# Image Classification using Graph Neural Network and Multiscale Wavelet Superpixels

Varun Vasudevan[a,1], Maxime Bassenne[b,1], Md Tauhidul Islam[b,*], and Lei Xing[b]

[a]Institute for Computational and Mathematical Engineering, Stanford University, Stanford, CA-94305, USA

[b]Department of Radiation Oncology, Stanford University, Stanford, CA-94305, USA

*Corresponding author. Email: tauhid@stanford.edu

## Abstract

Prior studies using graph neural networks (GNNs) for image classification have focused on graphs generated from a regular grid of pixels or similar-sized superpixels. In the latter, a single target number of superpixels is defined for an entire dataset irrespective of differences across images and their intrinsic multiscale structure. On the contrary, this study investigates image classification using graphs generated from an image-specific number of multiscale superpixels. We propose WaveMesh, a new wavelet-based superpixeling algorithm, where the number and sizes of superpixels in an image are systematically computed based on its content. WaveMesh superpixel graphs are structurally different from similar-sized superpixel graphs. We use SplineCNN, a state-of-the-art network for image graph classification, to compare WaveMesh and similar-sized superpixels. Using SplineCNN, we perform extensive experiments on three benchmark datasets under three local-pooling settings: 1) no pooling, 2) GraclusPool, and 3) WavePool, a novel spatially heterogeneous pooling scheme tailored to WaveMesh superpixels. Our experiments demonstrate that SplineCNN learns from multiscale WaveMesh superpixels on-par with similar-sized superpixels. In all WaveMesh experiments, GraclusPool performs poorer than no pooling / WavePool, indicating that poor choice of pooling can result in inferior performance while learning from multiscale superpixels.

## 1  Introduction

Convolutional neural networks (CNNs) achieve the best performance on various image classification tasks. CNNs learn to classify images from a regular pixel-grid representation of the image. Two limitations of this approach are: 1) Although not all pixels provide an equal amount of new information, by design, the filters in the first layer of a CNN operate on each pixel from top-left to bottom-right in the same way; 2) CNNs require images to be of the same size. Therefore, images are typically resized to a prescribed size before feeding into a CNN. In applications that use standard CNN architectures or pre-trained models on a new image classification dataset, the images are typically uniformly downsampled to meet the input size requirements of the architecture being used. Uniform downsampling may be suboptimal as real data naturally exhibits spatial and multiscale heterogeneity. Few studies
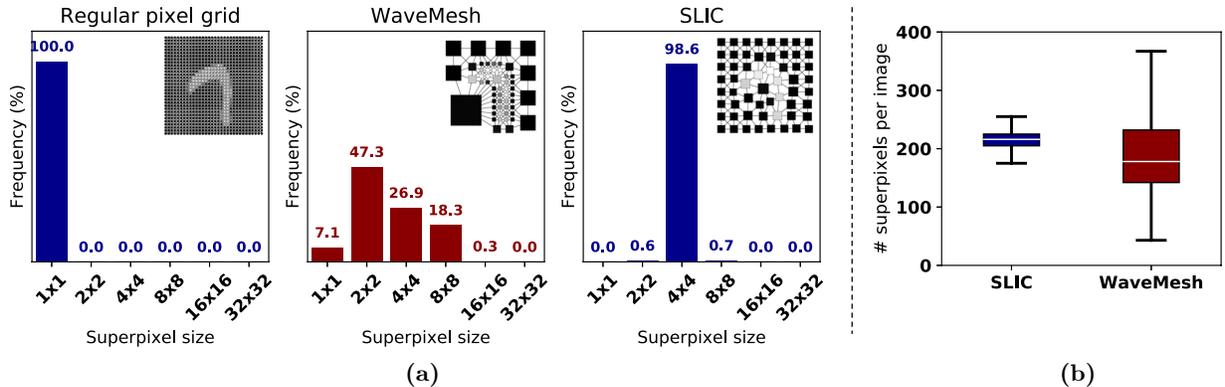
---

[1]Equal contribution

Figure 1: (a) Average distribution of superpixel size averaged across MNIST training dataset for different superpixel representation: none (left), WaveMesh (center), and similar-sized SLIC superpixels (right). In each panel, an insert shows the graph representation of a single sample for illustration. Size of a node in the graph is proportional to the superpixel size. SLIC superpixels are not cubic yet the x-axis binning is chosen to match other plots. (b) Boxplots of the # superpixels per image for CIFAR-10 training dataset.

have explored the impact of input image resolution on model performance [1], despite its recognized importance [2].

In contrast to CNNs, graph neural networks (GNN) learn from a graph representation of the image. Several studies have shown promise in classifying images from graphs using GNNs [3–9]. Unlike CNNs, the GNNs used in these studies do not require input graphs to have the same size/structure (e.g., number of nodes and edges can be different). However, these studies have been restricted to graphs that represent either a regular grid of pixels or similar-sized superpixels. In the latter, a single target number of superpixels is defined for an entire dataset irrespective of differences across images and their intrinsic multiscale structure.

In summary, while GNNs do not impose any restrictions on the size of superpixels in an image, prior studies have not systematically explored classifying images from graphs that represent multiscale superpixels. In this paper, we fill this gap by investigating image classification using a multiscale superpixel representation that can be considered as in between the regular-grid and similar-sized superpixel representations as shown in Figure 1. Our contributions are as follows.

- We present WaveMesh, an algorithm to superpixel (compress) images. WaveMesh is based on the quadtree representation of the wavelet transform. Our sample-specific method leads to non-uniformly distributed multiscale superpixels. The algorithm systematically computes the number and size of superpixels in an image based on the image content. WaveMesh requires at most one tunable parameter. WaveMesh superpixels allow us to rethink the process of downsampling (superpixeling) images to a fixed size.

- We propose WavePool, a spatially heterogeneous pooling method tailored to WaveMesh superpixels. WavePool preserves spatial structure leading to visually interpretable intermediate graphs. WavePool generalizes the classical pooling employed in CNNs and easily integrates with existing GNNs.

2

- We compare the performance of WaveMesh and similar-sized superpixels using SplineCNN, a state-of-the-art network for image graph classification [4]. Using SplineCNN, we perform extensive experiments on three benchmark datasets under three local-pooling settings: no pooling, GraclusPool, and WavePool.

# 2 Related work

**Superpixeling.** Grouping pixels to form superpixels was proposed as a preprocessing mechanism that preserves most of the structure necessary for image segmentation [10]. For a detailed review and evaluation of various state-of-the-art superpixeling algorithms see [11] and [12]. Few of them are ERS, SLIC, SEEDS, MSS, ERGC, LSC, ETPS, and SCALP. These algorithms generate similar-sized superpixels, and were originally developed and evaluated in the context of image segmentation and not image classification.

