

Intention-aware Feature Propagation Network for Interactive Segmentation

Chuyu Zhang^{1,3,4}, Chuanyang Hu¹, Yongfei Liu^{1,3,4}, and Xuming He^{1,2}

¹ ShanghaiTech University

² Shanghai Engineering Research Center of Intelligent Vision and Imaging

³ Shanghai Institute of Microsystem and Information Technology, Chinese Academy of Sciences

⁴ University of Chinese Academy of Sciences

{zhangchy2, hucy3, liuyf3, hexm}@shanghaitech.edu.cn

Abstract. We aim to tackle the problem of point-based interactive segmentation, in which two key challenges are to infer user’s intention correctly and to propagate the user-provided annotations to unlabeled regions efficiently. To address those challenges, we propose a novel intention-aware feature propagation strategy that performs explicit user intention estimation and learns an efficient click-augmented feature representation for high-resolution foreground segmentation. Specifically, we develop a coarse-to-fine sparse propagation network for each interactive segmentation step, which consists of a coarse-level network for more effective tracking of user’s interest, and a fine-level network for zooming to the target object and performing fine-level segmentation. Moreover, we design a new sparse graph network module for both levels to enable efficient long-range propagation of click information. Extensive experiments show that our method surpasses the previous state-of-the-art methods on all popular benchmarks, demonstrating its efficacy.

1 Introduction

Interactive image segmentation plays a vital role in a broad range of human-in-the-loop vision tasks, such as image editing [8], medical image analysis [32] and dense image annotation [30]. There has been a long history of interactive segmentation in vision literature, in which a variety of interaction strategies have been explored, including points [36,30], scribbles [3,2], and bounding boxes [27]. In this work, we mainly focus on the point-based interactive segmentation that only provides point clicks to indicate foreground or background on image, which typically requires less effort from human annotators.

A key characteristic of interactive segmentation is the diversity in the regions of user’s interest. Unlike semantic segmentation with a predefined label space, each instance in an interactive segmentation task may produce different foreground regions according to the user’s intent. Consequently, the main challenges of point-based interactive segmentation are to figure out the user intention [31] and to propagate the user annotation (i.e., clicks) to other unlabeled pixels [7] based on a limited number of user interactions.

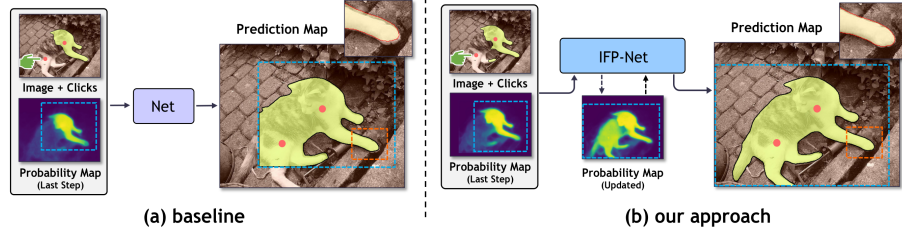


Fig. 1: Comparison of baseline with our approach. The blue box and orange box represent estimated user intention and zoomed-in details. The red points represent positive user clicks. The green finger represents user’s current click action. (a) The baseline prediction misses a part of the cat due to inaccurately estimated intention. (b) Our approach captures complete user intention by intention update and preserves more accurate object boundary.

Most previous work focus on modeling each individual interactive step and mainly rely on the sparsely-annotated clicks to infer the user intention [30,7,21]. Such a memoryless strategy tends to be less efficient in capturing foreground as it ignores predictions from previous steps. To alleviate this, recent work [31] takes the last-step prediction as an additional input to zoom into inferred target regions, leading to more effective foreground estimation. However, such a simple fusion method may produce erroneous information propagation due to potential inaccurate mask predictions in previous steps (as shown in Fig. 1). For the challenge of label propagation, previous methods typically utilize stacked convolutions to propagate user annotation in an implicit manner [36,21]. This has a limited capacity in capturing long-range dependency and often leads to incomplete foreground masks. More recently, Chen et al. [7] adopt fully-connected graph networks [33] to facilitate information propagation from both global and local perspective. Nonetheless, the fully-connected graph network can only cope with relatively low-resolution feature map due to its high-computation complexity, which can result in inaccurate foreground boundaries. Moreover, both propagation stages are susceptible to noisy affinity estimation that relies on previous predictions or color similarity.

To address the aforementioned limitations, we propose a novel intention-aware feature propagation strategy for interactive image segmentation. Our main ideas are two-folds: 1) to *learn a click-augmented feature representation* based on a sparse graph neural network (GNN), which allows efficient long-range information propagation at high spatial resolution, and thus enables us to generate more accurate object boundaries; and 2) to *improve user intention estimation* by integrating previous foreground prediction and current user clicks at each step, which can better localize target regions and hence results in an effective zoom-in mechanism for coping with scale variation.

To this end, we develop a coarse-to-fine sparse propagation network for each interactive segmentation step, consisting of a coarse-level network for foreground estimation and a fine-level network for detailed segmentation. Specifically, at each step, our coarse-level network first generates a refined foreground mask based on the entire image, user clicks and the last-step mask prediction. The resulting mask, which encodes the user intention more accurately, is used to zoom in to a selected region of interest. We then use the fine-level network to perform foreground segmentation at a fine resolution on the

zoomed-in region, of which the output is remapped onto the full image and sent to the next step. In each cascade stage, we design a new sparse graph network to propagate the user click information to the unlabeled region in a non-local and yet efficient manner. In particular, our sparse GNNs compute a set of click-augmented feature representations at high spatial resolution in linear complexity, which are highly efficient and can preserve more detailed information for foreground mask predictions.

We adopt a stage-wise training strategy for our network cascade which trains the coarse- and fine-level networks sequentially by simulating interactive steps [31]. We conduct extensive experiments on GrabCut, Berkeley, DAVIS, COCO, SBD, and Pascal3D+ datasets, and the results demonstrate that our method achieves the state-of-the-art performance. To summarize, our contribution is three-folds:

- We propose an effective intention-aware feature propagation strategy for interactive image segmentation, which employs a cascaded network for estimating user intention and performing long-range information propagation.
- We develop a novel click propagation strategy based on two sparse GNNs, capable of capturing long-range dependency on high-resolution feature maps and generating more accurate object boundaries.
- Our method achieves the state-of-the-art results on most public benchmarks, demonstrating the effectiveness of our design.

