

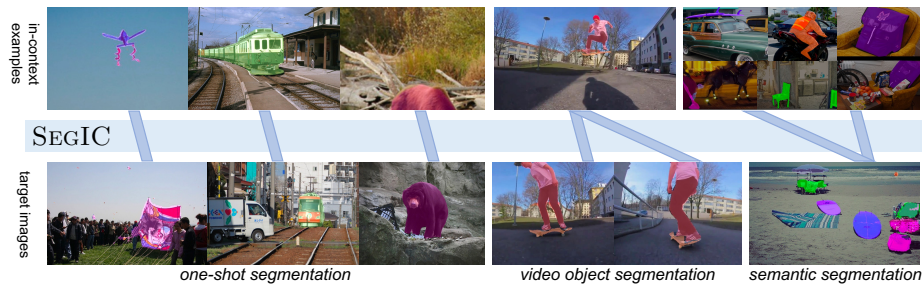
# SEGIC: Unleashing the Emergent Correspondence for In-Context Segmentation

Lingchen Meng<sup>1,2</sup> Shiyi Lan<sup>3</sup> Hengduo Li<sup>4</sup>  
Jose M. Alvarez<sup>3</sup> Zuxuan Wu<sup>1,2†</sup> Yu-Gang Jiang<sup>1,2</sup>

<sup>1</sup>Shanghai Key Lab of Intell. Info. Processing, School of CS, Fudan University

<sup>2</sup>Shanghai Collaborative Innovation Center of Intelligent Visual Computing

<sup>3</sup>NVIDIA <sup>4</sup>University of Maryland



**Fig. 1: Qualitative results of SEGIC.** SEGIC segments target images (the bottom row) according to a few labeled example images (top row, linked by  $\searrow$  in the figure), termed as “in-context segmentation”. SEGIC unifies various segmentation tasks via different types of in-context samples, including those annotated with one mask per sample (one-shot segmentation), annotated with a few masks per sample (video object segmentation), and the combination of annotated samples (semantic segmentation)

**Abstract.** In-context segmentation aims at segmenting novel images using a few labeled example images, termed as “in-context examples”, exploring content similarities between examples and the target. The resulting models can be generalized seamlessly to novel segmentation tasks, significantly reducing the labeling and training costs compared with conventional pipelines. However, in-context segmentation is more challenging than classic ones requiring the model to learn segmentation rules conditioned on a few samples. Unlike previous work with ad-hoc or non-end-to-end designs, we propose SEGIC, an end-to-end **segment-in-context** framework built upon a single vision foundation model (VFM). In particular, SEGIC leverages the emergent correspondence within VFM to capture dense relationships between target images and in-context samples. As such, information from in-context samples is then extracted into three types of instructions, *i.e.* geometric, visual, and meta instructions, serving as explicit conditions for the final mask prediction. SEGIC is a straightforward yet effective approach that yields state-of-the-art performance on one-shot segmentation benchmarks. Notably, SEGIC

<sup>†</sup> Corresponding author.

can be easily generalized to diverse tasks, including video object segmentation and open-vocabulary segmentation. Code will be available at <https://github.com/MengLcool/SEGIC>.

**Keywords:** In-context learning · Segmentation generalist · Vision foundation model

## 1 Introduction

Modern advancements in deep learning have established a routine process for addressing visual perception challenges, typically involving data collection, model training, and deployment. While this pipeline is highly effective, it invariably demands additional effort in data acquisition and model tuning when adapting to new domains. Although researchers have been seeking to learn generic representations with pre-training, the resulting models have to be fine-tuned on the target domain for improved performance.

In contrast, the success of large language models (LLMs) [6, 50, 53, 54] in Natural Language Processing (NLP) offers an alternative approach. These models are trained on vast datasets, handling various NLP tasks through next-token prediction guided by prompts [50]. A key strength of LLMs is their ability to learn from a few examples, a process known as in-context learning (ICL). This enables them to adapt to various tasks with a small and varied set of instructions without requiring extensive fine-tuning or retraining [6, 50]. The success of ICL in NLP highlights the potential for applying similar strategies in visual perception tasks.

While appealing, ICL in vision is particularly challenging as vision tasks are significantly different regarding inputs (2D/3D), outputs (one-hot labels/bounding boxes), and specialized architectures. Recent advances in vision generalist models [32, 74, 75] suggest different levels of segmentation tasks, *i.e.* instance, semantic, and video, can be unified within the same output space. This motivates us to explore ICL using segmentation as a testbed and investigate whether current vision models can be easily generalized. While there are a few previous attempts on ICL for segmentation, they either have fallen short in performance due to implicit modeling [62, 63] or have employed heavy and non-end-to-end pipelines [38, 59, 71], which are less effective and efficient.

At the heart of ICL for NLP tasks is mining the relationships among different words and then propagating labels from a few task-specific question-answer pairs, namely in-context samples to the target one [2, 29, 65]. We argue that in vision tasks, the similar entity that facilitates label propagation from in-context samples to novel samples is establishing dense correspondences between images. Although dense visual correspondences are difficult to obtain before the era of foundation models, recent studies [8, 59] have shown that high-quality correspondence emerges in visual foundation models (VFMs) [8, 49, 52, 56].

In light of this, we introduce SEGIC, an end-to-end **segment-in-context** framework without the need for sophisticated handcrafted prompt design. Specif-

ically, our framework is built upon a single frozen vision foundation model followed by a lightweight mask decoder. We leverage the emergent correspondence of the VFM to establish dense correspondences across target images and in-context samples. Based on that, we extract in-context information into three types of instructions: geometric, visual, and meta instructions. By explicitly utilizing these instructions, our model demonstrates remarkable generalization capabilities with low training costs across diverse segmentation tasks, as evidenced by extensive qualitative and quantitative experiments. We summarize our contributions in threefold:

- We introduce SEGIC, a simple yet effective in-context segmentation framework, exploring the strong emergent correspondence encoded in VFMs.
- We design geometric, visual, and meta instructions that explicitly transfer knowledge from in-context samples to the target to facilitate in-context segmentation without tuning the parameters of vision foundation models.
- SEGIC demonstrates state-of-the-art performance on COCO-20<sup>i</sup>, FSS-1000 and recent LVIS-92<sup>i</sup>. Moreover, we conduct a comprehensive study on vision foundation models of various pre-text tasks, model sizes, and pre-training data. Furthermore, SEGIC achieves competitive performance on novel tasks including video object segmentation and open-vocabulary segmentation, without ever seeing their training data.

## 2 Related Work

**Vision foundation models.** Recent years have witnessed great progress in large-scale vision pre-training [3,8,10,17,18,27,52], serving as the cornerstone for high-capacity foundation models. These pre-training approaches can be broadly categorized into two directions: vision-only pretraining [3, 8, 10, 17, 18, 49] and vision-language pre-training [24,33,52,56,69]. For vision-only pre-training, models aim to distinguish image/patch-level entities from different views [8, 10, 18] or reconstruct masked images [3, 18] from raw images. In vision-language pre-training, models strive to align cross-modal features into a unified visual-semantic space [24, 33, 52, 69], showcasing great open-set performance due to the transferable capabilities of language. Unlike these approaches that perform pre-training in an unsupervised or weakly supervised manner, SAM [27] is pre-trained on a huge amount of labeled segmentation data with precise prompts about locations. In this paper, we conduct extensive experiments using three types of pre-trained models as backbones to explore their potential for in-context segmentation. Interestingly, we observe that models with higher zero-shot semantic and geometric correspondence performance are more likely to be effectively utilized in our SEGIC framework for in-context segmentation.

