

Image-Text Co-Decomposition for Text-Supervised Semantic Segmentation

Ji-Jia Wu¹ Andy Chia-Hao Chang² Chieh-Yu Chuang² Chun-Pei Chen² Yu-Lun Liu²
Min-Hung Chen³ Hou-Ning Hu⁴ Yung-Yu Chuang¹ Yen-Yu Lin²

¹National Taiwan University ²National Yang Ming Chiao Tung University
³NVIDIA ⁴MediaTek

Abstract

This paper addresses text-supervised semantic segmentation, aiming to learn a model capable of segmenting arbitrary visual concepts within images by using only image-text pairs without dense annotations. Existing methods have demonstrated that contrastive learning on image-text pairs effectively aligns visual segments with the meanings of texts. We notice that there is a discrepancy between text alignment and semantic segmentation: A text often consists of multiple semantic concepts, whereas semantic segmentation strives to create semantically homogeneous segments. To address this issue, we propose a novel framework, Image-Text Co-Decomposition (CoDe), where the paired image and text are jointly decomposed into a set of image regions and a set of word segments, respectively, and contrastive learning is developed to enforce region-word alignment. To work with a vision-language model, we present a prompt learning mechanism that derives an extra representation to highlight an image segment or a word segment of interest, with which more effective features can be extracted from that segment. Comprehensive experimental results demonstrate that our method performs favorably against existing text-supervised semantic segmentation methods on six benchmark datasets. The code is available at <https://github.com/072jiajia/image-text-co-decomposition>.

1. Introduction

Semantic segmentation is essential to various applications [11, 15, 50] in computer vision but is hindered by several critical challenges. First, the expensive cost of acquiring pixel-level annotations limits the applicability of fully supervised semantic segmentation methods. Second, most existing methods [40, 43, 52] are developed to work on pre-defined categories and leave themselves inapplicable to rare or unseen concepts described by free-form text. To address these obstacles, a new research direction has emerged in vision-language models, referred to as *text-supervised semantic segmentation* [5, 28, 44–46, 49]. This task devel-

ops segmentation models capable of assigning labels across large vocabularies of concepts and supporting semantic segmentation model training without pixel-wise annotations.

Fig. 1 compares existing methods for text-supervised semantic segmentation by grouping their cross-domain alignment mechanisms into three categories, including *image-text*, *region-text*, and *region-word* alignment. Despite the differences, most of these methods compensate for the lack of pixel-wise annotations on broad semantic concepts by exploring abundant image-text pairs on the internet. The textual descriptions bring extensive knowledge across diverse categories. Thus, existing methods typically apply a vision-language model such as CLIP [34] to textual descriptions to acquire the semantic context of the corresponding images for segmentation model learning.

The image-text alignment is widely adopted in the literature *e.g.* [28, 44, 45]. As depicted in Fig. 1a, methods of this group derive an image encoder and a text encoder by aligning them in a joint embedding space. They then use their proposed zero-shot transfer techniques to enable the two encoders to predict segmentation output. Despite the simplicity, they introduce unfavorable discrepancies between the training and testing phases since we aim to match the semantic features from the text to the corresponding image segments rather than the whole image during testing.

To mitigate this issue, the region-text alignment is explored. As shown in Fig. 1b, methods of this group such as [5] utilize a pre-trained visual-language model to derive an additional image segmenter that discovers concepts described by the text. They enforce the consistency between the segmented region and the text but suffer from the discrepancy between the region-text alignment and semantic segmentation: A text may consist of multiple concepts, such as *pub*, *night*, and *car* in Fig. 1b, while semantic segmentation aims to identify regions of the same concept.

To address the aforementioned issues in the image-text and region-text alignments, we propose a novel framework, Image-Text Co-Decomposition (CoDe), to achieve *region-word alignment*. As illustrated in Fig. 1c, we utilize a

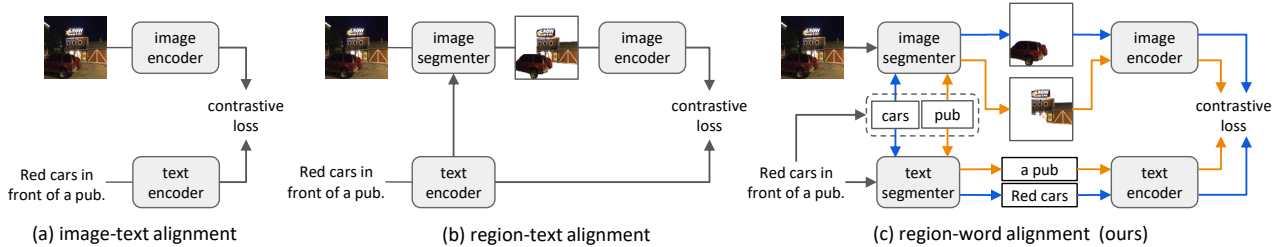


Figure 1. Existing methods perform text-supervised semantic segmentation by learning either (a) image-text alignment or (b) region-text alignment. This paper presents (c) region-word alignment via image-text co-decomposition, where the image and the text are decomposed into object regions and word segments, respectively, while contrastive learning is used to establish cross-modal correspondences between these image and word segments.

visual-language model to construct an image segmenter and a text segmenter: The former decomposes an image into image segments, while the latter decomposes a text into word segments. In addition, there exist one-to-one correspondence between image and word segments. This way, the discrepancy between training and testing is alleviated since each image segment is derived from a single concept given by the corresponding word segment.