2 Related works

2.1 Interactive segmentation

Interactive image segmentation has attracted much attention in computer vision research, and a variety of interaction strategies have been studied, which are based on bounding boxes, scribbles, or points. While the bounding-box based methods [27,37,34] can localize the target object quickly, and the scribble based methods [11,3,12,2] provide richer user-input cues, they often involve more user interactions. By contrast, the point-based, where a user provides points to indicate foreground or background on the image, requires less effort from human annotators [36,30,31,7,4,13]. Consequently, we mainly focus on the point-based methods in the discussion below.

Many point-based works have emerged since Xu et al [36] first propose a CNN-based method, which can be largely grouped into two categories. While one trend focuses on annotating object boundaries [5,17,1,22], most deep learning methods perform region-based segmentation, aiming to better leverage user clicks in each interactive step. In particular, Liew et al [19] attempt to refine local regions based on pairs of positive and negative clicks. Majumder et al [24] generate a content-aware guidance map for exploiting the hierarchical structural information in the image. Lin et al [21] argue that the first click is more important than others and design a first-click attention mechanism. Hao et al [13] improve the usage of interactive information from user clicks with edge-guided flow. To better adapt to test cases, recent work [15,30] develop a backpropagating refinement scheme to correct the mislabeled user clicks in the test time. Despite their promising performances, these approaches typically concentrate on a single interactive step and hence are inefficient in capturing user intention. To remedy this, Sofiiuk

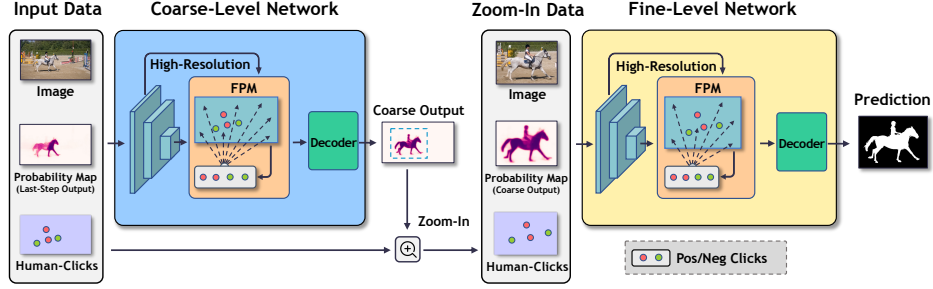


Fig. 2: The overall architecture of intention-aware feature propagation network. Our pipeline consists of a coarse-level and fine-level network. In each interactive step, we first update foreground estimation at image level by coarse-level network (blue box). The updated foreground probability map uses to zoom in onto user interested region. Subsequently, the fine-level network (yellow box) performs foreground segmentation at zoomed-in level. In each network, feature propagation module (FPM) is used to propagate user-provided information to the unlabeled regions.

et al [31] utilize the last-step prediction to facilitate localizing the target region, which however tends to produce erroneous information propagation due to inaccuracy in previous predictions. By contrast, we propose a coarse-to-fine strategy that better tracks the user intention from the coarse-level and perform effective segmentation based on both global and local image cues.

2.2 Graph neural network

Graph neural networks (GNN) [10,28] have been widely adopted to model long-range dependencies and there exists a large body of literature on this research topic [35]. However, only a few works utilize GNNs in the task of interactive segmentation. In particular, based on the non-local networks [33], Chen et al[7] propose a conditional diffusion network for interactive segmentation, which performs non-local feature propagation on the global convolutional features and local-level pixels using color similarity. This mixed strategy tends to suffer from inaccurate foreground boundaries due to the low-resolution deep feature map and/or the noisy graph affinity estimated based on previous foreground masks or color similarity. In contrast, we propose a simple sparse attention-based non-local graph network to propagate the click information, which can be applied to a high-resolution feature map and generate more accurate masks.

3 Method

In this section, we introduce our intention-aware feature propagation method for point-based interactive image segmentation. To this end, we first present the problem setup of interactive segmentation and outline our method in Sec. 3.1. Subsequently, we detail the design of two main components, including intention-aware framework (IAF) and feature propagation module (FPM) in Sec. 3.2 and Sec. 3.3, respectively. Finally, we formally describe the training strategy of our feature propagation network in Sec. 3.4.

3.1 Problem Setup and Method Overview

The goal of interactive segmentation is to correctly infer the region of user’s interest and segment the target object with as few clicks as possible. This sequential estimation task is typically converted into a series of foreground segmentation problems, each of which aims to output a foreground mask as accurate as possible given the current set of user’s clicks in an interactive step.

In this work, we consider a learning strategy that trains a neural network to predict the foreground mask in each step [36]. Formally, we assume a set of image-click pairs is generated as our training data \mathcal{D} (See Sec. 3.4 for more detail). It comprises data tuples $\{(I_i, U_i, Y_i)\}_{i=1}^{|\mathcal{D}|}$ where I_i indicates an input image, $U_i = \{(u_{i,j}, l_{i,j})\}_{j=1}^{M_i}$ is a set of user clicks represented by their pixel indices $u_{i,j}$ and labels $l_{i,j} \in \{\text{pos}, \text{neg}\}$, M_i is the number of clicks and Y_i denotes the groundtruth mask. Our goal is to learn a deep network \mathcal{M}_θ based on \mathcal{D} , which can be formulated as follows:

$$\min_{\theta} \mathbb{E}_{(I_i, U_i, Y_i) \sim \mathcal{D}} \mathcal{L}(\mathcal{M}_\theta(I_i, U_i), Y_i) \quad (1)$$

where \mathcal{L} denotes the loss function, and θ indicates the network parameters.

To tackle the problem in Eq (1), we focus on two key aspects of click-guided foreground mask prediction: 1) to capture the user’s intent effectively and 2) to efficiently propagate the click labels to unlabeled regions. To this end, we develop a novel intention-aware feature propagation network, which performs a coarse-to-fine foreground refinement at each interaction as shown in Fig.2. Specifically, our network consists of two main modules: a coarse-level network in the first stage which updates the foreground mask on the full image based on the last-step prediction and the up-to-date click inputs, and a subsequent fine-level network which zooms into the foreground region adaptively and performs segmentation at a fine resolution. Both stages employ newly-designed sparse GNNs to facilitate propagation of the user’s input information. Below, we will first introduce the overall intention-aware pipeline, followed by the details of our feature propagation module.

