

---

# From Spatial to Spectral: An Efficient, Frequency-Guided Representation Learner for Small Object Detection

---

Anonymous Authors<sup>1</sup>

## Abstract

This document provides a basic paper template and submission guidelines. Abstracts must be a single paragraph, ideally between 4–6 sentences long. Gross violations will trigger corrections at the camera-ready phase.

## 1. Introduction

## 2. Related Work

We review prior work from three angles that are most relevant to our goal: (i) efficient detector architectures, (ii) small object detection strategies, and (iii) frequency-domain modeling for dense prediction.

### 2.1. Efficient Detector Architectures

Real-time detection has been driven by architectural efficiency in backbones, feature pyramids, and heads. One-stage YOLO-style detectors optimize the accuracy–latency trade-off through carefully designed blocks and multi-scale prediction, with recent variants continuing to improve both speed and accuracy (Wang et al., 2024; Khanam & Hussain, 2024). Lightweight enhancements for challenging regimes (e.g., cluttered scenes) often rely on stronger feature aggregation or multi-kernel perception to increase representational diversity while keeping inference efficient (Xiao et al., 2025).

In parallel, Transformer-based detectors seek end-to-end set prediction by removing hand-crafted components such as anchors (Carion et al., 2020). Subsequent work improves the practicality of DETR-like models via more efficient attention and training strategies, enabling competitive performance under constrained budgets (Zhao et al., 2024; Zhang et al., 2022). Despite these advances, both CNN and Transformer-based detectors still face a common ten-

sion for tiny/small objects: improving fine-detail sensitivity typically increases computation, memory, or architectural intrusion, making it difficult to deploy a uniformly effective solution across detector families.

### 2.2. Small Object Detection

Small objects are inherently information-limited: they occupy few pixels, induce weak feature responses, and are easily suppressed by downsampling and coarse fusion. Early two-stage and one-stage frameworks (e.g., Faster R-CNN and SSD) already revealed the difficulty of preserving small-object cues under feature hierarchy and stride growth (Ren et al., 2015; Liu et al., 2016). A large body of work improves small-object performance by strengthening multi-scale feature fusion (e.g., FPN and its variants) (Lin et al., 2017), introducing additional pyramid levels, and designing attention or alignment modules to enhance small-scale features.

Recent methods increasingly emphasize *detail-aware* feature enrichment. For example, HS-FPN highlights tiny objects by generating high-frequency responses as mask weights and complements this with explicit spatial dependency modeling (Shi et al., 2025). Context modeling (e.g., large receptive fields or multi-kernel designs) also helps disambiguate tiny objects from background clutter (Wang et al., 2025; Xiao et al., 2025). However, many of these approaches focus on either spatial fusion or receptive-field engineering, while the *mechanism of how fine details are suppressed and should be reconstructed* is often left implicit, and portability across heterogeneous detector designs is not always validated.

### 2.3. Frequency-Domain Modeling for Dense Prediction

Frequency-domain analysis offers a complementary lens to understand and manipulate representation learning. A line of work uses Fourier transforms to achieve efficient global interactions. GFNet replaces quadratic self-attention with frequency-domain filtering (FFT–filtering–IFFT), yielding log-linear complexity while maintaining global receptive fields (Rao et al., 2023). Other work links common architectural operations to spectral decomposition: FcaNet interprets channel attention as a frequency-domain compres-

<sup>1</sup>Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

055 sion process and generalizes global pooling to multi-spectral  
056 channel attention (?).

057 More recently, frequency-aware modules have been ex-  
058 plored for dense prediction. FDConv observes that can-  
059 didate dynamic convolution kernels often have highly similar  
060 frequency responses, and proposes constructing frequency-  
061 diverse weights by allocating parameters to disjoint Fourier  
062 indices, together with frequency-band/spatial modulation  
063 (Chen et al., 2025). Frequency-aware fusion is also studied:  
064 FreqFusion explicitly introduces adaptive low-pass/high-  
065 pass filtering to improve feature consistency and boundary  
066 sharpness during upsampling and fusion (?). Wavelet-based  
067 approaches provide multi-resolution decomposition with  
068 partial spatial localization; WTConv performs convolutions  
069 in wavelet sub-bands to scale receptive fields efficiently and  
070 can be used as a drop-in layer in CNNs (Finder et al., 2025).  
071

