

Dual-Granularity Semantic Prompting for Language Guidance Infrared Small Target Detection

Zixuan Wang Haoran Sun Jiaming Lu Wenxuan Wang Zhongling Huang
Dingwen Zhang Xuelin Qian Junwei Han

School of Automation, Northwestern Polytechnical University
Xi'an, Shaanxi 710072, China

Abstract

Infrared small target detection remains challenging due to limited feature representation and severe background interference, resulting in sub-optimal performance. While recent CLIP-inspired methods attempt to leverage textual guidance for detection, they are hindered by inaccurate text descriptions and reliance on manual annotations. To overcome these limitations, we propose DGSPNet, an end-to-end language prompt-driven framework. Our approach integrates dual-granularity semantic prompts: coarse-grained textual priors (e.g., “infrared image”, “small target”) and fine-grained personalized semantic descriptions derived through visual-to-textual mapping within the image space. This design not only facilitates learning fine-grained semantic information but also can inherently leverage language prompts during inference without relying on any annotation requirements. By fully leveraging the precision and conciseness of text descriptions, we further introduce a text-guide channel attention (TGCA) mechanism and text-guide spatial attention (TGSA) mechanism that enhances the model’s sensitivity to potential targets across both low- and high-level feature spaces. Extensive experiments demonstrate that our method significantly improves detection accuracy and achieves state-of-the-art performance on three benchmark datasets.

1. Introduction

Infrared small target detection (IRSTD) aims to accurately localize tiny objects in infrared imagery captured by aerial and satellite platforms [1]. However, these targets often lack distinctive structural or textural features due to thermal diffusion, making them extremely difficult to differentiate from complex backgrounds. This challenge is further compounded by factors such as sensor noise, atmospheric interference, and highly dynamic scenes. Despite these dif-

ficulties,IRSTD remains crucial for mission-critical applications, including early-warning military systems, civil infrastructure monitoring, disaster response, and search-and-rescue operations. Consequently, the growing demand for robust and reliable detection has driven extensive research efforts in this field [19, 39].

Existing infrared small target detection (IRSTD) approaches, such as SCTrans [35] and DCFR-Net [8], primarily enhance target saliency by leveraging spatial-channel transformer mechanisms or detail-preserving convolutional networks. While these techniques improve feature representation and emphasize potential target regions, they predominantly rely on visual cues extracted from the image domain. This dependence becomes a major limitation in infrared imagery, where targets exhibit extremely low signal-to-noise ratios, weak contrast against complex backgrounds, and minimal structural or textural features [25]. Consequently, these methods often struggle to achieve robust detection performance in cluttered or highly dynamic scenes, underscoring the need for more adaptive and context-aware solutions.

To address these challenges, several pioneering studies have investigated visual-language models (VLMs) [13], leveraging text descriptions as semantic priors to guide target localization in infrared imagery. Although these text-guided approaches show considerable promise by introducing high-level contextual cues, they suffer from two critical limitations: **(1) Annotation overhead.** Generating precise text labels (e.g., “small aerial target in a sky-ground background”) requires domain expertise and extensive manual annotation, while automatically producing such descriptions using VLMs demands substantial computational resources. Both approaches impose prohibitive costs, severely limiting their scalability for large-scaleIRSTD deployment. **(2) Description discrepancy.** Text annotations—whether manually crafted or VLM-generated—often fail to accurately capture infrared-specific

properties. For example, low-intensity thermal signatures may be misinterpreted as noise, or target features may be confused with background clutter. Such semantic inaccuracies can propagate through the detection pipeline, ultimately degrading model performance.

To address these limitations, we propose an end-to-end framework for language-guided infrared small target detection, centered on **Dual-Granularity Semantic Prompts**. Specifically, (i) coarse-grained textual priors (e.g., “infrared image,” “small target”) provide general domain-aware context, while (ii) fine-grained, instance-specific semantic descriptions are generated via a visual-to-textual mapping within the image, capturing target-specific cues. This dual design enriches the feature space with complementary semantic cues, facilitating more discriminative representations of small targets.

Moreover, this framework naturally supports the use of language prompts during inference without requiring manual annotations, significantly reducing deployment overhead. To further exploit these semantic cues, we introduce a **Text-Guide Channel and Spatial Attention** mechanism, which integrates refined textual semantics into the model’s attention module. The mechanisms guide the model to concentrate on semantically relevant regions while suppressing background noise, thereby enhancing target saliency and improving localization precision. Our contributions can be summarized as follows:

- We propose a novel dual-granularity semantic prompting strategy that combines coarse-grained textual priors with fine-grained personalized semantic descriptions via visual-to-textual mapping, enhancing feature representation for small targets.
- We design a Text-Guide Channel and Spatial Attention mechanism that integrates refined semantic cues into the model, effectively directing focus toward semantically relevant regions while suppressing background interference, thereby enhancing target saliency and localization accuracy.
- Extensive experiments on benchmarkIRSTD datasets demonstrate that our framework significantly outperforms state-of-the-art methods, achieving superior detection in complex and dynamic infrared scenarios.

2. Related Work

2.1. Infrared Small Target Detection

Early infrared small target detection (IRSTD) methods were heavily reliant on handcrafted priors designed to enhance faint targets against noisy thermal backgrounds. Techniques such as filter-based methods, including Top-Hat and Laplacian-of-Gaussian (LoG) transformations, were commonly used to highlight target features by emphasizing differences in intensity [7, 41]. Additionally, local contrast-

based approaches that exploited regional saliency [2, 5] and sparse representation models for decomposing target and background signals [27] were also frequently applied. While these classical methods were computationally efficient, they proved highly sensitive to background clutter and often struggled in complex infrared scenes featuring non-uniform noise or dynamic structures, where distinguishing small targets from the background became difficult. With the advent of deep learning, convolutional neural networks (CNNs) revolutionizedIRSTD by offering more robust detection. Methods like ACUNet, IRNet, andIRSTDNet [6, 20, 43] utilized techniques such as multiscale fusion, encoder-decoder consistency, and detail-preserving designs to improve accuracy in detecting small targets. However, despite their advancements, CNN-based models still largely rely on low-level texture cues and lack strong semantic reasoning, leading to poor generalization when faced with ambiguous or unseen conditions.