3.2 Intention-Aware Framework (IAF)

We now present our intention-aware propagation framework which aims to better infer the region of user’s interest and produce a foreground segmentation with accurate object boundaries. To achieve this goal, we adopt a coarse-to-fine strategy and develop a network cascade architecture for each interactive step as shown in Fig.2.

Our network cascade first use a coarse-level network to fuse the last-step foreground map with the current user clicks and to update the foreground estimation at the image level. Formally, at each step t , we denote the input image as I , the current-step user click set as U^t and the last-step network prediction as P^{t-1} . The coarse-level network takes those three inputs and generates an updated foreground probability map for the entire image, which can be formulated as follows:

$$P_c^t = \mathcal{M}_c(I, P^{t-1}, F_{enc}(U^t); \theta_c) \quad (2)$$

where \mathcal{M}_c and θ_c denote the coarse-level network and its parameters respectively, and P_c^t is the output probability map; F_{enc} is an encoding function for representing user’s

click inputs, which generates two image-sized heatmaps for the positive and negative clicks respectively. Here we adopt a common encoding strategy [36,4] in which each click is represented by a click-centered Gaussian kernel or a disk transform function with a fixed radius.

Subsequently, a fine-level network as the second module, zooms in onto the foreground region and performs segmentation with a fine-level feature representation. Specifically, according to the image-level foreground estimation, we adaptively crop and resize the image, the click encoding maps and the foreground map, which are then fed into the fine-level network. The fine-level network generates a foreground probability map for the zoomed-in region and map it back to the full image view. This second-stage process can be written as follows,

$$\tilde{I}, \tilde{P}_c^t, \tilde{U}^t = \text{ZoomIn}(I, P_c^t, U^t) \quad (3)$$

$$P^t = \mathcal{M}_f(\tilde{I}, \tilde{P}_c^t, F_{enc}(\tilde{U}^t); \theta_f) \quad (4)$$

where $\tilde{I}, \tilde{P}_c^t, \tilde{U}^t$ denotes the image, foreground map and user clicks after the zoom-in step, P^t is the output of the fine-level network \mathcal{M}_f and θ_f are its parameters.

The function *ZoomIn* determines the region for cropping according to the coarse-level foreground prediction P_g^t . In particular, we first derive a bounding box of the target from binarizing P_g^t , and then expand the box by an adaptive margin. The margin size is inversely proportional to the size of the initial box so that large objects have a relatively tight bounding box which excludes distracting background while small objects take a relatively loose bounding box for including more context cues.

Our coarse-to-fine framework provides several key benefits for inferring foreground regions. First, it begins with an update of the foreground probability in the full image view, leading to a more accurate localization of the target region. This enables us to zoom in on small objects to extract more details and track the user’s intention in a more reliable manner. Moreover, the updated foreground provides a high-quality initial labeling for the fine-level network, which is essential for refining the segmentation in local regions.

3.3 Feature Propagation Module (FPM)

Both coarse- and fine-level network aim to propagate the user-provided information to unlabeled regions in their own input images. Due to the sparsity of user’s clicks, this typically requires modeling long-range feature relations across the image plane. To achieve efficient information propagation, we introduce a new graph neural network module, which augments a base CNN segmentation network for predicting the foreground mask. In contrast to previous non-local design [33,7], our graph network module are built on a sparse graph topology, which enables us to compute a user-input-aware representation on a high-resolution feature map and hence produces detailed foreground segmentation with accurate boundaries.

Specifically, we first use a CNN-based segmentation network to compute a stack of feature maps, from which a higher-level and a lower-level map are selected. The higher-level feature map, denoted as \mathcal{F} , typically encodes more semantics but has a

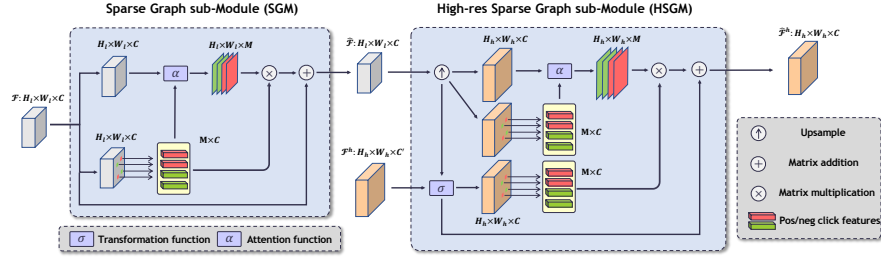


Fig. 3: The structure of sparse graph sub-module (SGM) and high-resolution sparse graph sub-module (HSGM) in FPM. H and W represent the height and width of the feature map. C and C' represent the channel of higher-level and lower-level feature maps, respectively. M is the number of user clicks. For simplicity, we ignore the reshape operation.

low resolution while the lower-level map, denoted as \mathcal{F}^h , has a high resolution and preserves more object boundary cues⁵. In order to better exploit both the higher- and lower-level features, we employ a two-stage design: a sparse graph network first augments the higher-level feature map with the user-click features, which is further integrated with the lower-level features in a high-res sparse graph network at the second stage. Below we describe the details of those two graph networks in turn.

Sparse Graph sub-Module (SGM): Our sparse graph submodule performs feature propagation on the higher-level feature maps to disseminate the user-click information to all features. Formally, we represent the higher-level feature map \mathcal{F} as $\{f_n\}_{n=1}^{H_l \times W_l}$, where H_l and W_l are the height and width of the low-res feature map respectively. At each location, $f_n \in \mathbb{R}^c$ is a c -channel feature vector and n is the spatial location index. To build the sparse graph, we select the feature vectors at the location of user clicks in \mathcal{U} , denoted as $\mathcal{F}_u = \{f_{u_i}\}_{i=1}^M$, and connect them to each location on the feature map.

