# A Survey on Segment Anything Model (SAM): Vision Foundation Model Meets Prompt Engineering

**Chaoning Zhang***
Kyung Hee University

**Fachrina Dewi Puspitasari**
KAIST

**Sheng Zheng**
Beijing Institute of Technology

**Chenghao Li**
KAIST

**Yu Qiao**
Kyung Hee University

**Taegoo Kang**
Kyung Hee University

**Xinru Shan**
Microsoft STCA

**Chenshuang Zhang**
KAIST

**Caiyan Qin**
Harbin Institute of Technology

**Francois Rameau**
State University of New York at Korea

**Lik-Hang Lee**
Hong Kong Polytechnic University

**Sung-Ho Bae**
Kyung Hee University

**Choong Seon Hong**
Kyung Hee University

July 4, 2023

## Abstract

Segment anything model (SAM) developed by Meta AI Research has recently attracted significant attention. Trained on a large segmentation dataset of over 1 billion masks, SAM is capable of segmenting any object on a certain image. In the original SAM work, the authors turned to zero-short transfer tasks (like edge detection) for evaluating the performance of SAM. Recently, numerous works have attempted to investigate the performance of SAM in various scenarios to recognize and segment objects. Moreover, numerous projects have emerged to show the versatility of SAM as a foundation model by combining it with other models, like Grounding DINO, Stable Diffusion, ChatGPT, etc. With the relevant papers and projects increasing exponentially, it is challenging for the readers to catch up with the development of SAM. To this end, this work conducts the first yet comprehensive survey on SAM. This is an ongoing project, and we intend to update the manuscript on a regular basis. Therefore, readers are welcome to contact us if they complete new works related to SAM so that we can include them in our next version.

## 1 Introduction

ChatGPT Zhang et al. [2023a] has revolutionized our perceptions of AI, gaining significant attention and interest across the world. It marks a breakthrough in generative AI (AIGC, a.k.a Artificial intelligence generated content) Zhang et al. [2023b], for which foundation models Bommasani et al. [2021] have played a significant role. The large language model has achieved significant performance in language tasks, leading to a new paradigm in various NLP areas. In the vision field, multiple works Radford et al. [2021], Jia et al. [2021], Yuan et al. [2021] have attempted to learn an image encoder together with a text encoder with contrastive learning He et al. [2020], Qiao et al. [2023a], Zhang et al. [2022a]. The resulting image encoder can be perceived as a vision foundation model. Another form of training a vision foundation model is through self-supervised learning, like masked autoencoder Zhang et al. [2022b]. However, such vision foundation models often require finetuning before they can be used for downstream tasks.

Very recently, Meta Research team has released the "Segment Anything" project Kirillov et al. [2023], where a model termed Segment Anything Model (SAM) has been proposed. An overall view of the "Segment Anything" project is

---

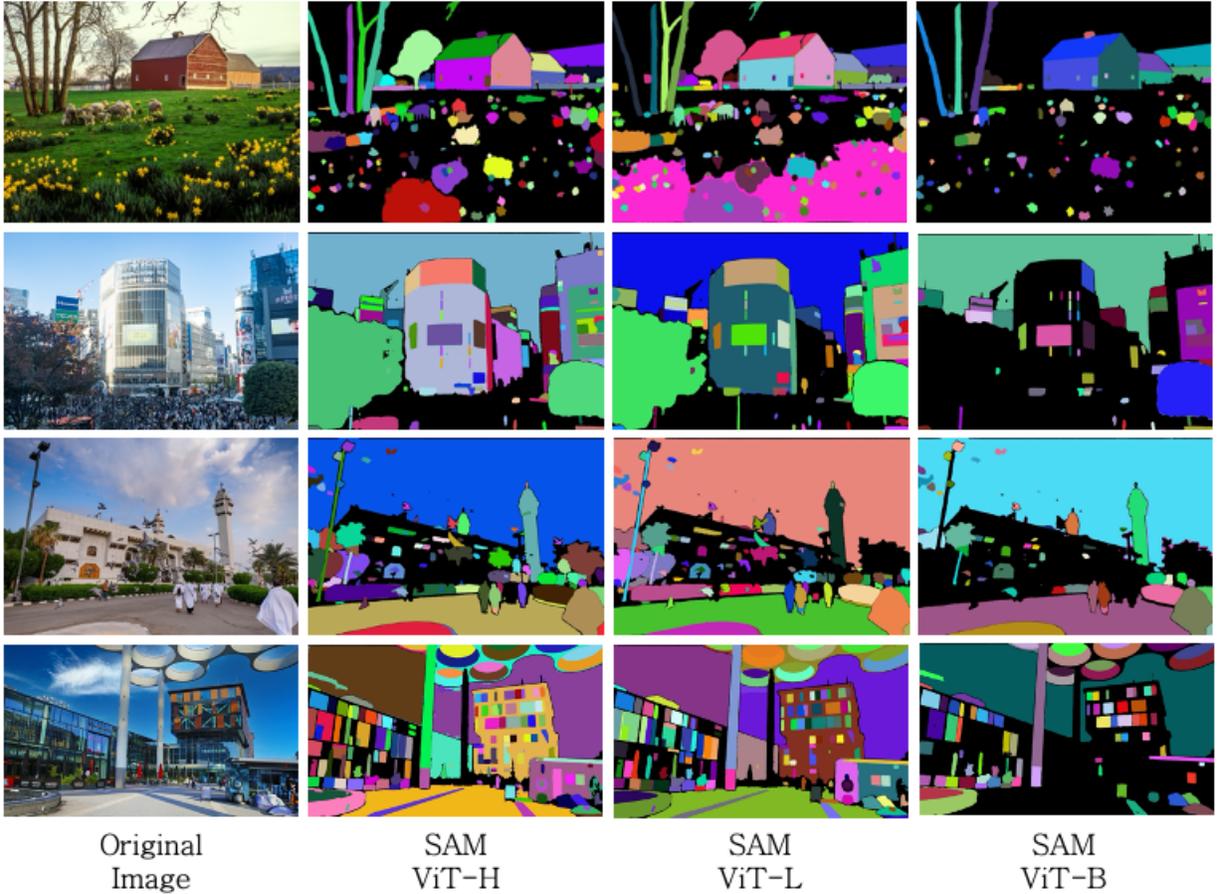*You are welcome to contact us through chaoningzhang1990@gmail.com

Figure 1: Segment Anything results depending on the model size.



(a) **Task:** promptable segmentation      (b) **Model:** Segment Anything Model (**SAM**)      (c) **Data:** data engine (top) & dataset (bottom)
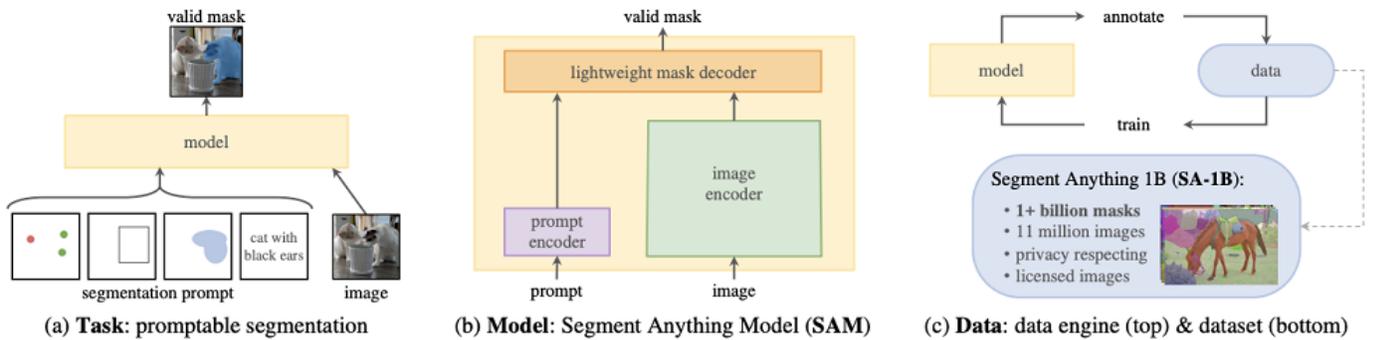
Figure 2: Segment anything project (figure obtained from Kirillov et al. [2023]).

shown in Figure 1. It is worth mentioning that SAM performs promptable segmentation, which differs from semantic segmentation in two ways: (1) the generated masks of SAM have no labels; (2) SAM relies on prompts. In other words, SAM only cut-outs the objects in the image without assigning labels (see Figure 1), and which object get cutout depends on the given prompt. Given the so-called prompt engineering, SAM has shown remarkable zero-shot transfer performance without the need for finetuning, making many believe that SAM is like GPT-3 Brown et al. [2020] moment for computer vision. SAM is trained on SA-1B, which contains more than 1B masks from 11M images making it the largest segmentation dataset ever released.

**Label prediction *v.s.* mask prediction.** Conceptually, semantic segmentation can be perceived as a combination of mask prediction and label prediction. The success of the "Segment Anything" project shows that these two sub-tasks

can be decoupled and SAM exclusively solves the first one. Without label prediction as existing image segmentation tasks (like instance segmentation and panoptic segmentation), the task that SAM solves might seem to be a trivial task at first sight. However, it actually solves a fundamental task in computer vision that contributes to the development of the vision foundation model. For maximizing the generalization to unseen distribution, a vision foundation model needs to be trained with a sufficiently large yet diverse dataset. When the dataset size and diversity increase, the object category and label have open-vocabulary nature, which makes it impossible to pre-determine a fixed list of labels beforehand.

**Vision foundation model meets promptable segmentation.** To overcome the above issue, SAM Kirillov et al. [2023] opts for the task of prompt-based mask prediction (a.k.a. promptable segmentation), where the role of prompt behaves like attention. When the human eye understands the world, it often focuses on a certain object while perceiving its surrounding area as a background. Given numerous objects in the same scene, without an attention mechanism, the human eye cannot make sense of them. Moreover, the human eye can identify and segment the object of interest even though the observer has never seen a similar object in the past. For example, a baby who first sees a samoyed dog will track the movement of the dog even though it never understands what a dog is. In other words, vision understanding mainly relies on the object mask instead of its corresponding label. Overall, the task of promptable segmentation well mimics how the human eye understands the world. The SAM trained on the promptable segmentation constitutes a vision foundation model that is not only generalizable to unseen distributions, but also compatible with other models for realizing more demanding tasks.
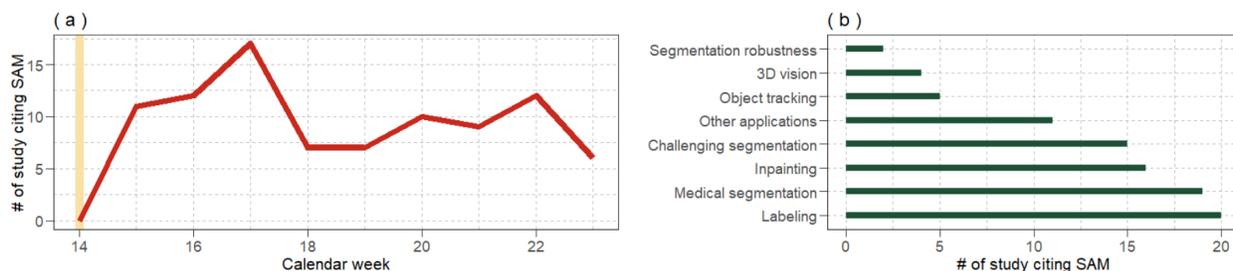


Figure 3: Volume of research work citing SAM (a) since its publication (indicated by the yellow vertical line). Number of studies in each topic related to the implementation of SAM (b).

