

Segment Everything Everywhere All at Once

Xueyan Zou^{*§2}, Jianwei Yang^{*†1}, Hao Zhang^{*‡}, Feng Li^{*‡}, Linjie Li[†], Jianfeng Gao^{¶†}, Yong Jae Lee^{¶†§}

[§] University of Wisconsin-Madison [†] Microsoft Research at Redmond [‡] HKUST [¶] Microsoft Cloud & AI

^{*}Equal Contribution [¶] Equal Advisory Contribution 1. Project Lead 2. Main Technical Contribution

{xueyan,yongjaelee}@cs.wisc.edu {jianwyan,jfgao,linjli}@microsoft.com {hzhangcx,fliay}@connect.ust.hk

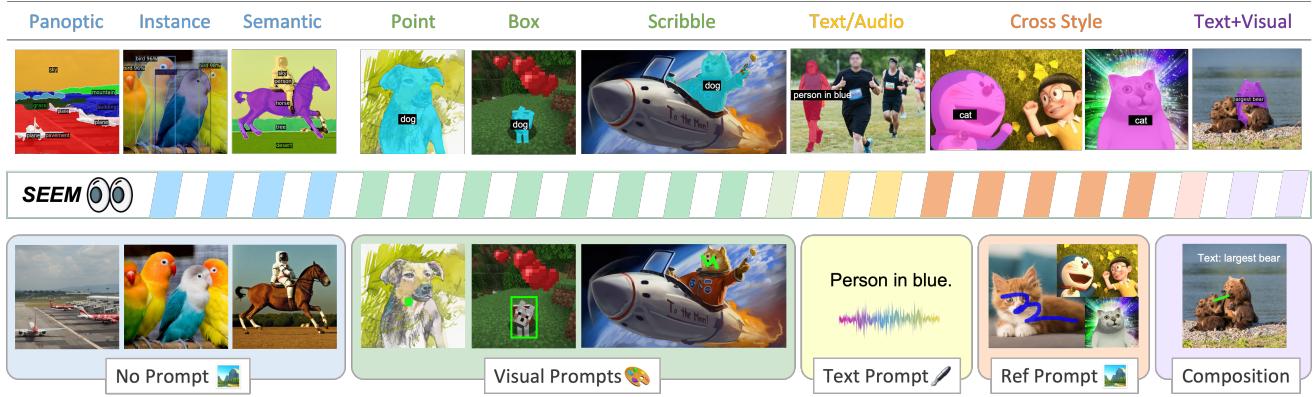


Figure 1: Our model, *SEEM*, can perform any segmentation task, such as semantic, instance, and panoptic segmentation, in open-set scenarios when no prompt is provided. Moreover, *SEEM* supports visual, textual, and referring region prompts in any arbitrary combination, allowing for versatile and interactive referring segmentation.

Abstract

Despite the growing demand for interactive AI systems, there have been few comprehensive studies on human-AI interaction in visual understanding e.g. segmentation. Inspired by the development of prompt-based universal interfaces for LLMs, this paper presents *SEEM*, a promptable, interactive model for Segmenting Everything Everywhere all at once in an image. *SEEM* has four desiderata: i) *Versatility* by introducing a versatile prompting engine for different types of prompts, including points, boxes, scribbles, masks, texts, and referred regions of another image; ii) *Compositionality* by learning a joint visual-semantic space for visual and textual prompts to compose queries on the fly for inference as shown in Fig. 1; iii) *Interactivity* by incorporating learnable memory prompts to retain dialog history information via mask-guided cross-attention; and iv) *Semantic-awareness* by using a text encoder to encode text queries and mask labels for open-vocabulary segmentation. A comprehensive empirical study is performed to validate the effectiveness of *SEEM* on various segmentation tasks. *SEEM* shows a strong capability of generalizing to unseen user intents as it learned to compose prompts of different types in a unified representation space. In addition, *SEEM* can

efficiently handle multiple rounds of interactions with a lightweight prompt decoder. The *SEEM* code and demo are available at <https://github.com/UX-Decoder/Segment-Everything-Everywhere-All-At-Once>.

1. Introduction

The success of Large Language Models (LLMs) such as ChatGPT [38] have shown the importance of modern AI models in interacting with humans, and have provided a glimpse of AGI [3]. The ability to interact with humans requires a user-friendly interface that can take as many types of human inputs as possible and generate responses that humans can easily understand. In NLP, such a universal interaction interface has emerged and evolved for a while from early models like GPT [2] and T5 [41], to some more advanced techniques like prompting [45, 63, 26] and chain-of-thought [48, 23, 44]. In the image generation area, several recent works attempt to combine text prompts with other types like sketches or layouts, to capture user intents more precisely [42, 39, 58], compose novel prompts [62, 27], and support multi-round human-AI interactions [1, 49].

Although interactive image segmentation has a long history [28, 14, 53, 7, 32], segmentation models that can interact with humans via *universal interfaces* that can take

multiple types of prompts (*e.g.*, texts, clicks, images) as input has not been well-explored. Existing ways of interactive segmentation are either using spatial hints like clicks or scribbles [7, 32] or referring segmentation using language [16, 36, 17, 65]. Most recently, a few works combine the textual prompts with box prompts [61]. However, they can only take one or two prompt types, which are far to reach the demands in real-world applications. A concurrent work SAM [22] supports multiple prompts. However, ours is very different from SAM as shown in Fig. 2. For example, SAM only supports limited interaction types like points and boxes and does not support high-level semantic tasks since it does not output semantic labels. Note that although we do not report edge detection results, our model can support it by simply converting masks to edges.

In this work, we advocate *a universal interface for segmenting everything everywhere with multi-modal prompts*. To achieve this goal, we propose a new prompting scheme that has four important properties, *versatility*, *compositionality*, *interactivity*, and *semantic-awareness*. For versatility, we propose to encode points, masks, text, boxes, and even a referred region of another image that are seemingly heterogeneous into prompts in the same joint visual-semantic space. As such, our model can deal with any combination of the input prompts, leading to strong compositionality. To enable interactivity, we further introduce memory prompts for condensing the previous segmentation information followed by communication with other prompts. As for semantic awareness, our model gives an open-set semantic to any output segmentation.