**Unified vision downstream models.** Unified vision downstream models have recently drawn significant attention due to their generalization capabilities and flexibilities. Unlike previous specialized vision models designed for specific datasets [9, 19, 40, 55, 64], vision generalists are tailored to handle multiple datasets [15, 43,

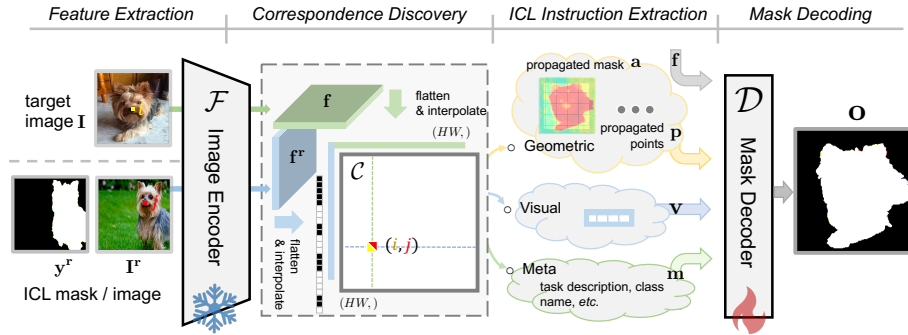
[44, 73] and a wide range of tasks [23, 74, 75] within a single yet unified model. Recently, many studies [23, 27, 62, 63, 74, 75] have focused on developing techniques that unify segmentation tasks. In a similar spirit, our work also builds upon a unified output space for segmentation tasks. However, our goal is different—we aim to perform in-context segmentation that allows a model to effectively segment novel entities conditioned on a few samples.

**Visual correspondence.** Establishing visual correspondences between different images is vital for various computer vision tasks. Traditionally, computing correspondences is based on hand-designed features like SIFT [42] and SURF [5]. Deep models [26, 30, 61] can also learn correspondence in a supervised manner. More recently, it has been shown features from large foundation models encode dense visual correspondence clues [3, 8, 17, 18, 49]. In this study, we discover a profound connection between correspondence and in-context segmentation—correspondence acts as explicit guidance, linking the target image with in-context images, thus facilitating label propagation for in-context segmentation.

**In-context learning.** For the first time, GPT-3 [6] introduces a new learning paradigm known as in-context learning, which unifies various NLP tasks as text completion or question-answering tasks using provided prompts and examples. This approach enables language models to handle various tasks, including novel ones, by leveraging task examples, without requiring re-training or fine-tuning. Recent studies [1, 4, 38, 62, 63] explore this mechanism in vision tasks. Painter [62] and SegGPT [63] aim to achieve in-context segmentation through in-painting. They build upon the Mask Image Modeling (MIM) framework [17] to concatenate images and predictions into a  $2 \times 2$  mosaic and make predictions by recovering the masked areas. In their pipelines, the vision backbone serves as both an image encoder and a mask decoder, which incurs significant computational costs. Moreover, they struggle to effectively leverage pre-trained models due to input shifts, leading to increased convergence challenges. Other approaches [38, 71] attempt using in-context segmentation via prompting SAM [27]. They build upon cross-image correspondences between in-context examples and target images by additional pre-trained models to generate prompts for SAM. However, these methods employ a two-stage pipeline, introducing redundancy and repeated computations. Consequently, if the model encounters limitations in one stage, it negatively impacts the final performance. In this work, we build an end-to-end in-context segmentation framework, leveraging the emergent correspondence using a single vision foundation model.

### 3 Approach

In-context learning equips a model with the ability to learn from example images, namely “in-context examples”, as humans, which has demonstrated great potential in NLP tasks [6, 50]. This process is akin to how humans intuitively grasp and replicate complex patterns from just a few guides, rapidly generalizing to new examples. In this paper, our goal is to establish an end-to-end in-context



**Fig. 2: Architecture overview.** SEGIC is built upon a **frozen** vision foundation model, consisting of four stages: (1) feature extraction; (2) correspondence discovery (Section § 3.1); (3) in-context instruction extraction (Section § 3.2); (4) mask decoding (Section § 3.3).

segmentation framework, enabling a model to effectively segment novel entities beyond the training set, such as video objects, with low training costs.

Formally, given a target image  $\mathbf{I} \in \mathbb{R}^{3 \times H \times W}$  where  $H$  and  $W$  represent the height and width, the goal is to segment a binary mask  $\mathbf{y} \in \mathbb{R}^{\{0,1\}^{H \times W}}$  conditioned on  $K$  in-context examples  $\{\mathbf{x}^r\}^K = \{(\mathbf{I}^r, \mathbf{y}^r)\}^K$ , where  $\mathbf{I}^r, \mathbf{y}^r$  are the image and ground-truth mask of an in-context example.

Our approach consists of four stages: feature extraction, correspondence discovery, in-context instruction extraction, and mask decoding, as shown in Figure 2. The feature extraction stage follows the common practice in visual perception tasks where we use a pre-trained visual foundation model to obtain the feature maps of both in-context and target images. Subsequently, we compute the dense correspondences based on their semantics and geometry (Section § 3.1), providing explicit guidance of the segmentation rules. Then, we extract in-context instructions based on the in-context samples and the dense correspondences (Section § 3.2). Finally, SEGIC uses a lightweight in-context mask decoder  $\mathcal{D}$  to produce segmentation masks for the desired target, leveraging the extracted information from in-context examples (Section § 3.3).

### 3.1 Dense Correspondence Discovery

To establish the relations between the target image  $\mathbf{I}$  and the reference image  $\mathbf{I}^r$  (considering one reference example for brevity), we extract dense cross-image correspondences at the pixel level. For this purpose, we leverage pre-trained VFMs due to their powerful generalization capability and emergent correspondence properties. Specifically, we first extract the visual features of  $\mathbf{I}$  and  $\mathbf{I}^r$  with a vision foundation model, then apply the cosine distance function to compute the patch-level distance map between the two images as shown in Figure 2. We further obtain the pixel-level correspondences by interpolating the patch-level

distance map to the size of the original image as follows:

$$\begin{aligned} \mathbf{f}, \mathbf{f}^r &= \mathcal{F}(\mathbf{I}), \mathcal{F}(\mathbf{I}^r) \\ \mathcal{C} &= \text{Upsample}(\text{Dist}(\mathbf{f}, \mathbf{f}^r)), \end{aligned} \quad (1)$$

where  $\mathcal{F}$  indicates the vision foundation model we use and  $\mathbf{f}, \mathbf{f}^r \in \mathbb{R}^{c \times h \times w}$  are the extracted features of the target and the reference images, respectively;  $c, h, w$  indicate the dimensions of the feature maps.  $\text{Dist}$  denotes the distance function (for which we use cosine distance) and  $\text{Upsample}$  is the interpolation function.  $\mathcal{C} \in \mathbb{R}^{HW \times HW}$  denotes the calculated dense correspondences between the target image and the in-context image.

### 3.2 In-context Instruction Extraction

After obtaining dense correspondences, the question becomes how to utilize the in-context information, including in-context examples and dense correspondences, as instructions to guide the segmentation process. Here we extract in-context instructions based on ideas from NLP tasks [29, 50]. Ideally, the representation for in-context information should clearly articulate how segmentation should be executed on the target image while being concise and efficient for effective segmentation. To this end, we decouple and encode the in-context information into three individual in-context instructions as shown in Figure 2: 1) geometric instructions; 2) visual instructions; and 3) meta instructions, each of which will be elaborated below:

**Geometric instructions** aim to provide a coarse location estimation of the target mask. Although we only have the mask annotation for the reference image, we can propagate the label from the reference to the target to obtain a propagated mask by exploring the dense correspondences  $\mathcal{C}$ :

$$\begin{aligned} \mathbf{a} &= \mathbf{y}^r / \|\mathbf{y}^r\| \cdot \mathcal{C} \in \mathbb{R}^{H \times W} \\ \mathbf{p} &= \text{PE}(\text{Topk}(\mathbf{a})), \end{aligned} \quad (2)$$

where  $\mathbf{a}$  represents the propagated coarse mask for the target image. As shown in Equation (2), by performing matrix multiplication, we gather and average the dense correspondences based on the positive points in the reference mask. This process is analogous to propagating labels from the reference image to the target image through dense correspondences, making  $\mathbf{a}$  a form of dense geometric instruction. Additionally, based on  $\mathbf{a}$ , we employ Topk (as referred to in Equation (2)) to select the top-k points with the highest values in  $\mathbf{a}$ , indicating the locations most likely relevant to the target mask. We further encode the 2D coordinates into high-dimensional vectors using cosine positional encoding PE as in [58, 60], resulting in a type of sparse geometric instruction. Overall,  $\mathbf{a}$  and  $\mathbf{p}$  together provide geometric information extracted from in-context examples.