072 While these spectral methods demonstrate that frequency-  
073 domain techniques can be integrated into modern architec-  
074 tures, existing designs are often *task- or component-specific*  
075 (e.g., classification backbones, fusion-only modules, or spe-  
076 cific convolution families), and do not provide a unified,  
077 plug-and-play operator that can be instantiated across *back-*  
078 *bone, neck, and head* and generalize across both CNN-  
079 and Transformer-style detectors. Our work fills this gap by  
080 introducing a decomposition–reconstruction operator that  
081 preserves and re-synthesizes discriminative spectral compo-  
082 nents with minimal overhead, and systematically validating  
083 its cross-architecture generality.  
084

### 085 3. Method

#### 086 3.1. Overall Framework

087 Small objects occupy few pixels and are therefore domi-  
088 nated by high spatial frequencies (edges and fine textures).  
089 However, in modern detectors these components are pro-  
090 gressively attenuated: backbone downsampling behaves as an  
091 implicit low-pass operator, neck fusion/upsampling further  
092 smooths boundaries, and the head’s regression gradients  
093 weaken once boundary evidence is diluted. Our key obser-  
094 vation is that these three failure modes can be addressed by  
095 a single frequency-domain principle: *decompose features*  
096 *into low/high-frequency components, selectively enhance*  
097 *them, and reconstruct (or inject) the enhanced signal.*

098 **Unified decompose–enhance–reconstruct (DER) opera-  
099 tor.** Given a feature tensor  $\mathbf{X} \in \mathbb{R}^{C \times H \times W}$ , we define

$$\begin{aligned} 100 \quad & (\mathbf{X}_L, \mathbf{X}_H) = \mathcal{D}(\mathbf{X}), \\ 101 \quad & \mathbf{X}_L^+ = \mathcal{E}_L(\mathbf{X}_L), \quad \mathbf{X}_H^+ = \mathcal{E}_H(\mathbf{X}_H), \quad (1) \\ 102 \quad & \mathbf{X}^+ = \mathcal{R}(\mathbf{X}_L^+, \mathbf{X}_H^+), \end{aligned}$$

103 where  $\mathcal{D}$  extracts low-/high-frequency evidence (wavelet  
104 subbands or oriented band-pass responses),  $\mathcal{E}_L$  and  $\mathcal{E}_H$  are  
105

106 lightweight enhancement functions, and  $\mathcal{R}$  reconstructs or  
107 injects the enhanced signal back to the spatial feature stream.  
108

109 This DER operator is instantiated at three locations with  
110 complementary roles along the detector pipeline and forms  
111 a coherent frequency flow: WDG performs DER early so  
112 that boundary-relevant  $\mathbf{X}_H$  is not prematurely suppressed  
113 and can be propagated to the neck; LGE/LGE-W applies  
114 DER again *before* multi-scale fusion so that the enhanced  
115 high-frequency residual survives aggregation; FDHead fi-  
116 nally converts the resulting high-frequency energy into a  
117 regression gain, making the localization head sensitive to  
118 boundary evidence that has been preserved and re-amplified  
119 upstream.

120 **DER instantiations across backbone/neck/head.** In our  
121 implementation, the three modules correspond to Eq. (1)  
122 with different choices of  $\mathcal{D}$ ,  $\mathcal{E}_{L/H}$ , and  $\mathcal{R}$ :

$$\begin{aligned} 123 \quad & \text{WDG: } \mathcal{D} = \text{DWT}, \quad \mathcal{E}_L = f_{\text{cdc}}, \quad \mathcal{E}_H = \text{Id}, \quad \mathcal{R} = \text{IDWT}; \\ 124 \quad & \text{LGE/LGE-W: } \mathcal{D} = \text{LGF} \text{ (or DWT)}, \quad \mathcal{E}_H : (13), \quad \mathcal{R} : (14), \text{ (or (15));} \\ 125 \quad & \text{FDHead: } \mathcal{D} = \text{DWT}, \quad \mathcal{R} : (18) \rightarrow (19). \end{aligned} \quad (2)$$

126 This shared DER view clarifies why the three changes are  
127 not independent: they progressively preserve, re-inject, and  
128 exploit the same high-frequency evidence across backbone–  
129 neck–head.