With the proposed prompting scheme, we build a segment-everything-everywhere model called ***SEEM*** following a simple Transformer encoder-decoder architecture with an extra text encoder [65, 61]. In ***SEEM***, the decoding process behaves similarly to generative LLMs but with multimodality-in-multimodality-out. All queries are taken as prompts and fed into the decoder, and the image and text encoder are used as the prompt encoder to encode all types of queries. Concretely, we encode all spatial queries such as points, boxes, and scribbles into *visual prompts* by pooling their corresponding visual features from the image encoder, while using the text encoder to convert text queries to *textual prompts*. This way the visual and textual prompts are always aligned with each other. When adopting an exemplar image segment as the query, we encode the image with the same image encoder and pool the image features accordingly. To the end, prompts of all 5 different types are mapped to the *joint visual-semantic space*, to enable unseen user prompts via zero-shot adaptation. By training on different segmentation tasks, our model has the ability to deal with various prompts. Moreover, different types of prompts can help each other via cross-attention among prompts. Therefore, we can composite prompts to obtain

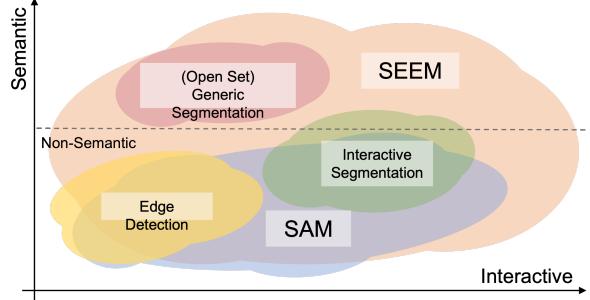


Figure 2: Comparison with concurrent work SAM [22]. Our approach has richer context on both interaction approaches (*e.g.* Referring region of an example image) as well as semantic spaces (support class aware segmentation).

better segmentation results. To the end, we build a segmentation interface with a single pre-trained model that can segment every object with semantics (everything), cover every pixel in the image (everywhere), and support all compositions of prompts (all at once).

In addition to the strong generalization capability, ***SEEM*** is also efficient to run. We take the prompts as input to the decoder. Therefore, when doing multi-round interactions with humans, our model only needs to run the feature extractor once at the beginning. In each iteration, we only need to run the lightweight decoder again with new prompts. When deploying the model, we can run the feature extractor which is usually heavy on a server and the relatively lightweight decoder on the user’s machine to reduce the network delay in multiple remote calls.

In summary, our contributions are threefolds:

- We design a unified prompting scheme that can encode various user intents into prompts in a *joint visual-semantic space* which possesses properties of versatility, compositionality, interactivity, and semantic awareness, leading to zero-shot capabilities in generalizing to unseen prompts for segmentation.
- We build ***SEEM***, a universal and interactive segmentation interface, by integrating the newly designed prompting mechanism into a lightweight decoder for *all* segmentation tasks, taking various types of prompts and combinations as inputs.
- We conduct experiments and visualizations to show that our model has strong performance on many segmentation tasks including closed-set and open-set panoptic segmentation, interactive segmentation, referring segmentation, and segmentation tasks with combined prompts.

2. Related Work

Close-set segmentation. Segmentation of visual concepts has been a persistent challenge in the field of computer vision, as evidenced by the extensive literature on

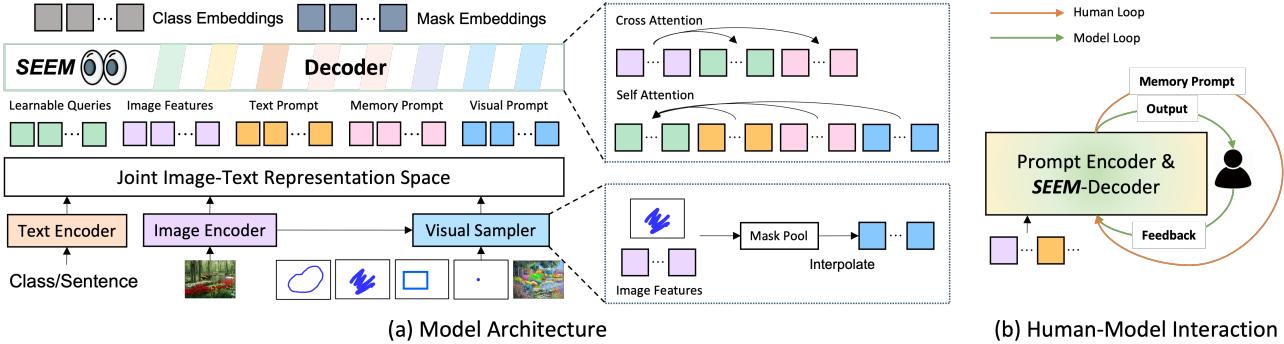


Figure 3: (a) The left part is an overview of the model. First, features and prompts are encoded by its corresponding encoder or sampler to a *joint visual-semantic space*. And learnable queries are randomly initialized. *SEEM* Decoder takes in queries, features, and prompts as input and outputs both class and mask predictions. The right part is the details of *SEEM* Decoder and visual sampler. (b) shows multiple rounds of interactions. Each round contains a human loop and a model loop. In the human loop, human receive the mask output of last iteration and give positive or negative feedback for the next round decoding through visual prompt. In the model loop, the model receives and updates memory prompt for future predictions.

the subject [12, 11, 66, 37]. Generic segmentation techniques encompass several subtasks, including instance segmentation, semantic segmentation, and panoptic segmentation [15, 5, 21], each focused on a different semantic level. For example, semantic segmentation aims to identify and label each pixel within an image based on its corresponding semantic class [6, 8, 34]. On the other hand, instance segmentation involves grouping pixels that belong to the same semantic class into separate object instances [25, 15]. Recently, the Detection Transformer (DETR)[4], a model based on the Transformer[47] architecture, has made significant advances in segmentation [29, 8, 25, 60, 19] tasks. However, these approaches cannot recognize objects absence in the training set, which constrains the model to a limited vocabulary size.

