorithm for two very popular classes of Bayes wavelet regression models—(i) those that model each wavelet coefficient independently given  $\mathcal{T}$  (Theorem 1); and (ii) those that induce a hidden Markov model (HMM) for incorporating dependency among the wavelet coefficients given  $\mathcal{T}$  (Theorem 2).

**Theorem 1.** *Suppose  $\mathcal{T} \sim \text{RRDP}(\boldsymbol{\lambda})$  and given the Haar DWT under  $\mathcal{T}$ , one models the wavelet coefficients independently, i.e.,  $(w_{j,k}, z_{j,k}) \stackrel{\text{ind}}{\sim} p_{j,k}(w, z | \boldsymbol{\phi})$  for all  $(j, k)$ , where  $\boldsymbol{\phi}$  represents the hyperparameters of the Bayesian wavelet regression model. Then the marginal posterior of  $\mathcal{T}$  is still an RRDP. Specifically,  $\mathcal{T} | \mathbf{y} \sim \text{RRDP}(\tilde{\boldsymbol{\lambda}})$  where the posterior selection probability mapping  $\tilde{\boldsymbol{\lambda}}$  is given as*

$$\tilde{\lambda}_d(A) = \lambda_d(A) M_d(A) \Phi(A_l^{(d)}) \Phi(A_r^{(d)}) / \Phi(A)$$

for any non-atomic  $A \in \mathcal{A}$  where  $M_d(A)$  is the marginal likelihood contribution from the wavelet coefficient on node  $A$  if it is a node in  $\mathcal{T}$  and divided in dimension  $d$ , i.e.,  $M_d(A) = \int p_{j,k}(w_d(A), z | \boldsymbol{\phi}) dz$  and  $\Phi : \mathcal{A} \rightarrow [0, \infty)$  is a mapping defined recursively (i.e., its value on  $A$  depends on its values on  $A$ ’s children) as

$$\Phi(A) = \sum_{d \in \mathcal{D}(A)} \lambda_d(A) M_d(A) \Phi(A_l^{(d)}) \Phi(A_r^{(d)})$$

if  $A$  is not atomic, and  $\Phi(A) = 1$  if  $A$  is atomic.

Remark:  $\Phi(\Omega)$  is the overall marginal likelihood. It is a function of the hyperparameters  $\boldsymbol{\phi}$ , and can be used for specifying the hyperparameters  $\boldsymbol{\phi}$  in an empirical Bayes strategy using maximum marginal likelihood estimation (MMLLE).

**Theorem 2.** *Suppose  $\mathcal{T} \sim \text{RRDP}(\boldsymbol{\lambda})$  and given  $\mathcal{T}$  under a Haar DWT, one models the wavelet coefficients conditionally independently given a set of latent state variables  $\mathcal{S} = \{S_{j,k} : j = 0, 1, 2, \dots, J, k = 0, 1, \dots, 2^j - 1\}$*

$$(w_{j,k}, z_{j,k}) | S_{j,k} = s \stackrel{\text{ind}}{\sim} p_{j,k}^{(s)}(w, z | \boldsymbol{\phi}) \quad \text{for all } (j, k)$$

where  $S_{j,k} \in \{1, 2, \dots, K\}$  is a latent state variable associated with  $(j, k)$ . Also, suppose the collection of all latent variables is modeled as a top-down Markov tree (MT) with transition kernel  $\boldsymbol{\rho}$ ,  $\mathcal{S} \sim \text{MT}(\boldsymbol{\rho})$ , i.e.,

$$P(S_{j,k} = s' | S_{j-1, \lfloor k/2 \rfloor} = s) = \rho_j(s, s')$$

where  $\rho_j(\cdot, \cdot)$  is the transition kernel of the Markov model which is allowed to be different over  $j$ . Then the joint marginal posterior of  $(\mathcal{T}, \mathcal{S})$  can be specified fully as the following sequential generative process. Suppose  $\mathcal{T}^0, \mathcal{T}^1, \dots, \mathcal{T}^j$  and the latent variables up to level  $j - 1$  have been generated. (To begin, we have  $j = 0$  and  $\mathcal{T}^0 = \{\Omega\}$ .) Then the state variables in level  $j$ , are generated from the following posterior transition probabilities

$$\begin{aligned} & P(S_{j,k} = s' | S_{j-1, \lfloor k/2 \rfloor} = s, \mathcal{T}^{(j)}, \mathbf{y}) \\ &= \rho_j(s, s') \sum_{d \in \mathcal{D}(A)} \lambda_d(A) M_d^{(s')} (A) \Phi_{s'}(A_l^{(d)}) \Phi_{s'}(A_r^{(d)}) / \Phi_s(A), \end{aligned}$$

where  $A$  is the node  $A_{j,k}$  in  $\mathcal{T}^j$ . Given  $S_{j,k} = s'$ , suppose  $j < J$ , then  $\mathcal{T}^{j+1}$  is generated by drawing  $D_{j,k}$  from a multinomial with probabilities  $\boldsymbol{\lambda}(A)$  such that

$$\begin{aligned} & P(D_{j,k} = d | S_{j,k} = s', \mathcal{T}^{(j)}, \mathbf{y}) \\ &= \frac{\lambda_d(A) M_d^{(s')} (A) \Phi_{s'}(A_l^{(d)}) \Phi_{s'}(A_r^{(d)})}{\sum_{d' \in \mathcal{D}(A)} \lambda_{d'}(A) M_{d'}^{(s')} (A) \Phi_{s'}(A_l^{(d')}) \Phi_{s'}(A_r^{(d')})}, \end{aligned}$$

where  $M_d^{(s)}(A)$  is the marginal likelihood contribution from the wavelet coefficient on node  $A$  if it is a node in

$\mathcal{T}$ , is divided in dimension  $d$  in  $\mathcal{T}$ , and its latent state is  $s$ . That is,  $M_d^{(s)}(A) = \int p_{j,k}^{(s)}(w_d(A), z | \phi) dz$  and  $\Phi = (\Phi_1, \Phi_2, \dots, \Phi_K) : \mathcal{A} \rightarrow [0, \infty)^K$  is a vector-valued mapping defined recursively as  $\Phi_s(A) = \sum_{s'} \rho_j(s, s') \sum_{d \in \mathcal{D}(A)} \lambda_d(A) M_d^{(s')} (A) \Phi_{s'}(A_i^{(d)}) \Phi_{s'}(A_r^{(d)})$  if  $A$  is not atomic, and  $\Phi_s(A) = 1$  if  $A$  is atomic, for all  $s \in \{1, 2, \dots, K\}$ , where  $j$  is the level of  $A$ .

Once the marginal posterior of  $\mathcal{T}$  is computed through Theorem 1 or Theorem 2, the full joint posterior is available as the conditional posterior of the rest of our model given  $\mathcal{T}$  is available for common Bayesian wavelet regressions. (More details are given in Section 2.2.2.) Then standard Bayesian inference can proceed.

In particular, one can draw samples for  $(\mathcal{T}, \mathcal{S})$  from their marginal posterior given in Theorem 2. Then given  $(\mathcal{T}, \mathcal{S})$ , one can further sample  $\mathbf{z}$  from the conditional posterior corresponding to the chosen wavelet regression model, and Bayesian inference can proceed in the usual manner. For example, one can obtain posterior samples of the underlying function  $\mathbf{f}$  by first drawing samples

$$(\mathcal{T}^{(1)}, \mathcal{S}^{(1)}, \mathbf{z}^{(1)}), (\mathcal{T}^{(2)}, \mathcal{S}^{(2)}, \mathbf{z}^{(2)}), \dots, (\mathcal{T}^{(B)}, \mathcal{S}^{(B)}, \mathbf{z}^{(B)}).$$

Then for the  $b$ th draw, we can compute the corresponding function  $\mathbf{f}^{(b)}$  using the inverse DWT

$$\mathbf{f}^{(b)} = \pi_{\mathcal{T}^{(b)}}^{-1} \left( W^{-1} \mathbf{z}^{(b)} \right),$$

where  $\pi_{\mathcal{T}}^{-1}$  denotes the inverse permutation corresponding to an RDP  $\mathcal{T}$ . Based on the posterior samples of  $\mathbf{f}$ , we can construct pointwise credible bands and estimate the posterior mean  $\mathbb{E}(\mathbf{f} | \mathbf{y})$ . We can apply Rao-Blackwellization and obtain the following estimate for the posterior mean

$$\mathbb{E}(\mathbf{f} | \mathbf{y}) \approx \frac{1}{B} \sum_{b=1}^B \pi_{\mathcal{T}^{(b)}}^{-1} \left( W^{-1} \mathbb{E}(\mathbf{z}^{(b)} | \mathcal{T}^{(b)}, \mathbf{y}) \right).$$

For several popular Bayesian wavelet regression models, the posterior mean can actually be computed analytically through message passing (MP) without posterior sampling when the Haar basis is adopted. We next turn to briefly reviewing these wavelet regression models in Section 2.2.2, and defer the MP algorithm (Theorem 3) to Supplementary Materials.

### 2.2.2 Examples of compatible Bayesian wavelet regression models

So far we have kept the description of the Bayesian wavelet regression model general, using generic notations such as  $p(w_{j,k}, z_{j,k} | \phi)$  and  $p(w_{j,k}, z_{j,k} | S_{j,k}, \phi)$  without spelling out the details. Next we describe some of the most popular Bayesian wavelet regression models. They indeed take these general forms and therefore our framework is applicable to them.

A popular class of Bayesian wavelet regression models for achieving adaptive shrinkage of  $\mathbf{z}$  utilize the so-called spike-and-slab prior, which introduces a latent binary random variable  $S_{j,k}$  for each  $(j, k)$  such that

$$z_{j,k} | S_{j,k} \stackrel{\text{ind}}{\sim} (1 - S_{j,k}) \delta_0(z_{j,k}) + S_{j,k} \gamma(z_{j,k} | \tau_j, \sigma)$$

where  $\delta_0(\cdot)$  is a point mass at 0, and  $\gamma(\cdot | \tau_j, \sigma)$  is a fixed unimodal symmetric density that possibly depends on  $\sigma$  and another scale parameter  $\tau_j$ . A common choice of  $\gamma(\cdot | \tau_j, \sigma)$  is the normal distribution with mean 0 and variance  $\tau_j \sigma^2$ , denoted by  $\phi(\cdot | 0, \sqrt{\tau_j} \sigma)$ , while heavy-tailed priors including the Laplace and quasi-Cauchy distributions [25] also enjoy desirable theoretical properties. Specifically, the function  $\gamma(x | \tau_j, \sigma)$  is

$$\gamma(x | \tau_j, \sigma) = a \exp(-a|x/\sigma|) / (2\sigma)$$

for Laplace priors where  $a = \sqrt{2/\tau_j}$ , and

$$\gamma(x | \tau_j, \sigma) = (2\pi)^{-1/2} \{1 - |x/\sigma| \cdot \tilde{\Phi}(|x/\sigma|) / \phi(x/\sigma)\} / \sigma$$

for quasi-Cauchy priors with  $\tilde{\Phi}(x) = \int_x^\infty \phi(t | 0, 1) dt$ .

Many authors [12], [13], [15], [26] adopt independent priors on the latent shrinkage state variable  $S_{j,k}$

$$S_{j,k} \stackrel{\text{ind}}{\sim} \text{Bern}(\rho_{j,k}).$$

One way to specify  $\rho = \{\rho_{j,k}, 0 \leq k < 2^j, 0 \leq j \leq J-1\}$  that properly controls for multiplicity is  $\rho_{j,k} \propto 2^{-j}$ . The specification of  $\tau = \{\tau_j, 0 \leq j \leq J-1\}$  of course depends on the choice of  $\gamma(\cdot | \tau_j, \sigma)$ . For instance, if one uses  $\tau_j = 2^{-\alpha j} \tau_0$  for the normal and Laplace prior, this leads to the reduced parameter  $\tau = (\alpha, \tau_0)$ . One can use  $\tau_j \equiv 1$  for the quasi-Cauchy prior. Other authors [11], [27] show that introducing Markov dependency into the latent shrinkage states can substantially improve inference by allowing effective borrowing of information across the location and scale.

