

TF-SSD: A Strong Pipeline via Synergic Mask Filter for Training-free Co-salient Object Detection

Supplementary Material

001 A. Overview

002 This supplementary material provides additional details to
003 complement our main paper. We elaborate on our imple-
004 mentation details in Sec. B, present more quantitative and
005 qualitative results in Sec. C and Sec. D, and discuss the lim-
006 itations of our method in Sec. E.

007 B. Additional implementation details

008 To handle challenging cases, several mechanisms are nec-
009 essary to ensure robustness in practice. Due to space lim-
010 itations in the main paper, we describe these details below.

011 B.1. Fallback mechanism in ISF

012 Our Intra-image Saliency Filter (ISF) relies on DINO’s at-
013 tention maps to identify salient objects from SAM [3] pro-
014 posals. However, in some challenging scenarios, QMG may
015 fail to segment salient masks, where all masks yield low
016 saliency scores $s_{n,t}^{sal}$. To ensure robustness, we employ a
017 fallback mechanism that directly generates masks from at-
018 tention maps to avoid these bad cases.

019 In particular, if the maximum saliency score $\max_t(s_{n,t}^{sal})$
020 is lower than a threshold τ_{fb} , the attention map \mathcal{A}_n is bina-
021 rized to obtain a mask $m_{n,t}^{fb}$ that captures the salient region,
022 formulated as:

$$023 \quad m_{n,t}^{fb}(x, y) = \begin{cases} 1 & \text{if } \mathcal{A}_n(x, y) > \tau_{attn}, \\ 0 & \text{otherwise,} \end{cases} \quad (1)$$

024 where τ_{attn} is the binarization threshold for \mathcal{A}_n . This
025 single fallback mask $m_{n,t}^{fb}$ will replace the low-quality
026 masks. This operation ensures that each image retains at
027 least one salient candidate mask for subsequent processing.
028 In our implementation, we set $\tau_{fb} = 0.05$, and τ_{attn} is dy-
029 namically set to a threshold that retains the top 50% of the
030 highest attention values for each map \mathcal{A}_n .

031 B.2. Merging of multiple co-salient objects in IPS

032 Our Inter-image Prototype Selector (IPS) is designed to
033 identify the target co-salient object in each image. How-
034 ever, individual images may contain multiple instances of
035 the co-salient object, and SAM often segments them sepa-
036 rately. To obtain accurate results for CoSOD, we have to
037 merge these separate masks into a single one.

038 To address such cases, we employ a dual-verification
039 mechanism with two thresholds: a semantic similarity
040 threshold τ_{sem} and a consistency difference threshold τ_{diff} .

041 For each image I_n with a selected mask \hat{m}_n , we examine
042 the remaining masks in $\mathcal{M}_n^{salient}$ to identify extra co-salient
043 objects. A candidate mask m_j is considered an extra co-
044 salient object if it satisfies the following two conditions:

- 045 1) **Semantic similarity:** The candidate mask is semanti-
046 cally similar to the primary mask, measured by pairwise
047 similarity $\langle p_{\hat{m}_n}, p_j \rangle \geq \tau_{sem}$.
- 048 2) **Cross-image consistency:** The candidate mask ex-
049 hibits comparable cross-image consistency, measured by
050 $|S_j^{co} - S_{\hat{m}_n}^{co}| < \tau_{diff}$.

051 The first condition ensures semantic coherence within
052 the image, while the second verifies that the candidate mask
053 represents an object that appears consistently across the
054 group. Verified masks are merged with the primary mask
055 to form the final segmentation mask. In our implementa-
056 tion, τ_{sem} is dynamically set as the top 80% of intra-image
057 pairwise similarities, and τ_{diff} is set to 0.1.