The proposed CoDe framework comprises four components: an image segmenter, a text segmenter, a region-word alignment module, and a prompt learning module. We randomly select nouns in the text. For each selected noun *e.g.*, “car”, the image segmenter identifies the image segment that matches the noun, *i.e.*, the region of the car, while the text segmenter discovers the corresponding word segment, *i.e.*, “red cars.” The region-word alignment is developed to enforce the consensus between the image and word segments. To better work with a vision-language model, we present a prompt learning module to derive an extra representation, enabling more effective feature extraction.

The main contributions of this work are as follows:

- We propose a new framework, Image-Text Co-Decomposition (CoDe), to learn the region-word alignment for eliminating train-test and image-text discrepancies, facilitating text-supervised semantic segmentation.
- We propose a prompt learning method to address domain shift issues arising from blank areas during the highlighting process and enhance the alignment between highlighted regions and highlighted words.
- Our method effectively carries out zero-shot semantic segmentation and performs favorably against the state-of-the-art methods on six benchmark datasets.

2. Related Works

2.1. Open-Vocabulary Semantic Segmentation

Open-vocabulary semantic segmentation focuses on segmenting any concepts within images, even those unseen during training, based solely on textual descriptions. Its three important branches are discussed as follows:

Semi-supervised setting with mask-annotations. Methods of this group such as [16, 17, 24, 26, 31, 47] learn from dense annotations to produce high-quality segmentation masks, and then utilize image-text pairs and pre-trained vision-language models to extend the segmentation capability to a larger target vocabulary. Despite the remarkable results, these methods are hindered by their reliance on costly dense annotations, posing a challenge in cases where such annotations are difficult to obtain.

Training-free methods. Another line of research *e.g.* [38, 41, 54] makes the most of large pre-trained models for open-vocabulary segmentation without training. MaskCLIP [54] introduces a modification to the final layer of the CLIP image encoder, yielding dense feature maps that could be employed as initial segmentation maps for further refinement. ReCo [38] constructs an image archive and makes use of retrieval and co-segmentation to identify co-occurrence regions among a specific category. Although these methods eliminate the process of training, the results exhibit significant room for improvement, which shows the need for additional supervision to accomplish this task.

Text-supervised semantic segmentation. It strikes a balance between the two aforementioned branches. Methods of this group are discussed in detail in the following because our method belongs to this group.

2.2. Text-Supervised Semantic Segmentation

Text-supervised semantic segmentation [4, 5, 28, 33, 36, 42, 44–46, 49] decomposes an image into semantic regions according to text descriptions. Unlike semi-supervised methods relying on a few images with mask annotations during training, methods of this group aim to learn semantic masks solely from text-based guidance. We roughly divide existing methods into two categories based on their cross-modal alignment between the image and text domains.

Image-text alignment. These methods train an image encoder alongside a text encoder to align pairs of image and text in a joint embedding space. They use zero-shot transfer to enable the encoders to produce segmentation results. GroupViT [45] introduces a bottom-up approach within

Transformers, grouping image patches into regions and utilizing object semantics derived from texts to guide training. SimSeg [49] further introduces a pretraining method that densely aligns visual and language representations, enabling the trained image encoder to generate segmentation masks in a zero-shot manner.

Region-text alignment. Another line of research targets at aligning the embedding of a region, instead of the whole image, with text descriptions. For instance, TCL [5] learns to segment specific regions within an image while ensuring consistency between the segmented region and the original text. It enables the model to segment the relevant region described in the text.

These methods for text-supervised semantic segmentation have shown that employing vision-language models and contrastive learning on image-text pairs enables aligning visual concepts with the meaning of the whole text. We notice that a text is usually a mix of multiple semantic concepts, but semantic segmentation aims to discover semantically homogeneous segments. To address this issue, inspired by the region-word matching techniques [8, 21, 22, 39] for cross-modal retrieval, we introduce image-text co-decomposition, where the image and the text are decomposed into image and word segments, respectively, and contrastive learning is adopted to enforce cross-modal consensus between these image and word segments. It turns out that image-text co-decomposition results in consistent performance gains on multiple benchmarks.

2.3. Prompt Tuning for Vision-Language Models

Emerged from natural language processing [23, 25, 27], prompt tuning focuses on parameter-efficient adaptation of large pre-trained models to new tasks. In computer vision [20, 55–57], pioneering work such as CoOp [55, 56] incorporates learnable tokens into the CLIP text encoder, enhancing the classification task performance. Recent studies *e.g.* [12, 14, 35] leverage prompt tuning in the text modality for extending CLIP’s capabilities to various applications such as detection and segmentation tasks. Notably, prompt learning methods are also applicable to the visual domain. VPT [19] employs prompt tuning in the visual modality by inserting learnable vectors into Vision Transformers. Further studies [18, 26] explore tuning methods that directly incorporate learnable prompts into the input image within the RGB domain to address downstream tasks.

Drawing inspiration from the success of these methods, our method leverages the capabilities of prompt tuning on segment feature extraction in both the visual and text domains. Prompt learning is beneficial in this work when applying contrastive learning to the visual and textual features extracted by a vision-language model.

3. Methodology

In this section, we first provide an overview of our method for image-text co-decomposition and define the notations in Sec. 3.1. Then, we specify the three major modules of our method, including 1) the image-text co-segmentation module in Sec. 3.2, 2) the region-word highlighting module in Sec. 3.3, and 3) the region-word alignment module in Sec. 3.4. These modules work harmoniously to address the region-word alignment for text-supervised semantic segmentation and enhance model performances. Finally, implementation details are given in Sec. 3.5.

