

# VLTP: Vision-Language Guided Token Pruning for Task-Oriented Segmentation

Hanning Chen<sup>1\*</sup> Yang Ni<sup>1</sup> Wenjun Huang<sup>1</sup> Yezi Liu<sup>1</sup> SungHeon Jeong<sup>1</sup>  
 Fei Wen<sup>2</sup> Nathaniel Bastian<sup>3</sup> Hugo Latapie<sup>4†</sup> Mohsen Imani<sup>1</sup>

<sup>1</sup> University of California, Irvine, CA, USA

<sup>2</sup> Texas A&M University, College Station, TX, USA

<sup>3</sup> United States Military Academy, West Point, USA

<sup>4</sup> Cisco, San Jose, CA, USA

{hanningc, m.imani}@uci.edu

## Abstract

*Vision Transformers (ViTs) have emerged as the backbone of many segmentation models, consistently achieving state-of-the-art (SOTA) performance. However, their success comes at a significant computational cost. Image token pruning is one of the most effective strategies to address this complexity. However, previous approaches fall short when applied to more complex task-oriented segmentation (TOS), where the class of each image patch is not predefined but dependent on the specific input task. This work introduces the *Vision Language Guided Token Pruning* (VLTP), a novel token pruning mechanism that can accelerate ViT-based segmentation models, particularly for TOS guided by multi-modal large language model (MLLM). We argue that ViT does not need to process every image token through all of its layers—only the tokens related to reasoning tasks are necessary. We design a new pruning decoder to take both image tokens and vision-language guidance as input to predict the relevance of each image token to the task. Only image tokens with high relevance are passed to deeper layers of the ViT. Experiments show that the VLTP framework reduces the computational costs of ViT by approximately 25% without performance degradation and by around 40% with only a 1% performance drop. The code associated with this study can be found at [this URL](#).*

## 1. Introduction

ViTs [9] and multi-head self-attention (MHSA) [45] have significantly advanced the development of computer vision, where ViTs have been widely used as backbone models for various tasks [21], including image classification [4] and segmentation [22]. Compared to previous convolutional neural networks (CNN) [13, 15, 16, 54],

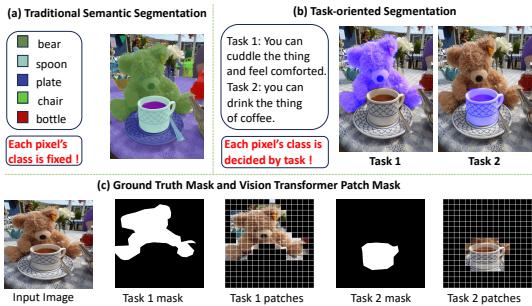


Figure 1. (a) Semantic segmentation example. (b) TOS example. (c) For the same image, the segmentation mask and corresponding image patches change when the input task changes.

ViT-based models excel in global reasoning through pairwise token attention. However, the success of ViT models comes with a significant computational overhead, which limits their deployment in resource-constrained environments, such as edge computing. This challenge is especially pronounced in segmentation tasks. Since segmentation requires detailed pixel-wise predictions, they often use high-resolution images, further increasing the number of image patches and leading to a quadratic increase in computation.

To improve the efficiency of ViT models, several techniques have been proposed [5]. Among the most notable methods is image token pruning [31, 44], which effectively accelerates model inference. These approaches aim to reduce the number of image tokens involved in ViT layers such as the MHSA computation, thus significantly decreasing the computational load due to the inherent quadratic complexity [45]. For example, the DToP method [44] aims to predict the semantic class of image tokens in the early layers of ViT, and those with high confidence can be omitted from computations in deeper ViT layers. Alternatively, CTS [31] aims to combine image tokens that share the same semantic class, thereby reducing the total token count. As illustrated in Fig. 1(a), these techniques

\*This work was done during an Cisco internship.

†Internship manager

work well for traditional segmentation tasks. However, they are less effective for more complex vision language tasks such as task-oriented segmentation (**TOS**) [38, 40], where the class of each pixel varies depending on the specified task. As seen in Fig. 1(b), TOS involves finding an item’s mask based on the query task. As the task changes, so does the ground truth mask, meaning that each pixel’s class also changes, as shown in Fig. 1(c). Directly applying previous methods [31, 44] would lead to a significant loss of visual information and thus low-quality masks during the reasoning stage, because these methods focus on the relevance of image patches to static semantic class, without considering a specific task. A further drawback of earlier approaches is their incompatibility with the latest trending MLLMs [28] since their pruning strategies did not account for external reasoning guidance. Recently, numerous frameworks [24, 33, 37, 56] have been proposed to guide vision models in performing vision language tasks using MLLM, such as Referring Expression Segmentation (RES) [53]. The strong reasoning capability of MLLM makes it promising for solving TOS tasks. Therefore, it is necessary to develop a new token pruning mechanism to accelerate ViT-based segmentation models, especially for TOS tasks guided by MLLM.

In this research, we present VLTP (**Vision-Language Guided Token Pruning**), a new ViT token pruning technique for TOS. The core concept of VLTP is that **only tokens relevant to reasoning tasks should be processed by all ViT layers**. We design a novel prune decoder and insert it at multiple selected layers of ViT to provide flexible multi-stage pruning. The prune decoder predicts the relevance of image tokens to the reasoning tasks, leveraging the guidance from MLLM. Based on the relevance predictions from the prune decoder, VLTP retains only the image tokens with high relation scores and forwards them to the following layers of ViT. By removing unrelated tokens that participated in ViT computation, we effectively accelerate the ViT model while preserving its high reasoning efficacy. To mitigate mispredictions by the prune decoder (e.g., due to pruning in the early stages of ViT), image tokens dropped are reactivated in the next pruning stage for re-evaluation of token relevance as well as in the mask generation. The contributions of the work are summarized as follows:

- We propose a token pruning strategy aimed at TOS with ViT-based models. By leveraging guidance from MLLM, our approach removes image tokens with low task relevance from later computation, thereby forcing the model to concentrate on important image tokens pertinent to the task.
- To the best of our knowledge, this is the first ViT token pruning method that takes MLLM guidance into consideration. Unlike previous works that focus on tradi-

tional vision-only tasks (such as image classification and segmentation), this work is also the first to investigate ViT pruning in vision-language reasoning tasks.

- Experiments on visual reasoning tasks show that VLTP reduces the computational costs of ViT by approximately 25% without any performance degradation and by around 40% with only a 1% performance drop. Guided by MLLM, our VLTP-integrated segmentation model outperforms the SOTA methods, achieving a +2.5% improvement in mIoU.

