

SpecDETR: A Transformer-based Hyperspectral Point Object Detection Network

Zhaoxu Li, Wei An, Gaowei Guo, Longguang Wang, Yingqian Wang, and Zaiping Lin

arXiv:2405.10148v3 [cs.CV] 25 May 2025

Abstract—Hyperspectral target detection (HTD) aims to identify specific materials based on spectral information in hyperspectral imagery and can detect extremely small-sized objects, some of which occupy a smaller than one-pixel area. However, existing HTD methods are developed based on per-pixel binary classification, neglecting the three-dimensional cube structure of hyperspectral images (HSIs) that integrates both spatial and spectral dimensions. The synergistic existence of spatial and spectral features in HSIs enable objects to simultaneously exhibit both, yet the per-pixel HTD framework limits the joint expression of these features. In this paper, we rethink HTD from the perspective of spatial-spectral synergistic representation and propose hyperspectral point object detection as an innovative task framework. We introduce SpecDETR, the first specialized network for hyperspectral multi-class point object detection, which eliminates dependence on pre-trained backbone networks commonly required by vision-based object detectors. SpecDETR uses a multi-layer Transformer encoder with self-excited subpixel-scale attention modules to directly extract deep spatial-spectral joint features from hyperspectral cubes. During feature extraction, we introduce a self-excited mechanism to enhance object features through self-excited amplification, thereby accelerating network convergence. Additionally, SpecDETR regards point object detection as a one-to-many set prediction problem, thereby achieving a concise and efficient DETR decoder that surpasses the state-of-the-art (SOTA) DETR decoder. We develop a simulated hyperSpectral Point Object Detection benchmark termed SPOD, and for the first time, evaluate and compare the performance of visual object detection networks and HTD methods on hyperspectral point object detection. Extensive experiments demonstrate that our proposed SpecDETR outperforms SOTA visual object detection networks and HTD methods. Our code and dataset are available at <https://github.com/ZhaoxuLi123/SpecDETR>.

Index Terms—Hyperspectral target detection, point object detection, Detection Transformer.

I. INTRODUCTION

Hyperspectral imagery (HSI) has gained significant attention in Earth observation due to its rich spectral information in recent decades. HSI usually provides hundreds of narrow spectral bands, allowing for the land cover differentiation at the pixel level. With advancements in hyperspectral classification [1], [2], unmixing [3], anomaly detection [4]–[6], change detection [7], biomass estimation [8], [9] and target detection [10], [11], HSI has assumed a crucial role in various remote

Zhaoxu Li, Wei An, Gaowei Guo, Yingqian Wang, and Zaiping Lin are with the College of Electronic Science and Technology, National University of Defense Technology, Changsha 410073, China. Longguang Wang is with the Aviation University of Air Force, Changchun 130010, China.(email: lizhaoxu@nudt.edu.cn; anwei@nudt.edu.cn; guogaowei22@nudt.edu.cn; wangyingqian16@nudt.edu.cn; linzaiping@nudt.edu.cn; wanglongguang15@nudt.edu.cn)

Corresponding author: Wei An

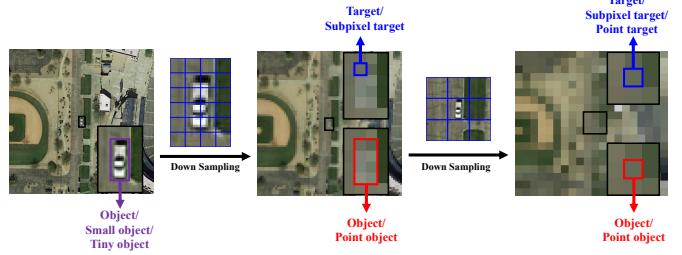


Fig. 1. Illustration of the concept of point object.

sensing applications such as agricultural survey, mineral exploration, and greenhouse gas detection. Hyperspectral target detection (HTD), which utilizes prior spectral information to locate specific materials, can identify extremely small objects that are undetectable in conventional optical images, with some objects areas being even smaller than a single pixel area.

HTD is generally regarded as a per-pixel binary classification problem based on prior spectral information, where each pixel in the HSI is categorized as either a target (containing the target spectrum) or background (lacking the target spectrum) based on known target spectra. In recent years, significant advancements have been made in the field of HTD. Feng et al. [12] introduced a cross-domain few-shot learning method tailored for HTD tasks, addressing the challenge of scarce target prior information. Chang [13] delved into the mathematical principles of Constrained Energy Minimization (CEM) widely used in sub-pixel HTD and proposed several novel CEM generalizations. These HSI methods treat the 3D hyperspectral cube as independent spectral vectors. Notably, some HTD approaches incorporate spatial characteristics. For instance, Feng et al. [14] proposed a coarse-to-fine HTD algorithm based on low-rank tensor decomposition, leveraging the spatial-temporal properties of HSI to extract pure background information. But these HTD methods still operate within a per-pixel classification framework, utilizing only the prior spatial information of targets, and fail to fully represent the spatial-spectral joint features of objects. The unique attribute of hyperspectral images, which integrate spatial and spectral features, allows objects to exhibit both simultaneously. We believe that such a per-pixel HTD framework can restrict the representation of instance-level spatial-spectral joint features of objects, leading to limitations in detection capability. On one hand, the current HTD framework primarily utilize prior target spectral information, often neglecting the extraction and utilization of instance-level object spatial characteristics. On the other hand, existing HTD methods typically output detection scores for

each pixel, failing to directly provide instance-level prediction results. The success of visual object detection networks such as Faster R-CNN [15], YOLO [16], and Detection Transformer (DETR) [17] offers a new perspective. Compared to the per-pixel binary classification framework in HTD, visual object detection networks can directly yield instance-level category and location predictions. Furthermore, HTD relies exclusively on prior spectra and test images for detection, while visual object detection networks can effectively learn robust, high-level, instance-level features from large-scale training images annotated with instance-level object information. During training, the learning frameworks of visual object detection networks can focus on difficult-to-classify samples, such as easily confused objects categories and background regions similar to objects. Therefore, in this paper, we introduce visual object detection into hyperspectral target detection and expand hyperspectral target detection to hyperspectral point object detection. Considering the differences between HTD and visual object detection, we will redefine the concepts of target and object, providing a clear definition of point object. Furthermore, we will elucidate the relationship between spatial information and spectral information, thereby clarifying the motivation behind our proposed framework for point object detection.

We have observed differences in similar concepts between the fields of visual object detection and HTD. In Fig. 1, we illustrate the distinctions among the object, the target in HTD, and the point object of focus in our paper using a vehicle object as an example. Object detection aims to locate and classify all pixels belonging to a single instance as a whole. In the object detection field, small objects or tiny objects are typically defined by the pixel number and can still be localized and classified based on their spatial information. For instance, the COCO dataset [18] and the SODA dataset [19] consider objects with a bounding box (bbox) area equal to or smaller than 1024 pixels as small objects, which still contain certain texture and shape detail information. In the field of HTD, most literature defines a target as a pixel containing spectral features of interest, and a subpixel target as a pixel whose target spectral proportion, also termed as target abundance, is less than 1. Some literature [20]–[22] specifically uses the terms point target or subpixel target to refer to instances that occupy only a single pixel. However, these works still employ the per-pixel classification framework and pixel-level evaluation metrics, essentially adhering to the pixel concept. We focus on shifting from pixel-based to instance-based concept, treating target pixels belonging to the same instance as a collective unit. We refer to these instances, which lack spatial detail features and are composed of very few pixels, as point objects. The definition of point objects is not based on pixel number, owing to the presence of many mixed pixels within their composition. Any objects whose sizes are close to or smaller than the spatial resolution or ground sample distance (GSD) of the image, and whose shape information alone is insufficient to support their classification, can be considered as point objects.

In the field of computer vision (CV), researchers typically employ feature extraction networks to extract spatial features as shape and texture from RGB images for visual downstream

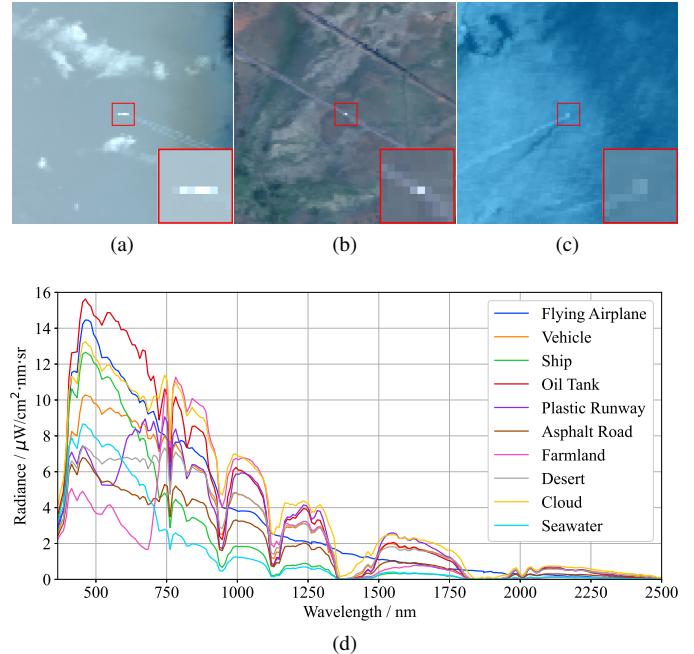


Fig. 2. Real-world point objects in HSI collected by AVIRIS. (a) Flying airplane. (b) Vehicle. (c) Ship. (d) Spectral radiance curves of the shown point objects and other ground objects.

tasks such as classification, object detection, and instance segmentation. In contrast, hyperspectral imaging systems generate data cubes containing hundreds of spectral bands, prompting HTD methods to prioritize spectral information utilization. However, hyperspectral images are inherently spatial-spectral integrated, forming a three-dimensional cube that encompasses both spatial and spectral dimensions. Fig. 2 illustrates the pseudo-color images and spectral curves of three real-world point objects: an airplane in flight, a car on the road, and a ship on the sea surface. The limited spatial resolution results in deficient shape and texture details for these objects, as shown in Fig. 2(a-c). While spatial information permits coarse localization of these point objects, it proves insufficient for reliable classification. Fig. 2(d) presents the spectral curves of the three point objects and other common ground objects, revealing distinct differences among their spectral signatures due to physical mechanisms. For instance, the shown airplane, owing to its high altitude and proximity to the sensor, exhibits less spectral radiation attenuation in the two atmospheric absorption bands at 1.4 μm and 1.8 μm , resulting in significantly higher radiation values at these bands compared to those of other curves. This observation motivates our investigation: Can we effectively extract both spatial and spectral features of point objects? We believe that current object detection frameworks inherently can represent spatial-spectral joint features. The spatial features manifested in single channel imagery or RGB imagery fundamentally arise from spectral radiance variations across pixel locations. As shown in Fig. 2(a-c), the observed spatial characteristics essentially represent a limited spatial-spectral joint representation constrained by limited spectral dimensionality. From this perspective, current visual detection networks already demonstrate foundational capacity for

spatial-spectral joint feature learning from hyperspectral cube data. Nevertheless, employing the object detection framework to achieve end-to-end localization and classification of subpixel-level point objects remains, to our knowledge, an unexplored challenge—the focal point of this paper.

Inspired by the current DETR-like detectors, we propose the first specialized network for hyperspectral point object detection, SpecDETR. We treat the spectral features of each pixel in the HSI as a token serving as the fundamental operation unit in Transformer [23]. In existing visual object detection networks, the pretrained feature extraction network, referred to as the backbone, is used to extract deep features from input RGB images, which are then fed into the detection head for object localization and classification. Unlike these visual object detection networks, we bypass the backbone and directly utilize a Transformer encoder to extract deep features from spectral tokens. Building upon the deformable attention module [24], we propose a self-excited subpixel-scale attention module (S2A module) for the extraction of spatial-spectral joint features of point objects, which efficiently achieves subpixel-scale deformable sampling while enabling self-excited amplification of object features. Furthermore, based on the DINO [25] decoder, we design a concise decoder framework tailored for point object detection, integrating one-to-many label assignment and the end-to-end advantages of DETR. With the same attention operations and decoder layer number, our decoder achieves superior point object detection accuracy with fewer parameters than the DINO decoder. To advance the development of hyperspectral point object detection, we developed a hyperSpectral multi-class **Point Object Detection** dataset, termed **SPOD**. Using the SPOD dataset, we thoroughly evaluated the mainstream visual object detection networks and HTD methods. Our SpecDETR outperforms state-of-the-art (SOTA) visual object detection networks and HTD methods on the SPOD dataset. The contributions of this paper can be summarized as follows:

- 1) We extend hyperspectral target detection to hyperspectral point object detection and propose the first specialized hyperspectral point object detection network, SpecDETR.
- 2) We regard SpecDETR as a one-to-many set prediction problem and develop a concise and efficient DETR decoder framework that surpasses the current SOTA DETR decoder framework on point object detection.
- 3) We develop a hyperspectral point object detection benchmark and evaluate visual object detection networks and HTD methods on the hyperspectral point object detection task for the first time.

The rest of this paper is organized as follows. Section II reviews the fields related to our work. Section III describes the details of SpecDETR. The datasets and experiments are presented in Section IV and Section V, respectively, followed by the conclusion in Section VI.

II. RELATED WORK

A. Hyperspectral Target Detection

HTD aims to locate targets based on the spectral characteristics in HSIs [26]. Classic HTD methods can be categorized into matching methods [27], statistical analysis methods [28]–[30], and subspace-based methods [31]–[33]. While these methods demonstrate effective target-background separation in homogeneous scenarios, they have limitations in handling complex scenarios. In response to these limitations, machine learning-based techniques such as sparse representation [34]–[40] and kernel methods [41]–[46] have gained popularity due to their ability to handle complex data and achieve impressive performance. Recently, many deep learning methods have been specifically designed to address the characteristics of HTD. Traditional HTD typically provide only a few, or even a single, prior spectral curve along with a test image, which results in a small sample number and may lead to network overfitting. Some methods [47]–[49] employ the data generation technique to enhance sample diversity. Others, like MLSN [50] and TLNSS [51], introduce transfer learning to mitigate the problem of insufficient training samples. In addition, some method incorporate prior spectra into their network structure by adopting siamese networks for contrastive learning, including TSCNTD [52], STTD [53], AGMS [49], and CS-TTD [54]. While other methods design networks based on physical models, such as JSPEN [55] and IRN [56]. Given that HTD involves the processing of one-dimensional spectral curves, Transformer, which originated in the NLP domain, have gained popularity in HTD research. Notable examples include STTD, CS-TTD, TSTTD [47], HTDFormer [57], MCLT [58], and SGT [59]. However, current HTD methods are based on a per-pixel binary classification framework which restricts the representation of instance-level object features. In this paper, we explore the feasibility of adopting an instance-level object detection framework as an alternative to the current spectral-level HTD framework.

B. Object Detection

In recent years, the field of visual object detection has made significant progress driven by deep learning, encompassing two-stage detectors, single-stage detectors, and DETR-like detectors. Two-stage detectors, such as Faster R-CNN [15], provide more precise object detection through a two-step process but are computationally intensive, whereas single-stage detectors like YOLO [16] offer faster, albeit slightly less accurate, detection. Carion et al. [17] proposed an end-to-end object detector based on Transformer [23], named DETR, which reformulates object detection as a one-to-one set prediction problem. Deformable DETR [24] significantly improves the training efficiency of DETR by introducing deformable attention mechanisms and multi-scale feature fusion. DAB-DETR [60] decomposes the object query into content and positional queries and introduces explicit positional representation. DN-DETR [61] and Group DETR [62] introduce denoising training and grouped queries, respectively, to address the instability of bipartite matching and accelerate model convergence. Building on these works, DINO [25], through

the introduction of contrastive denoising training, mixed query selection, and look forward twice techniques, surpasses CNN detectors on the visual object detection benchmark COCO dataset [18], achieving the best results. CO-DETR [63] further improves detection accuracy by introducing more dense supervision signals to enhance the discriminability of encoder features. However, these DETR-like detectors are typically computationally expensive. RT-DETR [64] achieves real-time detection performance comparable to that of the YOLO series. YOLOv10 [65] combines the label assignment strategies of DETR with CNN architecture, achieving an optimal balance between detection accuracy and inference speed. These visual object detection networks are designed for objects with distinct shape and texture details. In this work, we aim to further optimize and customize the DETR framework for the point object detection task.

C. Hyperspectral Object Detection

Recently, object detection has been introduced to hyperspectral image processing, showing promising results in fields such as camouflage object detection [66], vehicle detection [67]–[69], and methane gas detection [70]. For instance, S2ADet [69] employs a dual-backbone network architecture, where one backbone processes color images while the other operates on three-channel images derived from hyperspectral data through Principal Component Analysis, extracting spatial and spectral semantic information, respectively. MethaneMapper [70] uses two ResNet networks to extract features from visible and shortwave infrared bands, respectively. These works rely on visual backbone networks pre-trained on large RGB datasets. Similarly, recent popular hyperspectral object tracking networks also follow this approach [71]–[74]. SEE-Net uses pre-trained backbone networks to extract features from different pseudo-color image combinations [71]. SiamHYPER adds a hyperspectral awareness module to the RGB-based Siamese tracker, extracting features through dual-inputs of RGB and hyperspectral data [72]. HA-Net also employs dual-inputs of RGB and hyperspectral data to extract highly discriminative semantic features [73].

On one hand, the aforementioned researches primarily focus on large-sized objects and are not well-suited to address the challenge of subpixel object detection, which is a critical concern in the traditional HTD task. On the other hand, these works still add spectral information extraction designs on top of pre-trained backbone networks based on RGB inputs. In this paper, we propose the hyperspectral point object detection task, aiming to tackle the problem of detecting extremely small objects at the subpixel scale within the object detection framework. Furthermore, we explore whether the current object detection framework can directly extract spatial-spectral joint features from hyperspectral cube data.

III. METHODOLOGY

A. Task Definition

Traditional HTD methods use a prior target spectral library to perform binary classification on pixels. Our work extends

visual object detection to the hyperspectral point object detection task, aiming to achieve object-level detection of point objects. Our network learns object-level features on a set of training HSIs with bounding box (bbox) annotations of point objects. If annotations include category labels, the trained network can predict the categories of point objects. In contrast, the HTD methods always treat all target prior spectra as the same category and neglect classification ability.

For the hyperspectral point object detection task, we adopt the GT label format of the common object detection such as the COCO dataset [18]. For each hyperspectral image in both the training and testing sets, the GT label comprises the true bbox information and category label for each object. Each object is accompanied by a bbox, which precisely indicates the object's location in the image. This rectangular box is typically represented by four parameters: the coordinates of the top-left corner (x, y), and the width (w) and height (h) of the box, or alternatively, the center coordinates (x, y), width (w), and height (h). These two forms of representation are usually equivalent. The GT label for an object also includes a category label.

B. Overview

As shown in Fig. 3, our SpecDETR is built upon DINO [25] and removes the backbone network of DINO. It only consists of a multi-layer Transformer encoder, a multi-layer Transformer decoder, and multiple prediction heads. Given an HSI, we use a linear layer to transform the channel number and treat each pixel as a token. Then, we feed them into the Transformer encoder to extract deep features. Based on the output of the multi-layer encoder, a candidate anchor box is generated at each pixel location, from which a fixed number of anchor boxes are selected. For each selected anchor box, an initial object query vector is generated. The initial object queries and their corresponding anchor boxes are then fed into the multi-layer decoder and updated layer-by-layer. The output of SpecDETR consist of refined bboxes and classification results generated by updated object queries. Following DINO, SpecDETR introduces an additional denoising branch alongside the original decoder pipeline which is termed the matching branch. Unlike current DETR-like detectors, SpecDETR regards point object detection as a one-to-many set prediction problem, and accordingly, it has been specifically refined. We propose Center-Shifting Contrastive DeNoising (CCDN) training to replace Contrastive DeNoising (CDN) training in DINO. In the matching branch, we use a hybrid one-to-many label assigner instead of a Hungarian matching-based one-to-one label assignment in current DETR-like detectors. Based on the one-to-many matching, we simplify the decoder by removing cross-attention between object queries and initializing object queries as non-learnable, uniform vectors. Additionally, SpecDETR uses non-maximum suppression (NMS) to remove overlapping prediction boxes. The loss function of SpecDETR is consistent with that of DINO. We will elaborate on the details of our improvements in the following sections.

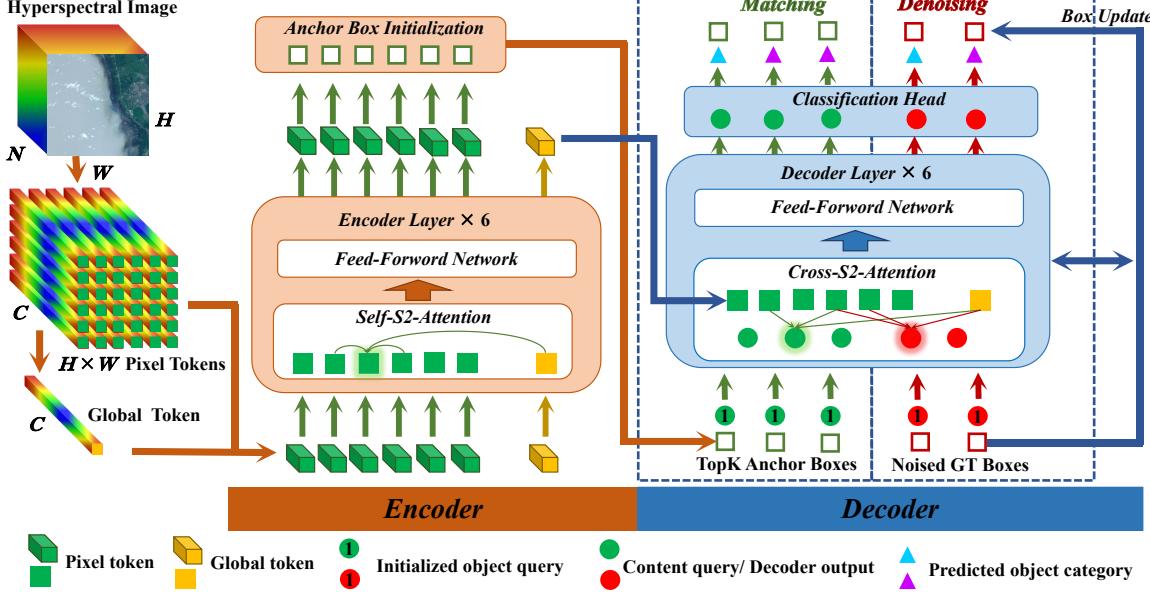


Fig. 3. An overview of our SpecDETR. For clear visualization, we omit positional embedding and only display the classification prediction head of the final decoder layer.

C. Data Tokenization

Given an HSI data cube $\mathbf{X} \in \mathbb{R}^{H \times W \times N}$, where H , W , and N represent the width, height, and band number, we preprocess \mathbf{X} before feeding it into the encoder by:

$$\mathbf{P}_0 = \text{LN} \left(\text{linear} \left(\frac{\mathbf{X}}{V} \right) \right), \quad (1)$$

where $\text{linear}(\cdot)$ represents a linear layer, $\text{LN}(\cdot)$ represents a layer normalization layer [75], and V is a constant. The output data cube $\mathbf{P}_0 \in \mathbb{R}^{H \times W \times C}$ has an initialized feature vector at each pixel position, serving as a pixel token for the following Transformer, where C represents the dimension of tokens. Subsequently, a initialized global token \mathbf{g}_0 is initialized by:

$$\mathbf{g}_0 = \frac{1}{H \times W} \sum_{i=1}^{H \times W} \mathbf{p}_{0,i}, \quad (2)$$

where $\mathbf{p}_{0,i}$ represents the i -th initialized pixel token.

Normalization is a crucial preprocessing step in image processing. For general vision tasks, 8-bit quantized RGB images are typically Z-score normalized for each band before being input into neural networks. For HSI processing tasks, HSIs are usually subjected to Max-Min normalization. However, both normalization methods are unsuitable for our hyperspectral point object detection framework. HSIs possess higher quantization bits and more bands than conventional RGB images, with each band exhibiting distinct data distributions. Some bands have minor standard deviations, resulting in excessively large values after Z-score normalization, which can adversely affect network convergence. Furthermore, traditional HTD work primarily validates on datasets containing only a single HSI, whereas our proposed framework is designed for datasets with a vast number of HSIs. On one hand, the numerical range varies between HSIs, and performing Max-Min normalization on each HSI individually can result in

inconsistencies in the spectral numerical characteristics of HSIs. On the other hand, actual hyperspectral radiance images sometimes contain pixels with extremely high values; after Max-Min normalization, most image regions are concentrated within a very small numerical range, causing quantization loss of spectral information. Therefore, in Eq. (1), we normalize all HSIs within the same dataset by dividing by the same constant V . We adopt different settings for V depending on whether the data is radiance or reflectance. For the SanDiego and Gulfport datasets used in our experiments, which are reflectance data subjected to radiometric correction, the theoretical numerical range is $[0,1]$, and V is set to 1. The numerical range of radiance data is influenced by various factors, including band range, acquisition scenario, unit system, and producer processing methods. For radiance data, we employ a lenient criterion for selecting V : ensuring that most values of the normalized HSIs within the same dataset are smaller than 1. The other two datasets used in our experiments, SPOD and AVON, are based on radiance data collected by different hyperspectral sensor in different data campaigns. For the SPOD dataset, we set V to 3000; for the AVON dataset, we set V to 5000. This setting can also be automated by setting V to the value corresponding to the 60%, 80%, or other specified percentiles of all values in the dataset, sorted in ascending order.

D. Transformer Encoder

Deformable DETR and its successors, such as DN-DETR, DINO, and Co-DETR, employ a deformable Transformer encoder on the multi-scale feature maps outputted by the backbone network to further refine features. The deformable Transformer encoder replaces the Transformer attention module in the original DETR with a deformable multi-scale attention module, which focuses solely on a small subset of key sampling points around reference points, serving as a pre-filter

to highlight critical elements among all feature map pixels. For the point objects of interest in this paper, given their extremely small size, they only interact with nearby local background pixels, while distant background pixels have negligible impact. The deformable attention module, as a local sparse sampling operator, is particularly well-suited for feature extraction of point objects. Therefore, we propose the self-excited subpixel-scale attention module (S2A module) based on the deformable attention module for the extraction of spatial-spectral joint features of point objects. Furthermore, following the DETR-like detectors, we cascade multiple Transformer encoder layers equipped with the S2A module to achieve a progressive refinement of object and background features from shallow to deep levels. The input to the multi-layer Transformer encoder consists of \mathbf{P}_0 from Eq. (1) and \mathbf{g}_0 from Eq. (2), with the set of all initialized tokens denoted as \mathbf{F}_0 . For the j -th Transformer encoder layer, it receives the output \mathbf{F}_{j-1} from the previous encoder layer and processes it as follows:

$$\tilde{\mathbf{F}}_j = \text{LN}(\text{Self-S2A}(\mathbf{F}_{j-1}) + \mathbf{F}_{j-1}), \quad (3)$$

$$\mathbf{F}_j = \text{LN}\left(\text{FFN}(\tilde{\mathbf{F}}_j) + \tilde{\mathbf{F}}_j\right), \quad (4)$$

where S2A (\cdot) represents the S2A module, and FFN (\cdot) denotes the feed-forward network, comprising linear layers and ReLU activation layers, also known as the multi-layer perceptron (MLP). During the classic Transformer encoder, the query and value elements of the attention module are identical, hence the attention module in the encoder is termed a self-attention module. Adhering to this convention, we refer to the S2A module in Eq. (3) as the self-S2A module.