3.1. Method Overview

Image-text co-decomposition enables text-supervised segmenters to learn region-word consensus when segmenting an image X^v with a paired text X^t . Our method aims to jointly learn an image segmenter F^v and a text segmenter F^t with solely the supervision from a set of K image-text pairs, $D = \{X_k^v, X_k^t\}_{k=1}^K$, where no annotations are given. In addition, we optimize two learnable prompts, including a region prompt P^v and a word prompt P^t , to alleviate the unfavorable effect of blank embeddings caused by applying a vision-language model to masked images or texts for feature extraction.

Fig. 2 illustrates the pipeline of our method, consisting of three modules, including the image-text co-segmentation, region-word highlighting, and region-word alignment modules. For an input image-text pair (X^v, X^t) , we initiate the process by randomly selecting a noun N , *e.g.*, *balloon* in the figure, from the text X^t using the noun selector [2]. This selected noun serves as a query. We take the query N along with the image X^v as input to the image segmenter F^v to generate the region mask M^v showing the estimated object region specified by the query. Similarly, a text segmenter F^t takes the query N and the text X^t as input and estimates the word mask M^t indicating the associated word segment.

Subsequently, we apply the region mask M^v to the image X^v to crop the estimated object region. For the estimated background, *i.e.*, the region outside the mask M^v , we crop the corresponding region from the learned region prompt P^v . The highlighted image H^v is yielded by combining the cropped object and background regions. Similarly, the highlighted text H^t is generated by combining the text X^t inside the mask M^t and the word prompt P^t outside the mask M^t . We extract features from the highlighted image and text by using the image encoder E^v and the text encoder E^t of CLIP [34], respectively. The procedure is repeated for each image-text pair and each selected noun. It follows that the region-word alignment is accomplished by contrastive learning [7]. Four loss functions, including \mathcal{L}_{kg} , $\mathcal{L}_{\text{seg}}^v$, $\mathcal{L}_{\text{seg}}^t$, and \mathcal{L}_{hcl} , are used for network optimization, and will be elaborated in the following.

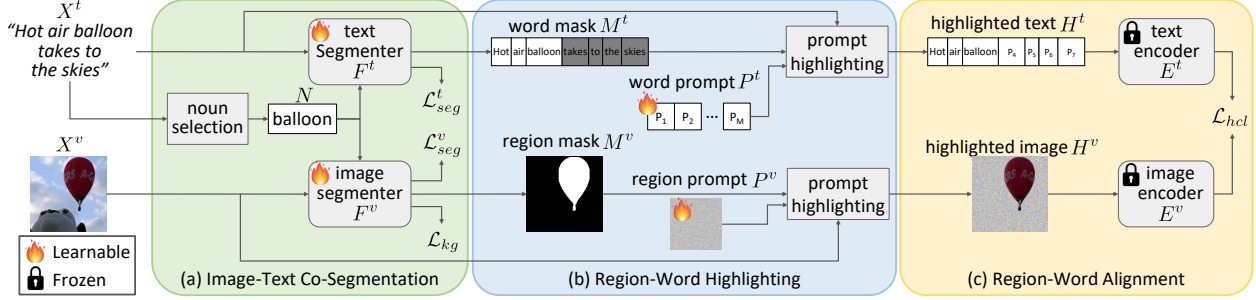


Figure 2. **Training pipeline of our method for image-text co-decomposition.** Our method consists of three major modules, including (a) the image-text co-segmentation module where the image and text segmenters estimate the region and word masks according to a selected noun, respectively, (b) the region-word highlighting module where the estimated masks together with two learnable prompts produce the highlighted image and text, and (c) the region-word alignment module where contrastive learning is applied to the embedded object regions and word segments to accomplish region-word alignment.

3.2. Image-Text Co-Segmentation

The image-text co-segmentation module comprises a noun selector, an image segmenter, and a text segmenter, as shown in Fig. 2a. Taking an image-text pair (X^v, X^t) as input, this module aims at jointly identifying an object region in image X^v and its accompanying word segment in text X^t according to a randomly selected noun.

To begin with, we employ the noun selector [2], which takes the text X^t as input and extracts a set of J nouns, $\{N_j\}_{j=1}^J$, in X^t . For each noun N_j , we carry out *region mask generation*, where the image segmenter F^v predicts a region mask M^v specifying the area in image X^v relevant to noun N_j . A similar task *word mask generation* is performed by the text segmenter F^t , which seeks a word mask M^t matching noun N_j . The tasks of region and word mask generation are depicted as follows.

Region mask generation. The image segmenter F^v takes image X^v and noun N_j as input. It encodes the image into a pixel-wise embedding $\mathbf{x}^v \in \mathbb{R}^{H \times W \times C}$, where $H \times W$ is the image resolution and C is the channel dimension. We also compute the noun embedding $\mathbf{n}_j \in \mathbb{R}^C$ for noun N_j . The image segmenter generates a region mask $M^v \in \mathbb{R}^{H \times W}$ by performing the dot product between the noun embedding \mathbf{n}_j and every location of the image embedding \mathbf{x}^v .

