

SPIN Road Mapper: Extracting Roads from Aerial Images via Spatial and Interaction Space Graph Reasoning for Autonomous Driving

Wele Gedara Chaminda Bandara, Jeya Maria Jose Valanarasu, and Vishal M. Patel

Abstract—Road extraction is an essential step in building autonomous navigation systems. Detecting road segments is challenging as they are of varying widths, bifurcated throughout the image, and are often occluded by terrain, cloud, or other weather conditions. Using just convolution neural networks (ConvNets) for this problem is not effective as it is inefficient at capturing distant dependencies between road segments in the image which is essential to extract road connectivity. To this end, we propose a Spatial and Interaction Space Graph Reasoning (SPIN) module which when plugged into a ConvNet performs reasoning over graphs constructed on spatial and interaction spaces projected from the feature maps. Reasoning over spatial space extracts dependencies between different spatial regions and other contextual information. Reasoning over a projected interaction space helps in appropriate delineation of roads from other topographies present in the image. Thus, SPIN extracts long-range dependencies between road segments and effectively delineates roads from other semantics. We also introduce a SPIN pyramid which performs SPIN graph reasoning across multiple scales to extract multi-scale features. We propose a network based on stacked hourglass modules and SPIN pyramid for road segmentation which achieves better performance compared to existing methods. Moreover, our method is computationally efficient and significantly boosts the convergence speed during training, making it feasible for applying on large-scale high-resolution aerial images. Code available at: https://github.com/wgcbn/SPIN_RoadMapper.git.

I. INTRODUCTION

Among all the topographic objects found in aerial images, road is one of the essential topographic features with numerous applications ranging from automatic navigation and guidance systems. Extraction of roads from aerial images helps to understand the connectivity between places and thus aid in automating navigation, disaster mitigation, and controlling traffic. Furthermore, road detection helps to determine the drivable areas for autonomous vehicles so that motion planning algorithms can be constrained on drivable roads. In addition, most of the algorithms designed for road boundary extraction and curb detection are based on road segmentation maps as the primary step [1], [2]. The extraction of road boundaries and curbs can be used to further improve the safety of autonomous driving [3], [4].

Classical methods for road segmentation involve geometric-stochastic models [5], [6], line network extraction [7], and snakes [8]. There have also been works that consider the problem of road extraction as a problem of graph extraction from images [9], [10]. Following the popularity of deep learning methods in computer vision [11],

Authors are with the Department of Electrical and Computer Engineering, The Johns Hopkins University, Baltimore, MD, USA. Emails: {wbandar1, jvalana1, vpatel136}@jhu.edu.

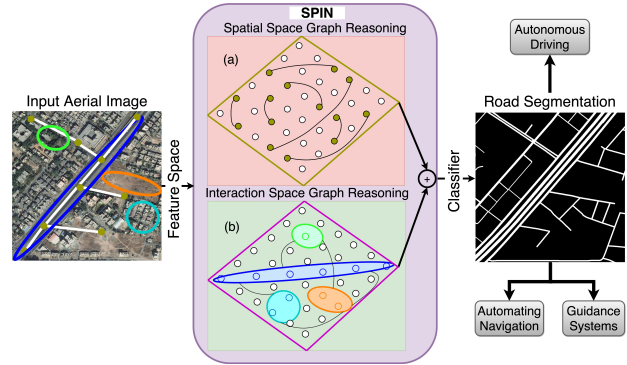


Fig. 1. An overview of our proposed method. We build graphs in two spaces: (a) spatial space and (b) a projected latent interaction graphs from feature maps. Graph reasoning in spatial space extracts connectivity between the road segments, whereas reasoning over interaction space delineates roads from other topographies. Nodes connected with lines in (a) denote how road segments are modeled to understand connectivity in the spatial space. Regions marked with different colors in (b) denote how different semantics are segregated for better road delineation in the interaction space.

[12], techniques involving ConvNets have been explored for automatic road extraction [2], [13], [14], [15]. These works pose road extraction as a problem of semantic segmentation where one tries to classify the pixels corresponding to the road from other semantics of the image.

Segmenting roads from aerial images is not a straightforward segmentation problem because roads appear at different scales in the image due to varying widths and certain road regions are often narrow and get occluded with respect to the terrain. Also, there exists some similarity of the road texture with respect to nearby regions and there are chances of occlusion due to clouds and various weather conditions. One major problem of using ConvNets directly for road segmentation is that they are not good at learning long-range dependencies due to their inherent inductive biases. In aerial images, the road structure is mostly branched throughout the image as road is a connected topography. Also, just using a ConvNet does not constrain the network to learn representations of connected road segments [16]. These issues make road segmentation from aerial images an open and challenging problem.

In this work, we focus on improving road segmentation by incorporating a global understanding of the image. Modeling dependencies and relations over regions in the image can help in understanding connectivity between the road segments. We note that transformer-based methods [17] are currently becoming popular for their property of extracting long-range dependencies. However, it is not feasible for applications on large-scale high-resolution remote sensing

datasets as it requires high compute power and significant training time. Thus, we propose using graph reasoning rather than just relying only on stacked convolutions or transformers to model global dependencies. Performing graph reasoning is light-weight and does not add on much to computation cost like transformers.