To address these shortcomings, Transformer-based architectures, such as ISFormer and DGFormer [11, 30], were introduced. These models employ self-attention mechanisms to capture long-range dependencies and global context, which enhances feature representation, especially in cluttered scenes. However, even these Transformer-based models remain limited by their reliance on purely visual information, which restricts their semantic expressiveness and adaptability to diverse imaging conditions. The limitations of both classical methods and deep learning models highlight the need for more comprehensive approaches that can incorporate multi-modal or semantic information to better adapt to the complexities of infrared target detection.

2.2. Language-guidedIRSTD

Vision-Language Models (VLMs), such as CLIP [23] and ALIGN [14], have made significant strides by aligning images and text in a shared embedding space using contrastive learning. This alignment has led to remarkable performance across a wide range of tasks, including open-vocabulary detection (e.g., ViLD [10], GLIP [16]), image captioning, visual question answering, and dense prediction through multimodal frameworks like OFA [29], BLIP-2 [17], and UniVL [44]. Moreover, generative models like GIT [31] treat vision-language tasks as conditional generation problems, while instruction-following detectors, such as LLMDet and DetGPT [21, 22], combine large language models with detection heads to enable natural-language-driven vision tasks. Despite their tremendous success in the realm of RGB imagery, VLMs face significant challenges when adapted to infrared data, primarily due to modality gaps, such as low signal-to-noise ratios, weak target contrast, and the distinct physics behind infrared imaging. These differences make it difficult for VLMs to effectively interpret and detect targets in infrared images, limiting their

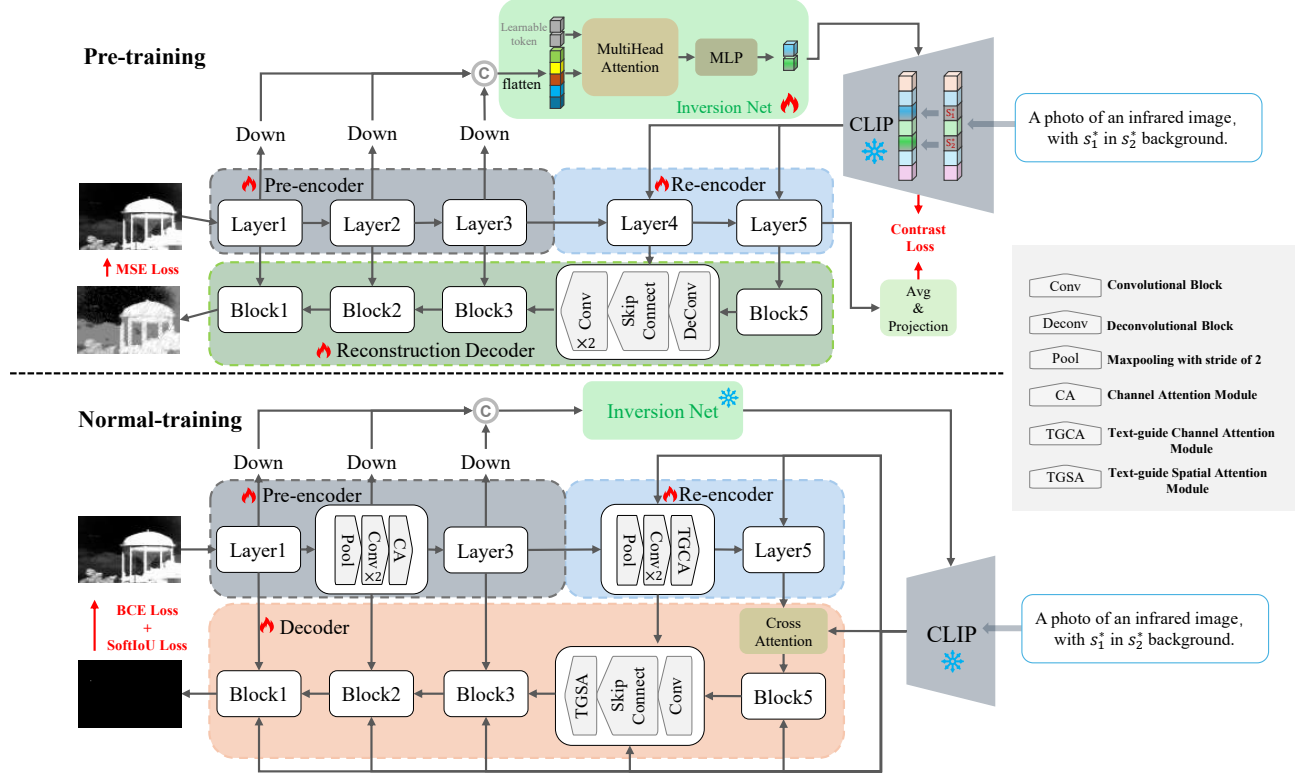


Figure 1. Overview of DGSPNet. The entire training process is divided into two phases: pre-training and normal training. During the pre-training phase, the network’s decoder is replaced with a reconstruction decoder, and contrastive loss is introduced to supervise the training of the inversion network. During the normal training phase, the weights of the inversion net are frozen. Each convolutional block consists of one convolution unit, a normalization layer, and a ReLU activation layer, and the deconvolutional block has the same composition.

performance in infrared small target detection (IRSTD).

To address this, recent efforts have sought to bridge the modality gap by exploring language-guided IRSTD. For example, Text-IRSTD [13] introduces fuzzy semantic prompts and a Progressive Cross-modal Semantic Interaction Decoder to dynamically align textual and visual features, improving target localization. However, these approaches still face challenges such as high annotation or computational overhead and an inability to effectively capture infrared-specific attributes, such as low-intensity thermal signals and background noise. These limitations often result in semantic mismatches between the textual descriptions and the infrared image content, which in turn degrade the overall detection performance. As a result, there remains a need for more efficient and precise methods that can bridge these gaps and fully leverage language-guided infrared detection.

3. Methodology

3.1. Overview

Figure 1 illustrates the schematic overview of the proposed DGSPNet, an end-to-end dual-granularity semantic prompt-

ing network for language-guided infrared small target detection. DGSPNet takes an infrared image as input and outputs a binary segmentation mask to localize small targets.

The framework consists of three principal components: *hierarchical image encoder*, *dual-granularity semantic prompt* and *prompt-driven semantic guidance*. The hierarchical image encoder is divided into a pre-encoder and a re-encoder, responsible for “low-to-middle-level feature extraction” and “text-guided high-level feature optimization”, respectively. The dual-granularity semantic prompt mechanism constructs multi-granularity textual descriptions that semantically align with the input infrared image, while the prompt-driven semantic guidance leverages these text prompts to dynamically guide the decoding process, enabling more accurate segmentation of small targets.