**Visual instructions** indicate the visual clues of the target entity. We utilize mask pooling [66] on reference image features to extract visual salient clues,  $\mathbf{v}$ :

$$\mathbf{v} = \mathbf{y}^r / \|\mathbf{y}^r\| \cdot \mathbf{f}^r. \quad (3)$$

By doing so, only information relevant to the reference mask, including low-level (texture, appearance) and high-level (semantic) clues, is used.

**Meta instructions** indicate other clues implicitly provided by in-context examples, such as task descriptions, class names, *etc.* We uniformly treat them as languages and encode them with a pre-trained language model, following [23,27,37].

$$\mathbf{m} = \mathcal{F}^t(\mathbf{meta}) \quad (4)$$

where  $\mathcal{F}^t$  is the pre-trained CLIP text encoder [52], **meta** indicates the meta information and **m** is the meta feature.

Finally, we use  $\mathbf{c} = \{\mathbf{a}, \mathbf{p}, \mathbf{v}, \mathbf{m}\}$  to denote the set of in-context instructions derived from reference samples, which can be further used for producing the segmentation mask, as will be introduced in section § 3.3.

### 3.3 Mask Decoding

In this section, we discuss how to predict the segmentation masks in target images following the aforementioned in-context instructions. In particular, we use a query-based mask decoder  $\mathcal{D}$  due to its high performance in segmentation tasks and flexibility [7, 11, 12]. Formally, the decoder network takes the target image feature  $\mathbf{f}$  and in-context instructions  $\mathbf{c}$  as input, as well as a learnable query  $\mathbf{q}$ , and outputs a mask  $\mathbf{o}$  as follows:

$$\mathbf{o} = \mathcal{D}(\mathbf{f}, \mathbf{c}; \mathbf{q}). \quad (5)$$

Unlike previous designs that use a set of object queries [7, 11, 12] for prediction, we only initialize one query since we just need to predict one mask in-context conditioned on the in-context instructions as shown in Figure 2.

To prepare these instructions for mask decoding, we first project them into the latent space used by the decoder. We categorize the in-context instructions into two types according to their spatial property: instructions with spatial shapes (*i.e.* the propagated coarse mask  $\mathbf{a}$ ) and without spatial shapes (*i.e.*  $\mathbf{p}, \mathbf{v}, \mathbf{m}$ ). For the spatial instructions, we employ a series of convolutional layers  $\mathcal{M}$  to encode it into the image feature space; for the non-spatial instructions, we use projection layers  $\mathcal{P}$  to project them into the query feature space:

$$\begin{aligned} \mathbf{q}^{\mathbf{p}}, \mathbf{q}^{\mathbf{v}}, \mathbf{q}^{\mathbf{m}} &= \mathcal{P}^{\mathbf{p}}(\mathbf{p}), \mathcal{P}^{\mathbf{v}}(\mathbf{v}), \mathcal{P}^{\mathbf{m}}(\mathbf{m}) \\ \mathbf{a}' &= \mathcal{M}(\mathbf{a}) \end{aligned} \quad (6)$$

where  $\mathcal{P}^{\mathbf{p}}, \mathcal{P}^{\mathbf{v}}, \mathcal{P}^{\mathbf{m}}$  indicate the projection layers for  $\mathbf{p}, \mathbf{v}, \mathbf{m}$ , respectively.  $\mathbf{a}'$ ,  $\mathbf{q}^{\mathbf{p}}$ ,  $\mathbf{q}^{\mathbf{v}}$  and  $\mathbf{q}^{\mathbf{m}}$  are the projected in-context instructions.

Furthermore, we inject the projected in-context instructions into the decoding stage to guide the decoder to segment in-context. Similarly, for the spatial features, we add them to image features, such that they are aware of the coarse mask produced by reference samples. For the non-spatial features, we

concatenate the initial query with them, which allows a deeper interaction via a self-attention mechanism in the decoder:

$$\begin{aligned} \mathbf{f}' &= \mathbf{f} + \mathbf{a}' \\ \mathbf{q}' &= \text{Concat}(\mathbf{q}, \mathbf{q}^1, \mathbf{q}^s, \mathbf{q}^m) \\ \mathbf{o} &= \mathcal{D}(\mathbf{f}'; \mathbf{q}') \end{aligned} \tag{7}$$

Finally, the mask prediction  $\mathbf{o}$  is produced based on image features  $\mathbf{f}'$  conditioned on instructions, and query features  $\mathbf{q}'$  as shown in Equation (7). For more details, please refer to our supplementary.

### 3.4 Training Pipeline

During training, we freeze all the parameters of the VFM and only leave the newly introduced mask decoder trainable. We employ a linear combination of a dice loss [57] and a binary cross-entropy loss for our mask loss:  $L_{mask} = \lambda_{ce}L_{ce} + \lambda_{dice}L_{dice}$ . It is worth noting that we calculate the segmentation loss on  $K$  selected points using importance sampling following [11, 28] instead of the whole image to save memory cost. To further improve the robustness toward noisy in-context samples, we introduce two strategies into our training recipe, namely “context reversion” and “negative entity augmentation”.

**Context reversion.** We artificially introduce noisy context during training to improve the robustness. To simulate situations where in-context examples are inaccurate, we propose “context reversion”: swapping the target and reference images—we use the prediction of the target image as an in-context example. The noisy context introduces randomness during training and hence can improve robustness.

**Negative entity augmentation.** In tasks like video object segmentation (VOS), in-context examples for different entities in the same image are mutually exclusive. Taking the case of video object segmentation in Figure 1 as an example, the person and the skateboard are exclusive in one image. Thus, entities that are not of interest can serve as negative samples, indicating that they are not relevant to the target. We augment the in-context instructions with these negative entities for a better result.

## 4 Experiments

### 4.1 Training Data

We train SEGIC on semantic and instance segmentation datasets. Unlike previous sophisticated task-specific designs for multiple datasets/tasks, our method offers a simple approach that distinguishes tasks by in-context examples.

**COCO** [35] is a widely used instance/semantic segmentation with 83K training samples and 41K samples of 80 categories. We use the 2014 split version to be consistent with the COCO-20<sup>i</sup> one-shot semantic segmentation benchmark.



**ADE20k** [72] is a widely used semantic segmentation dataset for 150 categories, with 20K training images.

**LVIS** [16] is a large instance segmentation dataset containing 1000+ categories with  $\sim$ 100K training images.

## 4.2 Training Details

We implement SEGIC in PyTorch and use 8 V100 GPUs for most of our experiments. We use a batch size of 32 (4 per GPU) in total, 1 point for point instruction, and 12544 points per mask to calculate mask loss following [11]. We merge the COCO instances for its semantic segmentation. Please refer to our supplementary material for the detailed training pipeline. Thanks to the frozen backbone, our method is extremely memory-efficient (with less than 10G memory cost in most of our experiments). SEGIC is a single unified model that is jointly trained on mixed datasets, while evaluated on various datasets, separately. We utilize DINOv2 [49] as the default vision foundation model, with a ViT-B for all ablations, and ViT-L/G for the main experiments. We employ an AdamW [41] optimizer with a weight decay of  $1e-4$ . We set an initial learning rate as  $1e-4$  and multiply 0.1 at the 10 epoch during training. Each dataset is sampled uniformly with 160K samples per epoch in the main experiments and 80K samples for the ablations. We perform data augmentations on target images and reference images respectively. We use large-scale jittering augmentation for semantic segmentation datasets, and normal data augmentations, including random resizing cropping, color jittering, and random horizontal flipping, for instance segmentation datasets. We random sample 1 mask per image for semantic segmentation datasets during training, while up to 10 masks for instance segmentation. The size of a single image is cropped/padded to  $896 \times 896$ .