Carrying out inference under WARP requires the conditional posterior of  $z_{j,k}$  given  $(\mathcal{T}, \mathcal{S})$ . For the above popular models, this posterior is given by

$$z_{j,k} | S_{j,k}, \mathbf{y} \stackrel{\text{ind}}{\sim} (1 - S_{j,k}) \delta_0(z_{j,k}) + S_{j,k} f_1(z_{j,k} | w_{j,k}, \tau_j, \sigma),$$

where  $f_1(z_{j,k} | w_{j,k}, \tau_j, \sigma) \propto \phi(w_{j,k} | z_{j,k}, \sigma) \cdot \gamma(z_{j,k} | \tau_j, \sigma)$ . The function  $f_1(z_{j,k} | w_{j,k}, \tau_j, \sigma)$  is analytically available if  $\gamma(\cdot | \tau_j, \sigma)$  is the density of normal, Laplace, or quasi-Cauchy distributions. For the normal prior where  $\gamma(\cdot | \tau_j, \sigma) = \phi(\cdot | 0, \sqrt{\tau_j} \sigma)$ ,  $f_1(\cdot | w_{j,k}, \tau_j, \sigma)$  is the density of  $N(w_{j,k} / (1 + \tau_j^{-1}), \sigma^2 / (1 + \tau_j^{-1}))$ . For Laplace and quasi-Cauchy priors, analytical forms of  $f_1(\cdot | \tau_j, \sigma)$  are available in [25, Sec. 2.3]. As it is often the mean corresponding to  $f_1$  that is needed for posterior estimation, we here give the closed forms of the means by integrating out  $z_{j,k}$  with respect to its posterior distribution. Let the corresponding mean function be  $\mu_1(w_{j,k}, \tau_j, \sigma)$ , which is given by

$$w_{j,k} / (1 + \tau_j^{-1})$$

for normal priors,

$$w_{j,k} - \sigma \frac{a \{e^{-aw_{j,k}/\sigma} \Phi(w_{j,k}/\sigma - a) - e^{aw_{j,k}/\sigma} \tilde{\Phi}(w_{j,k}/\sigma + a)\}}{e^{-aw_{j,k}/\sigma} \Phi(w_{j,k}/\sigma - a) + e^{aw_{j,k}/\sigma} \tilde{\Phi}(w_{j,k}/\sigma + a)}$$

for Laplace priors, and

$$w_{j,k} \left\{ 1 - \exp\left(-\frac{w_{j,k}^2}{2\sigma^2}\right) \right\}^{-1} - 2 \left(\frac{w_{j,k}}{\sigma}\right)^{-1}$$

for quasi-Cauchy priors.

For these wavelet regression models that adopt the spike-and-slab setup, by Theorem 2 we can derive a fully

conjugate posterior that takes the same form as the prior. In particular, for each  $A \in \mathcal{A}$ , under the normal prior for  $\gamma(\cdot | \tau_j, \sigma)$ , applying Theorem 2 shows that

- The marginal likelihood contribution from the data within node  $A$  if  $A$  is divided in dimension  $d$  is:

$$M_d^{(s)}(A) = \frac{1}{\sqrt{2\pi(1+s\tau_j)\sigma^2}} \exp\left\{-\frac{w_d(A)^2}{2\sigma^2(1+s\tau_j)}\right\}$$

for  $s = 0, 1$ .

- The posterior spike probability on  $A$  if  $A$  is divided in dimension  $d$  is:

$$\tilde{\rho}_d(A) = \rho(A)M_d^{(1)}(A)/M_d(A),$$

where  $M_d(A) = \rho(A)M_d^{(1)}(A) + (1 - \rho(A))M_d^{(0)}(A)$ .

In most practical problems, the variation in the function value within each partition block will eventually become negligible with respect to the noise level, and so further division within such homogeneous blocks will not improve statistical efficiency and could lead to overfitting. For example, in Figure 1 the partition in the upper left block (Level 3) along with its descendants is not necessary. Thus it is also desirable to incorporate adaptivity in the depth of the wavelet tree along each subbranch and allow it to be terminated earlier than reaching level  $J$  depending on how smooth the function is across the index space. This consideration is closely related to the idea of adaptive block shrinkage [28] in the frequentist wavelet regression analysis. Once there is little evidence for any interesting structure within a subset of the index space, then the function value within that subset can be shrunk to a constant. That is, the wavelet tree is “pruned” there. Remarkably, wavelet models with such pruning are also compatible with our WARP framework and can be readily achieved by introducing a pruning indicator to accompany  $S_{j,k}$ . We refer interested readers to Supplementary Materials for additional technical details on how to incorporate pruning.

For the Haar basis, the posterior mean  $E(\mathbf{f}|\mathbf{y})$  for the above wavelet models can be evaluated analytically through recursive message passing without any Monte Carlo sampling for Bayesian wavelet regression models that adopt the spike-and-slab setup along with optional pruning of the wavelet tree, which contains the models without optional pruning as special cases with zero pruning probabilities. For completeness, we describe this strategy in the Supplementary Materials and will use it to compute  $E(\mathbf{f}|\mathbf{y})$  in our numerical examples.

### 3 EXPERIMENTS

In this section, we conduct extensive experiments to evaluate the performance of our proposed framework in the image reconstruction task in terms of both estimation accuracy and computational scalability. Applications to other image processing tasks are discussed in Section 4. We compare WARP to a number of state-of-the-art wavelet, non-wavelet, and deep neural network-based methods available in the literature for denoising 2D images. We provide results on denoising 3D images in the Supplementary Materials.

Throughout these experiments we apply WARP with independent spike-and-slab Bayesian wavelet regression models under the Haar basis along with optional pruning.

Our prior specification is as follows:  $\rho(A) = \min(1, 2^{-\beta j} C)$  for  $A$  in the  $j$ th resolution (for  $j < J$ ),  $\tau_j = 2^{-\alpha j} \tau_0$ , and  $\eta(A) = \eta_0$  for all  $A$ ; we set  $\sigma^2$  to an estimate based on the finest scale wavelet coefficients [8]; all other parameters in  $\phi = (\alpha, \beta, \sigma^2, \tau, C, \eta_0)$  are estimated by maximizing the marginal likelihood (available in a closed form as  $\Phi(\Omega)$  from our recursive message passing algorithm) at a set of grid points. Supplementary Materials contain a sensitivity analysis showing that WARP is generally robust to the values of its hyperparameters. Therefore we recommend a grid search on a small set rather than a full optimization as the default tuning method. Gaussian noise with standard deviation  $\sigma$  is added to the true images and we apply all methods to the noisy observations for image reconstruction. For WARP, we use the posterior mean as the reconstructed image, which is analytically attainable through Theorem 3.

#### 3.1 Image reconstruction using ImageNet data

We use 100 test images randomly chosen from the ImageNet dataset [29] to evaluate selected methods in reconstructing images of various structures. ImageNet is originally used for large-scale visual recognition in the community of computer vision, and we here use its Fall 2011 release (consisting of 14,197,122 urls). We compare our method with eight existing wavelet and non-wavelet approaches with available software: 1-dimensional Haar denoising operated on vectorized observation [25] or 1D-Haar, translation-invariant 2D Haar estimation [14] or TI-2D-Haar, shape-adaptive Haar wavelets [30] or SHAH, adaptive weights smoothing [31] or AWS, Bayesian smoothing method using the Chinese restaurant process [32] or CRP, coarse-to-fine wedgelet [33] or Wedgelet, nonparametric Bayesian dictionary learning proposed by [34] or BPFA, and the conventional running median method or RM. We apply the cycle spinning technique to remove visual artifacts in image reconstruction [35], [36] for the methods of WARP, 1D-Haar, SHAH, AWS, CRP, Wedgelet and RM, by averaging 121 local shifts (a step size up to 5 pixels in each direction). TI-2D-Haar is translation invariant and BPFA includes cycle spinning based on patches, and thus no additional cycle spinning is needed for these two methods. For each method, we calculate the mean squared error (MSE) and mean absolute error (MAE) to measure its accuracy, and time each method based on one replication ran on MacBook Pro with 2.7 GHz Intel core i7 CPU and 16GB RAM. We implement the methods using publicly available code, either in R (1D-Haar, SHAH, and AWS) or Matlab (TI-2D-Haar, CRP, Wedgelet, BPFA, and RM). WARP is available in both R and Matlab, and we use the R version to time it.

Figure 2 presents the average MSEs and MAEs of all methods where  $\sigma$  varies from 0.1 to 0.7. We can first see that the proposed hierarchical adaptive partition improved the basic wavelet regression significantly (compare 1D-Haar vs. WARP) for all scenarios. In fact, WARP is uniformly the best method under both metrics for all scenarios, with the performance lead over other methods widening as the noise

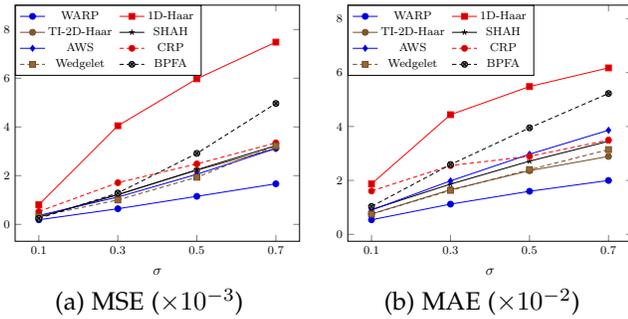


Fig. 2. Comparison of various methods based on 100 randomly selected  $512 \times 512$  images from ImageNet. The method of running median is off the chart (not plotted here). The maximum standard errors at each  $\sigma$  among all methods are  $(0.001, 0.042, 0.071, 0.058) \times 10^{-3}$  for MSE,  $(0.002, 0.062, 0.065, 0.058) \times 10^{-2}$  for MAE, respectively. The running time of each method in seconds is 7.2 (WARP), 76.9 (SHAH), 7.9 (AWS), 10.7 (CRP), 8.7 (Wedgelet),  $2.1 \times 10^3$  (BPFA), and less than 1 (1D-Haar, TI-2D-Haar, RM), based on one test image without cycle spinning at  $\sigma = 0.3$  including both tuning and estimation steps.

level increases. The sensitivity analysis in the Supplemen-

tary Materials indicates that the method of WARP is robust to hyperparameters and choices of  $\gamma$ .

WARP is computationally efficient, benefiting from the conjugacy of random recursive partition and closed form expression in Theorem 3. WARP is the fastest adaptive approach among SHAH, AWS, CRP, Wedgelet, and BPFA. (The computing times are given in the caption of Figure 2.) Section 3.2 further compares the scalability of selected methods using images of various sizes.

### 3.2 Scalability

Next we verify the linear complexity of the WARP framework using both 2D and 3D images. Usually there are various ways to tune each method, and we focus on the estimation step given tuning parameters for all methods to make a fair comparison. For WARP, one actually may choose the tuning parameters from a smaller image by downsampling without loss of much accuracy, in view of its insensitivity to hyperparameters (Section D in the Supplementary Materials).