Open-set segmentation. Referring segmentation models [16, 36, 17] targets segmentation from language descriptions, which is open-vocabulary by nature. However, due to the limited referring segmentation data [57, 35], the trained model often performs well on the target dataset but is difficult to extrapolate to real-world applications. Many open-vocabulary segmentation models have been proposed recently [24, 13, 18, 10, 43, 51], which employ large pre-trained vision-language models like CLIP [20] to transfer visual-semantic knowledge by freeze or fine-tune their weights. More recently, X-Decoder [65] proposed a unified approach to tackle various segmentation and vision-language tasks for open-vocabulary segmentation. To expand vocabulary size, OpenSeeD [61] proposes to use a large amount of detection data and a joint training method to improve segmentation. Additionally, ODISE [52] leverages a text-to-image diffusion model as the backbone for open-vocabulary segmentation. In contrast to prior works, our model explores the integration of interactive and open-

set segmentation for general usage.

Interactive segmentation. Interactive segmentation is segmenting objects by interactively taking user inputs. It has been a longstanding problem and achieved considerable progress [28, 14, 53, 32, 7, 22]. Generally, the interaction types can take various forms, such as clicks, boxes, polygons, and scribbles, among which click-based interaction models are the most prevalent. Concurrent to our work, SAM [22] proposed a promptable segmentation model trained on 11 million images, demonstrating strong zero-shot performance. It takes user interactions as prompts for general segmentation. Though SAM advances the vision progress, it produces segmentation without semantic meaning. In addition, the types of prompts are limited to points, boxes, and text.

3. Method

SEEM employs a generic encoder-decoder architecture, but specifically features a sophisticated interaction between queries and prompts, as shown in Fig. 3 (a). Given an input image $\mathbf{I} \in \mathcal{R}^{H \times W \times 3}$, an image encoder is first used to extract image features \mathbf{Z} . Then, seem-decoder predicts the masks \mathbf{M} and semantic concepts \mathbf{C} based on the query outputs \mathbf{O}_h^m (mask embeddings) and \mathbf{O}_h^c (class embeddings) that interacted with visual, text and memory prompts $\langle \mathbf{P}_t, \mathbf{P}_v, \mathbf{P}_m \rangle$:

$$\langle \mathbf{O}_h^m, \mathbf{O}_h^c \rangle = \text{Decoder}(\mathbf{Q}_h; \langle \mathbf{P}_t, \mathbf{P}_v, \mathbf{P}_m \rangle | \mathbf{Z}) \quad (1)$$

$$\mathbf{M} = \text{MaskPredictor}(\mathbf{O}_h^m) \quad (2)$$

$$\mathbf{C} = \text{ConceptClassifier}(\mathbf{O}_h^c) \quad (3)$$

where \mathbf{Q}_h is the learnable queries, and \mathbf{P}_t , \mathbf{P}_v , \mathbf{P}_m represent the text prompts, visual prompts, and memory

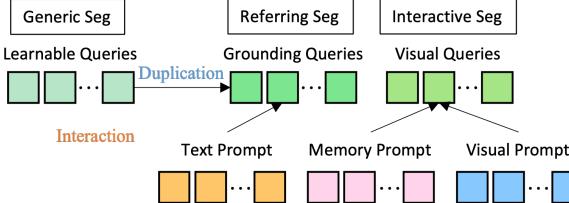


Figure 4: Queries and Prompt Interaction during **training**. Learnable queries are duplicated as grounding and visual queries with the same weights for referring and interactive segmentation.

prompts, respectively. During training, \mathbf{Q}_h is duplicated for generic segmentation, referring segmentation, and interactive segmentation as shown in Fig. 4. And the corresponding prompts are interacted with its queries through self-attention module. At inference time, the learnable queries are freely interacted with all prompts, thereby enabling zero-shot composition. Our design is inspired by the successful practice in X-Decoder [65]. However, we highlight the differences in Eq. (1), marked in red, which allows for a universal model for image segmentation with the following properties:

Versatile. In SEEM, we introduce visual prompts \mathbf{P}_v to handle *all* non-textual inputs, such as points, boxes, scribbles, and a referred region of another image. These non-textual queries are beneficial to disambiguate the user’s intents when textual prompts fail to identify the correct segment. For interactive segmentation, previous works either convert spatial queries to masks and feed them into the image backbone [32] or use the different prompt encoders for each input type (points, boxes) [22]. The first approach is too heavy in application because each interaction requires the image to go through the feature extractor. The second approach is hard to generalize to unseen prompts. To address these limitations, we propose a visual sampler (Fig. 3 (a)) to convert all kinds of non-textual queries to visual prompts that are lying in the same visual embedding space:

$$\mathbf{P}_v = \text{VisualSampler}(\mathbf{s}, \hat{\mathbf{Z}}) \quad (4)$$

where $\hat{\mathbf{Z}}$ is the feature maps extracted from either the target image (*i.e.*, $\hat{\mathbf{Z}} = \mathbf{Z}$) or a referred image, and $\mathbf{s} \in \{\text{points, box, scribbles, polygons}\}$ are the sampling locations specified by users. We first pooled the corresponding region from image feature through point sampling [8]. For all visual prompts, we interpolate at most 512 point feature vectors uniformly from the region specified by the prompt. When altering the input feature maps to the referred image, the sampling approach is the same where the sample image is the exemplar image. Another merit of our proposed method is that the visual prompts are naturally well-aligned with the textual prompts, as our model continuously learns

a common visual-semantic space through panoptic and referring segmentation.

Compositional. In practice, users may cast their intents using different or combined input types. Hence, a compositional approach to prompting is essential for real-world applications. However, we confront two issues during model training. First, the training data usually only cover a single type of interaction (*e.g.* none, textual, visual). Second, although we have used visual prompts to unify all non-textual prompts and align them with textual prompts, their embedding spaces remain inherently different. To address this, we propose matching different types of prompts with different outputs. Considering that visual prompts come from image features while textual prompts come from the text encoder, we select matched output indices for visual and textual prompts by matching them with the mask embeddings \mathbf{O}_h^m or class embeddings \mathbf{O}_h^c , respectively:

$$ID_v \leftarrow \text{Match}(\mathbf{O}_h^m \cdot \mathbf{P}_v + \text{IoU}_{mask}) \quad (5)$$

$$ID_t \leftarrow \text{Match}(\mathbf{O}_h^c \cdot \mathbf{P}_t + \text{IoU}_{mask}) \quad (6)$$

where IoU_{mask} is the IoUs between ground-truth masks and predicted masks. The proposed separate matching method turns out to outperform the approaches that only match with either \mathbf{O}_h^m or \mathbf{O}_h^c for all prompts.