**GNN for image graph classification.** Many studies have demonstrated the representational power and generalization ability of GNNs on image graph classification tasks using similar-sized superpixels. SplineCNN is a network for learning from irregularly structured data that builds on the work of MoNet [3], but uses a spline convolution kernel instead of Gaussian mixture model kernels [4]. Recognizing the importance of spatial and hierarchical structure inherent in images, Knyazev et al. model images as multigraphs that represent similar-sized superpixels computed at different user-defined scales, and then successfully train GNNs on the multigraphs [5]. Dwivedi et al. show that message passing graph convolution networks (GCN) outperform Weisfeiler-Lehman GNNs on MNIST and CIFAR-10 datasets [7].

**Local pooling.** Local pooling is used in GNNs to coarsen the graph by aggregating nodes within specified clusters [9]. Graclus is a kernel-based multilevel graph clustering algorithm that efficiently clusters nodes in large graphs without any eigenvector computation [13]. Graclus is used in many GNNs to obtain a clustering on the nodes, which the pooling operator then uses to coarsen the graph [3, 4, 9, 14]. Hereafter, we refer to pooling based on Graclus clustering as GraclusPool. Mesquita et al. show that convolutions play a leading role in the success of GNNs and not local pooling [9].

# 3 WaveMesh: Multiscale Wavelet Superpixel

The WaveMesh algorithm is broken down into its elementary steps below: 1) images are wavelet transformed, 2) images are filtered in wavelet space by thresholding the wavelet coefficients, and 3) the superpixel mesh is generated from the wavelet-filter mask. The algorithm is rooted in wavelet theory's seminal work [15, 16]. The particular way in which wavelets are used in this work is inspired by their related application in the physical sciences [17, 18].

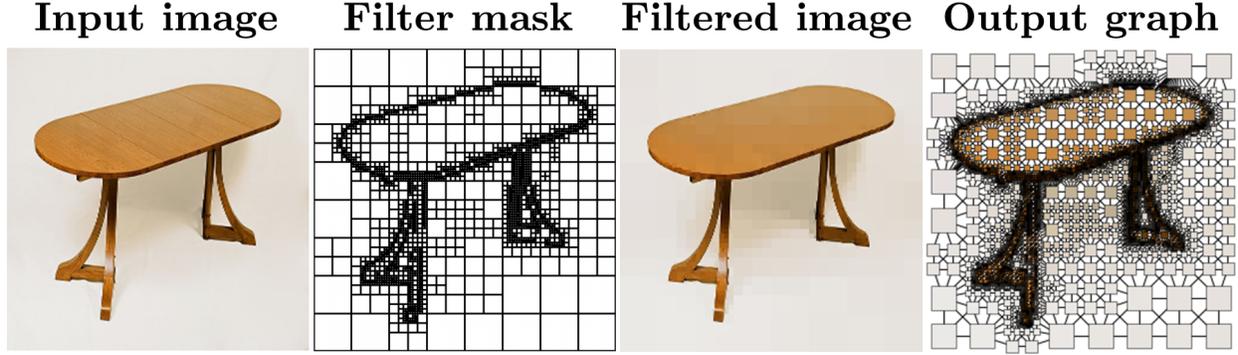| Input image | Filter mask | Filtered image | Output graph |

Figure 2: Filtering image in wavelet space generates a non-uniform superpixel mesh which is then represented as a graph. Input image is processed with the method described in section 3 with a threshold equal to five times the theoretical value.

## 3.1   Step 1: Wavelet Transform of the Input Image

Consider a two-dimensional (2D) image $I$ discretely described by its pixel values $I[\boldsymbol{x_0}]$ centered at locations $\boldsymbol{x_0} = 2^{-1}(i\Delta, j\Delta)$ of a $N{\times}N$ regular grid, where $\Delta$ is the inter-pixel spacing and $(i,j) = 1, 3, \ldots, 2N - 1$. A continuous wavelet representation of $I$ is $I(\boldsymbol{x}) = \sum_{\boldsymbol{x_0}} \widehat{I}^{(0)}[\boldsymbol{x_0}]\phi^{(0)}(\boldsymbol{x} - \boldsymbol{x_0})$, where $\boldsymbol{x}$ is the continuous pixel-space coordinate, and $\phi^0(\boldsymbol{x} - \boldsymbol{x_0})$ are scaling functions that form an orthonormal basis of low-pass filters centered at $\boldsymbol{x_0}$, with filter width $\Delta$. The scaling functions have unit energy $\langle \phi^0(\boldsymbol{x} - \boldsymbol{x_0})\phi^0(\boldsymbol{x} - \boldsymbol{x_0})\rangle = 1$, where the bracket operator $\langle y \rangle = 1/(N\Delta)^2 \int y(\boldsymbol{x})d\boldsymbol{x}$ denotes the global average for a general 2D continuous field $y(\boldsymbol{x})$. In practice, when dealing with discrete signals, $\widehat{I}^{(0)}[\boldsymbol{x_0}]$ cannot be computed exactly, since $I$ is only known at discrete points $\boldsymbol{x_0}$. Instead, it is numerically discretized and the approximation coefficients $\widehat{I}^{(0)}[\boldsymbol{x_0}]$ are estimated as an algebraic function of $I[\boldsymbol{x_0}]$. Assuming that $\phi^0(\boldsymbol{x} - \boldsymbol{x_0})$ decays fast away from $\boldsymbol{x} = \boldsymbol{x_0}$, we get $\widehat{I}^{(0)}[\boldsymbol{x_0}] = I[\boldsymbol{x_0}]/N$ [19]. This estimate for $\widehat{I}^{(0)}[\boldsymbol{x_0}]$ is the initialization stage of the recursive wavelet multiresolution algorithm (MRA) [15], which enables the computation of wavelet coefficients at coarser scales.

The decomposition of the finest-scale low-pass filter $\phi^0(\boldsymbol{x} - \boldsymbol{x_0})$ in terms of narrow-band wavelet filters $\psi^{(s,d)}(\boldsymbol{x} - \boldsymbol{x_s})$ with increasingly large filter width and a coarsest-scale scaling function $\phi^{(S)}(\boldsymbol{x} - \boldsymbol{x_S})$ yields the full wavelet-series expansion of $I$,

$$I(\boldsymbol{x}) = \sum_{s=1}^{S}\sum_{\boldsymbol{x_s}}\sum_{d=1}^{3} \widecheck{I}^{(s,d)}[\boldsymbol{x_s}]\psi^{(s,d)}(\boldsymbol{x} - \boldsymbol{x_s}) + \widehat{I}^{(S)}[\boldsymbol{x_S}]\phi^{(S)}(\boldsymbol{x} - \boldsymbol{x_S}). \tag{1}$$