Figure 3(a) shows that within a few weeks, numerous studies have been conducted in the computer vision community to investigate SAM from various perspectives (Figure 3(b)). Within this study, the two largest topics being explored are the integration of SAM into medical image segmentation, and the labeling of SAM's class-agnostic segmentation mask. Overall, we categorize our discussion of SAM-related study into two major topics: (1) *evaluation of SAM's segmentation performance and the endeavors to improve its performance* and (2) *review of how SAM is integrated with other foundation models in various vision- and non-vision-related tasks*. Given the increasing volume of work, it can be overwhelming for readers to catch up with the development of SAM. To this end, this work conducts a survey on SAM in the era of vision embarking on the path of Natural Language Processing (NLP) to embrace the foundation model.

## 2 Can SAM really segment anything in all scenarios?

As the title suggests, SAM Kirillov et al. [2023] is claimed to segment anything in the images. However, it remains unclear whether the SAM model can work well in real-world scenarios. Therefore, numerous works have been conducted recently to evaluate its performance in various scenarios, including medical images and beyond.

**Medical images.** Segmenting medical images to dissect abnormal tissue from others is often challenging because it requires expertise in the related pathology field. For this reason, current approaches in medical image segmentation still rely on data annotated by experts, which is then trained on various deep learning models. To alleviate the exhaustive resources needed for such training, practitioners seek to exploit the zero-shot segmentation capability of SAM. Nevertheless, raw segmentation outputs of SAM are found to be inferior to the accuracy (Dice score) of fully-supervised models, *e.g.*, U-Net He et al. [2023a]. Moreover, different organs produce different score discrepancies. For instance, the two largest gaps (70%) of accuracy between SAM and U-Net are observed in the segmentation for pancreas He et al. [2023a] and liver Hu and Li [2023] while the two smallest gaps (30%) are found from polyps Zhou et al. [2023a] and lung nodules He et al. [2023a]. This underperformance is potentially caused by the different nature of the object being segmented. Objects with more apparent boundaries produce higher segmentation accuracy than those with obscure lining Huang et al. [2023a]. For example, SAM can segment benign tumors easier than malignant tumors because the latter does not have a clear boundary as the effect of metastasis to the neighboring tissue Hu et al.

[2023a]. Nevertheless, despite its limitation against fully-supervised models, SAM segmentation still outperforms the non-deep learning segmentation methods, such as FSL BET Mohapatra et al. [2023]. SAM's performance is dependent on the types and magnitude of prompts fed to the prompt encoder. In general, automatic prompting yields unsatisfying segmentation results Huang et al. [2023a]. Meanwhile, box prompts produce higher segmentation accuracy compared to point prompts Cheng et al. [2023a], Wang et al. [2023a], Wald et al. [2023], Mattjie et al. [2023]. However, accuracy from point prompts can be increased by applying more points on the target object Cheng et al. [2023a], Hu and Li [2023], Mazurowski et al. [2023], Wald et al. [2023], Mattjie et al. [2023], Huang et al. [2023a]. Additionally, box prompts and point prompts can be combined to produce better accuracy, but it's not when it is applied simultaneously Huang et al. [2023a]. Improvement is apparent when the box prompt is applied in the initial prompting stage, whereas the point prompt is administered during the mask refinement stage Mattjie et al. [2023]. One of the causes of SAM's flaw in medical image segmentation is the imbalance proportion of medical image data that exists among SAM's training dataset because the domain-specific dataset is extremely scarce in open-world object segmentation tasks. For this reason, fine-tuning with domain-specific datasets can be a quick-fix solution for SAM's poor accuracy in medical image segmentation. To this end, such fine-tuning approaches have been evaluated on polyp colonoscopy Li et al. [2023a], skin lesions Hu et al. [2023b], and other medical image datasets Ma and Wang [2023]. On average, this strategy yields segmentation accuracy (Dice score) above 80%. More advanced accuracy improvement approaches are proposed by slightly modifying SAM's framework. Most of these modifications are performed by attaching domain-specific adapters, which role is to learn task-specific knowledge. Such an adapter can be fixed in between transformer layers of image encoder Qiu et al. [2023], Zhang and Liu [2023], Wu et al. [2023], between attention layers of mask decoder Wu et al. [2023], or even completely replacing both prompt encoder and mask decoder with task-specific head Qiu et al. [2023]. Another modification is proposed by decoupling the mask decoder into two modules which are responsible for handling IoU regression and mask learning respectively Gao et al. [2023]. Such a mechanism attempts to solve the coupling effect between image embedding and prompt token in the mask decoder that makes SAM segmentation output highly dependent on the prompt quality.
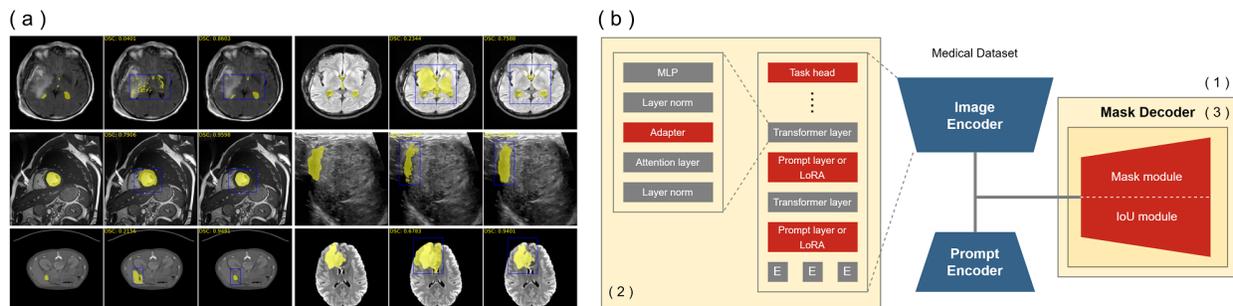


Figure 4: (a) From left to right, segmentation mask of ground truth, pre-trained SAM, and fine-tuned SAM Ma and Wang [2023]. (b) Best practices for improving SAM's segmentation in medical images; (1) fine-tuning mask decoder with a medical dataset, (2) fixing prompt layer or LoRA in between transformer layer or inside the transformer layer itself, and (3) decoupling mask decoder into two modules. (Blue: frozen modules, red: fine-tuned modules)

**Other scenarios.** Beyond medical images, many real-world segmentation tasks are exposed to challenging conditions that weaken the segmentation capability of SAM. Minuscule and slender objects Ji et al. [2023a,b], Ren et al. [2023], objects with obscure boundary Jie and Zhang [2023], Ji et al. [2023a,b], occluded objects in dense environments Yang et al. [2023a], camouflaged objects Tang et al. [2023], Ji et al. [2023a,b], and transparent objects Han et al. [2023], Ji et al. [2023b] are a few examples of which SAM segmentation outputs are rather inaccurate. For this reason, the raw output map of SAM cannot be used directly in object counting task Ma et al. [2023]. Hence, researchers introduce various approaches to overcome these challenges. The most straightforward approach is made by refining the input or output of SAM. For instance, an input image can be refined using image-level or pixel-level entropy filtering Guo et al. [2023], He et al. [2023b] or multi-augmentation He et al. [2023b], and prompts can be supplemented with additional high-quality tokens Ke et al. [2023]. Such practices are useful to correctly locate anchors to precisely prompt the target object. Additionally, instead of direct utilization, SAM's output masks can be refined, *e.g.*, by applying a mask filter or correcting the boundary Giannakis et al. [2023], Williams et al. [2023]. Further, a few approaches also attempt to fix adapters into SAM to enable task-specific learning Chen et al. [2023a], Julka and Granitzer [2023], Cao et al. [2023]. SAM also offers promising results in the semantic communication domain compared to traditional improvement approaches that are spectrum-limited and need high power consumption. Capitalizing on SAM's generalizability and promptability, Tariq et al. [2023] leverage SAM in their semantic communication framework to transmit only selected information at high broadcasting quality.
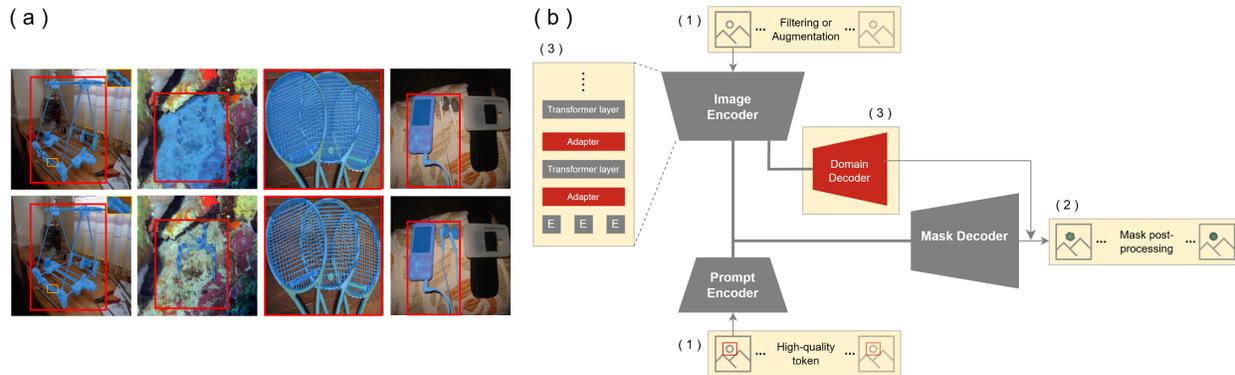
4

Figure 5: (a) From top to bottom, segmentation mask of pre-trained SAM and fine-tuned SAM Ke et al. [2023]. (b) Best practices for improving SAM's segmentation on challenging objects; (1) pre-processing of image and token, (2) post-processing mask output, and (3) fixing adapter in between transformer layer and fixing domain decoder.