A graph convolution [18] can extract dependencies over distant regions making it more meaningful for using it to understand road information on a global scale in aerial images. Graph convolutions have been explored for video recognition [19], semantic segmentation [20] and semi-supervised classification [18]. Unlike these works, we propose performing graph reasoning in two domains - spatial and interaction space. In graph reasoning over spatial space, we build a graph over the feature space to extract dependencies between different spatial regions in the input. As we operate on the original coordinate space, performing reasoning over the graph would help to extract rich contextual information for road segmentation. For graph reasoning over interaction space, we construct a new interaction space where we model semantics with similar information together. This causes different semantic objects of the aerial image like roads, buildings, clouds, trees, and other topographic features to be modeled into different spaces. Performing graph reasoning over this interaction space would help in appropriate delineation of roads from other topographies in the image. Combining both, we propose a stand alone Spatial and Interaction space (SPIN) graph reasoning module which performs reasoning in the spatial and interaction space projected from the feature maps. Fig 1 illustrates how SPIN module helps to make road segmentation better.

SPIN extracts long range dependencies between road segments and is effective at delineating roads from other semantics present in the image. When added to a base network, we show that it improves the segmentation performance by a reasonable amount. It has numerous other advantages as well. SPIN can be plugged easily into a ConvNet architecture after a convolutional block. As SPIN learns highly contextual information, it increases the convergence rate of the network by half saving a lot of training time. This property is highly useful for training ConvNets on large-scale high-resolution images like aerial images. Adding SPIN to a ConvNet is also computationally effective as it adds on only $0.03M$ parameters. Our proposed network consists of a feature extractor using residual blocks, stacked hourglass modules with skip connections for deep feature extraction and SPIN pyramid for graph reasoning. We analyze the effectiveness of our proposed method by conducting experiments on two large-scale road segmentation datasets - DeepGlobe [21] and Massachusetts Road [22] where we achieve a better performance than existing methods in the literature.

In summary, this paper makes the following contributions:

- We propose a new module - Spatial and Interaction Space Graph Reasoning (SPIN), which when plugged into a ConvNet performs reasoning over graphs constructed on spatial and interaction space projected from the feature maps.

- We propose a new network built using stacked hourglass modules and SPIN pyramid for road segmentation from aerial images.
- We conduct extensive experiments on large-scale road segmentation datasets where we achieve better performance than existing methods both qualitatively and quantitatively.
- Our SPIN module is highly computationally efficient and helps in fast network convergence which makes training ConvNets on high-resolution aerial images quick and effective.

II. RELATED WORK

Road segmentation: Road segmentation is a well-studied problem in which we classify each pixel in a given aerial image as “road” or “no road” [15]. Early research on road segmentation primarily relied on probabilistic models to enhance connectivity by combining contextual prior conditions such as road geometry [23], [24] and color intensity [25]. In [5], geometric probability models were used to represent road images, and then maximum likelihood estimation (MLE) was used to predict road pixels. In [26], a model based on high-order conditional random fields (CRF) was used to incorporate prior knowledge of roads. However, these traditional probabilistic methods require hand-designed features and complex optimization techniques [23].

One of the earliest attempts to automatically learn features for detecting roads in aerial images using expert labeled data was proposed in [1]. In this study, unsupervised learning methods such as Principal Component Analysis (PCA) was used to initialize the feature detectors. Later, with the introduction of ConvNets in deep learning, researchers have investigated various ConvNet architectures to efficiently extract roads from aerial images [27]. Among those, encoder-decoder based architectures are widely used due to its ability to capture relatively large spatial context [27], [15]. Examples of these include U-Net [28], LinkNet [29], ResNet18 [12] and multi-branch ConvNets [30], [31]. In addition to the architectural changes, researchers also investigated different types of loss functions to replace well-known binary cross-entropy loss (BCE) to further improve the quality of road proposals and to incorporate topological constraints. In [2], a differentiable IoU loss function was proposed and most of the later work on road segmentation then used it instead of the BCE loss or combined them together for obtaining improved performance. Instead of just formulating the road extraction as a binary segmentation task, [30] introduced a multi-task learning [32] approach where both segmentation and orientation of road line segments are used to improve the connectivity of the predicted road networks.

Graph convolutions: The main limitation of Fully Convolutional Networks (FCNs) is its limited receptive field [33]. To improve the receptive field of FCNs, researchers have proposed different solutions, such as adding pooling layers [33], dilated convolutions [34], depth-wise convolutions [35], etc. However, these methods generally learn relations implicitly and are computationally expensive [36]. Instead, graph

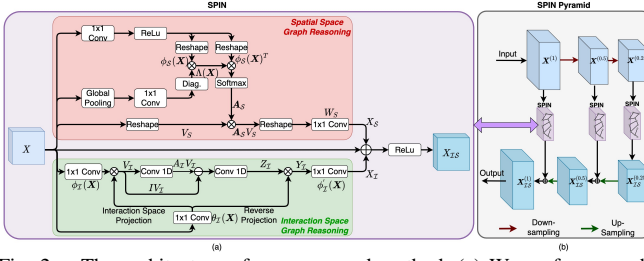


Fig. 2. The architecture of our proposed method. (a) We perform graph reasoning in both spatial and interaction space. (b) The proposed SPIN pyramid module which performs SPIN graph reasoning at multiple scales (1, 1/2, and 1/4) of original feature map to extract multi-scale long-range contextual information.