Here, $\widecheck{I}^{(s,d)}[\boldsymbol{x_s}] = \langle I(\boldsymbol{x})\psi^{(s,d)}(\boldsymbol{x} - \boldsymbol{x_s})\rangle$ and $\widehat{I}^{(S)}[\boldsymbol{x_S}] = \langle I(\boldsymbol{x})\phi^{(S)}(\boldsymbol{x} - \boldsymbol{x_S})\rangle$ are wavelet and approximation coefficients at scale $s$ and $S$, respectively, obtained from the orthonormality properties of the wavelet and scaling functions. In this formulation, $d = (1, 2, 3)$ is a wavelet directionality index, and $s = (1, 2 \ldots, S)$ is a scale exponent, with $S = \log_2 N$ the number of resolution levels allowed by the grid (5 for 32×32 images). Similarly, $\boldsymbol{x_s} = 2^{s-1}(i\Delta, j\Delta)$ is a scale-dependent wavelet grid of $(N/2^s){\times}(N/2^s)$ elements where the basis functions are

centered, with $i$, $j = 1, 3, \ldots, N/2^{s-1} - 1$. The wavelet coefficients represent the local fluctuations of $I$ centered at $\boldsymbol{x_s}$ at scale $s$, while the approximation coefficient is proportional to the global mean of $I$. At each scale, the filter width of the wavelets is $2^s\Delta$.

In this study, the 2D orthonormal basis functions $\psi^{(s,d)}(\boldsymbol{x} - \boldsymbol{x_s})$ are products of one-dimensional (1D) Haar wavelets [20]. The definition of 2D wavelets as multiplicative products of 1D wavelets is a particular choice that follows the MRA formulation [15]. Haar wavelets have a narrow spatial support that provides a high degree of spatial localization. However, they display large spectral leakage at high wavenumbers since infinite spectral and spatial resolutions cannot be simultaneously attained due to limitations imposed by the uncertainty principle [19]. Different boundary conditions can be assumed for the field $I$. We do not require such a choice in this study as we restrict ourselves to square images. However, the wavelet MRA framework is not limited to square inputs and can be generalized to rectangular inputs [19].

The definition of 2D wavelets as multiplicative products of 1D wavelets is a particular choice that follows the MRA formulation described by Mallat [15], in which, the multivariate wavelets are characterized by an isotropic scale and therefore render limited information about anisotropy in the image. A large number of alternative basis functions have been recently proposed for replacing traditional wavelets when analyzing multi-dimensional data that exhibit complex anisotropic structures such as filaments and sheets. These include, but are not limited to, curvelets, contourlets, and shearlets [21].

## 3.2 Step 2: Image Filtering in Wavelet Space

The second step decomposes $I$ as

$$I = I_> + I_\leq, \tag{2}$$

where the filtered $I_>$ and remainder $I_\leq$ components correspond to the highest and lowest energetic wavelet modes of $I$, respectively. By construction, these two components are not spatially cross-correlated, as implied by the orthogonality of the wavelets and by the filtering operation described below. Note that large wavelet coefficients are associated with large fluctuations within the corresponding region of the scale-dependent wavelet grid $\boldsymbol{x_s}$, these being markers of underlying coherent structures. Under the assumptions that $I_\leq$ is additive Gaussian white noise, Donoho and Johnstone described a wavelet-based algorithm that is optimal for achieving the target decomposition (2), since it minimizes the maximum $\mathbb{L}^2$-estimation error of $I_>$ [22]. $I_>$ is obtained by retaining only the wavelet coefficients $\breve{I}^{(s,d)}$ whose absolute values satisfy

$$\breve{I}_>^{(s,d)}(\boldsymbol{x_s}) = \begin{cases} \breve{I}^{(s,d)}(\boldsymbol{x_s}) & \text{if } |\breve{I}^{(s,d)}(\boldsymbol{x_s})| \geq T, \\ 0 & \text{otherwise,} \end{cases} \tag{3}$$

for all scales $s$, positions $\boldsymbol{x_s}$ and directions $d$. In (3), $T$ is a theoretical threshold defined as

$$T = \sqrt{2\sigma_{I_\leq}^2 \ln N^2}, \tag{4}$$

where $\sigma_{I_\leq}^2$ is the unknown variance of $I_\leq$. In this study, the iterative method of Azzalini et al. is employed, which converges to $T$ starting from a first iteration where $\sigma_{I_\leq}^2$ in (4)

is substituted by the variance $\sigma_I^2$ of the total image $I$ [23]. This iterative procedure does not introduce significant computational overhead, since only one wavelet transform is required independently of the number of iterations. The algorithm does not introduce any hyperparameter when the theoretical threshold value is used. Note that the threshold is image-dependent, thereby ensuring that the algorithm adapts the number of superpixels to each image appropriately. The above filtering operation is equivalent to applying a binary filter mask to wavelet coefficients, denoted as wavelet-filter mask below.

The iterative method is deemed as converged when the relative variation in the estimated threshold $T$ is less than $0.1\%$ across consecutive iterations. A maximum of $\mathcal{O}(10)$ iterations were required to obtain the results presented in this paper. The overall computational cost is $\mathcal{O}(n_i M)$, where $n_i$ is the number of iterations and $M$ is the number of pixels in the image [23]. In this work, we allow for further reduction in number of superpixels by varying the threshold $T$ to take larger values. Figure 2 illustrates the application of this wavelet filtering method on an RGB image (filtering is applied to each channel independently) [24]. Most of the structural and edge information is preserved at all scales. However, a drawback of the method is that the superpixel boundaries are necessarily regular and axis-aligned.

## 3.3   Step 3: Generating Superpixel Mesh from Wavelet-filter Mask

To generate superpixels for a given image, the final step is a grid adaptation based on the wavelet-filter mask described in subsection 3.2. The result is a non-uniform grid of multiscale superpixels adapted around regions of the image with high variability.

**Quadtree representation.** The algorithm is perhaps best understood by representing the wavelet coefficients in a quadtree [25], a tree data structure in which each node has exactly four children. A quadtree-based representation of wavelet coefficients was previously shown to be an efficient data structure for wavelet-based image compression [26, 27]. Here, the height of this quadtree equals the number of decomposition levels $S$ in the wavelet transform. Each vertex at a given level $s$ is associated with a triplet of wavelet coefficients $[\breve{I}^{(s,d=0)}(\boldsymbol{x_s}), \breve{I}^{(s,d=1)}(\boldsymbol{x_s}), \breve{I}^{(s,d=2)}(\boldsymbol{x_s})]$. All vertices from a given level correspond to wavelet coefficients across all locations at a given scale. The children vertices of a root vertex are the wavelet coefficients from that region in space at smaller scales. The quadtree representation of the wavelet coefficients of an $8\times8$ image is schematically represented in Figure 3(a). The number on each vertex indicates the scale, from the smallest scale $s{=}1$ associated with $2\times2$ pixel patches up to the largest scale $s{=}3$ associated with the entire $8\times8$ image. The pixel regions associated with each wavelet coefficient are delineated by solid lines in the three leftmost figures in Figure 3(c).