In this work, we use the image segmentation model in [5] to serve as the image segmenter F^v , and employ its corresponding loss, denoted by $\mathcal{L}_{\text{seg}}^v$ here, to help derive the image segmenter. This loss considers segment regularization and contrastive learning that can be directly applied to the segmentation results along with the noun embedding. We use the KgCoOp method [48] to obtain the noun embedding \mathbf{n}_j , as it avoids the pitfalls of improper prompt selection. It appends learnable context tokens to the noun, forming pseudo-sentences for optimal prompt tuning. The noun embedding loss \mathcal{L}_{kg} [48] is included to improve the accuracy of these embeddings, *i.e.*,

$$\mathcal{L}_{\text{kg}} = \|\mathbf{n}_j - \mathbf{n}'_j\|_2^2, \quad (1)$$

where the $\mathbf{n}'_j \in \mathbb{R}^C$ represents the knowledge-guided noun embedding generated from hand-crafted prompts such as “a photo of a N_j ” using the text encoder.

Word mask generation. The text segmenter F^t takes the text X^t and the noun N_j as input. For text feature extraction, we consider the CLIP text encoder appended with two learnable multi-head attention layers. With the resultant feature extractor \tilde{E}^t , the word-wise features of text X^t are obtained via $\mathbf{x}^t = \tilde{E}^t(X^t) \in \mathbb{R}^{L \times C}$, where L is the text length, *i.e.*, the number of word tokens. The word-specific logits $\ell_j = [\ell_{j,i}]_{i=1}^L \in \mathbb{R}^L$ for noun N_j are computed via

$$\ell_j = w \cdot \mathbf{x}^t \mathbf{n}_j + b, \quad (2)$$

where w and b are two learnable parameters, and $\mathbf{n}_j \in \mathbb{R}^C$ is the noun embedding.

Since each word in text X^t belongs to either one of the J word segments associated with nouns $\{N_j\}_{j=1}^J$ or none of them, the word mask $M^t = [m_i^t]_{i=1}^L \in \mathbb{R}^L$ for noun N_j is obtained by applying the softmax function over all the J noun-associated segments, *i.e.*,

$$m_i^t = \frac{\exp(\ell_{j,i})}{1 + \sum_{j'=1}^J \exp(\ell_{j',i})}, \quad \text{for } 1 \leq i \leq L, \quad (3)$$

where the additional 1 in the denominator is included for the case where word i does not belong to any noun-associated segments. The word mask M^t for noun N_j is produced.

According to the softmax function defined in Eq. 3, we get the probabilities of word i over $J + 1$ cases, namely belonging to one of the J noun-associated segments or none of them. We compile a pseudo label vector $\mathbf{p} = \{p_i\} \in \{0, 1\}^L$, where p_i takes value 1 if word i belonging to the j th noun-associated segment gets the highest probability, and 0 otherwise. We develop the text segmentation loss $\mathcal{L}_{\text{seg}}^t$, which is the cross-entropy loss on the word mask M^t with respect to the pseudo label vector \mathbf{p} , and can help learn the text segmenter F^t .

3.3. Region-Word Highlighting

We present a prompt learning method to reliably extract features from an image region or a word segment using a vision-language model. Specifically, we propose a region-highlighting prompt learning method and a word-highlighting prompt learning method, as shown in Fig. 2b.

Region highlighting prompt. When the region mask M^v is directly applied to the image X^v via $M^v * X^v$, where $*$ denotes the element-wise multiplication operation, it makes specific regions of the image being zeroed out, resulting in what we refer to as *blank areas*. When a pre-trained vision-language model like CLIP is applied to these areas, the domain distribution may shift due to the introduction of zero tokens, which are unseen in natural images. To mitigate this issue, we introduce a *region highlighting prompt*, which is a learnable, universal image representation, denoted by P^v . This representation is used alongside the original image in the process of feature extraction. The highlighted image is then obtained via

$$H^v = X^v * M^v + P^v * (1 - M^v). \quad (4)$$

In this way, the blank areas are filled with the corresponding areas of the region prompt P^v alleviating the unfavorable effect of domain shift.

Word highlighting prompt. A similar challenge arises in the text domain when applying the word mask M^t to text X^t . The resultant zero tokens in the masked part unintentionally carry meanings of specific words, leading to potential inaccuracies. To mitigate this issue, we introduce a *word highlighting prompt*, represented as a learnable, universal text representation P^t . The highlighted text H^t is obtained by

$$H^t = X^t * M^t + P^t * (1 - M^t). \quad (5)$$

Since the masked part is filled with content from P^t , the risk of including unexpected text meanings is reduced.

3.4. Region-Word Alignment

In the following, we describe how our method achieves region-word alignment. Our objective is to optimize mutual evidence between the highlighted object regions and the highlighted word segments, as illustrated in Fig. 2c.

Contrastive loss on highlighted region-word pairs. To achieve region-word alignment, we compute the highlighted region embedding e^v and highlighted word segment embedding e^t from the highlighted region-word pair by using the image and text encoders of CLIP by

$$e^v = E^v(H^v) \quad \text{and} \quad e^t = E^t(H^t), \quad (6)$$

where E^v and E^t are the CLIP image and text encoders, respectively.

