

Memory-Constrained Semantic Segmentation for Ultra-High Resolution UAV Imagery

Qi Li, Jiaxin Cai, Yuanlong Yu, Jason Gu, Jia Pan, Wenxi Liu

Abstract—Amidst the swift advancements in photography and sensor technologies, high-definition cameras have become commonplace in the deployment of Unmanned Aerial Vehicles (UAVs) for diverse operational purposes. Within the domain of UAV imagery analysis, the segmentation of ultra-high resolution images emerges as a substantial and intricate challenge, especially when grappling with the constraints imposed by GPU memory-restricted computational devices. This paper delves into the intricate problem of achieving efficient and effective segmentation of ultra-high resolution UAV imagery, while operating under stringent GPU memory limitation. The strategy of existing approaches is to downscale the images to achieve computationally efficient segmentation. However, this strategy tends to overlook smaller, thinner, and curvilinear regions. To address this problem, we propose a GPU memory-efficient and effective framework for local inference without accessing the context beyond local patches. In particular, we introduce a novel spatial-guided high-resolution query module, which predicts pixel-wise segmentation results with high quality only by querying nearest latent embeddings with the guidance of high-resolution information. Additionally, we present an efficient memory-based interaction scheme to correct potential semantic bias of the underlying high-resolution information by associating cross-image contextual semantics. For evaluation of our approach, we perform comprehensive experiments over public benchmarks and achieve superior performance under both conditions of small and large GPU memory usage limitations.

Index Terms—Ultra-high resolution image segmentation, implicit neural representation, memory module

I. INTRODUCTION

WITH the rapid progress of photography and sensor technologies, high-definition cameras have become commonplace in the deployment of Unmanned Aerial Vehicles (UAVs). Thus, there is a growing demand for ultra-high resolution (i.e., 2K, 4K, or even higher resolution) of UAV imagery for diverse applications, such as UAV localization [1], UAV detection [2], and agricultural monitoring [3], [4].

However, handling ultra-high resolution images will cost unaffordable computing resources, which is a formidable challenge for robotic systems with limited computation power [5]–[9]. The strategy of existing UAV imagery analysis approaches is to first downscale the ultra-high resolution images to achieve computationally efficient segmentation, but this strategy tends to overlook smaller, thinner, and curvilinear regions. More importantly, as input image resolutions continue to increase, this method demands a substantially greater amount of GPU memory, rendering it impractical for systems with limited GPU memory resources. Thus, is it possible to handle semantic segmentation of arbitrarily large images using limited GPU memory?

In this paper, we attempt to address such a new ultra-high resolution segmentation problem under limited GPU memory constraints. To resolve this problem, it is appropriate to resort to segmenting local patches instead of the entire ultra-high resolution image before merging all local segmentation results. However, this approach often leads to degraded performance, and thus additional context beyond the local patch needs to be introduced to address this concern in prior methods [10]–[12]. Unfortunately, the introduction of additional context also results in increased GPU memory usage. In this paper, we propose a novel method for performing local segmentation without relying on the context beyond local patches. This approach results in a GPU memory-efficient and effective framework of ultra-high resolution semantic segmentation.

In particular, we draw inspiration from the idea of implicit neural representation (INR) [13] and design an efficient spatial-guided high-resolution query module, enabling our model to infer high-quality pixel-wise segmentation results. In specific, our model queries the nearest latent embeddings of the spatial coordinates and the high-resolution spatial information as guidance, reducing the dependency on extra contextual information beyond the local patch to the largest extent. Moreover, we propose to guide the latent embeddings to supplement the details through high-resolution semantic masks in a more straightforward manner. However, the high-resolution spatial information tends to introduce semantic estimation bias during inference. To address this concern, we introduce a memory-based interaction scheme that efficiently facilitates the high-resolution semantics learning from compact cross-image contextual representation. Compared with previous memory-based schemes, our designed scheme adds only 1MB of GPU memory computational overhead thanks to its linear complexity.

Through comprehensive experiments, we demonstrate that our proposed model outperforms the state-of-the-art approaches under the condition of small GPU memory limits over Inria Aerial and DeepGlobe datasets. Besides, our model also can better trade-off segmentation performance, GPU memory usage, and computational overhead than the latest off-the-shelf methods in large GPU memory-limited systems.

Overall, the main contributions of our paper are:

- In this paper, we raise a new research problem on ultra-high image segmentation under a GPU memory-constrained condition. We propose a GPU memory-efficient and effective framework to handle such a challenging problem.
- A novel spatial-guided high-resolution query module is introduced to predict semantic masks of local image

patches without requiring additional contextual cues beyond the local region.

- We propose an efficient memory-based interaction scheme to address the issue of semantic bias arising from the local nature of image patches, which incorporates cross-image contextual information for high-resolution query, and introduces only a mild extra GPU memory overhead.
- Our model achieves superior performance in ultra-high resolution image segmentation, surpassing prior methods by a substantial margin, particularly under small GPU memory-limited conditions. Moreover, our approach offers a balanced trade-off between segmentation performance, GPU memory usage, and computational overhead under large GPU memory constraints.

II. RELATED WORKS

In this section, we survey the recent progress of ultra-high resolution semantic segmentation and introduce the related literatures on implicit neural representation and memory schemes.

A. Ultra-high Resolution Semantic Segmentation

Semantic segmentation is modeled as a dense prediction task, many works [14]–[20] based on convolutional neural network has achieved great success. In recent years, several methods [21]–[26] use Transformer architecture to conduct semantic segmentation tasks. However, most of the work is applied to ultra-high resolution images, which raises the trade-off between performance and GPU memory. To this end, Chen *et al.* [10] integrate global image and local patch each other in the deep layer to balance performance and GPU memory usage. Limited by the speed of global and local interaction, Wu *et al.* [27] design a classification network to choose important patches for the feature fusion. To further improve performance, Huynh *et al.* [11] progressively refine coarse segmentation results via a multi-stage processing pipeline. Li *et al.* [12] introduce a multi-scale locality-aware contextual correlation and the adaptive feature fusion scheme to strengthen local segmentation. These methods are patch-based approaches, which can save GPU memory but consume time due to multiple local segmentation. In the latest work, Guo *et al.* [28] leverage the shallow-deep network to directly process the full-scale ultra-high resolution images for accelerating the inference speed. In addition, some work [29], [30] generates high-quality semantic results by refining coarse segmentation maps from a pre-trained model. Comparing with the previous works, we focus on the systems constrained by limited GPU memory, by considering a better trade-off between accuracy, GPU memory, and speed.