Given the graph, we perform feature augmentation by passing messages from the click nodes \mathcal{F}_u to each feature location as follows (See Fig. 3 for illustration):

$$\hat{f}_n = f_n + \sum_{j=1}^M \alpha(f_n, f_{u_j}) W_c^T f_{u_j}, \quad \forall f_n \in \mathcal{F} \quad (5)$$

$$\alpha(f_n, f_{u_j}) = e^{\theta(f_n)^T \Phi(f_{u_j})} / Z_n(\mathcal{F}) \quad (6)$$

where \hat{f}_n is the updated feature, $W_c \in \mathbb{R}^{c \times c}$ is a weight matrix for feature transform, and α denotes an attention function in which θ and ϕ are linear transforms and $Z_n(\mathcal{F})$ is the normalization factor. Intuitively, the augmentation moves all the features towards

⁵ In this work, we adopt the DeeplabV3+[6] with a ResNet backbone as our base segmentation network. We select the feature map after the ASPP block as the higher-level \mathcal{F} , and the conv features after the first block in the ResNet as the lower-level \mathcal{F}^h . The size of \mathcal{F} and \mathcal{F}^h are 1/4 and 1/2 of the original image.

the click-annotated representations, which reduces in-class variation and improve foreground prediction. However, these augmented features, denoted as $\hat{\mathcal{F}}$, are built on the low-res feature map \mathcal{F} , which tends to produce coarse segmentation masks. To remedy this, we introduce a second graph network to refine them as below.

High-res Sparse Graph sub-Module (HSGM): The second graph network generates a click-augmented feature representation with a high spatial resolution. To this end, we integrate the first-stage output $\hat{\mathcal{F}}$ with the lower-level feature map \mathcal{F}^h , followed by another pass of click-to-feature propagation. Specifically, we denote the lower-level feature map as $\{f_n^h\}_{n=1}^{H_h \times W_h}$ where H_h and W_h are the height and width of the high-res feature map respectively. Similar to the SGM, we select the click-annotated feature, represented as $\mathcal{F}_u^h = \{f_{u_j}^h\}_{j=1}^M$, and link them to every feature location on \mathcal{F}^h .

Given the high-res graph, we first upsample the previous output $\hat{\mathcal{F}}$ so that it has the same spatial dimension as the lower-level feature map \mathcal{F}^h . We then perform feature integration and click-aware augmentation by a message passing as shown in Fig. 3. Formally, denoting the upsampled feature as $\{\hat{g}_n\}_{n=1}^{H_h \times W_h}$, we update the high-res feature representation as follows,

$$\hat{f}_n^h = \sigma(f_n^h \oplus \hat{g}_n) + \sum_{j=1}^M \alpha(\hat{g}_n, \hat{g}_{u_j}) W_f^T \sigma(f_{u_j}^h \oplus \hat{g}_{u_j}) \quad (7)$$

where \hat{f}_n^h denotes the high-res augmented feature, \oplus indicates feature concatenation, $\sigma(\cdot)$ is a transformation function and W_f is a weight matrix for feature transform. Note that we use the higher-level features $\{\hat{g}_n\}$ to compute the attention weights so that the information propagation is less susceptible to variations in the lower-level feature \mathcal{F}^h . Moreover, our sparse graphs only perform message passing from M selected nodes to N feature nodes, which can be computed efficiently with a complexity $\mathcal{O}(MN)$ where $M \ll N$ in the interaction.

Given the high-res augmented features $\{\hat{f}_n^h\}_{n=1}^{H_h \times W_h}$, we finally generate the foreground probability map by applying a two-layer conv-block followed by the Sigmoid function.

3.4 Model Training

To train our deep network for interactive segmentation, we first utilize a simulation process to generate a set of image-click pairs from a foreground segmentation dataset as in [31], which assumes the user clicks on the center of maximum error regions. Given the training dataset, we then adopt a stage-wise strategy to train our coarse-level and fine-level networks in turn.

Specifically, we first train the coarse-level network, which is then kept fixed during the second-stage training. The fine-level network is then initialized with the coarse-level model parameters and fine-tuned afterwards. Following [30], both networks employ the Normalized Focal Loss (NFL) [29] in their training objectives as shown in Eqn (1).

4 Experiments

In this section, we first depict the experiment setting and implementation details, then compare our model with existing works, followed by ablation study to validate each

component. Finally, we demonstrate some qualitative results to show the model efficacy, and analyze our model’s parameters and inference time.

4.1 Evaluation and Implementation Details

Datasets: We evaluate our method over a wide range of datasets including GrabCut, Berkeley, DAVIS, COCO, PASCAL VOC and SBD, by following the standard evaluation protocol.

GrabCut[27] is a typical interaction segmentation dataset, which contains 50 images with distinguishable foreground and background.

Berkeley[25] contains 96 images with 100 object masks from its test subset.

DAVIS[26] is originally introduced for video segmentation. Only 345 randomly sampled frames with finely labeled objects are used in our method by following [15].

COCO[20] is a typical semantic segmentation dataset, containing more complex scene and multiscale objects. Following [36], we split the dataset into COCO(seen) and COCO(unseen) according to their object class whether in PASCAL VOC or not. And finally, 10 images are sampled randomly for each category. For simplicity, we denote COCO(seen) and COCO(unseen) as COCO^s and COCO^u.

PASCAL VOC[9] contains 1449 images with 3427 object masks from its validation set. We ignore evaluating the precision of object boundaries since they are marked.

SBD[14] contains 6671 object masks for 2820 images.

Metric: To mimic the real user clicks in evaluation, we follow [36] to click the center of maximum error region to correct the output mask continuously. The interaction process will terminate when the IoU between prediction and groundtruth mask exceeds threshold τ , or reaching the maximum number of interactions. Due to large scene diversity in different datasets, we typically set τ as 85% or 90%, and set maximum number of interactions to 20 as previous works. The number of clicks (NoC) and number of failure (NoF) to meet the termination conditions is used for evaluation metric. For example, NoC@90 means the average number of clicks for test set is needed to reach 90% IoU under 20 maximum interactions, and NoF@90 means the number of failure case that doesn’t reach 90% IoU with 20 maximum interactions.