**Node tagging.** The vertices in the tree are tagged according to the filtering algorithm described in subsection 3.2. The tagged elements of the tree denoted by blue filled color in Figure 3(a,b) correspond to those with absolute values larger than the threshold $T$, and therefore correspond to locations in the image with important spatial variability. In the 2D case, tagging is applied if at least one of the 3 wavelet coefficients of $I$ per location is larger than the threshold. Additional tagging by green-filled color is applied to wavelet coefficients that are smaller than the threshold $T$ but that correspond to a spatial region with at least one tagged wavelet coefficient at a smaller scale. This corresponds to tagging all the ancestors
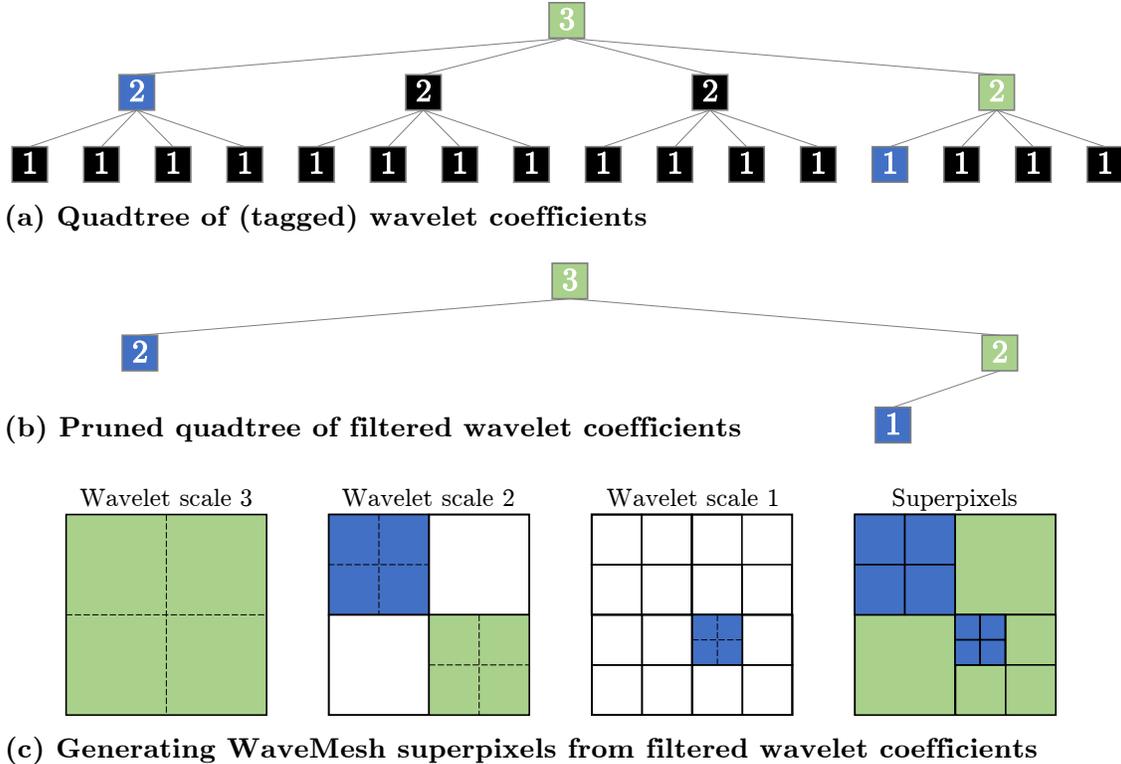
6

**(a) Quadtree of (tagged) wavelet coefficients**

**(b) Pruned quadtree of filtered wavelet coefficients**

**(c) Generating WaveMesh superpixels from filtered wavelet coefficients**

Figure 3: Illustration of the wavelet-based quadtree compression algorithm for an $8\times8$ image, along with the resulting adapted grid. Starting from the coarsest possible wavelet grid that contains just one superpixel, the algorithm adapts the grid by recursively splitting it. If the wavelet coefficient corresponding to a region is tagged (shown in blue), then that region is split into $2\times2$ superpixels.

of previously tagged vertices. This tagging procedure enforces cubic superpixels by ensuring that when there is a coherent structure at scale $s$ but not at a larger scale $s+1$, the wavelet coefficient at scale $s+1$ at that location are also tagged, hence triggering local grid refinement at level $s+1$. Non-tagged vertices are pruned as shown in Figure 3(b).

**Mesh generation.** Starting from the coarsest possible wavelet grid $\boldsymbol{x_s} = \boldsymbol{x_S}$ that contains just one superpixel, the algorithm adapts the grid by recursively splitting it as follows. If the wavelet coefficient corresponding to a region is tagged, then that region is split into $2\times2$ superpixels, which locally refines the grid. The algorithm is stopped otherwise. The same recursive loop is then applied to the refined superpixels. The final configuration of the adapted grid is obtained when none of the wavelet coefficients in any the superpixels are tagged. An example of final adapted grid is shown in Figure 3(c). The dashed lines correspond to the superpixel refinement due to the vertex being tagged. Adapted grid from a natural image is shown in Figure 2 where the superpixel mesh exhibits desired level of heterogeneity with multiscale refinement around edges. For RGB images, the most restrictive mesh is employed at every location and scale. In other words, tagging for the full image is applied if at least of the channels is tagged.
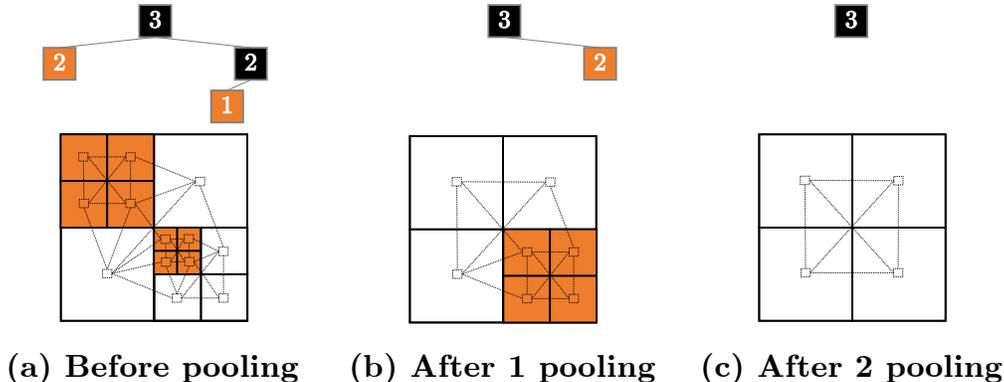
7

# 4    WavePool: Spatially heterogeneous pooling



(a) Before pooling        (b) After 1 pooling        (c) After 2 pooling

Figure 4: Illustration of WavePool from wavelet quadtree representation. Leaf nodes ($2\times2$ superpixels) are recursively pooled. In the lower panel, dashed squares and lines correspond to nodes and edges in the superpixel graph.