B. Implicit Neural Representation

In implicit neural representation (INR), the signals of the object and scene are maps from coordinates via a multi-layer perceptron (MLP) applied in modeling 3D reconstruction [31]–[36]. For example, Mildenhall *et al.* [35] present NeRF

that learns an implicit representation for a novel scene view using multiple image views. It can accurately capture the intricate details of the shape with a minimal amount of parameters. Later, INR is also used in video representation [37]–[39]. Chen *et al.* [37] propose a novel neural representation for videos that takes the time index as input and directly outputs the corresponding RGB frame. Recently, INR has also made some progress in 2D tasks [13], [30], [40], [41]. Chen *et al.* [13] represent natural and complex images in a continuous manner, which are trained in the image super-resolution task. Xu *et al.* [40] perform the spatial encoding in implicit functions and further introduce deep coordinate fusion and residual MLP architecture. Hu *et al.* [41] propose an alignment function using multi-level feature fusion for semantic segmentation. Among them, [41] is most related to our work. The key difference rests in that we utilize high-resolution spatial masks as guidance to the query module for the interpretable details of the ultra-high resolution images, while [41] only utilizes INR to perform multi-scale feature alignment.

C. Memory Scheme

Similar to the human brain, deep neural networks encode, store, and extract information via memory. In the vision tasks [42]–[48], it can capture cross-image information to serve the current image. Wang *et al.* [43] use a cross-batch memory mechanism to record and update the embeddings of past iterations for the collection of sufficient hard negative pairs. Xie *et al.* [46] relieve the ambiguity of similar objects by memorizing and tracking the regions of target objects. In [47], Kim *et al.* store the domain-agnostic categorical knowledge in the memory to achieve domain generalization for semantic segmentation. Jin *et al.* [48] set up a memory module to store the dataset-level distribution information of all classes and perform a coarse-to-fine iterative inference strategy in the memory. In our work, we introduce a memory-based interaction scheme that stores low-resolution semantic information to efficiently enhance the spatial semantics of the high-resolution image. Therefore, it is able to rectify spatial-guided semantic bias in our query module.

III. METHODOLOGY

In this section, we first describe the overall pipeline of our framework. Then, we elucidate the core components of our model, i.e., spatial-guided high-resolution query module and memory-based interaction scheme.

A. Framework Overview

In this paper, we introduce a semantic segmentation approach capable of efficiently processing ultra-high resolution images on systems with limited GPU memory.

In response to the GPU memory constraint, we follow the patch-based paradigm (e.g., [12]) in which we partition a large image into local patches, segment them subsequently, and then merge all the local segmentation results together. Formally, given an ultra-high resolution image \hat{I} with width \hat{W} and height \hat{H} , our approach partitions the image into N

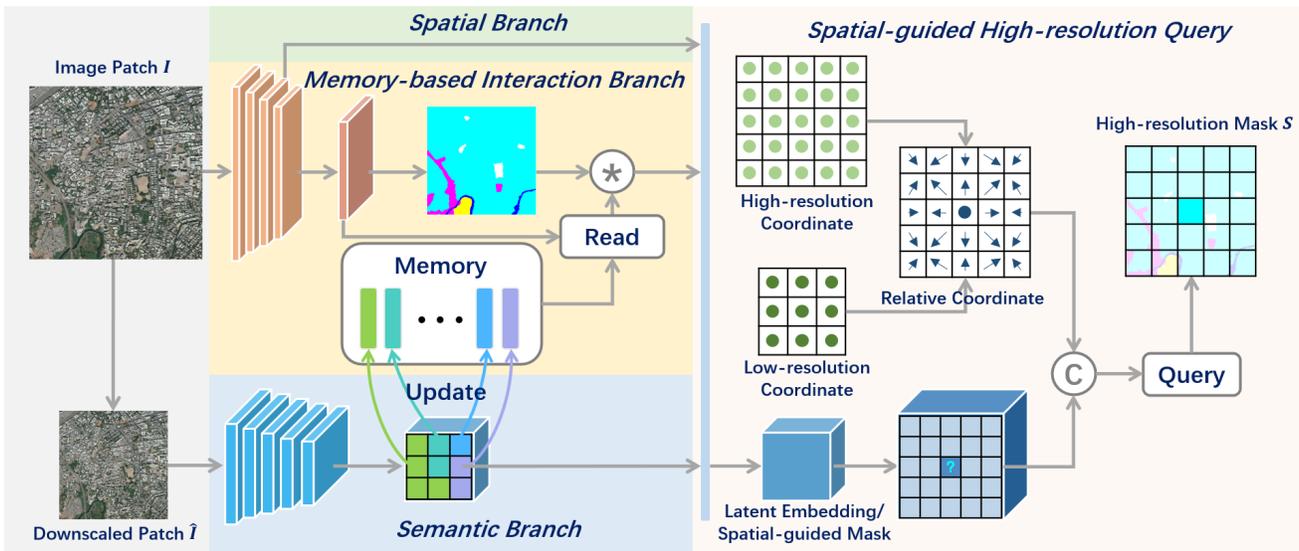


Fig. 1. The pipeline of our model. In specific, the input image is passed into a three-branch architecture to extract the low-resolution latent embedding and high-resolution spatial mask. They are jointly input to the spatial-guided high-resolution query module to obtain high-quality segmentation results. Besides, the memory-based interaction scheme is introduced to correct the semantic deviation of guide information.

overlapping local patches (denoted as I in the following) along both the row and column axes. Then, our model generates the segmentation results for each local patch, which are merged to produce an ultra-high resolution segmentation mask.

Given the possibly large resolution of the image patch I (e.g. 1280×1280), there still arises a necessity to mitigate GPU memory expenses. This, in turn, imposes constraints on the depth of the encoding network, which consequently limits its capability to extract an ample amount of semantic information. Therefore, we downscale I to a properly small resolution (e.g. 320×320) for extracting semantics. To this end, as illustrated in Fig. 1, our proposed framework is based on a three-branch architecture, comprising of the semantic branch, memory-based interaction branch, spatial branch, and spatial-guided high-resolution query module. Specifically, a partitioned image patch I and its downscaled version \hat{I} is fed into the encoders for feature extraction. The spatial branch and semantic branch obtain the visual features and semantic features, respectively. In addition, the memory-based interaction branch relies on an external memory bank to mitigate the bias in the guidance information, which associates the cross-image compact semantic representation with high-resolution spatial information for regularization. Inspired by [13], the spatial-guided high-resolution query module is dedicated to infer the high-resolution semantics in a pixel-by-pixel manner by querying the corresponding latent embeddings with the guidance of high-resolution structural information and low-resolution semantics. In the following, we will elaborate on the technical details of the two components.

