

G-SalAlignMamba: Geometry-Aware Vision Mamba for Dual-Modal Salient Object Detection

Author Name

Affiliation

email@example.com

Abstract

1 Recently, Visual State Space Models offer powerful global modeling for Dual-modal Salient Object Detection (SOD). However, they are still constrained by three inherent limitations: first, Mamba’s strict reliance on sequential ordering makes it highly vulnerable to *cross-modal geometric misalignment*, where spatial shifts disrupt token correspondence more severely compared with deep learning frameworks; second, general indiscriminate scanning treating all tokens equally may lead to *signal dilution*, where sparse foreground features are overwhelmed by massive background noise; third, conventional decoders rely on implicit upsampling, causing *boundary degradation* during resolution recovery. To address these challenges, we propose **G-SalAlignMamba**, a geometry-aware framework tailored for dual-modal SOD. We introduce **Geometry-Aware Encoding** with explicit alignment to correct spatial shifts, **Semantics-Informed Refinement** to prevent signal dilution by prioritizing foregrounds, and **Structure-Preserving Decoding** that integrates explicit alignment with unsupervised boundary refinement. Extensive experiments demonstrate that our proposed G-SalAlignMamba significantly outperforms state-of-the-art methods on RGB-D and RGB-T benchmarks. Furthermore, our framework exhibits favorable computational efficiency (30.41 FPS with 83.80M parameters), effectively breaking the bottleneck between global modeling capability and computational cost inherent in existing multi-modal networks.

1 Introduction

34 Dual-modal Salient Object Detection (SOD) aims to robustly identify visually attractive regions by integrating RGB images with auxiliary modalities (e.g., Depth or Thermal). Generally speaking, this field has primarily relied on Convolutional Neural Networks (CNNs) [Fan *et al.*, 2020a] or Transformer architectures [Liu *et al.*, 2021]. Although these models have achieved significant progress, CNNs are limited by

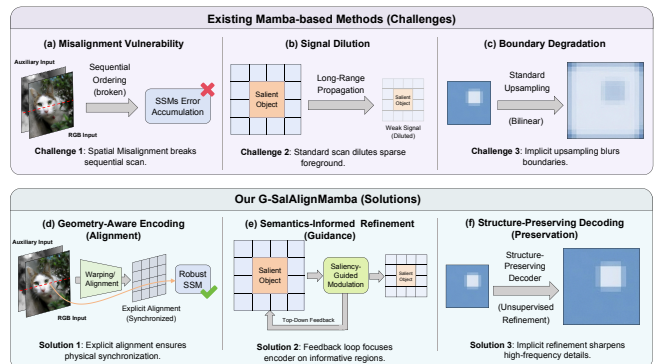


Figure 1: Top row illustrates three critical limitations in existing Mamba-based SOD: (a) Misalignment Vulnerability disrupts sequential state transitions; (b) Signal Dilution weakens sparse foreground signals during long-range propagation; (c) Boundary Degradation results from implicit upsampling. Bottom row presents our proposed G-SalAlignMamba solutions: (d) Geometry-Aware Encoding enforces physical synchronization; (e) Semantics-Informed Refinement prioritizes informative regions via feedback; (f) Structure-Preserving Decoding achieves sharp boundaries via implicit, unsupervised refinement.

local receptive fields that restrict their ability to model complex global contexts. For Transformers, despite possessing global receptive fields, face computational costs that grow quadratically with resolution. To break the stalemate between computational efficiency and global modeling, Visual State Space Models (SSMs), represented by VMamba [Liu *et al.*, 2024b], have recently garnered widespread attention for their ability to model long-range dependencies with linear computational complexity $\mathcal{O}(N)$ [Gu and Dao, 2024]. While promising, we argue that the direct application of VMamba to dual-modal SOD is severely constrained by three critical limitations rooted in current methodologies, as shown in Figure 1:

1. Misalignment Vulnerability in Mamba Architecture. Most pioneering Mamba-based fusion frameworks follow a simple “concatenate-then-scan” paradigm [He *et al.*, 2025], overlooking a fatal theoretical flaw—since Mamba strictly relies on a continuous sequential scanning order to update hidden states, it exhibits an extremely high sensitivity to cross-modal geometric misalignment. In dual-modal SOD scenarios, physical spatial shifts caused by parallax or

62 calibration errors directly disrupt the logical physical corre-
63 spondence of tokens within the scanning sequence. This Mis-
64 alignment Vulnerability renders existing simple fusion strate-
65 gies extremely fragile when handling spatial mismatches be-
66 tween modalities, causing the SSM model to construct incor-
67 rect causal relationships and accumulate errors.

68 **2. Signal Dilution Caused by General Scanning Mecha-**
69 **nisms.** General Vision Mamba backbones (e.g., ViM [Zhu *et*
70 *al.*, 2024], VMamba [Liu *et al.*, 2024b]) are largely designed
71 for image classification tasks, where feature distributions are
72 assumed to be relatively uniform. However, in SOD tasks,
73 salient objects typically occupy only a sparse fraction of the
74 image. General raster scanning mechanisms treat all tokens
75 indiscriminately, lacking a discriminative bias towards target
76 regions. In the recursive propagation mechanism of SSMs,
77 the current state depends on the accumulation of all histor-
78 ical information. This implies that sparse foreground signals
79 must traverse extremely long spatial sequences. During this
80 long-range propagation, effective information is prone to be-
81 ing “diluted” or overwhelmed by the vast majority of back-
82 ground noise. Without task-specific modulation, the Mamba
83 encoder easily fails to maintain the prominence of salient fea-
84 tures in deep layers.

85 **3. Boundary Degradation During Sequence Recon-**
86 **struction.** Existing Vision Mamba variants also face a “cou-
87 pling gap” between the SSMs’ sequence representation and
88 the integrity of the image’s 2D physical structure during res-
89 olution recovery. Current decoders in general Vision Mamba
90 models still adhere to traditional practices in vision tasks,
91 employing spatially-agnostic implicit upsampling operators
92 (e.g., bilinear interpolation) [Zhu *et al.*, 2024; Liu *et al.*,
93 2024b]. These operations lack awareness of Mamba’s scan-
94 ning characteristics and fail to utilize the geometric alignment
95 cues captured during the encoding stage [Huang *et al.*, 2021].
96 This dimensional projection distortion during the “sequence-
97 to-space” reconstruction not only triggers pixel-level Corre-
98 spondence Confusion, but also leads to blurred edges and
99 structural deterioration of salient targets—ultimately consti-
100 tuting a structural bottleneck that restricts SOD detection ac-
101 curacy [Qin *et al.*, 2019; Zhao *et al.*, 2019].

102 To address these interconnected issues, we propose **G-**
103 **SalAlignMamba**, a unified geometry-aware framework tai-
104 lored for dual-modal SOD. Guided by an *Alignment-*
105 *Guidance-Preservation* philosophy, we first introduce a
106 **Geometry-Aware Encoding** scheme with coarse-to-fine ex-
107 plicit alignment to enforce physical synchronization, resolv-
108 ing the misalignment vulnerability. Second, a **Semantics-**
109 **Informed Refinement** strategy is proposed to modulate the
110 scanning process via a top-down feedback loop, prevent-
111 ing signal dilution. Finally, a **Structure-Preserving Decod-**
112 **ing** scheme is presented for extending geometric constraints
113 to the upsampling stage, ensuring sharp boundary recovery
114 without explicit supervision.

