

---

# 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. FDCConv 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.

### 085 3. Method

#### 086 3.1. WDG

087 We introduce Wavelet-Difference Gate (WDG), a  
088 lightweight plug-and-play bottleneck that injects frequency-  
089 aware modulation into convolutional backbones. Given an  
090 input feature map  $\mathbf{x} \in \mathbb{R}^{C \times H \times W}$ , WDG first applies a  $1 \times 1$   
091 projection to hidden channels  $C' = \lfloor eC \rfloor$  (with expansion  
092 ratio  $e$ ) and then performs a 2D Haar discrete wavelet  
093 transform (DWT) to separate low- and high-frequency  
094 components. For simplicity, we describe the transform for  
095 even  $H, W$ ; in practice we align sizes by cropping/padding  
096 and restore the original resolution after reconstruction.  
097

098 **Projection and wavelet decomposition.** We first project  $\mathbf{x}$   
099 to a hidden space and decompose it into Haar subbands:

$$100 \quad \mathbf{x}' = f_{1 \times 1}(\mathbf{x}), \quad (1)$$

$$101 \quad (\mathbf{x}_{LL}, \mathbf{x}_{LH}, \mathbf{x}_{HL}, \mathbf{x}_{HH}) = \text{DWT}(\mathbf{x}').$$

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

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

$$110 \quad \mathbf{X}_{u,v}^{(c)} = \begin{pmatrix} \mathbf{x}_{2u,2v}^{'(c)} & \mathbf{x}_{2u,2v+1}^{'(c)} \\ \mathbf{x}_{2u+1,2v}^{'(c)} & \mathbf{x}_{2u+1,2v+1}^{'(c)} \end{pmatrix}. \quad (2)$$

111 Then Haar DWT and IDWT are given by

$$112 \quad \mathbf{S}_{u,v}^{(c)} = \frac{1}{2} \mathbf{H}_2 \mathbf{X}_{u,v}^{(c)} \mathbf{H}_2^\top, \\ 113 \quad \mathbf{X}_{u,v}^{(c)} = \frac{1}{2} \mathbf{H}_2^\top \mathbf{S}_{u,v}^{(c)} \mathbf{H}_2, \quad \mathbf{H}_2 = \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix}, \quad (3)$$

114 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-  
115 bands. This matrix form is exactly equivalent to the element-  
116 wise expressions used in our implementation.

117 **RepCDC for low-frequency refinement.** To enhance dis-  
118 criminative edges while keeping computation low, we refine  
119 the approximation subband at half resolution:

$$120 \quad \mathbf{y}_{LL} = f_{\text{cdc}}(\mathbf{x}_{LL}). \quad (4)$$

121 In our implementation,  $f_{\text{cdc}}$  is RepCDC followed by nor-  
122 malization and activation. RepCDC parameterizes a central-  
123 difference convolution by decreasing the center coefficient  
124 of a  $3 \times 3$  kernel with a learnable  $\theta$ . Concretely, the effective  
125 kernel is

$$126 \quad \mathbf{y}_{p,q}^{(o)} = \sum_c \sum_{i=-1}^1 \sum_{j=-1}^1 \mathbf{W}_{i,j}^{(o,c)} \mathbf{z}_{p+i,q+j}^{(c)} - \sum_c \boldsymbol{\theta}^{(o,c)} \mathbf{z}_{p,q}^{(c)}, \quad (5)$$

127 where  $\mathbf{z}$  denotes the input to RepCDC (e.g.,  $\mathbf{z} = \mathbf{x}_{LL}$ ), and  
128  $(p, q)$  indexes spatial locations. This expression is exactly  
129 equivalent to subtracting  $\theta$  from the center coefficient of a  
130  $3 \times 3$  kernel. During deployment, the resulting kernel is  
131 fused into a single standard convolution, so RepCDC incurs  
132 no extra inference branches. Operating on  $\mathbf{x}_{LL}$  reduces spa-  
133 tial cost by  $4 \times$  while strengthening edge sensitivity through  
134 the difference term.