**On the robustness of SAM.** Despite its generality to open-world imaging scenarios, SAM is still vulnerable to some image perturbation. Very recently, Zhang et al. [2023c] has investigated adversarial attacks on SAM. It shows that the SAM model is not robust against the attack of adversarial examples Szegedy et al. [2013], Goodfellow et al. [2015], Kurakin et al. [2017]. Specifically, the basic goal attack-SAM Zhang et al. [2023c] is set to remove the masks, for which the authors design a CLIP-MSE loss. The results show that the model is very vulnerable under white-box PGD attack Madry et al. [2018]. In the black-box setting, the model maintains robustness to some extent. On top of removing the original masks, the authors also experiment with generating new (target) masks, showing intriguing results. For details, interested readers can check Zhang et al. [2023c] for more details. Additionally, Huang et al. [2023b], Wang et al. [2023b] also examines SAM's robustness against fifteen image corruption of varying severity. They found that all corruptions, except blur-related ones, only slightly decreased (<5%) pixel accuracy and intersection over union metrics in various datasets. Additionally, Huang et al. [2023b] also examines SAM's robustness against fifteen image corruption of varying severity. They found that all corruptions, except blur-related ones, only slightly decreased (<5%) pixel accuracy and intersection over union metrics in various datasets. Nevertheless, Wang et al. [2023b] discovered that the decline in the above metric values in the perturbed image is larger in the segmentation task involving challenging objects, such as medical X-ray. Another recent work Qiao et al. [2023b] performs comprehensive evaluations on the robustness of SAM on corruption and beyond. Specifically, it interprets various corruptions as new styles and tests the SAM robustness against style transfer and common corruptions. Moreover, it investigates the SAM's robustness against local occlusion and adversarial perturbation. The results demonstrate that SAM has a moderate level of resilience against FGSM attacks, but not PGD attacks, even for perturbation with a very small magnitude Qiao et al. [2023b].

## 3  From *Segment Anything* to *X-anything*

The success of SAM for "*segment anything*" has motivated the community to investigate *X anything*. Specifically, SAM has been shown to be versatile in numerous projects when combined with other models to achieve impressive performance.

### 3.1  Label Anything

As the outcome of being trained on a massive dataset, SAM has a generalizable nature that also makes it class-agnostic because assigning a unique label to each segmented object is almost impossible. Nevertheless, labeling is one of the critical requirements in many computer vision tasks, *e.g.*, object detection and object classification. Thus, numerous studies have been initiated to exploit SAM's powerful segmentation while enabling it to produce a label for each mask.

*Labeling.* One of the earliest approaches in labeling SAM's segmentation mask is to combine Large-Language Model (LLM), Large-Vision Model(LVM), and Vision-Language Model (VLM) so they can complement each other's strengths Yu et al. [2023a] (Figure 6). This approach uses LLM, such as ChatGPT as input for AIGC models to generate images. This image is then passed through labeling processes by LVM. For instance, these LVM can be Bootstrapping Language-Image Pretraining (BLIP) Li et al. [2022], which converts the image to text, Grounding-DINO, which converts the text to visual prompt, and SAM which segments the image given the prompt. This process is done iteratively until all target objects can be labeled. This complementary approach can also be implemented in the 3D

scene understanding task Chen et al. [2023b], though some only at the semantic mask construction module of the entire 3D scene labeling pipeline Chen and Li [2023]. Another mutual association for labeling SAM's mask is also possible by combining SAM, which has strong boundary definition but lacks labels, with Class Activation Mapping (CAM) Zhou et al. [2015], which excels in semantic labeling but is poor at defining boundary Chen et al. [2023c]. The transfer labeling techniques open the possibility for SAM application in weakly-supervised segmentation tasks Sun et al. [2023a] and open-world object detection tasks He et al. [2023c]. Apparently, SAM-generated pseudo-label can exceed the accuracies of the earlier pseudo-labeling methods, with bounding box prompting giving the highest accuracy Jiang and Yang [2023]. Another simple approach to labeling SAM's mask is through voting between label-free SAM's mask and label-aware mask generated by Open-Vocabulary Semantic Segmentation (OVSeg) Liang et al. [2023] model. Depth mapping prior to SAM segmentation can also be applied in this approach to give SAM's mask richer geometric information that will be helpful during voting Cen et al. [2023a].

*Feature Matching.* Despite its class-free masks, SAM can be integrated into the semantic feature matching framework for it can automatically segment objects of the same class in the new image given the semantic feature of such objects in the input image. The key idea is to match the semantic features of two images and use the matching keypoints as a prompt for SAM to segment objects of the same class in the target image Liu et al. [2023a]. Further, SAM's mask outputs of this approach can be refined to alleviate their ambiguity by assigning weights to each mask Zhang et al. [2023d]. In 3D application, this semantic feature matching can also be used to find correspondence among similar shape objects, *e.g.*, human limbs and animal limbs Abdelreheem et al. [2023].

*Captioning.* The above approaches to label SAM's mask are extendable to the task of image captioning. Slightly different from labeling, captioning aims to give a description to an image based on the content understanding and the choice of user language style. For instance, a caption can be in the form of a short sentence or a whole paragraph. This task borrows the versatility of VLM to translate images into text and refine this text into a comprehensive paragraph using LLM. Here, SAM acts as the segmenter to select the specific regions where the users want their caption to emphasize Wang et al. [2023c]. Oftentimes, the captioning task can be complicated considering the accordance between word composition and spatial composition of the image. Thus, it is important to assess the caption density that evaluated these two aspects. In this task, SAM does not only function as a segmented but also as the informer of spatial composition of the image Doveh et al. [2023].

*Data Annotation.* Owing to SAM's capability to label anything as explained above, researchers further employ SAM in the automatic data annotation pipeline. This automation is helpful in computer vision task that deals with target object where the number of annotated data is sparse, *e.g.*, domain-specific application Wang et al. [2023d] such as medical image Zhang et al. [2023e], remote sensing Wang et al. [2023e], Yu et al. [2023b], Zhang et al. [2023f] and autonomous driving Zhou et al. [2023b] and specific imaging typeChen and Bai [2023].
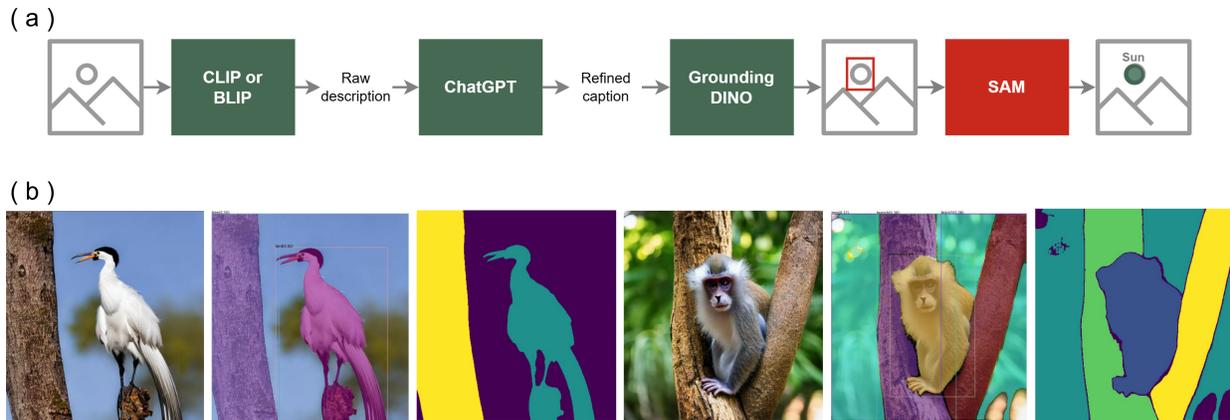


Figure 6: (a) General labeling pipeline for SAM's segmentation. (b) From left to right, original picture, label, and segmentation mask Yu et al. [2023a].

## 3.2 Inpaint Anything

Inpainting is the process of editing a part of an image, either by removing an existing object or adding a new instance to the frame. This process can be executed by traditional or deep-learning approaches. Particularly in the deep-learning

approach, selecting a target region in a frame is almost always necessary. Thus, this role is what SAM is expected to be capable of executing.

***2D Inpainting.*** To this end, numerous study that implements SAM in the image inpainting framework has been performed. The core idea of attaching SAM in their pipelines is to exploit SAM's class-agnostic segmentation ability on an input image which resulting masks are further processed by AIGC framework, *e.g.*, Stable Diffusion (SD) Rombach et al. [2022], or other deep-learning based inpainting frameworks (Figure 7). For instance, SAM mask outputs can be further nominated by the textual prompt to be processed by SD with the guidance of vision-language pertaining model, *e.g.*, Contrastive Language–Image Pre-training (CLIP) Radford et al. [2021], to generate final output as specified in the text instruction Xie et al. [2023], Levin and Fried [2023], Yu et al. [2023c], Liu et al. [2023b]. Here, the SAM segmentation task can be prompted by either visual prompts or textual prompts. However, since SAM reliability in handling text prompts is still weaker compared to visual prompts, researchers attempted to attach an intermediate module between prompt and SAM Wang et al. [2023f], Fang et al. [2023]. Grounded-SAM IDEA-Research [2023], which combines Grounding-DINO Liu et al. [2023c] and SAM, is one of the renowned solutions for such limitation. Other simpler applications feed SAM mask outputs into the autoencoder architecture with attention to producing customizable style transfer Liu et al. [2023d] or controlling the intensity of change at the selected region Jiang and Holz [2023]. SAM's agnostic nature to the fine-grained pixel information brings an advantage to the inpainting process in which input image resolution is low. Such insensitivity enables SAM to extract only global features that are useful when guiding image restoration tasks whose end goal is to enhance the image resolution. For this reason, SAM's mask outputs can be valuable prior to tuning the autoencoder network in the image/video restoration pipeline Xiao et al. [2023], Lu et al. [2023].

