

---

# Supplementary: Segment Everything Everywhere All at Once

---

Anonymous Author(s)

Affiliation

Address

email

## 1 A Clarification of Main Paper

2 In main.Table.1, we only marked SegGPT as concurrent work. However, #SAM is also a concurrent  
3 work to ours.

## 4 B Architecture details

5 **Human-Model Interaction** For the Human-Model Interaction depicted in main paper Fig.2.b, users  
6 have the flexibility to provide various inputs. These inputs can include positive/negative points,  
7 arbitrary shape strokes of the input image, or a referring image. These inputs are encoded as a visual  
8 prompt using the visual sampler. Additionally, a user can also input a description text for an object,  
9 which will be encoded through a text encoder. The visual and text prompts can be integrated to  
10 generate the final output. If the predicted mask does not meet the user’s expectations, the user can  
11 input corrections in any format, including both visual and text prompts. The previous mask is stored  
12 in a memory prompt. The visual prompt can indicate either a correct or incorrect region, marked as  
13 positive and negative respectively described in detail in the section below.

14 **Visual & Text Prompt Interaction** Taking inspiration from [1, 2], our SEEM decoder also employs  
15 a hierarchical structure. Our positive and negative prompts are pooled using a visual sampler, which  
16 utilizes features at different scales. These features align with image features through a cross-attention  
17 mechanism. This design choice helps bridge the gap in the embedding domain at each decoder  
18 layer, ensuring that the query, image, and prompt features are synchronized in the Joint Image Text  
19 Representation space. Detailed operation of the visual sampler is provided in the pseudo code below.

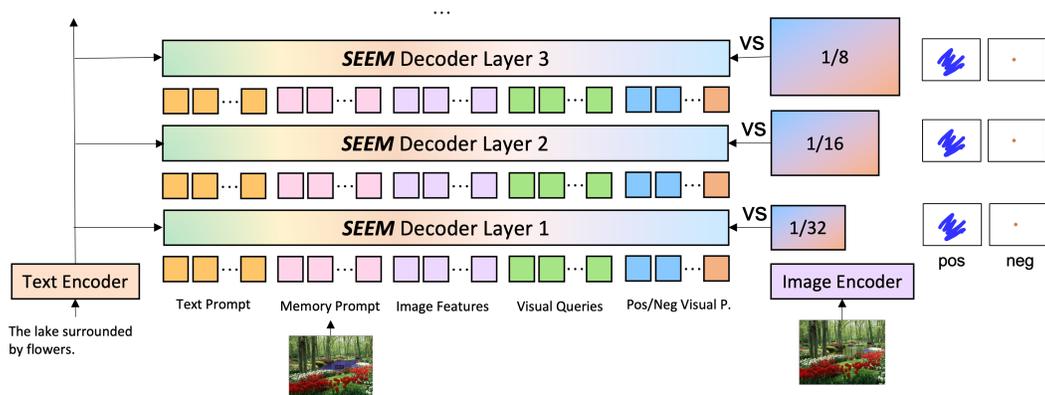


Figure 1: Detailed depiction of the Visual & Text Prompt Interaction in the SEEM Decoder. (VS denotes visual sampler – details are provided in Algorithm 1 in this document.)

---

**Algorithm 1** Pseudo code for visual sampler in SEEM.

---

```
# Inputs: img_f - image feature, [B,C,H,W]; pm/nm - pos/neg mask, [B,1,H,W]; max_len=512
# Functions: import torch.nn.functional as F; import torch
1 def Visual_Sampler(img_f, pm, nm, max_len):
2     pos_emb = img_f[pm]; neg_emb = img_f[nm];# Pool pos/neg features from image feature.
3     pos_emb = F.interpolate(pos_emb, min(max_len, len(pos_emb)), "nearest");# random selection.
4     neg_emb = F.interpolate(neg_emb, min(max_len, len(neg_emb)), "nearest");# random selection.
5     visual_prompt = torch.cat(pos_emb, neg_emb);# Formulate Visual Prompt.
6     visual_queries = torch.mean(pos_emb, dim=0, keepdims=True);# Formulate Queries Spatial.
```