The self-S2A module integrates two feature extraction operators based on deformable attention. The first is a subpixel-scale deformable sampling operator performed on the feature map \mathbf{P}_{j-1} at the original image resolution. In the traditional DETR framework, deep feature extraction is jointly accomplished by the pre-trained backbone network and the encoder. However, current backbone networks like ResNet are designed for general image understanding tasks, where object shape information is rich and sizes vary, with some objects even occupying the entire image. Consequently, these backbone networks extract feature maps at different downsampling ratios from the input image, while the encoder in current DETR-like detectors is used for further feature refinement. For the point objects focused on in this paper, the emphasis of current backbone networks on larger spatial features and the downsampling operations are not appropriate. Therefore, we directly perform deformable attention computation on each pixel token at the original image feature resolution, replacing the traditional DETR framework's pipeline of collaborative feature extraction between the backbone network and the encoder. The updated feature map retains the original resolution, thereby avoiding the loss of features for extremely small-sized objects due to downsampling operations. Deformable sampling is based on bilinear interpolation, which aligns well with the linear mixing model assumption of mixed pixels, enabling subpixel-scale deformable sampling on \mathbf{P}_{j-1} . In addition, we introduce a global token which combined with the subpixel-scale deformable sampling operator, and implements the self-

excited operator. Although the self-excited operator is also based on deformable attention, it functions to self-amplify object features rather than feature aggregation. We term this mechanism as the self-excited mechanism. The subpixel-scale deformable sampling operator and the self-excited operator represent two extremes: one computes deformable attention on the feature map at the original resolution, while the other does so on a single-pixel feature map. Consequently, the multi-scale deformable attention module from deformable DETR can be directly utilized for their implementation.

The schematic diagram of the self-S2A module is shown in Fig. 4. The number of deformable attention heads in the self-S2A module is denoted as M , with each attention head having K sampling points, and M divides the token dimension C evenly. Taking the i -th input pixel token $\mathbf{p}_{j-1,i}$ as the current query element, its 2D pixel position is set as the current reference point \mathbf{r} . $\mathbf{p}_{j-1,i}$, augmented with its position embedding, is passed through two independent linear layers to obtain $M \times K \times 2$ sampling point position offsets $\Delta\mathbf{r}$ relative to the reference point and attention weights \mathbf{A} , respectively:

$$\Delta\mathbf{r} = \text{linear}(\mathbf{p}_{j-1,i} + \text{pos}(\mathbf{p}_{j-1,i})) + \Delta\mathbf{r}_{\text{int}}, \quad (5)$$

$$\mathbf{A} = \text{linear}(\mathbf{p}_{j-1,i} + \text{pos}(\mathbf{p}_{j-1,i})), \quad (6)$$

where $\text{pos}(\mathbf{p}_{j-1,i})$ is the position encoding of the same dimension as \mathbf{r} , $\Delta\mathbf{r}_{\text{int}}$ represents the initialized position offsets, $\Delta\mathbf{r} \in \mathbb{R}^{2 \times M \times K \times 2}$, and $\mathbf{A} \in \mathbb{R}^{M \times K \times 2}$. $\Delta\mathbf{r}$ is split into the pixel token component $\Delta\mathbf{r}_P \in \mathbb{R}^{2 \times M \times K}$ and the global token component $\Delta\mathbf{r}_G \in \mathbb{R}^{M \times K \times 2}$. Similarly, \mathbf{A} is also split into $\mathbf{A}_P \in \mathbb{R}^{M \times K}$ and $\mathbf{A}_G \in \mathbb{R}^{M \times K}$. The \mathbf{P}_{j-1} and \mathbf{g}_{j-1} are fed into the same linear layer to generate the feature maps to be sampled (which can be understood as the value elements in the classical Transformer):

$$\hat{\mathbf{P}}_j, \hat{\mathbf{g}}_j = \text{linear}(\mathbf{P}_{j-1}, \mathbf{g}_{j-1}). \quad (7)$$

$\hat{\mathbf{P}}_j$ and $\hat{\mathbf{g}}_j$ are both evenly split along the channel dimension into M sub-feature maps to be sampled, resulting in a channel dimension of C/M . For the m -th attention head, the sampling value of the current query element $\mathbf{p}_{j-1,i}$ is:

$$\begin{aligned} \mathbf{v}^m &= \sum_{k=1}^K \left(A_P^{m,k} \cdot \text{bili}(\hat{\mathbf{P}}_j^m, \mathbf{r} + \Delta\mathbf{r}_P^{m,k}) \right) \\ &\quad + \sum_{k=1}^K \left(A_G^{m,k} \cdot \text{bili}(\hat{\mathbf{g}}_j^m, \mathbf{r} + \Delta\mathbf{r}_G^{m,k}) \right), \end{aligned} \quad (8)$$

where k is the index of the sampling point, m is the index of the attention head, i is the index of the encoder layer, $\mathbf{v}^m \in \mathbb{R}^{C/M}$, and $\text{bili}(\cdot)$ represents the sampling operation based on bilinear interpolation. The scalar attention weights $A_P^{m,k}$ and $A_G^{m,k}$ are normalized by $\sum_{k=1}^K (A_P^{m,k} + A_G^{m,k})$. After obtaining the sampled feature values from the M attention heads, the current query element $\mathbf{p}_{j-1,i}$ is updated as follows:

$$\tilde{\mathbf{p}}_{j,i} = \text{linear}(\mathbf{p}_{j-1,i} + \text{concat}(\mathbf{v}^1, \mathbf{v}^2, \dots, \mathbf{v}^M)), \quad (9)$$

where $\text{concat}(\cdot)$ denotes the concatenation operation along the channel dimension for the vectors. The same operations from Eq. (5) to Eq. (9) are applied to the remaining pixel tokens

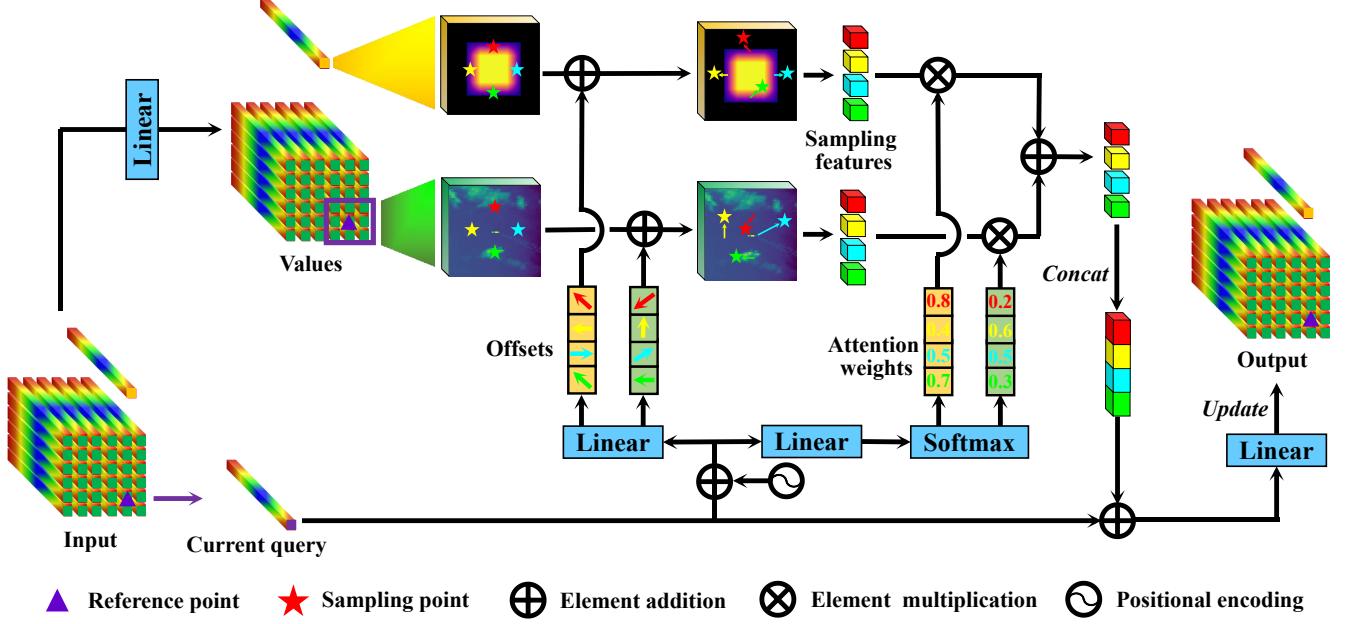


Fig. 4. A schematic diagram of the self-S2A module. For clear visualization, the number of sampling points for each attention head is set to 1. Different colors of the sampling points represent different attention heads.

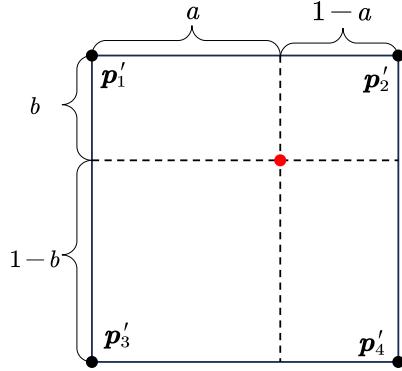


Fig. 5. A schematic diagram of bilinear interpolation.

and the global token \mathbf{g}_{j-1} , resulting in the updated $\tilde{\mathbf{P}}_j$ and $\tilde{\mathbf{g}}_j$. The union of these two sets constitutes the $\tilde{\mathbf{F}}_j$ in Eq. (3).

Compared to ordinary convolution operation or the ordinary Transformer attention module, the deformable sampling operation based on bilinear interpolation extends the feature sampling space from the pixel scale to the subpixel scale. In Eq. (8), the position range of the sampling points, $\mathbf{r} + \Delta\mathbf{r}_P^{m,k}$, includes components from the linear layer output, and therefore is within the real number domain. This approach allows the sampling point offsets to become learnable and to be smoothly updated during backpropagation, enabling the adaptive adjustment of the sampling point distribution during training. Additionally, sampling based on bilinear interpolation ensures effective extraction of spectral features at subpixel-scale locations. As shown in Fig. 5, $\mathbf{p}1'$, $\mathbf{p}2'$, $\mathbf{p}3'$, and $\mathbf{p}4'$ are the vectors at the integer coordinate points in $\hat{\mathbf{P}}_j^{m,k}$ closest to $\mathbf{r} + \Delta\mathbf{r}_P^{m,k}$. The horizontal and vertical distances from $\mathbf{r} + \Delta\mathbf{r}_P^{m,k}$ to $\mathbf{p}1'$ are a and b , respectively. The sampled

value by bilinear interpolation is given by:

$$\text{bili}(\hat{\mathbf{P}}_j^{m,k}, \mathbf{r} + \Delta\mathbf{r}_P^{m,k}) = a\mathbf{p}1' + (b - ab)\mathbf{p}2' + (a - ab)\mathbf{p}3' + (1 - a - b + ab)\mathbf{p}4'. \quad (10)$$

By denoting the coefficients of $\mathbf{p}1'$, $\mathbf{p}2'$, $\mathbf{p}3'$, and $\mathbf{p}4'$ as ε_1 , ε_2 , ε_3 , and ε_4 , respectively, the above equation can be rewritten in the form of the linear mixing model (LMM) [76]:

$$\text{bili}(\hat{\mathbf{P}}_j^{m,k}, \mathbf{r} + \Delta\mathbf{r}_P^{m,k}) = \sum_{z=1}^4 \varepsilon_z \mathbf{p}_z', \quad \text{s.t. } \begin{cases} \sum_{z=1}^4 \varepsilon_z = 1 \\ \varepsilon_z > 0 \end{cases}. \quad (11)$$

In remote sensing images, a pixel may contain multiple types of ground objects. The LMM handles this mixed phenomenon by assuming that the reflectance or radiance spectrum of a pixel is a linear combination of the spectra of different ground objects. The abundance of each spectral bases being non-negative and the sum of abundances equal to 1. Eq. (11) implies that the sub-pixel scale feature aggregation operator samples in a spectral feature space that is continuous in positional dimensions. Particularly for the first encoder layer, $\hat{\mathbf{P}}_1$ is generated by applying two fully connected operations along the feature dimension to the original spectral image, maintaining spatial features consistent with the original image. The sampling points can be viewed as mixed pixels generated by considering the nearest four spectra as spectral bases, with the feature proportion of each spectrum related to the spatial distance. The feature sampled at any given subpixel-scale position are capable of reflecting the spectral characteristics of neighboring pixels, thereby ensuring effective spectral feature extraction.

Fig. 6 illustrates the self-excited mechanism. The self-excited operator and the subpixel-scale deformable sampling operator are consistent in structure, but they implement dif-

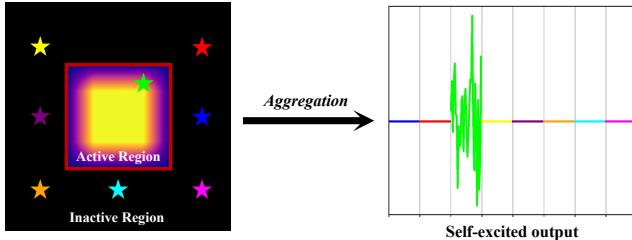


Fig. 6. A schematic diagram of the self-excited mechanism. For clear visualization, the number of sampling points for each attention head is set to 1. Different colors of the sampling points represent different attention heads.

ferent functions. This difference arises from the variation in the value elements to be sampled. The subpixel-scale deformable sampling operator treats the feature map as the value elements to be sampled, achieving spatial-spectral joint feature aggregation in the local region where the query element is located. In contrast, the self-excited operator uses the global token as the value element to be sampled, leveraging this single value element to amplify its own semantic information. In the deformable attention operation based on bilinear interpolation, the global token can also be transformed into a feature space that is continuous in positional dimensions. The global token fills the origin and its eight adjacent integer coordinate points, while other integer coordinates are filled with zero vectors. After bilinear interpolation, a non-zero region of 4×4 in size, centered at the origin, is generated, with all other regions being zero vectors. We refer to the non-zero square region as the active region, and the remaining area of the feature space as the inactive region. In deformable DETR and its successors, the number of attention heads is set to 8 by default. In the initialization for deformable attention, the initial positions of the k -th sampling point for the 8 attention heads are sequentially set to the four vertices and the midpoints of the four sides of a square with a size of $2k \times 2k$ centered at the reference point. The reference point is scaled to a square region of 2×2 in size centered at the origin. This ensures that the initial positions of the first sampling points for the 8 attention heads are located within the active region, while other sampling points are located in the inactive region. The final sampling positions are the initial positions plus the spatial offsets generated by the pixel token serving as the query element through linear projection. When all sampling points of an attention head are located in the inactive region, the output of this attention head is a zero vector. However, if any sampling point of an attention head is in the active region, that attention head is activated, and its output is non-zero. The output vector of the self-excited operator is the concatenation of the outputs from the 8 attention heads along the channel dimension. As illustrated in the right figure of Fig. 6, different combinations of activated attention heads generate special sampling outputs. As previously analyzed, the activation of different attention heads is determined by the pixel tokens of the query elements. Consequently, the semantic differences among various types of pixel markers can lead to different combinations of activated attention heads. The sampling output form of the self-excited operator accentuates these

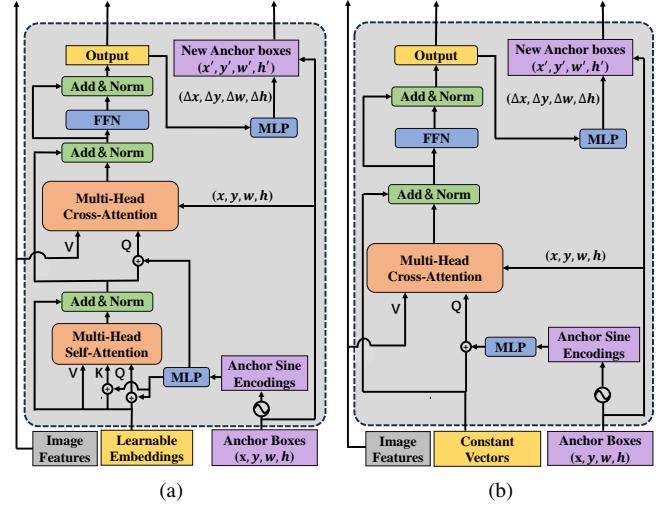


Fig. 7. Comparison of decoder layer structures. (a) DINO. (b) SpecDETR. Only the first layer is illustrated. The structure of the remaining layers is identical to that of the first layer, while the object queries are the output of the preceding layers.

semantic differences. The implementation of the self-excitation mechanism is straightforward and can directly utilize existing multi-scale deformable attention modules.

E. Transformer Decoder

In the design of the transformer decoder, we optimize based on DINO’s decoder, retaining the iterative bbox refinement and two-stage framework from deformable DETR, as well as the contrastive denoising training and a box prediction method named “look forward twice” from DINO. DETR-like detectors model object detection as a one-to-one set prediction problem, primarily to eliminate the need for traditional detectors’ NMS post-processing operation. General object detection datasets, such as COCO [18], contain numerous highly overlapping objects, while detectors prior to DETR can produce multiple redundant predictions for the same object. These two phenomena create a parameter setting dilemma for the IoU threshold of NMS, making the removal of NMS in DETR-like detectors particularly meaningful. However, our experiments reveal that DETR-like detectors still generate redundant predictions for point objects. Additionally, overlapping point objects, due to their lack of significant shape and texture details, cannot be distinguished like overlapping large objects in the COCO [18] dataset. This implies that overlapping point objects should be treated as a single object for detection. Therefore, we believe that removing NMS is not suitable for point objects, and we model point object detection as a one-to-many set prediction problem. Based on this, we optimize the decoder structure used by DINO, as shown in Fig. 7, and improve the contrastive denoising training and label assignment method. This section mainly introduces the forward propagation pipeline of the decoder.

The decoder is composed of multiple cascaded decoder layers, each of which uses object queries as query elements and the feature maps output by the encoder as value elements. These layers iteratively update the object queries and refine

the object anchor boxes. Let the number of encoder layers be E and the number of decoder layers be D . Initially, each pixel token output by the last encoder layer generates a candidate anchor box. For the i -th pixel token $\mathbf{p}_{E,i}$, the initial center position of the corresponding anchor box is set to the pixel location (p_{xi}, p_{yi}) , with both the initial width and height set to s . A normalized bbox vector $\mathbf{b}_{\text{int},i}$ is constructed as $[\frac{p_{xi}}{W}, \frac{p_{yi}}{H}, \frac{s}{W}, \frac{s}{H}]$. The pixel token $\mathbf{p}_{E,i}$ is fed into a bbox regression head to refine the initial bbox vector:

$$\mathbf{b}_{0,i} = \text{Sig}(\text{breg}_0(\mathbf{p}_{E,i}) + \text{InSig}(\mathbf{b}_{\text{int},i})) \quad (12)$$

where the bbox regression head $\text{breg}_0(\cdot)$ is a 3-layer MLP network, $\text{Sig}(\cdot)$ represents the Sigmoid activation function that maps the entire real number domain to the (0,1) interval, and $\text{InSig}(\cdot)$ represents the inverse of the Sigmoid activation function, mapping the (0,1) interval to the entire real number domain. Subsequently, $\mathbf{p}_{E,i}$ is fed into a classification head:

$$\mathbf{c}_{0,i} = \text{Sig}(\text{cls}_0(\mathbf{p}_{E,i})) \quad (13)$$

where the classification head $\text{cls}_0(\cdot)$ is a linear layer, and $\mathbf{c}_{0,i} \in \mathbb{R}^{N_{\text{cat}}}$ represents the category predictions, also known as object confidence scores, with N_{cat} being the number of object categories. After the Sigmoid activation function, the values of $\mathbf{c}_{0,i}$ range between (0,1), where values closer to 1 indicate a higher probability that the prediction belongs to the corresponding category. From all $H \times W$ candidate anchor boxes, the top Q_{match} boxes with the highest confidences are selected as the initial bboxes corresponding to the object queries of the decoder. During training, the denoising training branch additionally provides Q_{DN} noised GT boxes, the details of which will be introduced in the following text. Let the number of object queries for the decoder be Q ; during training, $Q = Q_{\text{match}} + Q_{\text{DN}}$, and during inference, $Q = Q_{\text{match}}$. DINO uses the index embedding to generate initial object queries, endowing them with learnability. Under one-to-one set prediction modeling, learnable and distinguishable object queries can enhance detection performance. However, in the context of point object detection with one-to-many set prediction modeling, we believe that distinguishable object queries is unnecessary. On one hand, we remove the self-attention computation between object queries, so they do not interfere with each other during forward pass. On the other hand, DINO aims for a single prediction per object, while our SpecDETR does not impose this constraint. Therefore, we initialize all object queries as all one vectors of dimension C .

Let d index the decoder layers and i index the object queries. In the d -th decoder layer, we take the object queries fed from the previous layer as the query elements and the output \mathbf{F}_E of the last encoder layer as the value elements to compute the cross-attention between the object queries and \mathbf{F}_E , thereby updating the object queries as follows:

$$\tilde{\mathbf{q}}_{d,i} = \text{LN}(\text{Cross-S2A}(\mathbf{q}_{d-1,i}, \mathbf{F}_E) + \mathbf{q}_{d-1,i}), \quad (14)$$

$$\mathbf{q}_{d,i} = \text{LN}(\text{MLP}(\tilde{\mathbf{q}}_{d,i}) + \tilde{\mathbf{q}}_{d,i}), \quad (15)$$

where $\mathbf{q}_{d,i} \in \mathbb{R}^C$, $d \in \{1, 2, \dots, D\}$, and $i \in \{1, 2, \dots, Q\}$. Compared to Eq. (3), the query elements and value elements of Eq. (14) are not the same, hence the S2A module in Eq. (14)

is termed a cross-S2A module. The cross-S2A module is used to capture the object category and boundary information near the anchor box $b_{d-1,i}$ corresponding to $\tilde{\mathbf{q}}_{d-1,i}$. Therefore, the reference point of the cross-S2A module is set to the center of the anchor box, and the initial positions of the sampling points are scaled within the anchor box, with the outermost sampling points placed on the box boundary. Other details are consistent with the self-S2A module.

The vector $\mathbf{q}_{d,i}$ is fed into a bbox regression head to update the bbox vector $\mathbf{b}_{d-1,i}$ as follows:

$$\hat{\mathbf{b}}_{d,i} = \text{Sig}(\text{breg}_d(\mathbf{q}_{d,i}) + \text{InSig}(\mathbf{b}_{d-1,i})), \quad (16)$$

$$\mathbf{b}_{d,i} = \text{detach}(\hat{\mathbf{b}}_{d,i}), \quad (17)$$

$$\mathbf{b}_{d,i}^{\text{pred}} = \text{Sig}(\text{breg}_d(\mathbf{q}_{d,i}) + \text{InSig}(\hat{\mathbf{b}}_{d-1,i})), \quad (18)$$

where $\text{breg}_d(\cdot)$ is a bbox regression head with the same structure as $\text{breg}_0(\cdot)$, and $\text{detach}(\cdot)$ represents the gradient detachment operation which can interrupt the backpropagation. Both $\mathbf{b}_{d,i}^{\text{pred}}$ and $\mathbf{b}_{d,i}$ are the bbox prediction for the i -th object query at the d -th decoder layer, with $\mathbf{b}_{d,i}$ used for bbox updates in the $d+1$ -th decoder layer and $\mathbf{b}_{d,i}^{\text{pred}}$ used to compute the loss function for the d -th decoder layer. Eq. (16)-Eq. (18) collectively implement DINO’s look forward twice method, which ensures that the parameter updates of $\text{breg}_d(\cdot)$ are influenced by both the bbox prediction error at the d -th decoder layer and the bbox prediction error at the $d+1$ -th decoder layer. $\mathbf{q}_{d,i}$ is then fed into a classification head to obtain the class prediction results corresponding to $\mathbf{b}_{d,i}^{\text{pred}}$:

$$\mathbf{c}_{d,i} = \text{Sig}(\text{cls}_d(\mathbf{q}_{d,i})) \quad (19)$$

where $\text{cls}_d(\cdot)$ is a classification head with the same structure as $\text{cls}_0(\cdot)$. The confidence scores for each class are combined with $\mathbf{b}_{d,i}^{\text{pred}}$, resulting in $Q \times N_C$ predictions. The refined bboxes and class confidence scores output from the last decoder layer serve as the prediction output of SpecDETR, while the prediction results from other layers are used only for loss function computation during the training phase. The output $Q \times N_C$ predictions can be reduced for redundancy through post-processing techniques such as NMS and confidence score filtering.

In the classical Transformer network, the decoder contains a self-attention module among query elements as well as a cross-attention module between query elements and encoder output elements. As shown in Fig. 7(a), the decoder in current DETR-like detectors still follows this pipeline to facilitate information interaction between object queries, thereby enabling low-redundancy predictions without the need for NMS post-processing. However, as discussed at the beginning of this section, point object detection is more suitably modeled as a one-to-many rather than a one-to-one set prediction problem, making the information interaction between object queries unnecessary. Therefore, SpecDETR removes the self-attention module between object queries, simplifying the processing pipeline and reducing computational complexity.

F. Hybrid Label Assigner

Group DETR, DN-DETR, and DINO have demonstrated that the one-to-one label assigner employed by the original DETR leads to a lack of supervisory signals during training, thereby resulting in slower convergence. As discussed above, removing NMS is not suitable for point objects, and point object detection should be modeled as a one-to-many set prediction problem. Therefore, we replace the Hungarian matching in current DETR-like detectors with a one-to-many label assigner in the matching branch. Traditional object detection networks often use some IoU-based label matching [15] as the one-to-many label assigners. However, due to the extremely small size of point objects, even the deviation of a few pixels can significantly reduce the IoU between the predicted bbox and the GT box. Therefore, we propose a hybrid one-to-many label assigner that combines the Hungarian matching in DETR [17] and Max IoU matching [15]. Initially, we use the Hungarian matching to forcibly assign a predicted bbox to each GT sample, a step we refer to as forced matching. Subsequently, we allocate an additional variable number of high-quality predicted bboxes to each GT sample based on IoU between the predicted boxes and GT boxes, a step termed dynamic matching.

Forced Matching. During training, when the bbox regression head is inaccurate or encounters challenging samples, the predicted bbox can deviate from the GT box. Given that point objects are only a few or even one pixel in size, pixel-level position deviations can result in extremely low IoU or even an IoU of zero. The Hungarian matching in DETR achieves bipartite matching between GT boxes and predicted bboxes based on a combined loss of the bbox GIoU loss [77], bbox L1 loss, and classification loss. Therefore, we first employ Hungarian matching to forcibly assign a positive sample to each GT box, ensuring the network can learn the hard samples.

Dynamic Matching. Following forced matching, we utilize the Max IoU assigner [15] to further allocate positive samples among the remaining anchor boxes. In addition to the IoU threshold, we introduce a number threshold T . If a GT box is assigned to more than T anchor boxes, we retain only the top T anchor boxes with the highest IoU as positive samples. We seek to ensure that predicted bboxes which closely match the GT boxes have higher confidence scores from the classification head and can be retained after NMS. Dynamic matching leverages both the IoU threshold and the number threshold to provide high-quality positive samples for the classification head. The introduction of a matching number threshold allows the actual matching IoU threshold to dynamically increase as the network converges, thereby continuously improving the quality of positive samples.

G. Center-Shifting Contrastive DeNoising Training

DN-DETR introduces a DeNoising (DN) training branch in the decoder to stabilize training and accelerate convergence. DINO improved DN training to Contrastive DeNoising (CDN) training, incorporating negative queries to reject useless anchors. CDN generates the anchor boxes of positive and negative queries by adding two levels of noise (position

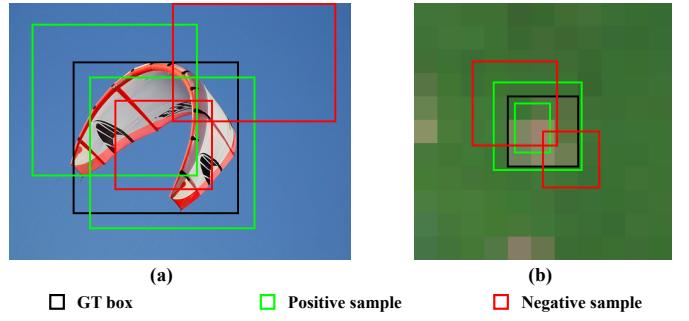


Fig. 8. Illustration of positive and negative positional queries of CDN and CCDN training. (a) CDN. (b) CCDN.