We adopt batch optimization for model training. Each batch has several triplets, each of which is composed of an image, its paired text, and a randomly selected noun from the text. Each triplet yields a region embedding and a word embedding via Eq. 6. Suppose that there are B triplets in this batch. We create a similarity matrix $S = [S_{i,j}] \in \mathbb{R}^{B \times B}$, where $S_{i,j}$ stores the cosine similarity between the i th region embedding e_i^v and the j th word segment embedding e_j^t . We adopt the symmetric version of InfoNCE loss to develop the highlighted region-word pair contrastive loss, which enhances the similarity of related region-word pairs while reducing it for unrelated pairs:

$$\mathcal{L}_{\text{hcl}} = -\frac{1}{2B} \sum_{i=1}^B \log \frac{\exp(S_{i,i}/\tau)}{\sum_{j=1}^B \exp(S_{i,j}/\tau)} - \frac{1}{2B} \sum_{i=1}^B \log \frac{\exp(S_{i,i}/\tau)}{\sum_{j=1}^B \exp(S_{j,i}/\tau)}, \quad (7)$$

where τ is a learnable temperature. Notably, even though nouns may be selected multiple times across image-caption pairs, the corresponding highlighted regions H^v and highlighted texts H^t vary, ensuring the effectiveness of the InfoNCE loss in precise region-word alignment.

Loss functions and optimization. In sum, the proposed network for image-text co-decomposition is optimized using a composite loss that combines the knowledge-guided, image segmentation, text segmentation, and highlighted region-word pair contrastive losses, defined as follows:

$$\mathcal{L} = \lambda_{\text{kg}} \mathcal{L}_{\text{kg}} + \lambda_{\text{seg}}^v \mathcal{L}_{\text{seg}}^v + \lambda_{\text{seg}}^t \mathcal{L}_{\text{seg}}^t + \lambda_{\text{hcl}} \mathcal{L}_{\text{hcl}}. \quad (8)$$

3.5. Implementation Details

We utilize NLTK’s [2] part-of-speech tagging algorithm for noun selection. For image segmentation, we utilize TCL’s image segmenter [5] to generate image masks, and we adopt the training loss in TCL, which relies solely on the image-caption pairs, to yield $\mathcal{L}_{\text{seg}}^v$. For text segmentation, we use a CLIP text encoder appended with two multi-head attention layers as the text segmenter \tilde{E}^t . Our model is trained on the CC3M and CC12M datasets. The resolution of input images is set to 224×224 . For each forward pass of an image-text pair, we randomly select 2 nouns from the text. The loss weights are set as follows: $\lambda_{\text{kg}} = 8.0$, $\lambda_{\text{seg}}^v = 1.0$, $\lambda_{\text{seg}}^t = 1.0$, and $\lambda_{\text{hcl}} = 0.5$ in the experiments. We train the model with a batch size of 64 on four NVIDIA 2080Ti GPUs and with a learning rate of 5×10^{-6} for a total of 50,000 iterations with 15,000 warmup steps and a cosine schedule. AdamW optimizer [29] is used with a weight decay of 0.05. To improve the quality of the predicted mask during the evaluation phase, we adopt the post-processing approach described in TCL [5], which uses pixel-adaptive mask refinement (PAMR) [1] for mask refinement.

Methods	Publication	Dataset	VOC	Context	Object	Stuff	City	ADE	Avg.
GroupViT [45]	CVPR 2022	CC3M+CC12M+YFCC14M	49.5	19.0	24.3	12.6	6.9	8.7	20.2
ViL-Seg [28]	ECCV 2022	CC12M	37.3	18.9	18.1				
ViewCo [36]	ICLR 2023	CC12M+YFCC14M	52.4	23.0	23.5				
OVSegmentor [46]	CVPR 2023	CC12M	53.8	20.4	25.1				
SimSeg [49]	CVPR 2023	CC3M+CC12M	<u>57.4</u>	26.2	29.7				
TCL [5]	CVPR 2023	CC3M+CC12M	55.0	<u>30.4</u>	<u>31.6</u>	<u>22.4</u>	<u>24.0</u>	<u>17.1</u>	<u>30.1</u>
SegCLIP [30]	ICML 2023	CC3M + COCO	52.6	24.7	26.5				
CoCu [44]	NeurIPS 2023	CC3M+CC12M+YFCC14M	51.4	23.6	22.7	15.2	22.1	12.3	24.6
PGSeg [51]	NeurIPS 2023	CC12M+RedCaps12M	53.2	23.8	28.7				
CoDe (Ours)	CVPR 2024	CC3M+CC12M	57.7	30.5	32.3	23.9	28.9	17.7	31.8

Table 1. **Text-supervised semantic segmentation performance comparison in terms of mIoU.** The proposed method is compared with nine SOTA methods on six popular semantic segmentation datasets: PASCAL VOC (VOC), PASCAL Context (Context), COCO-Object (Object), COCO-Stuff (Stuff), Cityscapes (City) and ADE20K (ADE). For each compared method, the dataset column lists its training datasets. Several methods used datasets in addition to CC3M and CC12M, such as YFCC14M, COCO and RedCaps12M. When applicable, we also provide an average mIoU across all six datasets. For each dataset, the best method is indicated by bold fonts, whereas the second best method is underlined.

4. Experiments

4.1. Datasets and Evaluation Settings

We utilize image-text datasets to train our proposed model and perform extensive experiments on six commonly used semantic segmentation benchmarks to validate our method.

Training datasets. We trained our model on two image-text datasets, Conceptual Captions 3M (CC3M) [37] and Conceptual 12M (CC12M) [6] containing 3M and 12M image-text pairs respectively. They have been widely adopted for training text-supervised semantic segmentation methods.