The proposed spatially heterogeneous pooling, WavePool, is best explained using the wavelet coefficient quadtree representation described in subsection 3.3. One WavePool operation involves aggregating all the leaf nodes of the wavelet quadtree. In the pixel domain, this step corresponds to merging patches of $2\times2$ superpixels into a parent superpixel, and aggregating the node features with a choice of pooling function (e.g. max). Figure 4 illustrates WavePool on a simple superpixel mesh and shows its effect on both the quadtree (upper panel) and region adjacency graph (lower panel) representation. In a region adjacency graph (RAG), nodes represent superpixel centroids, and edges connect neighboring superpixels. Note that GNNs are trained on RAGs. RAG is not a tree and should not be confused with the wavelet coefficient quadtree.

WavePool generalizes the classical CNN pooling operation. For a regular-pixel grid as in Figure 5, WavePool exactly matches the $2\times2$ pooling in CNN. However, this is not true with GraclusPool. Although more general than its CNN counterpart, WavePool is restricted to WaveMesh or more broadly to any quadtree based superpixels [28, 29], unlike GraclusPool.

# 5    Experimental setup

## 5.1    Datasets

To compare WaveMesh and similar-sized superpixels we perform experiments on three datasets: MNIST, Fashion-MNIST, and CIFAR-10 [30–32]. We represent superpixels by RAGs as shown in Figure 6, where mean intensity of superpixel is used as a node feature. Edges in the graph are directed with pseudo-coordinates as in [4].
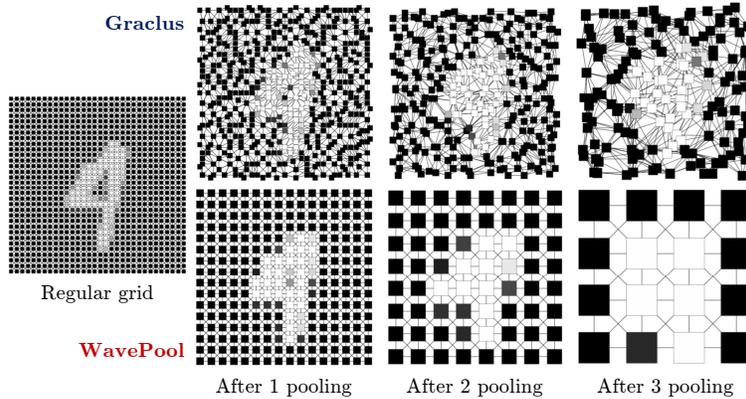
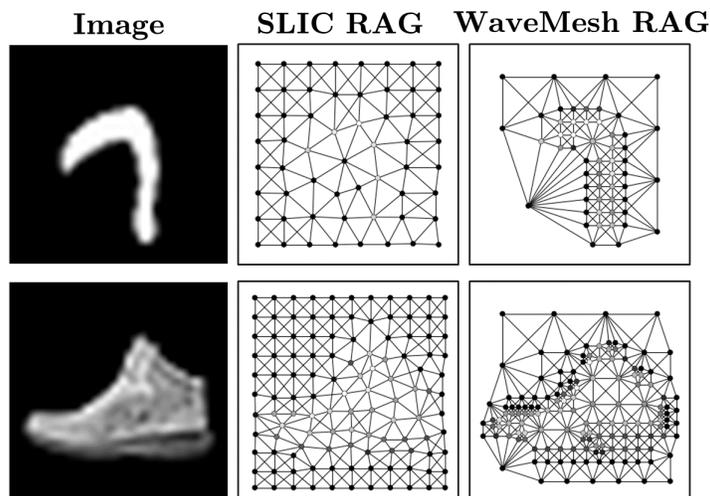Figure 5: WavePool vs. GraclusPool on a regular grid.



Figure 6: MNIST and Fashion-MNIST images with SLIC and WaveMesh RAG. Nodes represent superpixel centroids. SLIC and WaveMesh graphs are structurally different. In WaveMesh, there are more nodes along the object boundaries and fewer nodes in regions without much variation.

## 5.2   Model: SplineCNN

**Why SplineCNN?** We use SplineCNN [4], in all our experiments. SplineCNN is an ideal candidate for this study for the following reasons. First, it is a state-of-the-art GNN for image graph classification. Second, Fey et al. report that edge detecting patterns are learned by the kernels in SplineCNN when trained on superpixels. Third, the spline convolution (SConv) operator is a generalization of the convolution operator in CNNs with odd kernel size. This property of SConv operator combines nicely with WavePool that naturally collapses to classical CNN pooling on a regular pixel-grid.

    **SplineCNN configurations.** We conduct experiments on two SplineCNN configurations using the implementation available in PyTorch Geometric [33]. The configurations are:

1. SConv$((3,3),1,32) \rightarrow$ Pool $\rightarrow$ SConv$((3,3),32,64) \rightarrow$ Pool $\rightarrow$ Global mean pool $\rightarrow$

FC(128) $\rightarrow$ FC(10). Network has 30506 parameters.

2. SConv$((3,3),1,32) \rightarrow$ Pool $\rightarrow$ SConv$((3,3),32,64) \rightarrow$ Pool $\rightarrow$ SConv$((3,3),64,128) \rightarrow$ Pool $\rightarrow$ Global mean pool $\rightarrow$ FC(256) $\rightarrow$ FC(10). Network has 139178 parameters.

Here, FC denotes a fully-connected layer and Pool denotes a GraclusPool or WavePool layer or no pooling.

## 5.3 Comparison with SLIC

The SLIC superpixeling algorithm is based on $k$-means clustering [34]. We compare the performance of WaveMesh with SLIC. Many superpixeling algorithms have been proposed since SLIC, however it is still an ideal baseline for this study for the following reasons.

- SLIC is one among the six algorithms recommended by Stutz et al. after evaluating 28 state-of-the-art superpixeling algorithms [11]. All the recommended algorithms show superior performance in Boundary Recall, Undersegmentation Error, and Explained Variation (EV).

- Giraud et al. evaluates Achievable Segmentation Accuracy (ASA) of nine state-of-the-art superpixeling algorithms. SLIC is among the top six in ASA [12]. Also, the difference in ASA among top six algorithms is small.

- SLIC superpixels are used in all prior GNN studies [3–9]. Moreover, our goal is to compare multiscale WaveMesh superpixels with routinely used similar-sized superpixels for image classification, and not to find the best superpixeling algorithm for image classification using GNNs.