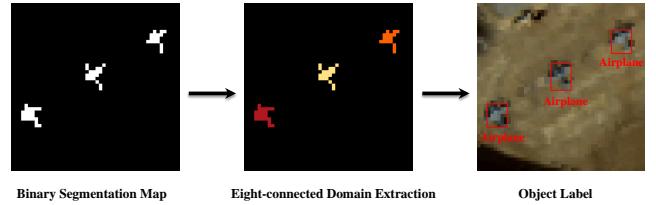


Fig. 9. The schematic diagram for generating object labels on binary segmentation maps.

offsets) to the corners of the GT boxes. However, CDN is designed for one-to-one matching, where the anchor box range of negative queries includes the positive queries, as shown in Fig. 8(a). However, negative samples with high overlap with GT boxes can disrupt the classification heads of our SpecDETR. Therefore, we adopt the query generation method of DN [61], consisting of center shifting and box scaling steps controlled by hyperparameters τ_1 and τ_2 , with $\tau_1 > 0$ and $\tau_2 > 0$. Our denoising approach adds two levels of shifting to the centers of the GT boxes to generate the anchor boxes of positive and negative queries, termed as Center-Shifting Contrastive DeNoising (CCDN).

Center Shifting. Given a GT box (x, y, w, h) , random noises (n_x, n_y) are added to the center (x, y) . For positive queries, the noise satisfies $-0.5\tau_1w < n_x < 0.5\tau_1w$ and $-0.5\tau_1h < n_y < 0.5\tau_1h$. For negative samples, n_x is randomly sampled within the range $[-\tau_1w, -0.5\tau_1w] \cup [0.5\tau_1w, \tau_1w]$, and n_y is randomly sampled within the range $[-\tau_1h, -0.5\tau_1h] \cup [0.5\tau_1h, \tau_1h]$.

Box Scaling. Regardless of positive or negative queries, the width and height are randomly sampled within the ranges $[\min(0, w - \tau_2w), w + \tau_2w]$ and $[\min(0, h - \tau_2h), h + \tau_2h]$, respectively.

IV. DATASETS AND EVALUATION METRICS

This section will introduce the hyperspectral multi-class point object detection dataset that we develop, as well as three public hyperspectral target detection datasets the Avon dataset, SanDiego dataset and Gulfport dataset. In addition, we will also discuss the construction of GT labels and evaluation criteria for the point object detection task.

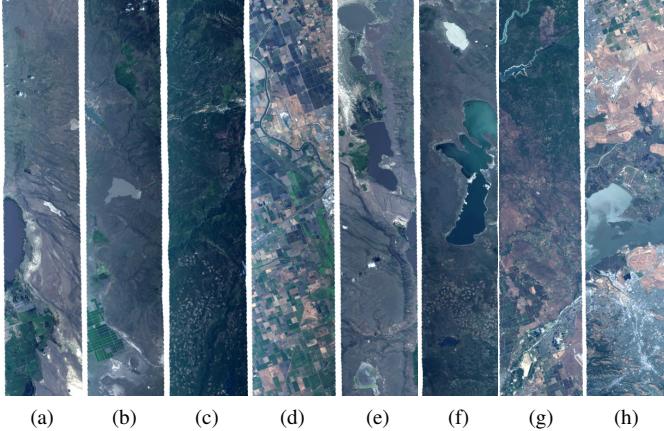


Fig. 10. Partial regions of AVIRIS flight lines used to construct the SPOD dataset. (a)-(d) The flight line f180601t01p00r06 for training data. (e)-(h) The flight line f180601t01p00r10 for test data.

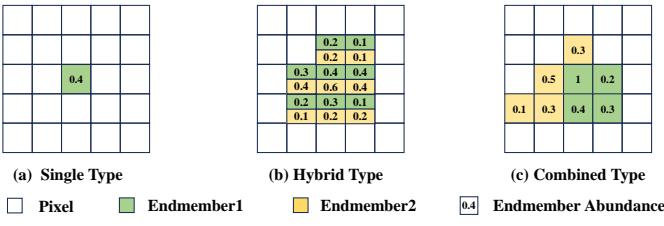


Fig. 11. Diagrams of the three types of object templates on the SPOD dataset.

A. Ground-truth Label Construction

Object detection datasets are typically annotated manually, where professional annotators view the images and draw the bounding boxes of the objects, while also marking the object categories. However, for point objects, many constituent pixels are mixed pixels, making manual annotation challenging to ensure the accuracy of mixed pixel labeling. Therefore, some real HTD datasets, such as the Avon and Gulfport datasets, provide rough GT positions obtained from GPS, which have significant errors. To accurately evaluate the performance of different methods on hyperspectral multi-class point object detection, we have constructed two datasets with precise labels: SPOD and Avon. These datasets provide 8 classes and 2 classes of point objects, respectively. The SPOD dataset is simulation-based, and precise GT labels can be directly obtained from the object simulation configuration. In contrast, the Avon dataset is a real-world dataset comprising 24 multi-pixel objects. Each object contains both pure object pixels and object pixels mixed with background. For the Avon dataset, we first perform target detection, followed by manual selection of object pixels for annotation. For each class of object in the Avon dataset, we selected a pure object spectrum and a mixed object spectrum as a prior spectrum to perform hyperspectral target detection, respectively. We then manually selected pure object pixels and mixed object pixels from the two HTD score maps, respectively. The pixels selected from both score maps were deemed valid object pixels, and they were utilized to create the binary pixel mask label map for the Avon dataset. Traditional HTD methods typically use this binary pixel mask

TABLE I
THE OBJECT SIMULATION SETUP ON THE SPOD DATASET.

Name	Type	Endmember components	Object pixels	Maximum abundance	Training samples	Test samples
C1	Single	Cardboard	1	0.05-0.2	122	686
C2	Single	Malachite	1	0.05-0.2	118	674
C3	Single	Fiberglass	1-2	0.2-1	136	633
C4	Single	NickelPowder	3-5	0.2-1	120	664
C5	Single	LutetiumOxide	3-5	0.2-1	108	688
C6	Hybrid	Kerogen, NylonWebbing	6-10	0.2-1	122	669
C7	Combined	PlasticSheet, OiledMarshPlant	6-10	0.2-1	106	686
C8	Combined	Verdigris, YtterbiumOxide, WoodBeam	11-16	0.2-1	123	629

as the GT labels. However, for the point object detection task, this binary segmentation map must be further transformed into bboxes and category labels, as illustrated in Fig. 9. On the binary segmentation map of the specified category, we utilized eight-connected domain analysis to cluster the object pixels, considering each eight-connected domain as an independent object entity. The maximum enclosing rectangle of the eight-connected domain is used as the bbox, and each bbox is assigned the corresponding object category.

The prior information required for the hyperspectral point object detection is fundamentally different from that of traditional HTD. Traditional HTD require only a few target spectra (often just one), whereas the hyperspectral point object detection necessitates numerous training images with object labels. Classic HTD datasets like the SanDiego dataset do not come with a prepared training image set. To enable our SpecDETR to be evaluated on the Avon, SanDiego, and Gulfport datasets, we employed a data simulation method from the SPOD dataset. Using target prior spectra and background images without objects, we generated training image sets for these three public HTD datasets. Since the objects in these three public HTD datasets are all single-material, we used only one prior spectrum per object category to generate the simulated training samples. Since the training sets for these three datasets are simulated, the GT labels are precise, ensuring effective training of the object detection network.

B. The SPOD Dataset

Due to the spectral mixing phenomenon, acquiring a large number of high-quality manually annotated point objects is extremely challenging. Therefore, we develop a simulated high-spectral multi-class point object detection benchmark based on observed real-world patterns to evaluate the accuracy performance of various methods in point object detection. As shown in Fig. 10, we selected two flight lines¹ captured by the AVIRIS sensor over California and Nevada, covering 21,000 km² of the Earth's surface with significant geographic variation. One flight line provides 150 128×128 training background images, while the other provides 500 128×128 testing background images. After removing the atmospheric absorption bands, we retained 150 bands with the band IDs [8-57, 66-79, 86-103, 123-145, 146-149, 173-214].

¹<https://aviris.jpl.nasa.gov/>

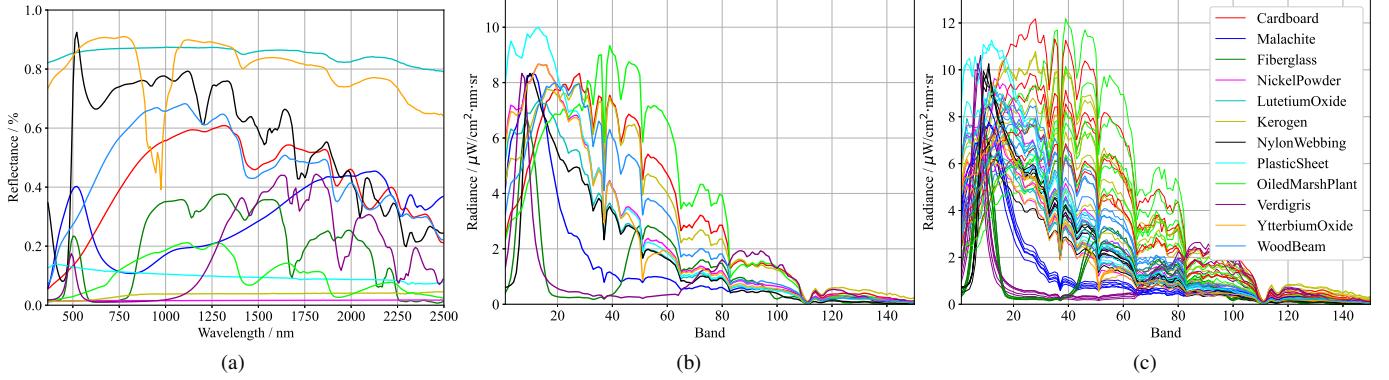


Fig. 12. Spectral curves of different endmembers. (a) Reflectance curves. (b) Radiance baselines. (c) Radiance curves after adding noise.

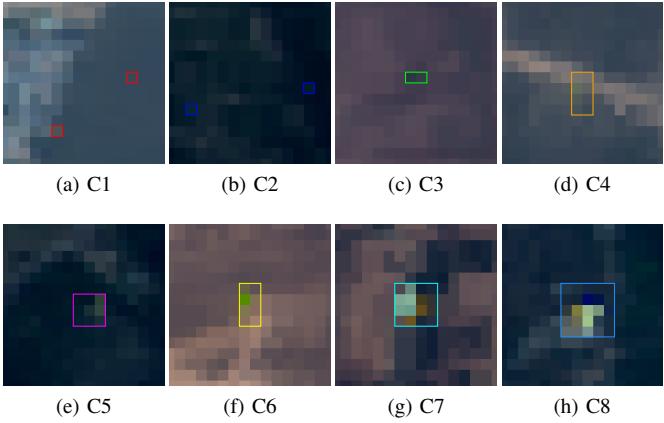


Fig. 13. Examples of objects on the SPOD dataset on pseudo-color images.

To ensure the realism of the simulation, we develop a simulation method based on the spatial characteristics of point objects and their spectral variation properties. This method is divided into three parts: object template generation, endmember spectrum generation, and simulated image synthesis. In the object template generation phase, we incorporate the point spread morphological characteristic of point objects and the co-existence of multiple endmember spectra. Additionally, since spectra of the same material can be affected by various factors such as illumination, we analyze the local spectral variation property and wide-area spectral variation property of water areas in AVIRIS images. These properties are used to generate simulated endmember spectra.

Object Template Generation: For each simulated object, we randomly select the pixel number N_P from the predefined range which is shown in Table I. Then, a clustered random combination of N_P pixels is generated to serve as a blank object template. Next, endmembers and abundances are assigned to each object pixel. Object pixels not adjacent to background pixels are considered pure pixels, with their object spectral abundance set to 1. Object pixels adjacent to background pixels are treated as mixed pixels, with their object spectral abundance randomly generated from the range (0.01, 1). The abundance of mixed pixels is assigned based on their distance to the object center, with closer pixels receiving

higher abundance values. Additionally, if all object pixels are mixed pixels, the abundance of the center pixel is randomly selected from the maximum abundance range, as specified in Table I. To account for the fact that real object spectra may contain multiple endmembers, we define three types of object templates: single-type, hybrid-type, and combined-type, as illustrated in Fig. 11. The single-type object template contains only one object endmember, so the object spectral abundance is equivalent to the endmember abundance. Hybrid-type and combined-type object templates contain two or more object endmembers. For the hybrid-type object template, each object pixel includes all types of object endmembers, and the object spectral abundance is randomly split into the abundances of each endmember. For the combined-type object template, each object pixel is assigned only one object endmember, and the object spectral abundance is equivalent to the endmember abundance. The combined-type object template represents a combination of multiple single-spectral objects.

Endmember Spectrum Generation: As shown in Fig. 11(a), we select 12 types of artificial object spectral reflectance curves r_t from the USGS spectral library [78] as initial object endmembers. Based on the radiance curve r_t of seawater and an average radiance curve s_w of a seawater region in the AVIRIS data, r_t is converted into an object endmember baseline s_t by:

$$s_t = \frac{M_t}{\max\left(\frac{r_t}{r_w} \circ s_w\right)} r_t \circ \frac{1}{r_w} \circ s_w, \quad (20)$$

where \circ denotes the pixel-wise product, M_t is the correction to ensure that the numerical range of the object endmember baseline is close to the HSIs to be synthesized, and M_t has slight differences for different endmember. Fig. 12(b) shows the spectral curves of endmember baselines. Then, we generate endmember spectra by adding noise into endmember baselines based on the statistical variation properties of water area spectra in AVIRIS data, as shown in Fig. 12(c). Spectra of the same object only exhibit the local fluctuation property, while different objects exhibit the wide-area variation property. Further details can be found in the Supplementary Material. **Simulated Image Synthesis:** For a simulated HSI to be synthesized, we first randomly select the number of object to be added from the range (1, 20). Each object is assigned a

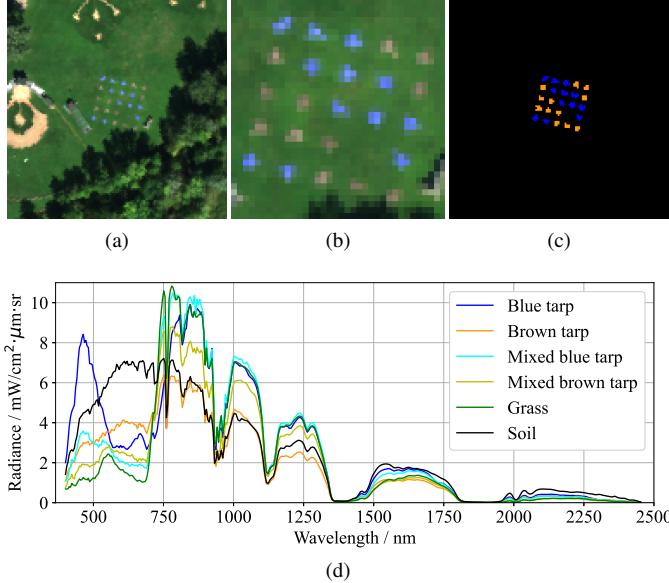


Fig. 14. The Avon dataset. (a) Pseudo-color image. (b) Zoomed-in image of the surrounding area of objects. (c) GT. (d) Radiance spectral curves of objects and background.



Fig. 15. Simulated training images of the Avon dataset.

category randomly chosen from the eight object types listed in Table I. For each object, we generate the corresponding object template and endmember spectra. The object template is then placed at a random location on the HSI, and the endmember spectra are injected into the HSI according to the object template. The injection process follows the linear mixing model [76], where the synthesized object spectrum is the abundance-weighted sum of the endmember spectra and the background pixel spectra.

As shown in Table I, we define eight types of point objects, including five single-type objects (labeled C1, C2, C3, C4, and C5), one hybrid-type object (labeled C6), and two combined-type objects (labeled C7 and C8). Among these, C1 and C2 are single-pixel objects with object spectral abundances smaller than 0.2. The categories and quantities of objects in each simulated image are randomly determined, with the maximum number of objects set to 20. Examples of each object type on pseudo-color images are shown in Fig. 13, where C1 and C2 are barely perceptible to the human eye. To increase the challenge, we average every 5 adjacent bands, reducing the 150 bands to 30 bands to reduce spectral information.

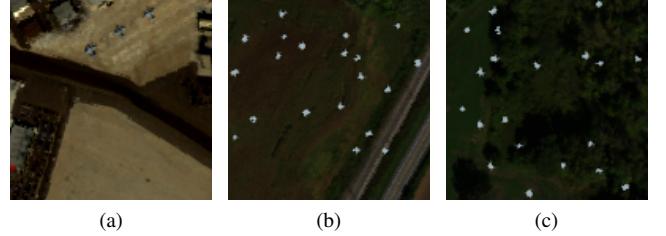


Fig. 16. The SanDiego dataset. (a) The test image. (b)(c) Simulated training images.

C. The Avon Dataset

The Avon dataset² was collected by a ProSpecTIR-VS sensor during the SHARE 2012 data campaign [79], providing hyperspectral images from multiple flights with different flight times. It has a spatial resolution of 1m and a spectral resolution of 5nm, with 360 bands between 400-2450nm. As shown in Fig. 14, we crop a test image of size 128×128 from the flight 0920-1701, which contains 12 blue tarps and 12 brown tarps. The spectra of the edge pixels of the tarps are a compromise between the pure tarp spectra and the grass background spectrum. Between 800-1200nm, the spectrum of brown tarps is similar to that of the soil background. We use the method described in Section IV-A to form the per-pixel annotations shown in Fig. 14(c) and convert them into the label format for point object detection. Additionally, we use the pure blue tarp spectrum and brown tarp spectrum from Fig. 14(d) as endmember spectra and cut background images from another flight 0920-1631, generating 150 training images. Each training image contains 10 simulated blue tarps and 10 simulated brown tarps, as shown in Fig. 15. To evaluate the impact of illumination on the performance of the proposed method, we extract the same area as the test image from four other flights and manually annotate the object labels.

D. The SanDiego Dataset

The SanDiego Dataset is captured by the AVIRIS over the SanDiego airport, with a size of 400400 pixels and a spatial resolution of 3.5 meters. As shown in Fig. 16(a), we extract a 128×128 test image focusing on 3 airplanes as objects. As depicted in Fig. 9, we convert the per-pixel annotations provided in literature [56] into the label format for point object detection. Since the SanDiego Dataset is the atmospherically corrected reflectance data, we use the reflectance data of the flight line 0920-1631 from the SHARE 2012 data campaign to generate 100 training images, each containing 10 airplanes, as illustrated in Fig. 16(b)(c). After removing bad bands and those not covered by the SHARE 2012 data, we retain 180 bands for the experiment.

E. The Gulfport Dataset

The Gulfport dataset³ was collected by a CASI-1500 sensor over the Gulf Park campus of the University of Southern

²<https://dirsapps.cis.rit.edu/share2012/>

³<https://github.com/GatorSense/MUUFLGulfport/>

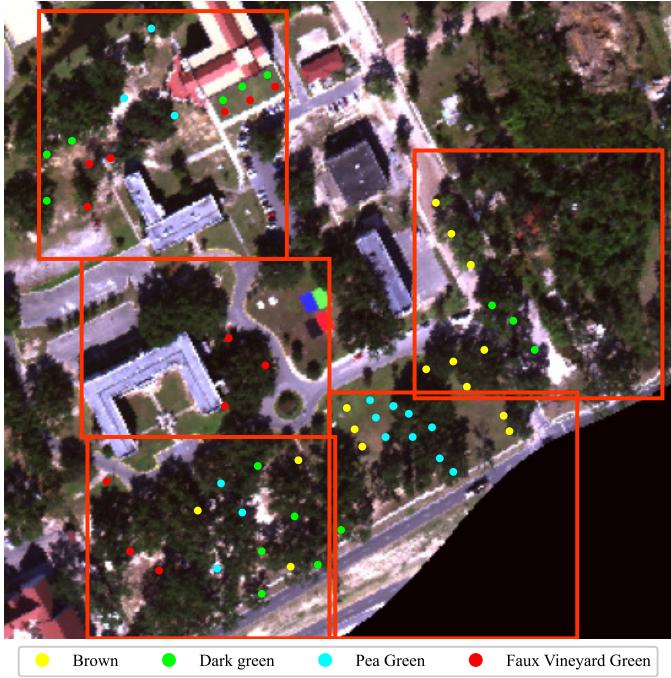


Fig. 17. The Gulfport Dataset. Red bboxes represent the cropped test images.

Mississippi, encompassing data from five flights. The images have a spatial resolution of 1 meter and consist of 72 bands between 367.7–1043.4 nm. This dataset includes 57 artificial objects made of cloth panels in four different colors: brown (15 panels), dark green (15 panels), pea green (15 panels), and faux vineyard green (12 panels). The object sizes include 3m×3m, 1m×1m, and 0.5m×0.5m. As depicted in Fig. 17, we crop five images from the flight 3 to serve as test images, each padded with zeros to a size of 128×128 pixels. As shown in Fig. 18, objects are classified into four confidence levels based on the ability of human annotators to identify them in the data: visible (Vis.), probably visible (Pro.Vis.), possibly visible (Pos.Vis.), and not visible (NotVis.), with counts of 9, 7, 8, and 33, respectively. The official GPS GT labels have a certain degree of positioning error, as can also be observed from Fig. 18. It is noted that there is a certain degree of position error in the GPS ground truth provided by the official sources, as can also be observed in Figure 3. Additionally, the flight 4 data where the object areas are masked is used to generate 200 training images. Each image contains five simulated panels of each type panels. The reference spectra for the simulated objects is the central spectra of the 3m×3m objects from the flight 3.

F. Evaluation Metrics

1) *Hyperspectral Point Object Detection*: In the evaluation of the general object detection task, the IoU is calculated as the area of overlap between the predicted bbox and the GT bbox, divided by the area of their union. A prediction is considered correct if the IoU between a predicted box and a GT box exceeds a specified threshold and the class prediction is accurate. Precision (Pr) and Recall (Re) are two fundamental metrics for evaluation. Pr represents the ratio of correctly

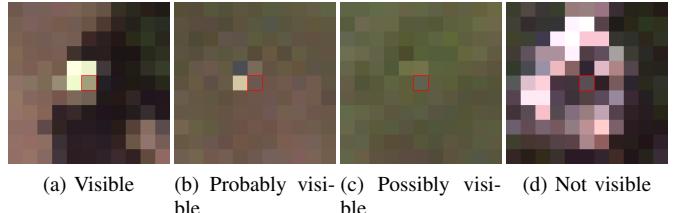


Fig. 18. Pea green panels at four confidence levels. Red bboxes represent the cropped test images.

detected objects to all detected objects, while Re indicates the ratio of correctly detected objects to the total number of GT objects.

Given the precise GT labels in the SPOD dataset and Avon dataset, we adopt the evaluation methodology of the COCO dataset [18] for these datasets, primarily reporting the average precision (AP) of each class, and mean AP (mAP) and mean average recall (mAR) across all classes. In the COCO evaluation, a fixed number of predictions are selected based on confidence scores for each test image, which we set to 100 in our experiments. For a specified object class, with a fixed IoU threshold, detection results are divided into positive and negative samples using different confidence thresholds. The precision and recall are calculated for each confidence threshold, and the precision-recall curve is plotted. The area under this curve represents the AP value for that class. The AP in COCO evaluation is typically the average of AP values over IoU thresholds ranging from 0.5 to 0.95, and the average recall (AR) is the average of recall values over the same IoU range. For the general object detection task, the object sizes are usually larger, and predictions with smaller IoU values tend to perform poorly in predicting object boundaries. However, for point objects, the ability to locate the object is far more important than predicting the object boundaries, and the bbox annotations of point object are prone to deviation, making smaller IoU thresholds also meaningful. Therefore, we also evaluate the mAP (mAP₂₅) and mean recall (mRe₂₅) at an IoU threshold of 0.25. In the SPOD dataset, subscripts C1 to C8 denote the eight distinct object classes, whereas in the Avon dataset, BL and BR signify the two tarp classes.

However, the GT labels for the test images in the SanDiego dataset and the Gulfport dataset are not precise. The labels for the SanDiego dataset are derived from widely used mask labels within the field. However, unlike the Avon dataset, the object edge pixels in the SanDiego dataset exhibit characteristics of nonlinear mixing. The training set for the SanDiego dataset is based on a linear mixing model simulation, therefore, the prediction boxes of the object detection networks tend to be smaller than the label boxes. For very small objects, even a one or two pixel deviation in the prediction box can lead to a significant drop in IoU. The labels for the Gulfport dataset are based on GPS point coordinates, which have a certain degree of location error and lack bbox information. Moreover, many objects in the Gulfport dataset are obscured, and may not have any information in the image, making manual annotation impossible. The COCO evaluation approach, which emphasizes

the fitting ability of bboxes, is not suitable for the SanDiego Dataset and the Gulfport dataset. To address this, we have fixed the criteria for correct predictions, no longer using the IoU traversal approach. For the SanDiego Dataset, a prediction is considered correct if the IoU between the predicted box and the GT box is greater than 0.25. For the Gulfport Dataset, we use a 5×5 rectangle centered on the pixel corresponding to the object's GPS location as the inner GT bbox, and a 9×9 rectangle as the outer GT bbox. A prediction is considered correct if the predicted box overlaps with the inner GT bbox and does not exceed the outer GT bbox. With the criteria for correct predictions fixed, we report the AP and Recall metrics for the SanDiego Dataset and the Gulfport Dataset.

2) *Hyperspectral Target Detection*: Current HTD methods typically produce detection score maps, from which segmentation result maps can be derived through thresholding. In contrast, the proposed SpecDETR provides bbox predictions for objects. To enable a direct comparison between SpecDETR and HTD methods, it is necessary to convert their predictions into a uniform format. Transforming from segmentation masks to instance bboxes is a well-posed problem, whereas the reverse process is ill-posed. Therefore, we convert the output of current HTD methods into bboxes for instance-level evaluation. Specifically, we normalize the score maps for each object class from every HTD compared method using min-max normalization, scaling them to the range [0, 1]. Subsequently, we perform binary segmentation on these score maps. Following the approach illustrated in Fig. 9, we convert the binary segmentation maps into predicted bboxes, setting the confidence score as the maximum score within the corresponding predicted bbox on the score maps. This transformation of HTD method predictions into object detection predictions allows for the application of object detection evaluation metrics.

Current HTD methods typically use the area under the ROC curve (AUC) [80] of the detection score maps as the evaluation metric, without considering the evaluation of segmentation maps. Consequently, many studies do not address the determination of segmentation thresholds. In the field of infrared small target detection, the segmentation IoU of the detection score map [81] is currently widely adopted, which is the ratio of overlapping pixels between the predicted target pixels and the GT target pixels to the number of pixels in their union. For the normalized score maps of HTD methods for each object class, we traverse the segmentation thresholds at a step size of 0.01, selecting the segmentation map with the highest IoU to maximize the effectiveness of HTD methods in object-level evaluation. Moreover, in our experiments on the SPOD and AVON datasets, we evaluate the HTD methods based on the mean AUC (mAUC) and mean segmentation IoU (mIoU) across all classes.

V. EXPERIMENTS

A. Comparison Experiments on the SPOD and Avon Dataset

1) *Implementation Details*: We use the MMDetection project [97] to implement SpecDETR. To ensure the diversity of generated samples within the denoising approach, the

parameters τ_1 and τ_2 for CCDN are set to 0.5 and 1.5, respectively, while the number of DN queries is set to 200 pairs. In the matching approach, to guarantee the quality of the assigned positive samples with the hybrid label assigner, the IoU threshold is set to 0.95, and the number threshold T is fixed at 9. The learning rate is initialized at 0.001, consistent with DINO, and is reduced by a factor of 0.1 at the 20th, 30th, and 90th epochs with the setting of 24, 36, and 100 training epochs, respectively. During the inference stage, we retain the top 300 prediction boxes with the highest confidence scores on each image and subsequently apply NMS to handle overlapping boxes. The IoU threshold for NMS is set to 0.01, indicating that the output predicted boxes are considered non-overlapping. All other parameters are kept identical to those used in DINO.

Given that the hyperspectral point object detection task developed in this paper spans the domains of HTD and visual object detection, the proposed SpecDETR is compared with both visual object detection networks and HTD methods on the SPOD dataset. Visual object detection networks typically use backbone networks to extract downsampled feature maps from input images. Directly inputting images of the original size without altering the network structure can be inefficient for extracting features of point objects, so we upscale the input HSIs by a factor of 4 by replicating each pixel 16 times before feeding them into the compared visual object detection networks. Additionally, the input channels of these networks are adjusted from 3 to match the band number of the input HSIs. The number of training epochs is set to 100 to ensure full convergence of the networks, while all other parameters are retained at their default settings. SpecDETR and the compared networks are trained with a batch size of 4 on an RTX 3090 GPU.