---

## 20 C Implementation details

### 21 C.1 Training Details

22 **Dataset Preparation.** For training, we employ the COCO dataset. Guided by the methodologies  
23 in [3, 4, 5], we replace overlapping masks in COCO and LVIS that have an Intersection over Union  
24 (IoU) value exceeding 0.7 with the corresponding LVIS mask. To emulate user inputs, we make use  
25 of the OpenCV library to create scribbles, strokes, polygons, and points. The shape and thickness  
26 of these elements are randomly adjusted during training, as is the number of masks. This process  
27 constitutes the mask initialization, with sample masks presented in Figure 2. Upon mask initialization,  
28 user input is simulated by utilizing the center points of false negative and false positive mask regions.  
29 Throughout training, the first N-1 iterations of interactive segmentation are performed without  
30 gradient updates. Here, the output mask serves as the input mask for the final layer, simulating human  
31 input, where N is a randomly selected number between 0 and 5.

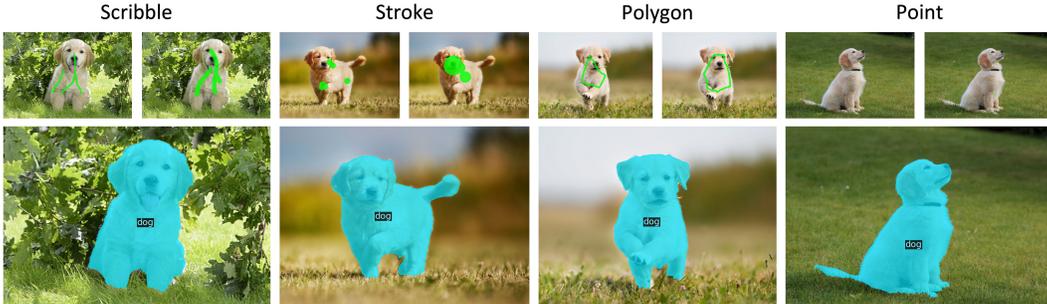


Figure 2: Visualization of mask initialization during training.

32 **Hyper-parameters.** In our study, we primarily follow the training hyperparameters set in [2]. We  
33 conducted training under two settings in the main paper: one that utilizes a vision language pretrained  
34 model starting from scratch, and another that employs an X-Decoder pretrained model. When training  
35 from a vision language model, we train all parameters with the exception of the language encoder.  
36 Conversely, when training from the X-Decoder checkpoint, we only train the SEEM-Decoder. The  
37 training duration differs between versions; for the tiny model, the first version requires approximately  
38 three days of training on 32 V100 GPUs, whereas the second version demands around one day and  
39 five hours on the same GPU configuration.

### 40 C.2 Evaluation Datasets for Interactive Segmentation

41 The evaluation datasets for interactive segmentation are selected with the primary principle of  
42 following existing works whenever possible. When this isn't feasible, we randomly sample 600  
43 images and corresponding masks from the dataset, ensuring a diverse vocabulary distribution. For  
44 the COCO dataset, we align with [6] to download the COCO mini dataset comprising around 600  
45 images. For the PASCAL VOC dataset, we adhere to [3, 7] and employ the original PASCAL VOC  
46 dataset. Similarly, for the DAVIS dataset, we procure the video frames as chosen in [6]. With regards  
47 to all input types assessed in Table 2 of the main paper, we follow the mask generation procedure  
48 illustrated in Figure 2.

### 49 C.3 Evaluation Metrics

50 For interactive segmentation, we use two metrics NOC and 1-IoU:

51 **NOC:** Number of Clicks (NOC) is commonly used in interactive segmentation. It calculate the  
52 average number of clicks required to output a high-quality mask which has an IoU above a predefined  
53 threshold. In some cases, a model cannot output masks that exceed the IoU threshold leading to  
54 infinite or very large NOC. Therefore, we also cap the largest number of clicks. For example 5-  
55 NOC85 denotes the NOC to achieve IoU=85 and the max possible NOC is 5.