## 2. Related Works

### 2.1. Image Segmentation

Image segmentation is the process of dividing an image into distinct regions or segments, each representing different objects or areas of interest. Traditional segmentation tasks include instance segmentation [2, 11, 46, 47] and semantic segmentation [3, 7, 35, 42, 48]. Instance segmentation identifies and separates individual objects in an image, assigning a unique label to each object instance. Semantic segmentation, on the other hand, classifies every pixel, partitioning the image into meaningful regions based on visual characteristics. These traditional segmentation tasks are typically image-based, where pixels are categorized into predefined classes, as shown in Fig. 1(a). More recently, segmentation tasks that incorporate both image and text inputs have emerged, such as RES [32] and affordance detection [8]. These tasks involve identifying object masks using a specified text query. Building on these approaches, more complex task-driven [40] and intention-driven [38] segmentation tasks have been introduced, focusing on locating appropriate items to accomplish a given task (or intention, which is a more abstract task). An example of Task-Oriented Segmentation (TOS) is shown in Fig. 1(b). These novel TOS tasks require models to perform not only visual recognition and scene comprehension but also reasoning.

### 2.2. Vision Transformer for Segmentation

ViT has become a popular backbone network for segmentation tasks [42, 48, 58]. Compared to CNN-based segmentation models [18], the global attention mechanism [45] enables ViTs to better understand complex scenes and achieve higher segmentation accuracy, in scenarios with intricate object relationships. The ViT-backbone also facilitates the integration of text embedding information, allowing the model to leverage additional semantic context for more accurate and complex segmentation task [22, 28, 55]. However, the computational cost of ViT limits its deployment, particularly in high-resolution image segmentation tasks. The need for deeper networks and numerous image tokens makes model compression essential.

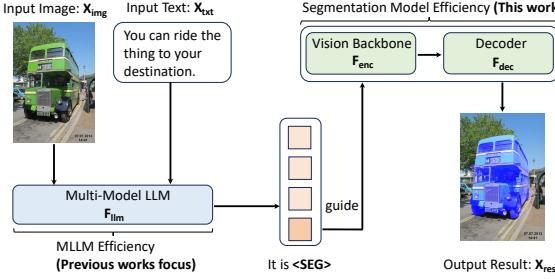


Figure 2. Multi-model LLM (MLLM) guide segmentation model for TOS.

(a) Comparison with Previous Works	
Previous works	1. MLLM quantization, Pruning, and KG distillation. 2. ViT pruning without LLM guidance.
This work	ViT pruning under LLM guidance.
(b) Latency Breakdown (single image and single text)	
MLLM (LLaVA 7B)	Segment Anything (SAM) ViT-H   Mask Decoder
71.0 ms	97.7 ms   8.9 ms
40%	55%   5%

Table 1. (a) Comparison with previous works. (b) MLLM guide segmentation model latency breakdown. We conduct the benchmark using an RTX A6000 GPU with the RIO dataset [38].

### 2.3. Token Pruning

Token pruning has been shown to be an effective mechanism for accelerating ViT [23, 30, 34, 39]. The key idea behind token pruning is the sparsification of image tokens. Since ViT’s attention computation latency and memory requirements are quadratic to the sequence length of the image tokens [20], reducing the number of tokens involved in computation has proven to be an effective method for accelerating ViT. Existing ViT token pruning techniques have been successful in image classification [36, 52] and traditional image segmentation tasks [30, 31, 44], but they are difficult to apply in TOS tasks. As depicted in Fig. 1(c), which highlights the complexity in TOS tasks, the label assigned to a pixel varies depending on the task. That is, the importance of pixels is ambiguous without extra information. Therefore, a new pruning mechanism is needed.

### 3. Task-oriented Segmentation

**Problem Definition:** The TOS task is formulated as follows: with an input image  $X_{img}$  and a task  $X_{txt}$ , we need to find the segmentation mask  $X_{res}$  for all suitable affordances.

$$X_{res} = \mathbf{F}(X_{img}, X_{txt}) \quad (1)$$

Here  $\mathbf{F}$  refers to the model to solve the task. As shown in Fig. 2, to solve the task “*You can ride the thing to your destination*”, we need to find the mask of the object “bus”. In this paper, we propose to leverage a nature

paradigm to solve this task. We first query an MLLM, such as LLaVA [28], to perform reasoning over the image based on the task and generate a special token  $\langle SEG \rangle$ .

$$\langle SEG \rangle = \mathbf{F}_{LLM}(X_{img}, X_{txt}) \quad (2)$$

Next, the special token is treated as a pointer to guide the segmentation model (such as Segment Anything (SAM) [22]) to find the mask.

$$X_{res} = \mathbf{F}_{seg}(X_{img}, \langle SEG \rangle) \quad (3)$$

, where  $\mathbf{F}_{seg}$  is the segmentation model.

**Comparison to Earlier TOS Solutions:** In addressing reasoning-intensive task-oriented object detection and segmentation tasks (e.g., COCO-Tasks [40] and RIO [38] datasets), two main types of solutions exist: single-stage and two-stage frameworks. Single-stage frameworks, such as MDETR [19] and Polyformer [29], directly combine image features and text features and perform end-to-end training of the entire system. The drawback of single-stage frameworks is their lack of logical reasoning capabilities. Consequently, recent works like CoTDet [43] and TaskCLIP [6] introduce two-stage frameworks that combine the reasoning abilities of LLMs with the object detection and segmentation capabilities of vision models. However, these frameworks require LLMs to explicitly generate task-dependent feature words and rationale. Sometimes the LLM provides only a vague indicator, such as a special token embedding, rather than a clear rationale. In contrast to these previous works, our two-step framework seamlessly integrates the robust reasoning abilities of MLLMs with the superior segmentation quality of vision foundation models.

**Nevertheless, a key challenge remains: the computational expense of two-stage framework is extremely high.** Previous research [41, 51, 59] has primarily focused on MLLM’s compression and pruning, largely overlooking the segmentation model, as depicted in Tab. 1(a). The segmentation model, as shown in Tab. 1(b), requires even more computational resources than MLLM. Following this, we will explore ways to enhance the efficiency of the segmentation model under the guidance of MLLM.

### 4. Vision-Language Guided Patch Pruning

As illustrated in Tab. 1(b), the ViT’s computation latency constitutes more than 90% of the total latency of the entire segmentation model. Thus, enhancing the ViT’s efficiency is essential for optimizing the segmentation model’s performance. This work introduces a vision-language guided token pruning method, which accelerates ViT when handling TOS. We detail the mechanism in the following sections.