Figure 4 (a) compares the scalability of selected methods

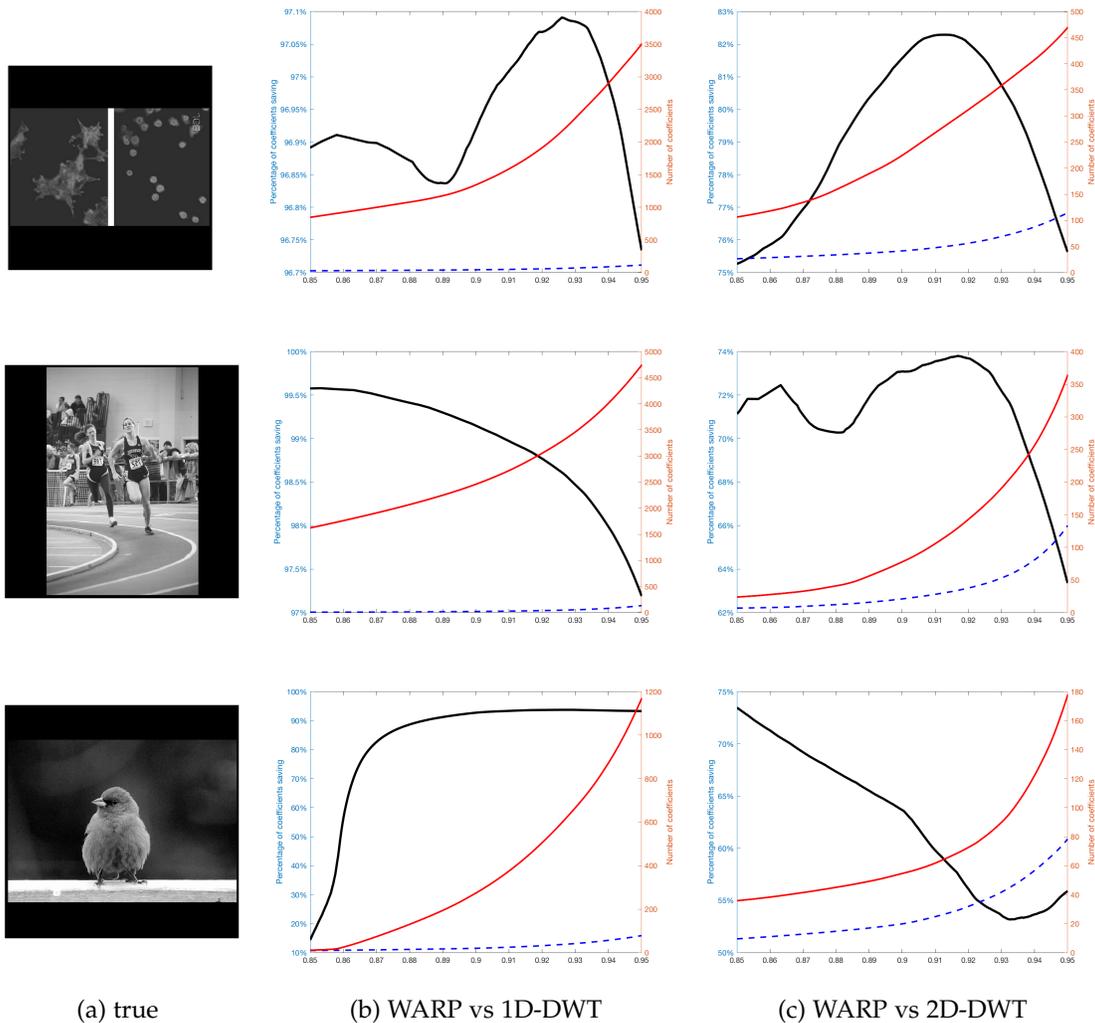


Fig. 3. Comparison of energy concentration for three methods—WARP, 1D Haar, and 2D Haar—on ImageNet images. Column (a) plots the true image, Column (b) compares WARP versus 1D DWT, and Column (c) compares WARP versus 2D DWT. In Columns (b) and (c), the red and blue lines correspond to the right  $y$  axis, plotting the number of coefficients to attain a specific energy level ( $x$  axis) by deterministic DWT and WARP, respectively. The black curve corresponds to the left  $y$  axis and is 100% less the ratio of the blue and red curves, indicating the percentage reduction in the number of wavelet coefficients to achieve the same sum of squares by WARP.

TABLE 1  
MSE ( $\times 10^{-3}$ ) of WARP and DnCNN on 12 widely used test images and BSD68.

$\sigma$		12 widely used test images												BSD68
		1	2	3	4	5	6	7	8	9	10	11	12	
0.2	WARP	2.89	1.42	2.49	4.07	3.65	3.39	3.17	1.69	4.33	2.55	2.50	2.59	3.75
	DnCNN	2.77	1.65	2.55	3.37	2.90	3.49	2.88	1.73	3.98	2.44	2.33	2.68	3.35
0.4	WARP	5.69	3.15	5.57	7.97	8.21	6.01	6.64	3.15	6.38	4.46	4.16	4.67	6.10
	DnCNN	15.64	14.14	15.39	15.93	15.85	16.56	15.53	13.28	15.37	14.57	13.49	14.69	14.90
0.6	WARP	8.23	4.31	8.12	10.86	12.26	8.36	9.57	4.40	7.96	6.06	5.66	6.22	7.88
	DnCNN	75.83	71.12	73.34	73.49	71.73	77.06	75.16	70.31	71.31	72.75	69.96	71.46	71.65

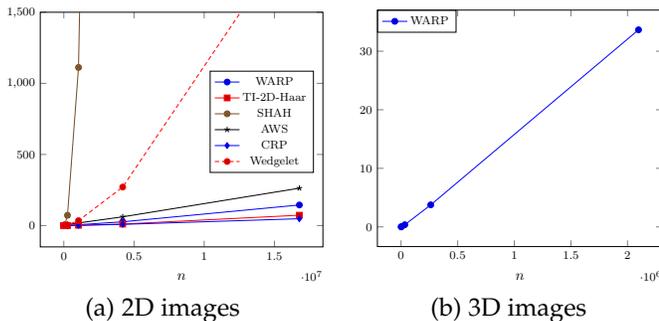


Fig. 4. Scalability of various methods using 2D and 3D images. Each line is the running time taken by the estimation step ( $y$ -axis) using the corresponding method versus the number of locations in the image ( $x$ -axis).

in Figure 2; we exclude 1D-Haar and RM as their reconstructions are highly inaccurate, and BPEFA as it scales poorly even at  $512 \times 512$  images. We can see that the empirical running time approximately follows a linear function of the number of locations. In fact, WARP takes only about 2 minutes for a large image of  $4096 \times 4096$  that contains 17 million pixels, and 5.3 seconds for an image of 1024 by 1024. Figure 4 suggests that Wedgelet and SHAH take quadratic time or even more, while TI-2D-Haar, AWS, and CRP takes linear time, but their performances are substantially inferior to that of WARP as shown in Figure 2. CRP seems to have a smaller slope than WARP, but it requires considerably longer tuning time than WARP according to the total running time with the tuning step included in the caption of Figure 2, at least based on its latest version of implementation to date.

It is worth noting that while many state-of-the-art methods designed for 2D images such as Wedgelet, TI-2D-Haar, and BPEFA require substantial modifications for a new dimensional setting, such as 3D images, the proposed WARP framework is directly applicable to  $m$ -dimensional data without modification, with the same linear scalability as suggested by Figure 4 (b).

### 3.3 Comparison with deep neural networks

In this section, we compare the proposed method WARP with convolutional neural networks. In particular, we apply WARP and the denoising convolutional neural networks (DnCNN) proposed in [37] to two popular benchmark datasets: the twelve widely used test images (Figure 7 in Supplementary Materials) and the BSD68 data [38] which contains 68 natural images from the Berkeley segmentation dataset. DnCNNs have been reported to achieve the state-of-the-art performance in various image processing tasks

[37]. We adopt a pre-trained model available in Matlab for DnCNN.

Table 1 reports the MSEs of WARP and DnCNN on the 12 widely used test images and BSD68 (averaged) at three noise levels when  $\sigma = 0.2, 0.4, 0.6$ . We can see that for light noise when  $\sigma = 0.2$ , WARP leads to smaller MSEs in five out of twelve images (i.e., Image 2, 3, 6, 8, 12) and gives comparable performance in other images. WARP gives uniformly smaller MSEs when  $\sigma = 0.4$  (intermediate noise) and constantly outperforms DnCNN by one order of magnitude when  $\sigma = 0.6$  (large noise), which is consistent with our observations in the ImageNet experiment. Besides the excellent performance of WARP, it is worth mentioning that unlike DnCNN which requires substantially more extensive pre-training and tuning, WARP does not require pre-training at all, and its small amount of tuning can be completely automated without user intervention. We do acknowledge that the performance of the pre-trained DnCNN might be improved with more extensive training.

## 4 ENHANCED ENERGY CONCENTRATION AND BEYOND 2D IMAGE RECONSTRUCTION

The excellent performance of WARP in image reconstruction suggests that the model is capable of identifying efficient representation of the underlying structure in a variety of real images as it is designed to achieve. This also suggests that extracting the underlying representation can potentially benefit a variety of other downstream processing tasks. In this section we first use a concept of “energy concentration” to examine how such efficiency is achieved and then discuss the potential applicability in other image processing tasks such as compression.

Energy (or information) concentration under wavelet transforms can be quantified by the number of wavelet coefficients needed to retain a given proportion of the sum of squares of the underlying function. An efficient wavelet representation will only need a small number of coefficients to capture most of the information contained in the function (as measured in terms of its sum of squares). Such a representation leads to high signal-to-noise ratios on a small number of coefficients that will facilitate all downstream processing tasks.

Next we compare energy concentration under WARP to that under classical 1D and 2D wavelet representations to quantify the improvement in energy concentration WARP achieves through adaptively identifying good permutations. To this end, we use the same ImageNet data as used in Section 3. For each image, we draw a sample from the posterior distribution of partition trees produced by WARP, and compute the number of coefficients required to attain a range

of energy levels (i.e., the total sum of squares) on a noisy observation at  $\sigma = 0.1$  and compare them to those required under standard 1D and 2D wavelet transforms. Figure 3 presents the numbers of wavelet coefficients required over the proportion of the sum of squares for three representative images.

Focusing on the proportion of the sum of squares from 0.85 to 0.95, we can see that the adaptive representation achieved by WARP requires substantially fewer numbers of wavelet coefficients (the red solid lines with scales on the right of each plot) to attain the same energy level than traditional 1D and 2D Haar DWT (the blue dashed lines with scales on the right of each plot). In Figure 3, we also plotted the percentage reduction in the number of coefficients (the black line with scales on the left of each plot) at each energy level. The largest coefficient saving of WARP is (80%, 70%, 70%) compared to 2D DWT, and this saving becomes (97%, 99%, 90%) when compared to 1D DWT. Enhanced energy concentration of WARP is observed in a wide range of test images in the database, and the extent of improvement in energy concentration varies according to the abundance of asymmetric structures present in the underlying image.

The improved energy concentration of WARP is expected to benefit a variety of downstream processing tasks beyond image denoising. For example, efficient image compression can be achieved using the posterior mode of the WARP model, which provides a sparse coding of the image. Coupling this idea with a pair of encoder and decoder, we introduce an algorithm for efficient image and video compression in a follow-up paper [39]. Interested readers may refer to that paper for additional numerical experiments involving a variety of datasets, including 2D ImageNet, 3D medical image, real-life YouTube videos, and surveillance videos, in which WARP-based compression substantially outperforms several state-of-the-art compression approaches.

## 5 APPLICATION TO RETINAL OPTICAL COHERENCE TOMOGRAPHY

We apply the proposed method to a dataset of optical coherence tomography (OCT) volumes. OCT provides a non-invasive imaging modality to visualize cross-sections of tissue layers at micrometer resolution, and thus is instrumental in various medical applications especially for the diagnosis and monitoring of patients with ocular diseases [40], [41], [42], [43]. The accurate interpretation of OCT images may require the involvement of both retina specialists and comprehensive ophthalmologists, and this task is compounded by heavily noised observations at a low signal-to-noise ratio due to sample-based speckle and detector noise [44], [45], [46]. Therefore, reconstruction of OCT images is necessary to improve both manual and automated OCT image analysis, and is increasingly important when OCT images are used to extract objective and quantitative assessment in ophthalmology which is touted as one advantage of OCT in clinical practice [42].