130 Formally, given a baseline detector with backbone  $\mathcal{B}$ , neck  
131  $\mathcal{N}$ , and head  $\mathcal{H}$ , we instantiate the three operators as

$$\begin{aligned} 132 \quad & \{\mathbf{C}_\ell\} = \mathcal{B}(\mathbf{I}), \quad \mathbf{C}'_\ell = \begin{cases} \mathcal{W}(\mathbf{C}_\ell), & \ell \in \mathcal{S}_{\mathcal{B}}, \\ \mathbf{C}_\ell, & \text{otherwise,} \end{cases} \\ 133 \quad & \{\mathbf{P}_\ell\} = \mathcal{N}(\{\mathbf{C}'_\ell\}), \quad \mathbf{P}'_\ell = \mathcal{E}(\mathbf{P}_\ell), \\ 134 \quad & \widehat{\mathbf{Y}} = \mathcal{H}_{\text{FD}}(\{\mathbf{P}'_\ell\}). \end{aligned} \quad (3)$$

135 Here  $\mathcal{W}$ ,  $\mathcal{E}$ , and  $\mathcal{H}_{\text{FD}}$  are the concrete DER instantiations  
136 for backbone, neck, and head, respectively.  $\mathcal{S}_{\mathcal{B}}$  denotes the  
137 set of backbone stages where WDG is inserted.

#### 138 3.2. Wavelet-Difference Gate (WDG)

139 We introduce Wavelet-Difference Gate (WDG), a  
140 lightweight plug-and-play bottleneck that injects frequency-  
141 aware modulation into convolutional backbones. Given an  
142 input feature map  $\mathbf{x} \in \mathbb{R}^{C \times H \times W}$ , WDG first applies a  $1 \times 1$   
143 projection to hidden channels  $C' = \lfloor eC \rfloor$  (with expansion  
144 ratio  $e$ ) and then performs a 2D Haar discrete wavelet  
145 transform (DWT) to separate low- and high-frequency  
146 components. For simplicity, we describe the transform for  
147 even  $H, W$ ; in practice we align sizes by cropping/padding  
148 and restore the original resolution after reconstruction.

149 **Projection and wavelet decomposition.** We first project  $\mathbf{x}$

110 to a hidden space and decompose it into Haar subbands:

$$\begin{aligned} \mathbf{x}' &= f_{1 \times 1}(\mathbf{x}), \\ (\mathbf{x}_{LL}, \mathbf{x}_{LH}, \mathbf{x}_{HL}, \mathbf{x}_{HH}) &= \text{DWT}(\mathbf{x}'). \end{aligned} \quad (4)$$

115 Here  $\mathbf{x}_{LL}$  is the low-frequency approximation, and  
 116  $\{\mathbf{x}_{LH}, \mathbf{x}_{HL}, \mathbf{x}_{HH}\}$  capture horizontal/vertical/diagonal  
 117 high-frequency details. This decomposition explicitly sepa-  
 118 rates coarse structures from fine details, enabling targeted  
 119 refinement for small objects.

120 For Haar DWT/IDWT, each spatial  $2 \times 2$  block is trans-  
 121 formed by a  $2 \times 2$  Haar matrix. For each channel  $c$  and  
 122 location  $(u, v)$ , define the local block  
 123

$$\mathbf{X}_{u,v}^{(c)} = \begin{pmatrix} \mathbf{x}_{2u,2v}^{'(c)} & \mathbf{x}_{2u,2v+1}^{'(c)} \\ \mathbf{x}_{2u+1,2v}^{'(c)} & \mathbf{x}_{2u+1,2v+1}^{'(c)} \end{pmatrix}. \quad (5)$$

128 Then Haar DWT and IDWT are given by

$$\begin{aligned} \mathbf{S}_{u,v}^{(c)} &= \frac{1}{2} \mathbf{H}_2 \mathbf{X}_{u,v}^{(c)} \mathbf{H}_2^\top, \\ \mathbf{X}_{u,v}^{(c)} &= \frac{1}{2} \mathbf{H}_2^\top \mathbf{S}_{u,v}^{(c)} \mathbf{H}_2, \quad \mathbf{H}_2 = \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix}, \end{aligned} \quad (6)$$