convolutions have the potential advantage of performing global reasoning on feature maps with explicit semantic meaning embedded in the graph structure. Due to this reason, many researchers have used graph convolutions in various computer vision tasks such as visual recognition [37], [38], semantic segmentation [20], [36], [39] and semi-supervised classification [18]. In [36], a graph reasoning module was proposed to capture multiple long-range contextual patterns of the original feature map through a data-dependent similarity matrix. In contrast to [36], [20] first transformed the original feature space into another latent coordinate space called interaction space and performed relational reasoning via graph convolution in the interaction space. In our proposed SPIN module, we combine the reasoning power of both spatial and interaction space graph reasoning by concatenating the individual outcomes. Further, we perform SPIN graph reasoning on different scales of the feature maps to learn multi-scale contextual relationships.

III. PROPOSED METHOD

A. SPIN graph reasoning

The proposed SPIN module is shown in Figure 2-(a). We build graphs in two spaces: spatial space and a projected latent interaction space from input feature maps. Then, graph reasoning is performed in spatial space to improve the connectivity between road segments and interaction space to delineate roads from other topographies. Assuming that spatial and interaction space graph reasoning provide different feature representations, we concatenate the output feature maps of both graph reasoning modules, as shown in Figure 2-(a), to extract rich global contextual information of road segments. In order to capture multi-scale context of input feature maps, we build SPIN pyramid by performing SPIN graph reasoning on different scales and then aggregate them as shown in Figure 2-(b).

In what follows, we elaborate on each block of the SPIN module in detail. Before that, we briefly review graph reasoning.

Graph reasoning: A graph $G = (\mathcal{V}, \mathcal{E}, \mathbf{A})$ is defined by its nodes \mathcal{V} , edges \mathcal{E} and similarity matrix \mathbf{A} that describes the similarity between each and every pixel (node) in the graph. Let $\mathbf{X} \in \mathbb{R}^{L \times C}$ denote the input feature map where C is the number of channels and $L = W \times H$. Here, W and H correspond to the width and height of \mathbf{X} . Standard 2D convolutions only share information among the positions

in a small neighborhood defined by the filter size. In order to achieve a large receptive field and to capture long-range dependencies among the pixels, ConvNet architectures stack multiple convolution layers which is highly inefficient. In contrast, a single graph convolution layer can extract long-range dependencies of input feature map very efficiently and effectively. Formally, the graph convolution is defined as [18],

$$\tilde{\mathbf{X}} = \sigma(\mathbf{A}\mathbf{X}\mathbf{W}), \quad (1)$$

where \mathbf{W} is the learnable weight matrix (usually modeled as a convolutional layer), $\sigma(\cdot)$ is the non-linear activation function (e.g. ReLU) and, \mathbf{A} and \mathbf{X} are the same as defined above. Note that the only difference between graph convolution and conventional convolution is that in graph convolution, we left-multiply the original feature map \mathbf{X} by the similarity matrix \mathbf{A} before doing the convolution operation.

With this background of graph reasoning, we now describe the two main building blocks of our proposed SPIN module: (1) Spatial space graph reasoning and, (2) Interaction space graph reasoning.

1) *Spatial space graph reasoning:* The overall procedure of spatial space graph reasoning is depicted in the red box in Figure 2-(a). As described in the previous section, the main intuition behind spatial space graph reasoning is to improve the connectivity between the predicted road segments. We first build a fully-connected graph in the spatial domain \mathcal{S} using the spatial similarity matrix \mathbf{A}_S and then perform spatial graph reasoning. We now describe the computation procedure of spatial graph reasoning in detail.

Computation of spatial similarity matrix \mathbf{A}_S : The first step of spatial graph reasoning is to compute the spatial similarity matrix $\mathbf{A}_S \in \mathbb{R}^{L \times L}$. There are different similarity metrics that have been proposed in the literature to calculate the similarity between two given pixels. The most popular are the Euclidean distance and the dot product. In our implementation, we use the dot product to compute the similarity matrix \mathbf{A}_S .

The similarity matrix \mathbf{A}_S for an input feature map \mathbf{X} can be represented as a multiplication of three transformations as follows:

$$\mathbf{A}_S = \text{Softmax}(\phi_S(\mathbf{X})\Lambda(\mathbf{X})\phi_S(\mathbf{X})^T), \quad (2)$$

where $\phi_S(\mathbf{X}) \in \mathbb{R}^{L \times M}$ is a linear transformation followed by ReLU non-linearity and $\Lambda(\mathbf{X}) \in \mathbb{R}^{M \times M}$ is the diagonal matrix. Note that M is the dimension of the intermediate feature map.

In this implementation, the linear transformation $\phi_S(\mathbf{X})$ is modeled using a 1×1 convolution layer that reduces input feature map dimension from C to M . The transformation $\Lambda(\mathbf{X})$ is represented by a global average pooling followed by a 1×1 convolution. Then we reshape the outputs $\phi_S(\mathbf{X})$ and $\Lambda(\mathbf{X})$ appropriately to perform matrix multiplication as shown in Figure 2-(a) to obtain the similarity matrix $\mathbf{A}_S \in \mathbb{R}^{L \times L}$.