We compare SpecDETR with current HTD methods on the SPOD dataset and the Avon dataset, both of which provide fine, per-pixel annotations for multi-class objects. Current HTD methods take a test image along with prior spectral signatures as input and typically treat the predictions as a single class. To obtain multi-class predictions, the compared methods [27], [28], [31]–[33], [36], [37], [42], [43], [47], [56], [96] sequentially use the prior spectra of different object classes for detection, yielding predictions for each object class individually. LSSA [40] detects sub-pixel targets by directly learning the target abundances, so we adapt LSSA into a multi-class detector by inputting all classes' prior spectra and simultaneously outputting abundance predictions for all classes. Moreover, IRN [56] and TSTTD [47] are the SOTA deep learning-based HTD methods, and we also make improvements to them. The modified IRN and TSTTD are trained on a set of training images and then directly infer on the test images. The original version of TSTTD generates training samples by mixing prior spectral signatures with background spectra, whereas our improved version trains directly on the target spectra from the image set, ensuring that TSTTD and SpecDETR are exposed to an equivalent level of object prior information. The SPOD dataset contains multi-spectral objects thus providing 20 prior spectra for each object class. In contrast, the Avon dataset consists of single-spectral objects,

TABLE II

PERFORMANCE COMPARISON OF SPECDETR AND COMPARED VISUAL OBJECT DETECTION NETWORKS ON THE SPOD DATASET. SIZE $\times 4$ INDICATES THAT THE INPUT IMAGES ARE ENLARGED BY A FACTOR OF 4, WHILE $\times 1$ REPRESENTS THE ORIGINAL SIZE.

Detector	Backbone	Epochs	Size	mAP↑	mAP ₂₅ ↑	mAR↑	mRe ₂₅ ↑	AP _{C1} ↑	AP _{C2} ↑	AP _{C3} ↑	AP _{C4} ↑	AP _{C5} ↑	AP _{C6} ↑	AP _{C7} ↑	AP _{C8} ↑
Faster R-CNN [15]	ResNet50 [82]	100	$\times 4$	0.197	0.377	0.245	0.419	0.000	0.000	0.000	0.035	0.026	0.537	0.430	0.550
Faster R-CNN [15]	RegNetX [83]	100	$\times 4$	0.227	0.379	0.267	0.399	0.000	0.000	0.000	0.042	0.043	0.631	0.522	0.578
Faster R-CNN [15]	ResNeSt50 [84]	100	$\times 4$	0.246	0.316	0.269	0.319	0.000	0.000	0.000	0.008	0.003	0.644	0.564	0.747
Faster R-CNN [15]	ResNeXt101 [85]	100	$\times 4$	0.220	0.368	0.253	0.376	0.000	0.000	0.000	0.016	0.010	0.618	0.518	0.596
Faster R-CNN [15]	HRNet [86]	100	$\times 4$	0.320	0.404	0.364	0.434	0.000	0.000	0.000	0.107	0.076	0.849	0.731	0.793
TOOD [87]	ResNeXt101 [85]	100	$\times 4$	0.304	0.464	0.401	0.570	0.000	0.000	0.000	0.181	0.194	0.743	0.648	0.663
CentripetalNet [88]	HourglassNet104 [89]	100	$\times 4$	0.695	0.829	0.840	0.956	0.831	0.888	0.915	0.373	0.367	0.810	0.655	0.725
CornerNet [90]	HourglassNet104 [89]	100	$\times 4$	0.626	0.736	0.855	0.969	0.797	0.751	0.855	0.328	0.308	0.768	0.554	0.644
RepPoints [91]	ResNet50 [82]	100	$\times 4$	0.207	0.691	0.346	0.934	0.043	0.143	0.269	0.071	0.073	0.372	0.265	0.417
RepPoints [91]	ResNeXt101 [85]	100	$\times 4$	0.485	0.806	0.635	0.961	0.373	0.561	0.632	0.242	0.253	0.658	0.508	0.649
RetinaNet [92]	EfficientNet [93]	100	$\times 4$	0.462	0.836	0.611	0.971	0.566	0.602	0.530	0.182	0.210	0.566	0.471	0.569
RetinaNet [92]	PVTv2-B3 [94]	100	$\times 4$	0.426	0.757	0.650	0.987	0.356	0.478	0.458	0.209	0.232	0.563	0.470	0.644
DeformableDETR [24]	ResNet50 [82]	100	$\times 4$	0.231	0.692	0.385	0.883	0.230	0.316	0.234	0.077	0.070	0.289	0.238	0.395
DINO [25]	ResNet50 [82]	100	$\times 4$	0.168	0.491	0.368	0.763	0.020	0.047	0.277	0.080	0.064	0.286	0.213	0.360
DINO [25]	Swin-L [95]	100	$\times 4$	0.757	0.852	0.909	0.983	0.915	0.912	0.951	0.483	0.497	0.847	0.728	0.721
SpecDETR	-	100	$\times 1$	0.856	0.938	0.897	0.955	0.963	0.969	0.970	0.698	0.648	0.905	0.844	0.850

TABLE III

COMPUTATIONAL EFFICIENCY COMPARISON OF SPECDETR AND COMPARED VISUAL OBJECT DETECTION NETWORKS ON THE SPOD DATASET.

Detector	Backbone	Image Size	Training time (min)	Inference Speed (FPS)	FLOPs (G)	Params (M)	GPUMem. (MB)
Faster R-CNN [15]	ResNet50 [82]	$\times 4$	8.8	40.3	68.8	41.5	309
Faster R-CNN [15]	RegNetX [83]	$\times 4$	9.0	33.2	57.7	31.6	270
Faster R-CNN [15]	ResNeSt50 [84]	$\times 4$	12.8	23.9	185.1	44.6	467
Faster R-CNN [15]	ResNeXt101 [85]	$\times 4$	14.6	26.3	128.4	99.4	528
Faster R-CNN [15]	HRNet [86]	$\times 4$	13.9	17.4	104.4	63.2	407
TOOD [87]	ResNeXt101 [85]	$\times 4$	17.3	11.2	114.3	97.7	485
CentripetalNet [88]	HourglassNet104 [89]	$\times 4$	30.2	14.7	501.3	205.9	1098
CornerNet [90]	HourglassNet104 [89]	$\times 4$	25.0	16.3	462.6	201.1	981
RepPoints [91]	ResNet50 [82]	$\times 4$	10.4	33.7	54.1	36.9	238
RepPoints [91]	ResNeXt101 [85]	$\times 4$	17.6	14.7	75.0	58.1	329
RetinaNet [92]	EfficientNet [93]	$\times 4$	9.0	30.5	36.1	18.5	215
RetinaNet [92]	PVTv2-B3 [94]	$\times 4$	13.4	22.3	71.3	52.4	313
DeformableDETR [24]	ResNet50 [82]	$\times 4$	13.8	24.6	58.7	41.2	252
DINO [25]	ResNet50 [82]	$\times 4$	17.3	21.3	86.3	47.6	301
DINO [25]	Swin-L [95]	$\times 4$	53.8	10.1	203.9	218.3	1168
SpecDETR	-	$\times 1$	17.8	24.7	139.7	16.1	287

with only 1 prior spectrum provided for each object class.

In the absence of further specifications, SpecDETR defaults to 100 training epochs on the SPOD dataset and 36 epochs on the Avon dataset. Since the comparison on the Avon dataset is exclusively with HTD methods, the output bounding boxes produced by SpecDETR are rounded to the nearest integer.

2) *Result Analysis:* Table II provides the performance comparison between SpecDETR and the compared visual object detection networks in point object detection. SpecDETR achieves the highest mAP of 0.856 after 100 epochs among all compared networks. SpecDETR exhibits remarkable detection capability for subpixel objects and achieves AP_{C1}, AP_{C2}, and AP_{C3} of 0.963, 0.969, and 0.97, respectively. Table II further demonstrates the viability of current visual object detection networks in directly extracting spatial-spectral joint

features from hyperspectral cube data. Upon substituting the input of these networks from a three-channel RGB image to the complete hyperspectral cube, they continue to operate effectively without any structural alterations. Notably, both C1 and C2, which occupy merely a single pixel and possess object abundances below 0.2, exhibit AP performance that is remarkably close to that of multi-pixel objects like C7 and C8 across several compared networks. This observation suggests that the features extracted by these networks indeed incorporate spectral information. We also observe that those backbones that can offer high-resolution shallow features are more valuable in point object detection. Followed by Faster R-CNN, HRNet, which achieves multi-scale feature interaction, outperforms ResNeXt. CentripetalNet and CornerNet with HourglassNet achieve 0.626 mAP and 0.691 mAP, respec-

TABLE IV
PERFORMANCE COMPARISON OF SPECDETR AND COMPARED HTD METHODS ON THE SPOD DATASET.

Method	mAUC↑	mIoU↑	mAP↑	mAP ₂₅ ↑	mAR↑	mRe ₂₅ ↑	AP _{C1} ↑	AP _{C2} ↑	AP _{C3} ↑	AP _{C4} ↑	AP _{C5} ↑	AP _{C6} ↑	AP _{C7} ↑	AP _{C8} ↑
ASD [31]	0.991	0.359	0.182	0.286	0.712	0.939	0.061	0.068	0.587	0.078	0.065	0.298	0.101	0.194
CEM [28]	0.981	0.166	0.040	0.122	0.262	0.672	0.011	0.062	0.215	0.009	0.005	0.008	0.010	0.002
OSP [32]	0.969	0.159	0.031	0.108	0.202	0.594	0.008	0.039	0.181	0.002	0.002	0.006	0.009	0.000
KOSP [42]	0.961	0.116	0.017	0.083	0.165	0.539	0.000	0.002	0.119	0.002	0.002	0.009	0.001	0.002
SMF [27]	0.904	0.039	0.003	0.016	0.078	0.301	0.000	0.000	0.020	0.000	0.000	0.000	0.000	0.000
KSMF [43]	0.891	0.041	0.003	0.015	0.059	0.206	0.000	0.000	0.024	0.000	0.000	0.000	0.000	0.000
TCIMF [33]	0.975	0.096	0.009	0.061	0.160	0.547	0.004	0.001	0.052	0.003	0.001	0.004	0.005	0.003
CR [96]	0.939	0.160	0.031	0.138	0.242	0.685	0.008	0.000	0.160	0.017	0.015	0.016	0.020	0.012
KSR [36]	0.877	0.140	0.015	0.091	0.133	0.530	0.008	0.001	0.094	0.002	0.003	0.006	0.004	0.001
KSRBBH [39]	0.785	0.144	0.007	0.087	0.105	0.466	0.000	0.000	0.000	0.007	0.007	0.028	0.008	0.010
LSSA [40]	0.702	0.119	0.041	0.093	0.163	0.302	0.001	0.002	0.272	0.001	0.000	0.051	0.000	0.000
IRN [56]	0.625	0.004	0.000	0.002	0.005	0.033	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
TSTTD [47]	0.964	0.056	0.044	0.057	0.745	0.918	0.005	0.022	0.207	0.017	0.018	0.020	0.013	0.049
SpecDETR	-	-	0.856	0.938	0.897	0.955	0.963	0.969	0.970	0.698	0.648	0.905	0.844	0.850

tively.

Fig. 19 presents some visualization results of SpecDETR. It is evident that SpecDETR demonstrates effective detection of objects that are challenging for human eyes to identify in the SPOD dataset. However, it is also shown that SpecDETR can accurately localize objects of C4 and C5 but sometimes provides incorrect class predictions, resulting in lower AP_{C4} and AP_{C5}. Enhancing the network’s ability to distinguish between hard-to-distinguish object categories remains a critical direction for future research.

Table III presents a quantitative evaluation of the computational efficiency performance of SpecDETR and compared visual object detection networks. The floating-point operations (FLOPs) are used to quantify the time complexity of the object detection networks, while the parameters (Params) and memory consumption are employed to assess the spatial complexity. Here, we report the maximum GPU memory occupied during the inference of a single image⁴. Additionally, we provide the training time with 100 training epochs and the inference speed which is the reciprocal of the single-image test time for each network. Owing to the removal of the backbone and the optimization of the decoder structure, SpecDETR has a relatively small number of Params, totaling only 16.1M, which is superior to the compared visual object detection networks. Although SpecDETR exhibits better performance in terms of spatial complexity, its performance in terms of time complexity is mediocre. Although the input images of SpecDETR are not upsampled as in other networks, SpecDETR achieves 24.7 FPS and 139.7 GFLOPs, whereas Faster R-CNN with ResNet50 achieves 40.3 FPS and 68.8 GFLOPs. This is primarily due to two reasons: first, SpecDETR is based on a Transformer architecture, which is inherently at a disadvantage in terms of computational speed compared to pure CNN-based networks; second, SpecDETR does not downsample the feature maps during the feature extraction process. Enhancing the computational efficiency of the proposed point object detection network could be a promising direction for future

⁴https://pytorch.org/docs/stable/generated/torch.cuda.max_memory_allocated.html

TABLE V
PERFORMANCE COMPARISON OF SPECDETR AND HTD METHODS ON THE AVON DATASET.

Method	mAUC↑	mIoU↑	mAP↑	mAP ₂₅ ↑	mAR↑	mRe ₂₅ ↑	AP _{BL} ↑	AP _{BR} ↑
ASD [31]	0.961	0.654	0.324	0.853	0.538	1.000	0.312	0.336
CEM [28]	0.967	0.460	0.178	0.566	0.471	1.000	0.218	0.139
OSP [32]	0.967	0.457	0.164	0.526	0.471	1.000	0.202	0.126
KOSP [42]	0.977	0.374	0.111	0.380	0.454	0.958	0.132	0.090
SMF [27]	0.957	0.155	0.047	0.240	0.350	0.958	0.002	0.092
KSMF [43]	0.965	0.211	0.068	0.229	0.463	0.958	0.064	0.073
TCIMF [33]	0.982	0.433	0.181	0.456	0.538	1.000	0.220	0.142
CR [96]	0.806	0.313	0.013	0.551	0.046	0.875	0.023	0.002
KSR [36]	0.944	0.269	0.003	0.329	0.008	0.750	0.004	0.001
KSRBBH [39]	0.868	0.495	0.181	0.578	0.483	1.000	0.148	0.213
LSSA [40]	0.923	0.388	0.097	0.407	0.292	0.833	0.187	0.008
IRN [56]	0.750	0.391	0.302	0.500	0.367	0.500	0.605	0.000
TSTTD [47]	0.951	0.463	0.312	0.570	0.579	1.000	0.560	0.065
SpecDETR	-	-	0.885	0.990	0.925	1.000	0.924	0.847

exploration.

On the Avon dataset, SpecDETR demonstrates superior performance in real-world single-spectral point object detection, achieving 0.885 mAP and 0.925 mAR. As illustrated in Tables IV and V, SpecDETR exhibits a significant advantage over the compared HTD methods in object-level evaluation on both the SPOD and Avon datasets. In the conversion process from HTD results to object detection outcomes, we apply the optimal segmentation threshold for each object class to ensure the recall of objects. Some HTD methods achieve high mAR₂₅, but perform poorly in terms of mAP₂₅. This is because current HTD methods primarily rely on the one-dimensional spectral information of objects, which tends to produce high detection responses in background regions with similar spectral characteristics or in other classes of similar objects. As shown in Fig. 20, most HTD methods generate numerous false alarms on soil, which has spectra similar to brown tarps. The poor performance in mAP₂₅ is actually due to some false alarms having higher detection responses than true positives. In contrast, SpecDETR can learn the spatial-spectral joint features of objects from a large number of training samples, which is superior in the dimension of object feature perception compared to existing HTD methods. As



Fig. 19. Visualization Result Examples of SpecDETR on the SPOD dataset. We filter out predictions with confidence scores below 0.2 and set the IoU threshold for GT matching to 0.25. Green boxes indicate true positive predictions are denoted, red boxes indicate false positive predictions, while yellow boxes or circles indicate undetected objects. Additionally, red circles represent cases where the localization is accurate, but the category prediction is incorrect, also classified as false negatives.

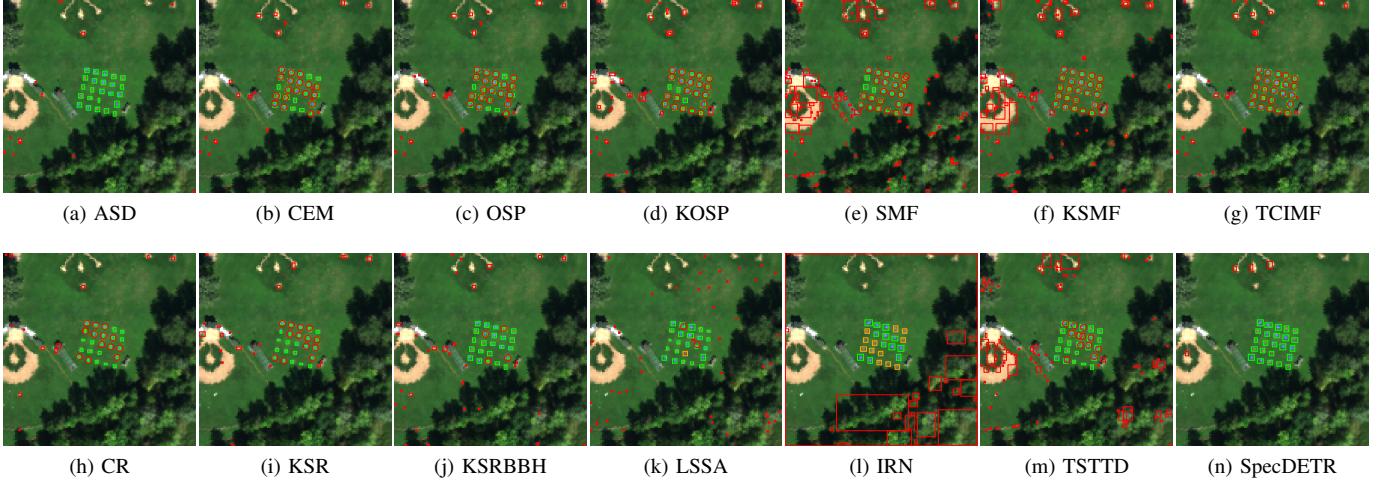


Fig. 20. Visualization results of SpecDETR and HTD methods on the Avon dataset. The confidence threshold is 0, the IoU threshold is 0.01, and the remaining visualization settings are consistent with Fig. 19.

TABLE VI

PERFORMANCE COMPARISON OF SPECDETR AND COMPARED VISUAL OBJECT DETECTION NETWORKS ON THE SANDIEGO AND GULFPORT DATASETS. THE BEST RESULTS ARE MARKED IN RED, AND THE SECOND-BEST RESULTS ARE MARKED IN BLUE. B., D.G., P.G., AND F.V.G. REPRESENT BROWN, DARK GREEN, PEA GREEN, AND FAUX VINEYARD GREEN CLOTH PANELS, RESPECTIVELY.

Method	SanDiego		Gulfport									
			Mean		B.		D.G.		P.G.		F.V.G.	
	AP	Re	AP	Re	AP	Re	AP	Re	AP	Re	AP	Re
CentripetalNet	0.776	1.000	0.178	0.213	0.208	0.200	0.192	0.267	0.139	0.133	0.171	0.250
CornerNet	0.750	1.000	0.160	0.271	0.203	0.267	0.069	0.133	0.242	0.267	0.125	0.417
DINO	0.210	1.000	0.174	0.408	0.220	0.333	0.232	0.400	0.166	0.400	0.077	0.500
SpecDETR	1.000	1.000	0.213	0.321	0.277	0.333	0.191	0.200	0.218	0.333	0.167	0.417



Fig. 21. Visualization results of SpecDETR and compared visual object detection networks on the SanDiego dataset. The confidence threshold is 0.09, object confidence scores are shown, and the remaining visualization settings are consistent with Fig. 19.

shown in Fig. 20(n), after NMS processing of the predictions, SpecDETR reports only 4 false alarms on the soil. These results demonstrate that the point object detection framework we propose is superior to the current per-pixel classification framework based on prior spectral information.

B. Experiments on the SanDiego and Gulfport Datasets

We select the three methods that performed best on the SPOD dataset, CentripetalNet, CornerNet, and DINO, for further comparative evaluation with SpecDETR on the SanDiego and Gulfport datasets. The training epochs for the SanDiego and Gulfport datasets are set to 12 and 24, respectively, with the learning rate reduced by a factor of 0.1 at the 10th and

20th epochs, respectively. All other settings are kept consistent with those used for the SPOD dataset.

Table VI presents the performance comparison of SpecDETR and the compared visual object detection networks on the SanDiego and Gulfport datasets. SpecDETR achieves an AP of 1.0 and a Re of 1.0 on the SanDiego dataset, while CentripetalNet, CornerNet, and DINO achieve AP values of 0.750, 0.776, and 0.210, respectively. Additionally, SpecDETR achieves the best mean AP of 0.213 and the second-best mean Recall of 0.321 on the Gulfport dataset. Fig. 21 shows the visualization results of SpecDETR and the compared visual object detection networks on the SanDiego dataset. Since the lowest confidence score for object prediction boxes among

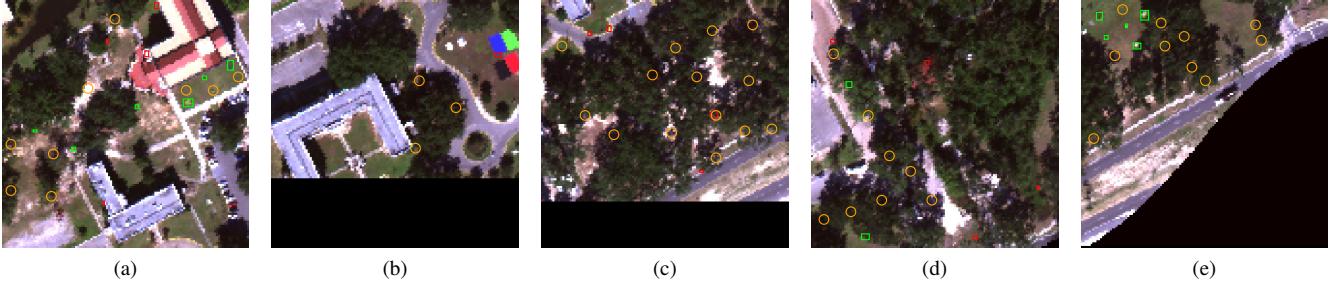


Fig. 22. Visualization results of SpecDETR on the Gulfport dataset. Visualization settings are consistent with Fig. 19.

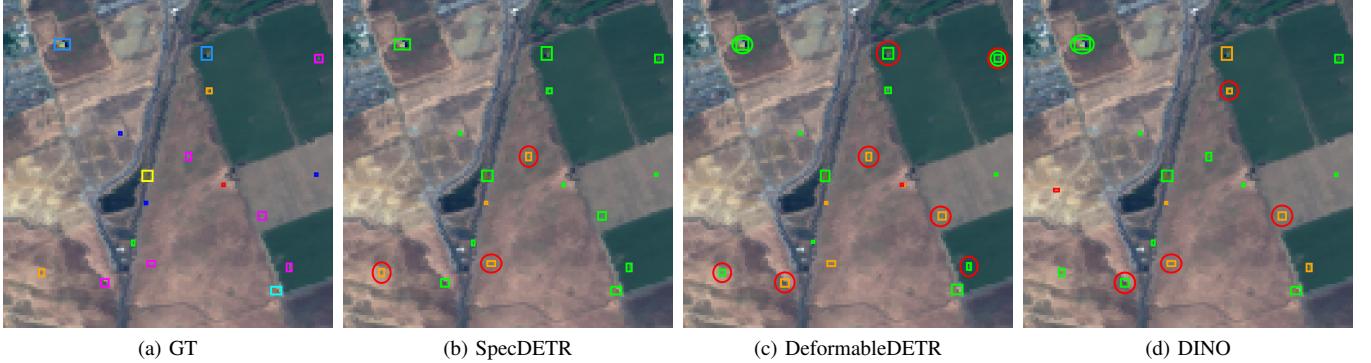


Fig. 23. Visual comparison of SpecDETR and other DETR-like detectors. (a) Ground truth labels, with color meanings consistent with Fig. 13. (b)-(d) Visualization of results, with visualization settings consistent with Fig. 19. Green circles indicate redundant true positives which are actually considered as false positives in the evaluation.

TABLE VII

THE PERFORMANCE OF SPECDETR ON THE GULFPORT DATASET UNDER DIFFERENT CONFIDENCE THRESHOLDS. NUM. REPRESENTS THE NUMBER OF PREDICTIONS AFTER CONFIDENCE FILTERING, AND TP STANDS FOR TRUE POSITIVES.

Confidence Threshold	Total				Vis.		Pro.Vis.		Pos.Vis.		NotVis.	
	Num.	TP	Pr	Re	TP	Re	TP	Re	TP	Re	TP	Re
0	380	18	0.047	0.316	9	1.000	4	0.571	3	0.375	2	0.061
0.05	97	15	0.155	0.263	9	1.000	4	0.571	1	0.125	1	0.030
0.1	52	15	0.289	0.263	9	1.000	4	0.571	1	0.125	1	0.030
0.2	26	13	0.500	0.228	8	0.889	4	0.571	1	0.125	0	0.000
0.3	17	12	0.706	0.211	7	0.778	4	0.571	1	0.125	0	0.000
0.4	13	9	0.692	0.158	6	0.667	3	0.429	0	0.000	0	0.000
0.5	9	6	0.667	0.105	4	0.444	2	0.286	0	0.000	0	0.000

the four methods is 0.092, Fig. 21 filters out prediction boxes with confidence scores below 0.09. Although all four methods can detect three airplanes, the other three comparison methods have false predictions with higher confidence scores than the object prediction boxes. Furthermore, despite the prediction boxes being filtered by confidence, DINO has many repeated predictions, indicating that the current DETR-like detectors' one-to-one prediction is not very suitable for hyperspectral point object detection. In contrast, SpecDETR successfully detects all three airplanes with confidence scores of 0.327, 0.314, and 0.351, respectively. Other predictions output by SpecDETR, with confidence scores below 0.05, are filtered out in Fig. 21.

Drawing inspiration from the task of moving object detection in satellite videos [98], Table VII presents the precision and recall performance of SpecDETR on the Gulfport

dataset under fixed confidence thresholds. At a confidence threshold of 0.1, SpecDETR detects 15 objects, achieving a Re of 0.26.3 and a Pr of 0.289. When the confidence threshold is raised to 0.3, SpecDETR detects 12 objects, with a Re of 0.706 and a Pr of 0.105. Fig. 22 visualizes the detection results of SpecDETR on the Gulfport dataset, where a large number of undetected objects were obscured by trees or shadows. Consequently, in Table VII, we analyze the detection performance of SpecDETR under different levels of object observational confidence. It is observed that the lower the object observational confidence level, the lower the Re of SpecDETR for detection. At a 0.1 confidence threshold, SpecDETR detects all 9 Vis. objects, 4 out of 7 Pro.Vis. objects, but only 1 each for the Pos.Vis. and NotVis. objects.

The predicted confidence scores for objects by SpecDETR are predominantly below 0.5, which can be attributed to the domain gap between the simulated training images and the real test images. Moreover, the extremely low Re for Pos.Vis. and NotVis. objects may be due, in part, to some objects being completely obscured. Additionally, it is possible that the linear mixing model on which our simulation is based does not accurately represent the actual imaging phenomena when objects are obscured by trees or shadows.

C. Model Analysis

1) *Robustness Analysis under Varying Lighting Conditions:* In the experiments on the Avon dataset in Section IV-A, the imaging time for the training data was at 11:31, while the test image was captured at 12:01. Additionally, we extract



Fig. 24. Visualization results of SpecDETR on the Avon dataset under varying lighting conditions. The confidence threshold is 0, and the remaining visualization settings are consistent with Fig. 19.

TABLE VIII

DETECTION PERFORMANCE OF SPECDETR ON THE AVON DATASET UNDER VARYING LIGHTING CONDITIONS.