134 where  $\mathbf{S}_{u,v}^{(c)} = \begin{pmatrix} \mathbf{x}_{LL,u,v}^{(c)} & \mathbf{x}_{LH,u,v}^{(c)} \\ \mathbf{x}_{HL,u,v}^{(c)} & \mathbf{x}_{HH,u,v}^{(c)} \end{pmatrix}$  collects the four sub-  
 135 bands. This matrix form is exactly equivalent to the element-  
 136 wise expressions used in our implementation.  
 137

138 **RepCDC for low-frequency refinement.** To enhance dis-  
 139 criminative edges while keeping computation low, we refine  
 140 the approximation subband at half resolution:  
 141

$$\mathbf{y}_{LL} = f_{\text{cdc}}(\mathbf{x}_{LL}). \quad (7)$$

144 In our implementation,  $f_{\text{cdc}}$  is RepCDC followed by nor-  
 145 malization and activation. RepCDC parameterizes a central-  
 146 difference convolution by decreasing the center coefficient  
 147 of a  $3 \times 3$  kernel with a learnable  $\theta$ . Concretely, the effective  
 148 kernel is  
 149

$$\mathbf{y}_{p,q}^{(o)} = \sum_c \sum_{i=-1}^1 \sum_{j=-1}^1 \mathbf{W}_{i,j}^{(o,c)} \mathbf{z}_{p+i,q+j}^{(c)} - \sum_c \theta^{(o,c)} \mathbf{z}_{p,q}^{(c)}, \quad (8)$$

150 where  $\mathbf{z}$  denotes the input to RepCDC (e.g.,  $\mathbf{z} = \mathbf{x}_{LL}$ ), and  
 151  $(p, q)$  indexes spatial locations. This expression is exactly  
 152 equivalent to subtracting  $\theta$  from the center coefficient of a  
 153  $3 \times 3$  kernel. During deployment, the resulting kernel is  
 154 fused into a single standard convolution, so RepCDC incurs  
 155 no extra inference branches. Operating on  $\mathbf{x}_{LL}$  reduces spa-  
 156 tial cost by  $4 \times$  while strengthening edge sensitivity through  
 157 the difference term.  
 158

159 **High-frequency gated modulation.** We use high-frequency  
 160 responses to predict a content-adaptive gate and modulate  
 161

the refined low-frequency feature:

$$\begin{aligned} \mathbf{g} &= \sigma(f_g(\text{Concat}(\mathbf{x}_{LH}, \mathbf{x}_{HL}, \mathbf{x}_{HH}))), \\ \tilde{\mathbf{x}}_{LL} &= \mathbf{y}_{LL} \odot (\mathbf{1} + \mathbf{g}). \end{aligned} \quad (9)$$

We use additive gating ( $\mathbf{1} + \mathbf{g}$ ) to emphasize informative regions without suppressing the overall magnitude of  $\mathbf{y}_{LL}$ .  $f_g$  is a  $1 \times 1$  convolution followed by normalization, and  $\text{Concat}(\cdot)$  denotes channel-wise concatenation. Since the gate is predicted from high-frequency subbands, it acts as a detail-aware selector that boosts regions with strong edge/texture cues.

**Reconstruction and residual output.** Finally, we preserve the original high-frequency subbands and reconstruct the feature via inverse Haar transform:

$$\begin{aligned} \hat{\mathbf{x}}' &= \text{IDWT}(\tilde{\mathbf{x}}_{LL}, \mathbf{x}_{LH}, \mathbf{x}_{HL}, \mathbf{x}_{HH}), \\ \mathbf{y} &= f_{1 \times 1}^{\text{out}}(\hat{\mathbf{x}}'). \end{aligned} \quad (10)$$

When input/output channels match, WDG uses a residual connection  $\mathbf{y} \leftarrow \mathbf{x} + \mathbf{y}$ . Since the wavelet-domain refinement operates on  $H/2 \times W/2$ , WDG adds only a small overhead and can be inserted as a generic bottleneck into different backbone designs. Preserving the original high-frequency subbands avoids over-smoothing and helps retain boundary sharpness after reconstruction.