We use the OCT data available at [http://people.duke.edu/~sf59/Fang\\_TMI\\_2013.htm](http://people.duke.edu/~sf59/Fang_TMI_2013.htm), acquired by a Bioptigen SDOCT system (Durham, NC, USA) at an axial resolution of  $\sim 4.5 \mu\text{m}$ . We apply the methods of TI-2D-Haar, SHAH,

TABLE 2  
Mean PSNR for 18 foveal images reconstructed by BRFOE, K-SVD, PGPD, BM3D, MSBTD, SSR, and WARP. Results for the methods other than WARP are from [46].

BRFOE	K-SVD	PGPD	BM3D	MSBTD	SSR	WARP
25.32	27.03	27.01	27.04	27.08	28.10	28.18

AWS, CRP, Wedgelet, BPFA, and WARP, to two noisy slices (plotted as “Obs.” in Figure 5). We also have access to a registered and averaged image by averaging 40 repeatedly sampled scans [46], which is referred to as the “noiseless” reference image and is used to compare the quality of reconstructed images. From the results in Figure 5, we clearly see that WARP gives the best global qualitative metric using MSE and MAE among all methods in comparison.

Visual comparison provides a detailed assessment of reconstructed images on local features that might be clinically relevant. For the first observation in Figure 5, we can see WARP distinguishes all layers well (the boxed region in the noiseless image), especially compared to TI-2D-Haar and AWS whose reconstructions are blurred across layers. For the second observation, we observe a separation of the posterior cortical vitreous from the internal limiting membrane in the noiseless image, which shows the potential to progress to vitreomacular traction (VMT) [47]. This separation becomes less clear if using TI-2D-Haar (especially the left proportion), although TI-2D-Haar gives MSE and MAE that are closer to WARP than the other methods. For both observations, there is still substantial noise left in the denoised images by SHAH, and AWS gives a reconstruction exhibiting undesirable patches. This study confirms that WARP is capable to denoise images while keeping important features present in the image, due to its ability to adapt to the geometry of the underlying structures.

We further compare WARP with a study conducted in [46], which considers another six method: BRFOE [48], K-SVD [49], PGPD [50], BM3D [51], MSBTD [52], and SSR [46]. These six methods have been applied to 18 foveal images from 18 subjects, using four slices nearby the original observation at various stages of their implementation. Although WARP does not require nearby information and can even process a 3D volume if such data exist, we apply WARP to the observation that averages the original observation and the four nearby slices only to make a fair comparison. In Table 2, we adopt the mean of peak signal-to-noise ratio (PSNR) for all methods to align with [46], which is calculated as  $-10 \log_{10}(\text{MSE})$  (noting that we rescale all observations and noiseless gray-scale images by 255). We can see that WARP gives the largest mean of PSNR, thus achieves excellent performances compared to a wide range of existing methods in this application setting. We choose the two subjects considered in Figure 5, and plot the reconstructed images by WARP utilizing the nearby four slices in Figure 6. It suggests that WARP even has an enhanced display compared to the “noiseless” image, especially in the lower half of the image.

## 6 DISCUSSION

We have introduced the WARP framework that uses random recursive partitioning to induce a prior on the permutations of the index space, thereby achieving efficient inference on multi-dimensional functions by converting it into a Bayesian model choice problem involving one-dimensional competitive generative models. While our approach is Bayesian, one may consider other methods such as frequentist adaptive partitioning and shrinkage methods that incorporate the same idea. We do find satisfying the fully principled probabilistic inferential recipes that arise under our approach.

The proposed framework WARP can be applied along with a wider range of Bayes wavelet regression models, including those that allow heterogeneous noise levels. If the error  $\epsilon$  in Model (1) has general covariance matrix  $\Sigma_\epsilon$ , it often still makes sense to assume that the covariance of the error  $u$  in the wavelet domain, i.e.  $W\Sigma_\epsilon W'$ , is diagonal, due to the so-called whitening property of wavelet transforms discussed in [53]. In this case, let  $\sigma_j^2 = \text{Var}(u_{j,k})$  for each  $j$ . Then one may estimate  $\sigma_j^2$  using a robust estimator of the scale based on  $\{w_{j,k}, 0 \leq k \leq 2^j - 1\}$  given a tree, for example, using the median absolute deviation of  $\{w_{j,k}, 0 \leq k \leq 2^j - 1\}$  rescaled by 0.6745. Alternatively, one can adopt a hyperprior on location-based unknown variance  $\sigma_j^2 \sim \text{IG}(\nu + 1, \nu\sigma_0^2)$ , which is an inverse gamma prior with shape  $\nu + 1$  and scale  $\nu\sigma_0^2$  (thus the prior mean is  $\sigma_0^2$ ). The hyperparameters  $(\nu, \sigma_0^2)$  are either specified by users or estimated using data, for instance, we may estimate  $\sigma_0^2$  by the median estimate based on the finest scale wavelet coefficients [8].

Finally, while we introduce the WARP framework in the context of image denoising, we believe that the adaptive wavelet representation is applicable to a wide range of other tasks involving multi-dimensional signal processing.

## ACKNOWLEDGMENTS

We are very grateful to an AE and three reviewers for providing extremely helpful comments and suggestions. We also thank Daniel Bourgeois for his help in porting our C++ code to R. Meng Li's research is partly supported by NSF grant DMS-2015569 and an ORAU Ralph E. Powe Junior Faculty Enhancement Award. Li Ma's research is partly supported by NSF grants DMS-1749789 and DMS-2013930.

## SUPPLEMENTARY MATERIALS

Supplementary materials contain Proposition 1 and its proof; descriptions of WARP with local block shrinkage; details of the recursive message passing algorithm; proofs of all theorems; a sensitivity analysis for the proposed framework; plots of the 12 widely used test images used in Section 3.3; and comparison of WARP and selected methods using experiments of 3D image reconstruction.

## REFERENCES

[1] T. Alasil, P. A. Keane, J. F. Updike, L. Dustin, Y. Ouyang, A. C. Walsh, and S. R. Sadda, "Relationship between optical coherence tomography retinal parameters and visual acuity in diabetic macular edema," *Ophthalmology*, vol. 117, no. 12, pp. 2379–2386, 2010.

[2] I. I. Bussell, G. Wollstein, and J. S. Schuman, "Oct for glaucoma diagnosis, screening and detection of glaucoma progression," *British Journal of Ophthalmology*, pp. bjophthalmol-2013, 2013.

[3] W. C. Huang, A. V. Cideciyan, A. J. Roman, A. Sumaroka, R. Sheplock, S. B. Schwartz, E. M. Stone, and S. G. Jacobson, "Inner and outer retinal changes in retinal degenerations associated with abca4 mutations," *Investigative ophthalmology & visual science*, vol. 55, no. 3, pp. 1810–1822, 2014.

[4] J. K. Sun, M. M. Lin, J. Lammer, S. Prager, R. Sarangi, P. S. Silva, and L. P. Aiello, "Disorganization of the retinal inner layers as a predictor of visual acuity in eyes with center-involved diabetic macular edema," *JAMA ophthalmology*, vol. 132, no. 11, pp. 1309–1316, 2014.

[5] R. Kafieh, H. Rabbani, F. Hajizadeh, M. D. Abramoff, and M. Sonka, "Thickness mapping of eleven retinal layers segmented using the diffusion maps method in normal eyes," *Journal of ophthalmology*, vol. 2015, 2015.

[6] A. Oishi, P. P. Fang, S. Thiele, F. G. Holz, and T. U. Krohne, "Longitudinal change of outer nuclear layer after retinal pigment epithelial tear secondary to age-related macular degeneration," *Retina*, vol. 38, no. 7, pp. 1331–1337, 2018.

[7] D. L. Donoho and I. M. Johnstone, "Ideal spatial adaptation by wavelet shrinkage," *Biometrika*, vol. 81, no. 3, pp. 425–455, 1994.

[8] ———, "Adapting to unknown smoothness via wavelet shrinkage," *J. Am. Statist. Ass.*, vol. 90, no. 432, pp. 1200–1224, 1995.

[9] S. Mallat, *A Wavelet Tour of Signal Processing: The Sparse Way*, 3rd ed. Academic press, 2008.

[10] F. Abramovich, T. Sapatinas, and B. W. Silverman, "Wavelet thresholding via a Bayesian approach," *J. R. Statist. Soc. B*, vol. 60, no. 4, pp. 725–749, 1998.

[11] M. S. Crouse, R. D. Nowak, and R. G. Baraniuk, "Wavelet-based statistical signal processing using hidden Markov models," *IEEE Transactions on Signal Processing*, vol. 46, no. 4, pp. 886–902, 1998.

[12] M. Clyde and E. I. George, "Flexible empirical Bayes estimation for wavelets," *J. R. Statist. Soc. B*, vol. 62, no. 4, pp. 681–698, 2000.

[13] P. J. Brown, T. Fearn, and M. Vannucci, "Bayesian wavelet regression on curves with application to a spectroscopic calibration problem," *J. Am. Statist. Ass.*, vol. 96, no. 454, pp. 398–408, jun 2001.

[14] R. Willett and R. Nowak, "Fast multiresolution photon-limited image reconstruction," in *IEEE International Symposium on Biomedical Imaging: Nano to Macro, 2004*, pp. 1192–1195.

[15] J. S. Morris and R. J. Carroll, "Wavelet-based functional mixed models," *J. R. Statist. Soc. B*, vol. 68, no. 2, pp. 179–199, apr 2006.

[16] D. L. Donoho, "Wedgelets: Nearly minimax estimation of edges," *Ann. Statist.*, vol. 27, no. 3, pp. 859–897, 1999.

[17] L. Jacques, L. Duval, C. Chau, and G. Peyré, "A panorama on multiscale geometric representations, intertwining spatial, directional and frequency selectivity," *Signal Processing*, vol. 91, no. 12, pp. 2699–2730, 2011.

[18] S. T. Ali, J.-P. Antoine, and J.-P. Gazeau, *Coherent States, Wavelets and Their Generalizations*, 2nd ed. Springer, 2014.

[19] A. E. Raftery, "Bayesian model selection in social research," *Sociological Methodology*, vol. 25, pp. 111–163, 1995.

[20] J. A. Hoeting, D. Madigan, A. E. Raftery, and C. T. Volinsky, "Bayesian model averaging: a tutorial," *Statist. Sci.*, vol. 14, no. 4, pp. 382–417, 11 1999.

[21] W. H. Wong and L. Ma, "Optional Pólya tree and Bayesian inference," *Ann. Statist.*, vol. 38, no. 3, pp. 1433–1459, 2010.

[22] L. Ma, "Adaptive testing of conditional association through recursive mixture modeling," *J. Am. Statist. Ass.*, vol. 108, no. 504, pp. 1493–1505, 2013.

[23] H. A. Chipman, E. I. George, and R. E. McCulloch, "Bayesian CART model search," *J. Am. Statist. Ass.*, vol. 93, no. 443, pp. 935–948, sep 1998.

[24] D. G. T. Denison, B. K. Mallick, and A. F. M. Smith, "A Bayesian CART algorithm," *Biometrika*, vol. 85, no. 2, pp. 363–377, 1998.

[25] I. M. Johnstone and B. W. Silverman, "Empirical Bayes selection of wavelet thresholds," *Ann. Statist.*, vol. 33, no. 4, pp. 1700–1752, 2005.

[26] H. A. Chipman, E. D. Kolaczyk, and R. E. McCulloch, "Adaptive Bayesian wavelet shrinkage," *J. Am. Statist. Ass.*, vol. 92, no. 440, pp. 1413–1421, dec 1997.

[27] L. Ma and J. Soriano, "Efficient functional ANOVA through wavelet-domain Markov groves," *J. Am. Statist. Ass.*, 2017, to appear, DOI:10.1080/01621459.2017.1286241.