058 C. Additional quantitative results

059 C.1. Different SAM area ratios

060 The area ratio threshold τ_{area} is a critical parameter in our
061 QMG that filters out excessively small masks (Eq. 2 in the
062 main paper). To evaluate its impact on CoSOD perfor-
063 mance, we conduct ablation experiments with four differ-
064 ent threshold values: 0.002, 0.005, 0.01, and 0.02 on the
065 CoCA [2] dataset. We employ two metrics: (1) **F-measure:**
066 the final CoSOD performance obtained using the complete
067 TF-SSD framework; (2) **Oracle F:** the upper bound perfor-
068 mance computed using selected oracle masks from QMG.

069 As shown in Tab. 1, when the threshold is too low
070 (0.005), excessive trivial masks are retained in the candi-
071 date set, introducing noise that degrades the subsequent ISF
072 and IPS. Meanwhile, the Oracle F score reflects the qual-
073 ity of the candidate pool. When the threshold is too high
074 (0.02), some valid co-salient objects are incorrectly filtered
075 out, which limits both the quality of the candidate pool and
076 the final performance. $\tau_{area} = 0.01$ strikes an optimal
077 balance, which effectively removes trivial masks while pre-
078 serving diverse co-salient object candidates.

079 C.2. Performance using other DINO backbones

080 Our main experiments use DINO ViT-B/8 [1] as the back-
081 bone for feature extraction in ISF. To provide a broader per-
082 formance assessment, we compare it with two other notable
083 backbones: (1) ViT-B/16 [1] from the same DINO frame-
084 work with a larger patch size; (2) ViT-B/14 [4] from the

Table 1. Ablation study on different area ratio thresholds τ_{area} on the CoCA dataset.

τ_{area}	F_{β}^{\max}	Oracle F_{β}^{\max}
0.002	0.545	0.696
0.005	0.601	0.722
0.01	0.686	0.768
0.02	0.653	0.737

085 more recent DINOv2 framework. We conduct experiments
086 on the CoCA dataset, which highlights the trade-offs be-
087 tween different DINO versions.

088 As shown in Tab. 2, DINO ViT-B/8 significantly out-
089 performs other configurations across all three metrics. A
090 smaller patch size enables our model to capture finer spatial
091 details, which is crucial for feature matching and prototype
092 selection in our ISF and IPS. In contrast, although DINOv2
093 ViT-B/14 utilizes more advanced pretraining strategies, its
094 larger patch size limits the perception of fine-grained infor-
095 mation. This result demonstrates that a smaller patch size is
096 more suitable for CoSOD that require precise localization
097 and fine-grained feature alignment.

Table 2. Performance comparison of different DINO backbones on the CoCA dataset.

Backbone	F_{β}^{\max}	E_{ξ}^{\max}	S_{α}
DINO ViT-B/16	0.556	0.742	0.697
DINOv2 ViT-B/14	0.607	0.783	0.708
DINO ViT-B/8	0.686	0.815	0.763

098 C.3. Ablation on hyperparameter settings

099 We conduct additional ablation studies to evaluate the im-
100 pact of hyperparameter settings.

101 **Quality score weights.** The balanced quality score $S_{n,t}^{ba}$
102 (Eq. 5 in the main paper) combines IoU prediction confi-
103 dence and area preference through weights α and β , where
104 $\alpha + \beta = 1$ to ensure normalized weighting between the
105 two terms. Tab. 3 illustrates the performance with different
106 weight settings on the CoCA dataset.

Table 3. Ablation on quality score weights on the CoCA dataset.

α	β	F_{β}^{\max}	E_{ξ}^{\max}	S_{α}
0.0	1.0	0.473	0.702	0.611
0.3	0.7	0.591	0.737	0.651
0.5	0.5	0.658	0.795	0.744
0.7	0.3	0.686	0.815	0.763
1.0	0.0	0.641	0.782	0.735