During training, the weights of the inversion net, which is designed for dual-granularity semantic prompt generation, are frozen. To obtain the weights of the inversion net, we elaborately design a reconstruction pre-training method. After pre-training, both the inversion net and the image encoder weights are loaded for formal training, providing a strong initialization that improves overall learning perfor-

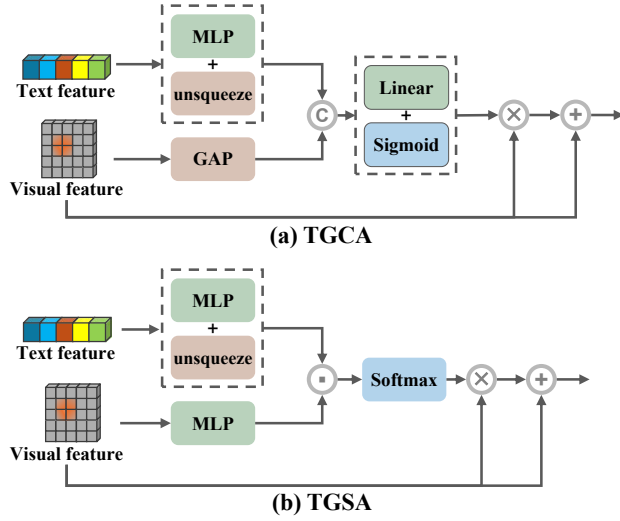


Figure 2. Illustration of TGCA and TGSA. Guided by text features, these modules generate channel and spatial attention weights, which are applied via multiplication and residual connections to enhance visual features.

mance.

In the following sections, we first introduce the hierarchical image encoder, which captures global contextual features while leveraging dual-granularity semantic information to emphasize target regions. Next, we describe the semantic prompt learning process, covering both coarse-grained priors (explicit descriptions) and fine-grained semantic cues (instance-specific, implicit descriptions). We then present the prompt-guided decoding strategy, which integrates attention mechanisms to model interactions between image features and text embeddings. Finally, we detail the training and inference procedures of DGSPNet.

3.2. Hierarchical Image Encoder

The hierarchical encoder comprises two stages: pre-encoding and re-encoding. The pre-encoding stage extracts basic visual features, such as edges, textures, and local brightness, providing a foundation for the inversion net to generate fine-grained semantic prompts. The re-encoding stage then leverages these dual-granularity semantic prompts to guide high-level feature optimization based on the pre-encoder outputs. To balance performance and efficiency, we adopt a U-Net architecture enhanced with attention mechanisms, enabling multi-resolution feature extraction and effective integration of semantic guidance.

Pre-encoder. Specifically, given an input infrared image $x \in \mathbb{R}^{H \times W}$, the pre-encoder consists of three convolutional blocks that iteratively extract multi-scale features $\{f^{(i)}\}_{i=1}^3 \in \mathbb{R}^{h_i \times w_i \times d_i}$, where h_i and w_i denote the height

and width of the i -th feature map and d_i is the corresponding feature dimension. With the exception of the first block, each block performs a downsampling operation on the feature maps, producing a final output $f^{(3)}$ with a spatial resolution reduced to $\frac{1}{4}$ of the original image.

For each pre-coding layer, assuming the input feature is f^{in} , max pooling operation is first applied to halve its spatial resolution, followed by two stacked 3×3 convolutional layers to extract deep features. To further enhance the representation ability of the encoder and focus on potential target regions, we integrate a channel attention layer at the end of each block to output the refined feature f^{out} . This layer combines global average pooling (GAP) and global max pooling (GMP) to compute channel-wise attention weights. By selectively emphasizing informative channels and suppressing redundant or noise signals, the encoder strengthens semantic cues relevant to the detection task while reducing the impact of complex background clutter. The overall process can be formulated as,

$$\begin{aligned} f^{\text{mid}} &= \gamma(\text{Conv}_{3 \times 3}(\gamma(\text{Conv}_{3 \times 3}(\text{MaxPool}(f^{\text{in}})))))) \\ f^{\text{out}} &= f^{\text{mid}} \cdot \sigma(\text{Conv}_{1 \times 1}(\text{GMP}(f^{\text{mid}}) + \text{GAP}(f^{\text{mid}}))) \end{aligned} \quad (1)$$

where $\text{Conv}_{k \times k}$ represents the convolution layer with $k \times k$ kernel size, γ and σ denote the rectified linear unit and sigmoid activation, respectively.

Re-encoder. Corresponding to the last two layers of the original UNet encoder, it takes the output features of the pre-encoder $f^{(3)}$ along with the dual-granularity semantic feature f^{text} as input, and outputs $\{f^{(i)}\}_{i=4}^5 \in \mathbb{R}^{h_i \times w_i \times d_i}$ by two re-encoder layers.

Like the layer-wise structure of the pre-encoder, each layer of the re-encoder comprises max-pooling, two convolutional layers, and a specialized channel attention mechanism. This specialized channel attention is referred to as Text-Guide Channel Attention(TGCA), which is designed to leverage semantic information as guidance, enabling the channel attention mechanism to further enhance or suppress specific regions within the image. So the re-encoder process can be illustrated as follows:

$$\begin{aligned} f^{\text{mid}} &= \gamma(\text{Conv}_{3 \times 3}(\gamma(\text{Conv}_{3 \times 3}(\text{MaxPool}(f^{\text{in}})))))) \\ f^{\text{out}} &= \text{TGCA}(f^{\text{mid}}, f^{\text{text}}) \end{aligned} \quad (2)$$

where TGCA guides the learning of channel attention weights for visual features through dual-granularity semantic text features. As shown in Figure 2(a), this module first generates text channel attention weights from text features via MLP, and adjusts their dimensionality through expansion to match the spatial dimensions of visual features. Meanwhile, global average pooling is applied to the visual features to obtain the channel attention weights inherent to the image itself. Subsequently, the text-guide channel attention weights and the image’s intrinsic channel attention

weights are concatenated fused by a linear layer to generate the final attention weights. After being normalized by the Sigmoid function, these weights are multiplied with the original visual features channel-wise. The result is then added to the original features via a residual connection, which ultimately outputs the text-guided enhanced visual feature representation. The specific formulas are as follows:

$$\begin{aligned} w^{\text{ch}} &= \sigma(\text{Fn}(\text{Concat}(\text{GAP}(f^{\text{mid}}), \text{MLP}(f^{\text{text}})))) \\ f^{\text{out}} &= f^{\text{mid}} + f^{\text{mid}} \cdot w^{\text{ch}} \end{aligned} \quad (3)$$

