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

Anonymous Authors¹

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-

¹Anonymous 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. Wavelet-Difference Gate (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 1×1
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 2×2 block is trans-
108 formed by a 2×2 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_{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 3×3 kernel with a learnable θ . 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 θ from the center coefficient of a
130 3×3 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 1×1 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. 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

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

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 angular term:@@

$$\begin{aligned}c_k &= \cos \phi_k, & s_k &= \sin \phi_k, \\ u' &= u c_k + v s_k, & v' &= -u s_k + v c_k, \\ r &= \sqrt{u'^2 + v'^2} + \varepsilon, & \theta &= \text{atan2}(v', u'), \\ \mathbf{g}_{s,k}(u, v) &= \exp\left(-\frac{\log^2(r/\rho_s)}{2 \log^2 2}\right) \cos \theta.\end{aligned}\quad (9)$$

where $\phi_k = k\pi/K$ and ρ_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

$$\mathbf{h}^{(c)} = \sum_{s=1}^S \sum_{k=1}^K \text{softmax}(\boldsymbol{\alpha})_s \text{softmax}(\boldsymbol{\beta})_k \mathbf{h}_{s,k}^{(c)}. \quad (10)$$

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

$$\mathbf{y} = \mathbf{x}_{\text{skip}} + f_{\text{mix}}(\sigma(\gamma) \mathbf{h}). \quad (11)$$

Here \mathbf{x}_{skip} is either the identity mapping (when channels match) or a 1×1 projection. In our implementation, f_{mix} is a 3×3 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. (8)–(11) but replaces f_{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:@@

$$\text{WTConv}(\mathbf{z}) = \mathcal{S}_0 \mathcal{D}_0(\mathbf{z}) + \text{IDWT}(\mathcal{S} \mathcal{D}_4(\text{DWT}(\mathbf{z}))), \quad (12)$$

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.

3.3. Frequency-Driven Head (FDHead)

We finally introduce Frequency-Driven Head (FDHead), a frequency-aware detection head that enhances localization sensitivity to fine boundaries while keeping the standard dense prediction interface. FDHead is instantiated over multi-scale feature maps $\{\mathbf{x}_i\}_{i=1}^N$ (e.g., $P2-P5$) and uses a shared prediction tower to reduce parameters.

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 merges several directional-difference filters into a single re-parameterized convolution, which biases the tower toward edge/contour responses important for small objects:

$$\mathbf{f}_i = \mathcal{T}(\mathbf{x}_i), \quad \mathcal{T} = \mathcal{T}_{\text{share}} \circ \mathcal{T}_{1 \times 1}. \quad (13)$$

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 $P2$). We split channels $\mathbf{f}_1 = [\mathbf{f}_a, \mathbf{f}_b]$ where \mathbf{f}_a keeps a small frequency-aware subset and \mathbf{f}_b bypasses unchanged. We then estimate high-frequency energy from Haar subbands