115 Our main contributions are summarized as follows:

- 116 • We first identify the intrinsic Misalignment Vulnerabil-
117 ity of Mamba’s sequential modeling. To address this, we
118 propose **G-SalAlignMamba**, pioneering a **Geometry-**
119 **Aware Encoding** scheme that enforces physical syn-
120 chronization to construct a robust foundation for multi-

modal SSMs. 121

- 122 • We propose a **Semantics-Informed Refinement** strat-
123 egy to counter the Signal Dilution inherent in generic
124 scanning. By constructing a top-down semantic feed-
125 back loop, we empower the encoder to prioritize
126 sparse salient signals, preventing them from being over-
127 whelmed by background noise.
- 128 • We design a **Structure-Preserving Decoding** scheme to
129 bridge the geometric gap in conventional upsampling.
130 This module enforces geometry-consistent reconstruc-
131 tion with an unsupervised boundary refinement mecha-
132 nism, effectively solving Boundary Degradation without
133 explicit supervision.
- 134 • Extensive experiments on RGB-D and RGB-T bench-
135 marks demonstrate that **G-SalAlignMamba** achieves
136 state-of-the-art performance, validating the necessity of
137 geometric constraints in Mamba-based fusion.

138 2 Related Work

139 2.1 Visual Mamba

140 The success of SSMs has driven their extension to the vi-
141 sual domain, where Vision Mamba (ViM) [Zhu *et al.*, 2024]
142 and VMamba [Liu *et al.*, 2024b] achieve efficient linear-
143 complexity modeling for fundamental recognition tasks. This
144 advantage has rapidly sparked widespread adoption across
145 dense prediction fields [Xing *et al.*, 2024; ?]. In the context
146 of SOD, recent pioneers like Samba [He *et al.*, 2025] have
147 also validated the potential of Mamba for maintaining spa-
148 tial continuity. However, current works primarily focus on
149 single-modal representations or simple fusion strategies, usu-
150 ally overlooking the inherent cross-modal spatial and geomet-
151 ric disparities in dual-modal settings. This limitation inspires
152 us to propose a geometry-aware framework that explicitly in-
153 corporates alignment mechanisms into the VMamba architec-
154 ture to suppress misalignment noise.

155 2.2 RGB-Depth Salient Object Detection (RGB-D 156 SOD)

157 To address RGB-only SOD limitations, many approaches
158 incorporate depth information, named RGB-D SOD, which
159 aims to enhance the model’s geometric and 3D structural un-
160 derstanding. Depth helps detect salient objects when col-
161 or/texture contrasts are weak. A representative work [Fu *et*
162 *al.*, 2021] used Siamese networks with joint learning and
163 densely-cooperative fusion, improving robustness. Another
164 method [Ji *et al.*, 2021] tackled depth noise by using depth
165 calibration followed by cross-modal fusion, boosting saliency
166 detection. However, existing methods still struggle with low-
167 quality depth maps and cross-modal spatial misalignment.

168 2.3 RGB-Thermal Salient Object Detection 169 (RGB-T SOD)

170 RGB-T SOD fuses RGB and thermal imagery, particularly
171 useful in low-light, nocturnal, or high-thermal-contrast en-
172 vironments. RGB-T methods utilize multi-stream networks,
173 attention-based fusion, and multi-scale decoding to integrate

174 RGB and thermal information [Pang *et al.*, 2023]. For ex-
 175 ample, [Zhou *et al.*, 2023] has shown that modeling pixel-
 176 wise positional relations and using Transformer-based en-
 177 coders and decoders leads to more accurate saliency masks
 178 in challenging thermal and visible scenarios. Thermal data
 179 significantly enhances performance when RGB data alone is
 180 insufficient. Nevertheless, they often fail to handle severe par-
 181 allax and thermal noise interference in complex scenes [Tu *et*
 182 *al.*, 2022a].

183 3 Method

184 To systematically resolve the issues of spatial misalignment
 185 and feature degradation in Mamba-based SOD, given paired
 186 RGB and auxiliary inputs $I_R, I_A \in \mathbb{R}^{3 \times H \times W}$, we pro-
 187 pose a geometry-aware dual-modal SOD framework to ad-
 188 dress cross-modal misalignment and signal dilution, as il-
 189 lustrated in Figure 2. The encoder employs a dual-branch
 190 Visual State-Space (VSS) backbone with lightweight Early
 191 Interaction to stabilize intermediate features, followed by a
 192 Coarse-to-Fine Explicit Alignment module that corrects spa-
 193 tial shifts via multi-scale displacement fields and a Comple-
 194 mentary Gated Fusion block to produce a deeply fused rep-
 195 resentation F_{fused} . The refinement module applies Semantics-
 196 Informed Refinement to shallow features, modulating them
 197 via top-down Saliency Priors derived from F_{fused} to prioritize
 198 foreground signals. Finally, the Structure-Preserving Decod-
 199 ing scheme progressively reconstructs the saliency map by
 200 integrating Context-Aware Alignment Upsampling with un-
 201 supervised Multi-Scale Boundary Refinement, ensuring pre-
 202 cise structural recovery and sharp edges essential for high-
 203 performance SOD.

204 3.1 Geometry-Aware Encoding: Cross-Modal 205 Alignment and Fusion

206 Hierarchical Feature Extraction with Early Interaction

207 To extract robust representations from RGB image I_R and
 208 auxiliary input I_A , we employ a dual-branch encoder based
 209 on VSS architecture. Different from CNNs, VSS captures
 210 long-range dependencies with linear complexity $\mathcal{O}(N)$.

211 We share parameters θ_{VSS} between branches to learn
 212 modality-agnostic geometric representations and reduce pa-
 213 rameter overhead as:

$$214 \mathbf{F}_R^{(i)}, \mathbf{F}_A^{(i)} = \text{VSS}(I_R, I_A; \theta_{\text{VSS}}), \quad i \in \{1, 2, 3, 4\}. \quad (1)$$

215 To prevent semantic divergence in deep layers, we intro-
 216 duce a lightweight Early Interaction mechanism. At stage
 217 $i = 2$, intermediate features are stabilized via window-based
 cross-attention as:

$$218 \mathbf{F}_R^{(2)} = \mathbf{F}_R^{(2)} + g_{\text{early}} \cdot \text{WindowCrossAttn}(\mathbf{F}_R^{(2)}, \mathbf{F}_A^{(2)}), \quad (2)$$

219 where $g_{\text{early}} = \sigma(\gamma_{\text{early}})$ is a learnable gate with γ_{early} (ini-
 220 tialized to weakly inject cross-modal information to prevent
 221 early degradation). This ensures that the subsequent align-
 222 ment module operates on features with compatible semantic
 levels.

223 Coarse-to-Fine Explicit Alignment

224 Implicit alignment in SOD often struggles with parallax-
 225 induced displacements, leading to blurred boundaries. Hence,
 226 we propose a **Coarse-to-Fine Explicit Alignment** strategy
 227 that predicts a pixel-wise displacement field $\Delta \mathbf{p}$ through
 228 multi-scale feature extraction and cross-attention fusion.

229 To capture misalignment at varying magnitudes, we con-
 230 struct a feature pyramid at three resolutions (original, 1/2,
 231 and 1/4) by projecting the stage-4 representations via scale-
 232 specific convolutions. For each scale $s \in \{0, 1, 2\}$, these
 233 multi-scale features are subsequently enhanced using cross-
 234 attention to yield the fused representation $\mathbf{F}_{\text{fused}}^{(s)}$. The final
 235 displacement field $\Delta \mathbf{p}$ is an adaptive aggregation of multi-
 236 scale predictions as:

$$\Delta \mathbf{p} = \sum_{s=0}^2 w_s \cdot \text{Predictor}_s(\mathbf{F}_{\text{fused}}^{(s)}), \quad \text{s.t.} \sum w_s = 1, \quad (3)$$

237 where w_s are learnable weights derived from Softmax. The
 238 auxiliary features are then explicitly aligned to the RGB co-
 239 ordinate system via differentiable warping \mathcal{W} :

$$\mathbf{F}_A^{\text{align}} = \mathcal{W}(\mathbf{F}_A^{(4)}, \Delta \mathbf{p}), \quad (4)$$

240 where $\Delta \mathbf{p}$ is bounded by a learnable maximum offset α_{max}
 241 through tanh normalization. The explicit rectification elim-
 242 inates geometric discrepancies, ensuring the spatial consis-
 243 tency critical for preserving sharp boundaries.