- [28] T. T. Cai, "Adaptive wavelet estimation: A block thresholding and oracle inequality approach," *Ann. Statist.*, vol. 27, no. 3, pp. 898–924, jun 1999.
- [29] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *2009 IEEE Conference on Computer Vision and Pattern Recognition*, 2009, pp. 248–255.
- [30] P. Fryzlewicz and C. Timmermans, "SHAH: SHape-Adaptive Haar wavelets for image processing," *Journal of Computational and Graphical Statistics*, vol. 25, no. 3, pp. 879–898, 2016.
- [31] J. Polzehl and V. G. Spokoiny, "Adaptive weights smoothing with applications to image restoration," *J. R. Statist. Soc. B*, vol. 62, no. 2, pp. 335–354, 2000.
- [32] M. Li and S. Ghosal, "Bayesian multiscale smoothing of Gaussian noised images," *Bayesian Analysis*, vol. 9, no. 3, pp. 733–758, 2014.
- [33] R. Castro, R. Willett, and R. Nowak, "Coarse-to-fine manifold learning," in *IEEE International Conference on Acoustics, Speech, and Signal Processing, 2004 (ICASSP'04)*, 2004.
- [34] M. Zhou, H. Chen, J. Paisley, L. Ren, L. Li, Z. Xing, D. Dunson, G. Sapiro, and L. Carin, "Nonparametric Bayesian dictionary learning for analysis of noisy and incomplete images," *IEEE Transactions on Image Processing*, vol. 21, no. 1, pp. 130–144, 2012.
- [35] R. Coifman and D. Donoho, "Translation-invariant de-noising," in *Wavelets and Statistics*, ser. Lecture Notes in Statistics, A. Antoniadis and G. Oppenheim, Eds. Springer New York, 1995, vol. 103, pp. 125–150.
- [36] M. Li and S. Ghosal, "Fast translation invariant multiscale image denoising," *IEEE Transactions on Image Processing*, vol. 24, no. 12, pp. 4876–4887, 2015.
- [37] K. Zhang, W. Zuo, Y. Chen, D. Meng, and L. Zhang, "Beyond a Gaussian denoiser: Residual learning of deep CNN for image denoising," *IEEE Transactions on Image Processing*, vol. 26, no. 7, pp. 3142–3155, 2017.
- [38] D. Martin, C. Fowlkes, D. Tal, and J. Malik, "A database of human segmented natural images and its application to evaluating segmentation algorithms and measuring ecological statistics," in *Proc. 8th Int'l Conf. Computer Vision*, vol. 2, July 2001, pp. 416–423.
- [39] R. Liu, M. Li, and L. Ma, "CARP: Compression Through Adaptive Recursive Partitioning for Multi-Dimensional Images," in *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, jun 2020, pp. 14 294–14 302.
- [40] D. Huang, E. A. Swanson, C. P. Lin, J. S. Schuman, W. G. Stinson, W. Chang, M. R. Hee, T. Flotte, K. Gregory, C. A. Puliafito *et al.*, "Optical coherence tomography," *Science*, vol. 254, no. 5035, pp. 1178–1181, 1991.
- [41] D. S. Grewal and A. P. Tanna, "Diagnosis of glaucoma and detection of glaucoma progression using spectral domain optical coherence tomography," *Current opinion in ophthalmology*, vol. 24, no. 2, pp. 150–161, 2013.
- [42] G. Virgili, F. Menchini, G. Casazza, R. Hogg, R. R. Das, X. Wang, and M. Michelessi, "Optical coherence tomography (oct) for detection of macular oedema in patients with diabetic retinopathy," *The Cochrane database of systematic reviews*, vol. 1, p. CD008081, 2015.
- [43] N. Cuenca, I. Ortuño-Lizarán, and I. Pinilla, "Cellular characterization of oct and outer retinal bands using specific immunohistochemistry markers and clinical implications," *Ophthalmology*, vol. 125, no. 3, pp. 407–422, 2018.
- [44] P. A. Keane, P. J. Patel, S. Liakopoulos, F. M. Heussen, S. R. Sadda, and A. Tufail, "Evaluation of age-related macular degeneration with optical coherence tomography," *Survey of ophthalmology*, vol. 57, no. 5, pp. 389–414, 2012.
- [45] F. Shi, X. Chen, H. Zhao, W. Zhu, D. Xiang, E. Gao, M. Sonka, and H. Chen, "Automated 3-d retinal layer segmentation of macular optical coherence tomography images with serous pigment epithelial detachments," *IEEE Trans. Med. Imaging*, vol. 34, no. 2, pp. 441–452, 2015.
- [46] L. Fang, S. Li, D. Cunefare, and S. Farsiu, "Segmentation Based Sparse Reconstruction of Optical Coherence Tomography Images," *IEEE Transactions on Medical Imaging*, vol. 36, no. 2, pp. 407–421, feb 2017.
- [47] J. S. Duker, P. K. Kaiser, S. Binder, M. D. De Smet, A. Gaudric, E. Reichel, S. R. Sadda, J. Sebag, R. F. Spaide, and P. Stalmans, "The international vitreomacular traction study group classification of vitreomacular adhesion, traction, and macular hole," *Ophthalmology*, vol. 120, no. 12, pp. 2611–2619, 2013.
- [48] Y. Weiss and W. T. Freeman, "What makes a good model of natural images?" in *2007 IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, 2007, pp. 1–8.
- [49] M. Elad and M. Aharon, "Image denoising via sparse and redundant representations over learned dictionaries," *IEEE Transactions on Image processing*, vol. 15, no. 12, pp. 3736–3745, 2006.
- [50] J. Xu, L. Zhang, W. Zuo, D. Zhang, and X. Feng, "Patch group based nonlocal self-similarity prior learning for image denoising," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 244–252.
- [51] K. Dabov, A. Foi, V. Katkovnik, and K. Egiazarian, "Image denoising by sparse 3-d transform-domain collaborative filtering," *IEEE Transactions on image processing*, vol. 16, no. 8, pp. 2080–2095, 2007.
- [52] L. Fang, S. Li, Q. Nie, J. A. Izatt, C. A. Toth, and S. Farsiu, "Sparsity based denoising of spectral domain optical coherence tomography images," *Biomedical optics express*, vol. 3, no. 5, pp. 927–942, 2012.
- [53] I. M. Johnstone and B. W. Silverman, "Wavelet threshold estimators for data with correlated noise," *J. R. Statist. Soc. B*, vol. 59, no. 2, pp. 319–351, 1997.
- [54] A. V. Aho and N. J. Sloane, "Some doubly exponential sequences," *Fibonacci Quart*, vol. 11, no. 4, pp. 429–437, 1973.
- [55] P. Mukherjee and P. Qiu, "3-d image denoising by local smoothing and nonparametric regression," *Technometrics*, vol. 53, no. 2, pp. 196–208, 2011.
- [56] P. Perona and J. Malik, "Scale-space and edge detection using anisotropic diffusion," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 12, no. 7, pp. 629–639, 1990.
- [57] L. Rudin, S. Osher, and E. Fatemi, "Nonlinear total variation based noise removal algorithms," *Physica D: Nonlinear Phenomena*, vol. 60, no. 1, pp. 259–268, 1992.
- [58] P. Coupé, P. Yger, S. Prima, P. Hellier, C. Kervrann, and C. Barillot, "An optimized blockwise nonlocal means denoising filter for 3-d magnetic resonance images," *IEEE Transactions on Medical Imaging*, vol. 27, no. 4, pp. 425–441, 2008.

Fig. 5. Two retinal OCT datasets (titled “Obs.”) and reconstructed images using TI-2D-Haar, SHAH, AWS, CRP, Wedgelet, BPFA, and WARP. The two metrics following each method are the MSE ( $\times 10^{-4}$ ) and MAE ( $\times 10^{-2}$ ) respectively. The “noiseless” reference is an registered and averaged image.

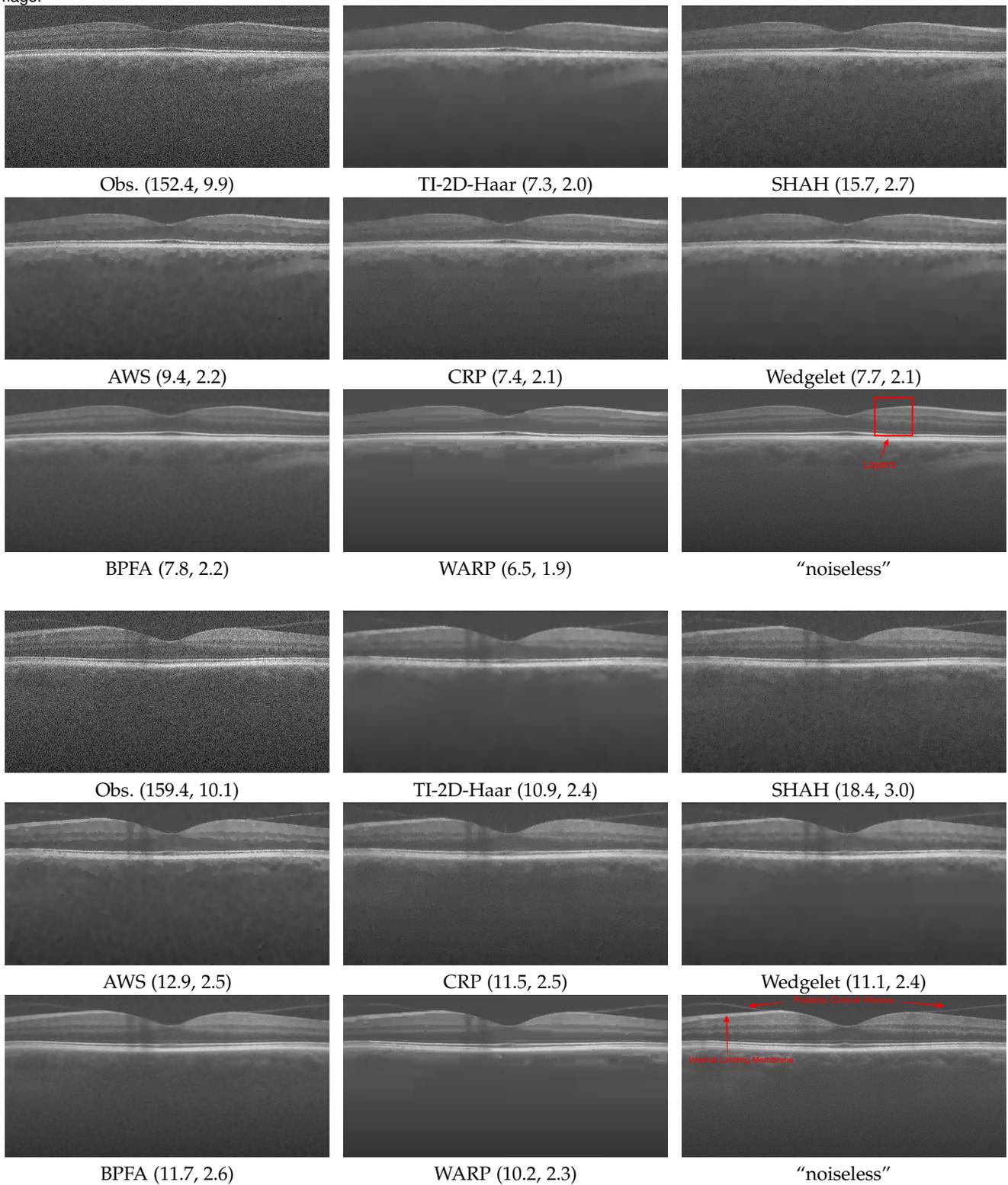
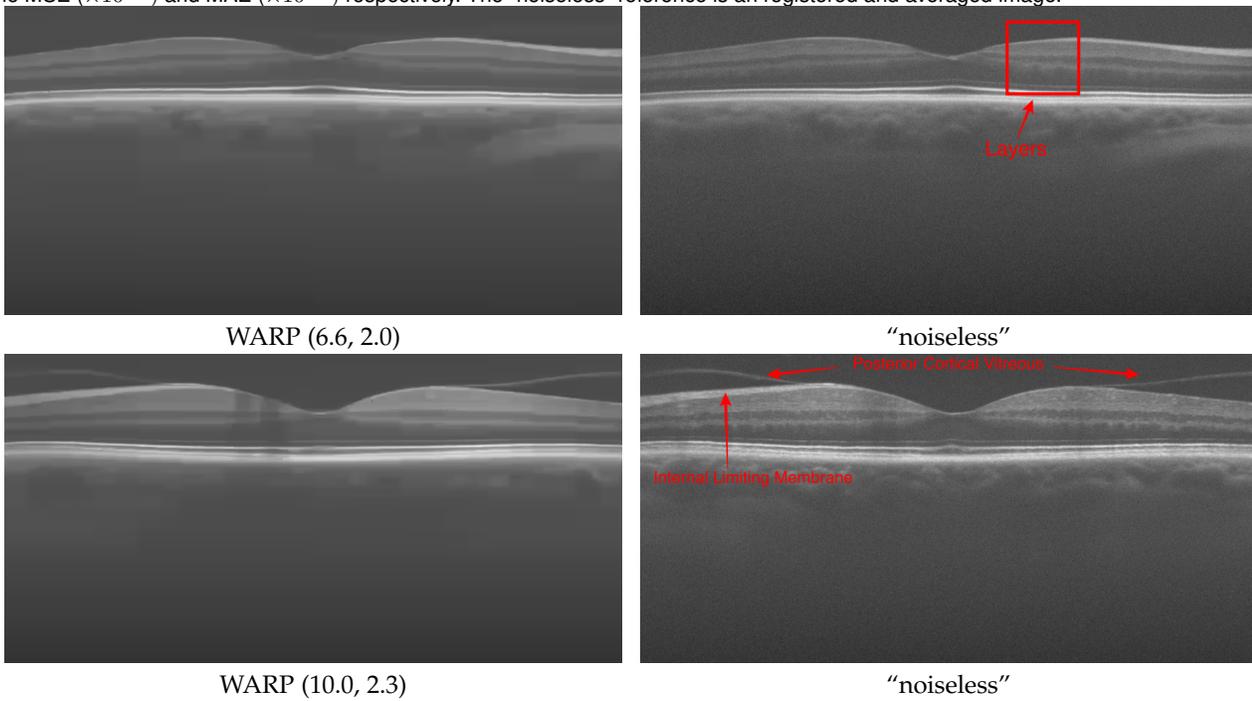