Due to the randomness of the training, we report the average NoC by running three times for the same experimental setting to obtain more stable results

Implementation Details: We adopt DeeplabV3+[6] as our coarse- and fine-level network. Besides, we follow [31] to utilize Conv1S to fuse click maps and foreground estimation, then sum the fused feature with image features at the output of the first convolutional block.

During training, we follow the same iterative sampling strategy in [31] to generate positive and negative clicks to alleviate the gap between training and inference stages. For clicks encoding, we adopt the disk encoding strategy proposed in [4] with a fixed radius of 5. Our network is trained over the SBD training set and COCO+LVIS dataset with common data augmentation approaches, including random cropping and scaling, horizontal flipping, and optimized by Adam optimizer with batch size 28. The image is resized with a size 320×480 after augmentation.

We adopt the stage-wise training strategy for our network. Specifically, for the first training stage, the coarse-level network is trained for 120 epochs and utilized to initialize the fine-level network. The learning rate is set as $5e-4$ and decays 10 times at 100th,

Method		GrabCut	Berkeley	COCO*	COCO ^u	PascalVOC	DAVIS	SBD
		NoC@90	NoC@90	NoC@85	NoC@85	NoC@85	NoC@85 / 90	NoC@85 / 90
DOS[36]	FCN	6.08	8.65	8.31	7.82	6.88	9.03 / 12.58	9.22 / 12.80
RIS[19]	-	5.00	6.03	5.98	6.44	5.12	- / -	6.03 / -
LD[18]	-	4.79	-	-	-	-	- / 9.57	- / -
ITIS[23]	-	5.60	-	-	-	3.80	- / -	- / -
BRS[15]	DenseNet	3.60	5.08	-	-	-	5.58 / 8.24	6.59 / 9.78
CMG[24]	FCN	3.58	5.60	5.40	6.10	3.62	- / -	- / -
CDNet[7]	ResNet-50	2.64	3.69	-	-	-	5.17 / 6.66	4.37 / 7.87
Ours	ResNet-50	<u>2.31</u>	<u>3.35</u>	<u>2.46</u>	<u>3.69</u>	<u>2.38</u>	<u>4.52 / 5.83</u>	<u>3.08 / 4.98</u>
IS+SA[16]	ResNet-101	3.07	4.94	4.08	5.01	3.18	5.16 / -	- / -
FCA [21]	ResNet-101	2.14	4.19	4.45	5.33	2.96	- / 7.90	- / -
f-BRS-B[30]	ResNet-101	2.72	4.57	-	-	-	5.04 / 7.41	4.81 / 7.73
RITM*[31]	ResNet-101	2.31	3.50	2.54	3.61	2.48	5.09 / 6.78	3.33 / 5.33
Ours	ResNet-101	2.15	<u>3.20</u>	<u>2.27</u>	<u>3.50</u>	<u>2.31</u>	<u>4.51 / 5.8</u>	2.98 / 4.83
RITM[31]	HRNet-18	2.04	3.22	2.40	3.61	2.51	4.94 / 6.71	3.39 / 5.43
Ours	HRNet-18	<u>1.75</u>	<u>2.86</u>	<u>2.23</u>	<u>3.00</u>	<u>2.49</u>	<u>4.68 / 6.01</u>	<u>3.33 / 5.25</u>
RITM[31] [†]	HRNet-18	1.54	2.26	-	-	2.28	4.36 / 5.74	3.80 / 6.06
Ours [†]	HRNet-18	1.68	2.12	2.17		2.40	4.03 / 5.22	<u>3.67 / 5.91</u>

Table 1: Comparison with SOTA. The bold means the best results across different backbones, and the underline means the best results under the same backbone. * means that we implement the results. [†] denotes the model are trained on COCO+LVIS, and others model are trained on SBD dataset. - denotes the results are not available. **The lower is the better.**

115th epochs. For the second stage training, we fix the parameters in coarse-level network and train our fine-level network for another 30 epochs with learning rate 5e-7 on the zoomed-in image.

4.2 Quantitative Results

We compare our framework with several existing approaches for evaluation, in terms of the following metrics: standard average number of clicks (NoC) and results with maximum 100 clicks.

Main results: As shown in Tab.1, we compare our method with previous approaches over a wide range of benchmarks covered by existing works. When trained on the SBD dataset, our method achieves the strongest performance over different datasets, and outperforms all existing approaches under the same backbone setting. For Berkeley and COCO datasets, our method achieves the best average number of clicks under HRNet-18 backbone, which only requires average **2.86** clicks to reach 90% IoU on Berkeley, **3.0** clicks to reach 85% IoU on COCO* datasets. The result on COCO^u shows that our framework can generalize to unseen classes greatly. Specifically, for the challenging SBD dataset and fine-grained DAVIS dataset, we achieve average **5.8** and **4.83** NoC@90 respectively on the ResNet-101 backbone, improve 0.98 and 0.5 compared

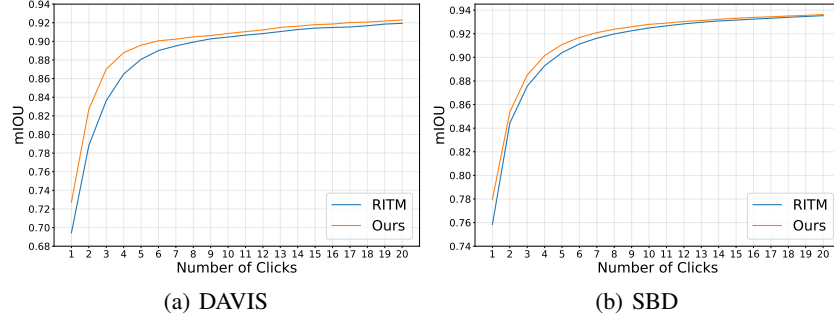


Fig. 4: mIOU varying over number of clicks k .

with RITM[31]. It indicates that our intention-aware framework with a click information propagation module can better segment objects with detailed boundaries. For GarbCut and PascalVOC datasets, we improve a little since the results are close to the upper bound. It is worth noting that the difficulty of interactive segmentation increases rapidly along with the number of clicks decreasing. Therefore, the improvements on all datasets are significant.