244 Complementary Gated Fusion

245 To address the failure of naive fusion in dual-modal SOD
 246 caused by indiscriminately treating noisy inputs, we propose
 247 a **Complementary Gated Fusion** block to balance shared
 248 and complementary cues:

249 1. **Commonality Path:** We integrate shared semantics
 250 using concatenation and Selective Scanning via the Concat
 251 Mamba Fusion Block:

$$\mathbf{F}_{\text{common}} = \text{ConcatMamba}(\mathbf{F}_R^{(4)}, \mathbf{F}_A^{\text{align}}). \quad (5)$$

252 This path concatenates RGB and aligned auxiliary features,
 253 then sequentially processes them through linear projection,
 254 reshape, depthwise convolution, Saliency-Guided SS2D (SG-
 255 SS2D) [He *et al.*, 2025], layer normalization (LN), and a final
 256 reshape to enhance feature representation.

257 2. **Difference-Aware Path:** To explicitly model comple-
 258 mentary signals, we compute feature differences between the
 259 RGB and aligned auxiliary modalities across multiple scales.
 260 The difference maps are encoded, gated via adaptive thresh-
 261 olds, and processed by lightweight SS2D blocks, before being
 262 upsampled and fused into the final difference representation
 263 \mathbf{F}_{diff} . A learnable gating mechanism dynamically balances
 264 these paths as:

$$\mathbf{F}_{\text{fused}} = (1 - g) \cdot \mathbf{F}_{\text{common}} + g \cdot \mathbf{F}_{\text{diff}} + \mathbf{F}_R^{(4)}, \quad (6)$$

265 where $g = \sigma(\gamma_{\text{gate}})$ with γ_{gate} (initialized to bias towards the
 266 commonality path). The residual connection with \mathbf{F}_R ensures
 267 gradient flow and feature preservation. Consequently, this
 268 mechanism adaptively exploits complementarity with sup-
 269 pressing noise, which could ensure robust SOD performance
 270 even with low-quality inputs.

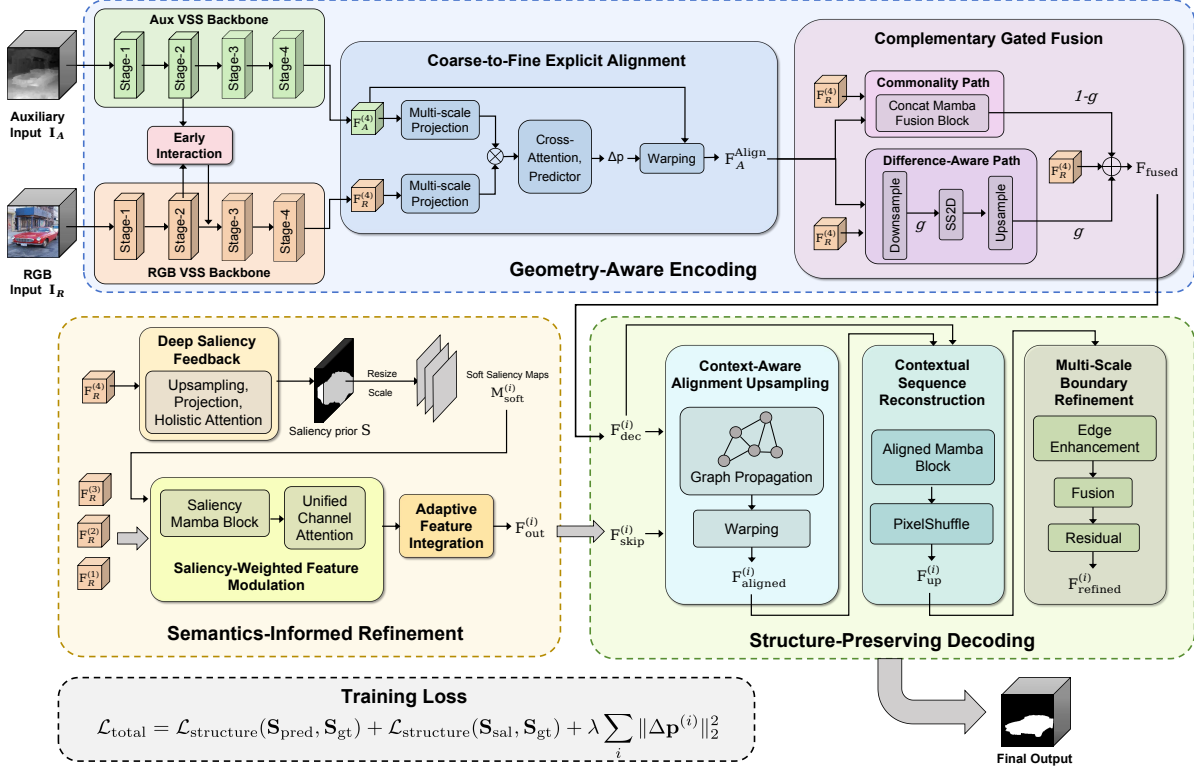


Figure 2: An overview of our proposed G-SalAlignMamba, including Geometry-Aware Encoding for dual-modal features alignment and fusion, Semantics-Informed Refinement to prevent signal dilution, and Structure-Preserving Decoding for geometry-consistent reconstruction.

3.2 Semantics-Informed Refinement: Mitigating Signal Dilution

In general Vision Mamba architectures, the sequential scanning mechanism treats all spatial tokens equally. However, salient objects often occupy a small fraction of the image in SOD tasks. Hence, general methods may lead to *signal dilution*, where strong foreground signals are overwhelmed by background tokens during state transitions. We address this issue via a feedback-driven refinement scheme.

Deep Saliency Feedback Mechanism

We first introduce a top-down feedback loop inspired by cognitive attention. We utilize the deeply fused representation $\mathbf{F}_R^{(4)}$ (from the final encoder stage) to generate a coarse **Saliency Prior** \mathbf{S} to guide shallow and high-resolution features.

Formally, we generate a raw prior \mathbf{S}_{raw} by projecting and upsampling the deep semantic features, followed by spatial smoothing via a Holistic Attention (HA) module to ensure coherence, ultimately yielding the final coarse saliency prior \mathbf{S} . The prior is then resized and transformed into a soft modulation mask via temperature scaling to guide the refinement:

$$\mathbf{S} = \text{HA}(\sigma(\mathbf{S}_{\text{raw}})) \in [0, 1]^{H \times W}, \quad (7)$$

where $\text{HA}(\cdot)$ denotes the Holistic Attention module that applies Gaussian smoothing for spatial coherence. This method leverages high-level semantic context to modulate low-level details, ensuring the enhancement module focuses computational resources on potential foreground regions before features propagate to the decoder.

Saliency-Weighted Feature Modulation

General state-space models' equal token treatment causes *feature dilution* in SOD tasks, where background often overwhelms sparse foregrounds. We propose a **Saliency-Weighted Modulation** strategy that utilizes a saliency-derived mask to modulate scanning, preserving foreground signals without altering sequence length.

For an input feature $\mathbf{F}_R^{(i)}$ at stage $i \in \{1, 2, 3\}$, the saliency prior \mathbf{S} is resized to match the specific spatial resolution (H_i, W_i) , yielding the mask $\mathbf{M}^{(i)}$. The mask is shaped using temperature scaling as:

$$\mathbf{M}_{\text{soft}}^{(i)} = \sigma \left(\frac{\mathbf{M}^{(i)}}{\tau} \right)^\gamma, \quad (8)$$

where τ is the temperature parameter and γ is the power scaling factor. Guided by the soft mask $\mathbf{M}_{\text{soft}}^{(i)}$, the input features $\mathbf{F}_R^{(i)}$ are processed by the Saliency Mamba Block—comprising linear projection, reshape, DWConv, SG-SS2D [He *et al.*, 2025], LN, and a final reshape to yield the refined representations.