Imaging time	Mean		Blue Tarp		Brown Tarp	
	mAP ₂₅	mRe ₂₅	AP ₂₅	Re ₂₅	AP ₂₅	Re ₂₅
12:01	0.990	1.000	1.000	1.000	0.981	1.000
12:06	0.974	1.000	1.000	1.000	0.949	1.000
13:44	0.981	1.000	1.000	1.000	0.962	1.000
13:51	0.987	1.000	1.000	1.000	0.974	1.000
13:57	1.000	1.000	1.000	1.000	1.000	1.000

TABLE IX

DETECTION PERFORMANCE OF SPECDETR ON THE AVON DATASET UNDER DIFFERENT NOISE LEVELS.

Noise level	Mean		Blue Tarp		Brown Tarp	
	mAP ₂₅	mRe ₂₅	AP ₂₅	Re ₂₅	AP ₂₅	Re ₂₅
SNR _{dB} =10	0.255	0.458	0.477	0.583	0.032	0.333
SNR _{dB} =15	0.416	0.708	0.752	0.750	0.079	0.667
SNR _{dB} =20	0.520	0.833	0.911	0.917	0.128	0.750
SNR _{dB} =25	0.481	0.917	0.842	0.917	0.120	0.917
SNR _{dB} =30	0.614	0.958	1.000	1.000	0.227	0.917
SNR _{dB} =35	0.761	1.000	1.000	1.000	0.523	1.000
w.o. noise	0.990	1.000	1.000	1.000	0.981	1.000

the same area as the test image from four other flights with imaging times at 12:06, 13:44, 13:51, and 13:57, respectively, to analyze and assess the robustness of SpecDETR under varying lighting conditions. Both the training and test images consist of uncorrected radiometric data, which are susceptible to solar illumination effects. The training and test images both consist of radiance data without atmospheric correction, which are influenced by solar illumination. As shown in Table VIII, SpecDETR achieves a mRe₂₅ of 1.0 and above 0.97 mRe₂₅ across the five test images with different lighting conditions, and it also attains a 1.0 AP₂₅ and 1.0 Re₂₅ for blue tarps. Figs. 20(n) and 24 present the visualization results, showing only a minimal number of false positives for brown tarps on soil. The robustness of SpecDETR to lighting conditions suggests that the differences in data distribution under various lighting conditions on the Avon dataset have a negligible impact on the point object detection task.

2) *Robustness Analysis under Different Noise Levels:* Following the classic work [99] in the field, we add additive

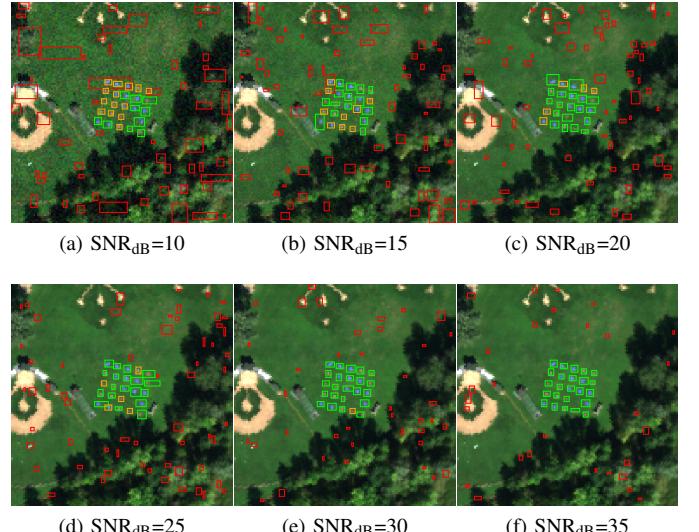


Fig. 25. Visualization results of SpecDETR on the Avon dataset under different noise levels. The visualization settings are consistent with Fig. 19.

Gaussian white noise to the Avon test image to assess the robustness of SpecDETR under different noise levels. As described in literature [99], zero-mean Gaussian noise with different variances are added to each spectral band, with the noise level quantified by SNR_{dB}, defined as:

$$\text{SNR}_{\text{dB}} = 10 \log_{10} \left(\frac{\sigma_{\text{signal},i}^2}{\sigma_{\text{noise},i}^2} \right), \text{ for } i = 1, 2, 3, \dots, C, \quad (21)$$

where i denotes the band index, $\sigma_{\text{signal},i}^2$ is the variance of the i -th band, and $\sigma_{\text{noise},i}^2$ is the variance of the added Gaussian noise. We set SNR_{dB} to 10 dB, 15 dB, 20 dB, 25 dB, 30 dB, and 35 dB, respectively, and add the corresponding $\sigma_{\text{noise},i}^2$ Gaussian noise into the Avon test image. Subsequently, we employ SpecDETR, trained on the noise-free training set, for direct inference without fine-tuning. Table IX provides a quantitative evaluation of the results, while Fig. 25 visualizes the predicted bboxes with confidence scores above 0.2. It can be observed that the detection performance is positively correlated with the magnitude of the added noise variance. Moreover, SpecDETR demonstrates superior noise robustness

for blue tarps compared to brown tarps, attributed to the greater spectral distinctiveness of blue tarps from the background. At SNR_{dB} of 10 dB, SpecDETR achieves a 0.477 AP₂₅ for blue tarps, which improves to 0.911 AP₂₅ when SNR_{dB} is 20 dB. However, at SNR_{dB} of 30 dB, SpecDETR attains only a 0.227 AP₂₅ for brown tarps. As depicted in Fig. 25, the introduction of Gaussian noise to the test image not only leads to false positives on soil but also on grass. This occurs because the data distribution of the noisy test images significantly diverges from that of the noise-free training images. The domain gap between training and test sets is a well-recognized challenge in the field of object detection and has given rise to the cross-domain few-shot object detection task [100].

3) *Analysis of the Necessity of Non-Maximum Suppression:* Fig. 23 presents the detection results of SpecDETR, Deformable DETR, and DINO on a test HSI from the SPOD dataset. Despite being end-to-end detectors, both Deformable DETR and DINO still exhibit redundant predictions, indicating the necessity of reintroducing NMS in SpecDETR.

4) *Feature Extraction in SpecDETR Encoder:* Following the conventional setup of DETR-like detectors, the number of SpecDETR encoder layers is set to 6 by default. Fig. 26 visualizes the feature maps output by each encoder layer of SpecDETR for Fig. 23(a). As the number of encoder layers increases, the background area with complex texture structure becomes progressively smoother, while objects emerge more distinctly from the background. Even for objects with extremely low abundance, C1 and C2, SpecDETR generates pronounced responses on the feature maps output by the final encoder layer. Moreover, Fig. 27 illustrates the feature maps output by the encoder of SpecDETR on various background test images, where the objects also exhibit clear features. This demonstrates that SpecDETR can effectively extract the correct object features rather than noise or proxy features. Additionally, in scenes with simpler background types such as farmland and woodland, the background regions in the feature maps are smooth and clean. However, in complex scenes like urban areas, the feature maps generate many isolated response points, which leads to some false alarms in the urban background regions as shown in Fig. 19.

5) *Analysis of the Self-S2A Module:* We visualize the sampling point distributions of the self-S2A module when using two object pixels and a background pixel as query elements in Fig. 23(a). The two object pixels belong to C1 and C8, respectively, with the former being a low-abundance single-pixel object and the latter a high-abundance pixel within a multi-pixel object. The selected background pixel shares the same background component class as the C1 pixel, resulting in similar spectral characteristics between the background pixel and the C1 pixel. Fig. 28 and Fig. 29 visualize the sampling point distributions of the subpixel-scale deformable sampling operator in the first and last encoder layers, respectively. As shown in Fig. 28, in the first encoder layer, the sampling point distributions for different types of query elements is largely consistent. Several sampling points are concentrated around the pixel location of the query element, while the remaining points are evenly distributed in the local spatial neighborhood. In contrast, Fig. 29 illustrates that in the last encoder layer, the

sampling point distribution varies significantly among different types of query elements. The sampling points for the C1 pixel are mainly focused around its own pixel location, whereas those for the C8 pixel are primarily distributed along the object boundary. This indicates that the pixel tokens updated in the penultimate encoder layer already possess distinct class information. Fig. 30 and Fig. 31 depict the sampling point distributions of the self-excited operator in the first and last encoder layers, respectively. As depicted in Fig. 30, in the first encoder layer, the attention heads for both the C1 pixel and the background pixel remain inactive, while the green attention head for the C8 pixel is activated. This is due to the lower object abundance of the C1 pixel compared to the higher abundance of the C8 pixel. However, as shown in Fig. 31, in the last encoder layer, the attention heads for the background pixel remain inactive, while both the C1 and C8 pixels have activated attention heads. This suggests that the preceding multiple encoder layers effectively extract the object features of C1. The output vectors are generated through concatenation of the sampling vectors from different attention heads, and the C1 and C8 pixels exhibit different combinations of activated attention heads. Consequently, the output vectors of the self-excited operator for the C1 and C8 pixels are significantly distinct, as illustrated in Fig. 32. Fig. 33 presents the sum of the absolute values across all channels of the output feature maps from these two sampling operators. In the output map of the self-excited operator, the background regions are either zero or close to zero, while high-response regions typically correspond to objects. This indicates that the self-excited operator helps object pixels to amplify their own features.

D. Ablation Experiments

1) *Network Architecture:* DINO stands as one of the most successful object detectors to date. Co-DETR [63], which inherits the DINO decoder structure, achieves the highest mAP of 0.66 on the COCO test-dev dataset⁵. Our experiments on the SPOD dataset also demonstrate that DINO, equipped with Swin-L as the backbone network, significantly outperforms other compared object detectors in terms of detection accuracy. SpecDETR builds upon DINO, specifically tailored for the point object detection task, by refining the feature extraction pattern and decoder structure. As shown in Table X, we compare the performance of different network module combinations on the SPOD dataset to validate the effectiveness of our improvements. DINO employs a backbone network to initially extract multi-scale feature maps from the input images, and then uses the Transformer encoder module for further feature extraction. In contrast, SpecDETR dispenses with the backbone network, directly utilizing the Transformer encoder to extract deep feature maps that match the input image size. While maintaining the DINO decoder unchanged, replacing Swin-L and the DINO encoder with the SpecDETR encoder results in an mAP increase from 0.391 to 0.663 under 24 training epochs, and from 0.757 to 0.849 under 100 training epochs. Additionally, the number of parameters decreases from

⁵<https://paperswithcode.com/sota/object-detection-on-coco>

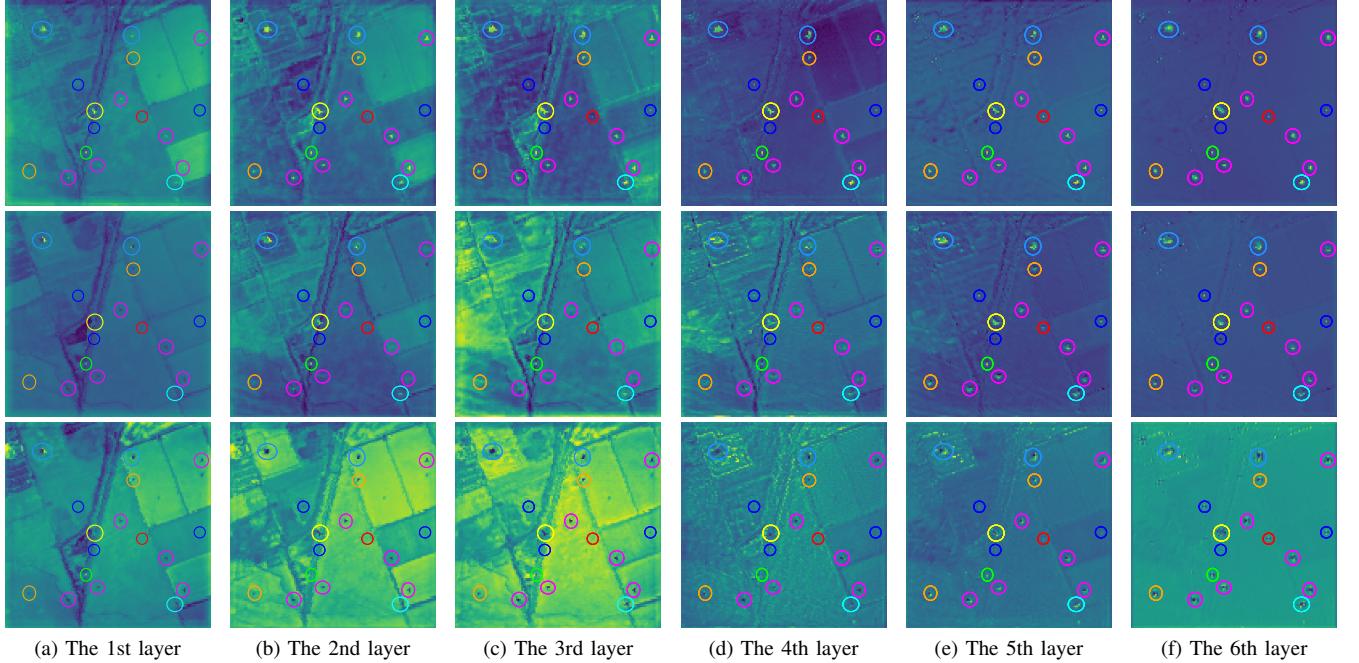


Fig. 26. Feature maps output by each encoder layer of SpecDETR on Fig. 23(a). From top to bottom, the channels 50, 100, and 150 of the output feature maps are displayed, with objects marked by circles. The meaning of the circle colors is consistent with Fig. 13.

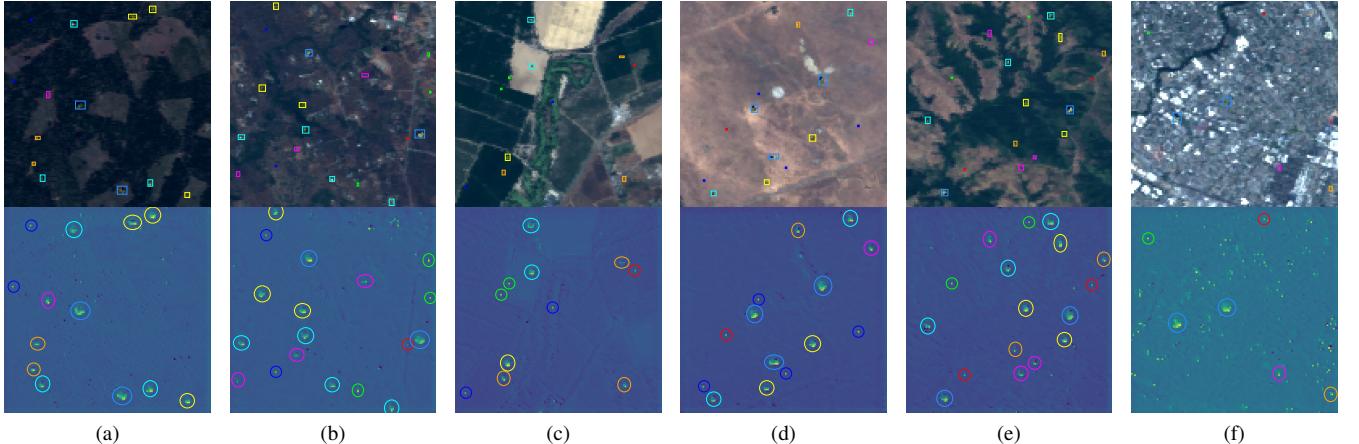


Fig. 27. Feature maps output by the encoder of SpecDETR on some test images from the SPOD dataset. The upper part shows the color images of the input, while the lower part displays the feature maps. The meaning of the circle and box colors is consistent with Fig. 13.

218.3M to 17.9M. However, due to the absence of downsampling in the feature maps extracted by the SpecDETR encoder, the advantage in FLOPs is not as pronounced as that in parameters. Furthermore, replacing the DINO decoder with the SpecDETR decoder also validates the superior performance of the SpecDETR decoder for the point object detection task through various quantitative metrics. Compared to the DINO decoder, the SpecDETR decoder achieves a 6.5% and 5.3% increase in mAP at 24 and 36 epoch settings, respectively. Meanwhile, there is a 2.7% reduction in FLOPs and a 10.6% decrease in parameters.

2) *Network Scale*: SpecDETR primarily consists of two modules: the Transformer encoder and decoder. In the default setting, we adopt the classical Transformer configuration,

setting both the encoder and decoder layers to 6. To explore the potential redundancy of each module, we compare the detection accuracy and complexity under different combinations of encoder and decoder layers in Table XI. When both the encoder and decoder layers are set to 2, SpecDETR achieves an mAP of 0.47 on the SPOD dataset under 36 training epochs. Increasing either the encoder or decoder layers results in improved detection accuracy for SpecDETR. The gain in detection accuracy from increasing the encoder layers is greater than that from increasing the decoder layers. SpecDETR with 2 encoder layers and 6 decoder layers achieves an mAP of 0.579, whereas SpecDETR with 6 encoder layers and 2 decoder layers achieves an mAP of 0.704. However, the gain in accuracy from increasing encoder layers comes at a cost. The

TABLE X
PERFORMANCE COMPARISON UNDER DIFFERENT NETWORK MODULE COMBINATIONS ON THE SPOD DATASET.

Image Size	Backbone	Encoder	Decoder	mAP			FLOPs (G)	Params (M)
				24 epochs	36 epochs	100 epochs		
$\times 4$	Swin-L	DINO	DINO	0.391	0.605	0.757	203.9	218.3
$\times 1$	-	SpecDETR	DINO	0.663	0.759	0.849	143.6	17.9
$\times 1$	-	SpecDETR	SpecDETR	0.706	0.799	0.856	139.7	16.1

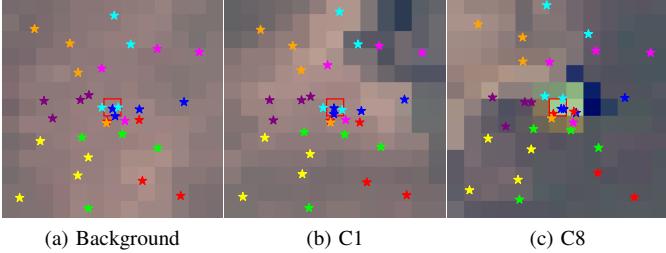


Fig. 28. Examples of sampling point distributions of the subpixel-scale deformable sampling operator in the first encoder layer of SpecDETR on Fig. 23(a). The red box indicates the position of the query element. Sampling points associated with different attention heads are distinguished by various colors.

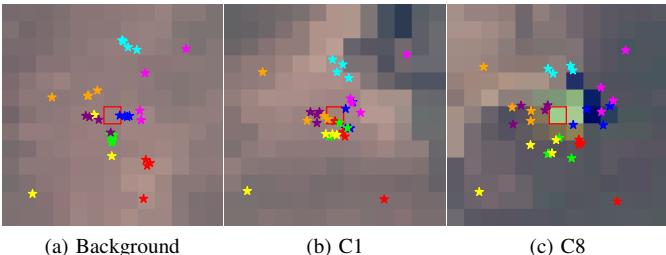


Fig. 29. Examples of sampling point distributions of the subpixel-scale deformable sampling operator in the last encoder layer of SpecDETR on Fig. 23(a).

increase in encoder layers results in an almost linear increase in FLOPs, while the increase in encoder layers has a smaller impact on FLOPs. This is because the number of queries in the encoder is equal to the number of image pixels, whereas in the decoder, it is equal to the number of anchor boxes, with the former being significantly larger than the latter. Additionally, increasing the number of decoder layers leads to a slightly higher increase in Params compared to encoder layers. This is due to the fact that while the Transformer compute operators in the SpecDETR encoder and decoder are largely consistent, each decoder layer includes an additional class prediction head and bbox regression head.

3) *Positive and Negative Sample Setting*: SpecDETR is a one-to-many set prediction problem, primarily reflected in the positive and negative sample setting, which includes the label assignment in the matching branch and the GT box noise addition in the denoising branch. As shown in Table XII, we investigate the impact of different positive and negative sample settings on the performance of SpecDETR. The label assignment considers one-to-one assignment (Forced Matching) and one-to-many assignment (Forced Matching + Dynamic Matching). For the denoising branch, we consider

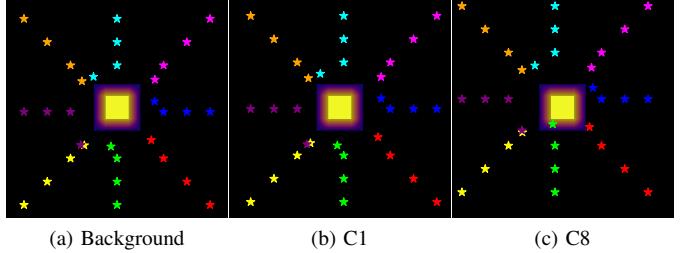


Fig. 30. Examples of sampling point distributions of the self-excited operator in the first encoder layer of SpecDETR on Fig. 23(a).

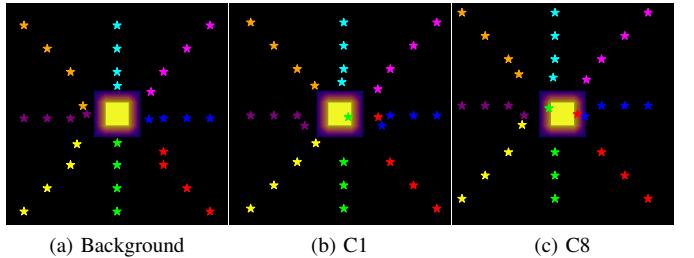


Fig. 31. Examples of sampling point distributions of the self-excited operator in the last encoder layer of SpecDETR on Fig. 23(a).

DN [61], CDN [25], our CCDN, and the case without denoising training. Under the three denoising ways, SpecDETR performs significantly better than the case without denoising training. This indicates that even when the matching part adopts one-to-many label assignment, SpecDETR still needs to rely on denoising training to accelerate convergence. Under one-to-many matching, CCDN outperforms CDN and DN. Additionally, under one-to-one label assignment, CDN is lower than DN by 2.2% and 1.7% in AP at 36 epochs and 100 epochs, respectively, because CDN generates negative samples that are highly similar to the GT, which interferes with the classification head of SpecDETR. Our one-to-many label assignment introduces more high-quality positive samples into the classification head from the matching branch. It improves detection performance regardless of whether DN, CDN, or CCDN is adopted.

4) *Attention Module*: As demonstrated in Table XIII, the removal of the self-excited operator leads to a significant decline in the detection performance of SpecDETR. Specifically, SpecDETR without the self-excited mechanism achieves only 0.228 mAP, 0.540 mAP, and 0.804 mAP at 24 epochs, 36 epochs, and 100 epochs, respectively, on the SPOD dataset. Furthermore, when SpecDETR operates without deformable sampling, the sampling points remain in their initial distribution and the self-excited mechanism is deactivated. Un-

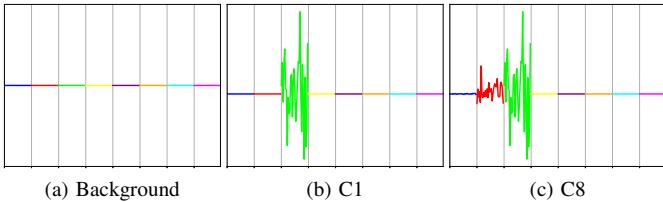


Fig. 32. Examples of sampling output of the self-excited operator in the last encoder layer of SpecDETR on Fig. 23(a).

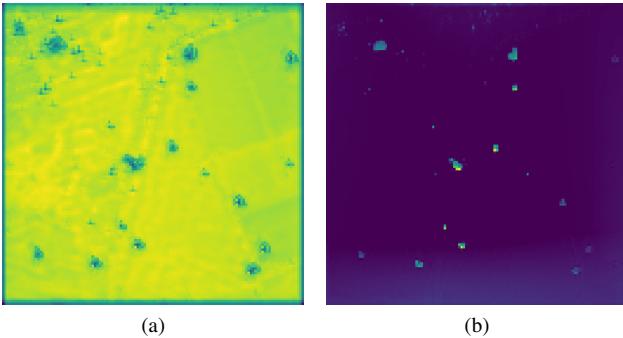


Fig. 33. Output feature maps of two deformable sampling operators from the S2A module in the last encoder layer of SpecDETR, corresponding to Fig. 23(a). (a) The subpixel-scale deformable sampling operator. (b) The self-excited operator. Both (a) and (b) are the sum of the absolute values across all channels of the feature maps.

der this configuration, SpecDETR achieves 0.652 mAP at the 24-epoch setting on the SPOD dataset. These results highlight two insights. First, the use of the subpixel-scale deformable sampling operator alone can disrupt the convergence of SpecDETR, indicating that deformable sampling, without proper mechanisms, may introduce instability. Second, the proposed self-excited mechanism plays a pivotal role in ensuring the convergence of SpecDETR, as it helps stabilize the training process when deformable sampling is employed. This underscores the importance of the self-excited mechanism in enhancing the robustness and performance of SpecDETR.

5) *Content Query Initialization*: As shown in Table XIV, we explore the impact of different content query initialization for SpecDETR. “Constant” represents initializing all content queries as all one vectors, “Random” represents initializing queries with random values each time, and “Embedding” represents the index embedding in DINO [25]. SpecDETR converges fastest with the constant initialization. At 24 epochs, the constant initialization achieved an mAP that is 9.1% higher than the embedding initialization and 15.4% higher than the random initialization. DINO is regarded as a one-to-one set prediction problem, requiring embedding initialization for content queries to differentiate them. We consider SpecDETR as a one-to-many set prediction problem, so the initialization of content queries should be consistent.

VI. CONCLUSION

In this paper, we extend the hyperspectral target detection task to the hyperspectral point object detection task, and propose the first specialized hyperspectral point object detection

TABLE XI
PERFORMANCE COMPARISON OF SPECDETR WITH DIFFERENT ENCODER AND DECODER LAYER SETTINGS ON THE SPOD DATASET UNDER 36 TRAINING EPOCHS.

Encoder Layers	Decoder Layers	mAP	FLOPs (G)	Params (M)
2	2	0.470	48.9	5.7
	4	0.577	54.0	8.4
	6	0.579	59.0	11.2
4	2	0.628	89.2	8.2
	4	0.703	94.3	10.9
	6	0.737	99.4	13.6
6	2	0.704	129.6	10.6
	4	0.750	134.7	13.4
	6	0.799	139.7	16.1

TABLE XII
PERFORMANCE COMPARISON OF SPECDETR WITH DIFFERENT POSITIVE
AND NEGATIVE SAMPLE SETTINGS ON THE SPOD DATASET.

DeNosing	Label Assigner		mAP	
	Forced Matching	Dynamic Matching	36epochs	100epochs
x	✓		0.611	0.796
	✓	✓	0.595	0.800
DN [61]	✓		0.773	0.836
	✓	✓	0.785	0.844
CDN [25]	✓		0.756	0.822
	✓	✓	0.791	0.846
CCDN	✓		0.780	0.829
	✓	✓	0.799	0.856

network, termed SpecDETR. Unlike the traditional HTD methods that rely on a per-pixel binary classification framework, our proposed SpecDETR is capable of learning instance-level spatial-spectral joint object feature representations from large-scale training images, and can directly provide end-to-end instance-level predictions of object location and classification. Extensive experiments verifies that our SpecDETR significantly outperforms both the SOTA visual object detection networks and HTD methods in hyperspectral object detection. However, our framework necessitates a substantial amount of labeled training data, which poses a challenge due to the difficulty in annotating point objects. To address this, we employ simulation techniques to construct synthetic training datasets for single-spectral point objects within three publicly HTD datasets. SpecDETR, trained on these simulation-based datasets, demonstrates robust detection performance for real-world multi-pixel single-spectral objects. Nevertheless, its performance diminishes when confronted with occlusions such as trees or shadows, or when additional noise is artificially introduced into the test images. We attribute this decline in performance to the domain gap between the training and test data. In the future, we plan to explore cross-domain few-shot object detection within the hyperspectral point object detection task, a direction we believe holds potential for mitigating the detection performance degradation caused by the domain gap and the dilemma of insufficient real point object training samples.

REFERENCES

- [1] M. Paoletti, J. Haut, J. Plaza, and A. Plaza, "A new deep convolutional neural network for fast hyperspectral image classification," *ISPRS*

TABLE XIII

PERFORMANCE COMPARISON OF SPECDETR WITH VARIANTS OF THE S2A MODULE ON THE SPOD DATASET. ✕ DENOTES CASES WHERE THE STRUCTURE EXISTS BUT FAILS TO PERFORM ITS INTENDED FUNCTION.