107 **Ideal size range.** The ideal size range $[r_{min}, r_{max}]$ (Eq. 4
108 in the main paper) defines the appropriate size for co-salient

objects. Tab. 4 presents results for different range set-
109 tings. The range $[0.15, 0.7]$ performs best, suggesting that
110 co-salient objects typically occupy 15%-70% of the image
111 area. More restrictive ranges (e.g. $[0.2, 0.5]$) limit the num-
112 ber of valid masks, while much wider ranges (e.g. $[0.05,$
113 $0.85]$) fail to effectively filter out background and large
114 masks.
115

Table 4. Ablation on ideal size range on the CoCA dataset.

r_{min}	r_{max}	F_{β}^{\max}	E_{ξ}^{\max}	S_{α}
0.05	0.85	0.661	0.798	0.747
0.1	0.6	0.678	0.808	0.756
0.15	0.7	0.686	0.815	0.763
0.2	0.5	0.619	0.741	0.669

Table 5. Ablation on overlap threshold on the CoCA dataset.

τ_{con}	F_{β}^{\max}	E_{ξ}^{\max}	S_{α}
0.5	0.652	0.791	0.741
0.7	0.673	0.805	0.754
0.85	0.686	0.815	0.763
0.95	0.679	0.796	0.756

Overlap threshold. The overlap threshold τ_{con} controls
116 how aggressively overlapping masks are filtered. As shown
117 in Tab. 5, $\tau_{con} = 0.85$ achieves the optimal performance.
118 Lower thresholds (e.g., 0.5) retain too many redundant
119 masks, while higher thresholds (e.g., 0.95) will incorrectly
120 remove valid masks that partially overlap with larger ob-
121 jects. They both lead to the reduced diversity of candidate
122 masks for subsequent processing.
123

D. Additional qualitative results

In this section, we present visualizations of several key
125 stages in our TF-SSD framework.
126

D.1. Visualization of SAM proposals

To demonstrate SAM’s segmentation capability and the ef-
128 fectiveness of our QMG module, Fig. 1 visualizes the top-5
129 mask proposals ranked by our quality score $S_{n,t}^{ba}$ (Eq. 5 in
130 the main paper). The masks with green borders represent
131 our final predictions after the complete TF-SSD pipeline.
132 It’s observed that SAM can generate mask proposals that
133 closely align with the GT, and our quality score $S_{n,t}^{ba}$ effec-
134 tively ranks these masks at the top positions. It validates the
135 effectiveness of SAM for mask generation and confirms that
136 QMG can successfully identify promising candidate masks
137 for subsequent processing.
138