Graph reasoning in spatial space: Once we have the similarity matrix \mathbf{A}_S , we can perform the spatial graph reasoning on input data according to the Eq. (1). First, we

reshape the input data appropriately and then we perform the matrix multiplication to obtain $A_S X$. Next, we multiply it by the trainable weight matrix W_S that is modeled as a 1×1 convolutional layer as shown in Figure 2-(a). Finally, we apply ReLU to obtain the spatial graph reasoned feature matrix X_S .

2) *Interaction space graph reasoning*: The overall procedure of interaction space graph reasoning is shown in the green box in Figure 2-(a). As we described earlier, the spatial space graph reasoning can improve the connectivity between predicted road segments. We now consider projecting the input feature space into another latent space, called the interaction space \mathcal{I} , where we try to delineate roads from other objects such as buildings, trees, vehicles, etc. Next, we build a graph that connects these features in the interaction space and performs a relational reasoning on the graph. After reasoning, the updated information is projected back to the original coordinate space. In what follows, we discuss these operations in detail.

Projection to interaction space: The first step is to project the original feature map X to the interaction space \mathcal{I} . This is done by the projection function $f(\cdot)$ such that the features $V_{\mathcal{I}} \in \mathcal{R}^{N \times S}$ in the interaction space are more friendly for global reasoning over disjoint and distant regions, where N is the number of nodes and S is the number of states.

In practice, we first reduce the dimension of the input feature X with the transformation $\theta_{\mathcal{I}}(X) \in \mathcal{R}^{L \times N}$ and formulate the projection function $\phi_{\mathcal{I}}(X) \in \mathcal{R}^{L \times S}$ as a linear combination of input X such that the new features can aggregate information from multiple regions. Concretely, the input feature X is projected as $V_{\mathcal{I}}$ in the interaction space \mathcal{I} through the projection function $\phi_{\mathcal{I}}(X)$ as follows:

$$V_{\mathcal{I}} = \theta_{\mathcal{I}}(X)^T \phi_{\mathcal{I}}(X). \quad (3)$$

We implement both functions $\phi_{\mathcal{I}}(\cdot)$ and $\theta_{\mathcal{I}}(\cdot)$ as 1×1 convolutional layer as shown in Figure 2-(a).

Graph reasoning in interaction space: After projecting the input feature space into interaction space, we build a fully-connected graph in the interaction space with the node similarity matrix $A_{\mathcal{I}} \in \mathcal{R}^{N \times N}$. The similarity matrix $A_{\mathcal{I}}$ is randomly initialized and learned during back propagation in contrast to the similarity matrix we defined for the spatial graph reasoning that is dependent on the input data. In addition, we use skip connection (i.e. identity matrix) that speeds up the optimization. Following Eq. (1), the graph convolution in the interaction space is formulated as:

$$Z_{\mathcal{I}} = A X W = ((I - A_{\mathcal{I}}) V_{\mathcal{I}}) W_{\mathcal{I}}, \quad (4)$$

where $W_{\mathcal{I}}$ is the trainable weight matrix. Here both matrices $W_{\mathcal{I}}$ and $A_{\mathcal{I}}$ are implemented as 1D convolution with kernel size of 1 as shown in Figure 2-(a).

Reverse projection to the original coordinate space: After graph reasoning in the interaction space, we project the output features $Z_{\mathcal{I}}$ to the original coordinate as:

$$Y_{\mathcal{I}} = \theta_{\mathcal{I}}(X)^T Z_{\mathcal{I}}, \quad (5)$$

$$X_{\mathcal{I}} = \phi'_{\mathcal{I}}(Y_{\mathcal{I}}). \quad (6)$$

We use the same projection matrix $\theta(X)$ to transform features to $Y_{\mathcal{I}} \in \mathcal{R}^{L \times S}$. Then we perform linear projection

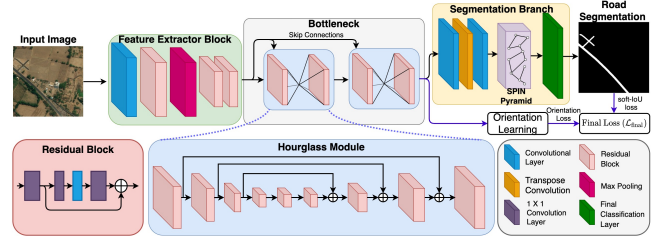


Fig. 3. Proposed network for road segmentation from aerial images.

$\phi'_{\mathcal{I}}(\cdot)$ using a 1×1 convolutional layer to transform $Y_{\mathcal{I}}$ into the original coordinate space. As a result we have the features $X_{\mathcal{I}}$ with feature dimension C at the original coordinate space.