***3D Inpainting.*** Although SAM is originally developed to handle 2D image, study shows that it can be integrated into 3D space. The key idea is to dissect 3D scene into multiple views projected onto 2D frames, and vice versa through Neural Radiance Field (NeRF) Mildenhall et al. [2021]. SAM's role in this 3D inpainting pipeline is to generate the target object's masks of multiple projected views to ensure the 3D rendering consistency, exploiting SAM's fine-grained agnostic nature Wang et al. [2023g]. Alternatively, SAM's manual segmentation can only be executed only at one single-view image while images of different views can be automatically prompted based on the 3D rendering result of the prior image Cen et al. [2023b]. Utilizing a similar segmentation idea, it is also possible to generate 3D objects only from a single view image by generating multiple view frames using AIGC framework Shen et al. [2023]. This development in 3D inpainting techniques set the starting point for its use in many applications, *e.g.*, robotic vision Lillrank et al. [2023]. Inpainting also allows users to click on an object to remove anything and then fill in anything or replace anything with a text prompt. In particular, the inpainting process can be employed in the pipeline of the diminished reality Mori et al. [2017], which is a critical step in various immersive environments of the Metaverse, e.g., removing unwanted objects before adding digital overlays on top of the physical surroundings Lee et al. [2023]. IA shows great advantages in simple operation in inpainting, making it highly user-friendly for application Open-vocabulary-Segment-Anything ngthanhtin [2023] combines OWL-ViT with Segment Anything, prompt by text and conduct inpainting via stable diffusion, which achieves a text interaction in inpainting.

***Matting*** . To properly execute the inpainting task, it is important to separate foreground from background accurately, which process is known as matting. Different from segmentation, the matting process is more intricate as it has to define an alpha-matte value that defines the opacity of the pixel in the form of trimap. For this reason, generating an alpha-matte mask is more resource-expensive than a normal segmentation mask. Thus, to alleviate this constraint, researchers employ SAM's segmentation power to generate accurate masks that will function as pseudo-trimap Yao et al. [2023]. Nevertheless, foreground-background separation can also be achieved through simple approaches, *e.g.*, enhancing the pixel quality of the foreground against a background of a target object choice. In this framework, SAM can be integrated to generate foreground segmentation mask Yang et al. [2023b].

### 3.3 Track Anything

Training an object tracking algorithm, particularly those that work with tracking-by-detection, requires ground-truth annotation work that is labor-intensive. This is due to the required accurate annotation at every frame to avoid the temporal accumulation of annotation errors. For this reason, the commonly desired approach in object tracking is to annotate using a bounding box due to its simplicity, *i.e.*, only demanding two corner points to draw. Meanwhile, more accurate annotations, *e.g.*, polygon segmentation or keypoints, are more desirable when handling inherent challenges in object detection, *e.g.*, occlusion and similarly-looking objects. Nevertheless, these annotations are more expensive to conduct due to the number of points that needs to be drawn for a single object.

***Tracking.*** To relieve this demanding annotation effort, researchers have explored the integration of SAM in object tracking tasks for its excellency in providing accurate segmentation. Nevertheless, given its original assignment in a
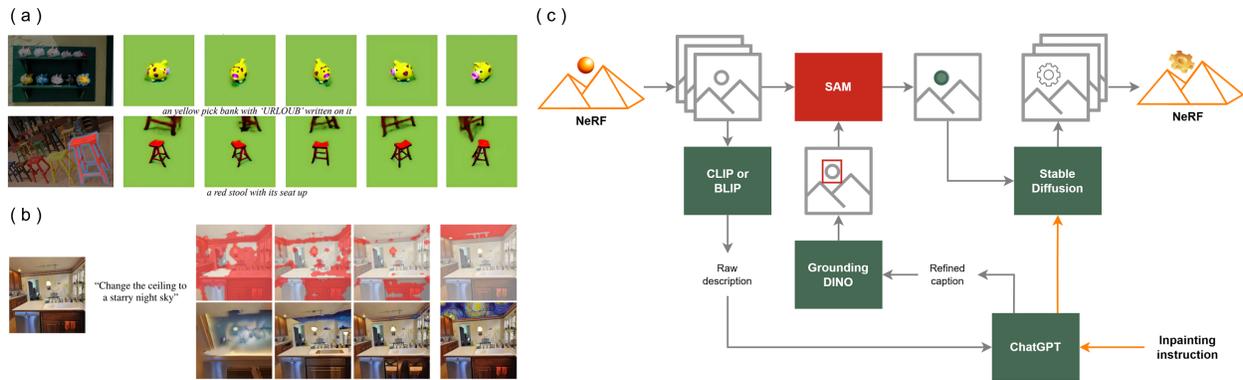
Figure 7: (a) 3D object reconstruction in NeRF from a single-view image Shen et al. [2023]. (b) 2D image inpainting given textual instruction, right-most image is generated using SAM-integrated pipeline Wang et al. [2023f]. (c) General pipeline for SAM-integrated inpainting.

single 2D image, SAM is potentially less versatile in handling frame changes in the temporal dimension He et al. [2023d]. To address this limitation, there have been few studies that integrate SAM with former video object segmentation (VOS) models, *e.g.*, XMem Cheng and Schwing [2022] or Decoupling Features in Hierarchical Propagation (DeAOT) Yang and Yang [2022] (Figure 8). Such combination gives advantages to each other, *i.e.*, VOS suppresses SAM temporal inconsistency while SAM refines VOS segmentation accuracy. Thus, this allows for a huge reduction of manual labor cost in segmenting video frames because manual annotation is only required at the initial frame and a few refinement annotations only at frames containing challenging scenarios Yang et al. [2023c]. These manual and refinement annotation efforts may also be done with textual instruction that is converted into visual annotation by Grounding-DINO Cheng et al. [2023b]. In an unsupervised video tracking setting, tracking is initially done without prior information by the video salient object tracking (VSOT) model. This model produces object trajectory based on its prediction on video frames. To produce accurate segmentation at every frame, SAM can be prompted by the predicted trajectory Zhang et al. [2023g]. Finally, SAM's implementation in object tracking tasks serves as the initiation to more advanced work, such as action recognition Wang et al. [2023h].
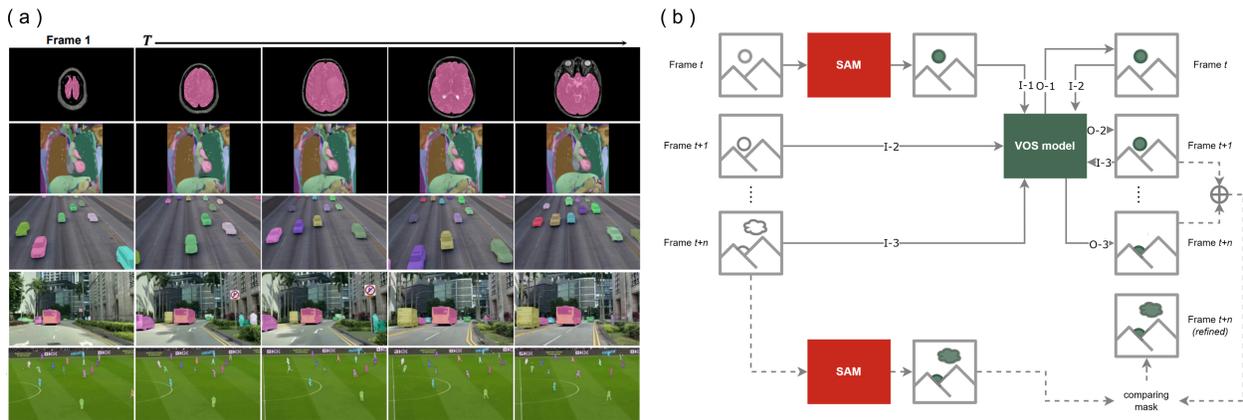


Figure 8: (a) Frame segmentation using video object segmentation (VOS) model Cheng et al. [2023b]. (b) General pipeline for SAM-integrated VOS (dashed line represents handling for newly appeared objects).

## 3.4 Anything in 3D

Previously, we have learned that inpainting and semantic labeling with SAM are applicable in 3D scenarios. Additionally, there are also other tasks that SAM is capable of guiding in 3D.

***Object Detection.*** One of the fundamental task in learning about 3D scenarios is object detection. The technique commonly used for this task is fundamentally similar to that of 2D images. The discrepancy lies in converting every annotation in 3D, *e.g.*, the bounding box becomes a bounding cube and the keypoint becomes a point cloud. Although it

is possible to conduct detection directly in 3D, it often requires enormous annotation effort. Therefore, a zero-shot technique is often desired, and this is where SAM is expected to take the role. Nonetheless, to accommodate SAM's segmentation capability, which is originally trained for 2D images, the 3D signal must be projected into 2D space before being processed by SAM. For instance, LiDAR point clouds can be projected into 2D bird-eye views, post-processed, and then fed into SAM together with mesh-grid prompts Zhang et al. [2023h]. The output mask can be further processed into a 2D box, which is then projected into 3D space to construct a bounding cube. This projection technique is also useful for assembling 3D segmentation masks in which a bi-directional merging technique can be used for assembling Yang et al. [2023d].

***Pose Estimation.*** Pose estimation is a common expansion from 3D object detection. Since this task requires the understanding of object movement with six degrees of freedom (6DoF), it is often desirable to observe directly in a 3D object. Nevertheless, the amount of 3D models available for training is often much more scarce than that of 2D images because building a proper 3D model is costly. Thus, again, zero-shot learning by applying SAM is preferable. Employing a similar 2D projection approach, 6DoF pose estimation further renders the pose using the rotation and translation elements in the feature descriptor of the SAM's segmentation output. Using this approach, SAM which only takes a few images of different views, is able to surpass the accuracy of the former 3D pose estimation model by halving their median pose error Fan et al. [2023]. Instead of rendering, hierarchical matching on object structure can also be applied to estimate the object pose, resulting in 10x faster runtime Chen et al. [2023d]. This technique utilizes SAM as a segmenter for multi-view images. The resulting mask is further combined with a depth image before the point features can be extracted. These features are then matched with the features extracted from the 3D model to generate the candidate pose for each single-view image (Figure 9).
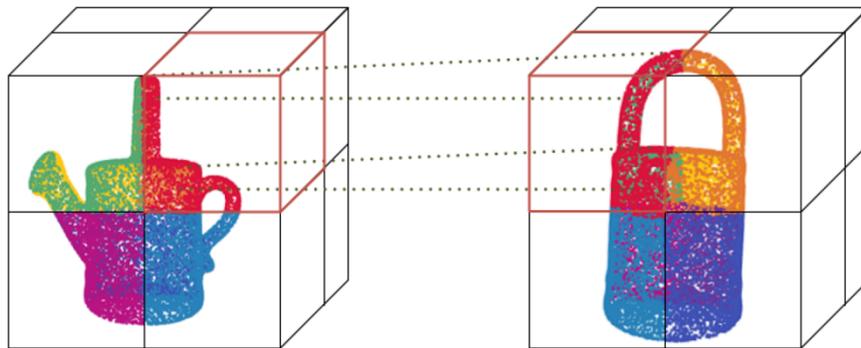


Figure 9: Illustration of how SAM's mask segmentations of multiple-view images are constructing 3D objects through point cloud matching Chen et al. [2023d].

## 3.5 Any-other-thing

The above X-anythings are the application of SAM in fundamental computer vision-related tasks. Nevertheless, segmentation can take care beyond those implementations. The generalizability of segmentation by SAM made it possible to extend to real-world scenarios. For instance, SAM can guide the audio-visual recognition task where it segments the target object that is most likely the source of audio. This task can be realized by fusing audio features and spatial features before feeding to the mask decoder of SAM Mo and Tian [2023]. SAM's fine-grained agnostic nature made it a helpful tool for eXplainable AI (XAI) because it is perceivable at the contextual level and less sensitive to the technical precision of the neural network (NN) model Sun et al. [2023b]. Utilizing this agnostic nature of SAM, another study developed an evaluation matrix for image translation tasks that compares spatial cosine similarity between semantic embedding of source and generated image, which is processed through SAM encoder Li et al. [2023b]. Further, another study also exploits this generalizability nature to perform calibration of LiDAR and camera in an autonomous driving scenario where SAM-generated mask is served as the container for LiDAR point cloud Luo et al. [2023]. Finally, many studies have proven that SAM is potentially reliably applied in many real-world applications, such as medical imaging Liu et al. [2023e], Wang et al. [2023i], civil construction Ahmadi et al. [2023], manufacturing Jain et al. [2023], molecular study Larsen et al. [2023], graph processing Jing et al. [2023], and communication Tariq et al. [2023].

# 4 A metric for evaluating SAM at everything mode

An important characteristic of SAM is that it can work in the mode of segmenting everything. The mode provides a straightforward way to visualize the quality of SAM. However, there is no metric to evaluate the performance of SAM in this mode. A major challenge in evaluating SAM at everything mode lies in that the predicted masks have no labels. In other words, the model only cut-outs the objects without assigning labels, therefore we call the everything mode as cut-out segmentation. Before we introduce our proposed metric for cut-out segmentation, we first summarize the metrics for existing image segmentation tasks, which are shown in Figure 10.
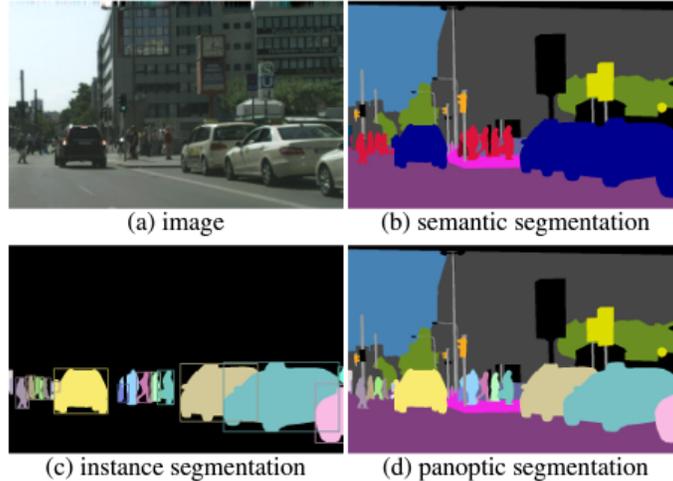


Figure 10: Illustration of semantic segmentation, instance segmentation, and panoptic segmentation. (figure obtained from Kirillov et al. [2019]).

## 4.1 Background on Image Segmentation

Segmentation is a crucial task in CV, which involves dividing an image into multiple regions based on semantic properties. It aims to extract meaningful information from images by identifying and separating the different objects or regions of interest. Segmentation task is a broad field:

**Semantic Segmentation.** In semantic segmentation, each pixel in an image is assigned a specific class label. A widely used metric to evaluate semantic segmentation is Mean Intersection over Union (mIoU) Han et al. [2023].

$$mIoU = \frac{1}{K+1} \sum_{i=0}^{K} \frac{TP_i}{TP_i + FP_i + FN_i} \tag{1}$$

**Instance Segmentation.** Instance segmentation is a task that not only assigns each pixel to a specific class label, but also distinguishes different instances of the same class. The metrics used to evaluate instance segmentation include Average Precision (AP), and Mean Average Precision (mAP) Henderson and Ferrari [2017].

$$AP_i = \frac{1}{m_i} \sum_{r=1}^{R} P(r)\Delta r \tag{2}$$

$$mAP = \frac{1}{m} \sum_{i=1}^{m} AP_i \tag{3}$$

**Panoptic Segmentation.** Panoptic segmentation combines both semantic and instance segmentation to provide a complete segmentation. Each pixel is labeled with either instance ID or semantic label, which depends on whether it belongs to an object instance or a background region. Panoptic Quality (PQ), Segmentation Quality (SQ), and Recognition Quality (RQ) are widely used to evaluate panoptic segmentation tasks Kirillov et al. [2019].

10

$$PQ = \underbrace{\frac{\sum_{(p,g) \in TP} IoU(p,g)}{|TP|}}_{\text{segmentation quality (SQ)}} \times \underbrace{\frac{|TP|}{|TP| + \frac{1}{2}|FP| + \frac{1}{2}|FN|}}_{\text{recognition quality (RQ)}} \qquad (4)$$
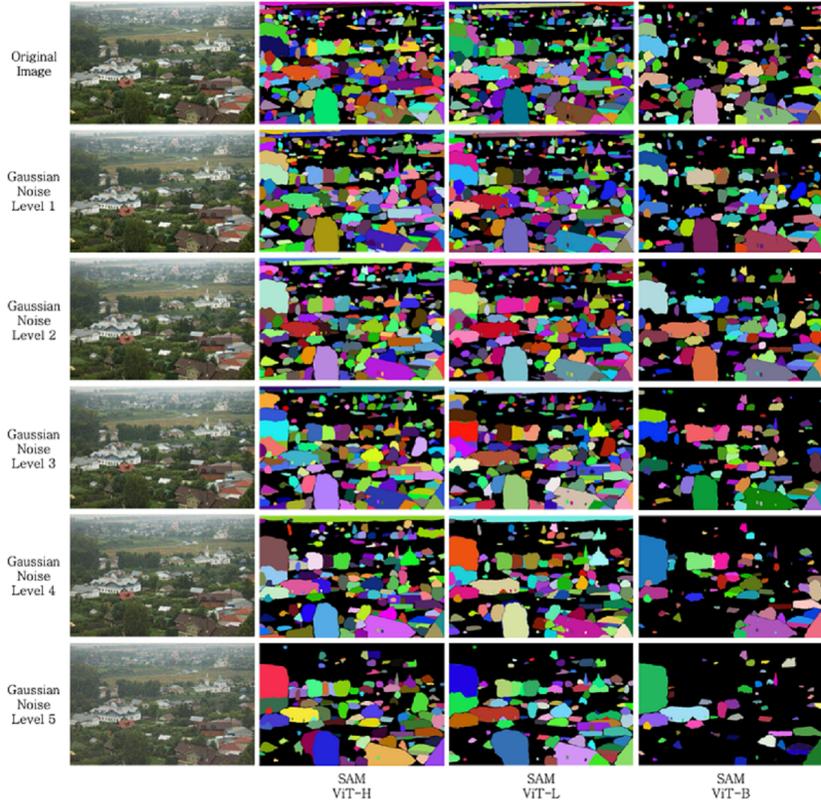


Figure 11: Segment Anything results depending on the model size, gaussian noise.

## 4.2 Greedy IoU: A Metric for Cut-out Segmentation

The above metrics cannot be directly applied to cut-out segmentation due to a lack of labels. Therefore, we propose a greedy search strategy to find the mask with the highest mIoU. Specifically, given the ground-truth mask, we seek the mask that has the highest IoU in the predicted cut-out segmentation maps. If the model fails to predict a valid mask, the resulting IoU is very low or zero. We then make an average over all the ground-truth masks to get a greedy mIoU. The algorithm is summarized in Algorithm 1. To show the use case of our proposed greedy mIoU, we take the generated output of SAM-ViT-H with clean samples as ground truth, and we compare and evaluate the performance of SAM-ViT-L and SAM-ViT-B. Moreover, the robustness of SAM models against Gaussian noise is also evaluated (see Figure 11). As we can see from Figure 11, SAM-ViT-L works comparably as SAM-ViT-L, while SAM-ViT-B often fails to detect and segment the objects. Moreover, as the severity level of Gaussian noise increases Hendrycks and Dietterich [2019], the performance of the SAM models also decreases, which is expected. A quantitative evaluation with our proposed greedy mIoU metric is shown in Table 1.

---
**Algorithm 1:** Greedy IoU algorithim

---
**Input:** Ground Truth Masks $M_{GT}$, Predicted Masks $M_P$
**Output:**
**for** $M_{GT}^i$ **in** $M_{GT}$, i = 1,2, ... K **do**
    **for** $M_P^j$ **in** $M_P$, j = 1,2, ... L **do**
        $iou_{ij}$ = $IoU(M_{GT}^i, M_P^j)$
    **end**
    $iou_i$ = max($iou_{i1}, iou_{i2}, ..., iou_{iL}$)
**end**
**return** *mean($iou_1, iou_2, ..., iou_K$)*

---

| Method | Clean | G.N Level 1 | G.N Level 2 | G.N Level 3 | G.N Level 4 | G.N Level 5 |
|---|---|---|---|---|---|---|
| SAM-ViT-B | 0.6247 | 0.5442 | 0.5335 | 0.5197 | 0.5000 | 0.4868 |
| SAM-ViT-L | 0.9157 | 0.8575 | 0.8492 | 0.8410 | 0.8300 | 0.8156 |

Table 1: Quantitative evaluation results based on our proposed greedy mIoU metric. G.N denotes Gaussian Noise.

## 5   Conclusion

Based on the task of promptable segmentation, the segment anything model (SAM) is the first vision foundation model that mimics the human eye to understand the world and its emergence has transformed the computer vision community. Our work conducts the first yet comprehensive survey on SAM. We hope our survey helps readers interested in SAM for performing their research.

## References

Chaoning Zhang, Chenshuang Zhang, Chenghao Li, Yu Qiao, Sheng Zheng, Sumit Kumar Dam, Mengchun Zhang, Jung Uk Kim, Seong Tae Kim, Jinwoo Choi, et al. One small step for generative ai, one giant leap for agi: A complete survey on chatgpt in aigc era. *arXiv preprint arXiv:2304.06488*, 2023a.

Chaoning Zhang, Chenshuang Zhang, Sheng Zheng, Yu Qiao, Chenghao Li, Mengchun Zhang, Sumit Kumar Dam, Chu Myaet Thwal, Ye Lin Tun, Le Luang Huy, et al. A complete survey on generative ai (aigc): Is chatgpt from gpt-4 to gpt-5 all you need? *arXiv preprint arXiv:2303.11717*, 2023b.

Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*, 2021.

Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *ICML*, 2021.

Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. In *ICML*, 2021.

Lu Yuan, Dongdong Chen, Yi-Ling Chen, Noel Codella, Xiyang Dai, Jianfeng Gao, Houdong Hu, Xuedong Huang, Boxin Li, Chunyuan Li, et al. Florence: A new foundation model for computer vision. *arXiv preprint arXiv:2111.11432*, 2021.

Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *CVPR*, 2020.

Yu Qiao, Md Munir, Apurba Adhikary, Huy Q Le, Avi Deb Raha, Chaoning Zhang, Choong Seon Hong, et al. Mp-fedcl: Multi-prototype federated contrastive learning for edge intelligence. *arXiv preprint arXiv:2304.01950*, 2023a.

Chaoning Zhang, Kang Zhang, Chenshuang Zhang, Trung X Pham, Chang D Yoo, and In So Kweon. How does simsiam avoid collapse without negative samples? a unified understanding with self-supervised contrastive learning. In *ICLR*, 2022a.

Chaoning Zhang, Chenshuang Zhang, Junha Song, John Seon Keun Yi, Kang Zhang, and In So Kweon. A survey on masked autoencoder for self-supervised learning in vision and beyond. *arXiv preprint arXiv:2208.00173*, 2022b.

Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. *arXiv preprint arXiv:2304.02643*, 2023.

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*, 2020.

Sheng He, Rina Bao, Jingpeng Li, P Ellen Grant, and Yangming Ou. Accuracy of segment-anything model (sam) in medical image segmentation tasks. *arXiv preprint arXiv:2304.09324*, 2023a.

Chuanfei Hu and Xinde Li. When sam meets medical images: An investigation of segment anything model (sam) on multi-phase liver tumor segmentation. *arXiv preprint arXiv:2304.08506*, 2023.

Tao Zhou, Yizhe Zhang, Yi Zhou, Ye Wu, and Chen Gong. Can sam segment polyps? *arXiv preprint arXiv:2304.07583*, 2023a.

Yuhao Huang, Xin Yang, Lian Liu, Han Zhou, Ao Chang, Xinrui Zhou, Rusi Chen, Junxuan Yu, Jiongquan Chen, Chaoyu Chen, Haozhe Chi, Xindi Hu, Deng-Ping Fan, Fajin Dong, and Dong Ni. Segment anything model for medical images?, 2023a.

Mingzhe Hu, Yuheng Li, and Xiaofeng Yang. Breastsam: A study of segment anything model for breast tumor detection in ultrasound images. *arXiv preprint arXiv:2305.12447*, 2023a.

Sovesh Mohapatra, Advait Gosai, and Gottfried Schlaug. Brain extraction comparing segment anything model (sam) and fsl brain extraction tool. *arXiv preprint arXiv:2304.04738*, 2023.

Dongjie Cheng, Ziyuan Qin, Zekun Jiang, Shaoting Zhang, Qicheng Lao, and Kang Li. Sam on medical images: A comprehensive study on three prompt modes, 2023a.

An Wang, Mobarakol Islam, Mengya Xu, Yang Zhang, and Hongliang Ren. Sam meets robotic surgery: An empirical study in robustness perspective, 2023a.

Tassilo Wald, Saikat Roy, Gregor Koehler, Nico Disch, Maximilian Rouven Rokuss, Julius Holzschuh, David Zimmerer, and Klaus Maier-Hein. Sam.md: Zero-shot medical image segmentation capabilities of the segment anything model. In *Medical Imaging with Deep Learning, short paper track*, 2023. URL `https://openreview.net/forum?id=iilLHaINUW`.

Christian Mattjie, Luis Vinicius de Moura, Rafaela Cappelari Ravazio, Lucas Silveira Kupssinskü, Otávio Parraga, Marcelo Mussi Delucis, and Rodrigo Coelho Barros. Zero-shot performance of the segment anything model (sam) in 2d medical imaging: A comprehensive evaluation and practical guidelines, 2023.

Maciej A. Mazurowski, Haoyu Dong, Hanxue Gu, Jichen Yang, Nicholas Konz, and Yixin Zhang. Segment anything model for medical image analysis: an experimental study, 2023.

Yuheng Li, Mingzhe Hu, and Xiaofeng Yang. Polyp-sam: Transfer sam for polyp segmentation, 2023a.

Mingzhe Hu, Yuheng Li, and Xiaofeng Yang. Skinsam: Empowering skin cancer segmentation with segment anything model, 2023b.

Jun Ma and Bo Wang. Segment anything in medical images. *arXiv preprint arXiv:2304.12306*, 2023.

Zhongxi Qiu, Yan Hu, Heng Li, and Jiang Liu. Learnable ophthalmology sam, 2023.

Kaidong Zhang and Dong Liu. Customized segment anything model for medical image segmentation, 2023.

Junde Wu, Yu Zhang, Rao Fu, Huihui Fang, Yuanpei Liu, Zhaowei Wang, Yanwu Xu, and Yueming Jin. Medical sam adapter: Adapting segment anything model for medical image segmentation, 2023.

Yifan Gao, Wei Xia, Dingdu Hu, and Xin Gao. Desam: Decoupling segment anything model for generalizable medical image segmentation. *arXiv preprint arXiv:2306.00499*, 2023.

Ge-Peng Ji, Deng-Ping Fan, Peng Xu, Ming-Ming Cheng, Bowen Zhou, and Luc Van Gool. Sam struggles in concealed scenes–empirical study on" segment anything". *arXiv preprint arXiv:2304.06022*, 2023a.

Wei Ji, Jingjing Li, Qi Bi, Wenbo Li, and Li Cheng. Segment anything is not always perfect: An investigation of sam on different real-world applications. *arXiv preprint arXiv:2304.05750*, 2023b.

Simiao Ren, Francesco Luzi, Saad Lahrichi, Kaleb Kassaw, Leslie M. Collins, Kyle Bradbury, and Jordan M. Malof. Segment anything, from space?, 2023.

Leiping Jie and Hui Zhang. When sam meets shadow detection. *arXiv preprint arXiv:2305.11513*, 2023.

Xiao Yang, Haixing Dai, Zihao Wu, Ramesh Bist, Sachin Subedi, Jin Sun, Guoyu Lu, Changying Li, Tianming Liu, and Lilong Chai. Sam for poultry science. *arXiv preprint arXiv:2305.10254*, 2023a.

Lv Tang, Haoke Xiao, and Bo Li. Can sam segment anything? when sam meets camouflaged object detection. *arXiv preprint arXiv:2304.04709*, 2023.

Dongsheng Han, Chaoning Zhang, Yu Qiao, Maryam Qamar, Yuna Jung, SeungKyu Lee, Sung-Ho Bae, and Choong Seon Hong. Segment anything model (sam) meets glass: Mirror and transparent objects cannot be easily detected. *arXiv preprint*, 2023.

Zhiheng Ma, Xiaopeng Hong, and Qinnan Shangguan. Can sam count anything? an empirical study on sam counting, 2023.

Xiaoyu Guo, Xiang Wei, Qi Su, Huiqin Zhao, and Shunli Zhan. Prompt what you need: Enhancing segmentation in rainy scenes with anchor-based prompting. *arXiv preprint arXiv:2305.03902*, 2023.

Chunming He, Kai Li, Yachao Zhang, Guoxia Xu, Longxiang Tang, Yulun Zhang, Zhenhua Guo, and Xiu Li. Weakly-supervised concealed object segmentation with sam-based pseudo labeling and multi-scale feature grouping. *arXiv preprint arXiv:2305.11003*, 2023b.

Lei Ke, Mingqiao Ye, Martin Danelljan, Yifan Liu, Yu-Wing Tai, Chi-Keung Tang, and Fisher Yu. Segment anything in high quality. *arXiv preprint arXiv:2306.01567*, 2023.

Iraklis Giannakis, Anshuman Bhardwaj, Lydia Sam, and Georgios Leontidis. Deep learning universal crater detection using segment anything model (sam), 2023.

Dominic Williams, Fraser MacFarlane, and Avril Britten. Leaf only sam: A segment anything pipeline for zero-shot automated leaf segmentation. *arXiv preprint arXiv:2305.09418*, 2023.

Tianrun Chen, Lanyun Zhu, Chaotao Ding, Runlong Cao, Shangzhan Zhang, Yan Wang, Zejian Li, Lingyun Sun, Papa Mao, and Ying Zang. Sam fails to segment anything?–sam-adapter: Adapting sam in underperformed scenes: Camouflage, shadow, and more. *arXiv preprint arXiv:2304.09148*, 2023a.

Sahib Julka and Michael Granitzer. Knowledge distillation with segment anything (sam) model for planetary geological mapping. *arXiv preprint arXiv:2305.07586*, 2023.

Yunkang Cao, Xiaohao Xu, Chen Sun, Yuqi Cheng, Zongwei Du, Liang Gao, and Weiming Shen. Segment any anomaly without training via hybrid prompt regularization, 2023.

Shehbaz Tariq, Brian Estadimas Arfeto, Chaoning Zhang, and Hyundong Shin. Segment anything meets semantic communication. *arXiv preprint arXiv:2306.02094*, 2023.

Chenshuang Zhang, Chaoning Zhang, Taegoo Kang, Donghun Kim, Sung-Ho Bae, and In So Kweon. Attack-sam: Towards evaluating adversarial robustness of segment anything model. *arXiv preprint*, 2023c.

Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*, 2013.

Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. In *ICLR*, 2015.

Alexey Kurakin, Ian Goodfellow, and Samy Bengio. Adversarial machine learning at scale. In *ICLR*, 2017.

Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. In *ICLR*, 2018.

Yihao Huang, Yue Cao, Tianlin Li, Felix Juefei-Xu, Di Lin, Ivor W Tsang, Yang Liu, and Qing Guo. On the robustness of segment anything. *arXiv preprint arXiv:2305.16220*, 2023b.

Yuqing Wang, Yun Zhao, and Linda Petzold. An empirical study on the robustness of the segment anything model (sam). *arXiv preprint arXiv:2305.06422*, 2023b.

Yu Qiao, Chaoning Zhang, Taegoo Kang, Donghun Kim, Shehbaz Tariq, Chenshuang Zhang, and Choong Seon Hong. Robustness of sam: Segment anything under corruptions and beyond. *arXiv preprint arXiv:2306.07713*, 2023b.

Qifan Yu, Juncheng Li, Wentao Ye, Siliang Tang, and Yueting Zhuang. Interactive data synthesis for systematic vision adaptation via llms-aigcs collaboration. *arXiv preprint arXiv:2305.12799*, 2023a.

Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *ICML*, pages 12888–12900. PMLR, 2022.