B. Spatial-guided High-resolution Query Module

High-quality local image segmentation relies on contextual cues beyond local patches. However, incorporating such context can increase computational overhead. To balance segmentation quality and computational efficiency, we propose a

novel spatial-guided high-resolution query module that queries the nearest latent embedding of a given spatial coordinate to obtain the corresponding semantic result without the need for additional context.

In concrete, we first define a query function f_θ (θ is learnable parameters) over the feature maps to achieve the high-resolution semantic mask S ($S \in \mathbb{R}^{C \times H \times W}$) where C denotes the number of semantic classes. Here, the feature maps are viewed as latent embeddings evenly distributed in spatial dimensions and we assign corresponding coordinates to them. Hence, the value at the coordinate x_q of the high-resolution semantic mask S can be queried as below:

$$S(x_q) = f_\theta(z^*, x_q - x^*), \quad (1)$$

where z^* is the nearest latent embedding from x_q and x^* is the low-resolution coordinate of the latent embedding z^* in the spatial domain. Given the relative coordinates, the high-resolution query function f_θ can query the nearest semantic result set of the latent embedding z^* .

To this end, we adopt a vanilla MLP as the query function shared by each latent embedding. However, the previous study [40] implies that neural networks are insensitive to high-frequency signals and are inclined to learn low-frequency signals. This may lead to undesirable artifacts for intricate ultra-high resolution images. Consequently, we encode the relative coordinates via a periodic function ϕ to enhance the capability of the network in the high-frequency domain. Thus, the encoding function of the relative coordinates (i.e., $x_q - x^*$) is defined as:

$$\phi(x_q - x^*) = [\omega_1 \sin(x_q - x^*), \omega_1 \cos(x_q - x^*), \dots, \omega_n \sin(x_q - x^*), \omega_n \cos(x_q - x^*)], \quad (2)$$

where $\omega_1, \dots, \omega_n$ are initially set to $2e^i$ ($i \in [1, \dots, n]$) as the frequency parameters that will be fine-tuned on the training

stage. As such, $\phi(\cdot)$ maps the relative coordinate to a $2n$ -dimensional positional information. Eq. 1 can be rewritten as:

$$S(x_q) = f_\theta(z^*, x_q - x^*, \phi(x_q - x^*)). \quad (3)$$

In general, the spatial resolution of the latent embedding is much smaller than that of a semantic map, which results in the loss of spatial details during the query. To this end, we propose to utilize the higher resolution masks to provide spatial guidance for the latent embeddings. Specifically, the high-resolution segmentation masks M_b ($M_b \in \mathbb{R}^{C \times \frac{H}{2} \times \frac{W}{2}}$) and M_l ($M_l \in \mathbb{R}^{C \times H \times W}$) are computed by the memory-based interaction branch and spatial branch, respectively, which modifies the high-resolution query function as below:

$$S(x_q) = f_\theta(z^*, x_q - x^*, \phi(x_q - x^*), m_b^*, x_q - x_b^*, \phi(x_q - x_b^*), m_l^*, x_q - x_l^*, \phi(x_q - x_l^*)), \quad (4)$$

where m_b^* and m_l^* denote the nearest mask values from x_q in mask M_b and M_l . x_b^* and x_l^* are the corresponding low-resolution coordinates of mask values m_b^* and m_l^* , separately. With the cues of high-resolution spatial information, the latent embeddings can better predict fine semantic results.

C. Memory-based Interaction Scheme

In our model, the high-resolution spatial mask M_b is computed by the last layer of the memory-based interaction branch. Since this branch is relatively shallow, its estimated high-resolution mask contains semantic bias that may interfere with spatial cues for the query function. To mitigate the negative impact, it requires contextual information for regularization. Without the need for additional large computation, we propose an efficient memory-based interaction scheme that adds semantics to the mask M_b . It involves an external memory bank \mathcal{M} ($\mathcal{M} \in \mathbb{R}^{D \times C}$) that stores the semantic features across images, where D denotes the feature dimension. To reduce the computation overhead to the largest extent, our proposed memory scheme costs linear complexity only.

As the preliminary step, we randomly select an image sample from the training set. Then, we calculate a mean vector of the semantic branch features for each category as the initial value of the memory bank, with the aid of the ground-truth segmentation maps. During training, the representation of each category c ($c \in [1, \dots, C]$) in the memory bank \mathcal{M} is updated by the moving average method in the t -th iteration.

$$\mathcal{M}_t^c = m \cdot \mathcal{M}_{t-1}^c + (1 - m) \cdot \varphi(\mathcal{R}_{t-1}), \quad (5)$$

where the momentum m is set as 0.9, and φ is a transform function. \mathcal{R}_t ($\mathcal{R}_t \in \mathbb{R}^{D \times \frac{H}{4} \times \frac{W}{4}}$) is the semantic branch features of the current sample in the t -th iteration. In φ , \mathcal{R}_{t-1} is permuted with the dimension $D \times N$ ($N = \frac{HW}{16}$).

The feature representation of each category in the memory bank can be denoted as \mathcal{R}^c ($\mathcal{R}^c \in \mathbb{R}^{D \times N^c}$) that stores the representation of the category c , where N^c is the number of the pixels labeled as the category c . \mathcal{GT} ($\mathcal{GT} \in \mathbb{R}^{D \times N}$) stores

the ground-truth category labels of \mathcal{R} . Next, we calculate the cosine similarity S^c ($S^c \in \mathbb{R}^{N^c}$) between \mathcal{R}^c and \mathcal{M}^c :

$$S^c = \frac{\mathcal{R}^c \cdot \mathcal{M}^c}{\|\mathcal{R}^c\|_2 \cdot \|\mathcal{M}^c\|_2}. \quad (6)$$

Finally, the transform function φ outputs:

$$\hat{\mathcal{R}}^c = \sum_{i=1}^{N^c} \frac{1 - S_i^c}{\sum_{j=1}^{N^c} (1 - S_j^c)} \cdot \mathcal{R}_i^c. \quad (7)$$

After updating the memory bank, we associate the memory bank as the compact cross-image semantic representation to high-resolution information to enhance the mask M_b in the semantic perspective. Specifically, we read the memory bank \mathcal{M} and extract the features \mathcal{F}_b ($\mathcal{F}_b \in \mathbb{R}^{D \times \frac{H}{2} \times \frac{W}{2}}$) from the memory-based interaction branch. Then, \mathcal{F}_b is permuted into the features with the dimension $D \times \frac{HW}{4}$. Thus, we calculate the relation \mathcal{W} :

$$\mathcal{W} = \text{Softmax}\left(\frac{\mathcal{M}^\top \otimes \mathcal{F}_b}{\sqrt{D}}\right), \quad (8)$$

where \otimes is matrix multiplication. The size of \mathcal{W} is $C \times \frac{HW}{4}$ and it is reshaped as $C \times \frac{H}{2} \times \frac{W}{2}$. Last, M_b is refined as follows:

$$M_b = (1 + \mathcal{W}) \cdot M_b. \quad (9)$$

Note that, compared with the previous memory scheme with quadratic complexity (e.g. [48]), our proposed memory scheme has a linear computational complexity of $\mathcal{O}(\frac{HWC}{4})$. This is because our scheme does not associate the global pixel memory for each pixel, which spares the self-attention computation, contributing significantly to a substantial increase in overall computational efficiency.