## 4.3 Main Results

For the main experiments, we compare SEGIC with other specialist/generalist methods on several benchmarks of different segmentation tasks. We use SEGIC of DINOv2 [49] of large and giant versions as the backbone for the main experiments.

**One-shot semantic segmentation.** To demonstrate the generalization capability from known categories to unknown ones, we evaluate SEGIC on one-shot semantic segmentation benchmarks. Following SegGPT [63], we evaluate SEGIC in two one-shot semantic segmentation settings: in-domain using COCO-20<sup>i</sup> [47] (the training set of COCO-20<sup>i</sup> is a subset of ours) and out-of-domain with FSS-1000 [34] (without seeing any training samples of FSS-1000). As depicted in Table 1, SEGIC has achieved state-of-the-art performance in both of these settings. To ensure a fair comparison, we report in-domain performance of specialist models for COCO-20<sup>i</sup> reported in [63]. Additionally, for FSS-1000, we highlight the performance trained on FSS-1000 in gray and zero-shot results in black. Notably, SEGIC outperforms previous generalist models by a significant margin

**Table 1: Main results on several segmentation benchmarks.** Param<sup>t</sup> indicates the number of trainable parameters. † indicates relying SAM; \* indicates that needs additional fine-tuning, # indicates that classes within images are known during inference; n/a indicates the model does not have capability for the task and - indicates that do not have reported number. For FSS-1000 and VOS datasets, we highlight the performance trained on corresponding datasets in gray and zero-shot results in black.

Method	Param <sup>t</sup>	one-shot segmentation			video object segmentation		semantic seg		open-vocab seg	
		COCO-20 <sup>i</sup>	FSS-1000	LVIS-92 <sup>i</sup>	DAVIS-17	YVOS-18	COCO ADE20k	PC-459	A-847	
		mean mIoU	mIoU	mean mIoU	$\mathcal{J}\&\mathcal{F}$	G	mIoU	mIoU	mIoU	mIoU
<i>few-shot seg specialist</i>										
HSNet (RN50) [45]	28M	41.7	86.5	17.4	n/a	n/a	n/a	n/a	n/a	n/a
VAT (RN50) [20]	52M	42.9	90.3	18.5	n/a	n/a	n/a	n/a	n/a	n/a
FPTrans (B) [70]	101M	56.5	-	-	n/a	n/a	n/a	n/a	n/a	n/a
<i>VOS specialist</i>										
AGAME (RN101) [25]	-	n/a	n/a	n/a	70.0	66.0	n/a	n/a	n/a	n/a
SWEM (RN50) [36]	58M	n/a	n/a	n/a	84.3	82.8	n/a	n/a	n/a	n/a
XMem (RN50) [13]	62M	n/a	n/a	n/a	87.7	86.1	n/a	n/a	n/a	n/a
<i>semantic seg specialist</i>										
MaskFormer (L) [12]	212M	n/a	n/a	n/a	n/a	n/a	64.8	54.1	n/a	n/a
Mask2Former (L) [11]	216M	n/a	n/a	n/a	n/a	n/a	67.4	56.1	n/a	n/a
MaskDINO (L) [31]	225M	n/a	n/a	n/a	n/a	n/a	-	56.6	n/a	n/a
<i>segmentation generalist</i>										
OneFormer (L) [23]	237M	n/a	n/a	n/a	n/a	n/a	67.4	57.7	-	-
UNINEXT (L) [68]	340M	n/a	n/a	n/a	77.2	78.1	-	-	-	-
X-decoder (L) [74]	341M	n/a	n/a	n/a	n/a	n/a	67.5	58.1	29.6	9.2
SEEM (L) [75]	341M	-	-	-	58.9	50.0	67.6	-	-	-
<i>in-context generalist</i>										
Painter (L) [62]	354M	33.1	61.7	10.5	34.6	24.1	-	49.9	-	-
SegGPT (L) [63]	354M	56.1	85.6	18.6	75.6	74.7	-	*39.6	-	-
PerSAM <sup>†</sup> (H) [71]	0	23.0	71.2	11.5	60.3	-	n/a	n/a	n/a	n/a
PerSAM-F <sup>†</sup> (H) [27]	2	23.5	75.6	12.3	71.9	-	n/a	n/a	n/a	n/a
Matcher <sup>†</sup> (H+G) [38]	0	52.7	87.0	33.0	79.5	-	n/a	n/a	n/a	n/a
SEGIC (L)	5M	76.1	86.8	44.6	71.4	62.7	#72.9	#55.5	#33.5	#18.9
SEGIC (G)	5M	74.5	88.4	47.8	73.7	65.4	#74.0	#59.0	#34.9	#20.1

(more than 20 of mean mIoU) on COCO-20<sup>i</sup> and achieves competitive results that are very close to specialist models on FSS-1000, even without ever being trained on it. Furthermore, we conduct experiments on LVIS-92<sup>i</sup> [38], a more challenging one-shot benchmark built upon LVIS [16]. On LVIS-92<sup>i</sup>, SEGIC surpasses Matcher, the previous SoTA, by a large margin (from 33.0 to 47.8).

Since we include the evaluation categories in the training process for COCO-20<sup>i</sup> in Table 1, for a rigorous comparison, we follow its training setting that trains and tests our model on 4 splits [47] separately, avoiding seeing categories for evaluation. As shown in Table 2, our method still achieves state-of-the-art performance on COCO-20<sup>i</sup>.

Furthermore, as shown in the last row of Table 2, by joint training on COCO-excluded datasets (including ADE20k, LVIS, and FSS-1000), there is a significant performance gain across all splits. This further demonstrates the effectiveness of our in-context generalization capabilities.

Overall, our best model surpasses all previous segmentation generalist models across all one-shot segmentation benchmarks, demonstrating its effectiveness.

**Table 2: Comparisons on one-shot COCO-20<sup>i</sup>.** ‡ indicates that is jointly trained on the COCO-excluded datasets.

Method	F0	F1	F2	F3	mean
HSNet [45]	37.2	44.1	42.4	41.3	41.2
VAT [20]	39.0	43.8	42.6	39.7	41.3
FPTrans [70]	44.4	48.9	50.6	44.0	47.0
MSANet [22]	47.8	57.4	48.7	50.5	51.1
SEGIC	55.8	54.7	52.4	51.4	53.6
SEGIC ‡	<b>62.3</b>	<b>62.5</b>	<b>63.3</b>	<b>60.9</b>	<b>62.3</b>

**Zero-shot video object segmentation.** Video object segmentation (VOS) aims to segment specific objects in video frames. In this work, we focus on the semi-supervised VOS setting [51, 67], where the masks that appeared first time are given as references. We evaluate SEGIC without any fine-tuning on video datasets to demonstrate our generalization capabilities. We choose two commonly used VOS datasets: DAVIS-17 [51] and YouTube-VOS-18 [67]. We export two metrics commonly used in VOS for evaluation: the  $\mathcal{J}\&\mathcal{F}$  score for DAVIS-17 and  $G$  score for YouTube-VOS-18, with their official evaluation servers or toolkits. As shown in Table 1, when compared to VOS specialist models, SEGIC achieves competitive performance on VOS benchmarks, even without seeing any training videos. Furthermore, in comparison to segmentation generalist models, SEGIC surpasses Painter [62], SEEM [75], and PerSAM [71] by a significant margin, and competes favorably with recent generalist models [38, 63]. Additionally, the VOS task pipeline in SEGIC is relatively simple. We do not use any test time augmentation tricks (TTA) used in [13]. Moreover, it does not involve dense feature interaction at the patch level, as in SegGPT, and does not require the use of a pre-trained SAM for segmentation.