The soft modulation ensures that the hidden state of the SSMS is primarily influenced by informative foreground signals, thereby preserving the integrity of salient features with suppressing background noise to yield the prior $\mathbf{F}_{\text{refined}}^{(i)}$. The saliency-weighted modulation adds negligible computational overhead, maintaining inference efficiency comparable to general scanning approaches.

322 Additionally, to enhance channel specificity, we integrate a
 323 unified channel attention module $\mathcal{A}_{\text{channel}}$ (combining SE [Hu
 324 *et al.*, 2018] and ECA mechanisms [Wang *et al.*, 2020])
 325 within the Saliency Mamba Block as:

$$\mathbf{F}_{\text{refined}}^{(i)} = \mathbf{F}_{\text{refined}}^{(i)} + \alpha \cdot \mathcal{A}_{\text{channel}}(\mathbf{F}_{\text{refined}}^{(i)}), \quad (9)$$

326 where α is a learnable balancing parameter, enhancing the
 327 signal-to-noise ratio for precise and complete SOD results.

328 Adaptive Feature Integration

329 Finally, to regulate the information flow to the decoder, we
 330 employ an Adaptive Feature Integration mechanism. The re-
 331 fined features $\mathbf{F}_{\text{refined}}^{(i)}$ are modulated by a learnable gating fac-
 332 tor $\sigma(\gamma_i)$. This method allows the network to dynamically
 333 determine the importance of the refined features, selectively
 334 accepting informative signals with suppressing any residual
 335 artifacts or noise before decoding, ultimately producing the
 336 integrated feature representation $\mathbf{F}_{\text{out}}^{(i)}$.

337 3.3 Structure-Preserving Decoding: Precise 338 Upsampling

339 General decoders often suffer from feature blurring during
 340 upsampling, which is particularly detrimental to SOD, where
 341 precise boundary delineation is critical for high evaluation
 342 metrics. To address this issue, we propose a structure-
 343 preserving scheme that enforces explicit alignment. The de-
 344 coder restores the resolution of the deepest encoder feature
 345 $\mathbf{F}_{\text{fused}}$ to the original scale.

346 Context-Aware Alignment Upsampling

347 In SOD tasks, spatial misalignment in skip connections often
 348 leads to boundary ghosting. The decoding process operates
 349 over three stages $i \in \{1, 2, 3\}$. At each stage, the module si-
 350 multaneously processes the current decoder feature $\mathbf{F}_{\text{dec}}^{(i)}$ and
 351 the skip-connected feature $\mathbf{F}_{\text{skip}}^{(i)}$ from the encoder’s shallow
 352 layers. Specifically, $\mathbf{F}_{\text{dec}}^{(i)} \in \mathbb{R}^{B \times C_i \times H_i \times W_i}$ denotes the fea-
 353 ture from the previous decoding stage (initialized as $\mathbf{F}_{\text{fused}}$ for
 354 $i = 1$). Correspondingly, $\mathbf{F}_{\text{skip}}^{(i)} \in \mathbb{R}^{B \times C_{\text{skip}}^{(i)} \times H_i \times W_i}$ represents
 355 the skip feature derived from encoder stage $(4 - i)$ (i.e., $\mathbf{F}_{\text{out}}^{(3)}$,
 356 $\mathbf{F}_{\text{out}}^{(2)}$, or $\mathbf{F}_{\text{out}}^{(1)}$ for $i = 1, 2, 3$ respectively).

357 To maintain feature correspondence with recovering de-
 358 tails, we introduce a **Context-Aware Alignment Upsam-**
 359 **pling** mechanism. We first project the concatenated decoder
 360 and skip features into node embeddings $\mathbf{N}^{(i)}$ and refine them
 361 via anisotropic graph propagation to capture boundary cues:

$$\mathbf{N}^{(i)} = \text{NodeEmbed}(\text{Concat}(\mathbf{F}_{\text{dec}}^{(i)}, \mathbf{F}_{\text{skip}}^{(i)})), \quad (10)$$

$$\mathbf{N}^{(i)} \leftarrow \mathbf{N}^{(i)} + \beta \cdot \sum_{d \in \mathcal{D}} \mathbf{w}_d \odot (\mathbf{N}_d - \mathbf{N}^{(i)}), \quad (11)$$

363 where β is a learnable diffusion strength, and $\mathcal{D} =$
 364 $\{\text{up, down, left, right}\}$ denotes the set of four cardinal neigh-
 365 bor directions used for anisotropic propagation. The refined
 366 nodes guide the prediction of a dense offset field $\Delta \mathbf{p}^{(i)}$,
 367 which aligns the skip features to the decoder’s coordinate sys-
 368 tem via differentiable warping \mathcal{W} :

$$\Delta \mathbf{p}^{(i)} = \tanh(\mathcal{R}_{\Delta}(\text{Concat}(\mathbf{F}_{\text{concat}}^{(i)}, \mathbf{N}_{\text{prop}}^{(i)}))) \cdot \alpha_{\text{max}}, \quad (12)$$

$$\mathbf{F}_{\text{aligned}}^{(i)} = \mathcal{W}(\mathbf{F}_{\text{skip}}^{(i)}, \text{Grid} + \Delta \mathbf{p}^{(i)}). \quad (13)$$

370 The explicit spatial shift ensures feature correspondence and
 371 preserves object structure during upsampling. Besides, this
 372 rectification guarantees the spatial precision required for
 373 high-quality SOD.

374 Contextual Sequence Reconstruction

375 In the decoding phase, the Contextual Sequence Reconstruc-
 376 tion block (AlignedMambaBlock) plays a vital role in recov-
 377 ering spatial details. It takes the upsampled decoder fea-
 378 ture $\mathbf{F}_{\text{dec}}^{(i)}$ and the explicitly aligned skip-connection feature
 379 $\mathbf{F}_{\text{aligned}}^{(i)}$ as inputs. Under the guidance of the resized saliency
 380 prior \mathbf{P}_{sal} , these features are fused to model global context
 381 and dependencies. The reconstructed sequence is then ex-
 382 panded and upsampled via a PixelShuffle operation to gener-
 383 ate the high-resolution feature $\mathbf{F}_{\text{up}}^{(i)}$ for the subsequent stage,
 384 aiming to preserve the semantic completeness in SOD tasks.

385 Multi-Scale Boundary Refinement

386 Preserving high-frequency details is critical for distinguish-
 387 ing salient foregrounds in SOD tasks. To avoid bound-
 388 ary degradation, we employ an unsupervised Multi-Scale
 389 Boundary Refinement module that implicitly sharpens high-
 390 frequency details without requiring explicit edge supervision.
 391 The module operates by extracting local boundary cues from
 392 $\mathbf{F}_{\text{up}}^{(i)}$ using parallel convolutional branches at three distinct
 393 scales ($k \in \{1, 2, 4\}$). These multi-scale boundary features
 394 are subsequently upsampled, concatenated, and fused back
 395 into the main feature stream via a residual connection. This
 396 process yields the final refined output $\mathbf{F}_{\text{refined}}^{(i)}$, characterized
 397 by sharper edges and clearer structural delineations.

398 Loss Functions

399 To stabilize training, we apply L2 regularization on the offset
 400 fields. The total loss combines the structure loss $\mathcal{L}_{\text{structure}}$ for
 401 both main and auxiliary outputs with the offset regularization:

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{structure}}(\mathbf{S}_{\text{pred}}, \mathbf{S}_{\text{gt}}) + \mathcal{L}_{\text{structure}}(\mathbf{S}_{\text{sal}}, \mathbf{S}_{\text{gt}}) + \lambda \sum_i \|\Delta \mathbf{p}^{(i)}\|_2^2, \quad (14)$$

402 where no explicit boundary supervision is required. This
 403 composite loss drives structural accuracy, ensuring robust
 404 convergence for SOD tasks.