Once we have the spatial and interaction space graph reasoning outputs, we combine them with the original input feature map and then apply ReLU non-linearity to get the final graph reasoned feature map X_{IS} . Mathematically, we can denote this as,

$$X_{IS} = \text{ReLU}(X_S + X + X_{\mathcal{I}}). \quad (7)$$

3) *SPIN pyramid*: We perform our SPIN graph reasoning at multiple scales to further increase the overall receptive field of the network and to improve long-range contextual information present in the intermediate feature maps. Concretely, we perform SPIN graph reasoning at different scales (1, 1/2, and 1/4) of the original feature map as shown in Figure 2-(b). In the results and discussion section, we conduct an ablation study to demonstrate the effect of spatial, interaction, and SPIN graph reasoning on the segmentation performance.

B. Network architecture

Feature extractor block: Operating the network at high-resolution (i.e. 256×256) requires a large GPU memory and computational power. Therefore, we employ a 7×7 convolutional layer with stride 2, followed by a residual block and a max pooling layer to bring it down to the size of 64×64 . We then add two subsequent residual modules before sending it to the hourglass module.

Bottleneck: Our proposed road segmentation network uses stack of two hourglass modules [31] at the bottleneck. The hourglass module captures information at different scales by cascading series of residual modules and max pooling layers. When the network reaches the lowest resolution, it performs bilinear upsampling and combines features across the same scales using skip connections. We feed forward the output of the bottleneck to the segmentation branch.

Segmentation branch: In the segmentation branch, we use a combination of convolution and transpose convolution layers to upsample the feature maps to the original scale. We then feed forward these feature maps to our proposed SPIN pyramid. To get the output segmentation map, we feed forward these graph reasoned feature maps from the SPIN pyramid to the final classification layer.

Orientation learning: For orientation learning, we adopt the same orientation learning technique described in *Batra et. al.* [30]. As shown in Figure 3, our road segmentation network is divided into two branches after the two hourglass

modules to support for both segmentation and orientation learning. The orientation learning task is formulated as a multi-class classification problem where, the orientation of each road-pixel is quantized into bins resulting in a total of 37 orientation classes. Please check the supplementary document for more details.

Loss function: The proposed road segmentation network utilizes predictions from intermediate feature maps to compute the loss at multiple scales: $((H/4, W/4), (H/2, W/2)$ and $(H, W))$ instead of computing it only at the original scale. This improves network’s ability to correctly predict road segments at multiple scales and helps to convergence faster. In this implementation, we make use of two loss functions: (1) segmentation loss, and (2) orientation loss. We use differentiable SoftIoU loss to compute the segmentation loss instead of using the BCE loss. The segmentation loss is computed at multiple scales. The road segmentation loss \mathcal{L}_{seg} is defined as follows,

$$\mathcal{L}_{\text{seg}} = \sum_s (1 - \text{SoftIoU}(Y_{\text{pred}}^s, Y_{\text{gt}}^s)), \quad (8)$$

where s denotes the scale having values $\{(H, W), (H/2, W/2), (H/4, W/4)\}$, Y_{pred}^s and Y_{gt}^s are the predicted and ground-truth segmentation maps at scale s , respectively. Similarly, we calculate the orientation loss at multiple scales. The orientation loss $\mathcal{L}_{\text{orient}}$ is defined as follows,

$$\mathcal{L}_{\text{orient}} = \sum_s \left(1 - \sum_{b=0}^{N_{\text{bins}}} O_b^s \log(\hat{O}_b^s) \right), \quad (9)$$

where N_{bins} is the number of bins in the quantized orientation, O_b^s and \hat{O}_b^s are the predicted and ground-truth orientation maps of orientation bin b and scale s , respectively. Finally, the overall loss function \mathcal{L} is defined as follows,

$$\mathcal{L}_{\text{final}} = \mathcal{L}_{\text{seg}} + \mathcal{L}_{\text{orient}}. \quad (10)$$

IV. EXPERIMENTAL SETTINGS

A. Datasets

Massachusetts road dataset: The Massachusetts Roads dataset [22] consists of train, validation and test sets with 1108, 14 and 49 images, respectively, each with a size of $1,500 \times 1,500$ pixels. Following [40], we fill the training images into size of 1536×1536 and then we crop each image into 512×512 patches with overlapping window of 256 pixels to make the training set. We observed that some parts of the images in the Massachusetts Road dataset are partially occluded and these images severely affect the performance of models. Hence, we removed these occluded images from the training set. Similarly, we crop each image in validation and test sets into 512×512 patches without any overlapping window. After these series of operations, the processed Massachusetts Road dataset contains 21782, 124, and 433 images with size of 512×512 , corresponding to the train, validation and test set, respectively.

DeepGlobe dataset: For the DeepGlobe dataset [21], we follow the same experimental and data preparation protocols mentioned in [30]. The DeepGlobe dataset consists of 6226 images with resolution of 1024×1024 . Following [30], we create splits of 4696 images for training and 1530 images

for testing. Then, we create the patches with 512×512 resolution with an overlapping window of 256 pixels and this results in total of 42264 images for training. Similarly, for the testing dataset also we create patches with resolution of 512×512 without any overlapping pixels and this results in total of 6116 images.