### 3.3. Log-Gabor Enhancer (LGE) and WTConv Variant (LGE-W)

We next improve the neck by introducing Log-Gabor Enhancer (LGE), a plug-and-play high-frequency refinement module applied to intermediate feature maps before multi-scale fusion. LGE is instantiated per feature level and is agnostic to the specific fusion topology (e.g., FPN/PAN/decoder-style aggregation).

**Log-Gabor filter bank (LGF).** Given a feature map  $\mathbf{x} \in \mathbb{R}^{C \times H \times W}$ , LGF applies a fixed Log-Gabor filter bank using depthwise convolutions. Let  $K$  and  $S$  denote the number of orientations and scales. For each channel  $c$ , orientation  $k$ , and scale  $s$ , we compute

$$\mathbf{h}_{s,k}^{(c)} = \mathbf{x}^{(c)} * \mathbf{g}_{s,k}, \quad (11)$$

where  $\mathbf{g}_{s,k}$  is a non-learnable Log-Gabor kernel and  $*$  is convolution. In our implementation,  $\mathbf{g}_{s,k}$  is instantiated in the spatial domain by rotating a centered coordinate grid and applying a log-normal radial envelope with a cosine

165 angular term:@@

$$\begin{aligned} c_k &= \cos \phi_k, & s_k &= \sin \phi_k, \\ u' &= u c_k + v s_k, & v' &= -u s_k + v c_k, \\ r &= \sqrt{u'^2 + v'^2} + \varepsilon, & \theta &= \text{atan2}(v', u'), \\ \mathbf{g}_{s,k}(u, v) &= \exp\left(-\frac{\log^2(r/\rho_s)}{2 \log^2 2}\right) \cos \theta. \end{aligned} \quad (12)$$

where  $\phi_k = k\pi/K$  and  $\rho_s$  is a fixed scale parameter. This produces a set of directional subband responses that explicitly emphasize edges and fine textures while introducing no additional learnable filter parameters.

**Learnable aggregation and residual enhancement (LGE).** LGE aggregates the subbands with learnable orientation/scale importance. Let  $\alpha \in \mathbb{R}^S$  and  $\beta \in \mathbb{R}^K$  be learnable logits; we obtain normalized weights by softmax and compute the high-frequency summary

$$\mathbf{h}^{(c)} = \sum_{s=1}^S \sum_{k=1}^K \text{softmax}(\alpha)_s \text{softmax}(\beta)_k \mathbf{h}_{s,k}^{(c)}. \quad (13)$$

We further apply a learnable global scale  $\gamma$  (implemented as a sigmoid-gated parameter) and a local mixing operator  $f_{\text{mix}}$ :

$$\mathbf{y} = \mathbf{x}_{\text{skip}} + f_{\text{mix}}(\sigma(\gamma) \mathbf{h}). \quad (14)$$

Here  $\mathbf{x}_{\text{skip}}$  is either the identity mapping (when channels match) or a  $1 \times 1$  projection. In our implementation,  $f_{\text{mix}}$  is a  $3 \times 3$  convolution (depthwise when  $C$  is preserved), so LGE adds only local mixing on top of fixed spectral decomposition while keeping a residual pathway.

**Wavelet variant (LGE-W).** LGE-W follows Eq. (11)–(14) but replaces  $f_{\text{mix}}$  with a wavelet-transform convolution (WTConv) when  $C$  is preserved. Using a fixed wavelet (Haar/db1), WTConv performs subband mixing in the wavelet domain and adds a lightweight depthwise branch:@@

$$\text{WTConv}(\mathbf{z}) = \mathcal{S}_0 \mathcal{D}_0(\mathbf{z}) + \text{IDWT}(\mathcal{S} \mathcal{D}_4(\text{DWT}(\mathbf{z}))), \quad (15)$$

where  $\mathcal{D}_0$  is a depthwise convolution in the spatial domain and  $\mathcal{D}_4$  denotes grouped depthwise convolutions applied over the four wavelet subbands. @@

### 3.4. Frequency-Driven Head (FDHead)

We finally introduce Frequency-Driven Head (FDHead), a frequency-aware detection head that improves small-object localization by injecting a boundary-sensitive prior into dense regression while preserving the standard anchor-free interface. FDHead is instantiated over multi-scale feature maps  $\{\mathbf{x}_i\}_{i=1}^N$  (e.g.,  $P2$ – $P5$ ) and shares most head parameters across levels to reduce capacity fragmentation.