3.3. Dual-Granularity Semantic Prompt

Since small targets in infrared images lack clear structural features, relying only on visual cues often leads to unclear or incorrect detections, especially in complex backgrounds. Previous studies[13] have shown that incorporating information from the word-embedding space can improve the expressiveness and understanding of image features. However, how to generate accurate and relevant text descriptions in a cost-effective way is an inherent limitation. Unlike prior approaches that depend on prompt engineering or large pre-trained language models, we propose a dual-granularity semantic prompt module, which leverages text inversion techniques to provide fine-grained attribute information that aligns with the visual content and is otherwise hard to obtain. Details of coarse- and fine-grained prompts are explained below.

Coarse-grained Textual Descriptions. The coarse-grained text description is designed to capture easily accessible prior knowledge about an input image, such as the modality (e.g., infrared image) or scene type (e.g., sky, suburb). This kind of information can be obtained directly from metadata or generated using predefined templates, avoiding the cost of additional annotations. Inspired by prior works[23], we create prompt templates tailored for infrared small target detection. For example,

“A photo of an infrared image, with targets in the sky/ground/ocean background.”

The prompts are encoded in natural language and processed by a text encoder to generate semantic embeddings $f^{\text{text}} \in \mathbb{R}^{l \times d_t}$, where l denotes the sequence length of text embeddings and d_t is the hidden dimension. These embeddings, combined with image features, are fed to the prompt-driven modules for joint reasoning. Coarse-grained prompts serve as global priors, providing initial semantic context and guiding attention to task-relevant regions.

Fine-grained Semantic Tokens. Coarse-grained descriptions are often too general to capture specific image details, such as scene layout (e.g., the sea surface) or target category (e.g., a ship). In practice, detailed information is more helpful for accurately locating infrared targets. For example, if the model is given the text prompt *“a ship on the sea surface”*, any prediction of a target in the sky can be con-

fidently treated as a false alarm and corrected. However, such detailed information cannot be predefined using fixed templates since they vary across different images.

Therefore, we introduce a set of semantic tokens $\{s_i^*\}_{i=1}^{n_t} \in \mathbb{R}^{1 \times d}$ that are dynamically generated for each image using a text inversion strategy. n_t denotes the number of tokens and d denotes the number of dimension of the text encoder. More concretely, these tokens serve as potential language representations and are learned by the inversion net after downsampling and fusing the features from the first three layers of the image encoder. This allows the model to directly extract detailed scene semantics and local target features from the image itself. Afterwards, the refined tokens are projected into the text embedding space and placed into the coarse-grained descriptions, effectively transforming them into implicit, image-specific text prompts. Taking $n = 2$ as an example,

“A photo of an infrared image, with s_1^ in s_2^* background.”*

Together, the coarse-grained and fine-grained prompts form a complementary dual-granularity prompt. The former offers stable, general prior knowledge to guide the learning process, while the latter injects personalized, image-specific semantics to dynamically adapt the detection process.

Here, we explain how this process is implemented. For the first three layers of features ($f^{(1)}$, $f^{(2)}$, and $f^{(3)}$) output by the image encoder, we use depthwise convolutions with a kernel size of 3 and strides of 16, 8, and 4, respectively, to achieve downsampling. Afterward, we fuse these three layers of features through concatenation and a 1×1 convolution to obtain image features f^{img} that contain both global semantic information and local detailed information:

$$f^{\text{img}} = \text{Conv}_{1 \times 1} \left(\text{Concat} \left(\text{DWConv}(f^{(1)}, f^{(2)}, f^{(3)}) \right) \right) \quad (4)$$

Containing rich spatial detailed features, f^{img} is subsequently fed into the inversion net, which is designed to map visual features from the pre-encoder into fine-grained semantic tokens, serving as a bridge between low-to-middle-level visual cues and text semantics. Given the fused visual features from the pre-encoder, $f^{\text{img}} \in \mathbb{R}^{n \times c \times h \times w}$ (where n is the batch size, c is the number of channel dimension), the first step is to reshape f^{img} into a sequence format $f_{\text{seq}}^{\text{img}} \in \mathbb{R}^{(h \times w) \times n \times c}$. After that, we introduce a set of learnable tokens, denoted as $T \in \mathbb{R}^{n_t \times n \times c}$ (where n_t is the number of fine-grained semantic tokens), and feed them into the multi-head attention as the query. The key and value, however, are the combined feature of T and $f_{\text{seq}}^{\text{img}}$. This design enables T to capture spatially discriminative information from visual features and interact semantically through self-attention. After passing through the multi-head attention, this set of learnable tokens are mapped to the semantic space via the MLP, ultimately resulting in tokens with fine-grained semantic information that can be embedded into the

text encoder—thereby completing the image-to-text inversion. The entire process is as follows:

$$\text{img}_{tokens} = \text{MLP} \left(\text{MHA} \left(T, (T, f_{\text{seq}}^{\text{img}}), (T, f_{\text{seq}}^{\text{img}}) \right) \right) \quad (5)$$

3.4. Prompt-driven Semantic Guidance

To effectively integrate textual semantics into the visual decoding process, we propose a Prompt-driven Semantic Guidance decoder, which adaptively modulates feature decoding based on language guidance. The decoder begins with a cross-attention mechanism that aligns the text prompting f^{text} with the high-level image features $f^{(5)}$ extracted from the final encoder stage. This initial interaction produces semantically guided features that establish a global alignment between the two modalities.