B. Implementation details

We use random crops of resolution 256×256 to train the network for the Massachusetts and DeepGlobe datasets. We use extensive data augmentation techniques such as image rotation, flipping, and mirroring. We use SGD optimizer with a batch size of 32, a momentum of 0.9 and a weight decaying of 0.0005. We use a step-learning rate scheduler with an initial learning rate of 0.01 where steps are scheduled at 50, 90, and 110. We reduce the learning rate by a factor of 0.1 at each step. We train the network for a total of 120 epochs. We implemented our model in PyTorch and used an NVIDIA Quadro RTX 8000 GPU for all of our experiments.

C. Performance metrics

For the DeepGlobe dataset accurate road segmentation masks are available and hence, we evaluate the quality of our road predictions using accurate road Intersection over Union (IoU^a) and $F1$ score. However, the groundtruth segmentation masks of Massachusetts road dataset have constant width and this will adversely affect the pure pixel based metrics. So, as proposed in [42] we also use relaxed IoU (IoU^r) with buffer size of 4 in our evaluations. Furthermore, we use Average Path Length Similarity (APLS) metric to measure the difference between ground truth and proposal graphs [43].

V. RESULTS

In this section, we compare the road segmentation performance of our SPIN Road Mapper with existing methods, quantitatively and qualitatively. In particular, we compare the performance of our method with that of Seg-Net [41], U-Net [28], LinkNet34 [29], HourGlass [31], Stack-HourGlass [31], and Batra *et al.* [30].

Quantitative results: The quantitative results are summarized in Table I. As can be seen from Table I, the proposed SPIN Road Mapper achieves the state-of-the-art (SOTA) results in terms of all the performance measures for the Massachusetts dataset. When considering the DeepGlobe dataset, our method achieves the SOTA results in terms of $F1$, IoU^r , and APLS. Furthermore, the improvement in terms of APLS metric is significant (+1.15% and +1.02% for Massachusetts and DeepGlobe datasets, respectively) and which confirms that the proposed SPIN module improves the connectivity of road segments by specially bringing gaps for occluded areas (see Fig. 5). These qualitative results confirms that the proposed SPIN module effectively captures the long-range dependencies of feature maps, and thereby helps final classifier to identify road pixels substantially well compared to the available ConvNet architectures for road segmentation.

TABLE I

A QUANTITATIVE COMPARISON OF OUR SPIN ROAD MAPPER WITH THE SOTA BASELINES IN TERMS OF F1 SCORE, IoU^r AND IoU^a .

Method	Massachusetts Road Dataset [22]						DeepGlobe Dataset [21]					
	Precision	Recall	F1	IoU^r	IoU^a	APLS	Precision	Recall	F1	IoU^r	IoU^a	APLS
Seg-Net [41]	77.34	79.84	78.57	64.71	58.59	57.76	69.48	72.97	71.19	55.26	49.20	58.55
U-Net [28]	82.46	84.34	83.39	71.51	60.97	61.33	73.55	74.98	74.26	59.06	55.02	61.23
LinkNet [29]	83.25	84.63	83.93	72.32	63.12	66.62	78.34	78.85	78.59	64.73	62.75	67.41
HourGlass [31]	81.26	81.86	81.56	68.86	61.37	65.37	79.43	80.14	79.78	66.34	60.71	65.33
Stack-HourGlass [31]	80.12	83.87	81.96	69.43	62.21	67.89	79.33	79.99	79.66	66.19	62.06	69.02
Batra <i>et al.</i> [30]	83.34	84.61	83.97	72.37	64.44	71.34	83.79	84.14	83.97	72.37	67.21	73.12
SPIN Road Mapper (ours)	83.90	85.06	84.47	73.12	65.24	72.49	84.14	84.50	84.32	72.89	67.02	74.14

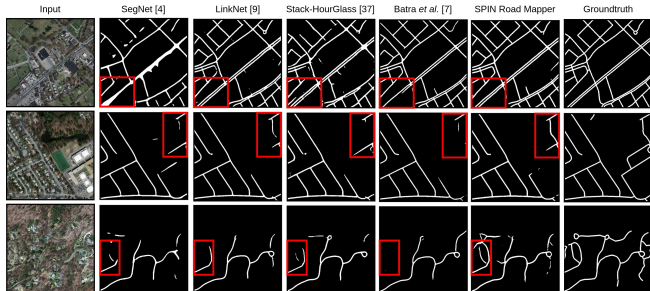


Fig. 4. A qualitative comparison between our SPIN Road Mapper and the SOTA methods.

Qualitative results: For the qualitative analysis, we visualize the predicted road maps from SegNet [41], LinkNet [29], Stack-HourGlass [21], Batra *et al.* [30], and our SPIN Road Mapper on the Massachusetts Road dataset in Figure 4. The red boxes in Figure 4 highlight the regions where our method performs better than the baseline methods. For example, consider the last row of Figure 4. Roads in the region highlighted by the red box are mostly covered by trees and buildings (as can be seen from the input aerial image), making it difficult for the baseline segmentation networks to correctly identify the presence of roads. In contrast, our method is able to predict most of the road segments due to its ability to capture long-range dependencies between road pixels through spatial graph reasoning, as well as its ability to delineate roads from surrounding structures through interaction space graph reasoning.