**Shared prediction tower.** For each level  $i$ , FDHead first aligns channels to a hidden width  $C_h$  (Conv+GroupNorm) and then applies a shared refinement stack (DEConv + depthwise–pointwise mixing). The DEConv block aggregates multiple directional-difference operators (center/adjacent/horizontal/vertical) and a standard kernel; at inference it can be written as a single convolution with merged weights:

$$\text{DEConv}(\mathbf{u}) = \varphi\left((\sum_m \mathbf{K}_m) * \mathbf{u} + \sum_m \mathbf{b}_m\right), \quad (16)$$

where  $m$  indexes the directional branches and  $\varphi(\cdot)$  denotes normalization and activation. This biases the shared tower toward contour-aware features that are beneficial for boundary-aligned regression.

$$\mathbf{f}_i = \mathcal{T}(\mathbf{x}_i), \quad \mathcal{T} = \mathcal{T}_{\text{share}} \circ \mathcal{T}_{1 \times 1}. \quad (17)$$

**P2 high-frequency gate.** Since the finest level ( $P2$ ) carries the most precise spatial details, FDHead applies a lightweight wavelet gate only on  $i = 1$  (corresponding to  $P2$ ). Let  $C_f$  be the gated channel width (set as a fraction of  $C_h$ ); we split channels  $\mathbf{f}_1 = [\mathbf{f}_a, \mathbf{f}_b]$  with  $\mathbf{f}_a \in \mathbb{R}^{C_f \times H \times W}$ . Using a fixed Haar transform, we estimate high-frequency energy as a softmax-weighted mixture of subband magnitudes and convert it to a channel-wise gain:

$$\begin{aligned} (\mathbf{f}_{LL}, \mathbf{f}_{LH}, \mathbf{f}_{HL}, \mathbf{f}_{HH}) &= \text{DWT}(\mathbf{f}_a), \\ \mathbf{w} &= \text{softmax}(\boldsymbol{\omega}), \\ \mathbf{h} &= w_{LH} |\mathbf{f}_{LH}| + w_{HL} |\mathbf{f}_{HL}| + w_{HH} |\mathbf{f}_{HH}|, \\ \mathbf{g} &= \text{Gate}(\text{AvgPool}(\mathbf{h})), \\ \tilde{\mathbf{f}}_a &= \mathbf{f}_a \odot (1 + \alpha \mathbf{g}). \end{aligned} \quad (18)$$

Here  $\boldsymbol{\omega}$  are learnable logits over  $\{LH, HL, HH\}$  and  $\alpha$  controls the gate strength. Gate( $\cdot$ ) is a squeeze-excitation style channel MLP (two  $1 \times 1$  convs with sigmoid output) driven by pooled high-frequency energy. We then form  $\tilde{\mathbf{f}}_1 = [\tilde{\mathbf{f}}_a, \mathbf{f}_b]$  and apply it only to the box branch: high-frequency energy is a direct proxy for boundary sharpness and thus improves offset estimation, while leaving the classification stream unchanged avoids over-fitting to textures and background clutter. For the remaining levels  $i > 1$ , we set  $\tilde{\mathbf{f}}_i = \mathbf{f}_i$ .

**Box/class prediction and decoding.** FDHead predicts per-location class logits and distributional box offsets (DFL) as

$$\mathbf{b}_i = \text{Scale}_i(\mathcal{H}_{\text{box}}(\tilde{\mathbf{f}}_i)), \quad \mathbf{p}_i = \mathcal{H}_{\text{cls}}(\mathbf{f}_i), \quad (19)$$

and decodes boxes by  $\widehat{\mathbf{B}} = \text{dist2bbox}(\text{DFL}(\mathbf{b}), \mathbf{A}) \cdot \mathbf{s}$  with anchors  $\mathbf{A}$  and strides  $\mathbf{s}$ . This design targets small objects by frequency-gating only the finest level while keeping the remaining head computation shared and lightweight.