IV. EXPERIMENTAL RESULTS

In this section, we conduct experiments over two public benchmarks to evaluate our problem and the capability of the proposed model. We first compare our model against the previous state-of-the-art methods and then perform ablation studies to delve into our model structure.

A. Datasets

Inria Aerial [49]. This dataset contains large resolution aerial images of five cities, ranging from dense metropolitan districts to alpine resorts. It contains 180 aerial images of 5000×5000 pixels with the binary mask for building/non-building areas. Following the protocol of [10], we split images into training, validation, and testing sets with 126, 27, and 27 images, respectively.

DeepGlobe [50]. This dataset provides 803 ultra-high resolution aerial images with 2448×2448 pixels. It contains seven classes of landscape regions, including urban, agriculture, rangeland, forest, water, barren, and unknown region not considered in the challenge. We split all images following [10], i.e., 454 training images, 207 validation images, and 142 testing images.

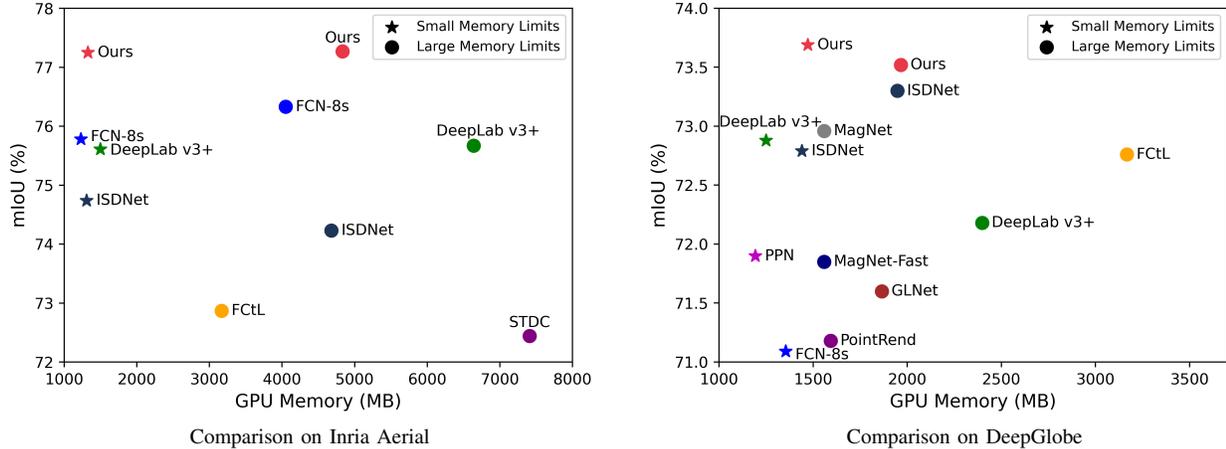


Fig. 2. Comparison of mIoU v.s. GPU Memory cost under small and large GPU memory constraints (denoted as stars and circles).

B. Implementation Details

Training details. We implement our framework using the mmSegmentation [51] toolbox on a workstation with a single NVIDIA RTX 3090 GPU. In particular, we adopt DeepLabv3 [52] with ResNet18 [16] as the encoder of the semantic branch and STDC [53] as the encoder of the memory-based interaction branch. During training, we optimize the parameters adapting Stochastic Gradient Descent (SGD) and set the batch size to 4 and 8 for Inria Aerial and DeepGlobe, respectively. The initial learning rate is set to 10^{-2} , which is decayed by a poly learning rate policy with the power of 0.9. In practice, it takes 40k and 80k iterations to converge our model for two datasets, respectively.

Inference details. During inference, we measure the GPU Memory and Frames-per-second (FPS) on an RTX 2080Ti GPU and adopt the same environment as [28] (i.e., CUDA 10.1, CuDNN 7.6.5, and Pytorch 1.6.0) for fair comparison.

C. Comparison with State-of-the-arts

In practice, segmenting an ultra-high resolution image often consumes an exceedingly large amount of GPU memory resources. As a consequence, for the robotic platforms like UAVs, they usually have limited computation resources and GPU memory for segmenting and analyzing the large images on board. To simulate the challenging situations, in the experiments, we compare the segmentation methods under small and large GPU memory limits. In this experiment, we set 1.5 GB as the small GPU memory limit. For large GPU memory limits, we set 7.5 GB for Inria Aerial and 3.5 GB for DeepGlobe, due to their different image resolutions. In the following, we conduct the comparison experiments and discuss the results.

Small memory limits. In general, there are two ways to segment the ultra-high resolution images: 1) segmenting downsampled global images (denoted as *global inference*); and 2) cropping, segmenting, and merging local patches (denoted as *local inference*). We first consider the methods that adopt the global inference strategy, which straightforwardly down-samples the whole input image I to meet the requirements of

small limited GPU memory. To this end, we retrain and test U-Net [54], FCN-8s [15], and DeepLab v3+ [19] on Inria Aerial and DeepGlobe datasets. As shown in Table I and Table II, these methods can hardly achieve satisfactory accuracy despite their high FPS, since compressing large images causes severe detail lost.

Thus, we adopt the local inference strategy for the segmentation methods to meet small GPU memory limits. Table. I and Table. II show our superior performance than other methods in the respective datasets. In specific, our model shows similar GPU memory usage and running time comparing to the state-of-the-art ISDNet [28], but elevates 2.50% mIoU on Inria Aerial and 0.87% mIoU on DeepGlobe. It is worth noting that our model achieves a significant advantage on Inria Aerial dataset, because our proposed spatial-guided high-resolution query module can depict delicate objects and this dataset contains a large number of small buildings. Besides, as a ResNet-50 based model, PPN [27] is also able to fulfill the requirement of limiting small GPU memory, but it underperforms our method and ISDNet. In Fig. 4, we illustrate several qualitative comparison results. We observe that our model has the ability to delineate fine regions (e.g., small urban and rivers) and stronger semantic discrimination (e.g., large forests and rangeland). This is attributed to the high-resolution query and memory-based semantic enhancement.