Ablation study: We conduct an ablation study to demonstrate the effect of spatial, interaction and SPIN graph reasoning on road segmentation. It can be seen from Table II that integrating spatial and interaction space graph reasoning to the ConvNet-based network results in increase road segmentation accuracy. Combining the spatial and interaction space graph reasoning together in SPIN pyramid results in further improvement over the individual components. In addition to the quantitative comparison, we also present a qualitative comparison in Figure 5 which clearly demonstrates how each graph reasoning technique improves the quality of road predictions. These experiments show that our proposed SPIN pyramid helps the network learn features with more global contextual information resulting in an improved performance.

Convergence: Figure 6 shows the training convergence plot of the proposed network with and without the SPIN module.

TABLE II

QUANTITATIVE RESULTS OF ABLATION STUDY.

Method	IoU^a	F1	APLS
ConvNet Only	83.97	66.58	73.01
ConvNet + Spatial GR	84.13	66.82	73.59
ConvNet + Interaction GR	84.12	66.76	73.52
ConvNet + SPIN GR	84.32	67.02	74.14

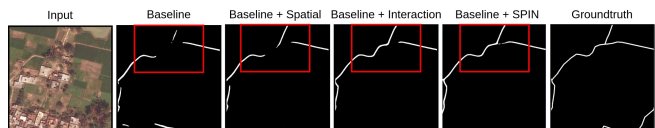


Fig. 5. A qualitative comparison for the ablation study.

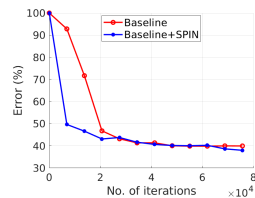


Fig. 6. The convergence characteristic for with and without the SPIN.

We can observe that adding the SPIN module helps achieve faster convergence. This leads to reduction in training time which is crucial for training ConvNets on large-scale high-resolution remote sensing datasets.

VI. CONCLUSION

We presented a Spatial and Interaction Space Graph Reasoning (SPIN) module that can be plugged into ConvNets to learn distant relationships between road segments in the feature space. Learning global dependencies are essential while extracting complex road topology from aerial images where most of the road segments are partially or completely occluded by trees, buildings, or clouds. The graph reasoning over the spatial space helps the network to extract more dependencies between different spatial regions and other contextual information whereas graph reasoning over a projected interaction space helps to delineate roads from surrounding objects. We conduct extensive experiments and compare the predicted road maps qualitatively and quantitatively with existing methods. We observe that our SPIN module helps convolutional networks to extract long-range dependencies and thereby improve the segmentation quality. SPIN is computationally light and also helps in faster convergence which are crucial while training ConvNets on large-scale high-resolutions datasets.