## 220 4. Experiment

### 221 4.1. Datasets and Metrics

222 We evaluate our framework on four benchmarks to demonstrate its robustness and cross-domain generalization: Vis-  
 223 Drone2019 (Du et al., 2019), TinyPerson (Yu et al., 2020),  
 224 UAVDT (Du et al., 2018), and DOTA v1 (Xia et al., 2018).  
 225 **VisDrone2019** is our primary benchmark and is particularly  
 226 challenging due to dense small objects and severe scale variation,  
 227 where most targets are smaller than  $50 \times 50$  pixels.  
 228

229 We report both accuracy and efficiency, including mAP<sub>50</sub>,  
 230 the number of parameters, GFLOPs, model size, and FPS.  
 231

### 232 4.2. Configuration

233 The experimental configuration is detailed in Table 1.

234 **Table 1. Configuration of Training and Testing Experiment**  
 235 **Environments.** Detailed hardware and software configuration  
 236 used for all experiments in this study.

237 Environment	238 Parameter
CPU	Intel(R) Xeon(R) Gold 5218R CPU @ 2.10GHz
GPU	NVIDIA A100-PCIE-40GB
VRAM	40 GB
RAM	46 GB
Operating System	Rocky Linux 8.5 (Green Obsidian)
Language	Python 3.10.14
Frame	PyTorch 2.1.0
CUDA Version	12.6

251 For YOLO-style architectures, models are trained for 300  
 252 epochs with an input resolution of  $640 \times 640$  and batch size  
 253 16, using SGD optimization. Unless otherwise specified,  
 254 Mosaic augmentation is enabled throughout training; we  
 255 use 4 dataloader workers and disable AMP.  
 256

## 257 5. Main Results

### 258 5.1. Alation

## 259 6. Analyses and Discussion

## 260 7. Conclusion

## 261 References

262 Carion, N., Massa, F., Synnaeve, G., Usunier, N., Kirillov,  
 263 A., and Zagoruyko, S. End-to-end object detection with  
 264 transformers. *arXiv preprint arXiv:2005.12872*, 2020.

265 Chen, L., Gu, L., Li, L., Yan, C., and Fu, Y. Frequency  
 266 dynamic convolution for dense image prediction. *arXiv*  
 267 preprint *arXiv:2503.18783*, 2025.

268 Du, D., Qi, Y., Yu, H., Yang, Y., Duan, K., Li, G., Zhang,  
 269

270 W., Huang, Q., and Tian, Q. The unmanned aerial vehicle  
 271 benchmark: Object detection and tracking. In *Proceedings*  
 272 of the European Conference on Computer Vision (ECCV), pp.  
 273 370–386, 2018.

274 Du, D., Zhu, P., Wen, L., Bian, X., Lin, H., Hu, Q., Peng, T.,  
 275 Zheng, J., Wang, X., Zhang, Y., et al. Visdrone-det2019:  
 276 The vision meets drone object detection in image chal-  
 277 lenge results. In *Proceedings of the IEEE/CVF Interna-*  
 278 *tional Conference on Computer Vision Workshops*, pp.  
 279 0–0, 2019.

280 Finder, S. E., Amoyal, R., Treister, E., and Freifeld, O.  
 281 Wavelet convolutions for large receptive fields. In  
 282 Leonardis, A., Ricci, E., Roth, S., Russakovsky, O., Sat-  
 283 tterer, T., and Varol, G. (eds.), *Computer Vision – ECCV*  
 284 2024, pp. 363–380, Cham, 2025. Springer Nature Switzer-  
 285 land. ISBN 978-3-031-72949-2.

286 Khanam, R. and Hussain, M. YOLOv11: An overview  
 287 of the key architectural enhancements. *arXiv preprint*  
 288 *arXiv:2410.17725*, 2024.

289 Lin, T.-Y., Dollár, P., Girshick, R., He, K., Hariharan, B.,  
 290 and Belongie, S. Feature pyramid networks for object  
 291 detection. In *Proceedings of the IEEE Conference on*  
 292 *Computer Vision and Pattern Recognition (CVPR)*, pp.  
 293 2117–2125, 2017.