The decoder follows a five-stage cascade structure. At each stage i , the input decoding feature $f_d^{(i)} \in \mathbb{R}^{\frac{h_i}{2} \times \frac{w_i}{2} \times d}$ is first upsampled via bilinear interpolation to match the spatial resolution of the corresponding encoder feature $f^{(i)} \in \mathbb{R}^{h_i \times w_i \times d_i}$. The features are then fused through element-wise addition, with the encoder feature first projected via a 1×1 convolution for dimension alignment,

$$f_{\text{fused}}^{(i)} = \text{Up}(f_d^{(i)}) + \text{Conv}_{1 \times 1}(f^{(i)}). \quad (6)$$

Next, the fused feature $f_{\text{fused}}^{(i)}$ is refined using a Text-Guided Spatial Attention (TGSA) block, which incorporates the [eot] token of the text embedding $f_{\text{text}}^{[\text{eot}]} \in \mathbb{R}^{1 \times d_t}$. As shown in figure 2(b), both the visual and textual features are projected into a shared latent space of dimension d_c via separate MLPs,

$$\begin{aligned} v^{(i)} &= \text{MLP}_{\text{vis}}^i(f_{\text{fused}}^{(i)}) \in \mathbb{R}^{h_i \times w_i \times d_c} \\ t^{(i)} &= \text{MLP}_{\text{txt}}^i(f_{\text{text}}^{[\text{eot}]}) \in \mathbb{R}^{1 \times d_c} \end{aligned} \quad (7)$$

The text vector $t^{(i)}$ is then broadcast across the spatial dimensions of $v^{(i)}$, and we perform the dot-product operation to measure the semantic similarity at each spatial location. A softmax normalization yields the attention weight,

$$w^{(i)} = \text{Softmax} \left(\frac{\langle v^{(i)}, t^{(i)} \rangle}{\sqrt{c}} \right) \in \mathbb{R}^{h_i \times w_i \times 1} \quad (8)$$

Finally, the attention map $w^{(i)}$ acts as a residual gating mechanism, selectively amplifying regions likely to contain targets while suppressing background noise:

$$f_d^{(i-1)} = \text{Conv}_{3 \times 3}(f_{\text{fused}}^{(i)} + w^{(i)} \cdot f_{\text{fused}}^{(i)}) \quad (9)$$

This process is repeated at each decoding stage, progressively refining the representation under the guidance of semantic prompts. At the final stage, the decoder outputs $f_d^{(1)}$, which is passed through a shared convolutional prediction head to generate the binary target mask.

3.5. Optimization and Inference

Reconstruction Pretraining. Since the weights of the inversion net are difficult to learn effectively during formal training and the random initial weights of the encoder lead to slow training convergence, we introduce the method of reconstruction pretraining. Its core goal is to obtain meaningful weights for the inversion net and, at the same time, acquire effective initial weights for the encoder—laying a solid parameter foundation for formal training and accelerating model convergence.

During the pretraining phase, we keep the structures of the encoder and inversion net of the original model unchanged, only modifying the decoder part: replacing the original decoder with a reconstruction decoder composed of transposed convolution layers. The core function of this reconstruction process is to obtain effective initial weights of the network while learning to reconstruct the input infrared image. Here express the pretraining loss function:

$$\mathcal{L} = \mathcal{L}_{\text{Contra}} \left(\text{sg}[f^{(5)}], f_{\text{text}}^{[\text{eot}]} \right) + \mathcal{L}_{\text{MSE}}(\text{out}, \text{input}) \quad (10)$$

where $\text{sg}[\cdot]$ is the stop-gradient operation. \mathcal{L}_{MSE} , the Mean Squared Error loss, is employed to compute the loss between the reconstructed results and the input image, thereby constraining the learning process of the entire network. Particularly, we introduce a contrastive loss for optimizing the inversion network. We calculate $\mathcal{L}_{\text{Contra}}$ between the average image feature $f^{(5)}$ and the text embedding $f_{\text{text}}^{[\text{eot}]}$, which maximizes their similarity and encourages the inversion network to learn a semantically faithful vision-to-language mapping.

Normal Training. DGSPNet generates a predicted target mask \hat{m} from an input infrared image and is trained using the corresponding ground-truth mask m as supervision. Specifically, for each input image, the training objective minimizes the following loss function:

$$\mathcal{L} = \lambda_1 \mathcal{L}_{\text{BCE}}(m, \hat{m}) + \lambda_2 \mathcal{L}_{\text{SoftIoU}}(m, \hat{m}) \quad (11)$$

where, \mathcal{L}_{BCE} denotes the binary cross-entropy loss to penalize pixel-level probability errors, $\mathcal{L}_{\text{SoftIoU}}$ means SoftIoU loss to optimize the segmentation error, and the coefficient λ_i balances different loss components.

4. Experiments

4.1. Experimental Setup

Datasets. For our experiments, we conducted evaluations on all three datasets, NUA-SIRST, NUDT-SIRST, and IRSTD-1K. Following the methodology of previous works, specifically [35], the training and testing splits for both NUA-SIRST and NUDT-SIRST are set to 5:5, meaning 50% of the data was used for training and the remaining

Table 1. Comparison with existing IRSTD approaches on the IRSTD-1K, NUDT-SIRST, and NUAAS-SIRST datasets. The evaluation metrics are IoU (10^{-2}), P_d (10^{-2}), and F_a (10^{-6}). The first and second best results are highlighted in **bold** and underlined, respectively.