Fig. 6. Reconstructed images using WARP based on the noisy observation and its four nearby slices. The two metrics following each method are the MSE ( $\times 10^{-4}$ ) and MAE ( $\times 10^{-2}$ ) respectively. The “noiseless” reference is an registered and averaged image.



# Supplementary Materials for “Learning Asymmetric and Local Features in Multi-Dimensional Data through Wavelets with Recursive Partitioning”

Supplementary materials contain (A) Proposition 1 and its proof, (B) descriptions of WARP with local block shrinkage, (C) details of the recursive message passing algorithm, (D) proofs of all theorems, (E) a sensitivity analysis for the proposed framework, (F) plots of the 12 widely used test images used in Section 3.3, and (G) comparison of WARP and selected methods using experiments of 3D image reconstruction.

## A CARDINALITY OF THE SPACE OF RDPs

**Proposition 1.** *The log cardinality of the tree space induced by RDPs is  $O(n)$  when  $m = 2$ .*

*Proof of Proposition 1.* Let  $c(a, b)$  be the cardinality of the tree space induced by RDPs for an  $2^a$  by  $2^b$  image. We can obtain the following recursive formula

$$c(a, b) = \begin{cases} c^2(a-1, b) + c^2(a, b-1), & \text{if } a, b \geq 1 \\ 1 & \text{if } a = 0 \text{ or } b = 0. \end{cases}$$

We assert that there exist two constants  $(k_1, k_2)$  such that  $k_2 \geq k_1 > 1$  and

$$c(a, b) \in \left[ \frac{1}{2} k_1^{2^{a+b}}, \frac{1}{2} k_2^{2^{a+b}} \right],$$

for any  $a \geq 1$  and  $b \geq 1$ .

First consider  $a = 1$  and  $b \geq 1$ . We have  $c(1, b) = c^2(1, b-1) + 1$  when  $b \geq 1$  and  $c(1, 0) = 1$  when  $b = 0$ . The quantity  $c(1, b)$  is actually the number of “strongly” binary trees of height  $\leq b$ , which possesses an analytical form

$$c(1, b) = \lfloor k^{2^b} \rfloor,$$

according to [54], where

$$k = \exp \left\{ \sum_{j=0}^{\infty} 2^{-j-1} \log(1 + c^{-2}(1, j)) \right\} \approx 1.503.$$

Letting  $k_1 = \sqrt{k}$  and  $k_2 = k$  and noting  $k^{2^b} \geq 2$  for all  $b \geq 1$ , we obtain that

$$\frac{1}{2} k_1^{2^{1+b}} = \frac{1}{2} k^{2^b} \leq k^{2^b} - 1 \leq \lfloor k^{2^b} \rfloor \leq k^{2^b} \leq \frac{1}{2} k^{2^{1+b}},$$

for all  $b \geq 1$ . Therefore, the assertion holds for all  $a = 1$  and  $b \geq 1$ . Since  $c(a, b) = c(b, a)$ , the assertion also holds for all  $a \geq 1$  and  $b = 1$ .

For any  $a \geq 1$  and  $b \geq 1$ , it is easy to verify that if the assertion holds for  $(a, b-1)$  and  $(a-1, b)$ , then it holds for  $(a, b)$ . We then complete the proof by induction.  $\square$

## B WARP WITH LOCAL BLOCK SHRINKAGE

Traditional wavelet analysis is done by fixing the maximum depth of the wavelet tree at  $J$ . That is, one partitions the index space all the way down to the finest level of “atomic” blocks. In most practical problems, once the blocks are small enough, the function value within the block becomes essentially constant with respect to the noise level, and so further division within such homogeneous blocks will be wasteful and will reduce statistical efficiency. For example, in Figure 1 the partition in the upper left block (Level 3) along with its descendants is not necessary. Thus it is often desirable to incorporate adaptivity in the depth of the wavelet tree and allow it to be terminated earlier than reaching level  $J$ . In practice the optimal maximum depth varies across  $\Omega$ . For example, some parts of an image may contain many interesting details, while the rest do not—e.g., an image of a painting hung on a gray wall. A high resolution will be needed to capture the details in the painting, but would be unnecessary and introduce additional variability in the estimation for the wall.

This consideration is closely related to the idea of adaptive block shrinkage [28] in the frequentist wavelet regression analysis. Once there is little evidence for any interesting structure within a subset of the index space, then the function value within that subset can be shrunk to a constant. That is, the wavelet tree is “pruned” there. Next we show that such pruning can be achieved in a hierarchical modeling manner, and the resulting Bayesian wavelet regression model is again compatible with our WARP framework.

To achieve such pruning, we introduce another set of latent variables  $\mathcal{R} = \{R_{j,k} : j = 0, 1, \dots, J-1, k = 0, 1, \dots, 2^j - 1\}$ , where  $R_{j,k} = 1$  indicates that the tree is pruned at node  $(j, k)$ . Next we describe a generative prior on  $\mathcal{R}$  that will blend well with the WARP framework. To start, let  $R_{0,0} \stackrel{\text{ind}}{\sim} \text{Bern}(\eta_{0,0})$  and for all  $j \geq 1$ , and

$$R_{j,k} | R_{j-1, \lfloor k/2 \rfloor} \stackrel{\text{ind}}{\sim} \begin{cases} \text{Bern}(\eta_{j,k}) & \text{if } R_{j-1, \lfloor k/2 \rfloor} = 0 \\ \text{Bern}(1) & \text{if } R_{j-1, \lfloor k/2 \rfloor} = 1. \end{cases}$$

That is, if a node’s parent has been pruned, then its children are also pruned by construction. From now on, we shall refer to this prior model on  $\mathcal{R}$  as an *optional pruning* (OP) model [22], which is specified by a set of *pruning probabilities*  $\eta_{j,k} \in [0, 1]$ . We write  $\mathcal{R} \sim \text{OP}(\eta)$ .

Given  $\mathcal{R}$ , we can modify our prior on  $\mathcal{S}$  to reflect the effect of pruning. For example, instead of an independent prior on  $\mathcal{S}$ , we can now generate them as follows

$$S_{j,k} | \mathcal{R} \stackrel{\text{ind}}{\sim} \begin{cases} \text{Bern}(\rho_{j,k}) & \text{if } R_{j,k} = 0 \\ \text{Bern}(0) & \text{if } R_{j,k} = 1. \end{cases}$$

That is, if the node has not been pruned, then we generate  $S_{j,k}$  from the independent Bernoulli as in the standard spike-and-slab setup, but if the node has been pruned, then by construction, we must have  $S_{j,k} = 0$  due to pruning.

It is often reasonable to specify the prior shrinkage and pruning probabilities as functions of the level in the RDP. That is,  $\rho_{j,k} = \rho_j$  and  $\eta_{j,k} = \eta_j$  for all  $k$ . In the node-specific notation,  $\rho(A) = \rho_j$  and  $\eta(A) = \eta_j$  for all  $j$ th node  $A \in \mathcal{A}$ . In this case, one can show that this joint model on  $(\mathcal{S}, \mathcal{R})$  is equivalent to a Markov tree model with three states defined in terms of the combinations of  $(S_{j,k}, R_{j,k}) = (1, 0), (0, 0)$ ,

or (0,1), and with the corresponding transition matrix for  $S_{j,k}$  given by

$$\rho_j = \begin{bmatrix} \rho_j(1 - \eta_j) & (1 - \rho_j)(1 - \eta_j) & \eta_j \\ \rho_j(1 - \eta_j) & (1 - \rho_j)(1 - \eta_j) & \eta_j \\ 0 & 0 & 1 \end{bmatrix}.$$

This allows us to derive the posterior from Theorem 2, and carry out inference accordingly. Specifically, for each  $A \in \mathcal{A}$ , let  $p_0(A)$  be the marginal likelihood contributed from the wavelet coefficients in  $A$  and its descendants if  $A$  is pruned, i.e.,

$$p_0(A) = \frac{1}{(\sqrt{2\pi\sigma^2})^{|A|-1}} \exp \left\{ -\frac{\sum_{x \in A} (y(x) - \bar{y}(A))^2}{2\sigma^2} \right\},$$

where  $\bar{y}(A) = \sum_{x \in A} y(x)/|A|$ . If  $A \in \mathcal{T}$ , the following maps are directly available from Theorem 2:

- The marginal likelihood contribution from the data within node  $A$  if  $A$  is divided in dimension  $d$ :

$$M_d(A) = \rho(A)M_d^{(1)}(A) + (1 - \rho(A))M_d^{(0)}(A);$$

- The posterior spike probability  $\tilde{\rho}_d$  of  $A$  if  $A$  is divided in dimension  $d$ :

$$\tilde{\rho}_d(A) = \rho(A)M_d^{(1)}(A)/M_d(A);$$

- The marginal likelihood from data on  $A$  and its descendants:  $\Psi(A) = (1 - \eta(A)) \sum_{d \in \mathcal{D}(A)} \lambda_d(A)M_d(A)\Psi(A_l^{(d)})\Psi(A_r^{(d)}) + \eta(A)p_0(A)$  if  $A$  is non-atomic;  $\Psi(A) = 1$  if  $A$  is atomic.
- The posterior probability of pruning  $A$ :

$$\tilde{\eta}(A) = \eta(A)p_0(A)/\Psi(A);$$

- The posterior probability for  $A$  to be divided in dimension  $d$  given  $A$  is not pruned:

$$\tilde{\lambda}_d(A) = (1 - \eta(A)) \frac{\lambda_d(A)M_d(A)\Psi(A_l^{(d)})\Psi(A_r^{(d)})}{\Psi(A) - \eta(A)p_0(A)}.$$

## C RECURSIVE MESSAGE PASSING ALGORITHM

For the Haar basis, the posterior mean  $E(\mathbf{f}|\mathbf{y})$  can be evaluated analytically through recursive message passing without any Monte Carlo sampling for Bayesian wavelet regression models that adopt the spike-and-slab setup along with optional pruning of the wavelet tree, which contains the models without optional pruning as special cases with zero pruning probabilities. We describe the strategy next and will use it to compute  $E(\mathbf{f}|\mathbf{y})$  in our numerical examples.