Self-excited operator	Deformable sampling	mAP		
		24epochs	36epochs	100epochs
✗	✓	0.228	0.540	0.804
✗	✗	0.652	0.748	0.842
✓	✓	0.706	0.799	0.856

TABLE XIV

AP PERFORMANCE COMPARISON OF SPECDETR WITH DIFFERENT CONTENT QUERY INITIALIZATION ON THE SPOD DATASET.

Content Query	24epochs	36epochs	100epochs
Embedding	0.647	0.771	0.845
Random	0.612	0.790	0.855
Constant	0.706	0.799	0.856

- Journal of Photogrammetry and Remote Sensing*, vol. 145, pp. 120–147, 2018, deep Learning RS Data.
- [2] R. Thoreau, L. Risser, V. Achard, B. Berthelot, and X. Briottet, “Toulouse hyperspectral data set: A benchmark data set to assess semi-supervised spectral representation learning and pixel-wise classification techniques,” *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 212, pp. 323–337, 2024.
 - [3] F. Zhu, Y. Wang, S. Xiang, B. Fan, and C. Pan, “Structured sparse method for hyperspectral unmixing,” *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 88, pp. 101–118, 2014.
 - [4] H. Su, Z. Wu, A.-X. Zhu, and Q. Du, “Low rank and collaborative representation for hyperspectral anomaly detection via robust dictionary construction,” *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 169, pp. 195–211, 2020.
 - [5] R. Zhao, B. Du, L. Zhang, and L. Zhang, “A robust background regression based score estimation algorithm for hyperspectral anomaly detection,” *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 122, pp. 126–144, 2016.
 - [6] Z. Li, Y. Wang, C. Xiao, Q. Ling, Z. Lin, and W. An, “You only train once: Learning a general anomaly enhancement network with random masks for hyperspectral anomaly detection,” *IEEE Transactions on Geoscience and Remote Sensing*, vol. 61, pp. 1–18, 2023.
 - [7] M. Hu, C. Wu, and L. Zhang, “Globalmind: Global multi-head interactive self-attention network for hyperspectral change detection,” *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 211, pp. 465–483, 2024.
 - [8] Y. Guo, K. Mokany, C. Ong, P. Moghadam, S. Ferrier, and S. R. Levick, “Plant species richness prediction from desis hyperspectral data: A comparison study on feature extraction procedures and regression models,” *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 196, pp. 120–133, 2023.
 - [9] T. Xu, F. Wang, Z. Shi, L. Xie, and X. Yao, “Dynamic estimation of rice aboveground biomass based on spectral and spatial information extracted from hyperspectral remote sensing images at different combinations of growth stages,” *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 202, pp. 169–183, 2023.
 - [10] C. Jiao, C. Chen, R. G. McGarvey, S. Bohlman, L. Jiao, and A. Zare, “Multiple instance hybrid estimator for hyperspectral target characterization and sub-pixel target detection,” *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 146, pp. 235–250, 2018.
 - [11] Y. Dong, L. Zhang, L. Zhang, and B. Du, “Maximum margin metric learning based target detection for hyperspectral images,” *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 108, pp. 138–150, 2015.
 - [12] S. Feng, X. Wang, R. Feng, F. Xiong, C. Zhao, W. Li, and R. Tao, “Transformer-based cross-domain few-shot learning for hyperspectral target detection,” *IEEE Transactions on Geoscience and Remote Sensing*, vol. 62, pp. 1–18, 2024.
 - [13] C.-I. Chang, “Constrained energy minimization (CEM) for hyperspectral target detection: Theory and generalizations,” *IEEE Transactions on Geoscience and Remote Sensing*, vol. 62, pp. 1–21, 2024.
 - [14] S. Feng, R. Feng, D. Wu, C. Zhao, W. Li, and R. Tao, “A coarse-to-fine hyperspectral target detection method based on low-rank tensor decomposition,” *IEEE Transactions on Geoscience and Remote Sensing*, vol. 61, pp. 1–13, 2023.
 - [15] S. Ren, K. He, R. Girshick, and J. Sun, “Faster R-CNN: Towards real-time object detection with region proposal networks,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 39, no. 6, pp. 1137–1149, 2016.
 - [16] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, “You only look once: Unified, real-time object detection,” in *IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 779–788.
 - [17] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko, “End-to-end object detection with transformers,” in *European Conference on Computer Vision*, 2020, pp. 213–229.
 - [18] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, “Microsoft COCO: Common objects in context,” in *European Conference on Computer Vision*, 2014, pp. 740–755.
 - [19] G. Cheng, X. Yuan, X. Yao, K. Yan, Q. Zeng, X. Xie, and J. Han, “Towards large-scale small object detection: Survey and benchmarks,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023.
 - [20] C. E. Caefer, M. S. Stefanou, E. D. Nielsen, A. P. Rizzuto, O. Raviv, and S. R. Rotman, “Analysis of false alarm distributions in the development and evaluation of hyperspectral point target detection algorithms,” *Optical Engineering*, vol. 46, no. 7, pp. 076402–076402, 2007.
 - [21] C. E. Caefer, J. Silverman, O. Orthal, D. Antonelli, Y. Sharoni, and S. R. Rotman, “Improved covariance matrices for point target detection in hyperspectral data,” *Optical Engineering*, vol. 47, no. 7, pp. 076402–076402, 2008.
 - [22] C. E. Caefer, S. R. Rotman, J. Silverman, and P. W. Yip, “Algorithms for point target detection in hyperspectral imagery,” in *Imaging Spectrometry VIII*, vol. 4816. SPIE, 2002, pp. 242–257.
 - [23] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, “Attention is all you need,” *Neural Information Processing Systems*, vol. 30, 2017.
 - [24] X. Zhu, W. Su, L. Lu, B. Li, X. Wang, and J. Dai, “Deformable DETR: Deformable transformers for end-to-end object detection,” *arXiv*, 2020.
 - [25] H. Zhang, F. Li, S. Liu, L. Zhang, H. Su, J. Zhu, L. M. Ni, and H.-Y. Shum, “DINO: DETR with improved denoising anchor boxes for end-to-end object detection,” *arXiv*, 2022.
 - [26] N. M. Nasrabadi, “Hyperspectral target detection : An overview of current and future challenges,” *IEEE Signal Processing Magazine*, vol. 31, no. 1, pp. 34–44, Dec. 2013.
 - [27] F. C. Robey, D. R. Fuhrmann, E. J. Kelly, and R. Nitzberg, “A CFAR adaptive matched filter detector,” *IEEE Transactions on Aerospace and Electronic Systems*, vol. 28, no. 1, pp. 208–216, Jan. 1992.
 - [28] J. C. Harsanyi, “Detection and classification of subpixel spectral signatures in hyperspectral image sequences,” Ph.D. dissertation, University of Maryland Baltimore County, Baltimore, MD, USA, 1993.
 - [29] S. Kraut, L. L. Scharf, and R. W. Butler, “The adaptive coherence estimator: A uniformly most-powerful-invariant adaptive detection statistic,” *IEEE Transactions on Signal Processing*, vol. 53, no. 2, pp. 427–438, 2005.
 - [30] B. Du, Y. Zhang, L. Zhang, and D. Tao, “Beyond the sparsity-based target detector: A hybrid sparsity and statistics based detector for hyperspectral images,” *IEEE Transactions on Image Processing*, vol. 25, no. 11, pp. 5345–5357, Nov. 2016.
 - [31] S. Kraut and L. L. Scharf, “The CFAR adaptive subspace detector is a scale-invariant GLRT,” *IEEE Transactions on Signal Processing*, vol. 47, no. 9, pp. 2538–2541, Sep. 1999.
 - [32] J. C. Harsanyi and C.-I. Chang, “Hyperspectral image classification and dimensionality reduction: An orthogonal subspace projection approach,” *IEEE Transactions on Geoscience and Remote Sensing*, vol. 32, no. 4, pp. 779–785, Jul. 1994.
 - [33] H. Ren and C.-I. Chang, “Target-constrained interference-minimized approach to subpixel target detection for hyperspectral images,” *Optical Engineering*, vol. 39, pp. 3138–3145, 2000.
 - [34] Y. Chen, N. M. Nasrabadi, and T. D. Tran, “Simultaneous joint sparsity model for target detection in hyperspectral imagery,” *IEEE Geoscience and Remote Sensing Letters*, vol. 8, no. 4, pp. 676–680, Jul. 2011.
 - [35] Y. Chen, N. M. Nasrabadi, and T. D. Tran, “Sparse representation for target detection in hyperspectral imagery,” *IEEE Journal of Selected Topics in Signal Processing*, vol. 5, no. 3, pp. 629–640, Jun. 2011.
 - [36] Y. Chen, N. M. Nasrabadi, and T. D. Tran, “Hyperspectral image classification via kernel sparse representation,” *IEEE Transactions on Geoscience and Remote Sensing*, vol. 51, no. 1, pp. 217–231, Jan. 2013.

- [37] Y. Zhang, B. Du, and L. Zhang, "A sparse representation-based binary hypothesis model for target detection in hyperspectral images," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 53, no. 3, pp. 1346–1354, Mar. 2015.
- [38] Q. Ling, Y. Guo, Z. Lin, L. Liu, and W. An, "A constrained sparse-representation-based binary hypothesis model for target detection in hyperspectral imagery," *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, vol. 12, no. 6, pp. 1933–1947, 2019.
- [39] Y. Zhang, L. Zhang, B. Du, and S. Wang, "A nonlinear sparse representation-based binary hypothesis model for hyperspectral target detection," *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, vol. 8, no. 6, pp. 2513–2522, Jun. 2015.
- [40] D. Zhu, B. Du, and L. Zhang, "Learning single spectral abundance for hyperspectral subpixel target detection," *IEEE Transactions on Neural Networks and Learning Systems*, 2023.
- [41] X. Jiao and C.-I. Chang, "Kernel-based constrained energy minimization (K-CEM)," in *Algorithms and Technologies for Multispectral, Hyperspectral, and Ultraspectral Imagery XIV*. SPIE, 2008, pp. 523–533.
- [42] H. Kwon and N. M. Nasrabadi, "Kernel orthogonal subspace projection for hyperspectral signal classification," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 43, no. 12, pp. 2952–2962, 2005.
- [43] H. Kwon and N. M. Nasrabadi, "A comparative analysis of kernel subspace target detectors for hyperspectral imagery," *EURASIP Journal on Advances in Signal Processing*, vol. 2007, pp. 1–13, 2006.
- [44] H. Kwon and N. M. Nasrabadi, "Kernel matched subspace detectors for hyperspectral target detection," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 28, no. 2, pp. 178–194, Feb. 2006.
- [45] T. Wang, B. Du, and L. Zhang, "A kernel-based target-constrained interference-minimized filter for hyperspectral sub-pixel target detection," *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, vol. 6, no. 2, pp. 626–637, Apr. 2013.
- [46] G. Camps-Valls, L. Gomez-Chova, J. Muñoz-Marí, J. Vila-Francés, and J. Calpe-Maravilla, "Composite kernels for hyperspectral image classification," *IEEE Geoscience and Remote Sensing Letters*, vol. 3, no. 1, pp. 93–97, Jan. 2006.
- [47] J. Jiao, Z. Gong, and P. Zhong, "Triplet spectral-wise transformer network for hyperspectral target detection," *IEEE Transactions on Geoscience and Remote Sensing*, 2023.
- [48] L. Sun, Z. Ma, and Y. Zhang, "Abal: Adaptive background latent space adversarial learning algorithm for hyperspectral target detection," *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 2023.
- [49] F. Luo, S. Shi, T. Guo, Y. Dong, L. Zhang, and B. Du, "Agms: Adversarial sample generation-based multi-scale siamese network for hyperspectral target detection," *IEEE Transactions on Geoscience and Remote Sensing*, 2024.
- [50] Y. Wang, X. Chen, F. Wang, M. Song, and C. Yu, "Meta-learning based hyperspectral target detection using siamese network," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1–13, 2022.
- [51] Y. Shi, H. Cui, Y. Yin, H. Song, Y. Li, and P. Gamba, "Transfer learning with nonlinear spectral synthesis for hyperspectral target detection," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 61, pp. 1–17, 2023.
- [52] D. Zhu, B. Du, and L. Zhang, "Two-stream convolutional networks for hyperspectral target detection," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 59, no. 8, pp. 6907–6921, 2020.
- [53] W. Rao, L. Gao, Y. Qu, X. Sun, B. Zhang, and J. Chanussot, "Siamese transformer network for hyperspectral image target detection," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1–19, 2022.
- [54] Q. Yang, X. Wang, L. Chen, Y. Zhou, and S. Qiao, "Cs-ttd: triplet transformer for compressive hyperspectral target detection," *IEEE Transactions on Geoscience and Remote Sensing*, 2024.
- [55] W. Dong, X. Wu, J. Qu, P. Gamba, S. Xiao, A. Vizziello, and Y. Li, "Deep spatial-spectral joint-sparse prior encoding network for hyperspectral target detection," *IEEE Transactions on Cybernetics*, 2024.
- [56] D. Shen, X. Ma, W. Kong, J. Liu, J. Wang, and H. Wang, "Hyper-spectral target detection based on interpretable representation network," *IEEE Transactions on Geoscience and Remote Sensing*, 2023.
- [57] Y. Li, H. Qin, and W. Xie, "Htdformer: Hyperspectral target detection based on transformer with distributed learning," *IEEE Transactions on Geoscience and Remote Sensing*, 2023.
- [58] Y. Wang, X. Chen, E. Zhao, C. Zhao, M. Song, and C. Yu, "An unsupervised momentum contrastive learning based transformer network for hyperspectral target detection," *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 2024.
- [59] X. Chen, M. Zhang, and Y. Liu, "Target detection with spectral graph contrast clustering assignment and spectral graph transformer in hyperspectral imagery," *IEEE Transactions on Geoscience and Remote Sensing*, 2024.
- [60] S. Liu, F. Li, H. Zhang, X. Yang, X. Qi, H. Su, J. Zhu, and L. Zhang, "DAB-DETR: Dynamic anchor boxes are better queries for DETR," *arXiv*, 2022.
- [61] F. Li, H. Zhang, S. Liu, J. Guo, L. M. Ni, and L. Zhang, "DN-DETR: Accelerate DETR training by introducing query denoising," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2022, pp. 13 619–13 627.
- [62] Q. Chen, X. Chen, J. Wang, S. Zhang, K. Yao, H. Feng, J. Han, E. Ding, G. Zeng, and J. Wang, "Group DETR: Fast DETR training with group-wise one-to-many assignment," in *International Conference on Computer Vision*, 2023, pp. 6633–6642.
- [63] Z. Zong, G. Song, and Y. Liu, "Dets with collaborative hybrid assignments training," in *International Conference on Computer Vision*, 2023, pp. 6748–6758.
- [64] Y. Zhao, W. Lv, S. Xu, J. Wei, G. Wang, Q. Dang, Y. Liu, and J. Chen, "Detrs beat yolos on real-time object detection," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2024, pp. 16 965–16 974.
- [65] A. Wang, H. Chen, L. Liu, K. Chen, Z. Lin, J. Han, and G. Ding, "Yolov10: Real-time end-to-end object detection," *arXiv*, 2024.
- [66] L. Yan, M. Zhao, X. Wang, Y. Zhang, and J. Chen, "Object detection in hyperspectral images," *IEEE Signal Processing Letters*, vol. 28, pp. 508–512, 2021.
- [67] G. Lee, J. Lee, J. Baek, H. Kim, and D. Cho, "Channel sampler in hyperspectral images for vehicle detection," *IEEE Geoscience and Remote Sensing Letters*, vol. 19, pp. 1–5, 2021.
- [68] A. Rangnekar, Z. Mulholland, A. Vodacek *et al.*, "Semi-supervised hyperspectral object detection challenge results-PBVS 2022," in *IEEE Conference on Computer Vision and Pattern Recognition Workshop*, 2022, pp. 390–398.
- [69] X. He, C. Tang, X. Liu, W. Zhang, K. Sun, and J. Xu, "Object detection in hyperspectral image via unified spectral-spatial feature aggregation," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 61, pp. 1–13, 2023.
- [70] S. Kumar, I. Arevalo, A. Iftekhar, and B. Manjunath, "Methanemap: Spectral absorption aware hyperspectral transformer for methane detection," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2023, pp. 17 609–17 618.
- [71] Z. Li, F. Xiong, J. Zhou, J. Lu, and Y. Qian, "Learning a deep ensemble network with band importance for hyperspectral object tracking," *IEEE Transactions on Image Processing*, vol. 32, pp. 2901–2914, 2023.
- [72] Z. Liu, X. Wang, Y. Zhong, M. Shu, and C. Sun, "SiamHYPER: Learning a Hyperspectral Object Tracker From an RGB-Based Tracker," *IEEE Transactions on Image Processing*, vol. 31, pp. 7116–7129, 2022.
- [73] Z. Liu, Y. Zhong, G. Ma, X. Wang, and L. Zhang, "A deep temporal-spectral-spatial anchor-free siamese tracking network for hyperspectral video object tracking," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 62, pp. 1–16, 2024.
- [74] Z. Li, G. Guo, X. He, Q. Xu, W. Wang, Q. Ling, Z. Lin, and W. An, "RawTrack: Toward Single Object Tracking on Mosaic Hyperspectral Raw Data," in *13th Workshop on Hyperspectral Imaging and Signal Processing: Evolution in Remote Sensing (WHISPERS)*, 2023, pp. 1–5.
- [75] J. L. Ba, J. R. Kiros, and G. E. Hinton, "Layer normalization," *arXiv*, 2016.
- [76] J. Liu and J. Zhang, "Spectral unmixing via compressive sensing," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 52, no. 11, pp. 7099–7110, Nov. 2014.
- [77] H. Rezatofighi, N. Tsoi, J. Gwak, A. Sadeghian, I. Reid, and S. Savarese, "Generalized intersection over union: A metric and a loss for bounding box regression," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 658–666.
- [78] R. Kokaly, R. Clark, G. Swayze, K. Livo, T. Hoefen, N. Pearson, R. Wise, W. Benzel, H. Lowers, R. Driscoll *et al.*, "USGS spectral library version 7 data: US geological survey data release," *United States Geological Survey (USGS): Reston, VA, USA*, vol. 61, 2017.
- [79] A. Giannandrea, N. Raqueno, D. W. Messinger, J. Faulring, J. P. Kerekes *et al.*, "The share 2012 data campaign," in *Algorithms and Technologies for Multispectral, Hyperspectral, and Ultraspectral Imagery XIX*. SPIE, 2013, pp. 94–108.

- [80] C.-I. Chang, “An effective evaluation tool for hyperspectral target detection: 3D receiver operating characteristic curve analysis,” *IEEE Transactions on Geoscience and Remote Sensing*, vol. 59, no. 6, pp. 5131–5153, 2020.
- [81] B. Li, C. Xiao, L. Wang, Y. Wang, Z. Lin, M. Li, W. An, and Y. Guo, “Dense nested attention network for infrared small target detection,” *IEEE Transactions on Image Processing*, vol. 32, pp. 1745–1758, 2022.
- [82] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 770–778.
- [83] I. Radosavovic, R. P. Kosaraju, R. Girshick, K. He, and P. Dollár, “Designing network design spaces,” in *IEEE Conference on Computer Vision and Pattern Recognition*, 2020, pp. 10428–10436.
- [84] H. Zhang, C. Wu, Z. Zhang, Y. Zhu, H. Lin, Z. Zhang, Y. Sun, T. He, J. Mueller, R. Manmatha *et al.*, “ResNeSt: Split-attention networks,” in *IEEE Conference on Computer Vision and Pattern Recognition*, 2022, pp. 2736–2746.
- [85] S. Xie, R. Girshick, P. Dollár, Z. Tu, and K. He, “Aggregated residual transformations for deep neural networks,” in *IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 1492–1500.
- [86] K. Sun, B. Xiao, D. Liu, and J. Wang, “Deep high-resolution representation learning for human pose estimation,” in *IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 5693–5703.
- [87] C. Feng, Y. Zhong, Y. Gao, M. R. Scott, and W. Huang, “TOOD: Task-aligned one-stage object detection,” in *International Conference on Computer Vision*, 2021, pp. 3490–3499.
- [88] Z. Dong, G. Li, Y. Liao, F. Wang, P. Ren, and C. Qian, “Centripetalnet: Pursuing high-quality keypoint pairs for object detection,” in *IEEE Conference on Computer Vision and Pattern Recognition*, 2020, pp. 10519–10528.
- [89] 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.
- [90] H. Law and J. Deng, “Cornernet: Detecting objects as paired keypoints,” in *European Conference on Computer Vision*, 2018, pp. 734–750.
- [91] Z. Yang, S. Liu, H. Hu, L. Wang, and S. Lin, “Reppoints: Point set representation for object detection,” in *IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 9657–9666.
- [92] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, “Focal loss for dense object detection,” in *IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 2980–2988.
- [93] M. Tan and Q. Le, “Efficientnet: Rethinking model scaling for convolutional neural networks,” in *International Conference on Machine Learning*, 2019, pp. 6105–6114.
- [94] W. Wang, E. Xie, X. Li, D.-P. Fan, K. Song, D. Liang, T. Lu, P. Luo, and L. Shao, “PVTv2: Improved baselines with pyramid vision transformer,” *Computational Visual Media*, vol. 8, no. 3, pp. 415–424, 2022.
- [95] 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 *International Conference on Computer Vision*, 2021, pp. 10012–10022.
- [96] W. Li and Q. Du, “Collaborative representation for hyperspectral anomaly detection,” *IEEE Transactions on Geoscience and Remote Sensing*, vol. 53, no. 3, pp. 1463–1474, Mar. 2015.
- [97] K. Chen, J. Wang, J. Pang, Y. Cao *et al.*, “MMDetection: Open mmlab detection toolbox and benchmark,” *arXiv*, 2019.
- [98] C. Xiao, W. An, Y. Zhang, Z. Su, M. Li, W. Sheng, M. Pietikinen, and L. Liu, “Highly efficient and unsupervised framework for moving object detection in satellite videos,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 46, no. 12, pp. 11532–11539, 2024.
- [99] L. Zhang, L. Zhang, D. Tao, and X. Huang, “Sparse transfer manifold embedding for hyperspectral target detection,” *IEEE Transactions on Geoscience and Remote Sensing*, vol. 52, no. 2, pp. 1030–1043, 2013.
- [100] Y. Fu, Y. Wang, Y. Pan, L. Huai, X. Qiu, Z. Shangguan, T. Liu, Y. Fu, L. Van Gool, and X. Jiang, “Cross-domain few-shot object detection via enhanced open-set object detector,” in *European Conference on Computer Vision*. Springer, 2024, pp. 247–264.

Supplementary Materials for SpecDETR: A Transformer-based Hyperspectral Point Object Detection Network

Zhaoxu Li, Wei An, Gaowei Guo, Longguang Wang, Yingqian Wang, and Zaiping Lin

I. SUPPLEMENTS TO THE SPOD DATASET

A. Spectral Fluctuation Model

To simulate spectral fluctuations more realistically, we modeled the spectral fluctuations of real hyperspectral data captured by the AVIRIS sensor. We chose these water areas as the modeling objects because land background pixels are often mixed, while water areas such as oceans and lakes have a single spectral composition.

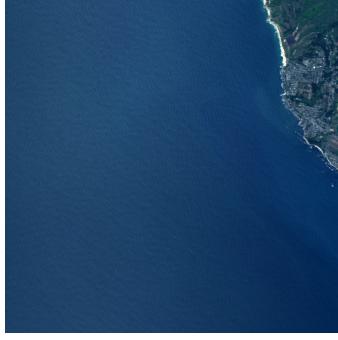


Fig. A. Golden Gate.

As shown in Fig. A, a 880×460 region of seawater in the Golden Gate, California, is selected as the reference area for analysis. First, the coefficient of variation of the radiances in each band is calculated for the reference area:

$$\gamma_i = \frac{\mu_i}{\sigma_i} \quad (1)$$

where γ_i represents the coefficient of variation of the radiances in the i -th band in the reference area, μ_i represents the mean of the radiances in the i -th band, and σ_i represents the standard deviation of the radiances in the i -th band.

Fig. B shows the coefficient curve of variation in the reference area with wavelength. It can be observed that except for the atmospheric absorption band at 1400nm, 1800nm, and 2500nm, the coefficient of variation in the reference area is continuous and smooth. For adjacent pixels of the same category on the same image, their imaging conditions, such as observation angle and lighting, can be considered consistent. To describe the fluctuations of these pixels of the same category in a small-scale area, we introduce a local fluctuation factor:

$$\alpha_i^j = \frac{s_i^j - \mu_i}{\mu_i} \quad (2)$$

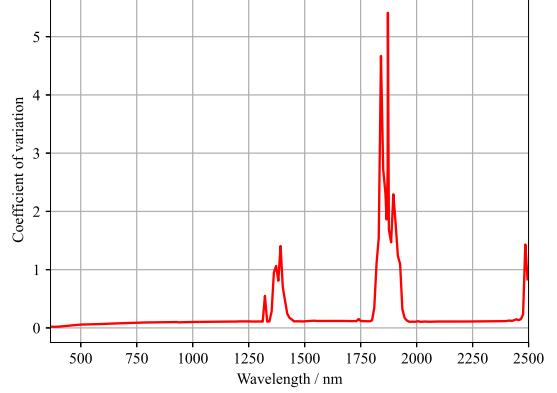


Fig. B. The coefficient of variation curve of the reference area.

where s_i^j represents the radiance value, and α_i^j represents the local fluctuation factor of the j -th pixel in the i -th band within the reference area.

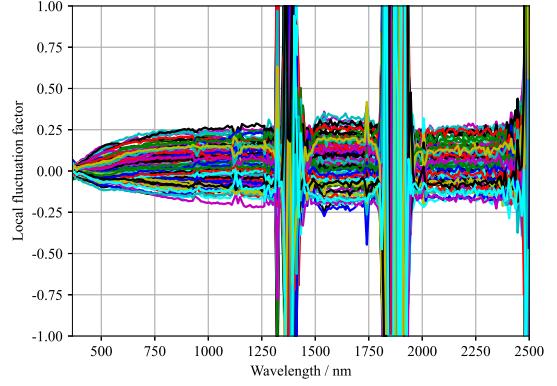


Fig. C. The local fluctuation factor curves in the reference area.

Fig. C shows the variation of the local fluctuation factor of different pixels in the reference area with wavelength. It can be observed that except for the atmospheric absorption bands at 1400nm, 1800nm, and 2500nm, the fluctuation level of the local fluctuation factor increases with increasing wavelength, similar to the coefficient curve of variation. Therefore, the standardized local fluctuation factor is defined as:

$$\bar{\alpha}_i^j = \frac{\alpha_i^j}{\gamma_i} \quad (3)$$

where $\bar{\alpha}_i^j$ represents the standardized local fluctuation factor of the j -th pixel in the i -th band in the reference area.

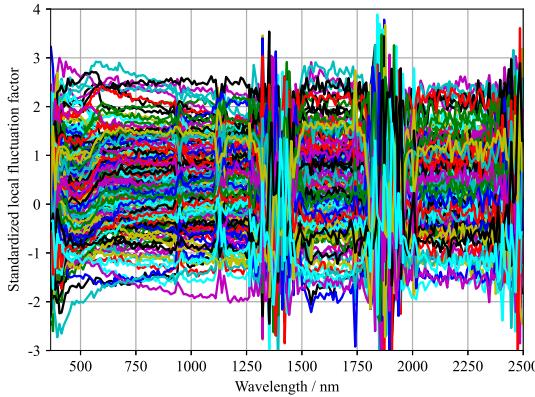


Fig. D. The standardized local fluctuation factor curves in the reference area.

Fig. D shows the standardized local fluctuation factor variation of different pixels in the reference area with wavelength. It can be observed that except for the atmospheric absorption bands, the curves of the standardized local fluctuation factor in the reference area are relatively stable in the visible to short-wave infrared spectral range. Hence, the standardized local fluctuation factor can be modeled as:

$$\bar{\alpha}_i^j = a^j + v_i^j \quad (4)$$

where a^j represents the baseline value of the standardized local fluctuation factor for the j -th pixel in the reference area, and v_i^j represents the noise of the standardized local fluctuation factor for the j -th pixel in the i -th band. The baseline value a^j can be approximated by the mean of the standardized local fluctuation factors in each band.