165 and produce a channel-wise gain:
 166
 167
 168

$$\begin{aligned}
 169 \quad & (\mathbf{f}_{LL}, \mathbf{f}_{LH}, \mathbf{f}_{HL}, \mathbf{f}_{HH}) = \text{DWT}(\mathbf{f}_a), \\
 170 \quad & \mathbf{w} = \text{softmax}(\boldsymbol{\omega}), \\
 171 \quad & \mathbf{h} = w_{LH} |\mathbf{f}_{LH}| + w_{HL} |\mathbf{f}_{HL}| + w_{HH} |\mathbf{f}_{HH}| \\
 172 \quad & \mathbf{g} = \text{Gate}(\text{AvgPool}(\mathbf{h})), \\
 173 \quad & \tilde{\mathbf{f}}_a = \mathbf{f}_a \odot (1 + \alpha \mathbf{g}). \\
 174 \quad & \\
 175 \quad & \text{where } \text{Gate}(\cdot) \text{ is a small channel MLP implemented by } 1 \times 1 \\
 176 \quad & \text{convs with sigmoid output. We then form } \tilde{\mathbf{f}}_1 = [\tilde{\mathbf{f}}_a, \mathbf{f}_b] \text{ and} \\
 177 \quad & \text{apply it only to the box branch: high-frequency cues mainly} \\
 178 \quad & \text{sharpen boundary-aligned localization, while keeping the} \\
 179 \quad & \text{classification stream unchanged avoids amplifying texture} \\
 180 \quad & \text{noise. For the remaining levels } i > 1, \text{ we set } \tilde{\mathbf{f}}_i = \mathbf{f}_i. \\
 181 \quad & \\
 182 \quad & \textbf{Box/class prediction and decoding.} FDHead predicts per-} \\
 183 \quad & \text{location class logits and distributional box offsets (DFL)} \\
 184 \quad & \text{as} \\
 185 \quad & \\
 186 \quad & \\
 187 \quad & \\
 188 \quad & \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 (15) \\
 189 \quad & \\
 190 \quad & \\
 191 \quad & \\
 192 \quad & \text{and decodes boxes by } \hat{\mathbf{B}} = \text{dist2bbox}(\text{DFL}(\mathbf{b}), \mathbf{A}) \cdot \mathbf{s} \text{ with} \\
 193 \quad & \text{anchors } \mathbf{A} \text{ and strides } \mathbf{s}. \text{ This design targets small objects} \\
 194 \quad & \text{by frequency-gating only the finest level while keeping the} \\
 195 \quad & \text{remaining head computation shared and lightweight.} \\
 196 \quad & \\
 197 \quad & \textbf{4. Experiment} \\
 198 \quad & \\
 199 \quad & \textbf{4.1. Datasets and Metrics} \\
 200 \quad & \\
 201 \quad & \text{We evaluate our framework on four benchmarks to demon-} \\
 202 \quad & \text{strate its robustness and cross-domain generalization: Vis-} \\
 203 \quad & \text{Drone2019 (Du et al., 2019), TinyPerson (Yu et al., 2020),} \\
 204 \quad & \text{UAVDT (Du et al., 2018), and DOTA v1 (Xia et al., 2018).} \\
 205 \quad & \textbf{VisDrone2019} is our primary benchmark and is particularly} \\
 206 \quad & \text{challenging due to dense small objects and severe scale var-} \\
 207 \quad & \text{iation, where most targets are smaller than } 50 \times 50 \text{ pixels.} \\
 208 \quad & \\
 209 \quad & \text{We report both accuracy and efficiency, including mAP}_{50}, \\
 210 \quad & \text{the number of parameters, GFLOPs, model size, and FPS.} \\
 211 \quad & \\
 212 \quad & \textbf{4.2. Configuration} \\
 213 \quad & \\
 214 \quad & \text{The experimental configuration is detailed in Table 1.} \\
 215 \quad & \\
 216 \quad & \text{For YOLO-style architectures, models are trained for 300} \\
 217 \quad & \text{epochs with an input resolution of } 640 \times 640 \text{ and batch size} \\
 218 \quad & 16, \text{ using SGD optimization. Unless otherwise specified,} \\
 219 \quad & \text{Mosaic augmentation is enabled throughout training; we} \\
 219 \quad & \text{use 4 dataloader workers and disable AMP.}$$

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

5. Main Results

6. Analyses and Discussion

7. Conclusion

References

- Carion, N., Massa, F., Synnaeve, G., Usunier, N., Kirillov, A., and Zagoruyko, S. End-to-end object detection with transformers. *arXiv preprint arXiv:2005.12872*, 2020.
- Chen, L., Gu, L., Li, L., Yan, C., and Fu, Y. Frequency dynamic convolution for dense image prediction. *arXiv preprint arXiv:2503.18783*, 2025.
- Du, D., Qi, Y., Yu, H., Yang, Y., Duan, K., Li, G., Zhang, W., Huang, Q., and Tian, Q. The unmanned aerial vehicle benchmark: Object detection and tracking. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 370–386, 2018.
- Du, D., Zhu, P., Wen, L., Bian, X., Lin, H., Hu, Q., Peng, T., Zheng, J., Wang, X., Zhang, Y., et al. Visdrone-det2019: The vision meets drone object detection in image challenge results. In *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops*, pp. 0–0, 2019.
- Finder, S. E., Amoyal, R., Treister, E., and Freifeld, O. Wavelet convolutions for large receptive fields. In Leonardis, A., Ricci, E., Roth, S., Russakovsky, O., SatTLer, T., and Varol, G. (eds.), *Computer Vision – ECCV 2024*, pp. 363–380, Cham, 2025. Springer Nature Switzerland. ISBN 978-3-031-72949-2.
- Khanam, R. and Hussain, M. YOLOv11: An overview of the key architectural enhancements. *arXiv preprint arXiv:2410.17725*, 2024.
- Lin, T.-Y., Dollár, P., Girshick, R., He, K., Hariharan, B., and Belongie, S. Feature pyramid networks for object