## 5.4 Evaluation

The two main objectives of our experiments are as follows. First, to understand the performance of WaveMesh superpixels under three different local-pooling settings, everything else being the same. Second, to understand how SplineCNN performs on SLIC and WaveMesh superpixels under the same network architecture and training settings. Together, these objectives systematically explore image classification using multiscale superpixels (WaveMesh) compared to single-scale superpixels (SLIC). Using the two SplineCNN configurations mentioned in subsection 5.2 we perform extensive experiments on WaveMesh and SLIC superpixels by varying the following. First, the number of superpixels. Second, changing the local-pooling: no pooling, GraclusPool, and WavePool.

The SplineCNN implementation in PyTorch Geometric uses Adam optimizer with an initial learning rate of 0.01, which is decreased by a factor of 10 after 15 and 25 epochs. Since the goal of our experiments is not to tune the best model for WaveMesh superpixels, we use the default hyperparameters from their implementation. We train the network for 30 epochs on MNIST and Fashion-MNIST and 75 epochs on CIFAR-10. The pooling function is max for both WavePool and GraclusPool. All experiments are repeated five times.

We also compare WaveMesh and SLIC superpixels on two traditional superpixel evaluation metrics: ASA and EV. While there are many superpixel evaluation metrics, we use ASA and EV because Giruad et al. recommends ASA to evaluate adherence to object boundaries, and EV to evaluate the color homogeneity within superpixels [12]. We calculate ASA and

Table 1: Results on MNIST superpixels. For each group of experiments, lowest value is marked in red. Experiments R1 and R2 are from [3] and [4]. Experiments R3–R5 from [7] are for models RingGNN, MoNet and GatedGCN. They report min and max value in **#Nodes**, and number of parameters in **Config**.

| # | SP | #Nodes | Config | Pool | Train acc (%) | Test acc (%) |
|---|----|--------|--------|------|---------------|--------------|
| 1 | WM | 238±50 | 1 | NP | 92.30±0.43 | 92.60±0.45 |
| 2 | WM | 238±50 | 1 | GR | 92.33±0.09 | 89.63±0.45 |
| 3 | WM | 238±50 | 1 | WP | 95.75±0.08 | 95.44±0.12 |
| 4 | WM | 238±50 | 2 | NP | 98.23±0.07 | 97.70±0.10 |
| 5 | WM | 238±50 | 2 | GR | 98.39± 0.05 | 96.80±0.11 |
| 6 | WM | 238±50 | 2 | WP | 99.68± 0.03 | 98.68±0.08 |
| 7 | SL | 241±5 | 1 | NP | 95.85±0.03 | 95.99±0.07 |
| 8 | SL | 241±5 | 1 | GR | 95.50±0.21 | 95.51±0.29 |
| 9 | SL | 241±5 | 2 | NP | 99.56±0.03 | 98.79±0.07 |
| 10 | SL | 241±5 | 2 | GR | 98.07±0.04 | 97.83±0.11 |
| 11 | WM | 57±12 | 1 | NP | 95.66±0.06 | 95.54±0.15 |
| 12 | WM | 57±12 | 1 | GR | 93.34±0.04 | 92.53±0.15 |
| 13 | WM | 57±12 | 1 | WP | 96.30±0.10 | 93.74±0.17 |
| 14 | WM | 57±12 | 2 | NP | 98.74±0.06 | 97.53±0.09 |
| 15 | WM | 57±12 | 2 | GR | 95.68±0.09 | 94.21±0.21 |
| 16 | WM | 57±12 | 2 | WP | 99.23±0.04 | 93.84±0.48 |
| 17 | SL | 59±2 | 1 | NP | 95.56±0.11 | 95.17±0.12 |
| 18 | SL | 59±2 | 1 | GR | 92.34±0.11 | 91.18±0.22 |
| 19 | SL | 59±2 | 2 | NP | 98.84±0.06 | 97.18±0.07 |
| 20 | SL | 59±2 | 2 | GR | 94.13±0.08 | 92.99±0.22 |
| R1 | SL | 75±0 | – | GR | – | 91.11 |
| R2 | SL | 75±0 | – | GR | – | 95.22 |
| R3 | SL | 40–75 | 105398 | – | 11.24±0.00 | 11.35±0.00 |
| R4 | SL | 40–75 | 104049 | – | 96.61±0.44 | 90.81±0.03 |
| R5 | SL | 40–75 | 104217 | – | 100.00±0.00 | 97.34±0.14 |

EV on the BSD300 dataset for SLIC and WaveMesh superpixels using the code provided by [12]. Default compactness value of 10 is used in the SLIC algorithm, and approximately 500 superpixels are generated for each image in the dataset [35].

# 6  Results and Discussion

**Format of Tables 1–3.** Classification results on MNIST, Fashion-MNIST, and CIFAR-10 graphs are reported in Tables 1–3. Experiment numbers starting with 'R' report results from prior studies. For brevity, we use the acronyms SP: superpixel, WM: WaveMesh, SL:

Table 2: Results on Fashion-MNIST superpixels. For each group of experiments, lowest value is marked in red. Experiment R1 is from [8].

| # | SP | #Nodes | Config | Pool | Train acc (%) | Test acc (%) |
|---|---|---|---|---|---|---|
| 1 | WM | 436±129 | 1 | NP | 80.34±0.42 | 79.60±0.44 |
| 2 | WM | 436±129 | 1 | GR | 80.36±0.39 | 65.35±2.94 |
| 3 | WM | 436±129 | 1 | WP | 85.77±0.18 | 76.60±0.83 |
| 4 | WM | 436±129 | 2 | NP | 86.86±0.15 | 85.71±0.16 |
| 5 | WM | 436±129 | 2 | GR | 85.40±0.10 | 75.69±1.47 |
| 6 | WM | 436±129 | 2 | WP | 92.58±0.07 | 83.66±1.49 |
| 7 | WM | 261±35 | 1 | NP | 82.54±0.20 | 81.61±0.24 |
| 8 | WM | 261±35 | 1 | GR | 81.32±0.13 | 76.75± 0.33 |
| 9 | WM | 261±35 | 1 | WP | 85.91±0.10 | 81.35±0.69 |
| 10 | WM | 261±35 | 2 | NP | 88.20±0.24 | 86.60±0.14 |
| 11 | WM | 261±35 | 2 | GR | 85.18±0.18 | 79.78±0.46 |
| 12 | WM | 261±35 | 2 | WP | 92.34±0.15 | 87.65±0.36 |
| 13 | SL | 259±7 | 1 | NP | 83.60±0.16 | 82.37±0.25 |
| 14 | SL | 259±7 | 1 | GR | 82.91±0.08 | 81.49±0.38 |
| 15 | SL | 259±7 | 2 | NP | 89.01±0.31 | 87.29±0.30 |
| 16 | SL | 259±7 | 2 | GR | 86.71±0.10 | 85.00±0.32 |
| 17 | WM | 134±22 | 1 | NP | 83.23±0.09 | 82.04±0.10 |
| 18 | WM | 134±22 | 1 | GR | 80.92±0.16 | 78.85±0.09 |
| 19 | WM | 134±22 | 1 | WP | 85.18±0.13 | 80.42±0.33 |
| 20 | WM | 134±22 | 2 | NP | 87.88±0.07 | 85.62±0.26 |
| 21 | WM | 134±22 | 2 | GR | 83.84±0.13 | 80.60±0.21 |
| 22 | WM | 134±22 | 2 | WP | 90.98±0.13 | 82.65±0.52 |
| 23 | SL | 118±4 | 1 | NP | 83.01±0.15 | 81.44±0.19 |
| 24 | SL | 118±4 | 1 | GR | 81.46±0.16 | 79.59±0.34 |
| 25 | SL | 118±4 | 2 | NP | 88.31±0.19 | 86.10±0.25 |
| 26 | SL | 118±4 | 2 | GR | 84.11±0.15 | 82.12±0.27 |
| R1 | SL | ≤ 75 | – | – | – | 83.07 |