After training, our model becomes familiar with all prompt types and supports a variety of composition ways, such as no prompts, one prompt type, or both visual and textual prompts using the same model and weights. *Particularly, the visual and textual prompts can be simply concatenated and fed to SEEM Decoder, even though it was never trained like that.*

Interactive. Interactive segmentation usually cannot be completed in one shot and requires multiple interaction rounds for refinement, similar to conversational agents like ChatGPT. In SEEM, we propose a new type of prompt called *memory prompts* \mathbf{P}_m and use them to convey the knowledge of the masks from the previous iteration to the current one. Unlike previous works that use a network to encode mask [32, 22], we introduce no extra module but simply a few memory prompts. These memory prompts are responsible for encoding the history information by using a mask-guided cross-attention layer [8]:

$$\mathbf{P}_m^l = \text{MaskedCrossAtt}(\mathbf{P}_m^{l-1}; \mathbf{M}_p | \mathbf{Z}) \quad (7)$$

where \mathbf{M}_p is the previous mask, \mathbf{Z} is the image feature map. So that, cross-attention only takes effect inside the regions specified by the previous mask. The updated memory prompts \mathbf{P}_m^l then interact with the other prompts via self-attention to convey the historical information for the current round. Notably, this design can be easily extended to support the simultaneous segmentation of multiple objects, which we leave for future study.

Table 1: **One suite of weights** for segmentation on a wide range of segmentation tasks. *SEEM* is the first that simultaneously supports generic segmentation, referring segmentation and interactive segmentation, as well as prompts compositions.

Method	Segmentation Data	Type	Generic Segmentation			Referring Segmentation			Interactive Segmentation			
			COCO	RefCOCOg	Pascal VOC	SBD	mIoU	AP50	NoC85	NoC90	NoC85	NoC90
Mask2Former (T) [8]	COCO (0.2M)	Segmentation	53.2	43.3	63.2	-	-	-	-	-	-	-
Mask2Former (B) [8]	COCO (0.2M)		56.4	46.3	67.1	-	-	-	-	-	-	-
Pano/SegFormer (B) [50]	COCO (0.2M)		55.4	*	*	-	-	-	-	-	-	-
LAVT (B) [56]	Ref-COCO (0.2M)		-	-	-	61.2	*	*	-	-	-	-
RITM (<T) [46]	COCO (0.2M)	Interactive	-	-	-	-	-	-	2.19	2.57	3.59	5.71
PseudoClick (<T) [33]	COCO (0.2M)		-	-	-	-	-	-	1.94	2.25	3.79	5.11
FocalClick (T) [7]	COCO (0.2M)		-	-	-	-	-	-	2.97	3.52	4.56	6.86
FocalClick (B) [7]	COCO (0.2M)		-	-	-	-	-	-	2.46	2.88	3.53	5.59
SimpleClick (B) [32]	COCO+LVIS (0.2M)		-	-	-	-	-	-	2.06	2.38	3.43	5.62
X-Decoder (T) [65]	COCO (0.2M)	Generalist	52.6	41.3	62.4	59.8	*	*	-	-	-	-
X-Decoder (B) [65]	COCO (0.2M)		56.2	45.8	66.0	64.5	*	*	-	-	-	-
SAM (B) [22]	SAM (11M)		-	-	-	-	-	-	3.30	4.20	6.50	9.76
<i>SEEM</i> (T)	COCO+LVIS (0.2M)		50.7	39.6	60.7	58.2	63.4	71.3	3.75	4.75	6.92	10.40
<i>SEEM</i> (B)	COCO+LVIS (0.2M)		56.1	46.6	65.2	63.1	68.2	76.8	2.99	3.89	5.93	9.23
<i>SEEM</i> (T)	COCO+LVIS (0.2M)	Composition	*	*	*	62.3	66.3	75.1	*	*	*	*
<i>SEEM</i> (B)	COCO+LVIS (0.2M)		*	*	*	65.7	69.8	77.0	*	*	*	*

Semantic-aware. Different from previous class-agnostic interactive segmentation works such as Simple Click [32] and the concurrent work SAM [22], our model gives semantic labels to masks from all kinds of prompt combinations in a zero-shot manner. Since our visual prompt features are aligned with textual features in a *joint visual-semantic space*. As shown in Fig. 4, semantic labels would be directly computed by \mathbf{O}_h^c (output of visual queries) and the vocabularies text embedding. Although **we did not train any semantic labels for interactive segmentation**, the calculated logits are well aligned, benefiting from the *joint visual-semantic space*.

4. Experiments

Datasets and Settings. *SEEM* is trained on three data types: panoptic segmentation, referring segmentation, and interactive segmentation. Panoptic and interactive segmentation are trained with COCO2017 [30] with panoptic segmentation annotations. Following [65], we exclude the validation set of Ref-COCOg [57] umd, resulting in 107K segmentation images in total. For referring segmentation we use a combination of Ref-COCO, Ref-COCOg and Ref-COCO+ for COCO image annotations. We evaluate all segmentation tasks covered by pertaining, including generic segmentation (instance/panoptic/semantic), referring segmentation, and interactive segmentation.

Implementation Details and Evaluation Metrics. Our model framework follows X-Decoder [65] except the decoder part that is composed of a vision backbone, a language backbone, an encoder, and seem-decoder. For the vision backbone, we use FocalIT [54] and DaViT-d3 (B) [9]. For the language encoder, we adopt a UniCL or Florence text encoder [55, 59]. Regarding evaluation metrics, for all segmentation tasks, we use standard metrics, such as PQ (Panoptic Quality) for panoptic segmentation, AP (Average

Table 2: Referring segmentation with text and visual prompts. *P* shorts for prompts and *Q* shorts for queries.