135 **High-frequency gated modulation.** We use high-frequency  
136 responses to predict a content-adaptive gate and modulate  
137 the refined low-frequency feature:

$$138 \quad \mathbf{g} = \sigma(f_g(\text{Concat}(\mathbf{x}_{LH}, \mathbf{x}_{HL}, \mathbf{x}_{HH}))), \quad (6)$$

$$139 \quad \tilde{\mathbf{x}}_{LL} = \mathbf{y}_{LL} \odot (1 + \mathbf{g}).$$

140 We use additive gating  $(1 + \mathbf{g})$  to emphasize informative  
141 regions without suppressing the overall magnitude of  $\mathbf{y}_{LL}$ .  
142  $f_g$  is a  $1 \times 1$  convolution followed by normalization, and  
143  $\text{Concat}(\cdot)$  denotes channel-wise concatenation. Since the  
144 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:

$$\begin{aligned}\hat{\mathbf{x}}' &= \text{IDWT}(\tilde{\mathbf{x}}_{LL}, \mathbf{x}_{LH}, \mathbf{x}_{HL}, \mathbf{x}_{HH}), \\ \mathbf{y} &= f_{1 \times 1}^{\text{out}}(\hat{\mathbf{x}}').\end{aligned}\quad (7)$$

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.

### 3.2. LGE and 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

$$\mathbf{h}_{s,k}^{(c)} = \mathbf{x}^{(c)} * \mathbf{g}_{s,k}, \quad (8)$$

where  $\mathbf{g}_{s,k}$  is a non-learnable Log-Gabor kernel and  $*$  is convolution. 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  $\boldsymbol{\alpha} \in \mathbb{R}^S$  and  $\boldsymbol{\beta} \in \mathbb{R}^K$  be learnable logits; we obtain normalized weights by softmax and compute the high-frequency summary

$$\mathbf{h}^{(c)} = \sum_{s=1}^S \sum_{k=1}^K \text{softmax}(\boldsymbol{\alpha})_s \text{softmax}(\boldsymbol{\beta})_k \mathbf{h}_{s,k}^{(c)}. \quad (9)$$

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

$$\mathbf{y} = \mathbf{x}_{\text{skip}} + f_{\text{mix}}(\sigma(\gamma) \mathbf{h}). \quad (10)$$

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), making

LGE a lightweight residual enhancer that strengthens detail sensitivity prior to subsequent fusion.

**Wavelet variant (LGE-W).** LGE-W follows the same LGF decomposition and aggregation in Eq. (8)–(10), but replaces  $f_{\text{mix}}$  with a wavelet-transform convolution (WTConv) when the input/output channels match. This variant injects multi-resolution wavelet-domain mixing with minimal architectural intrusion; when channel dimensions do not match, we fall back to the standard  $3 \times 3$  convolution for stability.

## 4. Experiment

### 4.1. Datasets and Metrics

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

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

### 4.2. Configuration

The experimental configuration is detailed in Table 1.