SLIC, NP: no pooling, GR: GraclusPool, and WP: WavePool. In each table, experiments are partitioned into groups (separated by a mid-rule). All experiments within a group are identical except for pooling. For each experiment, we report the mean and standard deviation values for number of nodes (superpixels), train and test accuracy. Lowest test accuracy within a group is highlighted in red.

**Number of superpixels.** Experiments 1–6 in Tables 1–3 uses WaveMesh superpixels obtained using the theoretical threshold $T$. The other WaveMesh experiments are on fewer superpixels obtained by scaling (increasing) $T$. To perform one-to-one comparison we generate approximately the same number of superpixels using the SLIC implementation in scikit-learn. We were unable to generate greater than $\approx$250 superpixels for FashionMNIST using SLIC.

**Effect of pooling.** In all WaveMesh experiments test accuracy is the lowest while using GraclusPool. Comparing WaveMesh+WavePool with WaveMesh+Graclus, the former performs significantly better in majority of the experiments. Therefore, cluster assignment has an effect on the performance of SplineCNN while using multiscale WaveMesh superpixels. This is unlike what was observed in [9] with MNIST single-scale SLIC superpixels on other GNN models. When comparing WaveMesh+NoPooling with WaveMesh+WavePool there is no clear winner. Similarly, SLIC+NoPooling is atleast as good as SLIC+Graclus. Therefore, while training a model on superpixels it is good to begin with a model that has just convolution layers.

**WaveMesh versus SLIC.** When comparing test accuracy in the absence of pooling WaveMesh is better than SLIC for CIFAR-10, better or same as SLIC for Fashion-MNIST, and there is no clear trend for MNIST. In the presence of pooling, WaveMesh+WavePool is better than or on-par with SLIC+Graclus. Recalling the objectives of this study, we conclude that the performance of SplineCNN with multiscale WaveMesh superpixels is just as good as SLIC superpixels.

**Performance with other GNNs.** Dwivedi et al. benchmark the performance of GNN models on MNIST and CIFAR-10 SLIC superpixels [7]. From their results, RingGNN and GatedGCN perform the worst and best, respectively. Therefore, it is clear that not all GNN models perform well on SLIC superpixels. While we have not trained these GNN models with WaveMesh superpixels, we expect a similar trend with WaveMesh. Results from their paper for RingGNN, MoNet and GatedGCN are reported in Table 1 (experiment R3–R5) and Table 3 (experiment R1–R3). Results for MoNet are shown because SplineCNN builds on the work of MoNet.

**ASA and EV.** Table 4 shows ASA and EV averaged across images in the BSD300 dataset. Giruad et al. recommends ASA to evaluate adherence to object boundaries, and EV to evaluate the color homogeneity within superpixels [12]. WaveMesh is comparable with SLIC on ASA and inferior on EV.

# 7 Conclusion

Prior GNN studies on image graph classification have been restricted to graphs that represent a regular grid or similar-sized SLIC superpixels. To fill this gap, we investigated image classification using multiscale superpixels and SplineCNN. We proposed 1) WaveMesh, a

Table 3: Results on CIFAR-10 superpixels. For each group of experiments, lowest value is marked in red. Experiments R1–R3 from [7] are for models RingGNN, MoNet and GatedGCN. They report min and max value in **#Nodes**, and number of parameters in **Config**.

| # | SP | #Nodes | Config | Pool | Train acc (%) | Test acc (%) |
|---|----|--------|--------|------|---------------|--------------|
| 1 | WM | 197±82 | 1 | NP | 51.18±0.15 | 50.59±0.17 |
| 2 | WM | 197±82 | 1 | GR | 52.63±0.36 | 43.36±0.72 |
| 3 | WM | 197±82 | 1 | WP | 55.04±0.21 | 52.58±0.21 |
| 4 | WM | 197±82 | 2 | NP | 61.52±0.36 | 58.33±0.40 |
| 5 | WM | 197±82 | 2 | GR | 60.28±0.18 | 50.42±0.27 |
| 6 | WM | 197±82 | 2 | WP | 70.25±0.30 | 56.89±0.31 |
| 7 | SL | 215±15 | 1 | NP | 48.37±0.24 | 47.25±0.21 |
| 8 | SL | 215±15 | 1 | GR | 50.96± 0.51 | 45.87±0.28 |
| 9 | SL | 215±15 | 2 | NP | 58.61±0.36 | 56.60±0.18 |
| 10 | SL | 215±15 | 2 | GR | 59.09± 0.20 | 50.69±0.45 |
| R1 | SL | 85–150 | 105165 | – | 19.56±16.40 | 19.30±16.12 |
| R2 | SL | 85–150 | 104229 | – | 65.92±2.52 | 54.66±0.52 |
| R3 | SL | 85–150 | 104357 | – | 94.55±1.02 | 67.31±0.31 |

Table 4: Performance on BSD300 dataset.

| Metric | SLIC | WaveMesh |
|--------|------|----------|
| Num superpixels | 475±20 | 515±129 |
| Achievable Segmentation Accuracy (ASA) | 0.967±0.014 | 0.950±0.025 |
| Explained Variation (EV) | 0.870±0.093 | 0.783±0.128 |

novel wavelet-based superpixeling algorithm, where the number and sizes of superpixels in an image are computed based on its content, and 2) WavePool, a novel spatially heterogeneous pooling scheme tailored to WaveMesh superpixels. Due to the multiscale nature of WaveMesh superpixels, their RAGs are structurally different from those of SLIC. Extensive experiments on benchmark datasets show that poor choice of local-pooling negatively affects the performance of SplineCNN while using WaveMesh superpixels. We also show that SplineCNN learns from multiscale WaveMesh superpixels on-par with SLIC superpixels under the same setting. Further investigation similar to [11] and [12] is required to rank superpixeling algorithms for image graph classification using popular GNNs.

# References

[1] C. F. Sabottke and B. M. Spieler. "The effect of image resolution on deep learning in radiography". In: *Radiology: Artificial Intelligence* 2.1 (2020), e190015.