Large memory limits. To evaluate the capability of our model under large GPU memory limits, we compare our model with BiSeNetV1 [55], STDC [53], PointRend [56], CascadePSP [29], GLNet [10], MagNet [11], FcTL [12], and ISDNet [28], in terms of mIoU, GPU Memory, and FPS. The comparison results are depicted in Table. III and Table. IV, in which we follow most of the results reported by [28]. As observed, our model is comparable to the state-of-the-art ISDNet on DeepGlobe, but we gain at least 3% improvement than it on Inria Aerial. This shows that our approach can further trade off performance, GPU memory, and speed.

Performance-memory trade-off. We quantify and visualize all the state-of-the-art methods based on the small and large GPU memory limits in Fig. 2. We can see that our

TABLE I

COMPARISON WITH STATE-OF-THE-ARTS WITH SMALL GPU MEMORY LIMITS ON INRIA AERIAL. * REPRESENTS OUR IMPLEMENTATION.

Model	Backbone	Inference	mIoU	Memory	FPS
FCN-8s* [15]	ResNet-18	Global	47.74	1412	12.72
DeepLab v3+* [19]	ResNet-18	Global	34.73	1532	16.58
FCN-8s* [15]	ResNet-18	Local	75.58	1228	1.23
DeepLab v3+* [19]	ResNet-18	Local	76.22	1496	1.10
ISDNet* [28]	ResNet-18	Local	74.75	1306	3.99
Ours	ResNet-18	Local	77.25	1324	3.55

TABLE II

COMPARISON WITH STATE-OF-THE-ARTS WITH SMALL GPU MEMORY LIMITS ON DEEPGLOBE. * REPRESENTS OUR IMPLEMENTATION.

Model	Backbone	Inference	mIoU	Memory	FPS
U-Net* [54]	U-Net	Global	20.61	1506	42.16
FCN-8s* [15]	ResNet-18	Global	60.41	1438	12.07
DeepLab v3+* [19]	ResNet-18	Global	52.43	1532	16.58
U-Net* [54]	U-Net	Local	69.72	1426	0.65
FCN-8s* [15]	ResNet-18	Local	71.09	1354	3.47
DeepLab v3+* [19]	ResNet-18	Local	72.53	1250	3.06
PPN [27]	ResNet-50	Local	71.90	1193	12.90
ISDNet* [28]	ResNet-18	Local	72.79	1440	11.47
Ours	ResNet-18	Local	73.66	1472	10.09

model can achieve optimal mIoU regardless of small or large GPU memory limits. More importantly, our model can obtain similar performances under small and large GPU memory limits (77.25% v.s. 77.27% on Inria Aerial and 73.66% v.s. 73.50% on DeepGlobe), indicating that our robustness to GPU memory constraints, whereas other methods are greatly affected due to narrow view of local patches. This still benefits from our query module, which predicts pixel categories only by the corresponding relative coordinates and the nearest latent embeddings.

D. Ablation Study

In the following, we conduct ablation studies on the two proposed modules. First, we demonstrate the effectiveness of spatial-guided high-resolution information and memory-based interaction scheme in our model. Next, we analyze different strategies for updating and reading the memory bank in our scheme. Last, we investigate the boundary case of GPU memory usage.

Model structure. First, we take the naive query module as our baseline, as shown in Table V, our model just costs 1284MB GPU memory with high FPS, but only gains 69.24% mIoU. To improve performance, we use STDC [53] to extract high-resolution spatial information M_b . It is able to boost the performance from 69.24% to 72.06%, which exceeds bilinear interpolation with 71.81% mIoU and demonstrates the effectiveness of the spatial mask guidance. Besides, the GPU memory increase is only 186MB and the total GPU memory is still less than 1.5GB. Based on this, we add memory bank \mathcal{M} and higher resolution information M_l , separately. It is observed that there is a great performance improvement, especially the guide of spatial information M_l (at least 1.1% boost). Importantly, the GPU memory and time hardly grow (about 1MB and less than 0.5 FPS) because both are linear operations. Next, with the aid of both, our model can achieve the best result (i.e., 73.66% mIoU and 1472MB GPU Memory). Besides, we

TABLE III

COMPARISON WITH STATE-OF-THE-ARTS WITH LARGE GPU MEMORY LIMITS ON INRIA AERIAL. * REPRESENTS OUR IMPLEMENTATION.

Model	Backbone	Inference	mIoU	Memory	FPS
FCN-8s* [15]	ResNet-18	Global	75.67	4050	1.96
DeepLab v3+* [19]	ResNet-18	Global	76.33	6638	1.75
STDC [53]	STDC	Global	72.44	7410	4.97
CascadePSP [29]	ResNet-50	Local	69.40	3236	0.03
GLNet [10]	ResNet-50	Local	71.20	2663	0.05
FCtL [12]	VGG-16	Local	72.87	3167	0.04
ISDNet [28]	ResNet-18	Global	74.23	4680	6.90
Ours	ResNet-18	Global	77.27	4834	5.53

TABLE IV

COMPARISON WITH STATE-OF-THE-ARTS WITH LARGE GPU MEMORY LIMITS ON DEEPGLOBE. * REPRESENTS OUR IMPLEMENTATION.

Model	Backbone	Inference	mIoU	Memory	FPS
U-Net* [54]	U-Net	Global	28.53	3511	9.34
FCN-8s* [15]	ResNet-18	Global	68.67	1890	7.98
DeepLab v3+* [19]	ResNet-18	Global	72.18	2398	7.22
BiSeNetV1 [55]	ResNet-18	Global	53.00	1801	14.20
STDC [53]	STDC	Global	70.30	2580	14.00
PointRend [56]	ResNet-50	Global	71.78	1593	6.25
CascadePSP [29]	ResNet-50	Local	68.50	3236	0.11
GLNet [10]	ResNet-50	Local	71.60	1865	0.17
MagNet-Fast [11]	ResNet-50	Local	71.85	1559	3.40
MagNet [11]	ResNet-50	Local	72.96	1559	0.80
FCtL [12]	VGG-16	Local	72.76	3167	0.13
ISDNet [28]	ResNet-18	Global	73.30	1948	27.70
Ours	ResNet-18	Global	73.50	1966	24.33

also try to enhance the semantics of higher resolution masks M_l using a new memory bank \mathcal{M}_l , but the result gets worse, we think that this mask mainly contains spatial details and the semantics are too weak to improve. Particularly, we can also adopt high-resolution features as the spatial information guide, that is implicit feature alignment (IFA) [41]. Table V shows IFA is inferior to our method by at least 1.4% mIoU and uses 270MB more GPU memory than ours. This largely reflects that our method is more effective and efficient.