405 4 Experiment

406 4.1 Datasets and Metrics

407 For RGB-D SOD, we use five benchmark datasets: NJUD
 408 [Ju *et al.*, 2014], NLPR [Peng *et al.*, 2014], SIP [Fan *et al.*,
 409 2020a], STERE [Niu *et al.*, 2012] and DUTLF-D [Piao *et al.*,
 410 2019]. For RGB-T SOD, we adopt seven datasets: VT5000
 411 [Tu *et al.*, 2022b], VT1000 [Tu *et al.*, 2019], VT821 [Wang
 412 *et al.*, 2018], un-VT5000 [Tu *et al.*, 2022a], un-VT1000 [Tu
 413 *et al.*, 2022a], un-VT821 [Tu *et al.*, 2022a] and UVT2000
 414 [Wang *et al.*, 2024]. For RGB-D and RGB-T SOD, we em-
 415 ploy three standard metrics: **E-measure** [Fan *et al.*, 2018]
 416 (E_m) jointly assesses pixel-level accuracy and image-level

Table 1: Quantitative comparison of our method against other recent advanced RGB-D SOD methods on five benchmark datasets. “-” indicates the result is not available. “↑” denotes that the larger value is better. The best two results are stressed in red and blue.

Method	Type	NJUD			NLPR			SIP			STERE			DUTLF-D			Params (M)	FPS
		$S_m \uparrow$	$F_m \uparrow$	$E_m \uparrow$	$S_m \uparrow$	$F_m \uparrow$	$E_m \uparrow$	$S_m \uparrow$	$F_m \uparrow$	$E_m \uparrow$	$S_m \uparrow$	$F_m \uparrow$	$E_m \uparrow$	$S_m \uparrow$	$F_m \uparrow$	$E_m \uparrow$		
BBSNet [Fan <i>et al.</i> , 2020b]	CNN	0.921	0.919	0.949	0.931	0.918	0.961	0.879	0.884	0.922	0.908	0.903	0.942	0.882	0.870	0.912	49.77	45.21
JL-DCF [Fu <i>et al.</i> , 2020]		0.877	0.892	0.941	0.931	0.918	0.965	0.885	0.894	0.931	0.900	0.895	0.942	0.894	0.891	0.927	143.52	22.15
SP-Net [Zhou <i>et al.</i> , 2021]		0.925	0.928	0.957	0.927	0.919	0.962	0.894	0.904	0.933	0.907	0.906	0.949	0.895	0.899	0.933	67.88	47.38
DCF [Ji <i>et al.</i> , 2021]		0.904	0.905	0.943	0.922	0.910	0.957	0.874	0.886	0.922	0.906	0.904	0.948	0.925	0.930	0.956	53.92	32.64
SPSN [Lee <i>et al.</i> , 2022]		0.918	0.921	0.952	0.923	0.912	0.960	0.892	0.900	0.936	0.907	0.902	0.945	-	-	-	-	-
SwinNet-B [Liu <i>et al.</i> , 2021]	Trans.	0.920	0.924	0.956	0.941	0.936	0.974	0.911	0.927	0.950	0.919	0.918	0.956	0.918	0.920	0.949	199.18	11.42
CATNet [Sun <i>et al.</i> , 2023]		0.932	0.937	0.960	0.938	0.934	0.971	0.910	0.928	0.951	0.920	0.922	0.958	0.952	0.958	0.975	262.73	23.09
VST-S++ [Liu <i>et al.</i> , 2024a]		0.928	0.928	0.957	0.935	0.925	0.964	0.904	0.918	0.946	0.921	0.916	0.954	0.945	0.950	0.969	143.15	27.85
CPNet [Hu <i>et al.</i> , 2024]		0.935	0.941	0.963	0.940	0.936	0.971	0.907	0.927	0.946	0.920	0.922	0.960	0.951	0.959	0.974	216.50	20.17
VSCoDe-S [Luo <i>et al.</i> , 2024]		0.944	0.949	0.970	0.941	0.932	0.968	0.924	0.942	0.958	0.931	0.928	0.958	0.960	0.967	0.980	74.72	30.56
DiMSOD [Zhang <i>et al.</i> , 2025]	Diff.	0.947	0.947	0.969	-	-	-	-	-	-	-	-	-	0.957	0.967	0.951	-	-
G-AlignSalMamba (Ours)	Mamba	0.956	0.948	0.983	0.947	0.932	0.978	0.929	0.941	0.955	0.935	0.932	0.958	0.977	0.954	0.981	83.80	30.41

Table 2: Quantitative comparison of our method against other recent advanced RGB-T SOD methods on seven benchmark datasets. “-” indicates the result is not available. “↑” denotes that the larger value is better. The best two results are stressed in red and blue.

Method	Type	VT5000			VT1000			VT821			un-VT5000			un-VT1000			un-VT821			UVT2000			Params (M)	FPS
		$S_m \uparrow$	$F_m \uparrow$	$E_m \uparrow$	$S_m \uparrow$	$F_m \uparrow$	$E_m \uparrow$	$S_m \uparrow$	$F_m \uparrow$	$E_m \uparrow$	$S_m \uparrow$	$F_m \uparrow$	$E_m \uparrow$	$S_m \uparrow$	$F_m \uparrow$	$E_m \uparrow$	$S_m \uparrow$	$F_m \uparrow$	$E_m \uparrow$	$S_m \uparrow$	$F_m \uparrow$	$E_m \uparrow$		
CSRNet [Huo <i>et al.</i> , 2021]	CNN	0.868	0.811	0.905	0.918	0.877	0.925	0.884	0.831	0.909	0.642	0.475	0.713	0.705	0.582	0.732	0.737	0.619	0.783	0.655	0.420	0.658	1.03	40.51
CGFNet [Wang <i>et al.</i> , 2021]		0.883	0.851	0.922	0.923	0.906	0.944	0.881	0.845	0.912	0.833	0.757	0.868	0.877	0.817	0.878	0.837	0.776	0.874	0.764	0.539	0.704	18.37	29.73
DCNet [Zhu <i>et al.</i> , 2025]		0.895	0.847	0.920	0.929	0.911	0.948	0.899	0.842	0.912	0.858	0.803	0.879	0.858	0.850	0.880	0.817	0.793	0.869	0.767	0.632	0.808	24.15	14.77
CAVER [Pang <i>et al.</i> , 2023]		0.892	0.841	0.924	0.936	0.903	0.945	0.891	0.839	0.924	0.850	0.805	0.893	0.873	0.838	0.881	0.818	0.780	0.856	0.786	0.616	0.782	93.84	16.32
SPNet [Zhang <i>et al.</i> , 2023]		0.914	0.880	0.948	0.941	0.925	0.954	0.913	0.873	0.936	0.900	0.848	0.929	0.931	0.902	0.938	0.894	0.833	0.910	0.765	0.558	0.733	113.59	29.93
SwinNet [Liu <i>et al.</i> , 2021]	Trans.	0.912	0.865	0.942	0.938	0.896	0.947	0.904	0.847	0.926	0.837	0.767	0.901	0.853	0.802	0.871	0.854	0.783	0.903	0.790	0.592	0.780	198.74	26.19
LSNet [Liu <i>et al.</i> , 2022]		0.877	0.825	0.915	0.925	0.885	0.935	0.878	0.825	0.911	0.847	0.767	0.892	0.868	0.797	0.875	0.842	0.754	0.875	0.763	0.527	0.711	4.64	88.83
SACNet [Wang <i>et al.</i> , 2024]		0.917	0.901	0.957	0.942	0.923	0.958	0.906	0.868	0.932	0.872	0.780	0.899	0.852	0.803	0.868	0.876	0.812	0.916	0.795	0.601	0.792	530.90	19.00
PCNet [Wang <i>et al.</i> , 2025]		0.920	0.899	0.956	0.943	0.924	0.958	0.915	0.879	0.941	0.879	0.861	0.936	0.922	0.904	0.947	0.893	0.869	0.936	0.819	0.686	0.851	-	-
DiMSOD [Zhang <i>et al.</i> , 2025]		Diff.	0.921	0.898	0.959	0.953	0.935	0.955	0.923	0.917	0.949	-	-	-	-	-	-	-	-	-	-	-	-	-
G-AlignSalMamba (Ours)	Mamba	0.939	0.913	0.956	0.941	0.939	0.961	0.921	0.919	0.940	0.926	0.849	0.931	0.935	0.910	0.913	0.871	0.909	0.930	0.905	0.704	0.788	83.80	30.41

417 consistency; **S-measure** [Fan *et al.*, 2017] (S_m) evaluates
 418 structural similarity; and **F-measure** [Achanta *et al.*, 2009]
 419 (F_m) balances precision and recall.