[2] P. Lakhani. "The Importance of Image Resolution in Building Deep Learning Models for Medical Imaging". In: *Radiology: Artificial Intelligence* 2.1 (2020), e190177.

[3] F. Monti, D. Boscaini, J. Masci, E. Rodola, J. Svoboda, and M. M. Bronstein. "Geometric deep learning on graphs and manifolds using mixture model cnns". In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2017, pp. 5115–5124.

[4] M. Fey, J. Eric Lenssen, F. Weichert, and H. Müller. "SplineCNN: Fast geometric deep learning with continuous B-spline kernels". In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2018, pp. 869–877.

[5] B. Knyazev, X. Lin, M. R. Amer, and G. W. Taylor. "Image classification with hierarchical multigraph networks". In: *arXiv preprint arXiv:1907.09000* (2019).

[6] B. Knyazev, G. W. Taylor, and M. Amer. "Understanding attention and generalization in graph neural networks". In: *Advances in Neural Information Processing Systems*. 2019, pp. 4202–4212.

[7] V. P. Dwivedi, C. K. Joshi, T. Laurent, Y. Bengio, and X. Bresson. "Benchmarking graph neural networks". In: *arXiv preprint arXiv:2003.00982* (2020).

[8] P. H. Avelar, A. R. Tavares, T. L. da Silveira, C. R. Jung, and L. C. Lamb. "Superpixel Image Classification with Graph Attention Networks". In: *arXiv preprint arXiv:2002.05544* (2020).

[9] D. Mesquita, A. H. Souza, and S. Kaski. "Rethinking pooling in graph neural networks". In: *arXiv preprint arXiv:2010.11418* (2020).

[10] X. Ren and J. Malik. "Learning a classification model for segmentation". In: *Proceedings Ninth IEEE International Conference on Computer Vision*. IEEE. 2003, p. 10.

[11] D. Stutz, A. Hermans, and B. Leibe. "Superpixels: An evaluation of the state-of-the-art". In: *Computer Vision and Image Understanding* 166 (2018), pp. 1–27.

[12] R. Giraud, V.-T. Ta, and N. Papadakis. "Evaluation Framework of Superpixel Methods with a Global Regularity Measure". In: *arXiv preprint arXiv:1903.07162* (2019).

[13] I. S. Dhillon, Y. Guan, and B. Kulis. "Weighted graph cuts without eigenvectors a multilevel approach". In: *IEEE transactions on pattern analysis and machine intelligence* 29.11 (2007), pp. 1944–1957.

[14] M. Defferrard, X. Bresson, and P. Vandergheynst. "Convolutional neural networks on graphs with fast localized spectral filtering". In: *Advances in neural information processing systems*. 2016, pp. 3844–3852.

[15] S. G. Mallat. "A theory for multiresolution signal decomposition: the wavelet representation". In: *IEEE T. Pattern Anal.* 11.7 (1989), pp. 674–693.

[16]  D. L. Donoho and I. M. Johnstone. "Ideal spatial adaptation by wavelet shrinkage". In: *biometrika* 81.3 (1994), pp. 425–455.

[17]  K. Schneider and O. V. Vasilyev. "Wavelet methods in computational fluid dynamics". In: *Annu. Rev. Fluid Mech.* 42 (2010), pp. 473–503.

[18]  M. Bassenne, P. Moin, and J. Urzay. "Wavelet multiresolution analysis of particle-laden turbulence". In: *Phys. Rev. Fluids* 3.8 (2018), p. 084304.

[19]  P. S. Addison. *The illustrated wavelet transform handbook: introductory theory and applications in science, engineering, medicine and finance.* CRC press, 2017.

[20]  C. Meneveau. "Analysis of turbulence in the orthonormal wavelet representation". In: *J. Fluid Mech.* 232 (1991), pp. 469–520.

[21]  G. Kutyniok and D. Labate. *Shearlets: Multiscale analysis for multivariate data.* Springer Science & Business Media, 2012.

[22]  D. L. Donoho and I. M. Johnstone. "Ideal spatial adaptation by wavelet shrinkage". In: *Biometrika* 81 (1994), pp. 425–455.

[23]  A. Azzalini, M. Farge, and K. Schneider. "Nonlinear wavelet thresholding: A recursive method to determine the optimal denoising threshold". In: *Applied and Computational Harmonic Analysis* 18.2 (2005), pp. 177–185.

[24]  M. Everingham, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman. "The pascal visual object classes (voc) challenge". In: *International journal of computer vision* 88.2 (2010), pp. 303–338.

[25]  R. A. Finkel and J. L. Bentley. "Quad trees a data structure for retrieval on composite keys". In: *Acta informatica* 4.1 (1974), pp. 1–9.

[26]  M. R. Banham and B. J. Sullivan. "A wavelet transform image coding technique with a quadtree structure". In: *[Proceedings] ICASSP-92: 1992 IEEE International Conference on Acoustics, Speech, and Signal Processing.* Vol. 4. IEEE. 1992, pp. 653–656.

[27]  M. B. Wakin, J. K. Romberg, H. Choi, and R. G. Baraniuk. "Geometric methods for wavelet-based image compression". In: *Wavelets: Applications in Signal and Image Processing X.* Vol. 5207. International Society for Optics and Photonics. 2003, pp. 507–520.

[28]  S. Tanimoto and T. Pavlidis. "A hierarchical data structure for picture processing". In: *Computer graphics and image processing* 4.2 (1975), pp. 104–119.

[29]  C. Zhang, G. Zhu, M. Chen, H. Chen, and C. Wu. "Image Segmentation Based on Multiscale Fast Spectral Clustering". In: *arXiv preprint arXiv:1812.04816* (2018).

[30]  Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. "Gradient-based learning applied to document recognition". In: *Proceedings of the IEEE* 86.11 (1998), pp. 2278–2324.

[31]  H. Xiao, K. Rasul, and R. Vollgraf. "Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms". In: *arXiv preprint arXiv:1708.07747* (2017).

[32]  A. Krizhevsky and G. Hinton. "Learning multiple layers of features from tiny images". In: *Technical Report* (2009).

[33]  M. Fey and J. E. Lenssen. "Fast Graph Representation Learning with PyTorch Geometric". In: *ICLR Workshop on Representation Learning on Graphs and Manifolds*. 2019.

[34]  R. Achanta, A. Shaji, K. Smith, A. Lucchi, P. Fua, and S. Süsstrunk. "SLIC superpixels compared to state-of-the-art superpixel methods". In: *IEEE transactions on pattern analysis and machine intelligence* 34.11 (2012), pp. 2274–2282.

[35]  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: *Proceedings Eighth IEEE International Conference on Computer Vision. ICCV 2001*. Vol. 2. IEEE. 2001, pp. 416–423.