Runnan Chen, Youquan Liu, Lingdong Kong, Nenglun Chen, Xinge Zhu, Yuexin Ma, Tongliang Liu, and Wenping Wang. Towards label-free scene understanding by vision foundation models. *arXiv preprint arXiv:2306.03899*, 2023b.

Zhimin Chen and Bing Li. Bridging the domain gap: Self-supervised 3d scene understanding with foundation models. *arXiv preprint arXiv:2305.08776*, 2023.

Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba. Learning deep features for discriminative localization, 2015.

Tianle Chen, Zheda Mai, Ruiwen Li, and Wei-lun Chao. Segment anything model (sam) enhanced pseudo labels for weakly supervised semantic segmentation. *arXiv preprint arXiv:2305.05803*, 2023c.

Weixuan Sun, Zheyuan Liu, Yanhao Zhang, Yiran Zhong, and Nick Barnes. An alternative to wsss? an empirical study of the segment anything model (sam) on weakly-supervised semantic segmentation problems, 2023a.

Yulin He, Wei Chen, Yusong Tan, and Siqi Wang. Usd: Unknown sensitive detector empowered by decoupled objectness and segment anything model. *arXiv preprint arXiv:2306.02275*, 2023c.

Peng-Tao Jiang and Yuqi Yang. Segment anything is a good pseudo-label generator for weakly supervised semantic segmentation, 2023.

Feng Liang, Bichen Wu, Xiaoliang Dai, Kunpeng Li, Yinan Zhao, Hang Zhang, Peizhao Zhang, Peter Vajda, and Diana Marculescu. Open-vocabulary semantic segmentation with mask-adapted clip, 2023.

Jun Cen, Yizheng Wu, Kewei Wang, Xingyi Li, Jingkang Yang, Yixuan Pei, Lingdong Kong, Ziwei Liu, and Qifeng Chen. Sad: Segment any rgbd. *arXiv preprint arXiv:2305.14207*, 2023a.

Yang Liu, Muzhi Zhu, Hengtao Li, Hao Chen, Xinlong Wang, and Chunhua Shen. Matcher: Segment anything with one shot using all-purpose feature matching. *arXiv preprint arXiv:2305.13310*, 2023a.

Renrui Zhang, Zhengkai Jiang, Ziyu Guo, Shilin Yan, Junting Pan, Hao Dong, Peng Gao, and Hongsheng Li. Personalize segment anything model with one shot. *arXiv preprint arXiv:2305.03048*, 2023d.

Ahmed Abdelreheem, Abdelrahman Eldesokey, Maks Ovsjanikov, and Peter Wonka. Zero-shot 3d shape correspondence. *arXiv preprint arXiv:2306.03253*, 2023.

Teng Wang, Jinrui Zhang, Junjie Fei, Yixiao Ge, Hao Zheng, Yunlong Tang, Zhe Li, Mingqi Gao, Shanshan Zhao, Ying Shan, et al. Caption anything: Interactive image description with diverse multimodal controls. *arXiv preprint arXiv:2305.02677*, 2023c.

Sivan Doveh, Assaf Arbelle, Sivan Harary, Amit Alfassy, Roei Herzig, Donghyun Kim, Raja Giryes, Rogerio Feris, Rameswar Panda, Shimon Ullman, et al. Dense and aligned captions (dac) promote compositional reasoning in vl models. *arXiv preprint arXiv:2305.19595*, 2023.

Rongsheg Wang, Yaofei Duan, and Yukun Li. Segment anything also detect anything. Technical report, EasyChair, 2023d.

Yizhe Zhang, Tao Zhou, Peixian Liang, and Danny Z Chen. Input augmentation with sam: Boosting medical image segmentation with segmentation foundation model. *arXiv preprint arXiv:2304.11332*, 2023e.

Di Wang, Jing Zhang, Bo Du, Dacheng Tao, and Liangpei Zhang. Scaling-up remote sensing segmentation dataset with segment anything model. *arXiv preprint arXiv:2305.02034*, 2023e.

Anzhu Yu, Wenjun Huang, Qing Xu, Qun Sun, Wenyue Guo, Song Ji, Bowei Wen, and Chunping Qiu. Sea ice extraction via remote sensed imagery: Algorithms, datasets, applications and challenges. *arXiv preprint arXiv:2306.00303*, 2023b.

Jielu Zhang, Zhongliang Zhou, Gengchen Mai, Lan Mu, Mengxuan Hu, and Sheng Li. Text2seg: Remote sensing image semantic segmentation via text-guided visual foundation models. *arXiv preprint arXiv:2304.10597*, 2023f.

Zhuyun Zhou, Zongwei Wu, Rémi Boutteau, Fan Yang, and Dominique Ginhac. Dsec-mos: Segment any moving object with moving ego vehicle, 2023b.

Junzhang Chen and Xiangzhi Bai. Learning to" segment anything" in thermal infrared images through knowledge distillation with a large scale dataset satir. *arXiv preprint arXiv:2304.07969*, 2023.

Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10684–10695, 2022.

Defeng Xie, Ruichen Wang, Jian Ma, Chen Chen, Haonan Lu, Dong Yang, Fobo Shi, and Xiaodong Lin. Edit everything: A text-guided generative system for images editing, 2023.

Eran Levin and Ohad Fried. Differential diffusion: Giving each pixel its strength. *arXiv preprint arXiv:2306.00950*, 2023.

Tao Yu, Runseng Feng, Ruoyu Feng, Jinming Liu, Xin Jin, Wenjun Zeng, and Zhibo Chen. Inpaint anything: Segment anything meets image inpainting. *arXiv preprint arXiv:2304.06790*, 2023c.

Zhaoyang Liu, Yinan He, Wenhai Wang, Weiyun Wang, Yi Wang, Shoufa Chen, Qinglong Zhang, Yang Yang, Qingyun Li, Jiashuo Yu, et al. Internchat: Solving vision-centric tasks by interacting with chatbots beyond language. *arXiv preprint arXiv:2305.05662*, 2023b.

Qian Wang, Biao Zhang, Michael Birsak, and Peter Wonka. Instructedit: Improving automatic masks for diffusion-based image editing with user instructions. *arXiv preprint arXiv:2305.18047*, 2023f.

Guian Fang, Zutao Jiang, Jianhua Han, Guangsong Lu, Hang Xu, and Xiaodan Liang. Boosting text-to-image diffusion models with fine-grained semantic rewards. *arXiv preprint arXiv:2305.19599*, 2023.

IDEA-Research. Grounded segment anything, 2023. URL `https://github.com/IDEA-Research/Grounded-Segment-Anything`. GitHub repository.

Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chunyuan Li, Jianwei Yang, Hang Su, Jun Zhu, et al. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. *arXiv preprint arXiv:2303.05499*, 2023c.

Songhua Liu, Jingwen Ye, and Xinchao Wang. Any-to-any style transfer. *arXiv preprint arXiv:2304.09728*, 2023d.

Jiaxi Jiang and Christian Holz. Restore anything pipeline: Segment anything meets image restoration. *arXiv preprint arXiv:2305.13093*, 2023.

Zeyu Xiao, Jiawang Bai, Zhihe Lu, and Zhiwei Xiong. A dive into sam prior in image restoration. *arXiv preprint arXiv:2305.13620*, 2023.

Zhihe Lu, Zeyu Xiao, Jiawang Bai, Zhiwei Xiong, and Xinchao Wang. Can sam boost video super-resolution? *arXiv preprint arXiv:2305.06524*, 2023.

Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM*, 65(1):99–106, 2021.

Dongqing Wang, Tong Zhang, Alaa Abboud, and Sabine Süsstrunk. Inpaintnerf360: Text-guided 3d inpainting on unbounded neural radiance fields. *arXiv preprint arXiv:2305.15094*, 2023g.

Jiazhong Cen, Zanwei Zhou, Jiemin Fang, Wei Shen, Lingxi Xie, Dongsheng Jiang, Xiaopeng Zhang, and Qi Tian. Segment anything in 3d with nerfs, 2023b.

Qiuhong Shen, Xingyi Yang, and Xinchao Wang. Anything-3d: Towards single-view anything reconstruction in the wild. *arXiv preprint arXiv:2304.10261*, 2023.

Dan Ogawa Lillrank, Shogo Akiyama, and Kai Arulkumaran. Zero-shot object manipulation with semantic 3d image augmentation for perceiver-actor. 2023.

Shohei Mori, Sei Ikeda, and H. Saito. A survey of diminished reality: Techniques for visually concealing, eliminating, and seeing through real objects. *IPSJ Transactions on Computer Vision and Applications*, 9:1–14, 2017.

Lik-Hang Lee, Pengyuan Zhou, Chaoning Zhang, and Simo Johannes Hosio. What if we have meta gpt? from content singularity to human-metaverse interaction in aigc era. *ArXiv*, abs/2304.07521, 2023.

ngthanhtin. owlvit segment anything, 2023. URL `https://github.com/ngthanhtin/owlvit_segment_anything`. GitHub repository.

Jingfeng Yao, Xinggang Wang, Lang Ye, and Wenyu Liu. Matte anything: Interactive natural image matting with segment anything models. *arXiv preprint arXiv:2306.04121*, 2023.

Lingfeng Yang, Yueze Wang, Xiang Li, Xinlong Wang, and Jian Yang. Fine-grained visual prompting. *arXiv preprint arXiv:2306.04356*, 2023b.

Haibin He, Jing Zhang, Mengyang Xu, Juhua Liu, Bo Du, and Dacheng Tao. Scalable mask annotation for video text spotting, 2023d.

Ho Kei Cheng and Alexander G. Schwing. Xmem: Long-term video object segmentation with an atkinson-shiffrin memory model, 2022.

Zongxin Yang and Yi Yang. Decoupling features in hierarchical propagation for video object segmentation. *NeurIPS*, 2022.

Jinyu Yang, Mingqi Gao, Zhe Li, Shang Gao, Fangjing Wang, and Feng Zheng. Track anything: Segment anything meets videos. *arXiv preprint arXiv:2304.11968*, 2023c.

Yangming Cheng, Liulei Li, Yuanyou Xu, Xiaodi Li, Zongxin Yang, Wenguan Wang, and Yi Yang. Segment and track anything. *arXiv preprint arXiv:2305.06558*, 2023b.

Zhenghao Zhang, Zhichao Wei, Shengfan Zhang, Zuozhuo Dai, and Siyu Zhu. Uvosam: A mask-free paradigm for unsupervised video object segmentation via segment anything model. *arXiv preprint arXiv:2305.12659*, 2023g.

Xijun Wang, Ruiqi Xian, Tianrui Guan, and Dinesh Manocha. Prompt learning for action recognition. *arXiv preprint arXiv:2305.12437*, 2023h.