#### 4.1. Preliminary: Segmentation Model

Refer to Fig. 2, a segmentation model  $\mathbf{F}_{seg}$  typically consists of two components: a vision backbone  $\mathbf{F}_{enc}$  and a

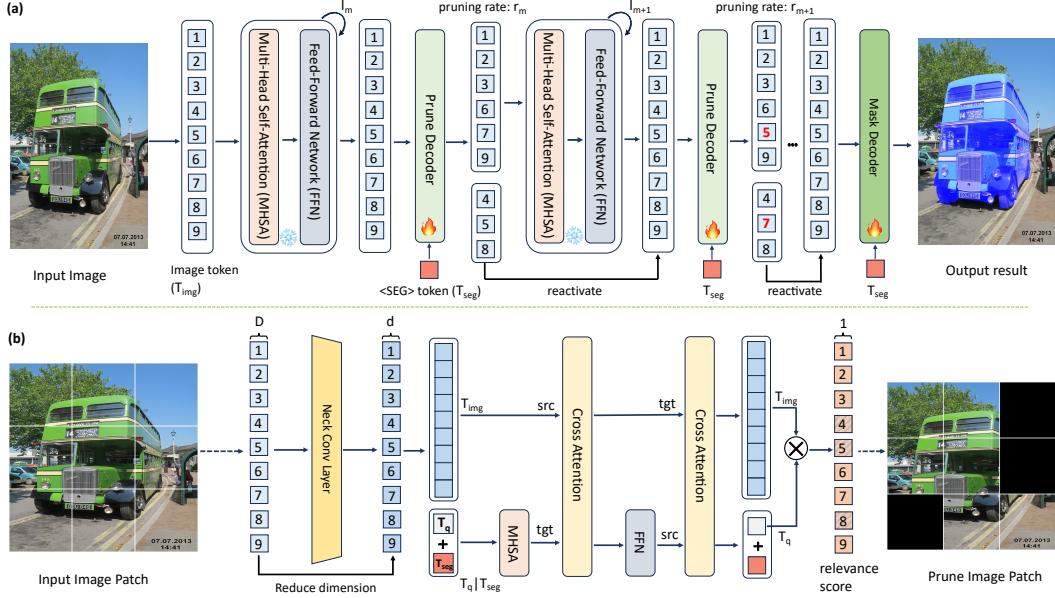


Figure 3. The illustration of vision-language guided token pruning (VLTP) framework. (a) ViT architecture with prune decoder. (b) Prune decoder model architecture. In this illustration, we consider an image with only 9 patch tokens.

mask decoder  $\mathbf{F}_{dec}$ . We illustrate this with SAM. The vision backbone in SAM is a ViT [9, 12]. The ViT extracts features from the input image and passes the feature map to the mask decoder. The mask decoder then predicts the target mask based on the feature map and prompt embeddings. Within our system, the prompt embedding is represented by the  $\langle \text{SEG} \rangle$  token embedding from MLLM, which serves as an ambiguous indicator towards the intended affordance to solve the task. As depicted in Tab. 1(b), given the large latency of vision backbone, the segmentation model will benefit significantly from our proposed VLTP.

A ViT splits an image ( $X_{img} \in R^{3 \times H \times W}$ ) into different patches via patch embedding layer. Suppose the embedding dimension is  $C$ , the patch size is  $P$ , and the image token sequence length is  $N$ , we have:

$$N = \frac{H \times W}{P^2} \quad (4)$$

ViTs are position-agnostic, so we add positional encoding to represent the spatial information of the image token. The resulting image token sequence is denoted as  $T_{img}^0 \in R^{N \times C}$ . We then pass the image token sequence into multiple repeated layers. Every ViT layer contains an MHSA, a feed-forward network (FFN), a layer normalization (LN) [1], and a residual connection [13]. The layers are indexed by  $l \in \{1, 2, \dots, L\}$ , and the output of each layer is marked as  $T_{img}^l$ .

$$\begin{aligned} T_{img}^l &= \text{MHSA}(\text{LN}(T_{img}^{l-1})) + T_{img}^{l-1} \\ T_{img}^l &= \text{FFN}(\text{LN}(T_{img}^l)) + T_{img}^l \end{aligned} \quad (5)$$

Given that the computation in Eq. (5) must be executed by  $L$  times, the computational burden of ViT is exceptionally high, especially when generating high-quality masks. The ViT model is particularly computation-intensive in terms of layer depth and embedding dimension. For example, SAM's ViT-H model comprises 32 layers with an embedding dimension of 1280. Traditional methods to enhance the ViT's efficiency include quantization [26], knowledge distillation [57], and model weight pruning [50]. Nonetheless, in this work, we present a novel orthogonal approach: **Not all image tokens need to traverse all ViT layers; we can eliminate the irrelevant tokens midway.**

## 4.2. Token Pruning Mechanism

To keep the image tokens that are relevant to reasoning tasks, we insert **prune decoder** between different ViT layers. Figure 3(a) presents the ViT model with prune decoder  $\mathbf{F}_{prune}$ . Here we divide the segmentation model's ViT into  $M$  pruning stages and put the prune decoder at the end of each stage. Each stage will have  $l_m$  ViT layers.  $\mathbf{F}_{prune}$  takes the image tokens  $T_{img}^{l_m}$  from the output of the  $l_m$  ViT layers and the  $\langle \text{SEG} \rangle$  token  $T_{seg}$  from the output of MLLM.

$$P_m = \mathbf{F}_{prune}(T_{img}^{l_m}, T_{seg}) \quad (6)$$

The  $P_{mn}$  ( $n \in \{1, 2, \dots, N\}$ ) is the relative score between  $n^{\text{th}}$  image token  $T_{img}^{l_m}[n]$  to the reasoning task. As shown in Fig. 3(a), after generating the relative score of each image token, we pass the image tokens with top  $(1 - r_m) \times N$  relative score to the next ViT layer while freezing the lowest

$r_m \times N$  tokens. Here “freeze” means the rest  $r_m \times N$  do not participate in the self-attention (Eq. (5)) until the next pruning stage. Therefore, after pruning, only  $1 - r_m$  image tokens participate in the computation which significantly reduces the computation of ViT. In Sec. 4.4, we discuss how we reuse those frozen tokens to improve the accuracy.