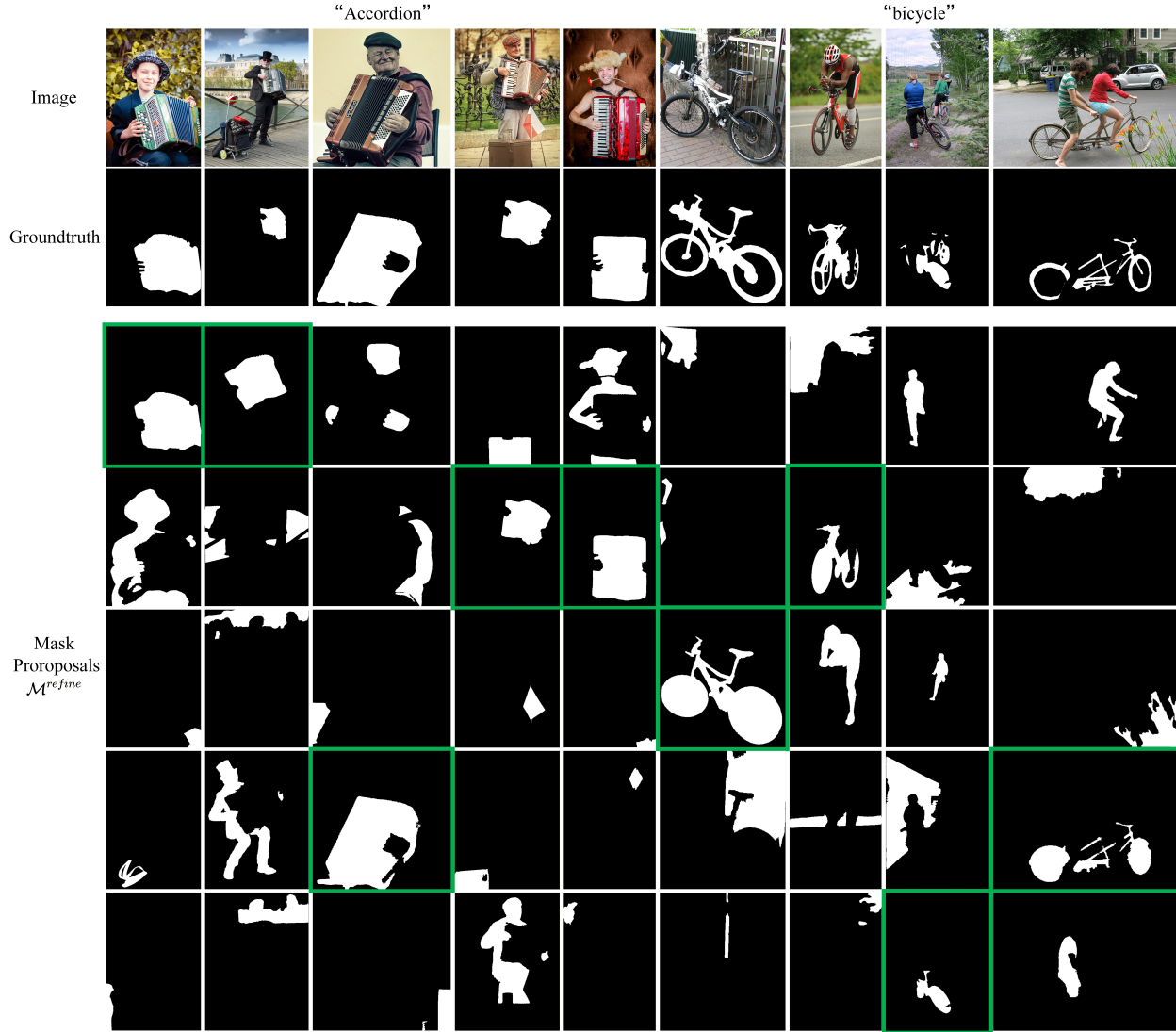


Figure 1. Visualization of SAM mask proposals for "Accordion" and "bicycle" image groups. For each group, we show the original images (row 1), ground truth (GT) masks (row 2), and the top-5 mask proposals ranked by quality score $S_{n,t}^{ba}$ (rows 3-7). Green borders indicate the final predictions selected by our TF-SSD pipeline.

139 D.2. Visualization of DINO attention maps

140 To illustrate the saliency-aware capability of DINO's atten-
 141 tion mechanism, Fig. 2 visualizes the attention maps \mathcal{A}_n of
 142 DINO's CLS-token. The attention maps can highlight the
 143 salient object regions that closely align with the GT masks.
 144 These observations validate the effectiveness of leveraging
 145 DINO attention maps to select salient mask proposals from
 146 the refined candidate set derived from our QMG.

147 E. Limitations

148 While our TF-SSD framework achieves strong performance
 149 on most CoSOD benchmarks, it still exhibits limitations in
 150 detecting small co-salient objects. As illustrated in Fig. 3,

the "moon" category contains predominantly small objects
 that occupy a relatively small portion of the image. Due to
 SAM's initial area-based filtering mechanism in our QMG,
 masks with small area ratios ($r_{n,t}^{area} < \tau_{area}$) are typically
 denoted as trivial objects that are filtered out. In addition,
 small objects also tend to receive lower quality scores $S_{n,t}^{ba}$
 from DINO's attention maps, which leads to wrong identi-
 fication for co-salient objects.

Our future work concentrates on addressing these bad
 cases in challenging scenarios. We will explore incorporat-
 ing semantic cues into an adaptive threshold that is modu-
 lated according to contexts of the image. We argue that it
 can identify the small salient objects from the trivial ones.