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

Environment	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

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

165 **5. Main Results**166 **6. Analyses and Discussion**167 **7. Conclusion**168 **References**169  
170 Carion, N., Massa, F., Synnaeve, G., Usunier, N., Kirillov,  
171 A., and Zagoruyko, S. End-to-end object detection with  
172 transformers. *arXiv preprint arXiv:2005.12872*, 2020.173  
174 Chen, L., Gu, L., Li, L., Yan, C., and Fu, Y. Frequency  
175 dynamic convolution for dense image prediction. *arXiv  
176 preprint arXiv:2503.18783*, 2025.177  
178 Du, D., Qi, Y., Yu, H., Yang, Y., Duan, K., Li, G., Zhang,  
179 W., Huang, Q., and Tian, Q. The unmanned aerial vehicle  
180 benchmark: Object detection and tracking. In *Proceedings  
181 of the European Conference on Computer Vision (ECCV)*, pp.  
182 370–386, 2018.183  
184 Du, D., Zhu, P., Wen, L., Bian, X., Lin, H., Hu, Q., Peng, T.,  
185 Zheng, J., Wang, X., Zhang, Y., et al. Visdrone-det2019:  
186 The vision meets drone object detection in image chal-  
187 lenge results. In *Proceedings of the IEEE/CVF Interna-  
188 tional Conference on Computer Vision Workshops*, pp.  
189 0–0, 2019.190  
191 Finder, S. E., Amoyal, R., Treister, E., and Freifeld, O.  
192 Wavelet convolutions for large receptive fields. In  
193 Leonardis, A., Ricci, E., Roth, S., Russakovsky, O., Sat-  
194 tterer, T., and Varol, G. (eds.), *Computer Vision – ECCV  
195 2024*, pp. 363–380, Cham, 2025. Springer Nature Switzer-  
196 land. ISBN 978-3-031-72949-2.197  
198 Khanam, R. and Hussain, M. YOLOv11: An overview  
199 of the key architectural enhancements. *arXiv preprint  
200 arXiv:2410.17725*, 2024.201  
202 Lin, T.-Y., Dollár, P., Girshick, R., He, K., Hariharan, B.,  
203 and Belongie, S. Feature pyramid networks for object  
204 detection. In *Proceedings of the IEEE Conference on  
205 Computer Vision and Pattern Recognition (CVPR)*, pp.  
206 2117–2125, 2017.207  
208 Liu, W., Anguelov, D., Erhan, D., Szegedy, C., Reed, S.,  
209 Fu, C.-Y., and Berg, A. C. SSD: Single shot multibox  
210 detector. In *Computer Vision–ECCV 2016: 14th Euro-  
211 pean Conference, Amsterdam, The Netherlands, October  
212 11–14, 2016, Proceedings, Part I*, pp. 21–37. Springer,  
213 2016.214  
215 Rao, Y., Zhao, W., Zhu, Z., Zhou, J., and Lu, J. GFNet:  
216 Global filter networks for visual recognition. *IEEE Trans-  
217 actions on Pattern Analysis and Machine Intelligence  
218 (TPAMI)*, 45(9):10960–10973, September 2023. doi:  
219 10.1109/TPAMI.2023.3263824.Ren, S., He, K., Girshick, R., and Sun, J. Faster R-CNN:  
Towards real-time object detection with region proposal  
networks. In *Advances in Neural Information Processing  
Systems (NeurIPS)*, volume 28, 2015.Shi, Z., Hu, J., Ren, J., Ye, H., Yuan, X., Ouyang, Y., He, J.,  
Ji, B., and Guo, J. HS-FPN: High frequency and spatial  
perception fpn for tiny object detection. *arXiv preprint  
arXiv:2412.10116*, 2025.Wang, A., Chen, H., Liu, L., Chen, K., Lin, Z., Han, J., and  
Ding, G. YOLOv10: Real-time end-to-end object  
detection. *arXiv preprint arXiv:2405.14458*, 2024.Wang, A., Chen, H., Lin, Z., Han, J., and Ding, G. LSNet:  
See large, focus small. In *Proceedings of the IEEE/CVF  
Conference on Computer Vision and Pattern Recognition  
(CVPR)*, 2025.Xia, G.-S., Bai, X., Ding, J., Zhu, Z., Belongie, S., Luo,  
J., Datcu, M., Pelillo, M., and Zhang, L. Dota: A large-  
scale dataset for object detection in aerial images. In  
*Proceedings of the IEEE conference on computer vision  
and pattern recognition*, pp. 3974–3983, 2018.Xiao, Y., Xu, T., Xin, Y., and Li, J. FBRT-YOLO: Faster and  
better for real-time aerial image detection. *arXiv preprint  
arXiv:2504.20670*, 2025.Yu, X., Gong, Y., Jiang, N., Ye, Q., and Han, Z. Scale  
match for tiny person detection. In *Proceedings of the  
IEEE/CVF Winter Conference on Applications of Com-  
puter Vision*, pp. 1257–1265, 2020.Zhang, H., Li, F., Liu, S., Zhang, L., Su, H., Zhu, J., Ni,  
L. M., and Shum, H.-Y. DINO: DETR with improved  
denoising anchor boxes for end-to-end object detection.  
*arXiv preprint arXiv:2203.03605*, 2022.Zhao, Y., Lv, W., Xu, S., Wei, J., Wang, G., Dang, Q., Liu,  
Y., and Chen, J. Detrs beat yolos on real-time object  
detection. *arXiv preprint arXiv:2304.08069*, 2024.

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

221   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  
222   page can be added. If you want, you can use an appendix like this one.  
223

224   The `\onecolumn` command above can be kept in place if you prefer a one-column appendix, or can be removed if you  
225   prefer a two-column appendix. Apart from this possible change, the style (font size, spacing, margins, page numbering, etc.)  
226   should be kept the same as the main body.  
227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

256

257

258

259

260

261

262

263

264

265

266

267

268

269

270

271

272

273

274