56 **1-IoU:** In most cases, our model outputs high-quality masks in one round of interaction. Therefore,  
57 we use 1-IoU to denote the IoU between the output mask and the GT mask with a single prompt.

58 For video object segmentation (VOS) we mainly report JF [8] metric which is the average of J (Region  
59 Similarity) and F (Contour Accuracy).

60 **Region Similarity J:** In order to quantify the similarity of region-based segmentation, specifically the  
61 extent of mislabeled pixels, we utilize the Jaccard index J. This index is defined as the intersection-  
62 over-union between the estimated segmentation and the ground-truth mask. The Jaccard index,  
63 introduced in PASCAL VOC2008 [9], has gained widespread acceptance due to its ability to provide  
64 intuitive and scale-invariant information about the number of mislabeled pixels. For a given output  
65 segmentation M and its corresponding ground-truth mask G, the Jaccard index is calculated as

$$66 J = \frac{|M \cap G|}{|M \cup G|}.$$

67 **Contour Accuracy F:** From a contour-based perspective, the set of closed contours  $c(M)$  can be  
68 seen as defining the spatial boundaries of the mask. To assess contour accuracy, one can compute  
69 precision and recall ( $P_c$  and  $R_c$ ) between the contour points of  $c(M)$  and  $c(G)$  using a robust bipartite  
70 graph matching approach. The F-measure (F) is considered a balanced metric between the two, and  
71 is defined as  $F = \frac{2P_c R_c}{P_c + R_c}$ . In our experiments, we adopt an efficient approximation of the bipartite  
72 matching using morphology operators.

## 73 D Experiments

### 74 D.1 Quantitative Results

75 **Ablation Study on Modality Composition** We rectify Main.Table.5 by utilizing our latest model that  
76 covers all backbone scales, from tiny to base and large. Please note that in the main paper’s Table.5,  
77 we inadvertently used an older configuration for the Tiny model that deviates from the configurations  
78 used in all other evaluations in the main paper. The results depicted in the table below elucidate  
79 several key points: (a) Merely ensembling the predicted visual and text masks does not generally  
80 enhance the performance of the referring segmentation. (b) Utilizing self-attention to amalgamate  
81 visual and text information (refer to Tentative Attention in Main.Figure.3) significantly improves  
82 performance across all backbone scale sizes in comparison to using solely text information, regardless  
83 of the output query we employ. (c) Generally speaking, supplementing text information with spatial  
84 data confers a larger improvement than the converse scenario, suggesting that spatial information  
85 tends to yield more accurate results than text information.

Table 1: The term ‘Text/Visual Prompt’ refers to the modality of information utilized in the study. ‘Output Query’ is indicative of the type of query employed to predict the output. ‘Composition Approach’ specifies the method through which text and visual information are integrated.

Text Prompt	Visual Prompt	Output Query	Composition Approach	Focal-Tiny			DavIt-Base			DavIt-Large		
				cIoU	mIoU	AP@50	cIoU	mIoU	AP@50	cIoU	mIoU	AP@50
Y	N	Text	N/A	58.4	63.4	71.6	63.0	68.2	76.7	62.4	67.6	75.3
Y	Y	All	Ensemble	63.0	60.0	66.9	69.3	66.6	74.3	68.9	65.5	72.7
Y	Y	Text	Self-Attn	66.5	69.6	78.8	75.0	76.9	86.3	73.2	76.5	85.9
N	Y	Visual	N/A	70.7	71.8	81.3	75.4	77.8	87.4	<b>75.2</b>	78.2	87.7
Y	Y	Visual	Self-Attn	<b>71.5</b>	<b>72.8</b>	<b>82.2</b>	<b>75.9</b>	<b>78.3</b>	<b>87.7</b>	74.9	<b>78.4</b>	<b>87.7</b>