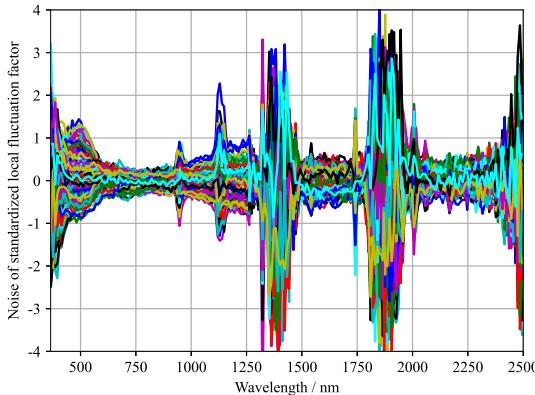


Fig. E. The noise curves of the standardized local fluctuation factor in the reference area.

Fig. E shows the noise curves of the standardized local fluctuation factor for different pixels in the reference area. It can be observed that the noise intensity of the standardized local fluctuation factor is greatly influenced by wavelength.

Fig. F displays the noise distribution of the standardized local fluctuation factor at 654nm in the reference area. It can be observed that at a specific wavelength, the noise of the standardized local fluctuation factor can be considered as Gaussian noise. Fig. G shows the standard deviation variation of the standardized local fluctuation factor in the reference area with wavelength.

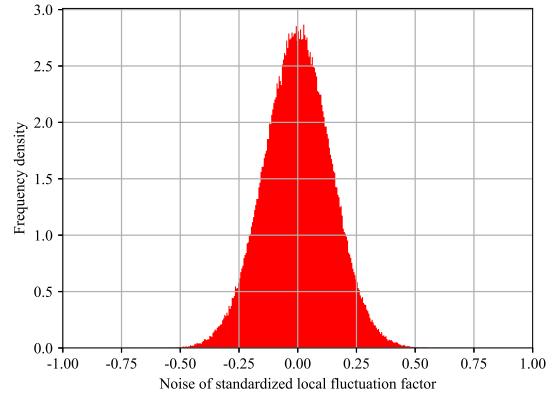


Fig. F. The noise distribution of the standardized local fluctuation factor at 654nm in the reference area.

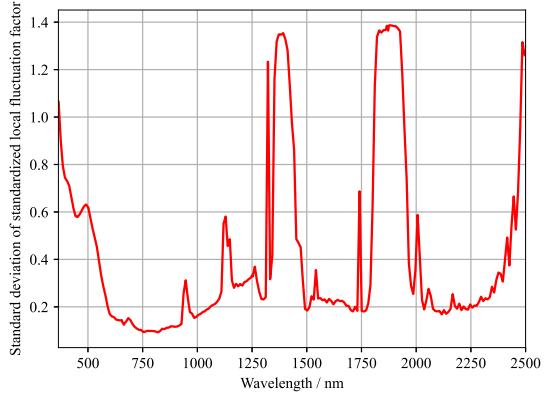


Fig. G. The standard deviation curves of the standardized local fluctuation factor in the reference area.

From Fig. H, it can be concluded that the baseline value of the standardized local fluctuation factor in the reference area can be approximated as a normal distribution with a mean of 0. Based on the above analysis, in a small-scale area on the same image, the spectra of adjacent pixels of the same category in non-atmospheric-absorption bands can be modeled numerically as:

$$s \approx \gamma \circ (a + v) \circ \bar{s} + \bar{s} \quad (5)$$

where \circ denotes element-wise multiplication, s represents the observed spectral curve of a certain pixel, \bar{s} represents the baseline spectral curve under the current observation conditions and can be approximated by the average radiance vector μ , s and \bar{s} are both $N \times 1$ column vectors, where N is the number of bands. a represents the baseline value of the standardized local fluctuation factor, which follows a normal distribution with mean 0 and standard deviation σ_a . v represents the noise of the standardized local fluctuation factor, which follows an N -dimensional normal distribution with mean 0 and standard deviation vector σ_v .

Based on the local fluctuation characteristics, we further analyzed the wide-area spectral fluctuation characteristics of six different water areas shown in Fig. I-A. These water areas are far apart with longer imaging time intervals, resulting in differences in imaging conditions such as observation distance, observation angle, and lighting. To describe the fluctuations

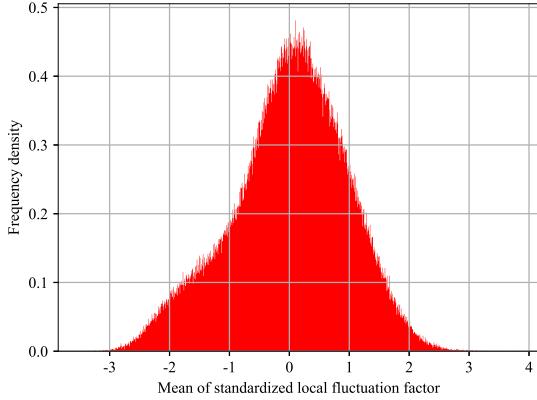


Fig. H. The baseline value distribution of the standardized local fluctuation factor for different pixels in the reference area.

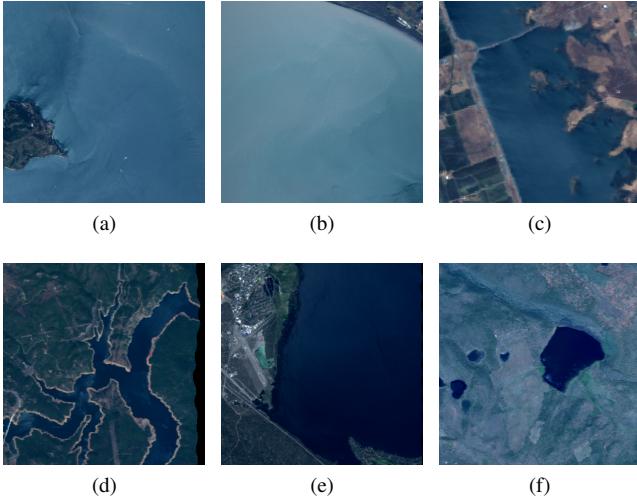


Fig. I. Six water areas near the Golden Gate. (a) San Francisco Bay. (b) San Pablo Bay. (c) Thermalito Afterbay. (d) West Branch Feather River. (e) Lake Almanor. (f) Big Jack Lake.

of spectra of the same category under different imaging conditions, we introduce a wide-area fluctuation factor:

$$\beta_i^k = \frac{\mu_i^k - \mu_i}{\mu_i} \quad (6)$$

where β_i^k represents the wide-area fluctuation factor of the k -th water area in the i -th band, μ_i^k represents the average radiance of the k -th water area in the i -th band, and μ_i represents the average radiance of the reference area in the i -th band.

Fig. I-A shows the variation of the wide-area fluctuation factor of different water areas with wavelength. It can be observed that, except for the atmospheric absorption bands at 1400nm, 1800nm, and 2500nm, the wide-area fluctuation factor exhibits a similar trend to the local fluctuation factor. Based on this, we define a standardized wide-area fluctuation factor as:

$$\bar{\beta}_i^k = \frac{\beta_i^k}{\gamma_i} \quad (7)$$

where $\bar{\beta}_i^k$ represents the standardized wide-area fluctuation factor of the k -th water area in the i -th band.

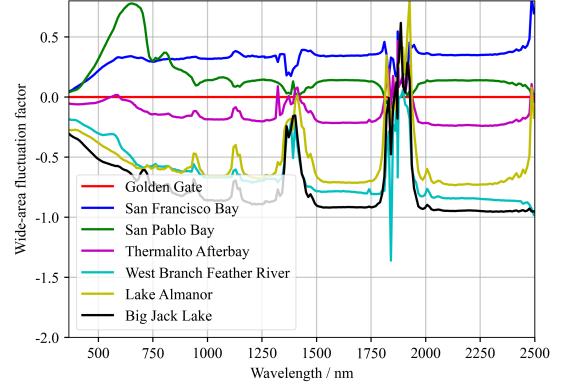


Fig. J. The wide-area fluctuation factor curves of different water areas.

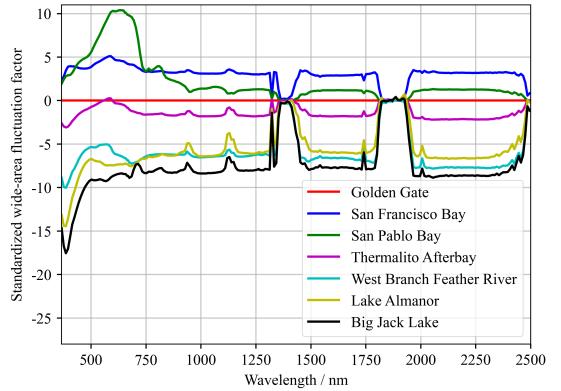


Fig. K. The standardized wide-area fluctuation factor curves of different water areas.

Fig. K shows the variation of the standardized wide-area fluctuation factor of different water areas with wavelength. The standardized wide-area fluctuation factor of different water areas remains stable on the near-infrared to short-wave infrared spectral range outside the atmospheric absorption bands. However, in the visible spectral range, some water areas exhibit more significant fluctuations in the wide-area fluctuation factor. This is because the reflectivity of water areas in the visible spectral range is greatly influenced by the chlorophyll content, which varies among different water areas. In the near-infrared and short-wave infrared spectral range, the influence of chlorophyll on spectral reflectivity is minimal, and these water areas have similar spectral reflection characteristics. By ignoring factors such as chlorophyll in the near-infrared to short-wave infrared spectral range outside the atmospheric absorption bands, the average spectral curve of the k -th water area can be approximated as:

$$\mu^k \approx b^k \gamma \circ \mu + \mu \quad (8)$$

where b^k represents the baseline value of the standardized wide-area fluctuation factor for the k -th water area.

Combining Eq. (5) and Eq. (8), the fluctuations of spectra of the same category in non-atmospheric-absorption bands can be modeled as:

$$s = (b\gamma + 1) \circ ((a + v) \circ \gamma + 1) \circ \bar{s} \quad (9)$$

where s represents the measured spectral curve, \bar{s} represents the average spectral curve of the reference area, γ represents the coefficient of variation curve of the reference area. b represents the standardized wide-area fluctuation factor, and a represents the baseline value of the standardized local fluctuation factor, which follows a normal distribution with mean 0 and standard deviation σ_a . v represents the noise of the standardized local fluctuation factor, which follows an N -dimensional normal distribution with mean 0 and standard deviation σ_v .

B. Supplements to Implementation Details

We use the AVIRIS raw data that has not been corrected with radiance gains to generate the SPOD dataset. To ensure the discriminability of simulated object spectra in the SPOD dataset, we select the spectral reflectance curves of 12 artificial materials from the USGS spectral library to generate simulated spectral radiance curves. The spectral reflectance curve refers to the ratio of the light flux reflected by an object to the light flux incident on the object. Accurately simulating radiance curves from spectral reflectance curves requires careful consideration of the entire process of light propagation in the atmosphere, including factors such as atmospheric absorption, atmospheric scattering, and surface reflection. We only need to validate the methods for hyperspectral object detection, so it is unnecessary to precisely approximate the actual radiance of specific land cover types for the simulated object spectra. Therefore, we simplify the simulation process by treating the transformation from spectral reflectance curves to radiance curves as a linear transformation. The simulated radiance curves can be obtained by :

$$s_t = \frac{s_w}{r_w} \circ r_t \quad (10)$$

where s_t represents the simulated object radiance curve, s_w represents the average radiance curve of the reference water area, r_w is the spectral reflectance curve of the seawater provided by the USGS spectral library, and r_t stands for the spectral reflectance curve of selected artificial materials provided by the USGS spectral library. As r_w is the water first surface reflection curve and has relatively low values, the values of s_t are excessively high. Therefore, linear scaling is required to map the numerical range of s_t to the normal range of the AVIRIS data. The simulated object spectra can be added to the hyperspectral images by :

$$s = \frac{M_t}{\max(s_t)} (b\gamma + 1) \circ ((a + v) \circ \gamma + 1) \circ s_t \quad (11)$$

where M_t represents the baseline of the maximum value of the simulated object spectra. Considering the numerical range of the AVIRIS data used in the SPOD dataset, The value range of M_t for each type of simulated spectrum is [2000, 3000], which corresponds to the distribution range of maximum digital numbers for most spectra in the AVIRIS data. γ , σ_a , and σ_v are the statistical values of the reference water area. It can be assumed that their observation conditions are consistent for different pixels of the same object, and b is the same for all of them. We set b to follow a uniform distribution in the range

[−0.3, 0.3], ensuring significant spectral fluctuations between the same type of objects.

II. SUPPLEMENTS TO EXPERIMENTS ON THE SPOD DATASET

Tables. G and H and serve as the supplements to Table II by providing AP and AR performance of more object detection networks [1], [2], [2]–[9] on the SPOD dataset, respectively. All object detection networks use parameter configurations of the COCO dataset [10] from the MMDetection framework [11], with variations limited to training epochs, image sizes, and input channels. Due to the domain gap between the SPOD and COCO datasets, many networks face challenges in detecting point objects in the SPOD dataset. However, our SpecDETR, as well as other networks such as DINO and CentripetalNet, demonstrate the feasibility of point object detection. In addition, SpecDETR performs the best across all APs. Despite not providing redundant predictions, which can significantly impact AR performance, the mAR of SpecDETR is second only to that of DINO.

III. EXPERIMENTAL SETUP FOR COMPARED HTD METHODS

The compared HTD methods, LSSA [12], IRN [13], and TSTTD [14], are adapted based on their official codes to suit the characteristics of our dataset. Among these, IRN and TSTTD are implemented in Python, while LSSA is implemented in MATLAB. The original LSSA takes a single target spectrum as input and outputs one abundance map, whereas the improved LSSA processes multiple target spectra simultaneously and output multiple abundance maps for different target classes. The original IRN and TSTTD utilize background pixels from a single test image for training, while the improved versions use background pixels from all training images. Additionally, the original TSTTD employs a single target prior spectrum for data augmentation, whereas the improved TSTTD directly uses all target pixels from the training images for training. Other parameters are retained at their default settings as provided in the official implementations.

We implement several compared HTD methods in Python, including ASD [15], CEM [16], OSP [17], KOSP [18], SMF [19], KSMF [20], TCIMF [21], CR [22], KSR [23], and KSRBBH [24]. These methods follow the classic HTD pipeline, where each type of object spectra is sequentially used as the prior spectra, and each pixel in the all test images is sequentially treated as the pixel to be tested. The background pixels are selected from a dual-window centered on the pixel to be tested. The dual-window sizes ($\omega_{in}, \omega_{out}$) are configured based on the size of each type of object, as detailed in Table A. We tune the other hyperparameters of these algorithms using the test image of the Avon dataset.

TABLE A
SIMULATION PARAMETER SETTINGS FOR TRAINING SETS OF THE AVON, SANDIEGO, AND GULFPORT DATASETS.

Dataset	SPOD Dataset								Avon Dataset	
	C1	C2	C3	C4	C5	C6	C7	C8	Blue tarp	Brown tarp
ω_{in}	1	1	3	11	11	13	13	13	5	5
ω_{out}	3	3	5	13	13	15	15	15	7	7

TABLE B
SIMULATION PARAMETER SETTINGS FOR TRAINING SETS OF THE AVON, SANDIEGO, AND GULFPORT DATASETS.

Dataset	Test data	Training data	Object pixels	Maximum abundance	Training images	Training samples
Avon	Avon Flight 0920-1701	Avon Flight 0920-1631	11-16	0.95-1	150	1500/1500
SanDiego	SanDiego	Avon Flight 0920-1631	11-20	0.95-1	100	2000
Gulfport	Gulfport Flight3	Gulfport Flight4	3-10	0.7-1	200	1000/1000/1000

IV. GENERATION OF SIMULATED TRAINING DATA FOR PUBLIC HTD DATASETS

Traditional HTD typically provides only a single test image and a few, or even just one, prior spectra. In contrast, the proposed hyperspectral point object detection framework requires a large number of annotated training images and training samples. To address this, we adopt the data simulation method from the SPOD dataset to generate simulated training sets for three public HTD datasets: Avon, SanDiego, and Gulfport. These datasets contain single-spectrum objects of 2, 1, and 4 classes, respectively. For each objects class, we extract a pure target spectrum from the test image for simulation. Table B presents the simulation parameter settings for the training sets of the three datasets. Taking the Avon dataset as an example, we crop 150 training images of size 128×128 from the flight line data numbered 0920-1631. Each training image contains 10 simulated blue tarps and 10 simulated brown tarps randomly placed. The number of pixels for each simulated object is randomly selected from the range (11, 16), and the pixels are combined in a random clustered pattern. The object pixels not adjacent to background pixels are considered pure pixels, with object spectral abundance set to 1. The object pixels adjacent to background pixels are treated as mixed pixels, with object spectral abundance randomly generated from the range (0.1, 1). The abundance of mixed pixels is assigned based on their distance to the object center, with closer pixels receiving higher abundance values. Additionally, if all the object pixels are mixed pixels, the abundance of the center pixel is randomly selected from the maximum abundance range (0.95, 1). For the Gulfport dataset, where objects are generally smaller and many consist entirely of mixed pixels, the pixel number range of simulated objects is set to (3, 10), and the maximum abundance range is adjusted to (0.3, 1).

Algorithm 1 Hyperspectral Point Object Detection Network

Input: Training hyperspectral image set D_{train} , test hyperspectral image set D_{test} , point object detection model M , learning rate η , number of epochs E , batch size B

Output: Trained model M^* , detection results R on the test set

- 1: **Training Phase:**
- 2: **for** each epoch $e = 1$ to E **do**
- 3: Shuffle the training dataset D_{train}
- 4: **for** each batch $b = 1$ to $\lceil |D_{train}|/B \rceil$ **do**
- 5: Sample a batch of training hyperspectral images X_b and corresponding labels Y_b from D_{train}
- 6: Forward pass: Compute model output $\hat{Y}_b = M(X_b)$
- 7: Compute loss $L_b = \text{Loss}(\hat{Y}_b, Y_b)$
- 8: Backward pass: Compute gradients ∇L_b
- 9: Update model parameters: $M \leftarrow M - \eta \nabla L_b$
- 10: **end for**
- 11: **end for**
- 12: Save the trained model $M^* = M$
- 13: **Inference Phase:**
- 14: **for** each test hyperspectral image $x \in D_{test}$ **do**
- 15: Forward pass: Compute model output $\hat{y} = M^*(x)$
- 16: Apply post-processing (e.g., non-maximum suppression) to \hat{y} to obtain detection result r
- 17: Add detection result r to the result set R
- 18: **end for**
- 19: **return** M^*, R

Algorithm 2 Forward Pass Process of SpecDETR

Input: Hyperspectral image cube $\mathbf{X} \in \mathbb{R}^{H \times W \times N}$, number of encoder layers E , number of decoder layers D , number of object queries for matching Q_{match} , number of object queries for denoising training Q_{DN} , normalization constant V

Output: Predicted bboxes \mathbf{B}^{pred} , class confidence scores \mathbf{C}^{pred}

1: **Data Tokenization:**

2: Initialize pixel tokens: $\mathbf{P}_0 = \text{LN}(\text{linear}(\mathbf{X}/V))$

3: Initialize global token: $\mathbf{g}_0 = \left(\sum_{i=1}^{H \times W} \mathbf{p}_{0,i} \right) / (H \times W)$

4: $\mathbf{F}_0 = \{\mathbf{P}_0, \mathbf{g}_0\}$

5: **Transformer Encoder:**

6: **for** $j = 1$ to E **do**

7: $\tilde{\mathbf{F}}_j = \text{LN}(\text{Self-S2A}(\mathbf{F}_{j-1}) + \mathbf{F}_{j-1})$

8: $\mathbf{F}_j = \text{LN}(\text{FFN}(\tilde{\mathbf{F}}_j) + \tilde{\mathbf{F}}_j)$

9: **end for**

10: **Transformer Decoder:**

11: Generate initial anchor boxes $\mathbf{b}_{0,i}$ and class scores $\mathbf{c}_{0,i}$:

12: **for** $i = 1$ to $H \times W$ **do**

13: $\mathbf{b}_{0,i} = \text{Sig}(\text{breg}_0(\mathbf{p}_{E,i}) + \text{InSig}(\mathbf{b}_{\text{int},i}))$

14: $\mathbf{c}_{0,i} = \text{Sig}(\text{cls}_0(\mathbf{p}_{E,i}))$

15: **end for**

16: Select top Q_{match} anchor boxes.

17: **if** training **then**

18: $Q = Q_{\text{match}} + Q_{\text{DN}}$

19: Generate noised GT boxes $\mathbf{b}_{0,i}$ for $i = Q_{\text{match}} + 1$ to Q

20: **else**

21: $Q = Q_{\text{match}}$

22: **end if**

23: Initialize object queries $\mathbf{q}_{0,i} = \mathbf{1}$ for $i = 1$ to Q

24: **for** $d = 1$ to D **do**

25: **for** $i = 1$ to Q **do**

26: $\tilde{\mathbf{q}}_{d,i} = \text{LN}(\text{Cross-S2A}(\mathbf{q}_{d-1,i}, \mathbf{F}_E) + \mathbf{q}_{d-1,i})$

27: $\mathbf{q}_{d,i} = \text{LN}(\text{MLP}(\tilde{\mathbf{q}}_{d,i}) + \tilde{\mathbf{q}}_{d,i})$

28: Update bbox $\mathbf{b}_{d,i}$ and class scores $\mathbf{c}_{d,i}$:

29: $\hat{\mathbf{b}}_{d,i} = \text{Sig}(\text{breg}_d(\mathbf{q}_{d,i}) + \text{InSig}(\mathbf{b}_{d-1,i}))$

30: $\mathbf{b}_{d,i} = \text{detach}(\hat{\mathbf{b}}_{d,i})$

31: $\mathbf{b}_{d,i}^{\text{pred}} = \text{Sig}(\text{breg}_d(\mathbf{q}_{d,i}) + \text{InSig}(\hat{\mathbf{b}}_{d-1,i}))$

32: $\mathbf{c}_{d,i} = \text{Sig}(\text{cls}_d(\mathbf{q}_{d,i}))$

33: **end for**

34: **end for**

35: **return** $\mathbf{B}^{\text{pred}} = \{\mathbf{b}_{D,i}^{\text{pred}}\}_{i=1}^Q, \mathbf{C}^{\text{pred}} = \{\mathbf{c}_{D,i}\}_{i=1}^Q$

Algorithm 3 Hybrid Label Assigner of SpecDETR**Input:** \mathcal{G} : Set of ground GT boxes, $\mathcal{G} = \{g_1, g_2, \dots, g_N\}$. \mathcal{P} : Set of predicted boxes, $\mathcal{P} = \{p_1, p_2, \dots, p_M\}$. τ_{IoU} : IoU threshold for dynamic matching. T : Maximum number of positive samples per GT box for dynamic matching.**Output:** \mathcal{A} : Set of assigned positive samples, $\mathcal{A} = \{(g_i, p_j)\}$ 1: **Initialize:** $\mathcal{A} \leftarrow \emptyset$.2: **Step 1: Forced Matching**3: Compute the combined loss $\mathcal{L}_{\text{comb}}$ for all (g_i, p_j) pairs: $\mathcal{L}_{\text{comb}}(g_i, p_j) = \mathcal{L}_{\text{GIoU}}(g_i, p_j) + \mathcal{L}_{\text{L1}}(g_i, p_j) + \mathcal{L}_{\text{cls}}(g_i, p_j)$

4: Perform bipartite matching using the Hungarian algorithm:

 $\mathcal{M}_{\text{forced}} \leftarrow \text{Hungarian}(\mathcal{L}_{\text{comb}})$

5: Assign one predicted box to each GT box:

 $\mathcal{A} \leftarrow \mathcal{A} \cup \mathcal{M}_{\text{forced}}$ 6: **Step 2: Dynamic Matching**7: Compute the IoU matrix IoU between \mathcal{G} and \mathcal{P} : $\text{IoU}_{i,j} = \text{IoU}(g_i, p_j)$ 8: **for** each $g_i \in \mathcal{G}$ **do**9: Find candidate predicted boxes with $\text{IoU}_{i,j} > \tau_{\text{IoU}}$: $\mathcal{C}_i \leftarrow \{p_j \mid \text{IoU}_{i,j} > \tau_{\text{IoU}}\}$ 10: Sort \mathcal{C}_i by IoU in descending order: $\mathcal{C}_i^{\text{sorted}} \leftarrow \text{Sort}(\mathcal{C}_i, \text{by } \text{IoU}_{i,j})$ 11: Retain the top T predicted boxes: $\mathcal{C}_i^{\text{top}} \leftarrow \mathcal{C}_i^{\text{sorted}}[:T]$

12: Assign additional positive samples:

 $\mathcal{A} \leftarrow \mathcal{A} \cup \{(g_i, p_j) \mid p_j \in \mathcal{C}_i^{\text{top}}\}$ 13: **end for**14: **return** \mathcal{A}

V. ALGORITHMIC PSEUDOCODE

To more clearly illustrate the differences between hyperspectral point object detection and the traditional HTD, we present the processing flow of the hyperspectral point object detection framework, as shown in Algorithm 1. Additionally, Algorithm 2 outlines the forward propagation process of SpecDETR on a single hyperspectral image, while Algorithm 3 introduces the steps of the hybrid label assigner in SpecDETR.

VI. COMPARISON OF SPECDETR AND SPARSE REPRESENTATION-BASED HTD METHODS FOR MULTI CLASSES

Classic sparse representation-based methods in HTD are also adept at simultaneously detecting sub-pixel targets across multiple classes, where different classes correspond to different atoms. Sparse representation-based HTD methods can be primarily categorized into two types. The first type involves performing two sparse reconstructions on the test pixels: one using only background spectra as atoms, and the other using target spectra or a combination of target and background spectra as atoms. The difference between the reconstruction errors from these two processes is used as the detection score. The second type, represented by LSSA [12], performs a single sparse reconstruction using a combination of target and background spectra as atoms, and regards the coefficients of the target atoms as the detection scores. The first type of sparse representation-based method cannot directly identify the class of the detected target; therefore, in the comparative experiments on the SPOD and Avon datasets, we conduct separate detection for each object class. To compare the performance of this type of methods with SpecDETR in multi-class target detection, we treat all objects in the SPOD dataset as a single class and re-run this type of HTD method on the SPOD dataset, and also re-evaluate SpecDETR under the single-class object setting. As shown in Table E, even when all objects are treated as the same class, this type of sparse representation method still underperforms SpecDETR in object-level evaluation. Furthermore, we improve LSSA to obtain prediction results for all object classes in a single run. However, in the comparative experiments on the SPOD and Avon datasets, LSSA also perform less effectively than SpecDETR in object-level evaluation.

Although sparse representation-based methods can utilize multi-class object spectral information, they still struggle to leverage multi-class object spatial information. For instance, in the SanDiego dataset, despite the low spatial resolution, aircraft still possess certain morphological information. On our developed SPOD dataset, C7 and C8 are combination objects consisting of two and three single-spectrum objects, respectively, meaning they actually have spatial characteristics. Current sparse representation-based HTD methods are designed from the perspective of prior target spectra, making it difficult for them to learn object spatial characteristics like SpecDETR.

VII. LOSS FUNCTION

SpecDETR uses the same loss function as DINO, which includes L1 loss and GIoU loss [25] for box regression, as

well as focal loss [26] for classification. For each layer of the decoder, bbox L1 loss, bbox GIoU loss, and classification loss are computed individually for both the denoising and matching branches. Furthermore, these losses are also calculated for the candidate anchor boxes derived from the encoder's output features. The total loss employed for backpropagation is a weighted sum of all these individual losses. We use the same loss coefficients as those used in DINO, setting 1.0 for classification loss, 5.0 for bbox L1 loss, 2.0 for GIoU loss.