Memory-based interaction scheme. In our memory scheme, we use semantic branch features to update the semantic memory, dubbed cross-branch interaction strategy. To demonstrate the necessity of this strategy, we perform ablation experiments with regard to it. Table VI shows performance barely changes (73.23% v.s. 73.25%) using the features of the memory-based branch, compared with no memory bank. For update mode, we ablate the effect of "Mean" mode (i.e., $\hat{\mathcal{R}}^c = \sum_{i=1}^{N^c} \frac{\mathcal{R}_i^c}{N^c}$). As shown in Table VI, this is only 0.06% mIoU improvement over no memory bank. Besides, we replace Concat instead of Softmax (i.e., $M_b = \text{Concat}(\frac{M^T \otimes \mathcal{F}_b}{\sqrt{D}}, M_b)$) for the reading mode, it leads to 0.08% mIoU degradation.

Memory usage boundary. To investigate the GPU memory usage boundary, we attempt to reduce the size of local patches in order to further decrease GPU memory consumption during the inference process. We perform this experiment on Inria Aerial with 5000x5000 pixel images. As shown in Fig. 3, our model and ISDNet achieve similar minimum GPU memory usage of approximately 940MB. However, even under the lowest-bound of GPU memory constraint, our model outperforms ISDNet by a significant margin, with comparatively less degradation in performance. In general, our model demonstrates superior performance over ISDNet, both under small and large GPU memory constraints.

TABLE V

EFFECTIVENESS OF OUR PROPOSED MODEL STRUCTURE.

Up-sampling	Spa. Inf.	Mem.	mIoU	Memory	FPS
Bilinear	-	-	71.81	1596	7.99
IFA [41]	-	-	72.23	1742	8.80
Ours	-	-	69.24	1284	19.51
Ours	M_b	-	72.06	1470	10.56
Ours	M_b	\mathcal{M}	72.53	1471	10.24
Ours	M_b+M_l	-	73.23	1471	10.13
Ours	M_b+M_l	\mathcal{M}	73.66	1472	10.09
Ours	M_b+M_l	$\mathcal{M}+\mathcal{M}_l$	73.09	1474	9.31

TABLE VI

DIFFERENT STRATEGIES FOR MEMORY-BASED INTERACTION SCHEME.

\mathcal{M}	Cross-Branch	Update	Read	mIoU
×	×	×	×	73.23
✓	×	Cosine	Softmax	73.25
✓	✓	Mean	Softmax	73.29
✓	✓	Cosine	Concat	73.15
✓	✓	Cosine	Softmax	73.66

V. CONCLUSION AND LIMITATIONS

In this paper, we propose an effective solution segmenting an ultra-high resolution image in a limited GPU memory system, which has practical value for robotic systems. In particular, we propose a spatial-guided high-resolution query module for local inference. Additionally, we also present a memory-based interaction scheme that efficiently enhances the semantics of high-resolution information by bridging cross-image contextual information. There are several limitations in our method. First, it can hardly eliminate the noise from high-resolution information, leading to bias in the query process. Moreover, how to deploy the memory-efficient model while preserving high FPS remains a big challenge.

REFERENCES

- [1] H. Goforth and S. Lucey, "Gps-denied uav localization using pre-existing satellite imagery," in *ICRA*. IEEE, 2019, pp. 2974–2980.
- [2] Y. Zheng, Z. Chen, D. Lv, Z. Li, Z. Lan, and S. Zhao, "Air-to-air visual detection of micro-uavs: An experimental evaluation of deep learning," *IEEE Robotics and automation letters*, vol. 6, no. 2, pp. 1020–1027, 2021.
- [3] J. Valente, L. Kooistra, and S. Mucher, "Fast classification of large germinated fields via high-resolution uav imagery," *IEEE Robotics and Automation Letters*, vol. 4, no. 4, pp. 3216–3223, 2019.
- [4] Z. Li and V. Isler, "Large scale image mosaic construction for agricultural applications," *IEEE Robotics and Automation Letters*, vol. 1, no. 1, pp. 295–302, 2016.
- [5] B. Ferrarini, M. J. Milford, K. D. McDonald-Maier, and S. Ehsan, "Binary neural networks for memory-efficient and effective visual place recognition in changing environments," *IEEE Transactions on Robotics*, vol. 38, no. 4, pp. 2617–2631, 2022.
- [6] L. Grossman and B. Plancher, "Just round: Quantized observation spaces enable memory efficient learning of dynamic locomotion," in *ICRA*. IEEE, 2023, pp. 3002–3007.
- [7] S. Gomez-Gonzalez, S. Prokudin, B. Scholkopf, and J. Peters, "Real time trajectory prediction using deep conditional generative models," *IEEE Robotics and Automation Letters*, vol. 5, no. 2, pp. 970–976, 2020.
- [8] P. Hu, F. Perazzi, F. C. Heilbron, O. Wang, Z. Lin, K. Saenko, and S. Sclaroff, "Real-time semantic segmentation with fast attention," *IEEE Robotics and Automation Letters*, vol. 6, no. 1, pp. 263–270, 2020.
- [9] Y. Sun, B. Pan, and Y. Fu, "Lightweight deep neural network for real-time instrument semantic segmentation in robot assisted minimally invasive surgery," *IEEE Robotics and Automation Letters*, vol. 6, no. 2, pp. 3870–3877, 2021.
- [10] W. Chen, Z. Jiang, Z. Wang, K. Cui, and X. Qian, "Collaborative global-local networks for memory-efficient segmentation of ultra-high resolution images," in *CVPR*, 2019, pp. 8924–8933.
- [11] C. Huynh, A. T. Tran, K. Luu, and M. Hoai, "Progressive semantic segmentation," in *CVPR*, 2021, pp. 16755–16764.

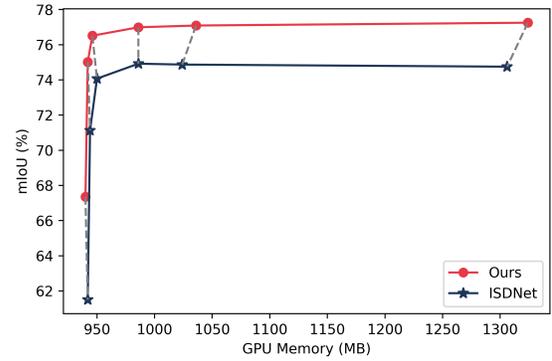


Fig. 3. By continuously reducing the size of local patches, we explore the minimum GPU memory cost of ISDNet and our approach pertaining to their corresponding segmentation performance.