Evaluation datasets. We used six zero-shot semantic segmentation benchmarks to validate the zero-shot transfer capability of our model on categories that were not specifically trained. As in previous work [5], the benchmarks can be categorized into two groups, with and without background classes. Benchmarks with a background generally label areas that do not belong to any predefined categories as “background,” which is usually removed by considering a probability threshold in text-supervised semantic segmentation. For this category, we use the validation split of the following datasets: PASCAL VOC 2012 [13], PASCAL Context [32], and COCO-Object [3]. They each contain 20, 59, and 80 foreground classes, respectively, with an additional background class. For the “without background category,” we evaluated our model with the validation split of COCO-Stuff [3], Cityscapes [10], and ADE20K [53] datasets. Each of them contains 171, 19, and 150 classes, respectively. In this category, all images are fully annotated, which is exceptionally challenging. Using datasets in this category, our model can be tested for its ability to recognize a variety of concepts. We employ mean intersection-over-union (mIoU) as our evaluation metric.

For zero-shot semantic segmentation evaluation, we rely solely on the image segmenter. The image segmenter pro-

cesses the input image in conjunction with class names from each dataset to produce segmentation predictions. In accordance with the settings of prior work [5], we adopt the class names provided by the default version of MMSegmentation [9] and adhere to its post-processing methodology.

4.2. Quantitative Comparisons

We compare the proposed method with nine text-supervised semantic segmentation methods on the six datasets. Tab. 1 reports the mIoU values. The numbers have been taken directly from the original papers. All methods were tested on the three datasets of the “with background class,” but only three methods (GroupViT [45], CoCu [44] and TCL [5]) were tested on the dataset of the “without background class.” For those three methods, we also report their average mIoU values across all six datasets. It is also worth noting that these methods use different combinations of training datasets, as indicated in the dataset column of Tab. 1.

Our method achieves the best performance in all six datasets, while TCL [5] and SimSeg [49] are the runners-up. In terms of average mIoU, our method (CoDe) achieves 31.8 whereas TCL achieves 30.1, resulting in a 5.65% improvement. The result demonstrates the effectiveness of our image-text co-decomposition method in addressing the alignment-level train-test discrepancy that exists in previous methods by directly learning the region-word alignment.

4.3. Qualitative Results

Visual comparison with existing methods. Fig. 3 visually compares the segmentation results of our methods and two runners-up, TCL [5] and SimSeg [49], on the PASCAL VOC, PASCAL Context, and COCO Object datasets.

This figure illustrates the fundamental benefit of our approach, which involves the direct learning of region-word

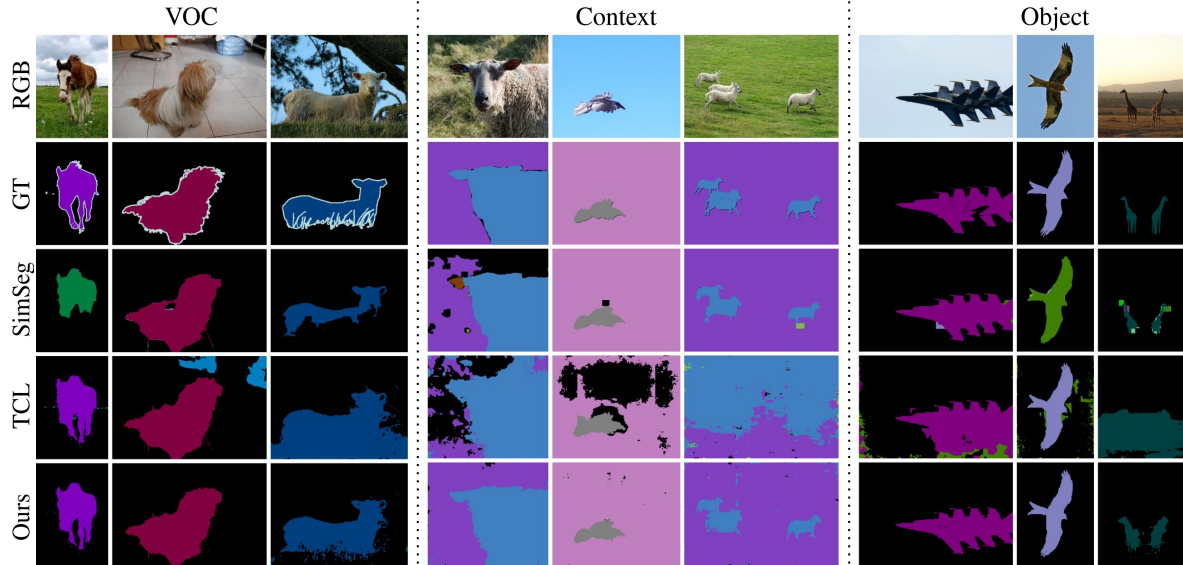


Figure 3. **Qualitative comparisons.** The proposed method is compared with the two most competitive methods, TCL [5] and SimSeg [49], on PASCAL VOC, PASCAL Context, and COCO Object datasets. Our method provides more precise object boundaries and effectively localizes objects within images without misclassification, leading to more accurate segmentation.

alignments. Our model effectively establishes a strong connection between object regions and word segments, allowing a better understanding of how objects are represented in images. Through this enhanced understanding, both segmentation quality and localization capabilities can be improved. As a result, our method provides more accurate classification and more precise masks than other methods.

The SimSeg[49] model, which learns from image-text alignments, occasionally assigns objects to the wrong classes. On the other hand, TCL [5], which is based on region-text alignment, produces coarser semantic masks. Accordingly, these observed limitations are most likely a result of the alignment-level discrepancy between the train and test, which may lead to suboptimal performance.