Text-P	Visual-P	Output-Q	Composition	Grounding cIoU	mIoU
✓	✓	All	Ensemble	68.1	71.8
✓	✗	Textual	-	57.3	62.3
✓	✓	Textual	Self-Attention	60.3	65
✗	✓	Visual	-	72.5	73.8
✓	✓	Visual	Self-Attention	72.9	74.5

Precision) for instance segmentation, and mIoU (mean Intersection over Union) for semantic segmentation. For interactive segmentation, we follow previous works [32, 31] to simulate user clicks by comparing the predicted segmentation with the ground-truth one in an automatic way. After one click on the image to generate the predicted mask, the next click will be placed at the center of the area with the largest segmentation error. We use the Number of Clicks (NoC) metric to evaluate interactive segmentation performance, which measures the number of clicks needed to achieve a certain Intersection over Union (IoU), i.e., 85%, and 90%, denoted as NoC@85, and NoC@90, respectively.

4.1. Interactive segmentation

In Table 1, we present a comparison of our model with the state-of-the-art interactive segmentation models. As a generalist model, our approach achieves comparable performance with RITM, SimpleClick, etc, and very similar performance compared with SAM [22] that is trained with $\times 50$ more segmentation data than ours. Notably, unlike existing interactive models, *SEEM* is the first interface that supports not only classical segmentation tasks but also a wide range of user input types, including text, points, scribbles, boxes, and images, providing strong compositional capabilities as shown in Tab. 2.

Table 3: **Ablation study** on interaction strategy. “Iter” denotes iterating segmentation for multiple rounds. “Negative” means adding negative points during interactive segmentation.

	Fix	Iter	Pos	Neg	PQ	COCO mAP	mIoU	Referring Segmentation cIoU	mIoU	AP@50	Pascal VOC NoC85	NoC90	SBD NoC85	NoC90
Baseline	✓	✗	✓	✗	51.1	39.8	62.1	57.4	63.1	71.0	4.52	5.64	7.61	11.04
- Scratch	✗	✗	✓	✗	49.0	38.7	59.9	55.6	61.3	69.6	4.19	5.35	7.64	12.21
+ Iter	✓	✓	✓	✗	50.8	39.5	60.9	58.0	63.1	71.3	4.20	5.44	7.86	11.30
+ Negative	✓	✓	✓	✓	50.6	39.5	61.2	56.6	62.7	70.9	4.12	5.23	6.81	10.10



Figure 5: Click for segmentation. *SEEM* supports arbitrary formats of clicks or scribbles by users. Moreover, it simultaneously gives the semantic label for the segmented mask, which is not possible in SAM [22].

4.2. Generic segmentation

With one suite of parameters pre-trained on all the segmentation tasks, we directly evaluate its performance on generic segmentation datasets. As shown in Table 1, *SEEM* maintains comparable panoptic, instance and semantic segmentation performance with strong baselines.

4.3. Referring segmentation

As shown in Table 1, compared with other referring segmentation and generalist models, *SEEM* achieves competitive performance. Notably as shown in Table 2, by adding a visual compositional prompt, referring segmentation performance are improved with a large margin with 5.7, 3.6, and 4.2 points, respectively under cIoU, mIoU, and AP50 metrics for the tiny model. And this gap retrains on the base model with improvements of 2.5, 1.5, and 0.4 points, respectively. Specifically, this number is computed by class embeddings \mathbf{O}_h^c (Output-Q-Textual). While the margin is even larger when computed with mask embeddings \mathbf{O}_h^m (Output-Q-Visual) as shown in Tab. 2. Further, we benchmark the vanilla composition (Ensemble) that directly combine visual and text mask on output probability.

4.4. Ablation Study

We conduct an ablation study of our model components in Tab. 3 on all the training segmentation tasks. The performance of generic segmentation decreases slightly when adding iteration and negative visual prompts. Furthermore, if we train the model from scratch, the performance of

generic segmentation declines even more. As expected, the referring segmentation performance drops when training from scratch. However, it further decreases when adding negative visual prompts. On the other hand, adding iterations for the interaction segmentation task slightly improves the grounding performance. The interactive segmentation performance gradually improves by adding iteration and negative visual prompts, while training from scratch surprisingly results in a slight improvement in performance on the Pascal VOC dataset.

4.5. Qualitative Results

We further qualitatively evaluate *SEEM*. Based on the proposed prompting scheme and decoder design, with the same suite of parameters, *SEEM* supports a wide range of visual input types.

Visual prompt interactive segmentation. In Fig. 5, we show the visualization of using *SEEM* to segment objects in an interactive way. The user can segment objects of interest by simply clicking or drawing a scribble. Taking these prompts, *SEEM* can simultaneously produce both masks and semantic labels for the objects. Note that our model is open-vocabulary, which empowers it to label unseen categories when given the candidate vocabulary (i.e., cheetah and butterfly in Fig. 5). When no vocabulary is given, *SEEM* can segment in a class-agnostic manner.

Text referring segmentation. We show the text referring to segmentation visualization results in Fig. 6. As expected, the results demonstrate that our model is semantic-aware of open-vocabulary concepts and attributes to understand



Figure 6: Text to mask or text referring segmentation. The referred text is shown on the masks. *SEEM* adapts to various types of input images in the domain of cartoons, movies, and games.