### 4.3. Prune Decoder Design and Training

Figure 3(b) presents the model architecture of  $\mathbf{F}_{prune}$ . Considering that the embedding dimension of image tokens (1280 for ViT-H) is generally higher than the  $\langle \text{SEG} \rangle$  embedding dimension (256 for LLaVa), we first use a neck convolution layer to reduce the image token embedding from  $D$  to  $d$ , where  $D$  is the ViT image token embedding dimension and  $d$  is the  $\langle \text{SEG} \rangle$  embedding dimension. We then concatenate the  $\langle \text{SEG} \rangle$  embedding vector  $T_{seg}$  with a trainable query token embedding  $T_q$  and perform self-attention of the concatenated vector.

$$T_{cat} = FFN(MHSA(T_q | T_{seg})) \quad (7)$$

The shape of  $T_{cat}$  is  $2 \times d$ , where  $T_{cat}[0]$  is the query token embedding and  $T_{cat}[1]$  is  $\langle \text{SEG} \rangle$  embedding. We then perform two-way cross attention between  $T_{cat}$  and  $T_{img}$ .

$$T'_{cat} = FFN(MHSA(tgt = T_{cat}, src = T_{img})) \quad (8)$$

$$T'_{img} = MHSA(tgt = T_{img}, src = T'_{cat}) \quad (9)$$

The self-attention of  $T_q | T_{seg}$  and cross attention between  $T_{cat}$  and  $T_{img}$  make the query token understand the relation between the image patch tokens and reasoning task. In the end, we perform inner-product between query embedding ( $T_q = T'_{cat}[0]$ ) and image tokens ( $T'_{img}$ ) to get the final image patch relation score:

$$P_s = T'_{img} \times T'^{\top}_{cat}[0] \quad (10)$$

Suppose the pruning rate is  $r_m$ , we get the pruning mask  $P_{mask}$  based on each image patch tokens’ score

$$P_{mask}[i] = \begin{cases} 1 & \text{if } P_s[i] \geq \text{Top}_{r_m}(P_s) \\ 0 & \text{otherwise} \end{cases} \quad (11)$$

Here we use the  $\text{Top}_{r_m}(P_s)$  as the threshold value, which is the  $(r_m \times N)^{\text{th}}$  highest score in  $P$ . As a result, roughly  $(1 - r_m) \times N$  of the image tokens will be retained. During the training stage, we use the Sigmoid function [10] to smooth out Eq. (11). To train the prune decoder, we freeze both MLLM and ViT and only train the prune decoder and mask decoder with the training loss  $L$ :

$$\begin{aligned} L = & \underbrace{CE(X_{res}, X_{gt}) + DICE(X_{res}, X_{gt})}_{L_X} \\ & + \underbrace{CE(P_{mask}, P_{gt}) + DICE(P_{mask}, P_{gt})}_{L_P} \end{aligned} \quad (12)$$

Here  $CE(\cdot, \cdot)$  and  $DICE(\cdot, \cdot)$  represent the cross entropy loss and dice loss functions [17] respectively. The training loss includes two parts: final reasoning segmentation mask loss  $L_X$  and pruning mask loss  $L_P$ . For the  $L_X$ , we follow previous work [22]. We set the ground truth of pruning masks as:

$$P_{gt}(\text{idx}) = \begin{cases} 1 & \text{if } \max(X_{gt}^{\text{patch}}(i, j)) > 0 \\ 0 & \text{otherwise} \end{cases}, \quad (13)$$

$$\text{where } \text{idx} = i \times \left(\frac{W}{P}\right) + j \quad \text{and} \quad \text{idx} \in [1, N] \quad (14)$$

Here the  $(i, j)$  is the patch index of the original ground truth mask  $X_{gt}$ . The ground truth pruning mask label  $P_{gt}$  at a specific index is assigned a value of 1 if the corresponding patch within the final mask includes at least one positive pixel; if not, it is assigned a value of 0.

### 4.4. Pruned Tokens Reactivation

As illustrated in Fig. 3, we reactivate the pruned image tokens in two scenarios. Firstly, the pruned tokens from stage  $m$  will be reactivated at stage  $m + 1$ . This process aids the prune decoder in re-selecting image tokens that are most pertinent to reasoning tasks, thereby correcting any misdrops from earlier stages. Secondly, following the completion of all ViT layers, we will merge the pruned image tokens with the remaining tokens pertinent to reasoning tasks, and feed the combined tokens (with shape  $N \times D$ ) into the mask decoder. This approach enhances the mask decoder’s ability to reconstruct the final segmentation mask.

## 5. Experiments

### 5.1. Dataset and Metrics

To verify the extensive applicability of VLTP for TOS tasks, we perform experiments on two datasets: RIO [38] and COCO-Tasks [40]. Both of these datasets are derived from MS COCO 2014 [27]. Specifically, RIO encompasses over 100 different tasks, and the entire dataset is categorized into two groups based on difficulty: *common* and *uncommon*. The COCO-Tasks dataset comprises 14 distinct tasks. Following the common convention, we utilize the mean intersection over union (mIoU) to assess the segmentation accuracy and the number of floating-point operations in giga (GFLOPs) to measure the model computational intensity. We perform all experiments using three NVIDIA RTX A6000 GPUs.

### 5.2. Vision Language Fintuning

We selected LLaVA 7B [28] to serve as MLLM. LLaVA processes both textual tasks and images as inputs to produce reasoning guidance  $\langle \text{SEG} \rangle$ . For the segmentation model, we

Table 2. Comparison with previous works on RIO [38] dataset.

Model	Common mIoU (%)	Uncommon mIoU (%)
MDETR [19]	44.14	22.03
TOIST [25]	45.07	19.41
Polyformer [29]	48.75	26.77
GROUNDHOG [56]	57.9	33.9
SAM ViT-H	<b>60.4</b>	<b>37.2</b>
SAM ViT-L	55.7	32.4
SAM ViT-B	50.9	27.9

Table 3. Comparison of training scheme. All results are reported based on RIO common dataset. ♦ represents the extra finetune epochs.

Method	GFLOPs	mIoU (%)
Baseline	2976	60.4
+VLTP@Direct (♦0)	2227	39.5
+VLTP@Finetune ( $F_{prune}$ ) (♦5)	2227	50.3
+VLTP@Finetune ( $F_{dec} + F_{prune}$ ) (♦9)	2227	<b>60.1</b>

opted for SAM, a well-known foundation model that enables prompt-guided segmentation. Following previous research [24, 37, 49], we perform end-to-end visual instruction fine-tuning on the entire system. For the LLaVA module, we utilize LoRA [14] for efficient fine-tuning, while for the SAM module, we freeze the ViT backbone and only train the mask decoder. Table 2 shows the reasoning segmentation accuracy of the entire system on the RIO dataset, along with a comparison to earlier studies. Compared to the SOTA method [56], guided by MLLM, the SAM with the ViT-H backbone model, exceeds it by 2.5% on RIO common parts and 3.3% on RIO uncommon parts. This result highlights the importance of integrating MLLM with the segmentation model in visual reasoning tasks.

### 5.3. VLTP Framework Setup

To enhance the computational efficiency of SAM ViT, we integrate the prune decoder at different ViT layers. For SAM ViT-H, there are a total of 32 layers, with every 8 layers containing 7 layers of window attention and 1 layer of global attention. To more effectively model the relational importance between image tokens and reasoning  $\langle \text{SEG} \rangle$ , we insert the prune decoder after the global attention layer. Consequently, for SAM ViT-H, we place the prune decoder at layers 8, 16, and 24. It's important to note that the prune decoder only involves query embedding self-attention and two-way query embedding-to-image patch embedding cross-attention. **Thus, its computation is merely 6 GFLOPs, significantly less than the computation of ViT-H, which exceeds 2900 GFLOPs.** Another crucial point is that there is only a single prune decoder

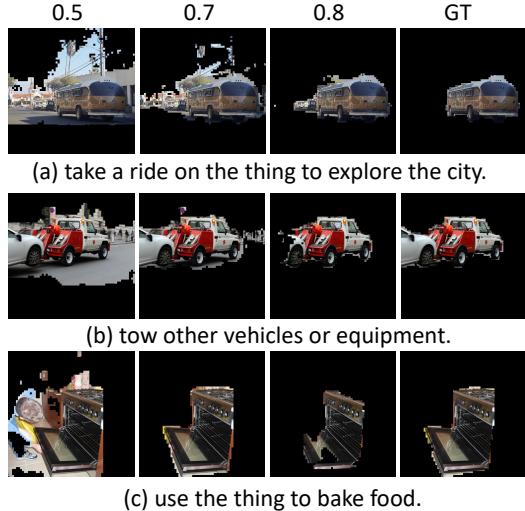


Figure 4. Visualization of VLTP image patch pruning for SAM ViT-H at layers 16 and 24. Three distinct pruning rates (0.5, 0.7, and 0.8) are illustrated alongside the ground truth (GT), task-related image patches.

model within the entire VLTP framework. However, due to the reactivation mechanism, the same prune decoder may be called at multiple layers.

Table 3 demonstrates the impact of different training strategies on final pruning outcomes. We assume a pruning rate of 50%, inserting the prune decoder at layers 16 and 24. This implies that starting from layer 16, only 50% of image tokens will be involved in the SAM ViT computation. The results indicate that when both  $F_{prune}$  and  $F_{dec}$  are fine-tuned simultaneously, VLTP reduces the SAM ViT's computation from 2976 GFLOPs to 2227 GFLOPs, yielding approximately a 25% reduction in GFLOPs with only a 0.3% drop in mIoU. Consequently, in the subsequent sections, both the prune decoder and mask decoder will be fine-tuned to achieve optimal pruning results.

### 5.4. Exploration of Pruning Rate and Position

Table 4 illustrates the effects of pruning as the pruning rate  $r_m$  increases. In this study, we set the pruning positions of SAM ViT at layers 16 and 24. Instead of finetuning the pruning and mask decoders from scratch, we base our experiments on a pre-trained model with a 50% pruning rate, as depicted in Tab. 3. Initially, we increase the pruning rate without any finetuning, and the results in Tab. 4 indicate an accuracy drop. Consequently, we subsequently fine-tune both the prune decoder and mask decoder together after adjusting the pruning rate. The experiment demonstrates that VLTP achieves approximately 35% GFLOPs with just a 0.8% mIoU reduction at a 70% pruning rate, and around 40% GFLOPs with only a 1% mIoU reduction at an 80% pruning rate, compared to the original SAM ViT-H. Additionally, we display the image patch-dropping visual results

Table 4. Pruning ratio exploration for SAM ViT-H. Here we assume the pruning position is at layer 16 and layer 24 and both layer’s pruning ratio is the same. All results are reported based on RIO common dataset.

Pruning Rate (%) Training scheme	50 -	70 zero shot	70 finetune	80 zero shot	80 finetune	90 zero shot	90 finetune
GFLOPs	2227	1930	<b>1930</b>	1782	<b>1782</b>	1636	1636
Common (%)	60.1	57.1	<b>59.6</b>	52.1	<b>59.4</b>	39.76	51.3
Uncommon (%)	37.5	35	<b>37.9</b>	32.2	<b>37.6</b>	19.01	28.4

Table 5. Investigation of SAM ViT-H pruning locations. The table presents the maximum accuracy and the most efficient outcome (measured in GFLOPs) for each pruning position combination. All outcomes are derived from the RIO common dataset. Each invocation of the prune decoder adds an additional 6 GFLOPs.

Position	Pruning Rate	GFLOPS	mIoU (%)
Baseline	{0}	2976	60.4
{8}	{20}	2529	59.8 (-0.6)
{8}	{40}	2083	53.1 (-7.3)
{8, 16}	{20, 40}	2232	58.8 (-1.6)
{8, 16}	{20, 60}	1783	53.7 (-6.7)
{16, 24}	{50, 50}	2227	60.1 (-0.3)
{16, 24}	{80, 80}	1782	58.9 (-1.2)
{8, 16, 24}	{20, 40, 40}	2232	58.9 (-1.5)
{8, 16, 24}	{20, 40, 60}	1853	53 (-7.4)
{8, 16, 24}	{50, 50, 50}	1860	53.3 (-7.1)

Table 6. Ablation study for effect’s of pruned tokens reevaluation. Here the model is SAM ViT-H and dataset is RIO common. Each invocation of the prune decoder adds an additional 6 GFLOPs.

Position	Pruning Rate	GFLOPS	mIoU (%)
{16}	{50}	2221	59.8
{16, 24}	{50, 50}	2227	60.1 (+0.3)
{16}	{70}	1923	58.2
{16, 24}	{70, 70}	1930	59.6 (+1.4)
{16}	{80}	1775	56.8
{16, 24}	{80, 80}	1782	59.4 (+2.6)

Table 7. A comparison of various pruning techniques. Each method is applied on SAM ViT-H. Random 50% indicates randomly dropping 50% of image tokens starting from layer 16.