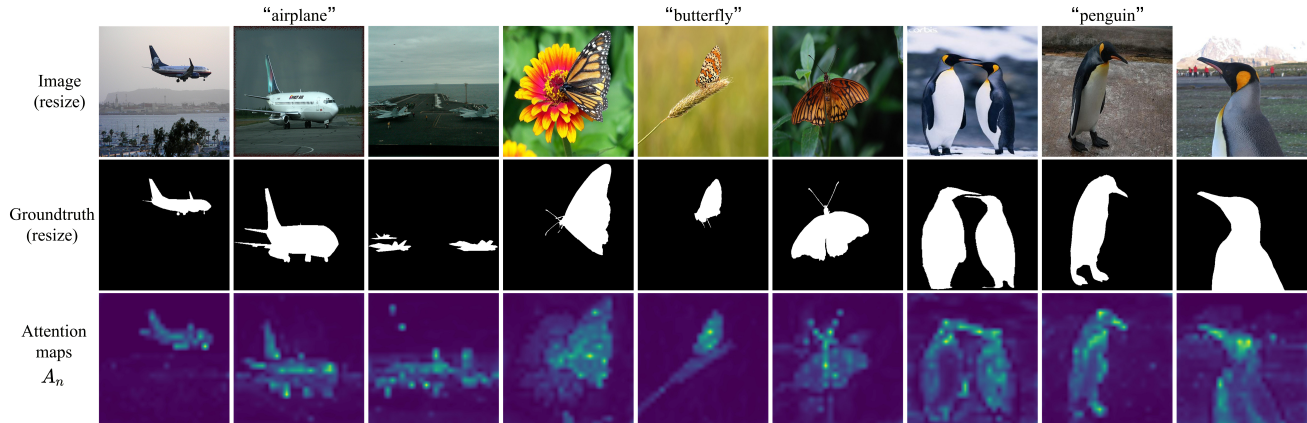


Figure 2. Visualization of DINO attention maps for three categories: "airplane", "butterfly", and "penguin". Row 1: Resized input images. Row 2: Resized GT masks. Row 3: DINO attention maps \mathcal{A}_n that naturally highlight salient objects.

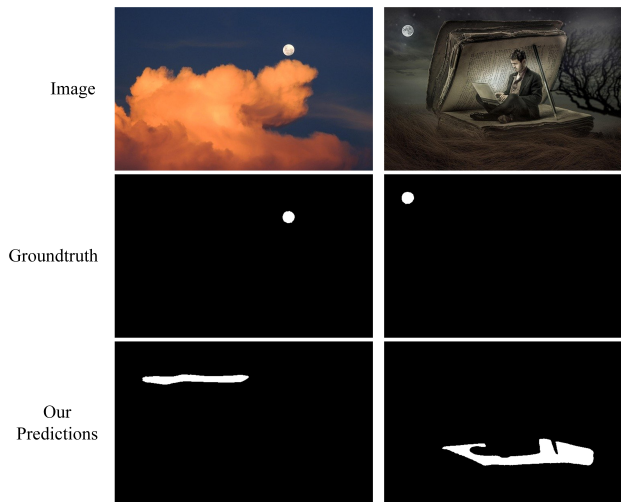


Figure 3. Limitation on small object detection. Visualizations of failed cases from the "moon" category.

Learning robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023. 1

180
181

164

References

165

166

167

168

169

170

171

172

173

174

175

176

177

178

179

- [1] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *ICCV*, pages 9650–9660, 2021. 1
- [2] Deng-Ping Fan, Zheng Lin, Ge-Peng Ji, Dingwen Zhang, Huazhu Fu, and Ming-Ming Cheng. Taking a deeper look at co-salient object detection. In *CVPR*, pages 2919–2929, 2020. 1
- [3] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *ICCV*, pages 4015–4026, 2023. 1
- [4] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: