# SUM: Saliency Unification through Mamba for Visual Attention Modeling

Alireza Hosseini[*,1]    Amirhossein Kazerouni[*,2,3,4]    Saeed Akhavan [1]
Michael Brudno [2,3,4]    Babak Taati [2,3,4]

[1] University of Tehran [2] University of Toronto [3] Vector Institute
[4] University Health Network

{arhosseini77, s.akhavan}@ut.ac.ir, {amirhossein, brudno}@cs.toronto.edu
babak.taati@uhn.ca

## Abstract

*Visual attention modeling, important for interpreting and prioritizing visual stimuli, plays a significant role in applications such as marketing, multimedia, and robotics. Traditional saliency prediction models, especially those based on Convolutional Neural Networks (CNNs) or Transformers, achieve notable success by leveraging large-scale annotated datasets. However, the current state-of-the-art (SOTA) models that use Transformers are computationally expensive. Additionally, separate models are often required for each image type, lacking a unified approach. In this paper, we propose* **S***aliency* **U***nification through* **M***amba* **(SUM)***, a novel approach that integrates the efficient long-range dependency modeling of Mamba with U-Net to provide a unified model for diverse image types. Using a novel Conditional Visual State Space (C-VSS) block, SUM dynamically adapts to various image types, including natural scenes, web pages, and commercial imagery, ensuring universal applicability across different data types. Our comprehensive evaluations across five benchmarks demonstrate that SUM seamlessly adapts to different visual characteristics and consistently outperforms existing models. These results position SUM as a versatile and powerful tool for advancing visual attention modeling, offering a robust solution universally applicable across different types of visual content. Our code and pretrained models are available at* https://github.com/Arhosseini77/SUM.

## 1. Introduction

Visual attention is a critical function of the human visual system, enabling the selection of the most relevant information in a visual scene [32]. Modeling of this mechanism, known as saliency prediction, plays pivotal roles in numerous applications such as marketing [22, 31], multi-

media [50], computer vision [52], and robotics [9].

Deep learning models have succeeded in saliency prediction, by exploiting large-scale annotated datasets [3, 30]. Typically, these models employ a pre-trained object recognition network for feature extraction [38], with the U-Net architecture as a popular choice. Most methods employ Convolutional Neural Networks (CNNs) to construct encoders and decoders for latent features, which generate visual saliency maps [6, 14, 23, 26, 36, 62]. Recurrent architectures, such as Long-Short Term Memory (LSTM) networks, are also sometimes used to model both local and long-range visual information [11, 46], enhancing the accuracy of saliency predictions. More recently, the use of Transformer-based models has led to significant improvements, achieving SOTA performance in saliency prediction by learning spatial long-range dependencies [13, 20, 22, 41, 48]. However, the computational demands of the standard self-attention mechanism used in these methods, which scales quadratically with image size, present a substantial challenge, especially for dense prediction tasks like saliency modeling.

Moreover, a significant limitation within current saliency prediction models lies in their design specificity for singular visual contexts. Saliency maps, and consequently the models that generate them, need to be adapted to the unique characteristics of different types of images. For example, in natural scenes, the visual attention of viewers may be driven largely by elements like color and movement, whereas in e-commerce images, textual information typically attracts more attention [29]. Similarly, in user interface (UI) designs, the upper-left quadrant often attracts more attention due to common eye movement patterns and a left-to-right viewing bias [31]. Although there are robust models tailored for specific datasets, such as those optimized for commercial imagery [22] or UIs [31], the research on the development of universally applicable models, that can effectively handle diverse requirements of various image types, remains limited. This gap underscores the necessity for a

---

model, which can be universally performant across all image types and saliency datasets, thus providing a more comprehensive solution to the field of saliency prediction.

To address the challenges outlined above, we leverage the capabilities of State Space Models (SSMs) [34] as used in Mamba [18, 47], and introduce a novel unified Mamba-U-Net-based model for visual saliency prediction. Models like Mamba capture long-distance dependencies with linear computational complexity. Inspired by these successes, we propose the **S**aliency **U**nification through **M**amba (**SUM**), which uses Mamba to efficiently capture long-range information. To ensure universal applicability across diverse image types, we incorporate a novel Conditional Visual State Space (C-VSS) block in our design. This component effectively separates the distributions of different data types, making the model robust across various modalities. It allows SUM to dynamically adapt to distinct visual characteristics found in natural scenes, e-commerce imagery, and UIs. Validation of SUM on six large-scale datasets across different visual contexts confirms its exceptional adaptability and strong performance, positioning it as a potent tool in the advancement of visual attention modeling. These attributes make SUM a valuable tool for a range of applications in visual saliency prediction. The **main contributions** of this work are summarized as follows:

- A novel efficient class-conditional unified model is proposed that employs Mamba to capture long-range visual information efficiently with linear computational complexity.
- A conditional component that dynamically adapts the model behavior at test time through the shift and scaling mechanisms, enhancing the adaptability of the model to various visual contexts.
- SUM is extensively evaluated on six diverse benchmark datasets, including natural scenes with gaze or mouse ground truth labels, web pages, and commercial images, consistently demonstrating superior or competitive performance against previous SOTA models.

## 2. Related work

**Saliency Prediction:** Saliency prediction models are designed to identify areas within an image or video that capture human visual attention. Initially inspired by biological insights, these models historically used contrasts in color, intensity, and orientation, or low-level, hand-designed features to mimic human visual perception. This approach was based on cues from studies of how humans prioritize visual information [17, 27, 55].

With the advent of deep learning and the availability of large-scale eye-tracking datasets [3, 28, 30], there has been a shift towards applying deep neural networks to the problem of saliency prediction. This shift was marked by significant improvements in the accuracy and reliability of saliency models [61]. Kummerer et al. [38] demonstrated that leveraging pretrained networks, originally designed for object recognition tasks, could enhance the performance of saliency prediction models. This insight paved the way for subsequent models such as EML-Net [26], DeepGaze II [39], and SALICON [25], which incorporated pretrained CNN encoders to enhance the prediction of saliency maps. Beyond the use of pretrained CNNs, researchers have explored various other network architectures for saliency prediction. These include fully convolutional networks (FCNs) [37], generative adversarial networks (GANs) [7, 51], and convolutional long short-term memory networks (ConvLSTM) [28]. Attention mechanism and Transformer models [60], which show remarkable success in various vision tasks [20], have also been applied to saliency prediction. Models like VGG-SSM [11] and TranSalNet [48], incorporate self-attention modules and transformer-based methods, respectively. These approaches highlight the growing interest in leveraging advanced architectures that go beyond traditional CNNs to improve saliency.

Saliency prediction has also expanded to cover diverse types of data beyond natural scenes, including commercial advertisements [29, 40, 42] and user interfaces [31, 58]. This diversification has led to specialized models that address unique dataset challenges. For instance, Kou et al.[35] proposed a method for integrating confidence scores in saliency predictions for advertising images, enhancing both robustness and performance. Similarly, Jiang et al. [29] introduced salient Swin-Transformers and incorporated text detection techniques into their models, demonstrating the potential of combining various data modalities to improve prediction accuracy. Following this trend, Hosseini et al. [22] proposed a model that combines pretrained CNNs and transformers with a text map detection module for advertising saliency prediction.

**Mamba:** Recent advancements in SSMs, particularly with the development of the Mamba model [18], have significantly changed the landscape of computational modeling. This model offers a promising alternative to traditional attention-based models. Introduced by Gu et al. [18], Mamba achieves linear computational complexity with respect to input size and is effective at capturing long-distance dependencies. This innovation has led to its broad application in fields such as language understanding and vision tasks [24, 47, 53, 54, 65, 67]. The development of vision-specific SSMs such as Vision Mamba [67] and Vmamba [47] has marked a significant step in SSM development. Notable examples include U-Mamba [49], which combines SSMs with CNNs for medical image segmentation. SegMamba [63] integrates SSMs in its encoder and uses a CNN-based decoder for 3D brain tumor segmentation. VM-UNet [57] explores a purely SSM-based ap-
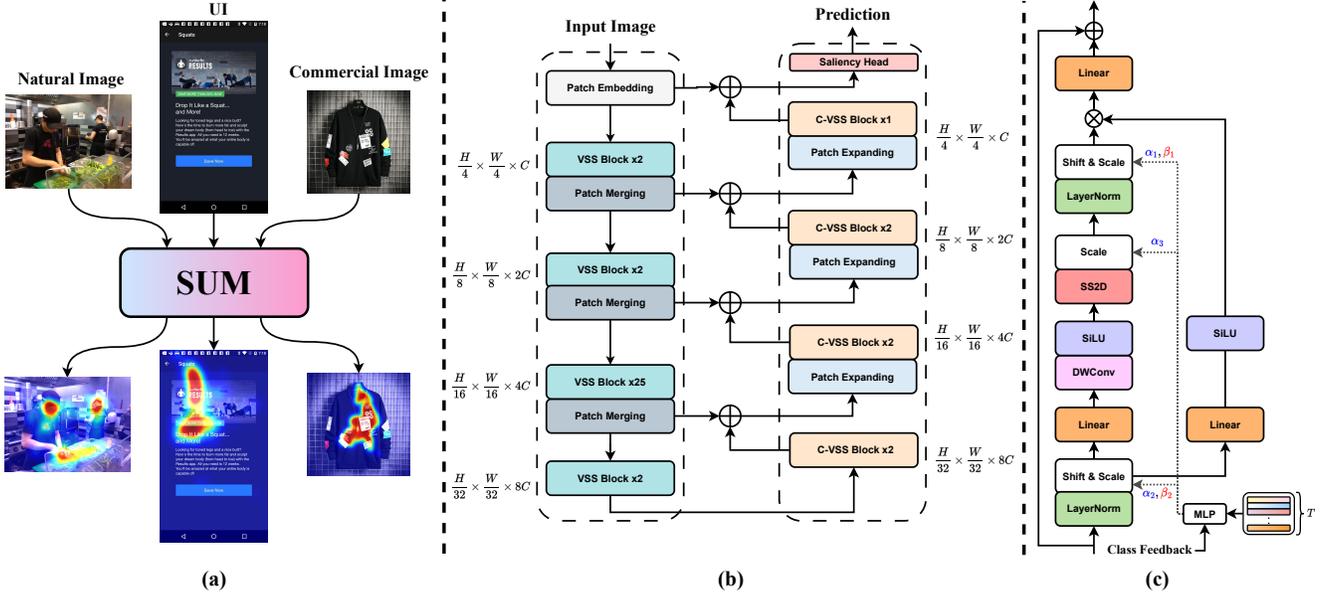
Figure 1. (a) Overview of our **SUM** model, (b) conditional-U-Net-based model for saliency prediction, and (c) C-VSS module.

proach in this area. Other models, like LightM-UNet [43], stand out for their efficiency, outperforming previous medical segmentation models with fewer parameters. Additionally, Mamba's versatility is demonstrated in video-based applications such as video medical segmentation [66] and understanding [8]. However, the use of Mamba in saliency map prediction is still largely unexplored.