Method	VLM support	GFLOPs	mIoU
Baseline	No	2976	60.4 (-0)
Random 50%	No	2297	38.2 (-22.2)
CTS [31]	No	2232	35.7 (-24.7)
DToP [44]	No	1892	36.3 (-24.1)
SViT [30]	No	1935	33.5 (-26.9)
<b>VLTP (this work)</b>	<b>Yes</b>	<b>1782</b>	<b>59.4 (-1.0)</b>

in Fig. 4. Under MLLM guidance, VLTP retains the most

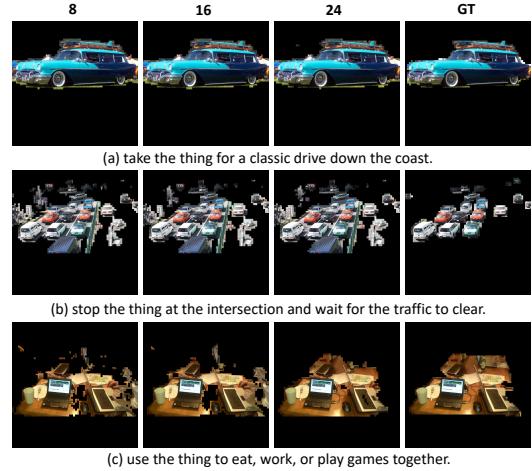


Figure 5. Visualization of VLTP image patch pruning for SAM ViT-H at layers 8, 16, and 24 along with the ground truth (GT). The pruning rate is 0.7.

essential image patches while freezing irrelevant ones, in comparison to the ground truth.

In Tab. 5, we illustrate the influence of different pruning positions on the ultimate pruning accuracy. The findings suggest that initiating pruning from layer 8 leads to a greater accuracy reduction compared to starting from layer 16, assuming the same level of ViT computation reduction. We attribute this to inadequate image feature extraction. Image patches that aren’t directly pertinent to reasoning tasks still contribute to the final mask generation. Consequently, removing these patches at an earlier layer can negatively impact the final mask generation. Figure 5 visualizes the patch dropping at various layer positions.

## 5.5. Ablation Study and Final Results

Table 6 illustrates the impact of reevaluating pruned image tokens. We compare the final reasoning segmentation accuracy between pruning only at layer 16 and pruning at both layer 16 and layer 24. Reassessing token relevance at layer 24 enhances the accuracy of the final reasoning. This is due to the prune decoder potentially making incorrect predictions about the relevance of image tokens to the reasoning task at earlier stages. Hence, it is crucial to reassess the importance of these tokens at deeper ViT layers. In Tab. 7,

Table 8. Main results on two different TOS benchmarks.

Model	pruning rate	RIO common		RIO uncommon		COCO-Tasks	
		mIoU(%)	GFLOPs	mIoU(%)	GFLOPs	mIoU(%)	GFLOPs
SAM ViT-H	-	60.4	2976	37.2	2976	44.6	2976
VLTP@Finetune	0.5	60.1	2227	37.5	2224	44.2	2219
VLTP@Finetune	0.8	59.4	1782	37.6	1774	43.2	1771
SAM ViT-L	-	55.7	1491	32.4	1491	40.1	1491
VLTP@Finetune	0.5	55.1	1118	32.1	1121	39.6	1127
VLTP@Finetune	0.8	54.3	895	31.4	901	38.7	905

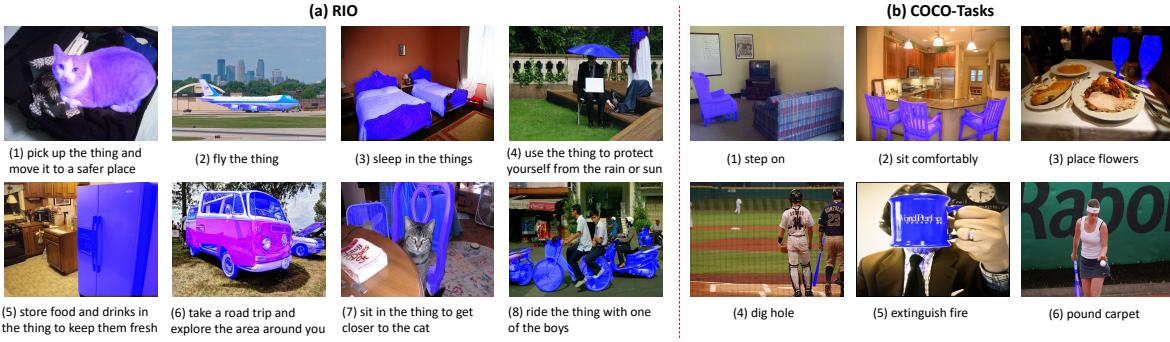


Figure 6. Visualized results. The segmentation results are predicted on RIO (left) and COCO-Tasks (right). The model is SAM ViT-H under the guidance of LLaVA 7B with VLTP@Finetune. We set the pruning position at layers 16 and 24 with the pruning rate as 50%.

we compare various pruning methods for task-oriented reasoning tasks. Although prior research [30,31,44] has shown effectiveness in pruning ViTs for traditional semantic or instance segmentation tasks, these methods are not suitable for TOS as they lack vision-language guidance. For instance, the policy network in CTS [31] is unable to accurately predict whether four image tokens belong to the same semantic class because each token’s class can change depending on the reasoning task. Consequently, the pruning networks proposed in DTOP [44] and SViT [30] are also ineffective as they do not incorporate reasoning guidance. VLTP addresses these limitations by integrating both ViT image tokens and MLLM reasoning guidance, thus achieving superior performance. In Tab. 8, we provide a summary of the results when applying VLTP to SAM ViT-H and SAM ViT-L for the RIO common, RIO uncommon, and COCO-Tasks datasets. VLTP effectively maintains task-relevant image tokens while significantly reducing SAM ViT’s computational costs. On average, VLTP reduces SAM ViT’s computation by 40% with only a 1% loss in mIoU. Figure 6 showcases our framework’s segmentation visualized results on the RIO and COCO-Task datasets.

## 6. Conclusion

In this work, we present a technique for vision language-guided token pruning in ViTs dedicated to TOS. Recogniz-

ing that not every image token needs to be processed by all layers of the ViT, we introduce a novel prune decoder to estimate the relevance of each image patch token to the reasoning task. After evaluation, we retain those image tokens with high relevance to the task and skip the rest. To mitigate incorrect predictions in the initial ViT layers, we also propose a reactivation mechanism that reassesses the relevance of all image tokens to the task. Comprehensive experimental results indicate that our method enhances model efficiency by significantly reducing the computational load of ViTs while producing state-of-the-art segmentation results.

## Acknowledgements

This work was supported in part by the DARPA Young Faculty Award, the National Science Foundation (NSF) under Grants #2127780, #2319198, #2321840, #2312517, and #2235472, the Semiconductor Research Corporation (SRC), the Office of Naval Research through the Young Investigator Program Award, and Grants #N00014-21-1-2225 and N00014-22-1-2067. Additionally, support was provided by the Air Force Office of Scientific Research under Award #FA9550-22-1-0253, along with generous gifts from Xilinx and Cisco.