294 Liu, W., Anguelov, D., Erhan, D., Szegedy, C., Reed, S.,  
 295 Fu, C.-Y., and Berg, A. C. SSD: Single shot multibox  
 296 detector. In *Computer Vision–ECCV 2016: 14th Euro-  
 297 pean Conference, Amsterdam, The Netherlands, October  
 298 11–14, 2016, Proceedings, Part I*, pp. 21–37. Springer,  
 299 2016.

300 Rao, Y., Zhao, W., Zhu, Z., Zhou, J., and Lu, J. GFNet:  
 301 Global filter networks for visual recognition. *IEEE Trans-  
 302 actions on Pattern Analysis and Machine Intelligence  
 303 (TPAMI)*, 45(9):10960–10973, September 2023. doi:  
 304 10.1109/TPAMI.2023.3263824.

305 Ren, S., He, K., Girshick, R., and Sun, J. Faster R-CNN:  
 306 Towards real-time object detection with region proposal  
 307 networks. In *Advances in Neural Information Processing  
 308 Systems (NeurIPS)*, volume 28, 2015.

309 Shi, Z., Hu, J., Ren, J., Ye, H., Yuan, X., Ouyang, Y., He, J.,  
 310 Ji, B., and Guo, J. HS-FPN: High frequency and spatial  
 311 perception fpn for tiny object detection. *arXiv preprint*  
 312 *arXiv:2412.10116*, 2025.

313 Wang, A., Chen, H., Liu, L., Chen, K., Lin, Z., Han, J.,  
 314 and Ding, G. YOLOv10: Real-time end-to-end object  
 315 detection. *arXiv preprint arXiv:2405.14458*, 2024.

- 275 Wang, A., Chen, H., Lin, Z., Han, J., and Ding, G. LSNet:  
276 See large, focus small. In *Proceedings of the IEEE/CVF*  
277 *Conference on Computer Vision and Pattern Recognition*  
278 (*CVPR*), 2025.
- 279 Xia, G.-S., Bai, X., Ding, J., Zhu, Z., Belongie, S., Luo,  
280 J., Datcu, M., Pelillo, M., and Zhang, L. Dota: A large-  
281 scale dataset for object detection in aerial images. In  
282 *Proceedings of the IEEE conference on computer vision*  
283 *and pattern recognition*, pp. 3974–3983, 2018.
- 285 Xiao, Y., Xu, T., Xin, Y., and Li, J. FBRT-YOLO: Faster and  
286 better for real-time aerial image detection. *arXiv preprint*  
287 *arXiv:2504.20670*, 2025.
- 289 Yu, X., Gong, Y., Jiang, N., Ye, Q., and Han, Z. Scale  
290 match for tiny person detection. In *Proceedings of the*  
291 *IEEE/CVF Winter Conference on Applications of Com-*  
292 *puter Vision*, pp. 1257–1265, 2020.
- 293 Zhang, H., Li, F., Liu, S., Zhang, L., Su, H., Zhu, J., Ni,  
294 L. M., and Shum, H.-Y. DINO: DETR with improved  
295 denoising anchor boxes for end-to-end object detection.  
296 *arXiv preprint arXiv:2203.03605*, 2022.
- 298 Zhao, Y., Lv, W., Xu, S., Wei, J., Wang, G., Dang, Q., Liu,  
299 Y., and Chen, J. Detrs beat yolos on real-time object  
300 detection. *arXiv preprint arXiv:2304.08069*, 2024.
- 301  
302  
303  
304  
305  
306  
307  
308  
309  
310  
311  
312  
313  
314  
315  
316  
317  
318  
319  
320  
321  
322  
323  
324  
325  
326  
327  
328  
329

330   **A. You *can* have an appendix here.**

331   You can have as much text here as you want. The main body must be at most 8 pages long. For the final version, one more  
332   page can be added. If you want, you can use an appendix like this one.  
333

334   The `\onecolumn` command above can be kept in place if you prefer a one-column appendix, or can be removed if you  
335   prefer a two-column appendix. Apart from this possible change, the style (font size, spacing, margins, page numbering, etc.)  
336   should be kept the same as the main body.  
337

338

339

340

341

342

343

344

345

346

347

348

349

350

351

352

353

354

355

356

357

358

359

360

361

362

363

364

365

366

367

368

369

370

371

372

373

374

375

376

377

378

379

380

381

382

383

384