For further comparisons, we train our model on the COCO+LVIS dataset. The results in the last two lines show that we still achieve a great improvement on fine-grained DAVIS. Meanwhile, same as RITM[31], we also observe a drop on SBD datasets compared with trained on SBD datasets. The improvement or drop on Berkeley, Pascal VOC, and GrabCut is slight because their performance are nearly saturated.

IOU@ k Analysis: We analysis the performance varying over number of clicks k . As shown in the Fig. 4, our method is consistently better than RITM, and about 2% higher in mIOU for the first five clicks, which validates the effectiveness of our method.

Success samples with different number of clicks: To further analysis our improvement, we report the distribution of success samples with different number of clicks on the challenging fine-grained DAVIS dataset, which contains 345 samples in total. As in Fig. 5, our approach can successfully segment 71.3% (246) samples within **5 clicks**, which greatly outperforms the RITM by 8.1%. It indicates that we can better utilize user-provided sparse click information. Also, we can know hard case in the tail drags down NoC metrics.

Results with maximum 100 clicks: Following[30], we report NoC with the maximum number of clicks limited to 100 on the DAVIS dataset with ResNet-50 backbone. In Tab. 2, it shows that our method improves a lot on NoC@90. And we improve 0.91 on NoC₁₀₀@90 metric compared with CDNet. This is because our sparse graph can be propagated on high-resolution feature maps to preserve more detail information, while CDNet’s cannot.

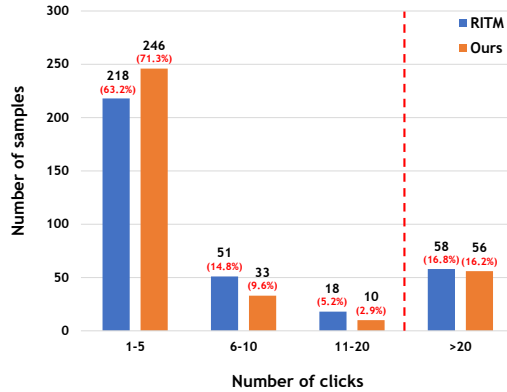


Fig. 5: The distribution of the number of clicks on DAVIS dataset. The experiments are on the ResNet-50 backbone. The left of red dotted line is the success samples and the right of red dotted line is the failed samples, which can't reach 90% IoU within 20 clicks. 1-5 means the number of samples that need at least one click and at most five clicks to reach 90% IoU.

Method	NoC@90	NoC ₁₀₀ @90	NoF ₁₀₀ @90
f-BRS[30]	78	20.70	50
CDNet[7]	65	18.59	48
Ours	56	17.68	46

Table 2: Analysis of 100 clicks.

Method	Components		DAVIS SBD	
	FPM	IAF		
Baseline	-	-	6.85	5.67
	✓	-	6.39	5.36
	-	✓	6.39	5.34
Ours	✓	✓	6.05	5.13

Table 3: Each component analysis.

4.3 Ablation study

In this section, we perform several ablative experiments to evaluate the effectiveness of each component in our approach. All ablation experiments are evaluated using NoC@90 metric on DAVIS and SBD datasets, which are more challenging than other datasets.

Baseline: We adopt previous state-of-the-art RITM[31] on ResNet-34 backbone as our baseline which achieves decent performance as shown in Tab.1.

Effectiveness of each component: The results in Tab.3 show that each component plays a critical role in the entire framework and contributes to the final results. For the *Feature propagation module (FPM)*, it achieves 0.46 and 0.31 NoC@90 improvements on DAVIS and SBD dataset respectively. The improvement on fine-grained DAVIS dataset reveals that the FPM can better capture accurate boundary details. Besides, the gain on SBD dataset indicates that the FPM can accurately segment objects with severe scale variation. Further, the performance is further improved on two datasets by incorporating our *Intention-aware framework (IAF)*, which indicates that the IAF can better estimate foreground in coarse-level and segment the object accurately in fine-level.

Sparse Graph Analysis: To validate the efficacy of sparse graph design in feature propagation module (FPM), we conduct more experiments to compare with feature diffusion

Method	Params(M)	Flops(G)	NoC ₂₀ @85	NoC ₂₀ @90
Baseline*	31.4	508.72	6.60	8.42
Baseline* + FDM [7]	31.42	513.82	5.40	7.64
Baseline* + FDM in both low & high res.	31.44	1510.16	✗	✗
Baseline* + FPM	31.5	531.42	5.05	7.17

Table 4: Graph design analysis on the DAVIS dataset. Low resolution denotes 1/4 feature map while high resolution represents and 1/2 scale feature map.

graph module (FDM) in CDNet [7], which is actually a fully-connected graph network. For a fair comparison, all experiments are built upon the Baseline* adopted in CDNet. As in Tab 4, we can observe our ‘Baseline* + FPM’ achieves significant improvement in both metrics. It is worth noting that our advanced sparse graph design makes it feasible to conduct message propagation in both low & high resolution feature maps. In addition, when applying FDM in both low & high scale feature maps, we find that the overall network consumes amount of computation (1510.16 GFLOPs v.s. 531.42 GFLOPs), which is unacceptable in real systems.

Feature Propagation Module (FPM): We conduct incremental ablation experiments to study the effectiveness of each component in the FPM on the baseline. As shown in Tab.5, the introduction of *Sparse Graph sub-Module (SGM)* brings 0.23 and 0.1 NoC@90 improvements on DAVIS and SBD respectively. It indicates that the feature propagation of SGM is helpful. Moreover, the introduction of HSGM improves performance further, which means that feature propagation on high-resolution does help to preserve more precise boundary information. The results show that each module in FPM is effective.

Specifically, we analyze the impact of the FPM module when there are all positive clicks. Since it is difficult to compare the case where both models have N positive points, we compare the case with only one positive point. In such case, Baseline+FPM improves mIOU on DAVIS (0.68 \rightarrow 0.71) compared to Baseline. Conceptually, if the user inputs N consecutive positive clicks, it indicates that the model tends to label the target object as background. In such cases, adding positive features to all points would shift the pixel representations towards the foreground side so that the model can predict more pixels as foreground, which could be beneficial.