For each  $A \in \mathcal{A}$ , let  $c(A)$  be the scale (father wavelet) coefficient on  $A$  if  $A \in \mathcal{T}$ , and let  $\varphi(A) = E(c(A)\mathbf{1}_{\{A \in \mathcal{T}\}}|\mathbf{y})$ . Note that  $E(\mathbf{f}|\mathbf{y})$  is given by  $\varphi(A)$  for all atomic  $A$ . To compute the mapping  $\varphi$ , we introduce two auxiliary mappings  $\psi_0(A) = P(A \in \mathcal{T}, R(A) = 0|\mathbf{y})$  and  $\varphi_0(A) = E(c(A)\mathbf{1}_{\{A \in \mathcal{T}, R(A)=0\}}|\mathbf{y})$ . Let  $\bar{A}^{(d)}$  denote the parent of  $A$  in  $\mathcal{T}$  if  $A$  is a child node after dividing its parent in the  $d$ th dimension, and let  $\mathcal{P}(A) \subset \{1, 2, \dots, m\}$  be the collection of dimensions of  $A$  that do not have full support  $[0, 2^{j_i} - 1]$ , i.e., those that have been partitioned at least once in previous levels. Using the tri-variate mapping  $(\phi_0, \varphi_0, \varphi) : \mathcal{A} \rightarrow \mathbb{R}^3$ .

**Theorem 3.** To initiate the recursion, for  $A = \Omega$ , we let  $\psi_0(A) = 1 - \tilde{\eta}(A)$ ,  $\varphi_0(A) = (1 - \tilde{\eta}(A))|A|/\sqrt{n}$ , and  $\varphi(A) = |A|/\sqrt{n}$ . Suppose we have evaluated these mappings up to level  $j - 1$ , for level  $j = 1, \dots, J$ , we have

$$\begin{aligned} \psi_0(A) &= \sum_{d \in \mathcal{P}(A)} \psi_0(\bar{A}^{(d)})\tilde{\lambda}_d(\bar{A}^{(d)})(1 - \tilde{\eta}(A)); \\ \varphi_0(A) &= (1 - \tilde{\eta}(A)) \cdot \sum_{d \in \mathcal{P}(A)} \frac{\tilde{\lambda}_d(\bar{A}^{(d)})}{\sqrt{2}} \left[ \varphi_0(\bar{A}^{(d)}) - \tilde{\rho}_d(\bar{A}^{(d)})\mu_1(w_d(\bar{A}^{(d)}))\psi_0(\bar{A}^{(d)})(-1)^{\mathbf{1}(A \text{ is the left child of } \bar{A}^{(d)})} \right]; \\ \varphi(A) &= \frac{\varphi_0(A)}{1 - \tilde{\eta}(A)} + \frac{1}{\sqrt{2}} \sum_{d \in \mathcal{P}(A)} \{\varphi(\bar{A}^{(d)}) - \varphi_0(\bar{A}^{(d)})\}\lambda_d(\bar{A}^{(d)}). \end{aligned}$$

**Remark:** Note that this recursion is top-down (from low to high resolutions), whereas that for computing  $\Phi$  is bottom-up (from high to low resolutions). The two-directional recursion shares the spirit of the forward-backward algorithm for HMMs.

Once we have computed the mapping  $(\varphi_0, \psi_0, \varphi) : \mathcal{A} \rightarrow \mathbb{R}^3$ , the posterior mean  $E(\mathbf{f}|\mathbf{y})$  is then given by  $\varphi$  applied on the atomic nodes. Note that this theorem applies to the special case with no pruning as well.

## D PROOFS OF THEOREMS

*Proof of Theorem 1.* Because Theorem 1 can be considered a special case with a single latent state, its proof follows immediately from the latter theorem, which we prove below.  $\square$

*Proof of Theorem 2.* First we verify that the mapping  $\Phi_s(A)$  is the marginal likelihood contributed from data with locations in  $A$ , given that  $A \in \mathcal{T}$  and that the latent state variable associated with the parent of  $A$  in  $\mathcal{T}$  is  $s$ . We show this by induction. First note that if  $A$  is atomic, then

$$\Phi_s(A) = P(\mathbf{y}(A) | A \in \mathcal{T}, S(A_p) = s) = 1$$

by design as there are no wavelet coefficients associated with atomic nodes. Now, suppose we have shown that  $\Phi_s(A) = P(\mathbf{y}(A) | A \in \mathcal{T}, S(A_p) = s)$  for all  $A$  with level higher than  $j$ . Then if  $A$  is of level  $j$ , it follows that

$$\begin{aligned} &P(\mathbf{y}(A) | A \in \mathcal{T}, S(A_p) = s) \\ &= \sum_{s'} \sum_d P(\mathbf{y}(A) | A \in \mathcal{T}, S(A) = s', S(A_p) = s, D(A) = d) \\ &\quad \times P(S(A) = s' | A \in \mathcal{T}, S(A_p) = s) \\ &\quad \times P(D(A) = d | A \in \mathcal{T}, S(A_p) = s) \\ &= \sum_{s'} \rho_j(s, s') \sum_{d \in \mathcal{D}(A)} \lambda_d M_d^{(s')} (A) \Psi_{s'}(A_l^{(d)}) \Psi_{s'}(A_r^{(d)}), \end{aligned}$$

which leads to the definition of  $\Phi_s(A)$  in Theorem 2.

Next let us derive the joint marginal posterior of  $(\mathcal{T}, S)$ . Note that

$$\begin{aligned} &P(S_{j,k} = s' | S_{j-1, \lfloor k/2 \rfloor} = s, \mathcal{T}^{(j)}, \mathbf{y}) \\ &= \frac{P(S_{j,k} = s', S_{j-1, \lfloor k/2 \rfloor} = s, \mathbf{y}(A) | \mathcal{T}^{(j)})}{P(S_{j-1, \lfloor k/2 \rfloor} = s, \mathbf{y}(A) | \mathcal{T}^{(j)})}. \end{aligned}$$

Now we have

$$\begin{aligned} & \mathbb{P}(S_{j,k} = s', D_{j,k} = d, \mathbf{y}(A) | \mathcal{T}^{(j)}, S_{j-1, \lfloor k/2 \rfloor} = s) \\ &= \rho_j(s, s') \lambda_d(A) M_d^{(s')} (A) \Phi_{s'}(A_l^{(d)}) \Phi_{s'}(A_r^{(d)}), \end{aligned}$$

which leads to

$$\begin{aligned} & \mathbb{P}(S_{j,k} = s', \mathbf{y}(A) | \mathcal{T}^{(j)}, S_{j-1, \lfloor k/2 \rfloor} = s) \\ &= \rho_j(s, s') \sum_d \lambda_d(A) M_d^{(s')} (A) \Phi_{s'}(A_l^{(d)}) \Phi_{s'}(A_r^{(d)}) \end{aligned}$$

and furthermore,

$$\begin{aligned} & \mathbb{P}(S_{j,k} = s' | S_{j-1, \lfloor k/2 \rfloor} = s, \mathcal{T}^{(j)}, \mathbf{y}) \\ &= \frac{\rho_j(s, s') \sum_d \lambda_d(A) M_d^{(s')} (A) \Phi_{s'}(A_l^{(d)}) \Phi_{s'}(A_r^{(d)})}{\sum_{s''} \rho_j(s, s'') \sum_d \lambda_d(A) M_d^{(s'')} (A) \Phi_{s''}(A_l^{(d)}) \Phi_{s''}(A_r^{(d)})}, \end{aligned}$$

where the denominator is just  $\Phi_s(A)$ .

Finally,

$$\begin{aligned} & \mathbb{P}(D_{j,k} = d | S_{j,k} = s', \mathcal{T}^{(j)}, \mathbf{y}) \\ &= \frac{\mathbb{P}(D_{j,k} = d, \mathbf{y}(A) | S_{j,k} = s', \mathcal{T}^{(j)})}{\mathbb{P}(\mathbf{y}(A) | S_{j,k} = s', \mathcal{T}^{(j)})} \\ &= \frac{\lambda_d(A) M_d^{(s')} (A) \Phi_{s'}(A_l^{(d)}) \Phi_{s'}(A_r^{(d)})}{\sum_{d'} \lambda_{d'}(A) M_{d'}^{(s')} (A) \Phi_{s'}(A_l^{(d')}) \Phi_{s'}(A_r^{(d')})}. \end{aligned}$$

This completes the proof.  $\square$

*Proof of Theorem 3.* We first obtain the recursive recipe for computing the maps  $(\psi_0, \varphi_0)$  following Theorem 1:

$$\begin{aligned} & \psi_0(A) \\ &= \sum_{d \in \mathcal{P}(A)} \mathbb{P}(\bar{A}^{(d)} \in \mathcal{T}, R(\bar{A}^{(d)}) = 0 | \mathbf{y}) \tilde{\lambda}_d(\bar{A}^{(d)}) (1 - \tilde{\eta}(A)) \\ &= \sum_{d \in \mathcal{P}(A)} \psi_0(\bar{A}^{(d)}) \tilde{\lambda}_d(\bar{A}^{(d)}) (1 - \tilde{\eta}(A)), \end{aligned}$$

and

$$\begin{aligned} & \varphi_0(A) = \mathbb{E}(c(A) \mathbf{1}_{\{A \in \mathcal{T}, R(A)=0\}} | \mathbf{y}) \\ &= \sum_{d \in \mathcal{P}(A)} \mathbb{E}(c(A) \mathbf{1}_{\{\bar{A}^{(d)} \in \mathcal{T}, D(\bar{A}^{(d)})=d, R(\bar{A}^{(d)})=0\}} | \mathbf{y}) \\ &= \sum_{d \in \mathcal{P}(A)} \mathbb{E}(c(A) | \bar{A}^{(d)} \in \mathcal{T}, D(\bar{A}^{(d)}) = d, R(\bar{A}^{(d)}) = 0, \mathbf{y}) \\ &\quad \times \mathbb{P}(\bar{A}^{(d)} \in \mathcal{T}, D(\bar{A}^{(d)}) = d, R(\bar{A}^{(d)}) = 0 | \mathbf{y}) \\ &= (1 - \tilde{\eta}(A)) \sum_{d \in \mathcal{P}(A)} \frac{\tilde{\lambda}_d(\bar{A}^{(d)})}{\sqrt{2}} \left[ \varphi_0(\bar{A}^{(d)}) - \right. \\ &\quad \left. \tilde{\rho}_d(\bar{A}^{(d)}) \mu_1(w_d(\bar{A}^{(d)})) \psi_0(\bar{A}^{(d)}) \cdot (-1)^{\mathbf{1}(A \text{ is the left child of } \bar{A}^{(d)})} \right]. \end{aligned} \quad (2)$$