**Generic semantic segmentation.** We also evaluate SEGIC on generic semantic segmentation benchmarks, which need to segment dataset-dependent pre-defined categories within each image. We adapt generic semantic segmentation to our in-context learning framework by gathering in-context examples from the training set and then performing segmentation in an in-context manner. We use two widely-used semantic segmentation datasets, COCO [35] and ADE20k [72], for evaluation. As illustrated in Table 1, compared to specialist and generalist models for semantic segmentation, SEGIC demonstrates strong performance under the settings that the classes are known in advance (marked with #). As shown in Table 1, our best model achieves 74.0 and 59.0 mIoU on COCO and ADE20k, respectively, surpassing the previous specialist/generalist models.

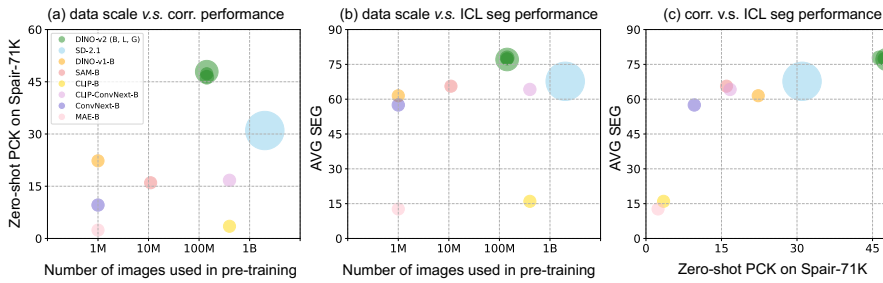
**Open-vocabulary semantic segmentation.** Similar to generic semantic segmentation, we can also extend SEGIC to open-vocabulary semantic segmentation. Furthermore, we could also utilize StableDiffusion [56] to synthesize images with coarse masks for those categories lacking in-context examples like [48]. Compared to X-decoder [74], our method obtains competitive and even better results on PC-459 and ADE-847.

#### 4.4 Rethinking Vision Foundation Model Pre-training

Existing works improve VFMs by scaling up data and model sizes, modifying the model architectures, and utilizing different pre-text tasks. It is unknown how much those factors affect the performance when the encoder is frozen. We delve deeper into the potential of various **frozen** pre-trained VFMs [8, 17, 21, 27, 39, 49, 52, 56]. For each VFM, we freeze it to extract image features and dense correspondences then map them into the hidden space for mask decoding, equally. Additionally, we assess the **zero-shot** semantic correspondence score to underscore the emerging downstream capabilities from different perspectives directly.

**Table 3: Ablation on foundation models as backbone.** HRes denotes high-resolution pre-training; MRes denotes multi-resolution pre-training or intrinsic support within the architecture, *e.g.* CNN and UNet; Task<sup>p</sup> denotes the pre-training task. ■ means that we only obtain a trivial performance.

Encoder	Arch	HRes	MRes	Task <sup>p</sup>	Data <sup>p</sup>	COCO-20 <sup>i</sup>	FSS-1000	DAVIS-17	AVG
DINOv2-B [49]	ViT	✓	✓	image/patch self-distillation	142M	74.6	87.4	68.4	76.8
DINOv2-L [49]	ViT	✓	✓	image/patch self-distillation	142M	75.4	86.9	68.9	77.1
DINOv2-G [49]	ViT	✓	✓	image/patch self-distillation	142M	75.7	87.5	70.1	77.8
DINOv1-B [8]	ViT		✓	image self-distillation	1M	53.5	76.9	54.2	61.5
SAM-B [27]	ViT	✓		interactive segmentation	11M	59.8	76.9	60.1	65.6
SD-2.1 [56]	UNet	✓	✓	text-to-image generation	2B	66.6	84.2	52.4	67.7
OpenCLIP-ConvNext-B [21]	CNN		✓	image-text alignment	400M	65.1	77.0	50.4	64.2
CLIP-ViT-B [52]	ViT			image-text alignment	400M	15.3	29.4	1.3	15.3
MAE-B [17]	ViT			mask image modeling	1M	12.0	23.7	2.3	12.6
ConvNext-B [39]	CNN		✓	image classification	1M	62.0	72.8	37.6	57.5



**Fig. 3: Performance on zero-shot semantic correspondence and in-context segmentation.** We use the dense feature similarity for correspondence estimation. The diameter of each bubble represents the number of parameters of each model.

More specifically, we evaluate a percent of correct keypoints score (PCK) on SPair-71k [46] following to [14] with cosine similarity between the dense feature maps of the two images for semantic correspondence. We observe that (1) high/multi-resolution is essential; (2) pre-text tasks are important; (3) model and data scales do not necessarily help for ICL segmentation; (4) zero-shot correspondence can imply the ICL segmentation capacity.

**High/multi-resolution is essential.** It is evident in Table 3 that the models (CLIP-ViT-B and MAE-B), lacking support for high-resolution or multi-resolution only exhibit trivial performance within our frozen-backbone framework. Despite the remarkable success of vision transformers, they encounter limitations in directly adapting to various input resolutions without fine-tuning due to fixed-length visual tokens and position embeddings. Unlike vision transformers, CNNs, with their stacked convolution design, demonstrate seamless adaptation of parameters to images of different resolutions, even only pre-trained on low-resolution images. To alleviate this limitation, DINO-v1/v2 enhance vanilla ViTs with multi-resolution inputs and position embedding interpolation during pre-training. Consequently, they achieve superior segmentation performance. Overall, the observations indicate that the high/multi-resolution capabilities of the pre-trained VFMs play a key role in in-context segmentation.

**Pre-text tasks are important.** Once the conditions for supporting high-resolution inputs are met, pre-text tasks employed during pre-training and the scale of data also impact the in-context segmentation capacities. Illustrated

in Figure 3, with the same model size and pre-training data, DINOv1 outperforms those classification-pre-trained models in both zero-shot semantic correspondence and in-context segmentation tasks by a large margin. This emphasizes the substantial influence of the pre-text task on downstream capabilities.

**Model and data scales do not directly aid ICL segmentation.** Equipping with larger-scale data and advanced self-distillation techniques, DINOv2 outperforms DINOv1 by a considerable margin. However, as shown in Table 3, scaling up the model size from base (B) to giant (G), is not a clear performance boost. Furthermore, compared to those pre-trained on than 10 times data, *e.g.* OpenCLIP-ConvNext-B and SD-2.1, DINOv2 still exhibits significant advantages. Those demonstrate that model and data scales do not directly help.

**Zero-shot correspondence implies the segmentation capacity.** We further study the relationship between correspondence and segmentation performance with different backbones. As shown in Figure 3, the performance on segmentation is proportional to correspondence, demonstrating a strong consistency between these two tasks. It further confirms our motivation to leverage the emergent correspondence for in-context segmentation. This insight further inspires us, suggesting that pre-training emphasizing inter/intra-image correspondence may lead to better in-context segmentation potentials, *e.g.* the image/patch-level discriminative self-distillation [49].

#### 4.5 Ablation Study

In the ablation study, we report the performance on COCO-20<sup>i</sup> (in-domain one-shot segmentation), FSS-1000 (out-of-domain one-shot segmentation) and DAVIS-17 (zero-shot video object segmentation) to investigate the in-domain convergence and out-of-domain generalization capability.