Visualization of image-text co-segmentation results.

Fig. 4 presents a visualization of the results obtained by our model. We denote regions and word segments associated with the different nouns in the corresponding colors. It demonstrates that our method effectively segments object regions within images based on various input nouns. It simultaneously segments corresponding word segments within the associated text, creating a harmonious alignment between the object region and the word segment.

The region-word alignment plays a pivotal role in our approach, serving as a supervisory signal for the model. By taking advantage of this alignment, our model not only performs visual localization but also captures correlations within the language domain. It indicates that our trained model possesses a more comprehensive understanding of the segmentation task.

C.	W.	R.	VOC	Context	Object	Stuff	City	ADE	Avg.
			54.4	27.6	32.7	22.5	25.0	16.6	29.8
✓			56.2	29.2	32.9	23.3	27.5	17.0	31.0
✓	✓		56.1	29.3	32.6	23.6	29.0	17.3	31.3
✓	✓	✓	57.7	30.5	32.3	23.9	28.9	17.7	31.8

Table 2. **Ablation study.** The baseline model is augmented with the image-text co-decomposition method (C.), the word highlighting prompt (W.), and the region highlighting prompt (R.), one at a time. We report the mIoU values of the resultant models on the six datasets and their averages.

4.4. Ablation Study

Contributions of individual components. The ablation study in Tab. 2 assesses the contribution of the proposed components, including the image-text co-decomposition method, the word highlighting prompt, and the region highlighting prompt. Without the co-decomposition method, our baseline model only trains the image segmenter, resulting in an average mIoU of 29.8. Afterward, each proposed component is added to the baseline model one at a time to verify its contribution. As a result of adding the image-text co-decomposition module alone, the average mIoU has been increased to 31.0. It suggests that the image-text co-decomposition method can achieve region-word alignment and enhance localization capability. The model is further enhanced with the addition of word highlighting prompts and image highlighting prompts, resulting in further performance improvement. It demonstrates that the highlighting prompt learning method enhances feature extraction and strengthens alignment between regions and words.

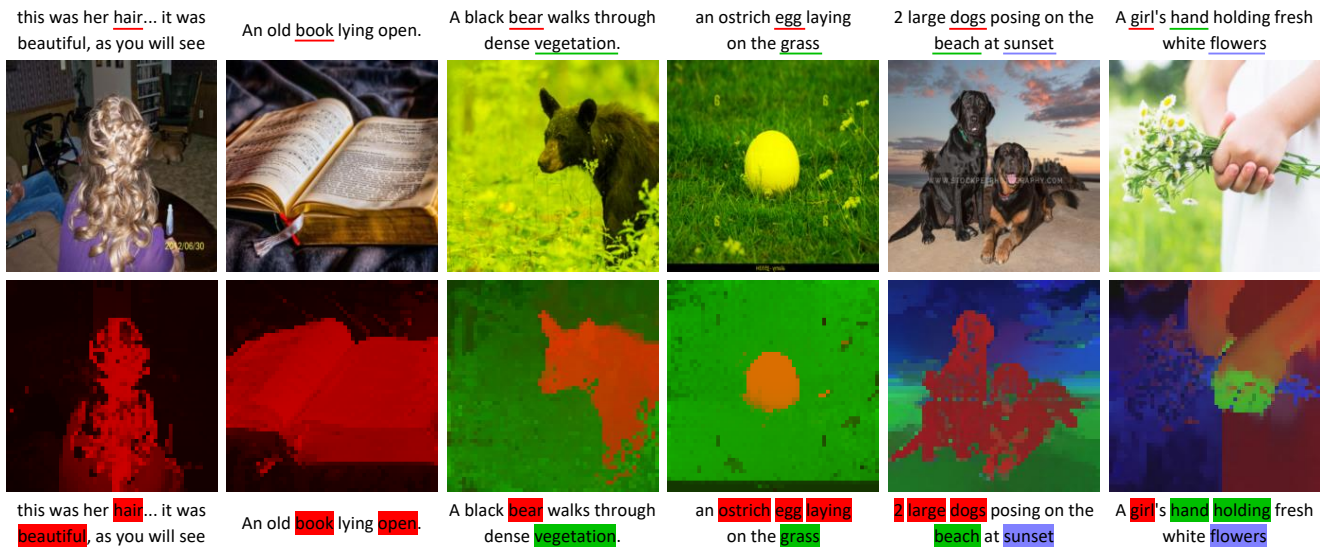


Figure 4. **Visualization of the results of our image-text co-decomposition method.** The first two rows display text and images, representing input image-text pairs. In each text, nouns are underlined with different colors. Our method uses these nouns as queries for performing image-text co-decomposition. Using our image-text co-decomposition method, the last two rows depict the method’s output, where regions and word segments associated with different nouns appear in corresponding colors.

λ_{hcl}	0.05	0.1	0.25	0.5	0.75	1.0
Avg.	30.6	31.2	31.7	31.8	31.5	30.8

Table 3. **Sensitivity analysis on the hyperparameter λ_{hcl} .** By varying λ_{hcl} , we examine the corresponding average mIoU values of all six datasets.