4.4 Visualization analysis

As shown in Fig.6, we visualize several samples to show the effectiveness of our method. In Fig.6(a), the blue rectangle represents the user intention estimated by the model. As we can see, the baseline focus on the incomplete object region due to the inaccurate user intention encoded in the last step prediction. Our method can correctly focus on the entire object area and obtain better segmentation results since our model updates user intention at coarse-level. In Fig.6(b), we observe our method can obtain more precise boundaries even with fewer user clicks. It shows the effectiveness of the coarse-to-fine strategy and FPM module, which help us to preserve more accurate boundary information. In Fig.6(c), the results indicate that our sparse graph neural network is able to capture more reliable long-range dependencies and make the prediction more complete.

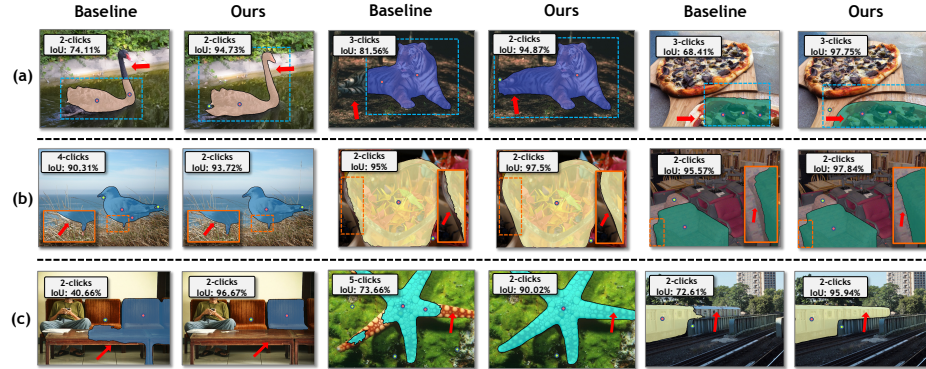


Fig. 6: **Visualization analysis:** The odd and even columns show the prediction result of the baseline[31] and our method. Row (a) indicates that our method maintain more complete user intention, Row (b) indicates that our method preserves more accurate boundary information, and Row (c) indicates that our method captures more reliable long-range dependencies

#	BS	SGM	HSGM	DAVIS	SBD
1	✓	-	-	6.85	5.67
2	✓	✓	-	6.62	5.57
3	✓	✓	✓	6.39	5.36

Table 5: Ablation study of FPM.

Method	Backbone	Params(M)	DAVIS	SBD
RITM*[31]	ResNet-50	31.41	6.74	5.39
RITM*[31]	ResNet-101	58.44	6.78	5.33
Ours	ResNet-50	64.70	5.83	4.98

Table 6: Parameter analysis. The parameters of our method are the sum of coarse and fine level networks.

4.5 Parameters and Inference time analysis

To study the influence of network parameters, we report the amount of parameters of our model and RITM [31], and compare their performance under the comparable parameter setting. As in Tab.6, we observe that the performance of RITM with ResNet-50 and ResNet-101 backbone shows that it is not easy to improve NoC by simply increasing the amount of parameters, but our stage-wise intention-aware framework with a sparse graph network can achieve significant improvement. The results indicate our improvements are not from more parameters but our delicately designed network structure.

Moreover, we test the inference time using the DAVIS dataset with ResNet-50 backbone on the Intel Xeon 2.20GHz CPU and a single NVIDIA Titan RTX GPU. The seconds per click (SPC) is 0.217, which is sufficiently fast for users and is much less than the average time cost of user interaction (3 seconds[4]) in each step.

In practical applications, we can save considerable user annotation time when there are billion images, approximately 4.17×10^5 hours per billion. In addition, as shown in Fig.5, when there are fewer human clicks, we significantly improves compared with baseline, which is very friendly to humans.

5 Conclusion

In this paper, we have developed a novel intention-aware feature propagation strategy for interactive image segmentation. Our method is capable of better inferring user intention and effectively propagating user-provided sparse annotations to the entire input image. To achieve this, we introduce a coarse-to-fine sparse propagation network consisting of a coarse-level and a fine-level sub-networks. Our coarse-level network is able to zoom in onto target regions more accurately and hence mitigates the impact of scale variation. Jointly with the fine-level network, they can capture long-range dependencies at a high spatial resolution with a newly-designed sparse GNN module and hence generate more accurate object boundaries. We evaluated our method on several public benchmarks, in which our approach outperforms the prior works by a sizable margin, achieving state-of-the-art performance.

References

1. Acuna, D., Ling, H., Kar, A., Fidler, S.: Efficient interactive annotation of segmentation datasets with polygon-rnn++. In: Proceedings of the IEEE conference on Computer Vision and Pattern Recognition. pp. 859–868 (2018) [3](#)
2. Agustsson, E., Uijlings, J.R., Ferrari, V.: Interactive full image segmentation by considering all regions jointly. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 11622–11631 (2019) [1](#), [3](#)
3. Bai, J., Wu, X.: Error-tolerant scribbles based interactive image segmentation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 392–399 (2014) [1](#), [3](#)
4. Benenson, R., Popov, S., Ferrari, V.: Large-scale interactive object segmentation with human annotators. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 11700–11709 (2019) [3](#), [6](#), [9](#), [14](#)
5. Castrejon, L., Kundu, K., Urtasun, R., Fidler, S.: Annotating object instances with a polygon-rnn. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 5230–5238 (2017) [3](#)
6. Chen, L.C., Zhu, Y., Papandreou, G., Schroff, F., Adam, H.: Encoder-decoder with atrous separable convolution for semantic image segmentation. In: Proceedings of the European Conference on Computer Vision (ECCV) (September 2018) [7](#), [9](#)
7. Chen, X., Zhao, Z., Yu, F., Zhang, Y., Duan, M.: Conditional diffusion for interactive segmentation. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 7345–7354 (2021) [1](#), [2](#), [3](#), [4](#), [6](#), [10](#), [12](#), [13](#)
8. Cheng, M.M., Zhang, F.L., Mitra, N.J., Huang, X., Hu, S.M.: Repfinder: finding approximately repeated scene elements for image editing. *ACM Transactions on Graphics (TOG)* **29**(4), 1–8 (2010) [1](#)
9. Everingham, M., Van Gool, L., Williams, C.K., Winn, J., Zisserman, A.: The pascal visual object classes (voc) challenge. *International journal of computer vision* **88**(2), 303–338 (2010) [9](#)
10. Gori, M., Monfardini, G., Scarselli, F.: A new model for learning in graph domains. In: Proceedings. 2005 IEEE International Joint Conference on Neural Networks, 2005. vol. 2, pp. 729–734. IEEE (2005) [4](#)
11. Grady, L.: Random walks for image segmentation. *IEEE transactions on pattern analysis and machine intelligence* **28**(11), 1768–1783 (2006) [3](#)
12. Gulshan, V., Rother, C., Criminisi, A., Blake, A., Zisserman, A.: Geodesic star convexity for interactive image segmentation. In: 2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition. pp. 3129–3136. IEEE (2010) [3](#)
13. Hao, Y., Liu, Y., Wu, Z., Han, L., Chen, Y., Chen, G., Chu, L., Tang, S., Yu, Z., Chen, Z., Lai, B.: Edgeflow: Achieving practical interactive segmentation with edge-guided flow. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV) Workshops. pp. 1551–1560 (October 2021) [3](#)
14. Hariharan, B., Arbeláez, P., Bourdev, L., Maji, S., Malik, J.: Semantic contours from inverse detectors. In: 2011 International Conference on Computer Vision. pp. 991–998. IEEE (2011) [9](#)
15. Jang, W.D., Kim, C.S.: Interactive image segmentation via backpropagating refinement scheme. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 5297–5306 (2019) [3](#), [9](#), [10](#)
16. Kontogianni, T., Gygli, M., Uijlings, J., Ferrari, V.: Continuous adaptation for interactive object segmentation by learning from corrections. In: European Conference on Computer Vision. pp. 579–596. Springer (2020) [10](#)