4.2 Implementation Details

420 Our method is implemented in PyTorch and trained on an
 421 NVIDIA GeForce RTX 4090 GPU. We train separate mod-
 422 els for the RGB-D SOD and RGB-T SOD tasks, respectively.
 423 Following previous works, we arrange the training sets for
 424 each task as follows: the training sets of NJUD, NLPR, and
 425 DUTLF-D for RGB-D SOD; the training sets of VT5000 and
 426 UVT2000 for RGB-T SOD.
 427

428 For the specific hyperparameter settings, we initialize the
 429 learnable gating factors γ_{early} and γ_{gate} to -2.0 , while the
 430 adaptive integration factor γ_i is initialized to 2.0 . The maxi-
 431 mum offset for explicit alignment is bounded by $\alpha_{\text{max}} = 3.0$.
 432 Both the temperature parameter τ and the power scaling fac-
 433 tor γ are set to 1.0 . The weight λ in the loss function is set to
 434 0.01 .

435 In the training process, we adopt the AdamW optimizer
 436 with an initial learning rate of 1×10^{-4} and a batch size of
 437 2. All input images are uniformly resized to 448×448 for
 438 training and testing. The model converges after 30 training
 439 epochs.

4.3 Model Comparison

440 **Quantitative Evaluation.** To demonstrate the universality of
 441 our unified framework in dual-modal SOD, we benchmark it
 442 against a comprehensive suite of 11 RGB-D and 10 RGB-T
 443 state-of-the-art methods, as shown in Table 1 and Table 2.
 444 Overall, our framework achieves consistently competitive or
 445 superior performance across all three metrics.
 446

447 For RGB-D (Table 1), our model gives the best S_m scores
 448 across all five datasets and consistently ranks top-2 in F_m and
 449 E_m . It significantly outperforms Transformer-based (e.g.,
 450 CATNet) and Diffusion-based (DiMSOD) counterparts, par-
 451 ticularly on challenging benchmarks (NJUD and DUTLF-D).
 452

453 For RGB-T (Table 2), our method also obtains state-of-
 454 the-art results on large-scale benchmarks. Specifically, it
 455 ranks first in both S_m and F_m on VT5000, un-VT1000, and
 456 UVT2000. Even under extreme unconstrained conditions
 457 (e.g., un-VT5000), it maintains a leading position with top-
 458 tier boundary accuracy (F_m) and structural consistency (E_m).
 459

460 Benefiting from Mamba’s linear complexity $\mathcal{O}(N)$, G-
 461 AlignSalMamba strikes a superior balance between perfor-
 462 mance and cost. It requires only 83.80M parameters, repre-
 463 senting a substantial reduction of approximately 61% to 68%
 464 compared to CPNet and CATNet. Furthermore, it achieves
 465 real-time inference at 30.41 FPS, significantly outpacing
 466 heavy Transformer baselines such as SwinNet-B (11.42 FPS).

Table 3: Ablation study of the proposed framework. The best results are highlighted in **bold**.

Idx	Encoder			Refinement			Decoder			NJUD			NLPR			VT5000			UVT2000		
	EI	CFEA	CGF	DSFM	SWM	AFI	CAAU	CSR	MBR	$S_m \uparrow$	$F_m \uparrow$	$E_m \uparrow$	$S_m \uparrow$	$F_m \uparrow$	$E_m \uparrow$	$S_m \uparrow$	$F_m \uparrow$	$E_m \uparrow$	$S_m \uparrow$	$F_m \uparrow$	$E_m \uparrow$
Base										0.852	0.903	0.911	0.842	0.892	0.909	0.850	0.835	0.910	0.847	0.849	0.909
E1	✓									0.864	0.913	0.915	0.858	0.902	0.911	0.856	0.847	0.911	0.849	0.856	0.917
E2		✓								0.869	0.917	0.917	0.866	0.903	0.912	0.860	0.852	0.914	0.862	0.857	0.921
E3			✓							0.879	0.921	0.921	0.888	0.905	0.916	0.877	0.859	0.916	0.880	0.867	0.925
E4	✓	✓								0.883	0.922	0.922	0.906	0.910	0.919	0.879	0.860	0.915	0.882	0.868	0.928
E5	✓	✓	✓							0.903	0.923	0.931	0.909	0.917	0.930	0.889	0.864	0.917	0.892	0.875	0.930
H1	✓	✓	✓	✓						0.905	0.924	0.933	0.909	0.917	0.931	0.895	0.870	0.918	0.893	0.876	0.932
H2	✓	✓	✓		✓					0.907	0.925	0.936	0.910	0.918	0.933	0.901	0.879	0.920	0.894	0.878	0.935
H3	✓	✓	✓			✓				0.911	0.922	0.932	0.910	0.920	0.930	0.911	0.901	0.921	0.891	0.880	0.931
H4	✓	✓	✓	✓	✓	✓				0.910	0.927	0.940	0.919	0.919	0.935	0.912	0.901	0.922	0.895	0.881	0.936
H5	✓	✓	✓	✓	✓	✓	✓			0.921	0.931	0.950	0.927	0.929	0.943	0.914	0.909	0.941	0.908	0.886	0.933
D1	✓	✓	✓	✓	✓	✓	✓			0.935	0.937	0.958	0.930	0.922	0.940	0.921	0.903	0.940	0.901	0.889	0.939
D2	✓	✓	✓	✓	✓	✓		✓		0.936	0.935	0.958	0.933	0.924	0.945	0.926	0.906	0.945	0.905	0.891	0.940
D3	✓	✓	✓	✓	✓	✓			✓	0.938	0.934	0.957	0.935	0.927	0.950	0.931	0.909	0.951	0.908	0.892	0.940
D4	✓	✓	✓	✓	✓	✓	✓	✓		0.939	0.936	0.959	0.938	0.929	0.951	0.936	0.912	0.955	0.912	0.894	0.941
D5	✓	✓	✓	✓	✓	✓	✓	✓	✓	0.955	0.945	0.963	0.947	0.930	0.975	0.939	0.923	0.960	0.915	0.902	0.965

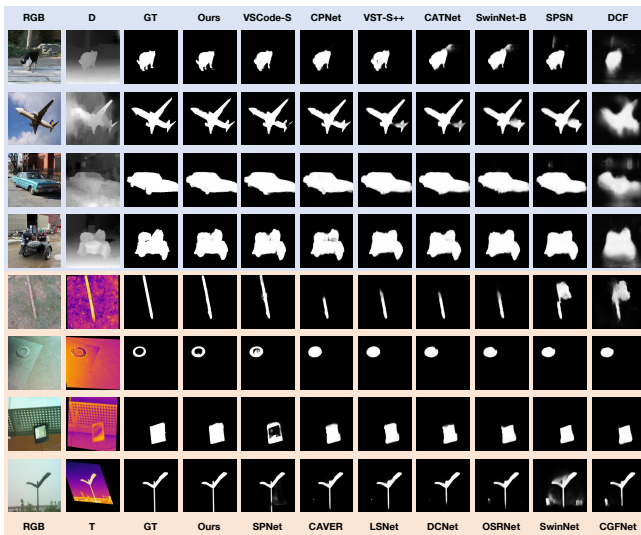


Figure 3: Visual comparison against other RGB-D SOD and RGB-T SOD methods on the NJUD and UVT5000 dataset.