**Table 4: Component ablation on in-context instructions.**  $\checkmark$  is the default setting for in-context instructions.

geometric	visual	meta	COCO-20 <sup>i</sup>	FSS-1000	DAVIS-17
$\checkmark$			68.8	87.3	65.8
	$\checkmark$		71.0	85.1	67.6
		$\checkmark$	73.7	63.7	21.6
			48.9	58.1	21.1
$\checkmark$	$\checkmark$	$\checkmark$	<b>74.6</b>	<b>87.4</b>	<b>68.4</b>

**Table 5: Ablations the training strategies.**  $\checkmark$  is the default training scheme in our main experiments.

reversion	negative	COCO-20 <sup>i</sup>	FSS-1000	DAVIS-17
		70.0	86.1	66.3
	$\checkmark$	72.9	86.9	65.0
$\checkmark$		73.1	86.5	64.3
$\checkmark$	$\checkmark$	<b>74.6</b>	<b>87.4</b>	<b>68.4</b>

**Ablations on in-context instructions.** We study the importance of each component of in-context instructions (*i.e.* geometric, visual, and meta instructions). We conduct ablations on different combinations of components. As shown in Table 4, we find that the geometric and visual instructions tend to help the performance for out-of-domain segmentation and meta instructions (class name and task description in the ablation) benefit in-domain performance: when each component is used individually, the geometric and visual instructions obtain the best results on FSS-1000 and DAVIS-17, respectively. Meanwhile, meta instructions achieve the best performance on COCO-20<sup>i</sup>, but with poor results on the other two datasets. Encouragingly, our method obtains the best performance

among three tasks when using all three prompts, further demonstrating that our model can effectively transfer knowledge from in-context samples with the proposed in-context instructions.

**Ablations on training strategies.** We investigate how the proposed training strategies affect performance. As shown in Table 5, since the two strategies, *i.e.* context reversion and negative entity augmentation, act as a form of “in-context augmentation”, they bring performance gains on COCO-20<sup>i</sup> and FSS-1000<sup>i</sup> with a slight drop on DAVIS-17. Furthermore, after combining them, our model achieves a consistent gain across all datasets, highlighting its efficacy.

**Table 6: Ablations on training data.**   is the default data combination.

<i>semantic segmentation</i>		<i>instance segmentation</i>		COCO-20 <sup>i</sup>	FSS-1000	DAVIS-17	AVG
COCO <sub>sem</sub>	ADE20k	COCO <sub>inst</sub>	LVIS				
✓				<b>76.3</b>	80.9	45.1	67.4
✓	✓			75.6	82.3	40.1	66.0
✓		✓		75.7	83.5	<b>70.1</b>	76.4
✓	✓	✓	✓	74.6	<b>87.4</b>	68.4	<b>76.8</b>

**Dataset Combination.** We further investigate the effectiveness of each dataset under our joint in-context training framework. We group the training data into two types: semantic segmentation (COCO<sub>sem</sub>, ADE20k) and instance segmentation (COCO<sub>inst</sub>, LVIS). As shown in Table 6, it is clear to see that when trained on COCO<sub>sem</sub> only, it achieves the best performance on COCO-20<sup>i</sup> but is weak on the other. Based on this, we can enrich the training data from two directions: (1) use extra semantic segmentation data and (2) introduce instance segmentation tasks into training. For the former direction, the performance on FSS-1000 increases after introducing ADE20k. For the latter, we still use COCO as the training set but introduce its instance-level annotation. It can be seen a clear performance boost on DAVIS-17  $\mathcal{J}\&\mathcal{F}$  score (45.1→70.1), since video object segmentation requires instance-level understanding. Finally, we further enrich our training data with LVIS [16], achieving the best overall performance.

## 5 Conclusion

We introduced SEGIC, an end-to-end in-context segmentation framework that leverages the emergent correspondence of a single frozen vision foundation model. By training on standard segmentation datasets, SEGIC achieved state-of-the-art performance on one-shot segmentation benchmarks. Impressively, SEGIC demonstrated competitive performance on novel tasks, providing a cost-effective training approach for universal segmentation. In this work, our primary focus is on utilizing one in-context example per entity. In the future, we plan to explore utilizing multiple in-context examples to enhance contextual information. Additionally, we aim to investigate our model’s potential in instance-level segmentation, such as open-world instance segmentation. We do not anticipate any undesirable ethical or social impacts.

**Acknowledgement** This project was supported by NSFC under Grant No. 62032006 and No. 62102092. We appreciate the valuable feedback from Xinlong Wang.

## References

1. Bai, Y., Geng, X., Mangalam, K., Bar, A., Yuille, A., Darrell, T., Malik, J., Efros, A.A.: Sequential modeling enables scalable learning for large vision models. arXiv preprint arXiv:2312.00785 (2023)
2. Balažević, I., Steiner, D., Parthasarathy, N., Arandjelović, R., Hénaff, O.J.: Towards in-context scene understanding. arXiv preprint arXiv:2306.01667 (2023)
3. Bao, H., Dong, L., Piao, S., Wei, F.: Beit: Bert pre-training of image transformers. In: ICLR (2021)
4. Bar, A., Gandelsman, Y., Darrell, T., Globerson, A., Efros, A.: Visual prompting via image inpainting. In: NeurIPS (2022)
5. Bay, H., Tuytelaars, T., Van Gool, L.: Surf: Speeded up robust features. In: ECCV (2006)
6. Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J.D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al.: Language models are few-shot learners. In: NeurIPS (2020)
7. Carion, N., Massa, F., Synnaeve, G., Usunier, N., Kirillov, A., Zagoruyko, S.: End-to-end object detection with transformers. In: ECCV (2020)
8. Caron, M., Touvron, H., Misra, I., Jégou, H., Mairal, J., Bojanowski, P., Joulin, A.: Emerging properties in self-supervised vision transformers. In: ICCV (2021)
9. Chen, L.C., Papandreou, G., Kokkinos, I., Murphy, K., Yuille, A.L.: Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. TPAMI (2017)
10. Chen, T., Kornblith, S., Norouzi, M., Hinton, G.: A simple framework for contrastive learning of visual representations. In: ICML (2020)
11. Cheng, B., Misra, I., Schwing, A.G., Kirillov, A., Girdhar, R.: Masked-attention mask transformer for universal image segmentation. In: CVPR (2022)
12. Cheng, B., Schwing, A., Kirillov, A.: Per-pixel classification is not all you need for semantic segmentation. In: NeurIPS (2021)
13. Cheng, H.K., Schwing, A.G.: Xmem: Long-term video object segmentation with an atkinson-shiffrin memory model. In: ECCV (2022)
14. Cho, S., Hong, S., Kim, S.: Cats++: Boosting cost aggregation with convolutions and transformers. TPAMI (2022)
15. Gu, X., Cui, Y., Huang, J., Rashwan, A., Yang, X., Zhou, X., Ghiasi, G., Kuo, W., Chen, H., Chen, L.C., et al.: Dataseg: Taming a universal multi-dataset multi-task segmentation model. arXiv preprint arXiv:2306.01736 (2023)
16. Gupta, A., Dollar, P., Girshick, R.: Lvis: A dataset for large vocabulary instance segmentation. In: CVPR (2019)
17. He, K., Chen, X., Xie, S., Li, Y., Dollár, P., Girshick, R.: Masked autoencoders are scalable vision learners. In: CVPR (2022)
18. He, K., Fan, H., Wu, Y., Xie, S., Girshick, R.: Momentum contrast for unsupervised visual representation learning. In: CVPR (2020)
19. He, K., Gkioxari, G., Dollár, P., Girshick, R.: Mask r-cnn. In: ICCV (2017)
20. Hong, S., Cho, S., Nam, J., Lin, S., Kim, S.: Cost aggregation with 4d convolutional swin transformer for few-shot segmentation. In: ECCV (2022)
21. Ilharco, G., Wortsman, M., Wightman, R., Gordon, C., Carlini, N., Taori, R., Dave, A., Shankar, V., Namkoong, H., Miller, J., Hajishirzi, H., Farhadi, A., Schmidt, L.: Openclip (2021)
22. Iqbal, E., Safarov, S., Bang, S.: Msanet: Multi-similarity and attention guidance for boosting few-shot segmentation. arXiv preprint arXiv:2206.09667 (2022)