- [12] Q. Li, W. Yang, W. Liu, Y. Yu, and S. He, "From contexts to locality: Ultra-high resolution image segmentation via locality-aware contextual correlation," in *ICCV*, 2021, pp. 7252–7261.
- [13] Y. Chen, S. Liu, and X. Wang, "Learning continuous image representation with local implicit image function," in *CVPR*, 2021, pp. 8628–8638.
- [14] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.
- [15] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *CVPR*, 2015, pp. 3431–3440.
- [16] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *CVPR*, 2016, pp. 770–778.
- [17] H. Zhao, J. Shi, X. Qi, X. Wang, and J. Jia, "Pyramid scene parsing network," in *CVPR*, 2017, pp. 2881–2890.
- [18] V. Badrinarayanan, A. Kendall, and R. Cipolla, "Segnet: A deep convolutional encoder-decoder architecture for image segmentation," *TPAMI*, vol. 39, no. 12, pp. 2481–2495, 2017.
- [19] L.-C. Chen, Y. Zhu, G. Papandreou, F. Schroff, and H. Adam, "Encoder-decoder with atrous separable convolution for semantic image segmentation," in *ECCV*, 2018, pp. 801–818.
- [20] J. Wang, K. Sun, T. Cheng, B. Jiang, C. Deng, Y. Zhao, D. Liu, Y. Mu, M. Tan, X. Wang *et al.*, "Deep high-resolution representation learning for visual recognition," *TPAMI*, vol. 43, no. 10, pp. 3349–3364, 2020.
- [21] S. Zheng, J. Lu, H. Zhao, X. Zhu, Z. Luo, Y. Wang, Y. Fu, J. Feng, T. Xiang, P. H. Torr *et al.*, "Rethinking semantic segmentation from a sequence-to-sequence perspective with transformers," in *CVPR*, 2021, pp. 6881–6890.
- [22] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo, "Swin transformer: Hierarchical vision transformer using shifted windows," in *ICCV*, 2021, pp. 10 012–10 022.
- [23] E. Xie, W. Wang, Z. Yu, A. Anandkumar, J. M. Alvarez, and P. Luo, "Segformer: Simple and efficient design for semantic segmentation with transformers," *NeurIPS*, vol. 34, pp. 12 077–12 090, 2021.
- [24] R. Strudel, R. Garcia, I. Laptev, and C. Schmid, "Segformer: Transformer for semantic segmentation," in *ICCV*, 2021, pp. 7262–7272.
- [25] B. Cheng, A. Schwing, and A. Kirillov, "Per-pixel classification is not all you need for semantic segmentation," *NeurIPS*, vol. 34, pp. 17 864–17 875, 2021.
- [26] B. Cheng, I. Misra, A. G. Schwing, A. Kirillov, and R. Girdhar, "Masked-attention mask transformer for universal image segmentation," in *CVPR*, 2022, pp. 1290–1299.
- [27] T. Wu, Z. Lei, B. Lin, C. Li, Y. Qu, and Y. Xie, "Patch proposal network for fast semantic segmentation of high-resolution images," in *AAAI*, vol. 34, no. 07, 2020, pp. 12 402–12 409.
- [28] S. Guo, L. Liu, Z. Gan, Y. Wang, W. Zhang, C. Wang, G. Jiang, W. Zhang, R. Yi, L. Ma *et al.*, "Isdnet: Integrating shallow and deep networks for efficient ultra-high resolution segmentation," in *CVPR*, 2022, pp. 4361–4370.
- [29] H. K. Cheng, J. Chung, Y.-W. Tai, and C.-K. Tang, "Cascadepsp: Toward class-agnostic and very high-resolution segmentation via global and local refinement," in *CVPR*, 2020, pp. 8890–8899.
- [30] T. Shen, Y. Zhang, L. Qi, J. Kuen, X. Xie, J. Wu, Z. Lin, and J. Jia, "High quality segmentation for ultra high-resolution images," in *CVPR*, 2022, pp. 1310–1319.
- [31] J. J. Park, P. Florence, J. Straub, R. Newcombe, and S. Lovegrove, "Deepsdf: Learning continuous signed distance functions for shape representation," in *CVPR*, 2019, pp. 165–174.