## References

- [1] JL Ba. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016. 4
- [2] Daniel Bolya, Chong Zhou, Fanyi Xiao, and Yong Jae Lee. Yolact: Real-time instance segmentation. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9157–9166, 2019. 2
- [3] Holger Caesar, Jasper Uijlings, and Vittorio Ferrari. Coco-stuff: Thing and stuff classes in context. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1209–1218, 2018. 2
- [4] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9650–9660, 2021. 1
- [5] Feiyang Chen, Ziqian Luo, Lisang Zhou, Xuetong Pan, and Ying Jiang. Comprehensive survey of model compression and speed up for vision transformers. *arXiv preprint arXiv:2404.10407*, 2024. 1
- [6] Hanning Chen, Wenjun Huang, Yang Ni, Sanggeon Yun, Fei Wen, Hugo Latapie, and Mohsen Imani. Taskclip: Extend large vision-language model for task oriented object detection. *arXiv preprint arXiv:2403.08108*, 2024. 3
- [7] Bowen Cheng, Alex Schwing, and Alexander Kirillov. Per-pixel classification is not all you need for semantic segmentation. *Advances in neural information processing systems*, 34:17864–17875, 2021. 2
- [8] Claudia Cuttano, Gabriele Rosi, Gabriele Trivigno, and Giuseppe Averta. What does clip know about peeling a banana? In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2238–2247, 2024. 2
- [9] Alexey DOSOVITSKIY. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. 1, 4
- [10] Jun Han and Claudio Moraga. The influence of the sigmoid function parameters on the speed of backpropagation learning. In *International workshop on artificial neural networks*, pages 195–201. Springer, 1995. 5
- [11] Junjie He, Pengyu Li, Yifeng Geng, and Xuansong Xie. Fastinst: A simple query-based model for real-time instance segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 23663–23672, 2023. 2
- [12] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 16000–16009, 2022. 4
- [13] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 1, 4
- [14] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021. 6
- [15] Wenjun Huang, Yang Ni, Arghavan Rezvani, SungHeon Jeong, Hanning Chen, Yezi Liu, Fei Wen, and Mohsen Imani. Recoverable anonymization for pose estimation: A privacy-enhancing approach. *arXiv preprint arXiv:2409.02715*, 2024. 1
- [16] Wenjun Huang, Arghavan Rezvani, Hanning Chen, Yang Ni, Sanggeon Yun, Sungheon Jeong, Guangyi Zhang, and Mohsen Imani. Intelligent sensing framework: Near-sensor machine learning for efficient data transmission. *IEEE Sensors Journal*, 2024. 1
- [17] Shruti Jadon. A survey of loss functions for semantic segmentation. In *2020 IEEE conference on computational intelligence in bioinformatics and computational biology (CIBCB)*, pages 1–7. IEEE, 2020. 5
- [18] Du Jiang, Gongfa Li, Chong Tan, Li Huang, Ying Sun, and Jianyi Kong. Semantic segmentation for multiscale target based on object recognition using the improved faster-rcnn model. *Future Generation Computer Systems*, 123:94–104, 2021. 2
- [19] Aishwarya Kamath, Mannat Singh, Yann LeCun, Gabriel Synnaeve, Ishan Misra, and Nicolas Carion. Mdetr-modulated detection for end-to-end multi-modal understanding. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 1780–1790, 2021. 3, 6
- [20] Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and François Fleuret. Transformers are rnns: Fast autoregressive transformers with linear attention. In *International conference on machine learning*, pages 5156–5165. PMLR, 2020. 3
- [21] Salman Khan, Muzammal Naseer, Munawar Hayat, Syed Waqas Zamir, Fahad Shahbaz Khan, and Mubarak Shah. Transformers in vision: A survey. *ACM computing surveys (CSUR)*, 54(10s):1–41, 2022. 1
- [22] 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. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4015–4026, 2023. 1, 2, 3, 5
- [23] Zhenglun Kong, Peiyan Dong, Xiaolong Ma, Xin Meng, Wei Niu, Mengshu Sun, Xuan Shen, Geng Yuan, Bin Ren, Hao Tang, et al. Spvit: Enabling faster vision transformers via latency-aware soft token pruning. In *European conference on computer vision*, pages 620–640. Springer, 2022. 3
- [24] Xin Lai, Zhuotao Tian, Yukang Chen, Yanwei Li, Yuhui Yuan, Shu Liu, and Jiaya Jia. Lisa: Reasoning segmentation via large language model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9579–9589, 2024. 2, 6
- [25] Pengfei Li, Beiwen Tian, Yongliang Shi, Xiaoxue Chen, Hao Zhao, Guyue Zhou, and Ya-Qin Zhang. Toist: Task oriented instance segmentation transformer with noun-pronoun distillation. *Advances in Neural Information Processing Systems*, 35:17597–17611, 2022. 6
- [26] Zhikai Li and Qingyi Gu. I-vit: Integer-only quantization for efficient vision transformer inference. In *Proceedings*