VIII. GENERALIZATION ANALYSIS

The training set and test set of the SPOD dataset we developed are both generated using the same simulation mechanism. The training set comprises 100 images, while the test set includes 500 images. Table. C presents the detection performance of SpecDETR on the training and test sets through cross-validation. As illustrated in Fig. L, SpecDETR has been sufficiently fitted after 100 training epochs. Regardless of whether the training set or test set is used as the training images, when the training images are used as validation images, the average metrics across all categories are essentially consistent and close to 1, with only slight differences in the AP metrics for each category. The mAP values are 0.979 and 0.980, mAP25 are 0.997 and 1.000, mAR are 0.984 and 0.986, and mRe25 are both 1.000. This phenomenon indicates that SpecDETR has a small bias on the SPOD Dataset. When validated on unseen images during training, the detection performance of SpecDETR decreases, showing a clear variance. Additionally, increasing the scale of the training data reduces the variance of SpecDETR and enhances the network's generalization capability. When trained on the training set and validated on the test set, SpecDETR's mAP drops to 0.856, and mAR drops to 0.897. Conversely, when trained on the test set and validated on the training set, SpecDETR's mAP is 0.922, and mAR is 0.942.

The test set of the Avon dataset only consists of a single image containing 24 real single-spectrum point objects, each with at least one pixel of high object spectral abundance, while the training set comprises 150 simulated images. Given that the test set of the Avon dataset comprises only one image, we solely analyze the performance of SpecDETR using the training set of the Avon Dataset as the training data. Table. D presents the detection performance of SpecDETR on the Avon dataset when validated on the training and test sets. The detection metrics of SpecDETR on the training set are nearly 1, indicating a small bias of SpecDETR on the AVON dataset. Considering the potential errors in manually annotated labels in the test set, we focus on the metrics mAP25 and mRe25 calculated at an IoU threshold of 0.25. SpecDETR achieves mAP25 of 0.990 and mRe25 of 1.000 on the test set, demonstrating good generalization capability of SpecDETR trained on simulated samples to real objects in the Avon dataset.

TABLE C
PERFORMANCE COMPARISON OF SPECDETR ON THE SPOD DATASET UNDER DIFFERENT TRAINING AND VALIDATION DATA COMBINATIONS.

Training	Validation	mAP↑	mAP ₂₅ ↑	mAR↑	mRe ₂₅ ↑	AP _{C1} ↑	AP _{C2} ↑	AP _{C3} ↑	AP _{C4} ↑	AP _{C5} ↑	AP _{C6} ↑	AP _{C7} ↑	AP _{C8} ↑
Training set	Training set	0.980	1.000	0.984	1.000	1.000	1.000	0.989	0.967	0.958	1.000	0.990	0.937
	Test set	0.856	0.938	0.897	0.955	0.963	0.969	0.970	0.698	0.648	0.905	0.844	0.850
Test set	Test set	0.979	0.997	0.986	1.000	0.990	1.000	1.000	0.975	0.959	0.977	0.967	0.967
	Training set	0.922	0.974	0.942	0.980	0.947	0.990	0.980	0.895	0.823	0.964	0.883	0.891

TABLE D
PERFORMANCE COMPARISON OF SPECDETR ON THE A DATASET UNDER DIFFERENT VALIDATION DATA.

Validation	mAP↑	mAP ₂₅ ↑	mAR↑	mRe ₂₅ ↑	AP _{BL} ↑	AP _{BR} ↑
Training set	0.979	0.998	0.991	1.000	0.987	0.971
Test set	0.885	0.990	0.925	1.000	0.924	0.847

TABLE E
PERFORMANCE COMPARISON OF SPECDETR AND SPARSE REPRESENTATION-BASED HTD METHODS ON THE SPOD DATASET UNDER THE SINGLE-CLASS OBJECT SETTING.

Method	AUC	IoU	AP	AP ₂₅	AR	Re ₂₅
CR [22]	0.943	0.446	0.031	0.346	0.144	0.566
KSR [23]	0.961	0.342	0.018	0.202	0.143	0.556
KSRBBH [24]	0.915	0.506	0.073	0.341	0.279	0.640
SpecDETR	-	-	0.907	0.989	0.930	0.994

IX. CONVERGENCE ANALYSIS

Figs. L and M present the training loss and mAP curves over epochs for our proposed SpecDETR and the current SOTA object detection network DINO on the SPOD dataset. Both SpecDETR and DINO are trained for 100 epochs, with the learning rate reduced to 0.1 of its original value at the 90th epoch. It is evident that during the first 50 epochs, SpecDETR exhibits a faster increase in mAP and a faster decrease in training loss compared to DINO. This indicates that SpecDETR has a superior convergence speed over DINO. In the latter 40 epochs, the training loss curves of SpecDETR and DINO are largely consistent, and the training loss remains essentially unchanged in the final few epochs, suggesting that both SpecDETR and DINO are well-fitted to the training set after training. However, the final mAP of SpecDETR is significantly higher than that of DINO, demonstrating that SpecDETR has better generalization and convergence properties than DINO.

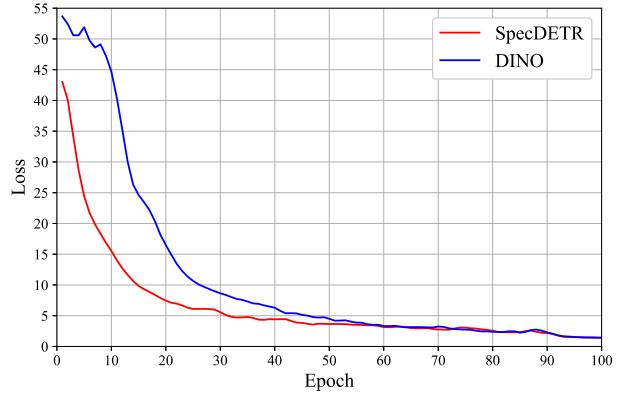


Fig. L. Training Loss convergence curves of SpecDETR and DINO on the SPOD dataset.

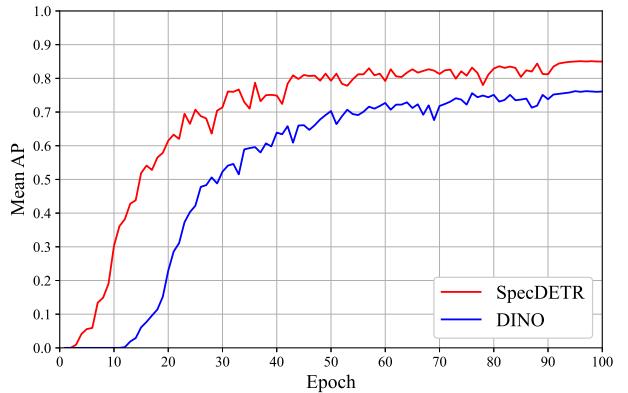


Fig. M. AP Convergence curves of SpecDETR and DINO on the SPOD dataset.

TABLE F
OBJECT-LEVEL DETECTION PERFORMANCE OF ASD UNDER DIFFERENT SEGMENTATION PIXEL NUMBER SETTINGS ON THE AVON DATASET.

Pixels	10	20	30	40	50	60	70	80	90	100
mAP	0.000	0.000	0.009	0.029	0.101	0.198	0.262	0.294	0.318	0.299
Pixels	110	120	130	140	150	160	170	180	190	200
mAP	0.252	0.245	0.219	0.201	0.187	0.172	0.160	0.153	0.146	0.143

X. DISCUSSION ON THE CONVERSION OF DETECTION SCORE MAPS TO OBJECT-LEVEL RESULTS FOR THE HTD METHODS

Current literature related to the HTD methods often evaluates the detection score map of a single detection image using ROC curves and AUC. However, our extended point object detection task employs an object-level evaluation approach, necessitating the binarization of the HTD methods' detection score maps and their conversion into object-level predictions. In our experiments, to maximize the HTD methods' performance on object-level evaluation metrics, we adopt the segmentation threshold corresponding to the maximum segmentation IoU across all images in the test set for each object class. Segmentation IoU is a widely used metric in segmentation tasks and has also become a benchmark evaluation metric for the infrared small target detection task. Another conversion approach is to directly select the top K pixels from the HTD methods' detection score maps. However, determining the appropriate number of pixels is challenging. A smaller number of pixels may result in missed objects or undersized object shapes, while a larger number may cause oversized object shapes or the merging of different objects. On the SPOD dataset, the number of pixels for different object classes varies within the same test image, and it also varies across different test images. A fixed number of pixels is difficult to apply simultaneously to different test images and object classes. Table. F presents the mAP of ASD under different segmentation pixel number settings on the Avon dataset. The Avon dataset consists of only one test image, and the pixel numbers for the two object classes are similar. Nevertheless, it is evident that the mAP is sensitive to the number of segmentation pixels. When the pixel number is 90, ASD's mAP reaches its best at 0.318, slightly lower than the 0.324 mAP under the optimal segmentation IoU criterion. When the pixel number is 70, the mAP drops to 0.262, and when the pixel number is 110, the mAP drops to 0.252. Therefore, selecting the top K pixels from the HTD methods' detection score maps is not well-suited for object-level evaluation.

TABLE G

AP PERFORMANCE COMPARISON OF SPECDETR AND MORE OBJECT DETECTION NETWORKS ON THE SPOD DATASET. DEFORMABLE DETR REPRESENTS THE ORIGINAL VERSION, DEFORMABLE DETR+ REPRESENTS DEFORMABLE DETR WITH ITERATIVE BOUNDING BOX REFINEMENT, AND DEFORMABLE DETR++ REPRESENTS TWO-STAGE DEFORMABLE DETR WITH ITERATIVE BOUNDING BOX REFINEMENT. TABLE II REPORTS THE RESULTS OF DEFORMABLEDETR++.

Detector	Bakebone	Epochs	Size	mAP↑	mAP ₂₅ ↑	mAP ₅₀ ↑	mAP ₇₅ ↑	AP _{C1} ↑	AP _{C2} ↑	AP _{C3} ↑	AP _{C4} ↑	AP _{C5} ↑	AP _{C6} ↑	AP _{C7} ↑	AP _{C8} ↑
Faster R-CNN [27]	ResNet50 [28]	100	×4	0.197	0.377	0.374	0.179	0.000	0.000	0.000	0.035	0.026	0.537	0.430	0.550
Faster R-CNN [27]	ResNet101 [28]	100	×4	0.207	0.368	0.365	0.209	0.000	0.000	0.000	0.027	0.015	0.587	0.450	0.577
Faster R-CNN [27]	Res2Net101 [29]	100	×4	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
Faster R-CNN [27]	RegNetX [30]	100	×4	0.227	0.379	0.378	0.242	0.000	0.000	0.000	0.042	0.043	0.631	0.522	0.578
Faster R-CNN [27]	ResNeSt50 [31]	100	×4	0.246	0.316	0.316	0.277	0.000	0.000	0.000	0.008	0.003	0.644	0.564	0.747
Faster R-CNN [27]	ResNeSt101 [31]	100	×4	0.183	0.253	0.252	0.216	0.000	0.000	0.000	0.006	0.001	0.386	0.392	0.681
Faster R-CNN [27]	ResNeXt101 [32]	100	×4	0.220	0.368	0.366	0.231	0.000	0.000	0.000	0.016	0.010	0.618	0.518	0.596
Faster R-CNN [27]	HRNet [33]	100	×4	0.320	0.404	0.402	0.345	0.000	0.000	0.000	0.107	0.076	0.849	0.731	0.793
PAA [1]	ResNet50 [28]	100	×4	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
PAA [1]	ResNet101 [28]	100	×4	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
TOOD [34]	ResNet50 [28]	100	×4	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
TOOD [34]	ResNeXt101 [32]	100	×4	0.304	0.464	0.440	0.303	0.000	0.000	0.000	0.181	0.194	0.743	0.648	0.663
FCOS [3]	ResNet50 [28]	100	×4	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
CenterNet [2]	ResNet50 [28]	100	×4	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
FreeAnchor [4]	ResNet50 [28]	100	×4	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
CentripetalNet [35]	HourglassNet104 [36]	100	×4	0.695	0.829	0.805	0.673	0.831	0.888	0.915	0.373	0.367	0.810	0.655	0.725
CornerNet [37]	HourglassNet104 [36]	100	×4	0.626	0.736	0.712	0.609	0.797	0.751	0.855	0.328	0.308	0.768	0.554	0.644
RepPoints [38]	ResNet50 [28]	100	×4	0.207	0.691	0.572	0.074	0.043	0.143	0.269	0.071	0.073	0.372	0.265	0.417
RepPoints [38]	ResNeXt101 [32]	100	×4	0.485	0.806	0.790	0.540	0.373	0.561	0.632	0.242	0.253	0.658	0.508	0.649
RetinaNet [38]	ResNet50 [28]	100	×4	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
RetinaNet [38]	EffcientNet [39]	100	×4	0.462	0.836	0.811	0.466	0.566	0.602	0.530	0.182	0.210	0.566	0.471	0.569
RetinaNet [38]	PVTv2-B3 [40]	100	×4	0.426	0.757	0.734	0.442	0.356	0.478	0.458	0.209	0.232	0.563	0.470	0.644
ATSS [5]	ResNet50 [28]	100	×4	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
ATSS [5]	Swin-L [41]	100	×4	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
Sparse R-CNN [6]	ResNet50 [28]	100	×4	0.011	0.109	0.048	0.001	0.000	0.000	0.000	0.003	0.003	0.015	0.023	0.042
Sparse R-CNN [6]	ResNet101 [28]	100	×4	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
DETR [7]	ResNet50 [28]	100	×4	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
ConditionalDETR [8]	ResNet50 [28]	100	×4	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
DAB-DETR [9]	ResNet50 [28]	100	×4	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
DeformableDETR [42]	ResNet50 [28]	100	×4	0.022	0.174	0.089	0.002	0.001	0.002	0.001	0.005	0.006	0.031	0.024	0.109
DeformableDETR+ [42]	ResNet50 [28]	100	×4	0.106	0.410	0.321	0.032	0.008	0.028	0.154	0.031	0.031	0.206	0.105	0.283
DeformableDETR++ [42]	ResNet50 [28]	100	×4	0.231	0.692	0.560	0.147	0.230	0.316	0.234	0.077	0.070	0.289	0.238	0.395
DINO [43]	ResNet50 [28]	100	×4	0.168	0.491	0.418	0.097	0.020	0.047	0.277	0.080	0.064	0.286	0.213	0.360
DINO [43]	Swin-L [41]	100	×4	0.757	0.852	0.842	0.764	0.915	0.912	0.951	0.483	0.497	0.847	0.728	0.721
SpecDETR	-	24	×1	0.706	0.877	0.873	0.734	0.946	0.966	0.937	0.360	0.357	0.757	0.645	0.682
		36	×1	0.799	0.907	0.903	0.821	0.958	0.965	0.961	0.562	0.501	0.869	0.777	0.799
		100	×1	0.856	0.938	0.930	0.863	0.963	0.969	0.970	0.698	0.648	0.905	0.844	0.850

TABLE H
AR PERFORMANCE COMPARISON OF SPECDETR AND MORE OBJECT DETECTION NETWORKS ON THE SPOD DATASET.

Detector	Bakebone	Epochs	Size	mAR↑	mAR ₂₅ ↑	mAR ₅₀ ↑	mAR ₇₅ ↑	AR _{C1} ↑	AR _{C2} ↑	AR _{C3} ↑	AR _{C4} ↑	AR _{C5} ↑	AR _{C6} ↑	AR _{C7} ↑	AR _{C8} ↑
Faster R-CNN [27]	ResNet50 [28]	100	×4	0.245	0.419	0.414	0.254	0.000	0.000	0.000	0.102	0.083	0.615	0.532	0.627
Faster R-CNN [27]	ResNet101 [28]	100	×4	0.248	0.393	0.390	0.272	0.000	0.000	0.000	0.081	0.046	0.660	0.553	0.645
Faster R-CNN [27]	Res2Net101 [29]	100	×4	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
Faster R-CNN [27]	RegNetX [30]	100	×4	0.267	0.399	0.398	0.304	0.000	0.000	0.000	0.089	0.090	0.699	0.607	0.654
Faster R-CNN [27]	ResNeSt50 [31]	100	×4	0.269	0.319	0.318	0.294	0.000	0.000	0.000	0.007	0.004	0.682	0.633	0.828
Faster R-CNN [27]	ResNeSt101 [31]	100	×4	0.201	0.252	0.252	0.229	0.000	0.000	0.000	0.003	0.001	0.410	0.437	0.760
Faster R-CNN [27]	ResNeXt101 [32]	100	×4	0.253	0.376	0.375	0.286	0.000	0.000	0.000	0.027	0.022	0.701	0.606	0.665
Faster R-CNN [27]	HRNet [33]	100	×4	0.364	0.434	0.430	0.386	0.000	0.000	0.000	0.185	0.155	0.899	0.818	0.859
PAA [1]	ResNet50 [28]	100	×4	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
PAA [1]	ResNet101 [28]	100	×4	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
TOOD [34]	ResNet50 [28]	100	×4	0.000	0.004	0.001	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.001	0.001
TOOD [34]	ResNeXt101 [32]	100	×4	0.401	0.570	0.532	0.391	0.000	0.000	0.000	0.429	0.442	0.816	0.756	0.764
FCOS [3]	ResNet50 [28]	100	×4	0.000	0.003	0.002	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.004
CenterNet [2]	ResNet50 [28]	100	×4	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
FreeAnchor [4]	ResNet50 [28]	100	×4	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
CentripetalNet [35]	HourglassNet104 [36]	100	×4	0.840	0.956	0.932	0.817	0.878	0.936	0.939	0.769	0.761	0.866	0.766	0.801
CornerNet [37]	HourglassNet104 [36]	100	×4	0.855	0.969	0.939	0.834	0.917	0.949	0.940	0.804	0.791	0.881	0.749	0.806
RepPoints [38]	ResNet50 [28]	100	×4	0.346	0.934	0.804	0.229	0.185	0.279	0.379	0.275	0.272	0.467	0.402	0.505
RepPoints [38]	ResNeXt101 [32]	100	×4	0.635	0.961	0.939	0.725	0.523	0.662	0.713	0.563	0.565	0.735	0.606	0.714
RetinaNet [38]	ResNet50 [28]	100	×4	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
RetinaNet [38]	EfficientNet [39]	100	×4	0.611	0.971	0.949	0.667	0.664	0.694	0.672	0.493	0.488	0.665	0.570	0.641
RetinaNet [38]	PVTv2-B3 [40]	100	×4	0.650	0.987	0.973	0.718	0.549	0.638	0.644	0.607	0.596	0.755	0.671	0.737
ATSS [5]	ResNet50 [28]	100	×4	0.001	0.004	0.003	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.005
ATSS [5]	Swin-L [41]	100	×4	0.001	0.003	0.002	0.001	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.004
Sparse R-CNN [6]	ResNet50 [28]	100	×4	0.091	0.477	0.292	0.032	0.000	0.000	0.000	0.071	0.068	0.178	0.195	0.215
Sparse R-CNN [6]	ResNet101 [28]	100	×4	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
DETR [7]	ResNet50 [28]	100	×4	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
ConditionalDETR [8]	ResNet50 [28]	100	×4	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
DAB-DETR [9]	ResNet50 [28]	100	×4	0.000	0.015	0.001	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.001
DeformableDETR [42]	ResNet50 [28]	100	×4	0.111	0.639	0.354	0.039	0.018	0.023	0.019	0.082	0.079	0.212	0.189	0.263
Deformable DETR+ [42]	ResNet50 [28]	100	×4	0.265	0.753	0.636	0.173	0.094	0.139	0.331	0.186	0.199	0.430	0.308	0.429
Deformable DETR++ [42]	ResNet50 [28]	100	×4	0.385	0.883	0.770	0.324	0.348	0.439	0.405	0.283	0.267	0.440	0.386	0.512
DINO [43]	ResNet50 [28]	100	×4	0.368	0.763	0.684	0.339	0.120	0.177	0.436	0.328	0.336	0.499	0.469	0.577
DINO [43]	Swin-L [41]	100	×4	0.909	0.983	0.976	0.914	0.948	0.942	0.973	0.869	0.853	0.933	0.888	0.865
SpecDETR	-	24	×1	0.784	0.924	0.919	0.821	0.963	0.975	0.955	0.560	0.491	0.821	0.742	0.768
		36	×1	0.854	0.937	0.933	0.873	0.974	0.970	0.975	0.696	0.602	0.909	0.841	0.867
		100	×1	0.897	0.955	0.948	0.902	0.975	0.974	0.981	0.766	0.742	0.945	0.892	0.903

REFERENCES

- [1] K. Kim and H. S. Lee, "Probabilistic anchor assignment with iou prediction for object detection," in *European Conference on Computer Vision*. Springer, 2020, pp. 355–371.
- [2] X. Zhou, D. Wang, and P. Krähenbühl, "Objects as points," *arXiv*, 2019.
- [3] Z. Tian, C. Shen, H. Chen, and T. He, "Fcos: Fully convolutional one-stage object detection," in *International Conference on Computer Vision*, 2019, pp. 9627–9636.
- [4] X. Zhang, F. Wan, C. Liu, R. Ji, and Q. Ye, "Freeanchor: Learning to match anchors for visual object detection," *Neural Information Processing Systems*, vol. 32, 2019.
- [5] S. Zhang, C. Chi, Y. Yao, Z. Lei, and S. Z. Li, "Bridging the gap between anchor-based and anchor-free detection via adaptive training sample selection," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2020, pp. 9759–9768.
- [6] P. Sun, R. Zhang, Y. Jiang, T. Kong, C. Xu, W. Zhan, M. Tomizuka, L. Li, Z. Yuan, C. Wang *et al.*, "Sparse r-cnn: End-to-end object detection with learnable proposals," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2021, pp. 14454–14463.
- [7] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko, "End-to-end object detection with transformers," in *European Conference on Computer Vision*, 2020, pp. 213–229.
- [8] D. Meng, X. Chen, Z. Fan, G. Zeng, H. Li, Y. Yuan, L. Sun, and J. Wang, "Conditional DETR for fast training convergence," in *International Conference on Computer Vision*, 2021, pp. 3651–3660.
- [9] S. Liu, F. Li, H. Zhang, X. Yang, X. Qi, H. Su, J. Zhu, and L. Zhang, "DAB-DETR: Dynamic anchor boxes are better queries for DETR," *arXiv*, 2022.
- [10] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft COCO: Common objects in context," in *European Conference on Computer Vision*, 2014, pp. 740–755.
- [11] K. Chen, J. Wang, J. Pang, Y. Cao *et al.*, "MMDetection: Open mmlab detection toolbox and benchmark," *arXiv*, 2019.
- [12] D. Zhu, B. Du, and L. Zhang, "Learning single spectral abundance for hyperspectral subpixel target detection," *IEEE Transactions on Neural Networks and Learning Systems*, 2023.
- [13] D. Shen, X. Ma, W. Kong, J. Liu, J. Wang, and H. Wang, "Hyperspectral target detection based on interpretable representation network," *IEEE Transactions on Geoscience and Remote Sensing*, 2023.
- [14] J. Jiao, Z. Gong, and P. Zhong, "Triplet spectral-wise transformer network for hyperspectral target detection," *IEEE Transactions on Geoscience and Remote Sensing*, 2023.
- [15] S. Kraut and L. L. Scharf, "The CFAR adaptive subspace detector is a scale-invariant GLRT," *IEEE Transactions on Signal Processing*, vol. 47, no. 9, pp. 2538–2541, Sep. 1999.
- [16] J. C. Harsanyi, "Detection and classification of subpixel spectral signatures in hyperspectral image sequences," Ph.D. dissertation, University of Maryland Baltimore County, Baltimore, MD, USA, 1993.
- [17] J. C. Harsanyi and C.-I. Chang, "Hyperspectral image classification and dimensionality reduction: An orthogonal subspace projection approach," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 32, no. 4, pp. 779–785, Jul. 1994.
- [18] H. Kwon and N. M. Nasrabadi, "Kernel orthogonal subspace projection for hyperspectral signal classification," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 43, no. 12, pp. 2952–2962, 2005.
- [19] F. C. Robey, D. R. Fuhrmann, E. J. Kelly, and R. Nitzberg, "A CFAR adaptive matched filter detector," *IEEE Transactions on Aerospace and Electronic Systems*, vol. 28, no. 1, pp. 208–216, Jan. 1992.
- [20] H. Kwon and N. M. Nasrabadi, "A comparative analysis of kernel subspace target detectors for hyperspectral imagery," *EURASIP Journal on Advances in Signal Processing*, vol. 2007, pp. 1–13, 2006.
- [21] H. Ren and C.-I. Chang, "Target-constrained interference-minimized approach to subpixel target detection for hyperspectral images," *Optical Engineering*, vol. 39, pp. 3138–3145, 2000.
- [22] W. Li and Q. Du, "Collaborative representation for hyperspectral anomaly detection," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 53, no. 3, pp. 1463–1474, Mar. 2015.
- [23] Y. Chen, N. M. Nasrabadi, and T. D. Tran, "Hyperspectral image classification via kernel sparse representation," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 51, no. 1, pp. 217–231, Jan. 2013.
- [24] Y. Zhang, L. Zhang, B. Du, and S. Wang, "A nonlinear sparse representation-based binary hypothesis model for hyperspectral target detection," *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, vol. 8, no. 6, pp. 2513–2522, Jun. 2015.
- [25] H. Rezatofighi, N. Tsoi, J. Gwak, A. Sadeghian, I. Reid, and S. Savarese, "Generalized intersection over union: A metric and a loss for bounding box regression," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 658–666.
- [26] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," in *International Conference on Computer Vision*, 2017, pp. 2999–3007.
- [27] S. Ren, K. He, R. Girshick, and J. Sun, "Faster R-CNN: Towards real-time object detection with region proposal networks," *IEEE transactions on pattern analysis and machine intelligence*, vol. 39, no. 6, pp. 1137–1149, 2016.
- [28] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 770–778.
- [29] S.-H. Gao, M.-M. Cheng, K. Zhao, X.-Y. Zhang, M.-H. Yang, and P. Torr, "Res2net: A new multi-scale backbone architecture," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 43, no. 2, pp. 652–662, 2019.
- [30] I. Radosavovic, R. P. Kosaraju, R. Girshick, K. He, and P. Dollár, "Designing network design spaces," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2020, pp. 10428–10436.
- [31] H. Zhang, C. Wu, Z. Zhang, Y. Zhu, H. Lin, Z. Zhang, Y. Sun, T. He, J. Mueller, R. Manmatha *et al.*, "ResNeSt: Split-attention networks," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2022, pp. 2736–2746.
- [32] S. Xie, R. Girshick, P. Dollár, Z. Tu, and K. He, "Aggregated residual transformations for deep neural networks," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 1492–1500.
- [33] K. Sun, B. Xiao, D. Liu, and J. Wang, "Deep high-resolution representation learning for human pose estimation," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 5693–5703.
- [34] C. Feng, Y. Zhong, Y. Gao, M. R. Scott, and W. Huang, "TOOD: Task-aligned one-stage object detection," in *International Conference on Computer Vision*, 2021, pp. 3490–3499.
- [35] Z. Dong, G. Li, Y. Liao, F. Wang, P. Ren, and C. Qian, "Centripetalnet: Pursuing high-quality keypoint pairs for object detection," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2020, pp. 10519–10528.
- [36] 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.
- [37] H. Law and J. Deng, "Cornernet: Detecting objects as paired keypoints," in *European Conference on Computer Vision*, 2018, pp. 734–750.
- [38] Z. Yang, S. Liu, H. Hu, L. Wang, and S. Lin, "Reppoints: Point set representation for object detection," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 9657–9666.
- [39] M. Tan and Q. Le, "Efficientnet: Rethinking model scaling for convolutional neural networks," in *International Conference on Machine Learning*, 2019, pp. 6105–6114.
- [40] W. Wang, E. Xie, X. Li, D.-P. Fan, K. Song, D. Liang, T. Lu, P. Luo, and L. Shao, "PVTv2: Improved baselines with pyramid vision transformer," *Computational Visual Media*, vol. 8, no. 3, pp. 415–424, 2022.
- [41] 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 *International Conference on Computer Vision*, 2021, pp. 10012–10022.
- [42] X. Zhu, W. Su, L. Lu, B. Li, X. Wang, and J. Dai, "Deformable DETR: Deformable transformers for end-to-end object detection," *arXiv*, 2020.
- [43] H. Zhang, F. Li, S. Liu, L. Zhang, H. Su, J. Zhu, L. M. Ni, and H.-Y. Shum, "DINO: DETR with improved denoising anchor boxes for end-to-end object detection," *arXiv*, 2022.