REFERENCES

- [1] V. Mnih and G. E. Hinton, "Learning to detect roads in high-resolution aerial images," in *European Conference on Computer Vision*. Springer, 2010, pp. 210–223.
- [2] G. Mátyus, W. Luo, and R. Urtasun, "Deeproadmapper: Extracting road topology from aerial images," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 3438–3446.
- [3] Z. Xu, Y. Sun, and M. Liu, "icurb: Imitation learning-based detection of road curbs using aerial images for autonomous driving," *IEEE Robotics and Automation Letters*, vol. 6, no. 2, pp. 1097–1104, 2021.
- [4] —, "Topo-boundary: A benchmark dataset on topological road-boundary detection using aerial images for autonomous driving," *arXiv preprint arXiv:2103.17119*, 2021.
- [5] M. Barzohar and D. B. Cooper, "Automatic finding of main roads in aerial images by using geometric-stochastic models and estimation," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 18, no. 7, pp. 707–721, 1996.
- [6] G. B. Vitor, A. C. Victorino, and J. V. Ferreira, "A probabilistic distribution approach for the classification of urban roads in complex environments," in *IEEE Workshop on International Conference on Robotics and Automation*, 2014.
- [7] D. Chai, W. Forstner, and F. Lafarge, "Recovering line-networks in images by junction-point processes," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2013, pp. 1894–1901.
- [8] I. Laptev, H. Mayer, T. Lindeberg, W. Eckstein, C. Steger, and A. Baumgartner, "Automatic extraction of roads from aerial images based on scale space and snakes," *Machine Vision and Applications*, vol. 12, no. 1, pp. 23–31, 2000.
- [9] J. Hu, A. Razdan, J. C. Femiani, M. Cui, and P. Wonka, "Road network extraction and intersection detection from aerial images by tracking road footprints," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 45, no. 12, pp. 4144–4157, 2007.
- [10] S. Hinz and A. Baumgartner, "Automatic extraction of urban road networks from multi-view aerial imagery," *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 58, no. 1-2, pp. 83–98, 2003.
- [11] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," *Advances in neural information processing systems*, vol. 25, pp. 1097–1105, 2012.
- [12] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [13] C. C. T. Mendes, V. Frémont, and D. F. Wolf, "Exploiting fully convolutional neural networks for fast road detection," in *2016 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2016, pp. 3174–3179.
- [14] D. Costea and M. Leordeanu, "Aerial image geolocalization from recognition and matching of roads and intersections," *arXiv preprint arXiv:1605.08323*, 2016.
- [15] F. Bastani, S. He, S. Abbar, M. Alizadeh, H. Balakrishnan, S. Chawla, S. Madden, and D. DeWitt, "Roadtracer: Automatic extraction of road networks from aerial images," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 4720–4728.
- [16] A. Mosinska, P. Marquez-Neila, M. Koziński, and P. Fua, "Beyond the pixel-wise loss for topology-aware delineation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 3136–3145.
- [17] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly et al., "An image is worth 16x16 words: Transformers for image recognition at scale," *arXiv preprint arXiv:2010.11929*, 2020.
- [18] T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," *arXiv preprint arXiv:1609.02907*, 2016.
- [19] X. Wang and A. Gupta, "Videos as space-time region graphs," in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 399–417.
- [20] Y. Chen, M. Rohrbach, Z. Yan, Y. Shuicheng, J. Feng, and Y. Kalantidis, "Graph-based global reasoning networks," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 433–442.
- [21] I. Demir, K. Koperski, D. Lindenbaum, G. Pang, J. Huang, S. Basu, F. Hughes, D. Tuia, and R. Raskar, "Deepglobe 2018: A challenge to parse the earth through satellite images," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2018, pp. 172–181.
- [22] V. Mnih, *Machine learning for aerial image labeling*. University of Toronto (Canada), 2013.
- [23] C. Unsalan and B. Sirmacek, "Road network detection using probabilistic and graph theoretical methods," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 50, no. 11, pp. 4441–4453, 2012.
- [24] R. Stoica, X. Descombes, and J. Zerubia, "A gibbs point process for road extraction from remotely sensed images," *International Journal of Computer Vision*, vol. 57, no. 2, pp. 121–136, 2004.
- [25] W. Wang, N. Yang, Y. Zhang, F. Wang, T. Cao, and P. Eklund, "A review of road extraction from remote sensing images," *Journal of Traffic and Transportation Engineering (English Edition)*, vol. 3, no. 3, pp. 271–282, 2016. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S2095756416301076>
- [26] J. D. Wegner, J. A. Montoya-Zegarra, and K. Schindler, "A higher-order crf model for road network extraction," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2013, pp. 1698–1705.
- [27] A. Abdollahi, B. Pradhan, N. Shukla, S. Chakraborty, and A. Alamri, "Deep learning approaches applied to remote sensing datasets for road extraction: A state-of-the-art review," *Remote Sensing*, vol. 12, no. 9, p. 1444, 2020.
- [28] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *International Conference on Medical image computing and computer-assisted intervention*. Springer, 2015, pp. 234–241.
- [29] A. Chaurasia and E. Culurciello, "Linknet: Exploiting encoder representations for efficient semantic segmentation," in *2017 IEEE Visual Communications and Image Processing (VCIP)*. IEEE, 2017, pp. 1–4.
- [30] A. Batra, S. Singh, G. Pang, S. Basu, C. Jawahar, and M. Paluri, "Improved road connectivity by joint learning of orientation and segmentation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 10385–10393.
- [31] A. Newell, K. Yang, and J. Deng, "Stacked hourglass networks for human pose estimation," in *European conference on computer vision*. Springer, 2016, pp. 483–499.
- [32] A. Kendall, Y. Gal, and R. Cipolla, "Multi-task learning using uncertainty to weigh losses for scene geometry and semantics," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 7482–7491.
- [33] A. Araujo, W. Norris, and J. Sim, "Computing receptive fields of convolutional neural networks," *Distill*, 2019, <https://distill.pub/2019/computing-receptive-fields>.
- [34] F. Yu and V. Koltun, "Multi-scale context aggregation by dilated convolutions," *arXiv preprint arXiv:1511.07122*, 2015.
- [35] A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, and H. Adam, "Mobilenets: Efficient convolutional neural networks for mobile vision applications," *arXiv preprint arXiv:1704.04861*, 2017.
- [36] X. Li, Y. Yang, Q. Zhao, T. Shen, Z. Lin, and H. Liu, "Spatial pyramid based graph reasoning for semantic segmentation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 8950–8959.
- [37] Y. Li and A. Gupta, "Beyond grids: Learning graph representations for visual recognition," in *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, 2018, pp. 9245–9255.
- [38] X. Liang, Z. Hu, H. Zhang, L. Lin, and E. P. Xing, "Symbolic graph reasoning meets convolutions," in *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, 2018, pp. 1858–1868.
- [39] X. Li, L. Zhang, A. You, M. Yang, K. Yang, and Y. Tong, "Global aggregation then local distribution in fully convolutional networks," *arXiv preprint arXiv:1909.07229*, 2019.
- [40] A. Wulamu, Z. Shi, D. Zhang, and Z. He, "Multiscale road extraction in remote sensing images," *Computational intelligence and neuroscience*, vol. 2019, 2019.
- [41] V. Badrinarayanan, A. Kendall, and R. Cipolla, "Segnet: A deep convolutional encoder-decoder architecture for image segmentation," *IEEE transactions on pattern analysis and machine intelligence*, vol. 39, no. 12, pp. 2481–2495, 2017.
- [42] V. Mnih and G. E. Hinton, "Learning to label aerial images from noisy data," in *Proceedings of the 29th International conference on machine learning (ICML-12)*, 2012, pp. 567–574.
- [43] A. Van Etten, "Spacenet road detection and routing challenge-part i," 2017.