We next derive the recursive formula for  $\varphi(A)$ . Let  $\varphi_1(A) = \mathbb{E}(c(A) \mathbf{1}_{\{A \in \mathcal{T}, R(A)=1\}} | \mathbf{y})$ , then we have  $\varphi(A) = \mathbb{E}(c(A) \mathbf{1}_{\{A \in \mathcal{T}\}} | \mathbf{y}) = \varphi_0(A) + \varphi_1(A)$ . Note that

$$\varphi(A) = \sum_{d \in \mathcal{P}(A)} \mathbb{E}(c(A) \mathbf{1}_{\{\bar{A}^{(d)} \in \mathcal{T}, D(\bar{A}^{(d)})=d\}} | \mathbf{y}), \quad (3)$$

and for each  $d \in \mathcal{P}(A)$ , we have

$$\begin{aligned} & \mathbb{E}(c(A) \mathbf{1}_{\{\bar{A}^{(d)} \in \mathcal{T}, D(\bar{A}^{(d)})=d\}} | \mathbf{y}) \\ &= \sum_{r=0,1} \mathbb{E}(c(A) \mathbf{1}_{\{\bar{A}^{(d)} \in \mathcal{T}, D(\bar{A}^{(d)})=d, R(\bar{A}^{(d)})=r\}} | \mathbf{y}) \\ &= \sum_{r=0,1} \mathbb{E}(c(A) | \bar{A}^{(d)} \in \mathcal{T}, D(\bar{A}^{(d)}) = d, R(\bar{A}^{(d)}) = r, \mathbf{y}) \\ &\quad \times \mathbb{P}(\bar{A}^{(d)} \in \mathcal{T}, D(\bar{A}^{(d)}) = d, R(\bar{A}^{(d)}) = r | \mathbf{y}). \end{aligned} \quad (4)$$

For the second term in (4), we have

$$\begin{aligned} & \mathbb{P}(\bar{A}^{(d)} \in \mathcal{T}, D(\bar{A}^{(d)}) = d, R(\bar{A}^{(d)}) = r | \mathbf{y}) \\ &= \mathbb{P}(D(\bar{A}^{(d)}) = d | \bar{A}^{(d)} \in \mathcal{T}, R(\bar{A}^{(d)}) = r, \mathbf{y}) \\ &\quad \times \mathbb{P}(\bar{A}^{(d)} \in \mathcal{T}, R(\bar{A}^{(d)}) = r | \mathbf{y}) \\ &= \tilde{\lambda}_d(\bar{A}^{(d)})^{1-r} \lambda_d(\bar{A}^{(d)})^r \psi_r(\bar{A}^{(d)}). \end{aligned}$$

For the first term in (4), it is easy to check that

$$\begin{aligned} & \mathbb{E}(c(A) | \bar{A}^{(d)} \in \mathcal{T}, D(\bar{A}^{(d)}) = d, R(\bar{A}^{(d)}) = r, \mathbf{y}) \\ &= \begin{cases} \frac{1}{\sqrt{2}} \left[ \frac{\varphi_0(\bar{A}^{(d)})}{\psi_0(\bar{A}^{(d)})} - \tilde{\rho}_d(\bar{A}^{(d)}) \mu_1(w_d(\bar{A}^{(d)})) \right. \\ \quad \left. \times (-1)^{\mathbf{1}(A \text{ is the left child of } \bar{A}^{(d)})} \right] & \text{if } r = 0 \\ \frac{1}{\sqrt{2}} \varphi_1(\bar{A}^{(d)}) / \psi_1(\bar{A}^{(d)}) & \text{if } r = 1, \end{cases} \end{aligned}$$

where we use the independence between  $c(A)$  and  $D(A)$  given  $A \in \mathcal{T}$ . Plugging the two terms into (4), we obtain that

$$\begin{aligned} & \mathbb{E}(c(A) \mathbf{1}_{\{\bar{A}^{(d)} \in \mathcal{T}, D(\bar{A}^{(d)})=d\}} | \mathbf{y}) \\ &= \frac{1}{\sqrt{2}} \left[ \varphi_0(\bar{A}^{(d)}) - \tilde{\rho}_d(\bar{A}^{(d)}) w_d(\bar{A}^{(d)}) / (1 + \tau_{j-1}^{-1}) \right. \\ &\quad \left. \times (-1)^{\mathbf{1}(A \text{ is the left child of } \bar{A}^{(d)})} \cdot \psi_0(\bar{A}^{(d)}) \right] \tilde{\lambda}_d(\bar{A}^{(d)}) \\ &\quad + \frac{1}{\sqrt{2}} \varphi_1(\bar{A}^{(d)}) \lambda_d(\bar{A}^{(d)}). \end{aligned} \quad (5)$$

Combining the result in (3) and (5), and comparing it with  $\varphi_0(A)$  in (2), we obtain that

$$\varphi(A) = \varphi_0(A) / (1 - \tilde{\eta}(A)) + \frac{1}{\sqrt{2}} \sum_{d \in \mathcal{P}(A)} \varphi_1(\bar{A}^{(d)}) \lambda_d(\bar{A}^{(d)}),$$

which concludes the proof by plugging in  $\varphi_1(\cdot) = \varphi(\cdot) - \varphi_0(\cdot)$ .  $\square$

## E SENSITIVITY ANALYSIS

In this section, we conduct a sensitivity analysis for the proposed WARP framework at various choices of hyperparameters.

We first implement the method of ‘‘WARP-full’’ which chooses  $\phi$  by a full optimization of the marginal likelihood using two simulated images  $(f_1, f_2)$  explicitly given in Section G. Recall that the row of WARP selects  $\phi$  at a limited set of grid points. Table 3 shows that the MSEs of WARP-full are almost identical to the row of WARP. This observation is consistent across many scenarios we have tested. Therefore, the method of WARP seems robust in terms of hyperparameters, and we shall recommend a

TABLE 3  
Average MSEs ( $\times 10^{-2}$ ) of WARP-full and WARP based on 100 replications under the setting of Table 5.

Method	$n = 64$				$n = 128$			
	$f = f_1$		$f = f_2$		$f = f_1$		$f = f_2$	
	$\sigma = 0.1$	0.2						
WARP-full	0.02	0.04	0.04	0.12	0.01	0.02	0.02	0.05
WARP	0.02	0.04	0.04	0.11	0.01	0.02	0.02	0.05

TABLE 4  
Sensitivity analysis of WARP when hyperparameters are selected differently using the Shepp-Logan phantom test image ( $256 \times 256$ ) in Matlab at various  $\sigma$ . The average MSEs ( $\times 10^{-2}$ ) are reported based on 5 replications.

$\tau$	$\eta$	0.1	0.3	0.5	0.7
function	constant	0.03	0.27	0.57	0.89
function	mix	0.03	0.27	0.58	0.88
function	full	0.03	0.27	0.57	0.87
mix	constant	0.03	0.26	0.57	0.94
mix	mix	0.03	0.26	0.57	0.91
mix	full	0.03	0.27	0.57	0.91
full	constant	0.03	0.27	0.58	0.86
full	mix	0.03	0.27	0.56	0.87
full	full	0.03	0.27	0.57	0.88

maximization over a small set of grid points as default. In addition, we investigate the performances of WARP at various choices of  $\gamma$  in  $\mathbb{B}_0$  including Laplace and quasi-Cauchy priors. We find out these  $\mathbb{B}_0$  lead to almost exactly the same MSEs as normal priors (results not shown here).

We further investigate the sensitivity of WARP by considering the following ways to select hyperparameters  $\tau$  and  $\eta$ :

- $\tau$ : “function” (we use  $\tau_j = 2^{-\alpha_j} \tau_0$  as in Section 3); “mix” (we use separate  $\tau_j$  only for the last two levels and a constant for other levels, therefore we have three free parameters for  $\tau$ ); “full” (we use separate  $\tau_j$ ’s for all levels  $j$ )
- $\eta$ : “constant” (we use  $\eta(A) = \eta_0$  for all  $A$  as in Section 3); “mix” (we use  $\eta_j$  for the last two levels and a constant for other levels, therefore we have three free parameters for  $\eta$ ); “full” (we use separate  $\eta_j$ ’s for all levels  $j$ ).

Table 4 shows that the MSEs only exhibit minimal differences across various combinations of tuning approaches. This confirms the previous findings that the proposed framework is not sensitive to hyperparameters.

## F 12 WIDELY USED TEST IMAGES

The 12 widely used test images used in Section 3.3 are plotted in Figure 7.

## G 3D IMAGES

Unlike WARP which is directly applicable to  $m$ -dimensional data for  $m > 2$ , other methods compared in 3.1 such as Wedgelet, TI-2D-Haar, and BPFA may require substantial modifications for a new dimensional setting. SHAH is conceptually applicable for 3D data, but the existing software takes hours to days in the tuning step for 3D images of



Fig. 7. The widely used 12 test images.

intermediate size while its performance in 2D settings is not among top two. Therefore, we compare WARP with RM and a collection of other approaches, including a 3D image denoising method via local smoothing and non-parametric regression (LSNR) proposed by [55], anisotropic diffusion (AD) method [56], total variation minimization (TV) method [57] and optimized non-local means (ONLM) method [58]. The TV method is modified by [55] by minimizing a 3D-version of the TV criterion. We adopt simulation settings in [55], which uses two artificial 3D images with the following true intensity functions:

$$f_1(x, y, z) = -(x - \frac{1}{2})^2 - (y - \frac{1}{2})^2 - (z - \frac{1}{2})^2 + \mathbf{1}_{\{(x, y, z) \in R_1 \cup R_2\}},$$

where  $R_1 = \{(x, y, z) : |x - \frac{1}{2}| \leq \frac{1}{4}, |y - \frac{1}{2}| \leq \frac{1}{4}, |z - \frac{1}{2}| \leq \frac{1}{4}\}$  and  $R_2 = \{(x, y, z) : (x - \frac{1}{2})^2 + (y - \frac{1}{2})^2 \leq 0.15^2, |z - \frac{1}{2}| \leq 0.35\}$ ;

$$f_2(x, y, z) = \frac{1}{4} \sin(2\pi(x + y + z) + 1) + \frac{1}{4} + \mathbf{1}_{\{(x, y, z) \in S_1 \cup S_2\}},$$

where  $S_1 = \{(x, y, z) : (x - \frac{1}{2})^2 + (y - \frac{1}{2})^2 \leq \frac{1}{4}(z - \frac{1}{2})^2, 0.2 \leq z \leq 0.5\}$  and  $S_2 = \{(x, y, z) : 0.2^2 \leq (x - \frac{1}{2})^2 + (y - \frac{1}{2})^2 + (z - \frac{1}{2})^2 \leq 0.4^2, z < 0.45\}$ .

Table 5 shows the comparison of various methods using MSE. It is worth mentioning that the numerical records for the other five methods to estimate  $f_1$  and  $f_2$  are from [55] as the code is not immediately available and the running time for some method such as LSNR can take hours to days (including the tuning step). WARP is uniformly the best approach among all the selected methods at least under the simulation setting.

TABLE 5  
 3D denoising for two images  $f_1, f_2$  in terms of MSE ( $\times 10^{-2}$ ). WARP uses  $5 \times 5 \times 5$  local shifts and are based on 100 replications. The mean of 100 MSEs is reported, and the maximum standard error is 0.00.

Method	$n = 64$				$n = 128$			
	$f = f_1$		$f = f_2$		$f = f_1$		$f = f_2$	
	$\sigma = 0.1$	0.2						
WARP	<b>0.02</b>	<b>0.04</b>	<b>0.04</b>	<b>0.11</b>	<b>0.01</b>	<b>0.02</b>	<b>0.02</b>	<b>0.05</b>
LSNR	0.03	0.08	0.06	0.13	<b>0.01</b>	<b>0.03</b>	<b>0.02</b>	0.06
TV	0.03	0.09	0.06	0.15	<b>0.01</b>	0.04	0.03	0.06
AD	0.06	0.35	0.07	0.38	0.03	0.20	0.04	0.22
ONLM	0.03	0.12	0.06	0.14	<b>0.01</b>	0.06	0.03	0.06
RM	0.22	0.33	0.11	0.26	0.08	0.19	0.06	0.14