86 **Ablation Study on Backbone Architecture.** To facilitate a fair comparison of the model backbone  
87 with SAM [10], SegGPT [11], and SimpleClick [4], we employ the SAM pretrained ViT backbone and  
88 present quantitative results in Table 2. Several conclusions can be drawn from these observations: (1)  
89 The utilization of the SAM pretrained ViT backbone tends to enhance the performance of interactive  
90 segmentation, while it may decrease the performance in the domains of generic and referring  
91 segmentation. This can be potentially attributed to the fact that the vision backbone is not trained

92 with any semantic labels. (2) As noted in the table, our ViT-L is trained with fix vision backbone  
 93 because of memory limitation. Compare with ViT-B generic segmentation result is extremely weak.  
 94 However surprisingly, the referring segmentation result is still comparable or even better.

Table 2: Ablation study on backbone architecture. ‘\*COCO+LVIS’ denotes that the model does not strictly adhere to this configuration, given that the SAM pretrained checkpoint is trained with the SAM dataset. We were unable to train the model until completion, thus we have only compared results for the same epoch. Also, ‘\*ViT-L’ denotes the backbone is fixed due to limited GPU memory. Full results will be posted in a future version.

Backbone	Segmentation Data	Epoch	COCO			Ref-COCOg				VOC	
			PQ	mAP	mIoU	cIoU	mIoU	AP@50	NoC@50	NoC@85	NoC@90
davit-d3	COCO+LVIS	27/50	50.6	41.6	61.4	49.2	57.9	64.7	1.54	3.45	4.43
ViT-B	SAM+COCO	27/50	48.5	40.3	56.8	50.1	59.0	66.1	1.49	2.99	3.91
davit-d5	COCO+LVIS	27/50	54.4	44.4	64.7	51.0	58.8	65.4	1.46	3.21	4.03
*ViT-L	SAM+COCO	27/50	44.4	38.0	52.3	51.5	59.8	66.6	1.42	2.91	3.71

## 95 D.2 Qualitative results

96 **Open Vocabulary Interactive Segmentation.** In Main.Figure.4, we evaluated interactive segmentation  
 97 within a closed vocabulary, identical to the COCO classes. In Figure 3, we examine this process  
 98 within an open vocabulary setting, manually providing candidate classes within the given image.  
 99 Remarkably, our model demonstrates adept performance on unseen classes. This capability is a  
 100 legacy of the X-Decoder and is powered by the Joint Image-Text Representation Space in SEEM.  
 101 Notably, certain items such as yellow corn, leaf, french fries, orange juice, and others, which have  
 102 never been trained, are successfully recognized by our model.



Figure 3: Visualization on Open Vocabulary Interactive Segmentation. (Best view with zoom in)

103 **Multi-Round Interactive Segmentation.** In the main paper, as seen in Figure.4, we evaluated the  
 104 single-shot result with an arbitrary input stroke provided by a human. However, in Figure 4, we assess  
 105 interactive segmentation in a multi-round format, illustrating the correction procedure of SEEM.  
 106 Images are grouped in sets of four following this sequence: Input Round1, Output1, Input Round2,  
 107 and Output2. Positive strokes are colored in green, while negative strokes are colored in blue. The  
 108 results clearly demonstrate that our negative tokens can effectively remove undesirable masks from  
 109 the previous round, and positive tokens are capable of inpainting missing parts from the previous  
 110 round.



Figure 4: Multi-Round Interactive Segmentation.

111 **Modality Composition.** As demonstrated in the main paper, specifically in Figures 4 and 5, our  
 112 model is capable of performing both interactive and referring segmentation. We have also shown  
 113 in the main paper’s Table 5, and in Table 1 of the supplementary material, that the composition of  
 114 spatial and textual information can enhance performance quantitatively. As illustrated in Figure 5, we  
 115 present visualization results that integrate spatial and textual content. Our spatial content efficiently  
 116 resolves ambiguities inherent in the text. For example, it helps determine which dog is playing with  
 117 the yellow ball, or identify what the left blue magical creature is.