- 220 detection. In *Proceedings of the IEEE Conference on*
221 *Computer Vision and Pattern Recognition (CVPR)*, pp.
222 2117–2125, 2017.
- 223 Liu, W., Anguelov, D., Erhan, D., Szegedy, C., Reed, S.,
224 Fu, C.-Y., and Berg, A. C. SSD: Single shot multibox
225 detector. In *Computer Vision–ECCV 2016: 14th Euro-*
226 *pean Conference, Amsterdam, The Netherlands, October*
227 *11–14, 2016, Proceedings, Part I*, pp. 21–37. Springer,
228 2016.
- 229 Rao, Y., Zhao, W., Zhu, Z., Zhou, J., and Lu, J. GFNet:
230 Global filter networks for visual recognition. *IEEE Trans-*
231 *actions on Pattern Analysis and Machine Intelligence*
232 (*TPAMI*), 45(9):10960–10973, September 2023. doi:
233 10.1109/TPAMI.2023.3263824.
- 234 Ren, S., He, K., Girshick, R., and Sun, J. Faster R-CNN:
235 Towards real-time object detection with region proposal
236 networks. In *Advances in Neural Information Processing*
237 *Systems (NeurIPS)*, volume 28, 2015.
- 238 Shi, Z., Hu, J., Ren, J., Ye, H., Yuan, X., Ouyang, Y., He, J.,
239 Ji, B., and Guo, J. HS-FPN: High frequency and spatial
240 perception fpn for tiny object detection. *arXiv preprint*
241 *arXiv:2412.10116*, 2025.
- 242 Wang, A., Chen, H., Liu, L., Chen, K., Lin, Z., Han, J.,
243 and Ding, G. YOLOv10: Real-time end-to-end object
244 detection. *arXiv preprint arXiv:2405.14458*, 2024.
- 245 Wang, A., Chen, H., Lin, Z., Han, J., and Ding, G. LSNet:
246 See large, focus small. In *Proceedings of the IEEE/CVF*
247 *Conference on Computer Vision and Pattern Recognition*
248 (*CVPR*), 2025.
- 249 Xia, G.-S., Bai, X., Ding, J., Zhu, Z., Belongie, S., Luo,
250 J., Datcu, M., Pelillo, M., and Zhang, L. Dota: A large-
251 scale dataset for object detection in aerial images. In
252 *Proceedings of the IEEE conference on computer vision*
253 *and pattern recognition*, pp. 3974–3983, 2018.
- 254 Xiao, Y., Xu, T., Xin, Y., and Li, J. FBRT-YOLO: Faster and
255 better for real-time aerial image detection. *arXiv preprint*
256 *arXiv:2504.20670*, 2025.
- 257 Yu, X., Gong, Y., Jiang, N., Ye, Q., and Han, Z. Scale
258 match for tiny person detection. In *Proceedings of the*
259 *IEEE/CVF Winter Conference on Applications of Com-*
260 *puter Vision*, pp. 1257–1265, 2020.
- 261 Zhang, H., Li, F., Liu, S., Zhang, L., Su, H., Zhu, J., Ni,
262 L. M., and Shum, H.-Y. DINO: DETR with improved
263 denoising anchor boxes for end-to-end object detection.
264 *arXiv preprint arXiv:2203.03605*, 2022.
- 265 Zhao, Y., Lv, W., Xu, S., Wei, J., Wang, G., Dang, Q., Liu,
266 Y., and Chen, J. Detrs beat yolos on real-time object
267 detection. *arXiv preprint arXiv:2304.08069*, 2024.

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

276
277 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
278 page can be added. If you want, you can use an appendix like this one.

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

282

283

284

285

286

287

288

289

290

291

292

293

294

295

296

297

298

299

300

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