23. Jain, J., Li, J., Chiu, M.T., Hassani, A., Orlov, N., Shi, H.: Oneformer: One transformer to rule universal image segmentation. In: CVPR (2023)
24. Jia, C., Yang, Y., Xia, Y., Chen, Y.T., Parekh, Z., Pham, H., Le, Q., Sung, Y.H., Li, Z., Duerig, T.: Scaling up visual and vision-language representation learning with noisy text supervision. In: ICML (2021)
25. Johnander, J., Danelljan, M., Brissman, E., Khan, F.S., Felsberg, M.: A generative appearance model for end-to-end video object segmentation. In: CVPR (2019)
26. Kim, S., Min, J., Cho, M.: Transformatcher: Match-to-match attention for semantic correspondence. In: CVPR (2022)
27. Kirillov, A., Mintun, E., Ravi, N., Mao, H., Rolland, C., Gustafson, L., Xiao, T., Whitehead, S., Berg, A.C., Lo, W.Y., et al.: Segment anything. arXiv preprint arXiv:2304.02643 (2023)
28. Kirillov, A., Wu, Y., He, K., Girshick, R.: Pointrend: Image segmentation as rendering. In: CVPR (2020)
29. Kossen, J., Rainforth, T., Gal, Y.: In-context learning in large language models learns label relationships but is not conventional learning. arXiv preprint arXiv:2307.12375 (2023)
30. Lee, J.Y., DeGol, J., Fragoso, V., Sinha, S.N.: Patchmatch-based neighborhood consensus for semantic correspondence. In: CVPR (2021)
31. Li, F., Zhang, H., Xu, H., Liu, S., Zhang, L., Ni, L.M., Shum, H.Y.: Mask dino: Towards a unified transformer-based framework for object detection and segmentation. In: CVPR (2023)
32. Li, H., Zhu, J., Jiang, X., Zhu, X., Li, H., Yuan, C., Wang, X., Qiao, Y., Wang, X., Wang, W., et al.: Uni-perceiver v2: A generalist model for large-scale vision and vision-language tasks. In: CVPR (2023)
33. Li, J., Selvaraju, R., Gotmare, A., Joty, S., Xiong, C., Hoi, S.C.H.: Align before fuse: Vision and language representation learning with momentum distillation. In: NeurIPS (2021)
34. Li, X., Wei, T., Chen, Y.P., Tai, Y.W., Tang, C.K.: Fss-1000: A 1000-class dataset for few-shot segmentation. In: CVPR (2020)
35. 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: ECCV (2014)
36. Lin, Z., Yang, T., Li, M., Wang, Z., Yuan, C., Jiang, W., Liu, W.: Swem: Towards real-time video object segmentation with sequential weighted expectation-maximization. In: CVPR (2022)
37. Liu, S., Zeng, Z., Ren, T., Li, F., Zhang, H., Yang, J., Li, C., Yang, J., Su, H., Zhu, J., et al.: Grounding dino: Marrying dino with grounded pre-training for open-set object detection. arXiv preprint arXiv:2303.05499 (2023)
38. Liu, Y., Zhu, M., Li, H., Chen, H., Wang, X., Shen, C.: Matcher: Segment anything with one shot using all-purpose feature matching. arXiv preprint arXiv:2305.13310 (2023)
39. Liu, Z., Mao, H., Wu, C.Y., Feichtenhofer, C., Darrell, T., Xie, S.: A convnet for the 2020s. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 11976–11986 (2022)
40. Long, J., Shelhamer, E., Darrell, T.: Fully convolutional networks for semantic segmentation. In: CVPR (2015)
41. Loshchilov, I., Hutter, F.: Decoupled weight decay regularization. arXiv preprint arXiv:1711.05101 (2017)
42. Lowe, D.G.: Distinctive image features from scale-invariant keypoints. IJCV (2004)



43. Meng, L., Dai, X., Chen, Y., Zhang, P., Chen, D., Liu, M., Wang, J., Wu, Z., Yuan, L., Jiang, Y.G.: Detection hub: Unifying object detection datasets via query adaptation on language embedding. In: CVPR (2023)
44. Meng, L., Dai, X., Yang, J., Chen, D., Chen, Y., Liu, M., Chen, Y.L., Wu, Z., Yuan, L., Jiang, Y.G.: Learning from rich semantics and coarse locations for long-tailed object detection. In: NeurIPS (2023)
45. Min, J., Kang, D., Cho, M.: Hypercorrelation squeeze for few-shot segmentation. In: CVPR (2021)
46. Min, J., Lee, J., Ponce, J., Cho, M.: Spair-71k: A large-scale benchmark for semantic correspondence. arXiv preprint arXiv:1908.10543 (2019)
47. Nguyen, K., Todorovic, S.: Feature weighting and boosting for few-shot segmentation. In: CVPR (2019)
48. Nguyen, Q.H., Vu, T.T., Tran, A.T., Nguyen, K.: Dataset diffusion: Diffusion-based synthetic data generation for pixel-level semantic segmentation. In: NeurIPS (2023)
49. Oquab, M., Darcet, T., Moutakanni, T., Vo, H., Szafraniec, M., Khalidov, V., Fernandez, P., Haziza, D., Massa, F., El-Nouby, A., et al.: Dinov2: Learning robust visual features without supervision. arXiv preprint arXiv:2304.07193 (2023)
50. Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C., Mishkin, P., Zhang, C., Agarwal, S., Slama, K., Ray, A., et al.: Training language models to follow instructions with human feedback. In: NeurIPS (2022)
51. Pont-Tuset, J., Perazzi, F., Caelles, S., Arbeláez, P., Sorkine-Hornung, A., Van Gool, L.: The 2017 davis challenge on video object segmentation. arXiv preprint arXiv:1704.00675 (2017)
52. Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al.: Learning transferable visual models from natural language supervision. In: ICML (2021)
53. Radford, A., Narasimhan, K., Salimans, T., Sutskever, I., et al.: Improving language understanding by generative pre-training. OpenAI blog (2018)
54. Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., Sutskever, I., et al.: Language models are unsupervised multitask learners. OpenAI blog (2019)
55. Ren, S., He, K., Girshick, R., Sun, J.: Faster r-cnn: Towards real-time object detection with region proposal networks. In: NeurIPS (2015)
56. Rombach, R., Blattmann, A., Lorenz, D., Esser, P., Ommer, B.: High-resolution image synthesis with latent diffusion models. In: CVPR (2022)
57. Sudre, C.H., Li, W., Vercauteren, T., Ourselin, S., Jorge Cardoso, M.: Generalised dice overlap as a deep learning loss function for highly unbalanced segmentations. In: Deep Learning in Medical Image Analysis and Multimodal Learning for Clinical Decision Support: Third International Workshop, DLMIA 2017, and 7th International Workshop, ML-CDS 2017, Held in Conjunction with MICCAI 2017, Québec City, QC, Canada, September 14, Proceedings 3 (2017)
58. Tancik, M., Srinivasan, P., Mildenhall, B., Fridovich-Keil, S., Raghavan, N., Singhal, U., Ramamoorthi, R., Barron, J., Ng, R.: Fourier features let networks learn high frequency functions in low dimensional domains. In: NeurIPS (2020)
59. Tang, L., Jia, M., Wang, Q., Phoo, C.P., Hariharan, B.: Emergent correspondence from image diffusion. arXiv preprint arXiv:2306.03881 (2023)
60. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, Ł., Polosukhin, I.: Attention is all you need. In: NeurIPS (2017)
61. Wang, X., Jabri, A., Efros, A.A.: Learning correspondence from the cycle-consistency of time. In: CVPR. pp. 2566–2576 (2019)
62. Wang, X., Wang, W., Cao, Y., Shen, C., Huang, T.: Images speak in images: A generalist painter for in-context visual learning. In: painter (2023)