METHOD	BACKBONE	SIZE	IRSTD-1K			NUDT-SIRST			NUAAS-SIRST		
			IoU \uparrow	P_d \uparrow	F_a \downarrow	IoU \uparrow	P_d \uparrow	F_a \downarrow	IoU \uparrow	P_d \uparrow	F_a \downarrow
IPI [9]	Traditional	256	27.92	81.37	16.18	17.76	74.49	41.23	25.67	84.63	16.67
PSTNN [36]	Traditional	256	24.57	71.99	35.26	14.85	66.13	44.17	30.30	72.80	48.99
MSLSTIPT [26]	Traditional	256	11.43	79.03	1524	8.342	47.40	888.1	10.30	82.13	1131
RKformer [37]	Hybird	512	64.12	93.27	18.65	92.25	96.86	6.580	77.24	99.11	1.580
IAANet [28]	Hybird	256	66.25	93.15	14.20	90.22	97.26	8.320	74.22	93.53	22.70
MTU-Net [33]	Hybird	256	66.11	93.27	36.80	74.85	93.97	46.95	74.78	93.54	22.36
SCTransNet [35]	Hybird	256	68.03	93.27	10.74	94.09	98.62	4.290	77.50	96.95	<u>13.92</u>
ACM [3]	CNN	256	59.23	93.27	65.28	61.12	93.12	55.22	68.93	91.63	15.23
ALCNet [4]	CNN	256	60.60	92.98	58.80	64.74	94.18	34.61	70.83	94.30	36.15
ISNet [39]	CNN	256	61.85	90.24	31.56	81.24	97.78	6.340	70.49	95.06	67.98
FC3-Net [38]	CNN	256	65.07	91.54	15.55	78.56	93.86	23.92	72.44	98.14	10.85
DNA-Net [15]	CNN	256	65.90	90.91	<u>12.24</u>	88.19	98.83	9.000	75.80	95.82	8.780
UIU-Net [34]	CNN	256	66.15	<u>93.98</u>	22.07	93.48	98.31	7.790	76.91	95.82	14.13
ISTDU [12]	CNN	256	66.36	<u>93.60</u>	53.10	89.55	97.67	13.44	75.52	96.58	14.54
RDIAN [24]	CNN	256	56.45	88.55	26.63	76.28	95.77	34.56	68.72	93.54	43.29
AGPCNet [42]	CNN	256	66.29	92.83	13.12	88.87	97.20	10.02	75.69	96.48	14.99
Dim2Clear [40]	CNN	256	66.34	93.75	20.93	81.37	96.23	9.170	77.29	<u>99.10</u>	<u>6.720</u>
MMLNet [18]	CNN	256	67.21	94.28	14.00	81.81	98.43	11.77	<u>78.71</u>	98.88	25.71
L2SKNet [32]	CNN	256	67.81	90.24	17.46	93.58	97.57	5.330	73.43	98.17	20.82
Text-IRSTD [13]	CNN	256	<u>69.57</u>	92.59	14.97	<u>95.84</u>	99.73	<u>1.032</u>	-	-	-
DGSPNet (<i>Ours</i>)	CNN	256	70.87	93.26	14.23	96.13	<u>99.15</u>	0.32	80.32	97.33	9.26

Table 2. Ablation study on prompt-driven modules to evaluate the necessity and synergy of sub-components (TGCA, Cross-Attention, TGSA). Experiments are conducted on IRSTD-1K and NUDT-SIRST datasets, with evaluation metrics including IoU(10^{-2}), P_d (10^{-2}), and F_a (10^{-6}).

PROMPT-DRIVEN MODULE VARIANT				IRSTD-1K			NUDT-SIRST		
Baseline	TGCA	Cross-Attention	TGSA	IoU \uparrow	P_d \uparrow	F_a \downarrow	IoU \uparrow	P_d \uparrow	F_a \downarrow
✓				65.31	91.91	17.72	94.38	98.15	8.66
✓	✓			66.99	92.25	20.89	94.43	98.83	2.69
✓	✓	✓		67.54	92.92	26.93	94.66	98.73	2.64
✓	✓		✓	67.87	92.59	27.97	95.27	99.04	2.57
✓		✓	✓	67.81	91.91	36.00	94.81	98.94	1.49
✓	✓	✓	✓	70.87	93.26	14.23	96.13	99.15	0.32

50% for testing. For IRSTD-1K, we used an 8:2 split, where 80% of the data was used for training and 20% for testing. This setup aligns with the standard practice in this field, enabling a fair comparison with prior research.

Implementation Details. Our experiments are conducted on the NVIDIA GeForce RTX 4090 GPU. All images are normalized and random cropped to a size of 256×256. During pre-training phase, the AdamW optimizer is adopted with a learning rate of 3e-4 on the inversion network and 1e-4 on other parameters of the model over 800 pretraining epochs. After pretraining, the encoder and inversion net weights are directly transferred to the formal training

model, while the reconstruction decoder weights are discarded. As to formal Training Phase, the Adam optimizer is used and a Polynomial Decay learning rate scheduler is applied, with an initial learning rate of 1e-4 and a decay power of 1.2, ensuring the learning rate gradually decreases over 800 training epochs with a batch size of 4. Additionally, we choose frozen weights CLIP-ViT-B/32 as the text encoder. We evaluate the proposed model using three complementary metrics: Intersection over Union (IoU), detection probability (P_d), and false-alarm rate (F_a). More explanations of these metrics can be found in the supplementary.

Table 3. Ablation study on text designs to evaluate the influence of the number of fine-grained semantic tokens (#Tokens) on model performance, implemented on the IRSTD-1K dataset. Various text variants (incorporating different quantities of fine-grained tokens such as s_1^* , s_2^*) are assessed, with evaluation metrics including IoU(10^{-2}), P_d (10^{-2}), and F_a (10^{-6}).

#Tokens	TEXT VARIANT	IoU \uparrow	P_d \uparrow	F_a \downarrow
0	A photo of an infrared image, with targets in the background.	66.41	91.92	22.43
1	A photo of an infrared image, with s_1^* in the background.	69.29	92.25	17.78
2	A photo of an infrared image, with s_1^* in s_2^* background.	70.87	93.26	14.23
3	A photo of an infrared image, with s_1^* s_2^* in s_3^* background.	67.37	93.26	16.45
4	A photo of an infrared image, with s_1^* s_2^* in s_3^* s_4^* background.	69.22	91.58	19.91

4.2. Comparisons with the State-of-the-Arts

Our method demonstrates a clear advantage over state-of-the-art competitors. Table 1 presents a detailed comparison of the quantitative results between our proposed method and other state-of-the-art approaches in IRSTD across three datasets. Our approach achieves the best or second-best performance across most metrics, highlighting the effectiveness of our design. Specifically, on IRSTD-1K, DGSPNet obtains an IoU of 70.87 (surpassing the second-best method, Text-IRSTD [13], with 69.57), a high P_d of 93.26, and an F_a of 14.23, excelling in both detection accuracy and false alarm control. On NUDT-SIRST, it achieves a remarkable IoU of 96.13, a P_d of 99.15, and an ultra-low F_a of 0.32, significantly outperforming other methods and demonstrating precise target localization and robust background suppression. On NUAA-SIRST, DGSPNet also delivers competitive results with an IoU of 80.32, leading in detection integrity and accuracy.

Notably, although DGSPNet adopts a CNN backbone, it outperforms methods that utilize hybrid architectures. For example, on NUDT-SIRST, SCTransNet [35] achieves an IoU of 94.09 and an F_a of 4.290, whereas DGSPNet improves the IoU to 96.13 and reduces the F_a to 0.32, revealing a substantial performance gap. This demonstrates that our reconstruction pre-training strategy, hierarchical encoder design, and dual-granularity semantic prompt module effectively exploit the strengths of CNNs in feature extraction and semantic fusion, while further extending the capability of traditional visual models through language-guided learning.