Qualitative Evaluation. Figure 3 compares our method with state-of-the-art approaches on RGB-D (top) and RGB-T (bottom) benchmarks. In RGB-D scenarios characterized by complex topologies (e.g., 2nd row) or noisy depth (e.g., 3rd row), our proposed G-SalAlignMamba produces sharper boundaries and clearer object delineation than counterparts. Similarly, in RGB-T cases featuring thin instances (5th, 8th rows) and small objects (6th row), our model effectively suppresses artifacts and preserves fine-grained structural details, maintaining robustness against thermal degradation where other methods often fail.

4.4 Ablation Study

We perform a comprehensive ablation study to evaluate the individual and combined contributions of each component in the proposed G-SalAlignMamba framework on NJUD, NLPR, VT5000, and UVT2000 datasets.

Geometry-Aware Encoding. We first construct a strong Baseline (Base) to validate the effectiveness of our proposed

modules. This baseline utilizes the pre-trained VMamba backbone with simple channel concatenation for multi-modal fusion and general bilinear interpolation for upsampling without any of our geometry-aware or refinement modules. Based on the above, we investigate the encoder components: *Early Interaction* (EI), *Coarse-to-Fine Explicit Alignment* (CFEA), and *Complementary Gated Fusion* (CGF). As shown in Table 3, we evaluate each module independently (E1–E3) to verify their standalone effectiveness relative to the baseline. We then test pairwise combinations (E4–E5), culminating in the full encoder (E6).

Semantics-Informed Refinement. For the refinement stage, we decouple the *Deep Saliency Feedback Mechanism* (DSFM) from the *Saliency-Weighted Modulation* (SWM) and evaluate them alongside the *Adaptive Feature Integration* (AFI). We report the performance of deploying these modules individually (H1–H3) and the gain from the feedback-modulation synergy (H4), demonstrating the necessity of the complete refinement scheme (H5).

Structure-Preserving Decoding. Finally, we examine the decoder by progressively adding *Context-Aware Alignment Upsampling* (CAAU), *Contextual Sequence Reconstruction* (CSR), and *Multi-Scale Boundary Refinement* (MBR). We assess the impact of each component in isolation (D1–D3) and their stepwise integration (D4–D5), confirming that the full decoder yields the best spatial alignment and boundary quality compared to general methods.

5 Conclusion

In this paper, we proposed G-SalAlignMamba, a geometry-aware framework for dual-modal SOD that targets three key issues: *Misalignment Vulnerability*, *Signal Dilution*, and *Boundary Degradation*. Despite the promising performance on RGB-D and RGB-T benchmarks, our proposed G-SalAlignMamba may still encounter challenges in scenarios with extreme occlusion or highly cluttered backgrounds, where the target features are severely obstructed in both modalities. Future work will focus on incorporating stronger semantic priors to improve robustness in these edge cases.

References

- [Achanta *et al.*, 2009] Radhakrishna Achanta, Sheila Hemami, Francisco Estrada, and Sabine Susstrunk. Frequency-tuned salient region detection. In *2009 IEEE conference on computer vision and pattern recognition*, pages 1597–1604. IEEE, 2009.
- [Fan *et al.*, 2017] Deng-Ping Fan, Ming-Ming Cheng, Yun Liu, Tao Li, and Ali Borji. Structure-measure: A new way to evaluate foreground maps. In *Proceedings of the IEEE international conference on computer vision*, pages 4548–4557, 2017.
- [Fan *et al.*, 2018] Deng-Ping Fan, Cheng Gong, Yang Cao, Bo Ren, Ming-Ming Cheng, and Ali Borji. Enhanced-alignment measure for binary foreground map evaluation. *arXiv preprint arXiv:1805.10421*, 2018.
- [Fan *et al.*, 2020a] Deng-Ping Fan, Zheng Lin, Zhao Zhang, Menglong Zhu, and Ming-Ming Cheng. Rethinking rgb-d salient object detection: Models, data sets, and large-scale benchmarks. *IEEE Transactions on neural networks and learning systems*, 32(5):2075–2089, 2020.
- [Fan *et al.*, 2020b] Deng-Ping Fan, Yingjie Zhai, Ali Borji, Jufeng Yang, and Ling Shao. Bbs-net: Rgb-d salient object detection with a bifurcated backbone strategy network. In *European conference on computer vision*, pages 275–292. Springer, 2020.
- [Fu *et al.*, 2020] Keren Fu, Deng-Ping Fan, Ge-Peng Ji, and Qijun Zhao. Jl-dcf: Joint learning and densely-cooperative fusion framework for rgb-d salient object detection. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 3052–3062, 2020.
- [Fu *et al.*, 2021] Keren Fu, Deng-Ping Fan, Ge-Peng Ji, Qijun Zhao, Jianbing Shen, and Ce Zhu. Siamese network for rgb-d salient object detection and beyond. *IEEE transactions on pattern analysis and machine intelligence*, 44(9):5541–5559, 2021.
- [Gu and Dao, 2024] Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. In *First conference on language modeling*, 2024.
- [He *et al.*, 2025] Jiahao He, Keren Fu, Xiaohong Liu, and Qijun Zhao. Samba: A unified mamba-based framework for general salient object detection. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 25314–25324, 2025.
- [Hu *et al.*, 2018] Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7132–7141, 2018.
- [Hu *et al.*, 2024] Xihang Hu, Fuming Sun, Jing Sun, Fasheng Wang, and Haojie Li. Cross-modal fusion and progressive decoding network for rgb-d salient object detection. *International Journal of Computer Vision*, 132(8):3067–3085, 2024.
- [Huang *et al.*, 2021] Shihua Huang, Zhichao Lu, Ran Cheng, and Cheng He. Fapn: Feature-aligned pyramid network for dense image prediction. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 864–873, 2021.
- [Huo *et al.*, 2021] Fushuo Huo, Xuegui Zhu, Lei Zhang, Qifeng Liu, and Yu Shu. Efficient context-guided stacked refinement network for rgb-t salient object detection. *IEEE Transactions on Circuits and Systems for Video Technology*, 32(5):3111–3124, 2021.
- [Ji *et al.*, 2021] Wei Ji, Jingjing Li, Shuang Yu, Miao Zhang, Yongri Piao, Shunyu Yao, Qi Bi, Kai Ma, Yefeng Zheng, Huchuan Lu, et al. Calibrated rgb-d salient object detection. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9471–9481, 2021.
- [Ju *et al.*, 2014] Ran Ju, Ling Ge, Wenjing Geng, Tongwei Ren, and Gangshan Wu. Depth saliency based on anisotropic center-surround difference. In *2014 IEEE international conference on image processing (ICIP)*, pages 1115–1119. IEEE, 2014.
- [Lee *et al.*, 2022] Minhyeok Lee, Chaewon Park, Suhwan Cho, and Sangyoun Lee. Spsn: Superpixel prototype sampling network for rgb-d salient object detection. In *European conference on computer vision*, pages 630–647. Springer, 2022.
- [Liu *et al.*, 2021] Zhengyi Liu, Yacheng Tan, Qian He, and Yun Xiao. Swinnet: Swin transformer drives edge-aware rgb-d and rgb-t salient object detection. *IEEE Transactions on Circuits and Systems for Video Technology*, 32(7):4486–4497, 2021.
- [Liu *et al.*, 2022] Biyuan Liu, Huaixin Chen, Zhixi Wang, Wenqiang Xie, and Lingyu Shuai. Lsnet: Extremely lightweight siamese network for change detection of remote sensing image. In *IGARSS 2022-2022 IEEE International Geoscience and Remote Sensing Symposium*, pages 2358–2361. IEEE, 2022.
- [Liu *et al.*, 2024a] Nian Liu, Ziyang Luo, Ni Zhang, and Junwei Han. Vst++: Efficient and stronger visual saliency transformer. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 46(11):7300–7316, 2024.
- [Liu *et al.*, 2024b] Yue Liu, Yunjie Tian, Yuzhong Zhao, Hongtian Yu, Lingxi Xie, Yaowei Wang, Qixiang Ye, Jianbin Jiao, and Yunfan Liu. Vmamba: Visual state space model. *Advances in neural information processing systems*, 37:103031–103063, 2024.
- [Luo *et al.*, 2024] Ziyang Luo, Nian Liu, Wangbo Zhao, Xuguang Yang, Dingwen Zhang, Deng-Ping Fan, Fahad Khan, and Junwei Han. Vscope: General visual salient and camouflaged object detection with 2d prompt learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 17169–17180, 2024.
- [Niu *et al.*, 2012] Yuzhen Niu, Yujie Geng, Xueqing Li, and Feng Liu. Leveraging stereopsis for saliency analysis. In *2012 IEEE conference on computer vision and pattern recognition*, pages 454–461. IEEE, 2012.
- [Pang *et al.*, 2023] Youwei Pang, Xiaoqi Zhao, Lihe Zhang, and Huchuan Lu. Caver: Cross-modal view-mixed trans-