Dingyuan Zhang, Dingkang Liang, Hongcheng Yang, Zhikang Zou, Xiaoqing Ye, Zhe Liu, and Xiang Bai. Sam3d: Zero-shot 3d object detection via segment anything model. *arXiv preprint arXiv:2306.02245*, 2023h.

Yunhan Yang, Xiaoyang Wu, Tong He, Hengshuang Zhao, and Xihui Liu. Sam3d: Segment anything in 3d scenes. *arXiv preprint arXiv:2306.03908*, 2023d.

Zhiwen Fan, Panwang Pan, Peihao Wang, Yifan Jiang, Dejia Xu, Hanwen Jiang, and Zhangyang Wang. Pope: 6-dof promptable pose estimation of any object, in any scene, with one reference. *arXiv preprint arXiv:2305.15727*, 2023.

Jianqiu Chen, Mingshan Sun, Tianpeng Bao, Rui Zhao, Liwei Wu, and Zhenyu He. 3d model-based zero-shot pose estimation pipeline. *arXiv preprint arXiv:2305.17934*, 2023d.

Shentong Mo and Yapeng Tian. Av-sam: Segment anything model meets audio-visual localization and segmentation. *arXiv preprint arXiv:2305.01836*, 2023.

Ao Sun, Pingchuan Ma, Yuanyuan Yuan, and Shuai Wang. Explain any concept: Segment anything meets concept-based explanation. *arXiv preprint arXiv:2305.10289*, 2023b.

Yunxiang Li, Meixu Chen, Wenxuan Yang, Kai Wang, Jun Ma, Alan C Bovik, and You Zhang. Samscore: A semantic structural similarity metric for image translation evaluation. *arXiv preprint arXiv:2305.15367*, 2023b.

Zhaotong Luo, Guohang Yan, and Yikang Li. Calib-anything: Zero-training lidar-camera extrinsic calibration method using segment anything. *arXiv preprint arXiv:2306.02656*, 2023.

Yihao Liu, Jiaming Zhang, Zhangcong She, Amir Kheradmand, and Mehran Armand. Samm (segment any medical model): A 3d slicer integration to sam. *arXiv preprint arXiv:2304.05622*, 2023e.

Bin Wang, Armstrong Aboah, Zheyuan Zhang, and Ulas Bagci. Gazesam: What you see is what you segment, 2023i.

Mohsen Ahmadi, Ahmad Gholizadeh Lonbar, Abbas Sharifi, Ali Tarlani Beris, Mohammadsadegh Nouri, and Amir Sharifzadeh Javidi. Application of segment anything model for civil infrastructure defect assessment, 2023.

Siddarth Jain, Radu Corcodel, Devesh K Jha, and Diego Romeres. Vision guided food assembly by robot teaching from target composition. 2023.

Rasmus Larsen, Torben L Villadsen, Jette K Mathiesen, Kirsten MØ Jensen, and Espen D Boejesen. Np-sam: Implementing the segment anything model for easy nanoparticle segmentation in electron microscopy images. 2023.

Yongcheng Jing, Xinchao Wang, and Dacheng Tao. Segment anything in non-euclidean domains: Challenges and opportunities. *arXiv preprint arXiv:2304.11595*, 2023.

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.

Paul Henderson and Vittorio Ferrari. End-to-end training of object class detectors for mean average precision. In *Computer Vision–ACCV 2016: 13th Asian Conference on Computer Vision, Taipei, Taiwan, November 20-24, 2016, Revised Selected Papers, Part V 13*, pages 198–213. Springer, 2017.

Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common corruptions and perturbations. *ICLR*, 2019.

Kevmo. magic-copy, 2023. URL `https://github.com/kevmo314/magic-copy`. GitHub repository.

Feizc. Iea, 2023. URL `https://github.com/feizc/IEA`. GitHub repository.

Gasvn. Editanything, 2023. URL `https://github.com/sail-sg/EditAnything`. GitHub repository.

Peize Sun, Shoufa Chen, and Ping Luo. Grounded segment anything: From objects to parts. `https://github.com/Cheems-Seminar/grounded-segment-any-parts`, 2023c.

Jiaqi Chen, Zeyu Yang, and Li Zhang. Semantic-segment-anything, 2023e. URL `https://github.com/fudan-zvg/Semantic-Segment-Anything`. GitHub repository.

Curt Park. segment anything with clip, 2023. URL `https://github.com/Curt-Park/segment-anything-with-clip`. GitHub repository.

Vietanhdev. Anylabeling, 2023. URL `https://github.com/vietanhdev/anylabeling`. GitHub repository.

RockeyCoss. Prompt-segment-anything, 2023. URL `https://github.com/RockeyCoss/Prompt-Segment-Anything`. GitHub repository.

Amine. Sam-medical-imaging, 2023. URL `https://github.com/amine0110/SAM-Medical-Imaging`. GitHub repository.

Karol. Segment anything model (sam) in napari, 2023. URL `https://github.com/MIC-DKFZ/napari-sam`. GitHub repository.

Jo Okuma. napari-segment-anything, 2023. URL `https://github.com/JoOkuma/napari-segment-anything`. GitHub repository.

Yukang Chen. 3d box segment anything, 2023. URL `https://github.com/dvlab-research/3D-Box-Segment-Anything`. GitHub repository.

# A  Appendix

Table 2: List of X-anything projects on GitHub that exploits SAM's segmentation potential.

| Function | Application and Key Idea |
|---|---|
| **Inpainting** | *Magic Copy* Kevmo [2023]<br>  Mate image interactively<br>*Image Editing Anything* Feizc [2023], *Edit Anything* Gasvn [2023]<br>  Inpaint image with AIGC using point and text prompt<br>*Open Vocabulary Segment Anything* ngthanhtin [2023]<br>  Inpaint image with AIGC using text prompt<br>*Grounded Segment Anything* Sun et al. [2023c]<br>  Segment object at fine-grained part level |
| **Labeling** | *Semantic Segment Anything* Chen et al. [2023e]<br>  Label SAM's mask through voting with semantic segmenter mask<br>*Segment Anything with CLIP* Park [2023]<br>  Label SAM mask output by calculating its similarity with CLIP<br>*AnyLabeling* Vietanhdev [2023]<br>  Annotate object interactively using YOLO and SAM<br>*Prompt Segment Anything* RockeyCoss [2023]<br>  Label SAM's mask<br>*SAM Medical Imaging* Amine [2023]<br>  Segment medical imaging using DICOM files<br>*SAM in Napari* Karol [2023], *Napari Segment Anything* Okuma [2023]<br>  Offline interactive medical image segmentation using Napari |
| **Tracking** | *Track Anything* Yang et al. [2023c], *SAM-Track* Cheng et al. [2023b]<br>  Track selected object in video frame with minimum annotation cost |
| **3D-Scene** | *3D-Box via SAM* Chen [2023]<br>  3D object detection using VoxelNeXT |

Table 3: List of X-anything study in computer vision that exploits SAM's segmentation potential.

| Function | Application and Key Idea |
|---|---|
| **Inpainting** | *Edit Everything* Xie et al. [2023], *Inpaint Anything* Yu et al. [2023c]<br>  Instruct editing target and style of AIGC using text prompt<br>*Differential Diffusion* Levin and Fried [2023]<br>  Control editing magnitude through text instruction in AIGC loop<br>*InternChat* Liu et al. [2023b]<br>  Integrate interactive text prompt with visual prompt for editing<br>*Any-to-Any Style Transfer* Liu et al. [2023d]<br>  Transfer style between images with attention on the target part segmented by SAM<br>*SPT Image Restoration* Xiao et al. [2023], *SAM-Guided Refinement Module* Lu et al. [2023]<br>  Tune image/video restoration by injecting SAM's masks as prior<br>*Restore Anything* Jiang and Holz [2023]<br>  Select restoration region by SAM's mask<br>*InpaintNeRF360* Wang et al. [2023g]<br>  Inpaint 3D object segmented by SAM using text prompt<br>*Anything-3D* Shen et al. [2023]<br>  Generate 3D object from single-view image segmented by SAM using text prompt<br>*Segment Anything in 3D* Cen et al. [2023b]<br>  Generate 3D object from single-view image through cross-view rendering and self-prompting<br>*3D Augmentation PerAct* Lillrank et al. [2023]<br>  Manipulate robotic vision using SAM's mask to select target object<br>*MatAny* Yao et al. [2023]<br>  Generate trimap for image matting from SAM's mask<br>*Fine-Grained Visual Prompting* Yang et al. [2023b]<br>  Separate foreground object from background for visual prompting task using SAM's mask<br>*InstructEdit* Wang et al. [2023f]<br>  Instruct inpainting task with Grounded-SAM prompted by text<br>*FineRewards* Fang et al. [2023]<br>  Refine prompt for text-to-image diffusion model through caption matching |
| **Labeling** | *ChatGenImage* Yu et al. [2023a]<br>  Label AIGC generated image iteratively by combining LLM, LVM, and VLM<br>*Segment Any RGBD* Cen et al. [2023a]<br>  Label SAM's mask through voting with OVSeg-generated mask<br>*Matcher* Liu et al. [2023a], *Personalized SAM* Zhang et al. [2023d], *SAM-3D* Abdelreheem et al. [2023]<br>  Transfer segmentation of same-class object using semantic feature matching keypoints as SAM prompt<br>*Cross-Modality Noise Supervision* Chen et al. [2023b]<br>  Label 3D scene using image captioning foundation model and SAM<br>*Bridge3D* Chen and Li [2023]<br>  Label 3D scene using image captioning foundation model and masked autoencoder<br>*Caption Anything* Wang et al. [2023c]<br>  Selective captioning on a specific region in an image using SAM<br>*Dense and Aligned Captions* Doveh et al. [2023]<br>  Evaluate the similarity between the word composition of the caption and image spatial composition |
| **Tracking** | *Track Anything* Yang et al. [2023c], *SAM-Track* Cheng et al. [2023b]<br>  Track selected object in video frame with minimum annotation cost<br>*UVOSAM* Zhang et al. [2023g]<br>  Track the object in a video without supervision and use the predicted trajectory to build a segmentation mask<br>*PLAR* Wang et al. [2023h]<br>  Recognize action of target object segmented by SAM |
| **3D-Scene** | *SAM3D* Zhang et al. [2023h], *SAM3D* Yang et al. [2023d]<br>  Zero-shot 3D segmentation task that is includes the projection into 2D space<br>*Promptable Object Pose Estimation* Fan et al. [2023]<br>  Estimate 3D object pose by rendering features from limited view images<br>*3D Model-Base Pose Estimation* Chen et al. [2023d]<br>  Estimate 3D object pose by hierarchically matching features from multi-view images and 3D model |