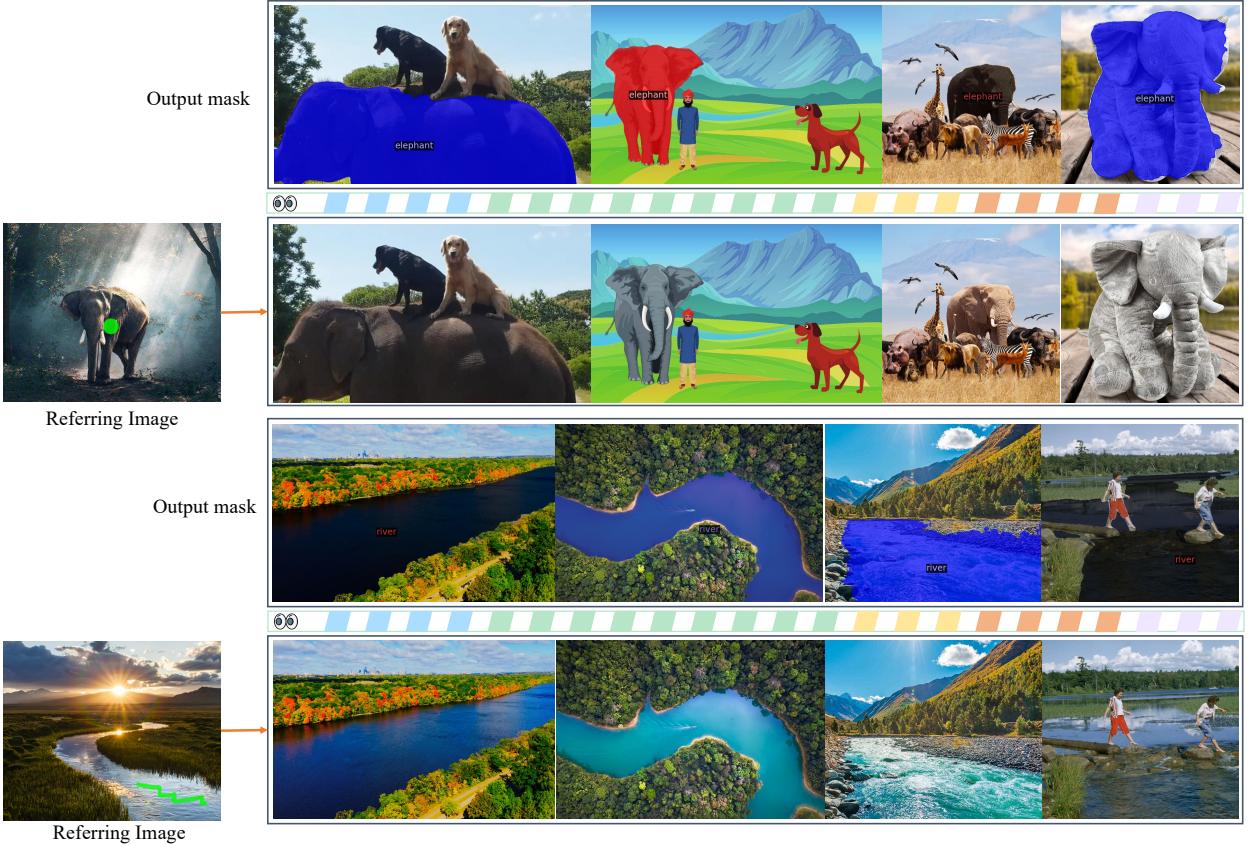


Figure 7: Visual referring segmentation with *SEEM*. Given a referring image with simple spatial hints, *SEEM* can segment the contents which are semantically similar in different target images.

language. In addition, *SEEM* is able to generalize to unseen scenarios like cartoons, movies, and games.

Visual referring segmentation. We visualize *SEEM* the segment with referred regions from another image in Fig.7. By simply drawing a click or scribble on one referring image, *SEEM* can take it as input and segment objects with similar semantics on other images. Notably, this referring

segmentation has a powerful generalization capability to images of other domains. For example, by referring to the elephant in the forest, another object of the same category can be segmented well under drastically different scenes like cartoons, plush toys, and grassland.

In Fig. 8, we further show the referring segmentation ability on the video object segmentation task in a zero-shot



Figure 8: Zero-shot video object segmentation using the first frame plus one stroke. From left to right, the videos are “parkour” and “horsejump-low” from DAVIS [40], video 101 from YouCook2 [64], and a teaser video from EA game “It takes two”. *SEEM* precisely segments referred objects even with significant appearance changes caused by blurring or intensive deformations.

manner. By referring to the objects in the first frame with scribbles, *SEEM* can precisely segment the corresponding objects in the following frames, even when the following objects change in appearance by blurring or intensive deformations.

5. Conclusion

In this paper, we have presented *SEEM*, which can segment everything (all semantics) everywhere all at once (all possible prompt composition). Apart from performing generic open-vocabulary segmentation, *SEEM* can interactively take different types of visual prompts from the user, including click, box, polygon, scribble, text, and referring image segmentation. These visual prompts are mapped into a *joint visual-semantic space* with a prompt encoder, which makes our model versatile to various prompts and can flexibly composite different prompts. Extensive experiments indicate that our model yields competitive perfor-

mance on several open-vocabulary and interactive segmentation benchmarks.

Limitations and future works. We would like to highlight that our *SEEM* is the first and preliminary step toward the universal and interactive interface for image segmentation. Due to the limited resources, our current model still has two main limitations. First, our model is trained on a small scale of segmentation data (mostly COCO [30]). Although it shows strong performance and generality across different datasets and settings, we believe its performance can be further improved by leveraging more training data and supervision. Second, our model does not support part-based segmentation since it was mainly trained with object-level mask annotations. However, given the versatility of our built pipeline, it can seamlessly learn from part-based segmentation data *without* any change to the model. In this sense, we believe the segmentation dataset developed by the concurrent work SAM [22] is a highly valuable resource.

Acknowledgement. We would like to express our gratitude to Lijuan Wang and Lei Zhang for their generous support and valuable suggestions from Jianfeng Wang, Hongyang Li and Zhenyuan Yang.