4.3. Ablation study

Component Effectiveness in the Prompt-Driven Module. To clarify the necessity of each sub-component in the prompt-driven semantic guidance module, we conduct ablation studies on Text-Guide Channel Attention (TGCA), Cross-Attention, and Text-Guide Spatial Attention (TGSA). By incrementally adding these components, we compare IoU, detection rate (P_d), and false alarm rate (F_a) on the IRSTD-1K and NUDT-SIRST datasets, with the results shown in Table 2. On IRSTD-1K, both Cross-Attention and

TGSA improve IoU and P_d , indicating enhanced localization accuracy and completeness, but they also introduce a higher F_a . In contrast, TGCA can strengthen text–visual semantic alignment at the channel level, effectively reducing F_a and providing a basis for false-alarm suppression. Only when all three components are integrated does the model achieve the best trade-off among the metrics, yielding optimal overall performance on both IRSTD-1K and NUDT-SIRST. These results verify the functional complementarity of the components and the soundness of our designs.

Effectiveness of Dual-Granularity Semantic Prompt.

Table 3 presents the performance of different text designs for the dual-granularity semantic prompts on the IRSTD-1K dataset. The results show that using two semantic tokens and constructing the prompt as “A photo of an infrared image, with s_1^* in s_2^* background” yields the best performance. Using only a single semantic token leads to suboptimal results. This suggests that while one token is generally sufficient to describe either target characteristics or background attributes, introducing too many tokens can cause semantic ambiguity and degrade the guidance quality.

5. Conclusion

We introduce DGSPNet, an advanced end-to-end framework designed to combine coarse-grained textual priors with fine-grained, instance-adaptive prompts to provide explicit linguistic guidance for infrared small-target detection. Unlike traditional models, DGSPNet enhances detection by aligning hierarchical visual features with dynamically distilled semantic cues, enabling the decoder to focus attention on faint targets while effectively minimizing background noise. This approach leads to state-of-the-art performance across three challenging infrared datasets, demonstrating the model’s robust ability to handle diverse and complex scenarios. Our experiments highlight the synergy between fixed semantic anchors, which provide context, and learnable image tokens, which adapt to specific target features. This combination not only improves detection accuracy and robustness but also stabilizes training through a simple gradient-isolation trick. Additionally, DGSPNet shows the broader potential of cross-modal prompting in low-SNR

detection tasks, positioning language as a lightweight, plug-and-play supervisory signal that can significantly enhance computer vision models, especially in resource-constrained or noisy environments.

References

- [1] CL Philip Chen, Hong Li, Yantao Wei, Tian Xia, and Yuan Yan Tang. A local contrast method for small infrared target detection. *IEEE transactions on geoscience and remote sensing*, 52(1):574–581, 2013. 1
- [2] Zhiwei Chen, Bin Wang, and Hongjun Zhou. Local contrast method based on directional information for infrared small target detection. *Optics Communications*, 313:123–130, 2014. 2
- [3] Yimian Dai, Yiquan Wu, Fei Zhou, and Kobus Barnard. Asymmetric contextual modulation for infrared small target detection. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pages 950–959, 2021. 7
- [4] Yimian Dai, Yiquan Wu, Fei Zhou, and Kobus Barnard. Attentional local contrast networks for infrared small target detection. *IEEE transactions on geoscience and remote sensing*, 59(11):9813–9824, 2021. 7
- [5] Sameer D Deshpande, Meng Joo Er, Yu Chan, and Soon Huat Lee. Max-mean and max-median filters for detection of small targets. *IEEE Transactions on Aerospace and Electronic Systems*, 35(1):64–73, 1999. 2
- [6] Yuan Ding, Zhi Ye, and Yuanjing Tang. Irnet: A two-stage network for infrared small target detection. *IEEE Access*, 8: 52513–52522, 2020. 2
- [7] Richard O Duda, Peter E Hart, and David G Stork. *Pattern Classification*. John Wiley & Sons, 2012. 2
- [8] Linyu Fan, Yingying Wang, Guoliang Hu, Feifei Li, Yuhang Dong, Hui Zheng, Changqing Lin, Yue Huang, and Xinghao Ding. Diffusion-based continuous feature representation for infrared small-dim target detection. *IEEE Transactions on Geoscience and Remote Sensing*, 62:1–17, 2024. 1
- [9] Chenqiang Gao, Deyu Meng, Yi Yang, Yongtao Wang, Xiaofang Zhou, and Alexander G Hauptmann. Infrared patch-image model for small target detection in a single image. *IEEE transactions on image processing*, 22(12):4996–5009, 2013. 7
- [10] Jiuxiang Gu and et al. Open-vocabulary object detection via vision and language knowledge distillation. *ICLR*, 2022. 2
- [11] Shuai Han, Bin Wang, and Yu Zhang. Dgformer: Dual-granularity transformer for infrared small target detection. *IEEE Transactions on Geoscience and Remote Sensing*, 2023. 2
- [12] Qingyu Hou, Liuwei Zhang, Fanjiao Tan, Yuyang Xi, Hao-liang Zheng, and Na Li. Istdu-net: Infrared small-target detection u-net. *IEEE Geoscience and Remote Sensing Letters*, 19:1–5, 2022. 7
- [13] Feng Huang, Shuyuan Zheng, Zhaobing Qiu, Huanxian Liu, Huanxin Bai, and Liqiong Chen. Text-irstd: Leveraging semantic text to promote infrared small target detection in complex scenes. *arXiv preprint arXiv:2503.07249*, 2025. 1, 3, 5, 7, 8
- [14] Chao Jia, Yinfei Yang, Ye Xia, et al. Scaling up visual and vision-language representation learning with noisy text supervision. *ICML*, 2021. 2
- [15] Boyang Li, Chao Xiao, Longguang Wang, Yingqian Wang, Zaiping Lin, Miao Li, Wei An, and Yulan Guo. Dense nested attention network for infrared small target detection. *IEEE Transactions on Image Processing*, 32:1745–1758, 2022. 7
- [16] Junnan Li and et al. Grounded language-image pretraining. *CVPR*, 2022. 2
- [17] Junnan Li and et al. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. *ICLR*, 2023. 2
- [18] Qiang Li, Wei Zhang, Wanxuan Lu, and Qi Wang. Multi-branch mutual-guiding learning for infrared small target detection. *IEEE Transactions on Geoscience and Remote Sensing*, 2025. 7
- [19] Fanzhao Lin, Kexin Bao, Yong Li, Dan Zeng, and Shim-ing Ge. Learning contrast-enhanced shape-biased representations for infrared small target detection. *IEEE Transactions on Image Processing*, 33:3047–3058, 2024. 1
- [20] Jian Liu, Liang Liu, and Junlin Yang. Acunet: Adaptive context u-net for infrared small target detection. *Sensors*, 21 (13):4563, 2021. 2
- [21] Yuxin Liu and et al. Llmdet: Instruction-tuned vision-language models for zero-shot object detection. *arXiv preprint arXiv:2309.00676*, 2023. 2
- [22] Ziyi Liu and et al. Detgpt: Detect what you want via reasoning with large language models. *arXiv preprint arXiv:2305.15062*, 2023. 2
- [23] Alec Radford, Jong Wook Kim, Jack Hallacy, et al. Learning transferable visual models from natural language supervision. *ICML*, 2021. 2, 5
- [24] Heng Sun, Junxiang Bai, Fan Yang, and Xiangzhi Bai. Receptive-field and direction induced attention network for infrared dim small target detection with a large-scale dataset irdst. *IEEE Transactions on Geoscience and Remote Sensing*, 61:1–13, 2023. 7
- [25] Haotian Sun, Qiuyu Jin, Jun Xu, and Linbo Tang. Infrared small-target detection based on multi-level local contrast measure. *Procedia Computer Science*, 221:549–556, 2023. 1
- [26] Yang Sun, Jungang Yang, and Wei An. Infrared dim and small target detection via multiple subspace learning and spatial-temporal patch-tensor model. *IEEE Transactions on Geoscience and Remote Sensing*, 59(5):3737–3752, 2021. 7
- [27] Fang Wang, Lei Zhang, and Xiaodan Liang. Infrared small target detection via simultaneous orthogonal matching pursuit and sparsity measure. *Infrared Physics & Technology*, 69:81–89, 2015. 2
- [28] Kewei Wang, Shuaiyuan Du, Chengxin Liu, and Zhiguo Cao. Interior attention-aware network for infrared small target detection. *IEEE Transactions on Geoscience and Remote Sensing*, 60:1–13, 2022. 7
- [29] Peng Wang and et al. Ofa: Unifying architectures, tasks, and modalities through a simple sequence-to-sequence learning framework. *ICML*, 2022. 2