**Unified Models:** Unified models in visual saliency prediction have made significant advancements in integrating image and video saliency within a single framework. The UNISAL model [14] is an example that addresses the integration of image and video saliency through domain adaptation. However, while UNISAL is a lightweight model, it is not a universal model for all image saliency datasets. It primarily relies on the Salicon Dataset [30] for image saliency prediction, and its performance on this dataset has been outperformed by other models over time. Furthermore, UNISAL's universal model does not include diverse image types, limiting its applicability. As another notable model, UniAR [41], focuses on image-based saliency prediction and incorporates a multimodal transformer to capture diverse human behaviors across various visual content and tasks. While it encompasses UI and natural scene images, it overlooks incorporating e-commercial images, which have become increasingly important in recent years [22, 29]. Additionally, UniAR's complexity is highlighted by its model size, with 848M parameters, making it computationally demanding and potentially limiting its practical use. Despite advancements in existing unified saliency models, there is still a significant gap in developing efficient, comprehensive models that effectively address real-world needs across di-

verse image types while maintaining manageable complexity.

## 3. Proposed Method

This section provides an overview of the proposed network architecture, as shown in Figure 1(a). Next, we revisit the concept of VSS as introduced by Liu et al. [47]. Building on this foundation, we introduce our novel C-VSS module and a conditional Mamba-U-Net-based model for visual saliency prediction.

### 3.1. Model Architecture

The architecture of **SUM**, as illustrated in Figure 1(b), adopts a U-Net configuration. The process initiates with an input image $X \in \mathbb{R}^{H \times W \times 3}$, with spatial dimensions $H$ and $W$, and 3 channels, which undergoes initial transformation via a patch embedding module, reducing its dimensions to $\frac{H}{4} \times \frac{W}{4} \times C$. The encoder module generates four hierarchical output representations. Each stage is followed by a downsampling layer, which reduces the spatial dimensions by half while simultaneously doubling the number of channels. Transitioning to the decoder, comprises four stages of C-VSS layers, with each stage incorporating two blocks, except for the final stage, which contains a single block. Patch-expanding layers are then applied to achieve resolution upsampling while also decreasing channel dimensions by a factor of 2. Finally, a linear layer is responsible for generating the ultimate output. Our SUM architecture uses the VMamba [47] weights pre-trained on ImageNet [12]. This pre-training accelerates the learning process, improves the model's ability to detect salient regions more accurately,

and ensures better generalization on diverse images.

## 3.2. Visual State Space (VSS)

Mamba [18] employs SSMs [19] to shift the complexity of attention from quadratic to linear in long-sequence modeling. This has proven particularly beneficial in vision tasks due to its higher accuracy, reduced computational load, and lower memory requirements [67]. However, adapting Mamba's inherently 1D, causal scanning for 2D images presents challenges due to its restricted receptive field and inability to process unscanned data effectively. To address these issues, VMamba [47] introduces the Cross-Scan Module, which employs bidirectional scanning along horizontal and vertical axes. This module expands the image into sequences of patches scanned in four directions, enabling each pixel to integrate information from all directions. Subsequently, these sequences are reassembled into the original 2D format to form a complete image. Termed the 2D-Selective-Scan (SS2D), this method enhances Mamba's functionality for 2D spatial processing, ensuring both local and global spatial relevance. Building upon these insights, we incorporate the VSS block as the fundamental unit in SUM. As shown in Figure 1(c), the VSS module can be formulated as:

$$\mathbf{X} = \text{LN}_1(\mathbf{F}),$$
$$\mathbf{Attention} = \text{LN}_2(\text{SS2D}(\text{SiLU}(\text{DW-Conv}(\text{Linear}(\mathbf{X}))))),$$
$$\mathbf{Output} = \text{Linear}(\text{SiLU}(\text{Linear}(\mathbf{X})) \otimes \mathbf{Attention}) + \mathbf{F},$$
(1)

where the input feature is denoted by $\mathbf{F} \in \mathbb{R}^{H' \times W' \times C'}$. The operator $\otimes$ denotes an element-wise product operation, LN represents LayerNorm, DW-Conv stands for depth-wise convolution, and SiLU [15] is an activation function.

## 3.3. Conditional Visual State Space (C-VSS)

We enhance the model's adaptability to diverse visual content by conditioning the VSS block in the decoder based on input type. This is crucial for predicting saliency maps effectively, as different content types inherently attract viewer attention in distinct ways. For instance, natural scenes may focus on color and movement, e-commerce images on textual information, and UI designs on specific layout patterns such as the upper-left quadrant. To address these variations, we implement modulation of the feature map through dynamic scaling and shifting operations that adjust feature activations based on the input type. The modulated feature map can be generally defined as:

$$\text{Modulated Feature Map} = \alpha_i \odot \mathbf{F} + \beta_i$$

where $\mathbf{F}$ denotes the original feature map. Here, $\alpha$ is a scaling factor, $\beta$ is a shifting factor, and $\odot$ is an element-wise multiplication.

To refine our model's ability to effectively handle different data types, we define $T = 4$ learnable tokens, where $D$ represents the dimensionality of each token. Each token is designated to capture distinct information about one of the following data categories: `Natural Scene-Mouse`, `Natural Scene-Eye`, `E-Commerce`, and `UI`. These tokens provide a more nuanced mechanism than a simple one-hot encoding of data types, enabling the model to adapt and learn detailed, type-specific information. We have allocated two tokens for the natural scene data because different methodologies are used in data collection for these categories, `eye` and `mouse`. Grouping them into a single token could potentially confuse the model during inference. As discussed in [59], mouse tracking data is less consistent and more scattered than eye tracking data, which does not fully align with eye tracking data distribution, particularly in terms of different contextual regions. Furthermore, while mouse tracking data can lead to acceptable outcomes for training existing models, it is less reliable for model selection and evaluation. Based on these insights and our experiments, we differentiate mouse and eye data of natural scenes.

Subsequently, the relevant token is fed into a Multi-Layer Perceptron (MLP) model to ensure that learning is conditioned on the specific characteristics of each data type. The MLP is composed of $K$ hidden layers and $p_1, p_2, \ldots, p_K$ features per layer. This MLP is designed to regress the parameters $\alpha_i$ and $\beta_i$, which modulate the model based on the diversity of inputs. The MLP, defined as $g(\mathbf{z}; \theta) : \mathbb{R}^{4 \times D} \to \mathbb{R}^{4 \times 5}$, outputs a matrix $\mathbf{Y}$, with each row representing one of four input tokens and generating five key parameters. These parameters include pairs and individual instances of $\alpha_i$ and $\beta_i$, specifically $\{(\alpha_1, \beta_1), (\alpha_3), (\alpha_2, \beta_2)\}$. An input label $L$ determines the selection of the relevant row from $\mathbf{Y}$, resulting in the output vector $\mathbf{S} = \mathbf{Y}_L$. This $1 \times 5$ vector contains modulation parameters finely tuned to the specifics of the designated input. These parameters are then integrated into the model to modify its behavior dynamically: $(\alpha_1, \beta_1)$ are used to shift and scale $\text{LN}_1$, $(\alpha_3)$ adjusts the scaling of the SS2D block to regulate feature intensity, and $(\alpha_2, \beta_2)$ shift and scale $\text{LN}_2$. This enables the MLP to precisely control the normalization and scaling within the model, thereby enhancing its performance and generalization across different visual content types.

## 3.4. Loss Function

Our model utilizes a composite loss function inspired by [2, 14, 48] in visual saliency prediction. This function integrates five distinct components, each designed to optimize the prediction accuracy of saliency maps by targeting different aspects of the saliency prediction task. The loss function

is formulated as:

$$\text{Loss} = \lambda_1 \cdot \mathcal{L}_{\text{KL}}(s^g, s) + \lambda_2 \cdot \mathcal{L}_{\text{CC}}(s^g, s) + \lambda_3 \cdot \mathcal{L}_{\text{SIM}}(s^g, s) \\ + \lambda_4 \cdot \mathcal{L}_{\text{NSS}}(f^g, s) + \lambda_5 \cdot \mathcal{L}_{\text{MSE}}(s^g, s) \tag{2}$$

where $s^g$ represents the ground truth saliency map, $f^g$ denotes the ground truth fixation map, and $s$ is the network's predicted saliency map. Each component of the loss function serves a specific purpose as defined in the following.

**Kullback-Leibler Divergence (KL):** KL divergence measures the dissimilarity between the predicted and ground truth distributions, providing a method to penalize the model when its predictions deviate significantly from the actual data distribution.

$$\mathcal{L}_{\text{KL}}(s^g, s) = \sum_{i=1}^{n} s_i^g \log \left( \epsilon + \frac{s_i}{s_i^g + \epsilon} \right), \tag{3}$$

where, the regularization constant $\epsilon$ is set to $2.2 \times 10^{-16}$.

**Linear Correlation Coefficient (CC):** The correlation coefficient assesses the linear relationship between the predicted and ground truth saliency maps. A higher correlation indicates that the model predictions align well with the ground truth trends, improving the reliability of the saliency maps.

$$\mathcal{L}_{\text{CC}}(s^g, s) = \frac{\text{cov}(s^g, s)}{\sigma(s^g) \cdot \sigma(s)}, \tag{4}$$

where, $\text{cov}(.)$ represents the covariance and $\sigma(.)$ denotes the standard deviation.

**Similarity (SIM):** SIM evaluates the overlap between the predicted and actual saliency maps, emphasizing the importance of accurately predicting the salient regions.

$$\mathcal{L}_{\text{SIM}}(s^g, s) = \sum_{i=1}^{n} \min(s_i^g, s_i) \tag{5}$$

**Normalized Scan-path Saliency (NSS):** NSS measures the correlation between the normalized predicted saliency map and the actual fixation points, highlighting the model's effectiveness at capturing human attention patterns.

$$\mathcal{L}_{\text{NSS}}(f^g, s) = \frac{1}{\sum_i (f_i^g)} \sum_i \left( \frac{s_i - \mu(s)}{\sigma(s)} \right) f_i^g \tag{6}$$

**Mean Squared Error (MSE):** This component calculates the mean squared error between the predicted and actual saliency maps, directly penalizing inaccuracies in the pixel-wise saliency values.

By adjusting the weighting coefficients $\lambda_i$ ($i = 1, \ldots, 5$), we aim to minimize dissimilarity metrics (KL, MSE) and maximize similarity metrics (CC, SIM, NSS). This strategy ensures that the model predicts accurate saliency maps and closely aligns with human visual attention patterns and saliency distributions.

Table 1. Comprehensive compilation of datasets used for training and testing.

| Dataset | Image domain | Acquisition Type | # Image | Image Resolution | # Training Sample |
|---|---|---|---|---|---|
| *Salicon* [30] | Natural scene | Mouse | 15,000 | $640 \times 480$ | 10,000 |
| *MIT1003* [33] | Natural scene | Eye | 1003 | Varied | 904 |
| *CAT2000* [3] | Natural scene | Eye | 2000 | $1080 \times 1920$ | 1600 |
| *OSIE* [64] | Natural scene | Eye | 700 | $800 \times 600$ | 500 |
| *U-EYE* [31] | Web page | Eye | 1979 | Varied | 1583 |
| *SalECI* [29] | E-Commercial | Eye | 972 | $720 \times 720$ | 871 |

## 4. Experiments

**Datasets:** We leverage six benchmark large-scale datasets for training and evaluating our models, as outlined in Table 9. This table presents a list of these datasets along with specific details about each.

**Evaluation Metrics:** To assess the accuracy of predicted saliency maps, we use two types of metrics: *location-based* and *distribution-based*, followed by [5]. Location-based metrics, such as NSS and AUC (Area under the ROC Curve), evaluate predictions using a binary fixation map image as ground truth and focus on specific salient locations. Distribution-based metrics, including CC (Correlation Coefficient), SIM (Similarity), and KLD (Kullback-Leibler Divergence), utilize a grayscale saliency map image to measure the similarity between predicted and actual distributions. Higher values generally indicate better performance for all metrics, except for KLD, where a value closer to zero signifies a more accurate prediction.

**Implementation Details:** Our model is implemented using the PyTorch framework and is trained on A40 with 48 GB memory for 15 epochs with an early stopping after 4 epochs. We optimize the network using the Adam optimizer. The learning rate is initially set to $1 \times 10^{-4}$, and we employ a learning rate scheduler that decreases the factor by 0.1 after every four epochs. The batch size is set to 16. Additionally, we resize all data and labels to a resolution of $256 \times 256$ and combine the training data from all six datasets for model training. The optimal values for the loss function weighting coefficients $\lambda_i$ are as follows: $\lambda_1 = 10$, $\lambda_2 = -2$, $\lambda_3 = -1$, $\lambda_4 = -1$, and $\lambda_5 = 5$. In addition, the MLP architecture in our implementation comprises three linear layers with widths of 128, 64, and 5, respectively, interleaved with GELU [21] activation functions. The number of tokens is set to $T = 4$ with each token having a dimensionality of $D = 128$.

### 4.1. Experiment Results

We conducted comprehensive testing of our universal model, SUM, across six different datasets, each benchmarked against state-of-the-art (SOTA) models for comparison. These datasets encompass a range of areas including natural scenes, user interfaces, and e-commerce. SUM consistently outperformed the best existing models across all datasets. In the 30 metrics presented in Table 2, SUM

Table 2. Saliency prediction performance across various datasets. * indicates that we have trained those models ourselves for fair comparison because results were not available for the corresponding dataset or the input image size was varied. † signifies that the results have been taken from the paper by Hosseini et al. [22], and the rest of the results are taken from their respective papers. For our model, we note the percentage (%) change in performance relative to the second-best result, or to the best result if ours is not the top performer.

| Dataset | Method | CC ↑ | KLD ↓ | AUC ↑ | SIM ↑ | NSS ↑ | # Parameters |
|---|---|---|---|---|---|---|---|
| *U-EYE* [31] (Web page) | SAM* [11] | 0.580 | 1.490 | 0.811 | 0.520 | 1.640 | 30M |
| | UMSI* [16] | 0.562 | 1.580 | 0.805 | 0.510 | 1.690 | 30M |
| | SAM++* [31] | 0.580 | 1.190 | 0.800 | 0.530 | 1.660 | 42M |
| | Transalnet* [48] | 0.696 | 0.616 | 0.839 | 0.598 | 1.601 | 72M |
| | UMSI++* [31] | 0.670 | 0.860 | 0.830 | 0.580 | 1.610 | 30M |
| | **SUM (Ours)** | **0.731** +5.03% | **0.544** −11.69% | **0.846** +0.83% | **0.630** +5.35% | **1.704** +0.83% | 57.5M |
| *SalECI* [29] (E-Commercial) | SSM† [11] | 0.720 | 0.599 | 0.830 | 0.611 | 1.396 | 42M |
| | DeepGaze IIE† [44] | 0.560 | 0.995 | 0.842 | 0.399 | 1.327 | 104M |
| | EML-NET† [31] | 0.510 | 1.220 | 0.807 | 0.536 | 1.232 | 47M |
| | Transalnet† [31] | 0.717 | 0.873 | 0.824 | 0.534 | 1.723 | 72M |
| | Temp-Sal† [2] | 0.719 | 0.712 | 0.813 | 0.629 | 1.768 | 242M |
| | SSwin Transformer [29] | 0.687 | 0.652 | 0.868 | 0.606 | 1.701 | 29M |
| | Hosseini et al. [22] | 0.750 | 0.578 | 0.892 | 0.645 | 1.890 | 66M |
| | **SUM (Ours)** | **0.789** +5.20% | **0.473** −18.17% | **0.899** +0.78% | **0.680** +5.43% | **2.012** +6.46% | 57.5M |
| *OSIE* [64] (Natural scene) | UMSI [16] | 0.746 | 0.513 | 0.856 | 0.631 | 1.788 | 30M |
| | EML-NET [26] | 0.717 | 0.537 | 0.854 | 0.619 | 1.737 | 47M |
| | SAM-ResNet [11] | 0.758 | 0.480 | 0.860 | 0.648 | 1.811 | 43M |
| | Chen et al. [10] | 0.761 | 0.506 | 0.860 | 0.652 | 1.840 | - |
| | Transalnet* [31] | 0.791 | 0.667 | 0.923 | 0.651 | 2.448 | 72M |
| | UniAR [41] | 0.754 | 0.547 | 0.867 | 0.647 | 1.842 | 848M |
| | **SUM (Ours)** | **0.861** +8.85% | **0.340** −29.17% | **0.924** +6.57% | **0.727** +11.5% | **3.416** +39.54% | 57.5M |
| *Salicon* [30] (Natural scene) | UniAR [41] | 0.901 | 0.215 | 0.870 | 0.792 | 1.947 | 848M |
| | SimpleNet [56] | 0.907 | 0.193 | 0.871 | 0.797 | 1.926 | 116M |
| | MDNSal [56] | 0.899 | 0.217 | 0.868 | 0.797 | 1.893 | - |
| | MSI-Net [36] | 0.899 | 0.307 | 0.865 | 0.784 | 1.931 | 20M |
| | GazeGAN[7] | 0.879 | 0.376 | 0.864 | 0.773 | 1.899 | - |
| | UNISAL[14] | 0.879 | 0.354 | 0.864 | 0.775 | 1.952 | 4M |
| | Transalnet* [48] | 0.89 | 0.220 | 0.867 | 0.783 | 1.924 | 72M |
| | DeepGaze IIE* [44] | 0.872 | 0.285 | 0.869 | 0.733 | **1.996** | 104M |
| | Temp-Sal* [2] | **0.911** | 0.195 | 0.869 | 0.800 | 1.967 | 242M |
| | **SUM (Ours)** | 0.909 −0.22% | **0.192** −1.54% | **0.876** +0.64% | **0.804** +0.50% | 1.981 −0.75% | 57.5M |
| *CAT2000* [3] (Natural scene) | FastSal [23] | 0.721 | 0.552 | 0.86 | 0.603 | 1.859 | 4M |
| | SAM-Resnet [11] | 0.87 | 0.670 | 0.878 | 0.739 | 2.411 | 43M |
| | MSI-Net* [36] | 0.866 | 0.428 | 0.881 | 0.730 | 2.355 | 20M |
| | DVA [62] | 0.861 | 0.449 | 0.878 | 0.734 | 2.345 | - |
| | UNISAL [14] | 0.842 | 0.530 | 0.876 | 0.721 | 2.257 | 4M |
| | MDNSal [56] | **0.889** | 0.293 | 0.878 | 0.751 | 2.329 | - |
| | Transalnet* [48] | 0.877 | 0.287 | 0.882 | 0.744 | 2.373 | 72M |
| | **SUM (Ours)** | 0.882 −0.79% | **0.270** −5.92% | **0.888** +0.68% | **0.754** +0.4% | **2.424** +0.54% | 57.5M |
| *MIT1003* [33] (Natural scene) | FastSal [23] | 0.590 | 1.036 | 0.875 | 0.478 | 2.008 | 4M |
| | SAM-Resnet [11] | 0.746 | 1.247 | 0.902 | 0.597 | 2.752 | 43M |
| | DVA [62] | 0.699 | 0.753 | 0.897 | 0.566 | 2.574 | - |
| | UNISAL [14] | 0.734 | 1.014 | 0.902 | 0.597 | 2.759 | 4M |
| | Transalnet* [48] | 0.722 | 0.660 | 0.903 | 0.592 | 2.631 | 72M |
| | **SUM (Ours)** | **0.768** +2.95% | **0.563** −14.7% | **0.913** +1.11% | **0.630** +5.53% | **2.839** +2.9% | 57.5M |

achieved SOTA results in 27 cases and secured second place in the other three. These results demonstrate that our model is highly effective and versatile across various types of data, setting a new standard for future advancements in the field. This consistency in performance underscores its robustness and capability to handle the diverse challenges presented
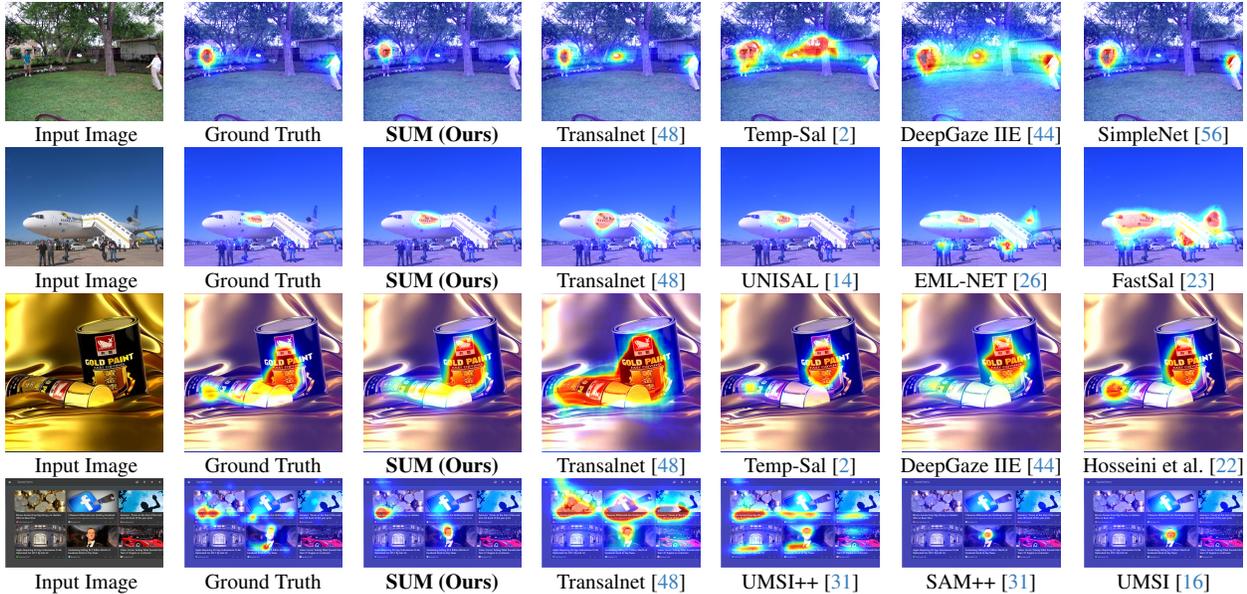
Figure 2. Comparative visualizations of saliency predictions across different data types. The first row depicts Natural Scene-Mouse data, the second row showcases Natural Scene-Eye data, the third row features E-commerce, and the fourth row displays UI. Each row highlights the model's performance in identifying salient features within these distinct categories.

by different datasets. Moreover, compared to counterparts like Transalnet [48], Temp-Sal [2], DeepGaze IIE [44], and UniAR [41], our model is relatively efficient. This efficiency underscores the advantages of our streamlined approach, which leverages Mamba's capabilities to develop a model that is efficient, robust, and universally applicable. Additionally, Figure 2 displays saliency prediction images selected from the validation sets, showing that our model's predictions are much closer to the ground truth than those of the SOTA models, further proving that SUM can more accurately predict human attention behavior.

## 5. Ablation Study

**Impact of different loss combinations :** We investigated the impact of different loss metric combinations on the model validation performance, as summarized in Table 3. Our approach involved normalizing each metric using a min-max scaling technique to ensure a balanced evaluation across different metrics. The score function, described in Equation 7, is specifically designed to maximize the beneficial metrics (CC, SIM, NSS) and minimize the detrimental metric (KL). The function's configuration is as follows:

$$\mathcal{F}_{\text{score}} = CC_{\text{scaled}} + SIM_{\text{scaled}} + NSS_{\text{scaled}} - KL_{\text{scaled}} \quad (7)$$

From the results in Table 3, it is evident that the inclusion of KL loss significantly impacts the model's performance, demonstrating its crucial role in defining saliency loss. When loss functions are used individually, the performance varies, with SIM typically showing higher values

for both CC and $\mathcal{F}$ scores, indicating its strong standalone impact on model saliency. Excluding MSE, which is less directly related to saliency, still results in high performance, but the highest scores are consistently observed when MSE is included, suggesting its underlying contribution to model robustness and generalization. The integration of all five loss functions results in the highest $\mathcal{F}$ scores. This combination not only balances the enhancement and suppression of features but also stabilizes the training process, as indicated by the highest scores of 2.853 and 2.836 for Salicon and all datasets, respectively.

**Impact of the C-VSS module:** We compared the impact of using a C-VSS module conditioned on three and four classes against a standard VSS, which serves as the unconditional setup. The three classes include Natural Scene, UI, and E-Commerce, in contrast to our model's broader categorization into four classes. As shown in Table 4, the C-VSS module significantly enhances performance across all evaluated datasets compared to the standard VSS. Notably, conditioning the model on four classes yields better results than limiting it to three. This suggests that the finer categorization in the four-class setup better aligns with the varied data characteristics, especially when dealing with different data acquisition setups, thereby improving the model's predictive accuracy and robustness.

**Impact of different prompt lengths:** We explored the influence of prompt length on model performance. We experimented with various prompt lengths—32, 96, 128, 256—to determine how they impact model behavior during the train-

Table 3. Evaluation of different combinations of loss functions on model performance.

| Loss Functions | | | | | Avg. Performance on Salicon [30] | | | | | Avg. Performance Across All Datasets | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| KL | CC | SIM | NSS | MSE | CC ↑ | KLD ↓ | NSS ↑ | SIM ↑ | $\mathcal{F}_{\text{Score}}$ ↑ | CC ↑ | KLD ↓ | NSS ↑ | SIM ↑ | $\mathcal{F}_{\text{Score}}$ ↑ |
| ✓ | ✗ | ✗ | ✗ | ✗ | 0.910 | 0.189 | 1.908 | 0.805 | 2.797 | 0.85 | 0.465 | 2.498 | 0.723 | 2.386 |
| ✗ | ✓ | ✗ | ✗ | ✗ | 0.907 | 0.732 | 1.926 | 0.787 | 1.634 | 0.851 | 1.08 | 2.532 | 0.7 | 1.218 |
| ✗ | ✗ | ✓ | ✗ | ✗ | **0.911** | 0.447 | 1.91 | **0.807** | 2.391 | 0.85 | 0.747 | 2.469 | **0.728** | 1.917 |
| ✗ | ✗ | ✗ | ✓ | ✗ | 0.834 | 0.765 | **2.044** | 0.721 | 0 | 0.804 | 1.072 | 2.614 | 0.658 | -0.079 |
| ✗ | ✗ | ✗ | ✗ | ✓ | 0.909 | 0.234 | 1.919 | 0.803 | 2.696 | 0.846 | 0.525 | 2.479 | 0.719 | 2.089 |
| ✓ | ✗ | ✓ | ✗ | ✗ | **0.911** | 0.196 | 1.928 | 0.806 | 2.833 | **0.852** | 0.465 | 2.337 | **0.728** | 1.972 |
| ✓ | ✗ | ✗ | ✓ | ✗ | 0.892 | 0.199 | 2.029 | 0.792 | 2.537 | 0.841 | 0.467 | 2.594 | 0.712 | 2.353 |
| ✓ | ✓ | ✗ | ✗ | ✗ | **0.911** | **0.185** | 1.191 | 0.805 | 1.977 | **0.852** | 0.453 | 2.515 | 0.720 | 2.46 |
| ✓ | ✗ | ✗ | ✗ | ✓ | 0.909 | 0.192 | 1.917 | 0.802 | 2.755 | 0.851 | 0.456 | 2.504 | 0.723 | 2.441 |
| ✗ | ✓ | ✓ | ✗ | ✗ | 0.910 | 0.531 | 1.921 | 0.802 | 2.188 | 0.85 | 0.871 | 2.503 | 0.721 | 1.733 |
| ✓ | ✓ | ✓ | ✗ | ✗ | 0.909 | 0.198 | 1.920 | 0.803 | 2.759 | **0.852** | 0.464 | 2.527 | 0.726 | 2.568 |
| ✓ | ✗ | ✓ | ✗ | ✓ | 0.909 | 0.192 | 1.919 | 0.799 | 2.722 | **0.852** | 0.461 | 2.514 | 0.726 | 2.53 |
| ✓ | ✗ | ✗ | ✓ | ✓ | 0.887 | 0.208 | 2.038 | 0.788 | 2.421 | 0.830 | 0.472 | **2.642** | 0.711 | 2.259 |
| ✓ | ✓ | ✗ | ✗ | ✓ | 0.910 | 0.188 | 1.914 | 0.803 | 2.783 | 0.851 | **0.447** | 2.511 | 0.722 | 2.464 |
| ✓ | ✓ | ✓ | ✓ | ✗ | 0.907 | 0.198 | 1.989 | 0.803 | 2.815 | 0.850 | 0.466 | 2.614 | 0.725 | 2.794 |
| ✓ | ✓ | ✓ | ✗ | ✓ | 0.905 | 0.208 | 1.920 | 0.798 | 2.632 | **0.852** | 0.457 | 2.510 | 0.720 | 2.437 |
| ✓ | ✓ | ✓ | ✓ | ✓ | 0.909 | 0.192 | 1.981 | 0.804 | **2.853** | **0.852** | 0.450 | 2.602 | 0.726 | **2.836** |

Table 4. Mean value and standard deviation of saliency prediction performance comparison of conditional VSS modules for three and four classes and standard VSS (no-condition) across all datasets.

| Dataset | Method | CC ↑ | KLD ↓ | AUC ↑ | SIM ↑ | NSS ↑ |
|---|---|---|---|---|---|---|
| *U-EYE* [31] (Web page) | No-condition | $0.725 \pm 0.035$ | $0.562 \pm 0.062$ | $0.845 \pm 0.012$ | $0.626 \pm 0.023$ | $1.689 \pm 0.121$ |
| | 3-class | $0.729 \pm 0.035$ | $0.551 \pm 0.057$ | $0.845 \pm 0.012$ | $0.628 \pm 0.012$ | $1.699 \pm 0.012$ |
| | 4-class | $\mathbf{0.731} \pm \mathbf{0.037}$ | $\mathbf{0.544} \pm \mathbf{0.057}$ | $\mathbf{0.846} \pm \mathbf{0.012}$ | $\mathbf{0.630} \pm \mathbf{0.023}$ | $\mathbf{1.704} \pm \mathbf{0.125}$ |
| *SalECI* [29] (E-Commercial) | No-condition | $0.783 \pm 0.046$ | $0.502 \pm 0.112$ | $0.898 \pm 0.014$ | $0.677 \pm 0.039$ | $\mathbf{2.017} \pm 0.168$ |
| | 3-class | $0.781 \pm 0.055$ | $0.505 \pm 0.131$ | $0.896 \pm 0.016$ | $0.678 \pm 0.047$ | $1.99 \pm 0.181$ |
| | 4-class | $\mathbf{0.789} \pm \mathbf{0.0453}$ | $\mathbf{0.473} \pm \mathbf{0.088}$ | $\mathbf{0.899} \pm \mathbf{0.012}$ | $\mathbf{0.680} \pm \mathbf{0.041}$ | $2.012 \pm 0.161$ |
| *OSIE* [64] (Natural scene) | No-condition | $0.842 \pm 0.033$ | $0.403 \pm 0.05$ | $0.918 \pm 0.009$ | $0.703 \pm 0.022$ | $3.18 \pm 0.32$ |
| | 3-class | $0.845 \pm 0.033$ | $0.395 \pm 0.049$ | $0.918 \pm 0.009$ | $0.706 \pm 0.022$ | $3.213 \pm 0.323$ |
| | 4-class | $\mathbf{0.861} \pm \mathbf{0.029}$ | $\mathbf{0.340} \pm \mathbf{0.050}$ | $\mathbf{0.924} \pm \mathbf{0.008}$ | $\mathbf{0.727} \pm \mathbf{0.02}$ | $\mathbf{3.416} \pm \mathbf{0.319}$ |
| *Salicon* [30] (Natural scene) | No-condition | $0.903 \pm 0.012$ | $0.206 \pm 0.028$ | $0.875 \pm 0.014$ | $0.798 \pm 0.013$ | $1.979 \pm 0.203$ |
| | 3-class | $0.904 \pm 0.011$ | $0.205 \pm 0.027$ | $0.875 \pm 0.014$ | $0.798 \pm 0.013$ | $1.981 \pm 0.204$ |
| | 4-class | $\mathbf{0.909} \pm \mathbf{0.011}$ | $\mathbf{0.192} \pm \mathbf{0.025}$ | $\mathbf{0.876} \pm \mathbf{0.014}$ | $\mathbf{0.804} \pm \mathbf{0.012}$ | $\mathbf{1.981} \pm \mathbf{0.201}$ |
| *CAT2000* [3] (Natural scene) | No-condition | $0.880 \pm 0.014$ | $0.272 \pm 0.022$ | $0.887 \pm 0.010$ | $0.752 \pm 0.010$ | $2.42 \pm 0.141$ |
| | 3-class | $0.881 \pm 0.016$ | $0.271 \pm 0.023$ | $\mathbf{0.888} \pm \mathbf{0.010}$ | $0.753 \pm 0.011$ | $\mathbf{2.424} \pm \mathbf{0.142}$ |
| | 4-class | $\mathbf{0.882} \pm \mathbf{0.0158}$ | $\mathbf{0.270} \pm \mathbf{0.026}$ | $\mathbf{0.888} \pm \mathbf{0.010}$ | $\mathbf{0.754} \pm \mathbf{0.011}$ | $\mathbf{2.424} \pm \mathbf{0.142}$ |
| *MIT1003* [33] (Natural scene) | No-condition | $0.737 \pm 0.035$ | $0.641 \pm 0.083$ | $0.908 \pm 0.010$ | $0.596 \pm 0.024$ | $2.648 \pm 0.255$ |
| | 3-class | $0.741 \pm 0.034$ | $0.636 \pm 0.077$ | $0.908 \pm 0.010$ | $0.597 \pm 0.023$ | $2.678 \pm 0.249$ |
| | 4-class | $\mathbf{0.768} \pm \mathbf{0.039}$ | $\mathbf{0.563} \pm \mathbf{0.075}$ | $\mathbf{0.913} \pm \mathbf{0.009}$ | $\mathbf{0.630} \pm \mathbf{0.027}$ | $\mathbf{2.839} \pm \mathbf{0.285}$ |

ing and validation phases. The results, detailed in Table 5, indicate that both shorter and longer prompt lengths contribute to fitting issues. Among the tested lengths, 128 demonstrated the most balanced and effective outcome.

**Comparison of Prompt vs. one-hot encoding:** In our experiments, we compared two approaches: one using generated prompts tailored to specific conditions, and another using a one-hot vector to represent class conditions. Our goal was to see how these methods influence the model's ability to handle different types of data. Table 6 illustrates the results of this comparison. Using the prompt-based approach, the model demonstrates higher performance across all metrics. This method helps the model better distinguish between the diverse data distributions in each domain, as

Table 5. Impact of prompt length on model performance.

| Prompt Length | Performance Metrics | | | | # Parameters |
|---|---|---|---|---|---|
| | CC ↑ | KL ↓ | NSS ↑ | SIM ↑ | |
| Salicon [30] | | | | | |
| 64 | **0.909** | 0.196 | 1.98 | **0.804** | 57.4M |
| 96 | **0.909** | **0.188** | 1.958 | 0.802 | 57.4M |
| **128** | **0.909** | 0.192 | **1.981** | **0.804** | 57.5M |
| 256 | 0.906 | 0.195 | 1.953 | 0.801 | 58M |
| Average Performance Across Datasets | | | | | |
| 64 | 0.849 | 0.463 | 2.601 | 0.725 | 57.4M |
| 96 | 0.847 | 0.455 | 2.567 | 0.722 | 57.4M |
| **128** | **0.852** | **0.450** | **2.602** | **0.726** | 57.5M |
| 256 | 0.850 | 0.456 | 2.558 | 0.723 | 58M |

Table 6. Prompt vs. One-Hot Encoding.

| Method | Performance Metrics | | | | # Parameters |
|---|---|---|---|---|---|
| | CC ↑ | KL ↓ | NSS ↑ | SIM ↑ | |
| Avg. Performance on Salicon [30] | | | | | |
| SUM - One-hot | $0.902_{\pm 0.012}$ | $0.201_{\pm 0.024}$ | $1.97_{\pm 0.972}$ | $0.795_{\pm 0.012}$ | 57.3M |
| **SUM - Prompt** | $\mathbf{0.909}_{\pm 0.011}$ | $\mathbf{0.192}_{\pm 0.025}$ | $\mathbf{1.981}_{\pm 0.201}$ | $\mathbf{0.804}_{\pm 0.012}$ | 57.5M |
| Avg. Performance Across Datasets | | | | | |
| SUM - One-hot | $0.843_{\pm 0.034}$ | $0.485_{\pm 0.046}$ | $2.583_{\pm 0.222}$ | $0.716_{\pm 0.023}$ | 57.3M |
| **SUM - Prompt** | $\mathbf{0.852}_{\pm 0.029}$ | $\mathbf{0.45}_{\pm 0.053}$ | $\mathbf{2.602}_{\pm 0.206}$ | $\mathbf{0.726}_{\pm 0.022}$ | 57.5M |

Table 7. Comparison of C-VSS placement in the proposed U-Net structure.

| Configuration | CC ↑ | KL ↓ | NSS ↑ | SIM ↑ | # Parameters |
|---|---|---|---|---|---|
| Avg. Performance on Salicon [30] | | | | | |
| Bottleneck | **0.909** | 0.195 | 1.97 | **0.804** | 57.37M |
| Decoder | **0.909** | **0.192** | **1.981** | **0.804** | 57.5M |
| All-Blocks | 0.907 | 0.198 | 1.975 | 0.801 | 58.5M |
| Avg. Performance Across Datasets | | | | | |
| Bottleneck | 0.847 | 0.466 | 2.581 | 0.724 | 57.37M |
| Decoder | 0.852 | **0.450** | **2.602** | **0.726** | 57.5M |
| All-Blocks | **0.854** | 0.458 | 2.601 | 0.724 | 58.5M |

opposed to the more straightforward one-hot vector method.

**Optimal C-VSS Placement in U-Net:** We evaluated the impact of deploying the C-VSS in different sections of our U-Net structure: solely in the bottleneck, across all blocks of the decoder, and in every block of both the encoder and decoder. Our objective was to ascertain the optimal placement of the C-VSS for enhancing model performance. As summarized in Table 7, incorporating C-VSS in the encoder, in addition to the decoder, tends to undermine the features in the encoder, leading to suboptimal performance. This observation suggests that integrating C-VSS throughout the entire U-Net may disrupt the model's ability to leverage its foundational pre-trained features effectively. Conversely, limiting the use of C-VSS to the bottleneck provides some benefits but does not fully capitalize on the potential enhancements the module offers. The most effective strategy, as indicated by our results, is employing C-VSS across all decoder blocks. This approach allows the model to better adapt to the unique characteristics of each input domain, resulting in superior performance metrics compared to the other configurations tested.

# 6. Conclusion

In this paper, we have presented **SUM**, a novel approach designed to address the limitations of traditional saliency prediction models. By integrating the Mamba architecture with U-Net and enhancing it with a Conditional Visual State Space (C-VSS) block, SUM adapts dynamically to various image types, making it universally applicable across diverse visual contexts. Our extensive evaluations across six benchmark datasets demonstrated SUM's superior performance, consistently outperforming existing models. The model excelled in both location-based and distribution-based metrics, proving its robustness and adaptability to use in real-world problems.

# References

[1] Hani Alers, Hantao Liu, Judith Redi, and Ingrid Heynderickx. Studying the effect of optimizing the image quality in saliency regions at the expense of background content. In *Image Quality and System Performance VII*, pages 59–67. SPIE, 2010. 1, 3

[2] Bahar Aydemir, Ludo Hoffstetter, Tong Zhang, Mathieu Salzmann, and Sabine Süsstrunk. Tempsal-uncovering temporal information for deep saliency prediction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6461–6470, 2023. 4, 6, 7

[3] Ali Borji and Laurent Itti. Cat2000: A large scale fixation dataset for boosting saliency research. *arXiv preprint arXiv:1505.03581*, 2015. 1, 2, 5, 6, 8

[4] Neil Bruce and John Tsotsos. Attention based on information maximization. *Journal of Vision*, 7(9):950–950, 2007. 1, 3

[5] Zoya Bylinskii et al. What do different evaluation metrics tell us about saliency models? *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 41(3):740–757, 2018. 5

[6] Zhaohui Che, Ali Borji, Guangtao Zhai, Xiongkuo Min, Guodong Guo, and Patrick Le Callet. How is gaze influenced by image transformations? dataset and model. *IEEE Transactions on Image Processing*, 29:2287–2300, 2019. 1

[7] Zhaohui Che, Ali Borji, Guangtao Zhai, Xiongkuo Min, Guodong Guo, and Patrick Le Callet. Gazegan: A generative adversarial saliency model based on invariance analysis of human gaze during scene free viewing. *arXiv preprint arXiv:1905.06803*, 2019. 2, 6

[8] Guo Chen, Yifei Huang, Jilan Xu, Baoqi Pei, Zhe Chen, Zhiqi Li, Jiahao Wang, Kunchang Li, Tong Lu, and Limin Wang. Video mamba suite: State space model as a versatile alternative for video understanding. *arXiv preprint arXiv:2403.09626*, 2024. 3

[9] Jiazhong Chen, Zongyi Li, Yi Jin, Dakai Ren, and Hefei Ling. Video saliency prediction via spatio-temporal reasoning. *Neurocomputing*, 462:59–68, 2021. 1

[10] Shi Chen, Nachiappan Valliappan, Shaolei Shen, Xinyu Ye, Kai Kohlhoff, and Junfeng He. Learning from unique perspectives: User-aware saliency modeling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2701–2710, 2023. 6

[11] Marcella Cornia, Lorenzo Baraldi, Giuseppe Serra, and Rita Cucchiara. Predicting human eye fixations via an lstm-based saliency attentive model. *IEEE Transactions on Image Processing*, 27(10):5142–5154, 2018. 1, 2, 6

[12] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009. 3

[13] Yasser Abdelaziz Dahou Djilali, Kevin McGuinness, and Noel O'Connor. Learning saliency from fixations. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 383–393, 2024. 1

[14] Richard Droste, Jianbo Jiao, and J Alison Noble. Unified image and video saliency modeling. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part V 16*, pages 419–435. Springer, 2020. 1, 3, 4, 6, 7

[15] Stefan Elfwing, Eiji Uchibe, and Kenji Doya. Sigmoid-weighted linear units for neural network function approximation in reinforcement learning. *Neural networks*, 107:3–11, 2018. 4

[16] Camilo Fosco, Vincent Casser, Amish Kumar Bedi, Peter O'Donovan, Aaron Hertzmann, and Zoya Bylinskii. Predicting visual importance across graphic design types. In *Proceedings of the 33rd Annual ACM Symposium on User Interface Software and Technology*, pages 249–260, 2020. 6, 7

[17] Stas Goferman, Lihi Zelnik-Manor, and Ayellet Tal. Context-aware saliency detection. *IEEE transactions on pattern analysis and machine intelligence*, 34(10):1915–1926, 2011. 2

[18] Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv preprint arXiv:2312.00752*, 2023. 2, 4

[19] Albert Gu, Karan Goel, and Christopher Ré. Efficiently modeling long sequences with structured state spaces. *arXiv preprint arXiv:2111.00396*, 2021. 4

[20] Kai Han, Yunhe Wang, Hanting Chen, Xinghao Chen, Jianyuan Guo, Zhenhua Liu, Yehui Tang, An Xiao, Chunjing Xu, Yixing Xu, et al. A survey on vision transformer. *IEEE transactions on pattern analysis and machine intelligence*, 45(1):87–110, 2022. 1, 2

[21] Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus). *arXiv preprint arXiv:1606.08415*, 2016. 5

[22] Alireza Hosseini, Kiana Hooshanfar, Pouria Omrani, Reza Toosi, Ramin Toosi, Zahra Ebrahimian, and Mohammad Ali Akhaee. Brand visibility in packaging: A deep learning approach for logo detection, saliency-map prediction, and logo placement analysis. *arXiv preprint arXiv:2403.02336*, 2024. 1, 2, 3, 6, 7

[23] Feiyan Hu and Kevin McGuinness. Fastsal: A computationally efficient network for visual saliency prediction. In *2020 25th International Conference on Pattern Recognition (ICPR)*, pages 9054–9061. IEEE, 2021. 1, 6, 7

[24] Vincent Tao Hu, Stefan Andreas Baumann, Ming Gui, Olga Grebenkova, Pingchuan Ma, Johannes Fischer, and Bjorn Ommer. Zigma: Zigzag mamba diffusion model. *arXiv preprint arXiv:2403.13802*, 2024. 2

[25] Xun Huang, Chengyao Shen, Xavier Boix, and Qi Zhao. Salicon: Reducing the semantic gap in saliency prediction by adapting deep neural networks. In *Proceedings of the IEEE international conference on computer vision*, pages 262–270, 2015. 2

[26] Sen Jia and Neil DB Bruce. Eml-net: An expandable multi-layer network for saliency prediction. *Image and vision computing*, 95:103887, 2020. 1, 2, 6, 7

[27] Lai Jiang, Mai Xu, Zhaoting Ye, and Zulin Wang. Image saliency detection with sparse representation of learnt texture atoms. In *Proceedings of the IEEE international conference on computer vision workshops*, pages 54–62, 2015. 2

[28] Lai Jiang, Mai Xu, Zulin Wang, and Leonid Sigal. Deepvs2. 0: A saliency-structured deep learning method for predicting dynamic visual attention. *International Journal of Computer Vision*, 129(1):203–224, 2021. 2

[29] Lai Jiang, Yifei Li, Shengxi Li, Mai Xu, Se Lei, Yichen Guo, and Bo Huang. Does text attract attention on e-commerce images: A novel saliency prediction dataset and method. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2088–2097, 2022. 1, 2, 3, 5, 6, 8

[30] Ming Jiang, Shengsheng Huang, Juanyong Duan, and Qi Zhao. Salicon: Saliency in context. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1072–1080, 2015. 1, 2, 3, 5, 6, 8, 9

[31] Yue Jiang, Luis A Leiva, Hamed Rezazadegan Tavakoli, Paul RB Houssel, Julia Kylmälä, and Antti Oulasvirta. Ueyes: Understanding visual saliency across user interface types. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*, pages 1–21, 2023. 1, 2, 5, 6, 7, 8

[32] John Jonides, David E Irwin, and Steven Yantis. Integrating visual information from successive fixations. *Science*, 215 (4529):192–194, 1982. 1

[33] Tilke Judd, Krista Ehinger, Frédo Durand, and Antonio Torralba. Learning to predict where humans look. In *2009 IEEE 12th international conference on computer vision*, pages 2106–2113. IEEE, 2009. 5, 6, 8

[34] Rudolph Emil Kalman. A new approach to linear filtering and prediction problems. 1960. 2

[35] Qiqi Kou, Ruihang Liu, Chen Lv, He Jiang, and Deqiang Cheng. Advertising image saliency prediction method based on score level fusion. *IEEE Access*, 11:8455–8466, 2023. 2

[36] Alexander Kroner, Mario Senden, Kurt Driessens, and Rainer Goebel. Contextual encoder–decoder network for visual saliency prediction. *Neural Networks*, 129:261–270, 2020. 1, 6

[37] Srinivas SS Kruthiventi, Kumar Ayush, and R Venkatesh Babu. Deepfix: A fully convolutional neural network for predicting human eye fixations. *IEEE Transactions on Image Processing*, 26(9):4446–4456, 2017. 2

[38] Matthias Kümmerer, Lucas Theis, and Matthias Bethge. Deep gaze i: Boosting saliency prediction with feature maps trained on imagenet. *arXiv preprint arXiv:1411.1045*, 2014. 1, 2

[39] Matthias Kümmerer, Thomas SA Wallis, and Matthias Bethge. Deepgaze ii: Reading fixations from deep

features trained on object recognition. *arXiv preprint arXiv:1610.01563*, 2016. 2

[40] Lucie Leveque and Hantao Liu. An eye-tracking database of video advertising. In *2019 IEEE International Conference on Image Processing (ICIP)*, pages 425–429. IEEE, 2019. 2

[41] Peizhao Li, Junfeng He, Gang Li, Rachit Bhargava, Shaolei Shen, Nachiappan Valliappan, Youwei Liang, Hongxiang Gu, Venky Ramachandran, Golnaz Farhadi, et al. Uniar: Unifying human attention and response prediction on visual content. *arXiv preprint arXiv:2312.10175*, 2023. 1, 3, 6, 7

[42] Song Liang, Ruihang Liu, and Jiansheng Qian. Fixation prediction for advertising images: Dataset and benchmark. *Journal of Visual Communication and Image Representation*, 81:103356, 2021. 2

[43] Weibin Liao, Yinghao Zhu, Xinyuan Wang, Cehngwei Pan, Yasha Wang, and Liantao Ma. Lightm-unet: Mamba assists in lightweight unet for medical image segmentation. *arXiv preprint arXiv:2403.05246*, 2024. 3

[44] Akis Linardos, Matthias Kümmerer, Ori Press, and Matthias Bethge. Deepgaze iie: Calibrated prediction in and out-of-domain for state-of-the-art saliency modeling. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 12919–12928, 2021. 6, 7

[45] Hantao Liu and Ingrid Heynderickx. Studying the added value of visual attention in objective image quality metrics based on eye movement data. In *2009 16th IEEE international conference on image processing (ICIP)*, pages 3097–3100. IEEE, 2009. 1, 3

[46] Nian Liu and Junwei Han. A deep spatial contextual long-term recurrent convolutional network for saliency detection. *IEEE Transactions on Image Processing*, 27(7):3264–3274, 2018. 1

[47] Yue Liu, Yunjie Tian, Yuzhong Zhao, Hongtian Yu, Lingxi Xie, Yaowei Wang, Qixiang Ye, and Yunfan Liu. Vmamba: Visual state space model. *arXiv preprint arXiv:2401.10166*, 2024. 2, 3, 4

[48] Jianxun Lou, Hanhe Lin, David Marshall, Dietmar Saupe, and Hantao Liu. Transalnet: Towards perceptually relevant visual saliency prediction. *Neurocomputing*, 494:455–467, 2022. 1, 2, 4, 6, 7

[49] Jun Ma, Feifei Li, and Bo Wang. U-mamba: Enhancing long-range dependency for biomedical image segmentation. *arXiv preprint arXiv:2401.04722*, 2024. 2

[50] Dipti Mishra, Satish Kumar Singh, Rajat Kumar Singh, and Divanshu Kedia. Multi-scale network (mssg-cnn) for joint image and saliency map learning-based compression. *Neurocomputing*, 460:95–105, 2021. 1

[51] Junting Pan, Cristian Canton Ferrer, Kevin McGuinness, Noel E O'Connor, Jordi Torres, Elisa Sayrol, and Xavier Giro-i Nieto. Salgan: Visual saliency prediction with generative adversarial networks. *arXiv preprint arXiv:1701.01081*, 2017. 2

[52] Yash Patel, Srikar Appalaraju, and R Manmatha. Saliency driven perceptual image compression. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pages 227–236, 2021. 1

[53] Xiaohuan Pei, Tao Huang, and Chang Xu. Efficientvmamba: Atrous selective scan for light weight visual mamba. *arXiv preprint arXiv:2403.09977*, 2024. 2

[54] Maciej Pióro, Kamil Ciebiera, Krystian Król, Jan Ludziejewski, and Sebastian Jaszczur. Moe-mamba: Efficient selective state space models with mixture of experts. *arXiv preprint arXiv:2401.04081*, 2024. 2

[55] Umesh Rajashekar, Ian Van Der Linde, Alan C Bovik, and Lawrence K Cormack. Gaffe: A gaze-attentive fixation finding engine. *IEEE transactions on image processing*, 17(4):564–573, 2008. 2

[56] Navyasri Reddy, Samyak Jain, Pradeep Yarlagadda, and Vineet Gandhi. Tidying deep saliency prediction architectures. In *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 10241–10247. IEEE, 2020. 6, 7

[57] Jiacheng Ruan and Suncheng Xiang. Vm-unet: Vision mamba unet for medical image segmentation. *arXiv preprint arXiv:2402.02491*, 2024. 2

[58] Chengyao Shen and Qi Zhao. Webpage saliency. In *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part VII 13*, pages 33–46. Springer, 2014. 2, 1, 3

[59] Hamed R Tavakoli, Fawad Ahmed, Ali Borji, and Jorma Laaksonen. Saliency revisited: Analysis of mouse movements versus fixations. In *Proceedings of the ieee conference on computer vision and pattern recognition*, pages 1774–1782, 2017. 4

[60] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017. 2

[61] Eleonora Vig, Michael Dorr, and David Cox. Large-scale optimization of hierarchical features for saliency prediction in natural images. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2798–2805, 2014. 2

[62] Wenguan Wang and Jianbing Shen. Deep visual attention prediction. *IEEE Transactions on Image Processing*, 27(5):2368–2378, 2017. 1, 6

[63] Zhaohu Xing, Tian Ye, Yijun Yang, Guang Liu, and Lei Zhu. Segmamba: Long-range sequential modeling mamba for 3d medical image segmentation. *arXiv preprint arXiv:2401.13560*, 2024. 2

[64] Juan Xu, Ming Jiang, Shuo Wang, Mohan S Kankanhalli, and Qi Zhao. Predicting human gaze beyond pixels. *Journal of vision*, 14(1):28–28, 2014. 5, 6, 8

[65] Chenhongyi Yang, Zehui Chen, Miguel Espinosa, Linus Ericsson, Zhenyu Wang, Jiaming Liu, and Elliot J Crowley. Plainmamba: Improving non-hierarchical mamba in visual recognition. *arXiv preprint arXiv:2403.17695*, 2024. 2

[66] Yijun Yang, Zhaohu Xing, and Lei Zhu. Vivim: a video vision mamba for medical video object segmentation. *arXiv preprint arXiv:2401.14168*, 2024. 3

[67] Lianghui Zhu, Bencheng Liao, Qian Zhang, Xinlong Wang, Wenyu Liu, and Xinggang Wang. Vision mamba: Efficient visual representation learning with bidirectional state space model. *arXiv preprint arXiv:2401.09417*, 2024. 2, 4

# SUM: Saliency Unification through Mamba for Visual Attention Modeling

## Supplementary Material

## A. Experimental Results

### A.1. Impact of different loss combinations

In addition, to provide additional details about the coefficients used for each loss combination, we conducted several experiments to determine the optimal coefficients for each combination. The best coefficients for each combination are depicted in Table 8.

### A.2. More visualization results

We have included an additional visualization of SUM's predictions in Figure 3. Compared to ground truths, SUM consistently delivers accurate predictions across various image types and datasets, underscoring its robustness and versatility in visual saliency modeling. Moreover, to further validate the robustness of our proposed method, we conducted comparative analyses using publicly available datasets that had not been previously seen, as detailed in Table Table 9. The performance, as depicted in Figure 4, notably remains consistent when applied to new and previously unseen datasets. This suggests that SUM adeptly identifies and highlights the salient features in images, maintaining close alignment with the ground truth data. Therefore, SUM can be reliably utilized in diverse real-world applications where accuracy in visual recognition is critical.

Table 8. loss weighting coefficients $\lambda_i$ ($i = 1, \ldots, 5$) as used in Table 3.

| KL | CC | SIM | NSS | MSE |
|----|----|-----|-----|-----|
| 1 | 0 | 0 | 0 | 0 |
| 0 | -1 | 0 | 0 | 0 |
| 0 | 0 | -1 | 0 | 0 |
| 0 | 0 | 0 | -1 | 0 |
| 0 | 0 | 0 | 0 | 1 |
| 10 | 0 | -3 | 0 | 0 |
| 10 | 0 | 0 | -3 | 0 |
| 10 | -3 | 0 | 0 | 0 |
| 10 | 0 | 0 | 0 | 5 |
| 0 | -2 | 0 | -1 | 0 |
| 10 | -2 | -1 | 0 | 0 |
| 10 | 0 | -3 | 0 | 5 |
| 10 | 0 | 0 | -3 | 5 |
| 10 | -3 | 0 | 0 | 5 |
| 10 | -2 | -1 | -1 | 0 |
| 10 | -2 | -1 | 0 | 5 |
| 10 | -2 | -1 | -1 | 5 |

Table 9. Details of unseen datasets used for quantitative analysis of SUM in Figure 4.

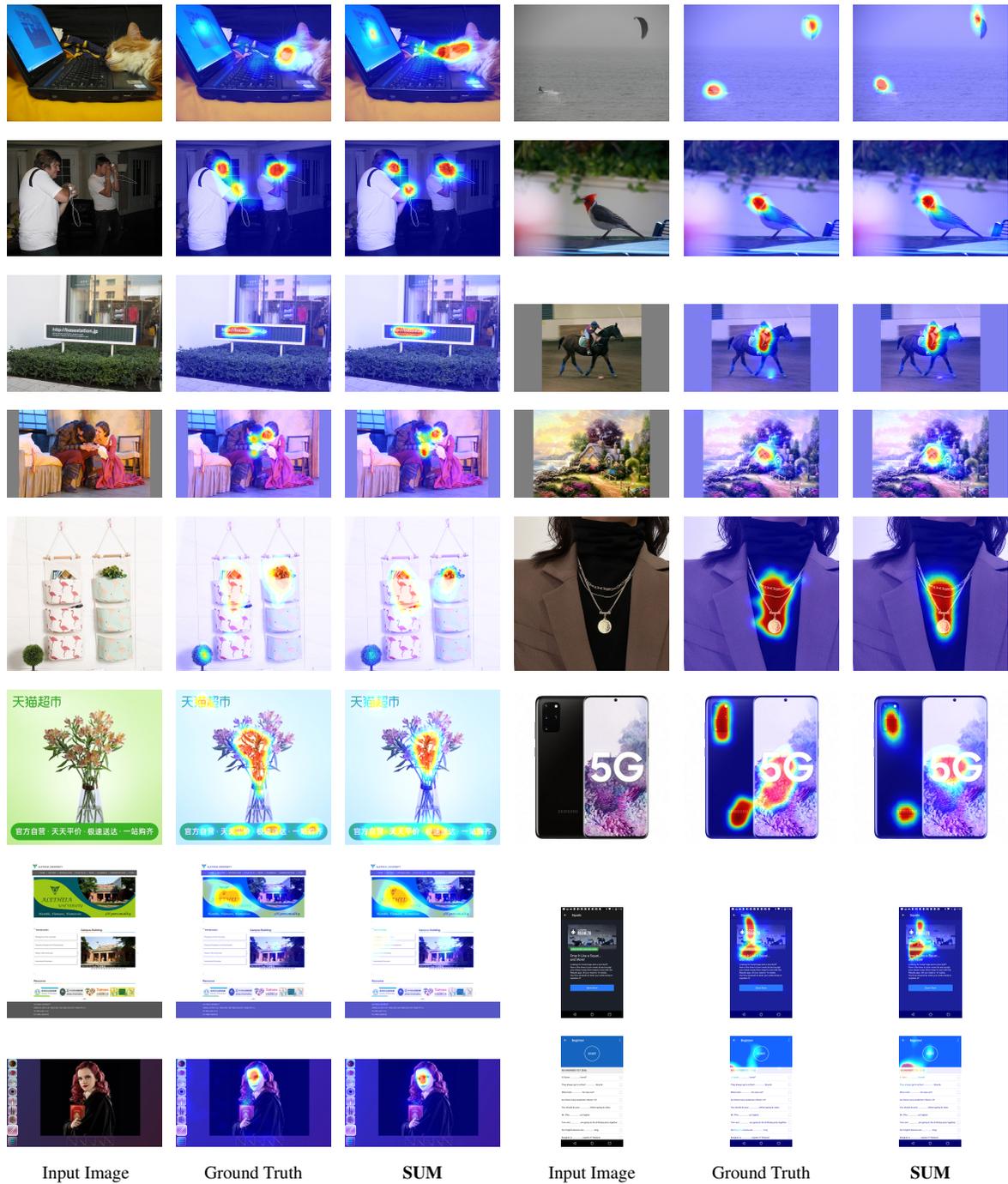| Dataset | Image domain | # Image | Image Resolution |
|---------|-------------|---------|------------------|
| *Toronto* [4] | Natural scene | 120 | $681 \times 511$ |
| *TUD Image Quality Database 1* [45] | Natural scene | 29 | $768 \times 512$ |
| *TUD Image Quality Database 2* [1] | Natural scene | 160 | $600 \times 600$ |
| *FIWI* [58] | Web page | 149 | $1360 \times 768$ |

Figure 3. Visualizations of SUM's predictions across different datasets. The first and second rows depict Natural Scene-Mouse data, while the third and fourth rows showcase Natural Scene-Eye data. The fifth and sixth rows present E-commerce data, and the seventh and eighth rows display UI data.
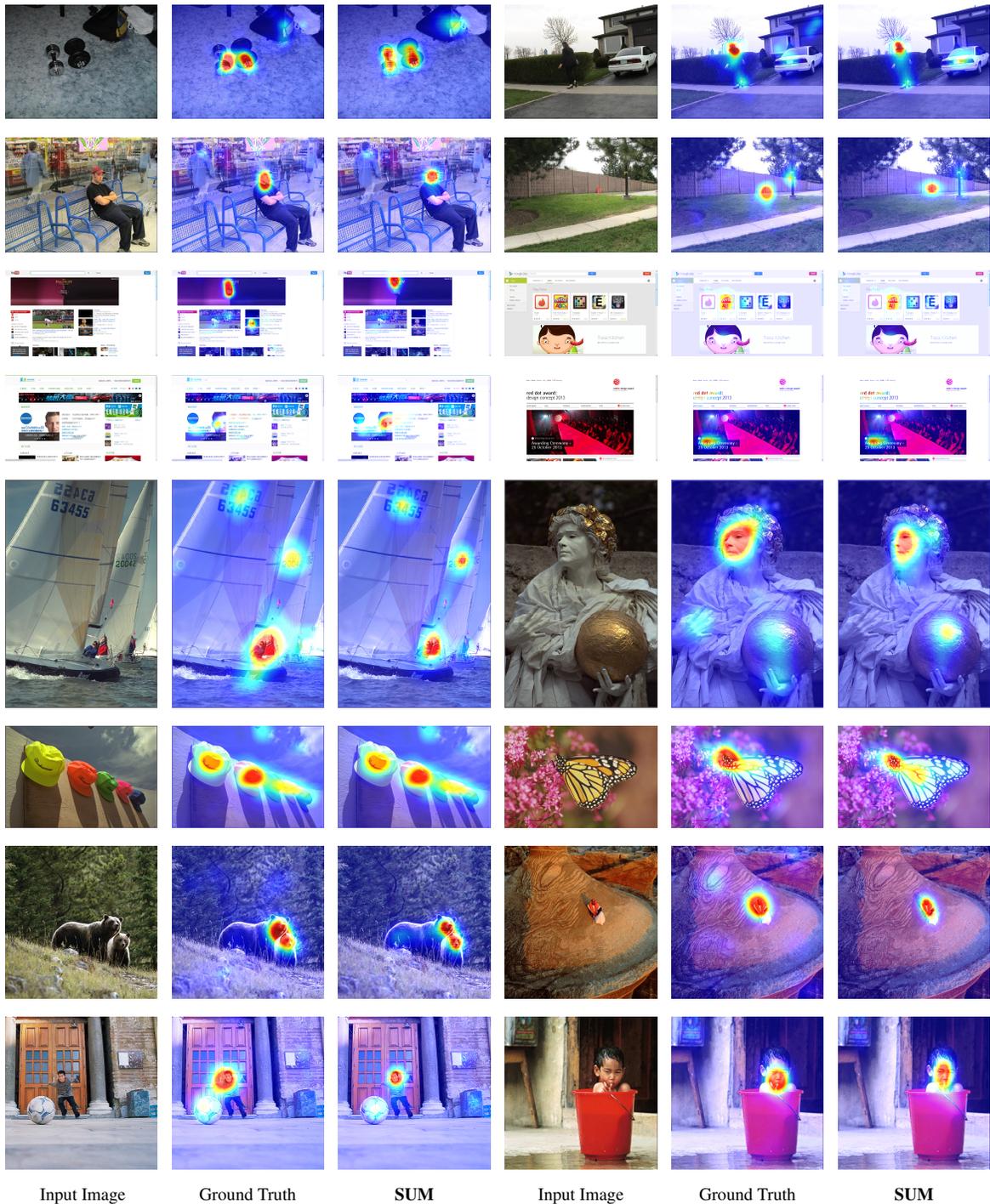
| Input Image | Ground Truth | **SUM** | Input Image | Ground Truth | **SUM** |

Figure 4. Visualizations of SUM's predictions across different datasets. The first and second rows showcase the Toronto dataset [4], while the third and fourth rows present the FIWI dataset [58]. The fifth and sixth rows display data from the TUD Image Quality Database 1 [45], and the seventh and eighth rows exhibit data from the TUD Image Quality Database 2 [1].