References

- [1] Tim Brooks, Aleksander Holynski, and Alexei A Efros. Instructpix2pix: Learning to follow image editing instructions. *arXiv preprint arXiv:2211.09800*, 2022.
- [2] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- [3] Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. Sparks of artificial general intelligence: Early experiments with gpt-4. *arXiv preprint arXiv:2303.12712*, 2023.
- [4] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *European Conference on Computer Vision*, pages 213–229. Springer, 2020.
- [5] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *IEEE transactions on pattern analysis and machine intelligence*, 40(4):834–848, 2017.
- [6] Liang-Chieh Chen, George Papandreou, Florian Schroff, and Hartwig Adam. Rethinking atrous convolution for semantic image segmentation. *arXiv preprint arXiv:1706.05587*, 2017.
- [7] Xi Chen, Zhiyan Zhao, Yilei Zhang, Manni Duan, Donglian Qi, and Hengshuang Zhao. Focalclick: towards practical interactive image segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1300–1309, 2022.
- [8] Bowen Cheng, Ishan Misra, Alexander G Schwing, Alexander Kirillov, and Rohit Girdhar. Masked-attention mask transformer for universal image segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1290–1299, 2022.
- [9] Mingyu Ding, Bin Xiao, Noel Codella, Ping Luo, Jingdong Wang, and Lu Yuan. Davit: Dual attention vision transformers. *arXiv preprint arXiv:2204.03645*, 2022.
- [10] Zheng Ding, Jieke Wang, and Zhuowen Tu. Open-vocabulary panoptic segmentation with maskclip. *arXiv preprint arXiv:2208.08984*, 2022.
- [11] Pedro F Felzenszwalb, Ross B Girshick, David McAllester, and Deva Ramanan. Object detection with discriminatively trained part-based models. *IEEE transactions on pattern analysis and machine intelligence*, 32(9):1627–1645, 2009.
- [12] King-Sun Fu and JK Mui. A survey on image segmentation. *Pattern recognition*, 13(1):3–16, 1981.
- [13] Golnaz Ghiasi, Xiuye Gu, Yin Cui, and Tsung-Yi Lin. Open-vocabulary image segmentation. *arXiv preprint arXiv:2112.12143*, 2021.
- [14] Leo Grady. Random walks for image segmentation. *IEEE transactions on pattern analysis and machine intelligence*, 28(11):1768–1783, 2006.
- [15] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 2961–2969, 2017.
- [16] Ronghang Hu, Marcus Rohrbach, and Trevor Darrell. Segmentation from natural language expressions. In *European Conference on Computer Vision*, pages 108–124. Springer, 2016.
- [17] Shijia Huang, Feng Li, Hao Zhang, Shilong Liu, Lei Zhang, and Liwei Wang. A unified mutual supervision framework for referring expression segmentation and generation. *arXiv preprint arXiv:2211.07919*, 2022.
- [18] Dat Huynh, Jason Kuen, Zhe Lin, Jiuxiang Gu, and Ehsan Elhamifar. Open-vocabulary instance segmentation via robust cross-modal pseudo-labeling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7020–7031, 2022.
- [19] Jitesh Jain, Jiachen Li, MangTik Chiu, Ali Hassani, Nikita Orlov, and Humphrey Shi. Oneformer: One transformer to rule universal image segmentation. *arXiv preprint arXiv:2211.06220*, 2022.
- [20] Chao Jia, Yinfai Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc V Le, Yunhsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. In *ICML*, 2021.
- [21] Alexander Kirillov, Kaiming He, Ross Girshick, Carsten Rother, and Piotr Dollár. Panoptic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9404–9413, 2019.
- [22] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and Ross Girshick. Segment anything, 2023.
- [23] Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large language models are zero-shot reasoners. *arXiv preprint arXiv:2205.11916*, 2022.
- [24] Boyi Li, Kilian Q Weinberger, Serge Belongie, Vladlen Koltun, and René Ranftl. Language-driven semantic segmentation. *arXiv preprint arXiv:2201.03546*, 2022.
- [25] Feng Li, Hao Zhang, Shilong Liu, Lei Zhang, Lionel M Ni, Heung-Yeung Shum, et al. Mask dino: Towards a unified transformer-based framework for object detection and segmentation. *arXiv preprint arXiv:2206.02777*, 2022.
- [26] Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation. *arXiv preprint arXiv:2101.00190*, 2021.
- [27] Yuheng Li, Haotian Liu, Qingyang Wu, Fangzhou Mu, Jianwei Yang, Jianfeng Gao, Chunyuan Li, and Yong Jae Lee. Gligen: Open-set grounded text-to-image generation. *arXiv preprint arXiv:2301.07093*, 2023.
- [28] Yin Li, Jian Sun, Chi-Keung Tang, and Heung-Yeung Shum. Lazy snapping. *ACM Transactions on Graphics (ToG)*, 23(3):303–308, 2004.
- [29] Zhiqi Li, Wenhui Wang, Enze Xie, Zhiding Yu, Anima Anandkumar, Jose M Alvarez, Ping Luo, and Tong Lu.