- [30] Yongtao Wang, Yifan Liu, and Liangpei Zhang. Isformer: Infrared small target detection via transformer-based feature enhancement. *Remote Sensing*, 14(18):4722, 2022. [2](#)
- [31] Zhe Wang and et al. Git: A generative image-to-text transformer for vision and language. *CVPR*, 2022. [2](#)
- [32] Fengyi Wu, Anran Liu, Tianfang Zhang, Luping Zhang, Junhai Luo, and Zhenming Peng. Saliency at the helm: Steering infrared small target detection with learnable kernels. *IEEE Transactions on Geoscience and Remote Sensing*, 2024. [7](#)
- [33] Tianhao Wu, Boyang Li, Yihang Luo, Yingqian Wang, Chao Xiao, Ting Liu, Jungang Yang, Wei An, and Yulan Guo. Mtnet: Multilevel transunet for space-based infrared tiny ship detection. *IEEE Transactions on Geoscience and Remote Sensing*, 61:1–15, 2023. [7](#)
- [34] Xin Wu, Danfeng Hong, and Jocelyn Chanussot. Uiu-net: U-net in u-net for infrared small object detection. *IEEE Transactions on Image Processing*, 32:364–376, 2022. [7](#)
- [35] Shuai Yuan, Hanlin Qin, Xiang Yan, Naveed Akhtar, and Ajmal Mian. Sctransnet: Spatial-channel cross transformer network for infrared small target detection. *IEEE Transactions on Geoscience and Remote Sensing*, 62:1–15, 2024. [1](#), [6](#), [7](#), [8](#)
- [36] Landan Zhang and Zhenming Peng. Infrared small target detection based on partial sum of the tensor nuclear norm. *Remote Sensing*, 11(4):382, 2019. [7](#)
- [37] Mingjin Zhang, Haichen Bai, Jing Zhang, Rui Zhang, Chaoyue Wang, Jie Guo, and Xinbo Gao. Rkformer: Runge-kutta transformer with random-connection attention for infrared small target detection. In *Proceedings of the 30th ACM International Conference on Multimedia*, pages 1730–1738, 2022. [7](#)
- [38] Mingjin Zhang, Ke Yue, Jing Zhang, Yunsong Li, and Xinbo Gao. Exploring feature compensation and cross-level correlation for infrared small target detection. In *Proceedings of the 30th ACM International Conference on Multimedia*, pages 1857–1865, 2022. [7](#)
- [39] Mingjin Zhang, Rui Zhang, Yuxiang Yang, Haichen Bai, Jing Zhang, and Jie Guo. Isnet: Shape matters for infrared small target detection. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 877–886, 2022. [1](#), [7](#)
- [40] Mingjin Zhang, Rui Zhang, Jing Zhang, Jie Guo, Yunsong Li, and Xinbo Gao. Dim2clear network for infrared small target detection. *IEEE Transactions on Geoscience and Remote Sensing*, 61:1–14, 2023. [7](#)
- [41] Shuang Zhang, Lei Zhang, and Xuanjing Mou. Infrared small target detection via non-local low-rank regularization and salient feature extraction. *Infrared Physics & Technology*, 77:364–377, 2016. [2](#)
- [42] Tianfang Zhang, Siying Cao, Tian Pu, and Zhenming Peng. Agpcnet: Attention-guided pyramid context networks for infrared small target detection. *arXiv preprint arXiv:2111.03580*, 2021. [7](#)
- [43] Yu Zhang, Xian Zhang, and Hao Sun. Infrared small target detection based on multiscale residual fusion network. In *ICIP*, pages 2508–2512, 2021. [2](#)
- [44] Liunian Harold Zhou and et al. Unified vision-language pre-training for image captioning and vqa. In *AAAI*, 2021. [2](#)