- of the IEEE/CVF International Conference on Computer Vision*, pages 17065–17075, 2023. 4
- [27] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part V 13*, pages 740–755. Springer, 2014. 5
- [28] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in neural information processing systems*, 36, 2024. 2, 3, 5
- [29] Jiang Liu, Hui Ding, Zhaowei Cai, Yuting Zhang, Ravi Kumar Satzoda, Vijay Mahadevan, and R Manmatha. Polyformer: Referring image segmentation as sequential polygon generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18653–18663, 2023. 3, 6
- [30] Yifei Liu, Mathias Gehrig, Nico Messikommer, Marco Caninic, and Davide Scaramuzza. Revisiting token pruning for object detection and instance segmentation. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 2658–2668, 2024. 3, 7, 8
- [31] Chenyang Lu, Daan de Geus, and Gijs Dubbelman. Content-aware token sharing for efficient semantic segmentation with vision transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 23631–23640, 2023. 1, 2, 3, 7, 8
- [32] Jiasen Lu, Vedanuj Goswami, Marcus Rohrbach, Devi Parikh, and Stefan Lee. 12-in-1: Multi-task vision and language representation learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10437–10446, 2020. 2
- [33] Chuoan Ma, Yi Jiang, Jiannan Wu, Zehuan Yuan, and Xiaojuan Qi. Groma: Localized visual tokenization for grounding multimodal large language models. *arXiv preprint arXiv:2404.13013*, 2024. 2
- [34] Tanvir Mahmud, Burhaneddin Yaman, Chun-Hao Liu, and Diana Marculescu. Papr: Training-free one-step patch pruning with lightweight convnets for faster inference. *arXiv preprint arXiv:2403.16020*, 2024. 3
- [35] Roozbeh Mottaghi, Xianjie Chen, Xiaobai Liu, Nam-Gyu Cho, Seong-Whan Lee, Sanja Fidler, Raquel Urtasun, and Alan Yuille. The role of context for object detection and semantic segmentation in the wild. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 891–898, 2014. 2
- [36] Bowen Pan, Rameswar Panda, Yifan Jiang, Zhangyang Wang, Rogerio Feris, and Aude Oliva. Ia-red<sup>2</sup>: Interpretability-aware redundancy reduction for vision transformers. *Advances in Neural Information Processing Systems*, 34:24898–24911, 2021. 3
- [37] Shengyi Qian, Weifeng Chen, Min Bai, Xiong Zhou, Zhiwen Tu, and Li Erran Li. Affordancellm: Grounding affordance from vision language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7587–7597, 2024. 2, 6
- [38] Mengxue Qu, Yu Wu, Wu Liu, Xiaodan Liang, Jingkuan Song, Yao Zhao, and Yunchao Wei. Rio: A benchmark for reasoning intention-oriented objects in open environments. *Advances in Neural Information Processing Systems*, 36, 2024. 2, 3, 5, 6
- [39] Yongming Rao, Wenliang Zhao, Benlin Liu, Jiwen Lu, Jie Zhou, and Cho-Jui Hsieh. Dynamicvit: Efficient vision transformers with dynamic token sparsification. *Advances in neural information processing systems*, 34:13937–13949, 2021. 3
- [40] Johann Sawatzky, Yaser Souri, Christian Grund, and Jurgen Gall. What object should i use?-task driven object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7605–7614, 2019. 2, 3, 5
- [41] Yuzhang Shang, Mu Cai, Bingxin Xu, Yong Jae Lee, and Yan Yan. Llava-prumerge: Adaptive token reduction for efficient large multimodal models. *arXiv preprint arXiv:2403.15388*, 2024. 3
- [42] Robin Strudel, Ricardo Garcia, Ivan Laptev, and Cordelia Schmid. Segmenter: Transformer for semantic segmentation. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 7262–7272, 2021. 2
- [43] Jiajin Tang, Ge Zheng, Jingyi Yu, and Sibei Yang. Cotdet: Affordance knowledge prompting for task driven object detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3068–3078, 2023. 3
- [44] Quan Tang, Bowen Zhang, Jiajun Liu, Fagui Liu, and Yifan Liu. Dynamic token pruning in plain vision transformers for semantic segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 777–786, 2023. 1, 2, 3, 7, 8
- [45] A Vaswani. Attention is all you need. *Advances in Neural Information Processing Systems*, 2017. 1, 2
- [46] Xudong Wang, Rohit Girdhar, Stella X Yu, and Ishan Misra. Cut and learn for unsupervised object detection and instance segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 3124–3134, 2023. 2
- [47] Yuqing Wang, Zhaoliang Xu, Xinlong Wang, Chunhua Shen, Baoshan Cheng, Hao Shen, and Huaxia Xia. End-to-end video instance segmentation with transformers. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8741–8750, 2021. 2
- [48] Enze Xie, Wenhui Wang, Zhiding Yu, Anima Anandkumar, Jose M Alvarez, and Ping Luo. Segformer: Simple and efficient design for semantic segmentation with transformers. *Advances in neural information processing systems*, 34:12077–12090, 2021. 2
- [49] Cilin Yan, Haochen Wang, Shilin Yan, Xiaolong Jiang, Yao Hu, Guoliang Kang, Weidi Xie, and Efstratios Gavves. Visa: Reasoning video object segmentation via large language models. *arXiv preprint arXiv:2407.11325*, 2024. 6
- [50] Huanrui Yang, Hongxu Yin, Maying Shen, Pavlo Molchanov, Hai Li, and Jan Kautz. Global vision transformer pruning with hessian-aware saliency. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 18547–18557, 2023. 4
- [51] Xubing Ye, Yukang Gan, Xiaoke Huang, Yixiao Ge, Ying Shan, and Yansong Tang. Voco-llama: Towards vision

- compression with large language models. *arXiv preprint arXiv:2406.12275*, 2024. 3
- [52] Hongxu Yin, Arash Vahdat, Jose M Alvarez, Arun Mallya, Jan Kautz, and Pavlo Molchanov. A-vit: Adaptive tokens for efficient vision transformer. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10809–10818, 2022. 3
- [53] Licheng Yu, Patrick Poirson, Shan Yang, Alexander C Berg, and Tamara L Berg. Modeling context in referring expressions. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part II 14*, pages 69–85. Springer, 2016. 2
- [54] Guangyi Zhang and Wenjun Huang. Advancing circuit transient response macromodeling: From conventional neural networks to siamese-lstm. *IEEE Transactions on Circuits and Systems I: Regular Papers*, 2024. 1
- [55] Haotian Zhang, Pengchuan Zhang, Xiaowei Hu, Yen-Chun Chen, Liunian Li, Xiyang Dai, Lijuan Wang, Lu Yuan, Jenq-Neng Hwang, and Jianfeng Gao. Glipv2: Unifying localization and vision-language understanding. *Advances in Neural Information Processing Systems*, 35:36067–36080, 2022. 2
- [56] Yichi Zhang, Ziqiao Ma, Xiaofeng Gao, Suhaila Shakiah, Qiaozhi Gao, and Joyce Chai. Groundhog: Grounding large language models to holistic segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 14227–14238, 2024. 2, 6
- [57] Zhuoyang Zhang, Han Cai, and Song Han. Efficientvit-sam: Accelerated segment anything model without performance loss. *arXiv preprint arXiv:2402.05008*, 2024. 4
- [58] Sixiao Zheng, Jiachen Lu, Hengshuang Zhao, Xiatian Zhu, Zekun Luo, Yabiao Wang, Yanwei Fu, Jianfeng Feng, Tao Xiang, Philip HS Torr, et al. Rethinking semantic segmentation from a sequence-to-sequence perspective with transformers. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6881–6890, 2021. 2
- [59] Baichuan Zhou, Ying Hu, Xi Weng, Junlong Jia, Jie Luo, Xien Liu, Ji Wu, and Lei Huang. Tinyllava: A framework of small-scale large multimodal models. *arXiv preprint arXiv:2402.14289*, 2024. 3