- 631 former for bi-modal salient object detection. *IEEE Transactions on Image Processing*, 32:892–904, 2023.
- 632
- 633 [Peng *et al.*, 2014] Houwen Peng, Bing Li, Weihua Xiong,
634 Weiming Hu, and Rongrong Ji. Rgb-d salient object
635 detection: A benchmark and algorithms. In *European conference on computer vision*, pages 92–109. Springer, 2014.
- 636
- 637 [Piao *et al.*, 2019] Yongri Piao, Wei Ji, Jingjing Li, Miao
638 Zhang, and Huchuan Lu. Depth-induced multi-scale recur-
639 rent attention network for saliency detection. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 7254–7263, 2019.
- 640
- 641
- 642 [Qin *et al.*, 2019] Xuebin Qin, Zichen Zhang, Chenyang
643 Huang, Chao Gao, Masood Dehghan, and Martin Jagersand.
644 Basnet: Boundary-aware salient object detection.
645 In *Proceedings of the IEEE/CVF conference on computer
646 vision and pattern recognition*, pages 7479–7489, 2019.
- 647
- 648 [Sun *et al.*, 2023] Fuming Sun, Peng Ren, Bowen Yin,
649 Fasheng Wang, and Haojie Li. Catnet: A cascaded and ag-
650 gregated transformer network for rgb-d salient object
651 detection. *IEEE Transactions on Multimedia*, 26:2249–2262,
2023.
- 652
- 653 [Tu *et al.*, 2019] Zhengzheng Tu, Tian Xia, Chenglong Li,
654 Xiaoxiao Wang, Yan Ma, and Jin Tang. Rgb-t image
655 saliency detection via collaborative graph learning. *IEEE
Transactions on Multimedia*, 22(1):160–173, 2019.
- 656
- 657 [Tu *et al.*, 2022a] Zhengzheng Tu, Zhun Li, Chenglong Li,
658 and Jin Tang. Weakly alignment-free rgb-t salient object
659 detection with deep correlation network. *IEEE Transac-
tions on Image Processing*, 31:3752–3764, 2022.
- 660
- 661 [Tu *et al.*, 2022b] Zhengzheng Tu, Yan Ma, Zhun Li, Chen-
662 glong Li, Jieming Xu, and Yongtao Liu. Rgb-t salient ob-
663 ject detection: A large-scale dataset and benchmark. *IEEE
Transactions on Multimedia*, 25:4163–4176, 2022.
- 664
- 665 [Wang *et al.*, 2018] Guizhao Wang, Chenglong Li, Yunpeng
666 Ma, Aihua Zheng, Jin Tang, and Bin Luo. Rgb-t saliency
667 detection benchmark: Dataset, baselines, analysis and a
668 novel approach. In *Chinese Conference on Image and
Graphics Technologies*, pages 359–369. Springer, 2018.
- 669
- 670 [Wang *et al.*, 2020] Qilong Wang, Banggu Wu, Pengfei Zhu,
671 Peihua Li, Wangmeng Zuo, and Qinghua Hu. Eca-net:
672 Efficient channel attention for deep convolutional neural
673 networks. In *Proceedings of the IEEE/CVF conference
674 on computer vision and pattern recognition*, pages 11534–
11542, 2020.
- 675
- 676 [Wang *et al.*, 2021] Jie Wang, Kechen Song, Yanqi Bao,
677 Liming Huang, and Yunhui Yan. Cgfnnet: Cross-guided
678 fusion network for rgb-t salient object detection. *IEEE
Transactions on Circuits and Systems for Video Technol-
679 ogy*, 32(5):2949–2961, 2021.
- 680
- 681 [Wang *et al.*, 2024] Kunpeng Wang, Danying Lin, Cheng-
682 long Li, Zhengzheng Tu, and Bin Luo. Alignment-free
683 rgb-t salient object detection: Semantics-guided asymmet-
684 ric correlation network and a unified benchmark. *IEEE
Transactions on Multimedia*, 26:10692–10707, 2024.
- [Wang *et al.*, 2025] Kunpeng Wang, Keke Chen, Chenglong
685 Li, Zhengzheng Tu, and Bin Luo. Alignment-free rgb-t
686 salient object detection: A large-scale dataset and progres-
687 sive correlation network. In *Proceedings of the AAAI Con-
688 ference on Artificial Intelligence*, volume 39, pages 7780–
689 7788, 2025.
- [Xing *et al.*, 2024] Zhaohu Xing, Tian Ye, Yijun Yang,
691 Guang Liu, and Lei Zhu. Segmamba: Long-range sequen-
692 tial modeling mamba for 3d medical image segmentation.
693 In *International conference on medical image comput-
694 ing and computer-assisted intervention*, pages 578–588.
695 Springer, 2024.
- [Zhang *et al.*, 2023] Zihao Zhang, Jie Wang, and Yahong
697 Han. Saliency prototype for rgb-d and rgb-t salient object
698 detection. In *Proceedings of the 31st ACM international
699 conference on multimedia*, pages 3696–3705, 2023.
- [Zhang *et al.*, 2025] Shuo Zhang, Jiaming Huang, Wenbing
701 Tang, Yan Wu, Terrence Hu, Xiaogang Xu, and Jing
702 Liu. Dimsod: A diffusion-based framework for multi-
703 modal salient object detection. In *Proceedings of the AAAI
704 Conference on Artificial Intelligence*, volume 39, pages
705 10103–10111, 2025.
- [Zhao *et al.*, 2019] Jia-Xing Zhao, Jiang-Jiang Liu, Deng-
707 Ping Fan, Yang Cao, Jufeng Yang, and Ming-Ming Cheng.
708 Egnnet: Edge guidance network for salient object detection.
709 In *Proceedings of the IEEE/CVF international conference
710 on computer vision*, pages 8779–8788, 2019.
- [Zhou *et al.*, 2021] Tao Zhou, Huazhu Fu, Geng Chen,
712 Yi Zhou, Deng-Ping Fan, and Ling Shao. Specificity-
713 preserving rgb-d saliency detection. In *Proceedings of the
714 IEEE/CVF international conference on computer vision*,
715 pages 4681–4691, 2021.
- [Zhou *et al.*, 2023] Heng Zhou, Chunna Tian, Zhenxi Zhang,
717 Chengyang Li, Yuxuan Ding, Yongqiang Xie, and
718 Zhongbo Li. Position-aware relation learning for rgb-
719 thermal salient object detection. *IEEE Transactions on
720 Image Processing*, 32:2593–2607, 2023.
- [Zhu *et al.*, 2024] Lianghui Zhu, Bencheng Liao, Qian
722 Zhang, Xinlong Wang, Wenyu Liu, and Xinggong Wang.
723 Vision mamba: Efficient visual representation learning
724 with bidirectional state space model. *arXiv preprint
725 arXiv:2401.09417*, 2024.
- [Zhu *et al.*, 2025] Jiayi Zhu, Xuebin Qin, and Abdulmoteleb
727 Elsaddik. Dc-net: Divide-and-conquer for salient object
728 detection. *Pattern Recognition*, 157:110903, 2025.
- 729