17. Le, H., Mai, L., Price, B., Cohen, S., Jin, H., Liu, F.: Interactive boundary prediction for object selection. In: Proceedings of the European Conference on Computer Vision (ECCV). pp. 18–33 (2018) [3](#)
18. Li, Z., Chen, Q., Koltun, V.: Interactive image segmentation with latent diversity. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 577–585 (2018) [10](#)
19. Liew, J., Wei, Y., Xiong, W., Ong, S.H., Feng, J.: Regional interactive image segmentation networks. In: 2017 IEEE international conference on computer vision (ICCV). pp. 2746–2754. IEEE Computer Society (2017) [3](#), [10](#)
20. Lin, T.Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L.: Microsoft coco: Common objects in context. In: European conference on computer vision. pp. 740–755. Springer (2014) [9](#)
21. Lin, Z., Zhang, Z., Chen, L.Z., Cheng, M.M., Lu, S.P.: Interactive image segmentation with first click attention. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 13339–13348 (2020) [2](#), [3](#), [10](#)
22. Ling, H., Gao, J., Kar, A., Chen, W., Fidler, S.: Fast interactive object annotation with curve-gcn. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 5257–5266 (2019) [3](#)
23. Mahadevan, S., Voigtlaender, P., Leibe, B.: Iteratively trained interactive segmentation. In: British Machine Vision Conference (BMVC) (2018) [10](#)
24. Majumder, S., Yao, A.: Content-aware multi-level guidance for interactive instance segmentation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 11602–11611 (2019) [3](#), [10](#)
25. McGuinness, K., O’connor, N.E.: A comparative evaluation of interactive segmentation algorithms. Pattern Recognition **43**(2), 434–444 (2010) [9](#)
26. Perazzi, F., Pont-Tuset, J., McWilliams, B., Van Gool, L., Gross, M., Sorkine-Hornung, A.: A benchmark dataset and evaluation methodology for video object segmentation. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 724–732 (2016) [9](#)
27. Rother, C., Kolmogorov, V., Blake, A.: ” grabcut” interactive foreground extraction using iterated graph cuts. ACM transactions on graphics (TOG) **23**(3), 309–314 (2004) [1](#), [3](#), [9](#)
28. Scarselli, F., Gori, M., Tsoi, A.C., Hagenbuchner, M., Monfardini, G.: The graph neural network model. IEEE transactions on neural networks **20**(1), 61–80 (2008) [4](#)
29. Sofiiuk, K., Barinova, O., Konushin, A.: Adaptis: Adaptive instance selection network. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 7355–7363 (2019) [8](#)
30. Sofiiuk, K., Petrov, I., Barinova, O., Konushin, A.: f-brs: Rethinking backpropagating refinement for interactive segmentation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 8623–8632 (2020) [1](#), [2](#), [3](#), [8](#), [10](#), [11](#), [12](#)
31. Sofiiuk, K., Petrov, I.A., Konushin, A.: Reviving iterative training with mask guidance for interactive segmentation. arXiv preprint arXiv:2102.06583 (2021) [1](#), [2](#), [3](#), [4](#), [8](#), [9](#), [10](#), [11](#), [12](#), [14](#)
32. Wang, G., Zuluaga, M.A., Li, W., Pratt, R., Patel, P.A., Aertsen, M., Doel, T., David, A.L., Deprest, J., Ourselin, S., et al.: Deepigeos: a deep interactive geodesic framework for medical image segmentation. IEEE transactions on pattern analysis and machine intelligence **41**(7), 1559–1572 (2018) [1](#)
33. Wang, X., Girshick, R., Gupta, A., He, K.: Non-local neural networks. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 7794–7803 (2018) [2](#), [4](#), [6](#)
34. Wu, J., Zhao, Y., Zhu, J.Y., Luo, S., Tu, Z.: Milcut: A sweeping line multiple instance learning paradigm for interactive image segmentation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 256–263 (2014) [3](#)

- 35. Wu, Z., Pan, S., Chen, F., Long, G., Zhang, C., Philip, S.Y.: A comprehensive survey on graph neural networks. *IEEE transactions on neural networks and learning systems* **32**(1), 4–24 (2020) [4](#)
- 36. Xu, N., Price, B., Cohen, S., Yang, J., Huang, T.S.: Deep interactive object selection. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. pp. 373–381 (2016) [1](#), [2](#), [3](#), [5](#), [6](#), [9](#), [10](#)
- 37. Zhang, S., Liew, J.H., Wei, Y., Wei, S., Zhao, Y.: Interactive object segmentation with inside-outside guidance. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 12234–12244 (2020) [3](#)