63. Wang, X., Zhang, X., Cao, Y., Wang, W., Shen, C., Huang, T.: Seggpt: Towards segmenting everything in context. In: CVPR (2023)
64. Xiao, T., Liu, Y., Zhou, B., Jiang, Y., Sun, J.: Unified perceptual parsing for scene understanding. In: ECCV (2018)
65. Xie, S.M., Raghunathan, A., Liang, P., Ma, T.: An explanation of in-context learning as implicit bayesian inference. In: ICLR (2021)
66. Xu, J., Liu, S., Vahdat, A., Byeon, W., Wang, X., De Mello, S.: Open-vocabulary panoptic segmentation with text-to-image diffusion models. In: CVPR (2023)
67. Xu, N., Yang, L., Fan, Y., Yue, D., Liang, Y., Yang, J., Huang, T.: Youtube-vos: A large-scale video object segmentation benchmark. arXiv preprint arXiv:1809.03327 (2018)
68. Yan, B., Jiang, Y., Wu, J., Wang, D., Luo, P., Yuan, Z., Lu, H.: Universal instance perception as object discovery and retrieval. In: CVPR. pp. 15325–15336 (2023)
69. Yu, J., Wang, Z., Vasudevan, V., Yeung, L., Seyedhosseini, M., Wu, Y.: Coca: Contrastive captioners are image-text foundation models. arXiv preprint arXiv:2205.01917 (2022)
70. Zhang, J.W., Sun, Y., Yang, Y., Chen, W.: Feature-proxy transformer for few-shot segmentation. *Advances in Neural Information Processing Systems* (2022)
71. Zhang, R., Jiang, Z., Guo, Z., Yan, S., Pan, J., Dong, H., Gao, P., Li, H.: Personalize segment anything model with one shot. arXiv preprint arXiv:2305.03048 (2023)
72. Zhou, B., Zhao, H., Puig, X., Xiao, T., Fidler, S., Barriuso, A., Torralla, A.: Semantic understanding of scenes through the ade20k dataset. *IJCV* (2019)
73. Zhou, X., Koltun, V., Krähenbühl, P.: Simple multi-dataset detection. In: CVPR (2022)
74. Zou, X., Dou, Z.Y., Yang, J., Gan, Z., Li, L., Li, C., Dai, X., Behl, H., Wang, J., Yuan, L., et al.: Generalized decoding for pixel, image, and language. In: CVPR (2023)
75. Zou, X., Yang, J., Zhang, H., Li, F., Li, L., Gao, J., Lee, Y.J.: Segment everything everywhere all at once. In: NeurIPS (2023)

## A Additional Implementation Details

In this section, we provide additional details of the mask decoding and training pipeline.

**Mask decoding.** Inspired by recent query-based segmentation methods [7, 11, 12, 27], we utilize a lightweight query-based mask decoder to effectively map the in-context enhanced image features and object query to an output mask. We employ a learnable object query that will be used for the decoder’s output. The workflow of SEGIC is concisely summarized in Pytorch-style pseudocode, as presented in Algorithm 1. Initially, the image feature and in-context instructions are projected into the same feature space for mask decoding. Subsequently, the projected in-context features are leveraged to enhance the image feature and object feature. The mask decoder then executes a multi-layered decoding process, structured as a four-step procedure within each layer, including (1) self-attention for the concatenated object query; (2) cross-attention from the image feature to the object query; (3) cross-attention from the query back to the image; and (4) calculating the mask. The mask calculation is performed using the image feature and the first element of the concatenated object feature, which corresponds to the position of the initial object query.

**Training pipeline.** In our main experiments, we adopt a mixed training scheme using both semantic segmentation datasets (COCO and ADE20k) and instance segmentation datasets (COCO and LVIS), as presented in Algorithm 2. We do not focus intensively on adjusting the dataset ratios, instead opting for uniform dataset-level sampling. For the segmentation datasets, we employ large-scale jittering (ranging from 0.1 to 2.0) for both the target image and in-context examples. These in-context examples are constructed based on the semantic class label of the target image, sampling one class per image during training. In the case of instance segmentation datasets, in-context examples are generated by applying two separate data augmentations to the same image. The instances from these differently augmented views then serve as mutual in-context examples. Our standard data augmentation techniques for this task include random resizing cropping (ranging from 0.3 to 1.0), random color jittering (with a 0.2 probability), and random horizontal flipping (with a 0.1 probability).

## B Additional Visualization

In this section, we provide more visualizations of the middle output and the predictions of SEGIC.

**Propagated mask.** As outlined in Section 3.2, the propagated mask  $\mathbf{a}$  is derived from a weighted mean of dense correspondences according to the ground-truth mask of in-context samples. To facilitate visualization, we first apply the sigmoid function to map  $\mathbf{a}$  into  $(0, 1)$ . Subsequently, this range is transformed into RGB space using the JET color map. As depicted in Figure 4, this process demonstrates that the propagated masks predominantly concentrate on the objects referenced in the in-context examples, providing strong guidance for the

subsequent mask decoding process. This observation further demonstrates the emerging potential of pre-trained vision foundation models in the realm of segmentation tasks.



**Fig. 4: Visualization of propagated masks.** We propagate labels from the in-context examples to the targets to obtain propagated masks by exploring the dense correspondences. We employ DINO-v2-large [49] for the visualization.

**More qualitative results on VOS.** We further provide more qualitative results on video object segmentation tasks (VOS). Note that SEGIC is never trained on video datasets and just treats VOS as in-context segmentation using the first frame as in-context examples. As shown in Figure 5, SEGIC well handles challenging scenarios, including (a) occlusions, (b) interwoven objects, and (c) small objects.

---

**Algorithm 1:** Pseudo code for SEGIC Mask Decoding.
 

---

```

# Inputs: Image Embedding f; In-context Instructions c = {a, p, v, m}
# Variables: Learnable Object Queries q
# Functions: Conv4ImgFeature(), Conv4ProgatedLabel(); Proj4Pos(), Proj4Vis(),
Porj4Meta(); QuerySelfAttn(), Query2ImgAttn(), Image2QueryAttn(), output()
1 def InContext_Enhancement(f, a, p, v, m):
    # Project image feature and in-context propagated mask into the hidden space for
    # mask decoding.
2     f', a' = Conv4ImgFeature(f), Conv4ProgatedLabel(a);
    # Enhance image feature with in-context propagated mask.
3     f' = f' + a';
    # Project other in-context instructions into the hidden space for mask decoding.
4     qp, qv, qm = Proj4Pos(p), Proj4Vis(v), Proj4Meta(m)
    # Enhance the object query by contacting with the hidden features of in-context
    # instructions.
5     q' = Concat(q, qp, qv, qm)

6 def Mask_Decoder(F, Q):
7     Q' = QuerySelfAttn(Q) # Query self-attention
8     Qo = Img2QueryAttn(Q', F) # Image-to-query cross-attention
9     Fo = Query2ImgAttn(F, Qo) # Query-to-image cross-attention
10    O = output(Fo, Qo[0]) # Compute mask

11 def forward(f, a, p, v, m):
12    f', q' = InContext_Enhancement(f, a, p, v, m) # Enhance image feature and object
    # query with in-context instructions.
13    Qo, Fo = q', f' # Initialize variables for mask decoding
14    for i in range(max_iter):
15        O, Qo, Fo = Mask_Decoder(Qo, Fo)
    
```

---



---

**Algorithm 2:** Pseudo code for training pipeline.
 

---

```

# training set: mixed dataset D = Dinst ∪ Dsem
1 def ICL_Preprocess(data):
2     It, yt, category = data
3     if task_type(data) == 'semantic':
4         Ir, yr = CategoryAwareSample(Dsem, category)
5         meta = 'a photo of a {category}.'
6     if task_type(data) == 'instance':
7         Ir, yr = DataAug(It, yt) # Individual data aug to build a different view as
            # reference
8         meta = 'please segment the instances.'

9 def train_epoch(model, D):
10    for data in D:
11        It, yt, Ir, yr, meta = ICL_Preprocess(data)
12        loss = model(It, yt, Ir, yr, meta)
    
```

---



**Fig. 5: Qualitative results on VOS.** SEGIC perform well on challenging scenarios in video object segmentation, including (a) occlusions, (b) interwoven objects, and (c) small objects.