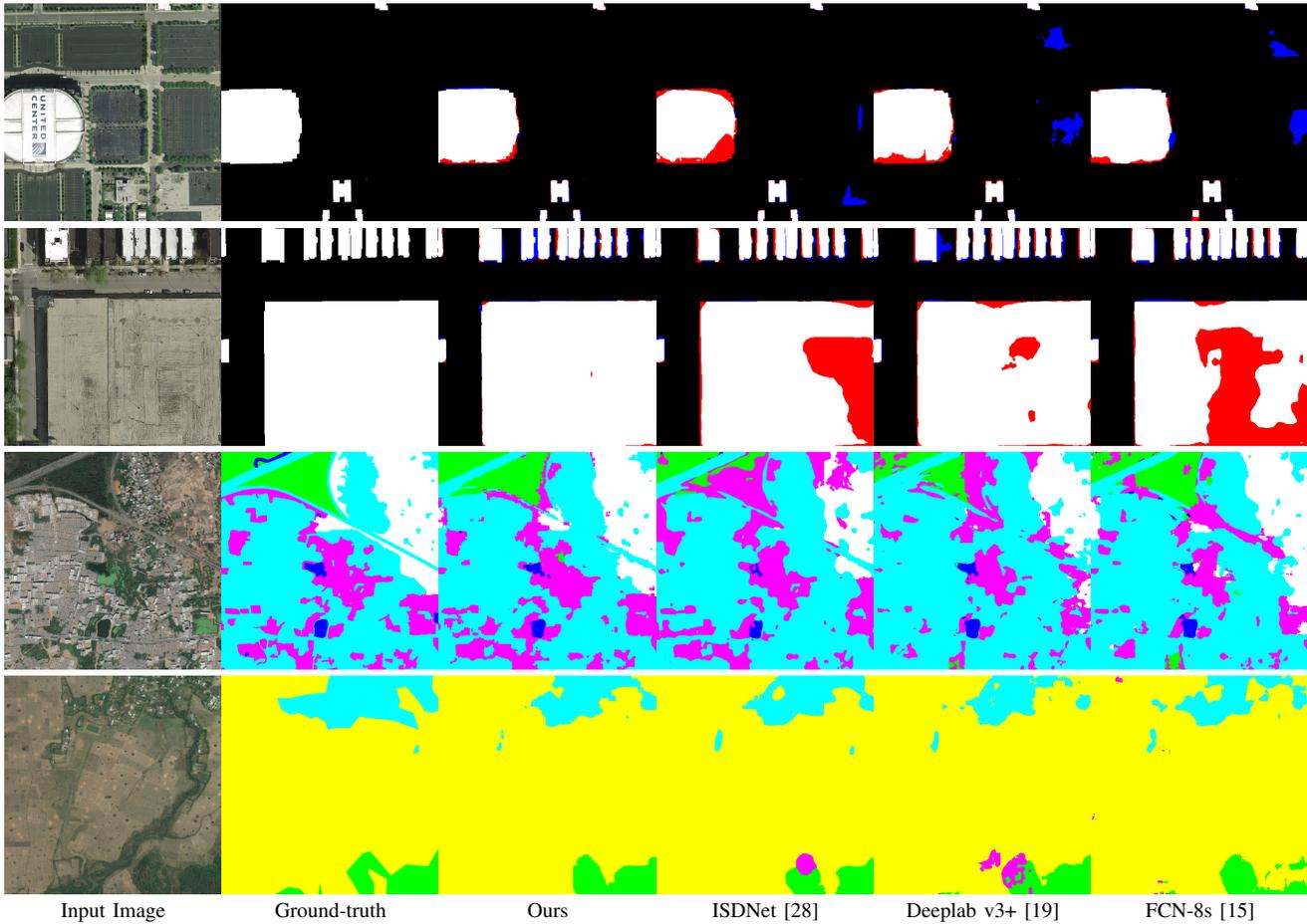


Fig. 4. On the first two rows, we show the representative comparison results from Inria Aerial, where white regions represent the estimated regions for buildings. We highlight the pixels which are the discrepancy between our estimation and ground truth. The blue and red pixels represent False Negatives and False Positives, respectively. On the last two rows, we show representative results from DeepGlobe, where cyan represents "urban", yellow represents "agriculture", purple represents "rangeland", green represents "forest", blue represents "water", and white represents "barren".

[32] V. Sitzmann, M. Zollhöfer, and G. Wetzstein, "Scene representation networks: Continuous 3d-structure-aware neural scene representations," *NeurIPS*, vol. 32, 2019.

[33] Z. Chen and H. Zhang, "Learning implicit fields for generative shape modeling," in *CVPR*, 2019, pp. 5939–5948.

[34] M. Niemeyer, L. Mescheder, M. Oechsle, and A. Geiger, "Differentiable volumetric rendering: Learning implicit 3d representations without 3d supervision," in *CVPR*, 2020, pp. 3504–3515.

[35] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng, "Nerf: Representing scenes as neural radiance fields for view synthesis," in *ECCV*, 2020.

[36] C. Jiang, A. Sud, A. Makadia, J. Huang, M. Nießner, T. Funkhouser *et al.*, "Local implicit grid representations for 3d scenes," in *CVPR*, 2020, pp. 6001–6010.

[37] H. Chen, B. He, H. Wang, Y. Ren, S. N. Lim, and A. Shrivastava, "Nerv: Neural representations for videos," *NeurIPS*, vol. 34, pp. 21 557–21 568, 2021.

[38] Z. Li, M. Wang, H. Pi, K. Xu, J. Mei, and Y. Liu, "E-nerf: Expedite neural video representation with disentangled spatial-temporal context," in *ECCV*. Springer, 2022, pp. 267–284.

[39] W. Shangguan, Y. Sun, W. Gan, and U. S. Kamilov, "Learning cross-video neural representations for high-quality frame interpolation," in *ECCV*. Springer, 2022, pp. 511–528.

[40] X. Xu, Z. Wang, and H. Shi, "Ultrasr: Spatial encoding is a missing key for implicit image function-based arbitrary-scale super-resolution," *arXiv preprint arXiv:2103.12716*, 2021.

[41] H. Hu, Y. Chen, J. Xu, S. Borse, H. Cai, F. Porikli, and X. Wang, "Learning implicit feature alignment function for semantic segmentation," in *ECCV*. Springer, 2022, pp. 487–505.

[42] S. W. Oh, J.-Y. Lee, N. Xu, and S. J. Kim, "Video object segmentation using space-time memory networks," in *ICCV*, 2019, pp. 9226–9235.

[43] X. Wang, H. Zhang, W. Huang, and M. R. Scott, "Cross-batch memory for embedding learning," in *CVPR*, 2020, pp. 6388–6397.

[44] Y. Chen, Y. Cao, H. Hu, and L. Wang, "Memory enhanced global-local aggregation for video object detection," in *CVPR*, 2020, pp. 10 337–10 346.

[45] L. Hu, P. Zhang, B. Zhang, P. Pan, Y. Xu, and R. Jin, "Learning position and target consistency for memory-based video object segmentation," in *CVPR*, 2021, pp. 4144–4154.

[46] H. Xie, H. Yao, S. Zhou, S. Zhang, and W. Sun, "Efficient regional memory network for video object segmentation," in *CVPR*, 2021, pp. 1286–1295.

[47] J. Kim, J. Lee, J. Park, D. Min, and K. Sohn, "Pin the memory: Learning to generalize semantic segmentation," in *CVPR*, 2022, pp. 4350–4360.

[48] Z. Jin, D. Yu, Z. Yuan, and L. Yu, "Mcibi++: Soft mining contextual information beyond image for semantic segmentation," *TPAMI*, 2022.

[49] E. Maggiori, Y. Tarabalka, G. Charpiat, and P. Alliez, "Can semantic labeling methods generalize to any city? the inria aerial image labeling benchmark," in *IGARSS*. IEEE, 2017, pp. 3226–3229.

[50] 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 *CVPRW*, 2018, pp. 172–181.

[51] M. Contributors, "Mmsegmentation: Openmmlab semantic segmentation toolbox and benchmark," 2020.

[52] L.-C. Chen, G. Papandreou, F. Schroff, and H. Adam, "Rethinking atrous convolution for semantic image segmentation," *arXiv preprint arXiv:1706.05587*, 2017.

[53] M. Fan, S. Lai, J. Huang, X. Wei, Z. Chai, J. Luo, and X. Wei, "Rethinking bisenet for real-time semantic segmentation," in *CVPR*, 2021, pp. 9716–9725.

[54] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *MICCAI*. Springer, 2015, pp. 234–241.

[55] C. Yu, J. Wang, C. Peng, C. Gao, G. Yu, and N. Sang, "Bisenet: Bilateral segmentation network for real-time semantic segmentation," in *ECCV*, 2018, pp. 325–341.

[56] A. Kirillov, Y. Wu, K. He, and R. Girshick, "Pointrend: Image segmentation as rendering," in *CVPR*, 2020, pp. 9799–9808.