- Panoptic segformer: Delving deeper into panoptic segmentation with transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1280–1289, 2022.
- [30] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *ECCV*, 2014.
- [31] Zheng Lin, Zheng-Peng Duan, Zhao Zhang, Chun-Le Guo, and Ming-Ming Cheng. Focuscut: Diving into a focus view in interactive segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2637–2646, 2022.
- [32] Qin Liu, Zhenlin Xu, Gedas Bertasius, and Marc Niethammer. Simpleclick: Interactive image segmentation with simple vision transformers. *arXiv preprint arXiv:2210.11006*, 2022.
- [33] Qin Liu, Meng Zheng, Benjamin Planche, Srikrishna Karanam, Terrence Chen, Marc Niethammer, and Ziyan Wu. Pseudoclick: Interactive image segmentation with click imitation. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part VI*, pages 728–745. Springer, 2022.
- [34] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3431–3440, 2015.
- [35] Junhua Mao, Jonathan Huang, Alexander Toshev, Oana Camburu, Alan L Yuille, and Kevin Murphy. Generation and comprehension of unambiguous object descriptions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 11–20, 2016.
- [36] Edgar Margffoy-Tuay, Juan C Pérez, Emilio Botero, and Pablo Arbeláez. Dynamic multimodal instance segmentation guided by natural language queries. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 630–645, 2018.
- [37] Shervin Minaee, Yuri Y Boykov, Fatih Porikli, Antonio J Plaza, Nasser Kehtarnavaz, and Demetri Terzopoulos. Image segmentation using deep learning: A survey. *IEEE transactions on pattern analysis and machine intelligence*, 2021.
- [38] OpenAI. Gpt-4 technical report, 2023.
- [39] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744, 2022.
- [40] Federico Perazzi, Jordi Pont-Tuset, Brian McWilliams, Luc Van Gool, Markus Gross, and Alexander Sorkine-Hornung. A benchmark dataset and evaluation methodology for video object segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 724–732, 2016.
- [41] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research*, 21(1):5485–5551, 2020.
- [42] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International Conference on Machine Learning*, pages 8821–8831. PMLR, 2021.
- [43] Yongming Rao, Wenliang Zhao, Guangyi Chen, Yansong Tang, Zheng Zhu, Guan Huang, Jie Zhou, and Jiwen Lu. Denseclip: Language-guided dense prediction with context-aware prompting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18082–18091, 2022.
- [44] Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. Toolformer: Language models can teach themselves to use tools. *arXiv preprint arXiv:2302.04761*, 2023.
- [45] Taylor Shin, Yasaman Razeghi, Robert L Logan IV, Eric Wallace, and Sameer Singh. Autoprompt: Eliciting knowledge from language models with automatically generated prompts. *arXiv preprint arXiv:2010.15980*, 2020.
- [46] Konstantin Sofiiuk, Ilya A Petrov, and Anton Konushin. Reviving iterative training with mask guidance for interactive segmentation. In *2022 IEEE International Conference on Image Processing (ICIP)*, pages 3141–3145. IEEE, 2022.
- [47] 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.
- [48] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed Chi, Quoc Le, and Denny Zhou. Chain of thought prompting elicits reasoning in large language models. *arXiv preprint arXiv:2201.11903*, 2022.
- [49] Chenfei Wu, Shengming Yin, Weizhen Qi, Xiaodong Wang, Zecheng Tang, and Nan Duan. Visual chatgpt: Talking, drawing and editing with visual foundation models. *arXiv preprint arXiv:2303.04671*, 2023.
- [50] Enze Xie, Wenhui Wang, Zhiding Yu, Anima Anandkumar, Jose M Alvarez, and Ping Luo. Segformer: Simple and efficient design for semantic segmentation with transformers. *arXiv preprint arXiv:2105.15203*, 2021.
- [51] Jiarui Xu, Shalini De Mello, Sifei Liu, Wonmin Byeon, Thomas Breuel, Jan Kautz, and Xiaolong Wang. Groupvit: Semantic segmentation emerges from text supervision. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18134–18144, 2022.
- [52] Jiarui Xu, Sifei Liu, Arash Vahdat, Wonmin Byeon, Xiaolong Wang, and Shalini De Mello. Open-vocabulary panoptic segmentation with text-to-image diffusion models. *arXiv preprint arXiv:2303.04803*, 2023.
- [53] Ning Xu, Brian Price, Scott Cohen, Jimei Yang, and Thomas S Huang. Deep interactive object selection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 373–381, 2016.
- [54] Jianwei Yang, Chunyuan Li, and Jianfeng Gao. Focal modulation networks. *arXiv preprint arXiv:2203.11926*, 2022.

- [55] Jianwei Yang, Chunyuan Li, Pengchuan Zhang, Bin Xiao, Ce Liu, Lu Yuan, and Jianfeng Gao. Unified contrastive learning in image-text-label space. In *CVPR*, 2022.
- [56] Zhao Yang, Jiaqi Wang, Yansong Tang, Kai Chen, Hengshuang Zhao, and Philip HS Torr. Latv: Language-aware vision transformer for referring image segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18155–18165, 2022.
- [57] Licheng Yu, Patrick Poirson, Shan Yang, Alexander C Berg, and Tamara L Berg. Modeling context in referring expressions. In *European Conference on Computer Vision*, pages 69–85. Springer, 2016.
- [58] Zhentao Yu, Zixin Yin, Deyu Zhou, Duomin Wang, Finn Wong, and Baoyuan Wang. Talking head generation with probabilistic audio-to-visual diffusion priors. *arXiv preprint arXiv:2212.04248*, 2022.
- [59] Lu Yuan, Dongdong Chen, Yi-Ling Chen, Noel Codella, Xiyang Dai, Jianfeng Gao, Houdong Hu, Xuedong Huang, Boxin Li, Chunyuan Li, Ce Liu, Mengchen Liu, Zicheng Liu, Yumao Lu, Yu Shi, Lijuan Wang, Jianfeng Wang, Bin Xiao, Zhen Xiao, Jianwei Yang, Michael Zeng, Luowei Zhou, and Pengchuan Zhang. Florence: A new foundation model for computer vision. *arXiv preprint arXiv:2111.11432*, 2021.
- [60] Hao Zhang, Feng Li, Huazhe Xu, Shijia Huang, Shilong Liu, Lionel M Ni, and Lei Zhang. Mp-former: Mask-piloted transformer for image segmentation. *arXiv preprint arXiv:2303.07336*, 2023.
- [61] Hao Zhang, Feng Li, Xueyan Zou, Shilong Liu, Chunyuan Li, Jianfeng Gao, Jianwei Yang, and Lei Zhang. A simple framework for open-vocabulary segmentation and detection. *arXiv preprint arXiv:2303.08131*, 2023.
- [62] Lvmin Zhang and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. *arXiv preprint arXiv:2302.05543*, 2023.
- [63] Zihao Zhao, Eric Wallace, Shi Feng, Dan Klein, and Sameer Singh. Calibrate before use: Improving few-shot performance of language models. In *International Conference on Machine Learning*, pages 12697–12706. PMLR, 2021.
- [64] Luowei Zhou, Chenliang Xu, and Jason Corso. Towards automatic learning of procedures from web instructional videos. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018.
- [65] Xueyan Zou, Zi-Yi Dou, Jianwei Yang, Zhe Gan, Linjie Li, Chunyuan Li, Xiyang Dai, Harkirat Behl, Jianfeng Wang, Lu Yuan, et al. Generalized decoding for pixel, image, and language. *arXiv preprint arXiv:2212.11270*, 2022.
- [66] Zhengxia Zou, Zhenwei Shi, Yuhong Guo, and Jieping Ye. Object detection in 20 years: A survey. *arXiv preprint arXiv:1905.05055*, 2019.