Figure 5: Modality Composition. (best viewed with zoom in.)

118 **Comparison of Interactive Segmentation with Concurrent Work.** While we have showcased  
 119 numerous capabilities of the SEEM model independently, here we examine how it fares against  
 120 concurrent works on interactive segmentation. We contrast our work with both SAM and SegGPT.  
 121 As evident in the column of SAM output, the model seems to struggle with either identifying the  
 122 specific object (e.g., which cat the user is pointing to) or capturing all components of an object (e.g.,  
 123 missing legs of a chair). Additionally, as SegGPT is trained with referring images, it appears that the  
 124 model may falter when partial masks are provided. In summary, our model holds the advantage in  
 125 identifying objects, and it also exhibits a proficiency in classifying instances, an aspect that seems to  
 126 be absent in the other concurrent methods.



Figure 6: Comparison of Interactive Segmentation with concurrent work.

127

## 128 References

- 129 [1] Bowen Cheng, Ishan Misra, Alexander G Schwing, Alexander Kirillov, and Rohit Girdhar.  
 130 Masked-attention mask transformer for universal image segmentation. In *Proceedings of the*  
 131 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1290–1299, 2022.
- 132 [2] Xueyan Zou, Zi-Yi Dou, Jianwei Yang, Zhe Gan, Linjie Li, Chunyuan Li, Xiyang Dai, Harkirat  
 133 Behl, Jianfeng Wang, Lu Yuan, et al. Generalized decoding for pixel, image, and language.  
 134 *arXiv preprint arXiv:2212.11270*, 2022.
- 135 [3] Konstantin Sofiiuk, Ilya A Petrov, and Anton Konushin. Reviving iterative training with mask  
 136 guidance for interactive segmentation. In *2022 IEEE International Conference on Image*  
 137 *Processing (ICIP)*, pages 3141–3145. IEEE, 2022.

- 138 [4] Hao Zhang, Feng Li, Xueyan Zou, Shilong Liu, Chunyuan Li, Jianfeng Gao, Jianwei Yang,  
139 and Lei Zhang. A simple framework for open-vocabulary segmentation and detection. *arXiv*  
140 *preprint arXiv:2303.08131*, 2023.
- 141 [5] Xi Chen, Zhiyan Zhao, Yilei Zhang, Manni Duan, Donglian Qi, and Hengshuang Zhao. Fo-  
142 calclick: towards practical interactive image segmentation. In *Proceedings of the IEEE/CVF*  
143 *Conference on Computer Vision and Pattern Recognition*, pages 1300–1309, 2022.
- 144 [6] Konstantin Sofiiuk, Iliia A. Petrov, and Anton Konushin. Reviving iterative training with mask  
145 guidance for interactive segmentation, 2021.
- 146 [7] Qin Liu, Zhenlin Xu, Gedas Bertasius, and Marc Niethammer. Simpleclick: Interactive image  
147 segmentation with simple vision transformers. *arXiv preprint arXiv:2210.11006*, 2022.
- 148 [8] Federico Perazzi, Jordi Pont-Tuset, Brian McWilliams, Luc Van Gool, Markus Gross, and  
149 Alexander Sorkine-Hornung. A benchmark dataset and evaluation methodology for video  
150 object segmentation. In *Proceedings of the IEEE conference on computer vision and pattern*  
151 *recognition*, pages 724–732, 2016.
- 152 [9] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman.  
153 The pascal visual object classes (voc) challenge. *International journal of computer vision*,  
154 88:303–338, 2010.
- 155 [10] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson,  
156 Tete Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and Ross Girshick.  
157 Segment anything, 2023.
- 158 [11] Xinlong Wang, Xiaosong Zhang, Yue Cao, Wen Wang, Chunhua Shen, and Tiejun Huang.  
159 Seggpt: Segmenting everything in context. *arXiv preprint arXiv:2304.03284*, 2023.