Hyperparameter sensitivity analysis. Tab. 3 investigates the impact of the loss weight for the highlighted region-word pair contrastive loss, denoted as λ_{hcl} in Eq. (8). We observe that, when we apply the highlighted region-word pair contrastive loss in our training phase, the performance consistently outperforms our baseline model. The method is robust to the parameter to some degree as it achieves reasonable performance for a wide range of values. When λ_{hcl} is set to 0.5, our model achieves a peak score of 31.8. It is evident from these results that the image-text co-decomposition method is superior to the image-text decomposition method for achieving region-word alignment.

Effectiveness of jointly decomposing text. We validate the effectiveness of decomposing text by assessing the performance enhancement achieved by generating word masks, as opposed to simply using extracted nouns. This experiment is conducted by modifying the calculation of \mathcal{L}_{hcl} . Instead of using word segment embeddings as mentioned in Sec. 3.4, we opt to compute the similarity matrix S using region embeddings with the embeddings of *individual nouns*. The average mIoU across all benchmarks is

30.2%, which is below our method’s 31.8%. This indicates the benefits of using word segments encompassing extra words associated with each noun. The contextual information encoded in these additional words can serve as valuable supervisory signals, thereby improving performance.

5. Conclusions

We propose Image-Text Co-Decomposition (CoDe) to address cross-domain alignment discrepancies in the existing methods for text-supervised semantic segmentation. First, our method decomposes image-text pairs into corresponding regions and word segments to enforce the region-word alignment. CoDe, underpinned by contrastive learning, alleviates the train-test discrepancy by unifying image-text and region-text alignments to region-word alignment. Then, we introduce a region-highlighting prompt learning method to enhance feature extraction on masked images or texts for precise region-word alignment. Moreover, CoDe surpasses state-of-the-art methods in zero-shot semantic segmentation across six benchmark datasets. This novel approach opens new possibilities for research in vision-language models and their broader applications in computer vision.

6. Acknowledgement

This work was supported in part by the National Science and Technology Council (NSTC) under grants 112-2221-E-A49-090-MY3, 111-2628-E-A49-025-MY3, 112-2634-F-002-005, 112-2634-F-002-006, and 110-2221-E-002-124-MY3, and NTU under grants 112L9009. This work was funded in part by MediaTek and NVIDIA.

References

- [1] Nikita Araslanov and Stefan Roth. Single-stage semantic segmentation from image labels. In *CVPR*, 2020. 5
- [2] Steven Bird, Ewan Klein, and Edward Loper. *Natural language processing with Python: analyzing text with the natural language toolkit*. 2009. 3, 4, 5
- [3] Holger Caesar, Jasper Uijlings, and Vittorio Ferrari. Coco-stuff: Thing and stuff classes in context. In *CVPR*, 2018. 6
- [4] Kaixin Cai, Pengzhen Ren, Yi Zhu, Hang Xu, Jianzhuang Liu, Changlin Li, Guangrun Wang, and Xiaodan Liang. Mixreorg: Cross-modal mixed patch reorganization is a good mask learner for open-world semantic segmentation. In *ICCV*, 2023. 2
- [5] Junbum Cha, Jonghwan Mun, and Byungseok Roh. Learning to generate text-grounded mask for open-world semantic segmentation from only image-text pairs. In *CVPR*, 2023. 1, 2, 3, 4, 5, 6, 7
- [6] Soravit Changpinyo, Piyush Sharma, Nan Ding, and Radu Soricut. Conceptual 12M: Pushing web-scale image-text pre-training to recognize long-tail visual concepts. In *CVPR*, 2021. 6
- [7] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *ICML*, 2020. 3
- [8] Sanghyuk Chun, Seong Joon Oh, Rafael Sampaio De Rezende, Yannis Kalantidis, and Diane Larlus. Probabilistic embeddings for cross-modal retrieval. In *CVPR*, 2021. 3
- [9] MMSegmentation Contributors. MMSegmentation: Openmmlab semantic segmentation toolbox and benchmark, 2020. 6
- [10] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding. In *CVPR*, 2016. 6
- [11] Foivos I Diakogiannis, François Waldner, Peter Caccetta, and Chen Wu. Resunet-a: A deep learning framework for semantic segmentation of remotely sensed data. *ISPRS*, 2020. 1
- [12] Yu Du, Fangyun Wei, Zihe Zhang, Miaojing Shi, Yue Gao, and Guoqi Li. Learning to prompt for open-vocabulary object detection with vision-language model. In *CVPR*, 2022. 3
- [13] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *IJCV*, 2010. 6, 1, 2
- [14] Chengjian Feng, Yujie Zhong, Zequn Jie, Xiangxiang Chu, Haibing Ren, Xiaolin Wei, Weidi Xie, and Lin Ma. Promptdet: Towards open-vocabulary detection using uncurated images. In *ECCV*, 2022. 3
- [15] Di Feng, Christian Haase-Schütz, Lars Rosenbaum, Heinz Hertlein, Claudius Glaeser, Fabian Timm, Werner Wiesbeck, and Klaus Dietmayer. Deep multi-modal object detection and semantic segmentation for autonomous driving: Datasets, methods, and challenges. *IEEE Trans. Intell. Transp. Syst.*, 2020. 1
- [16] Golnaz Ghiasi, Xiuye Gu, Yin Cui, and Tsung-Yi Lin. Scaling open-vocabulary image segmentation with image-level labels. In *ECCV*, 2022. 2
- [17] Cong Han, Yujie Zhong, Dengjie Li, Kai Han, and Lin Ma. Open-vocabulary semantic segmentation with decoupled one-pass network. In *ICCV*, 2023. 2
- [18] Qidong Huang, Xiaoyi Dong, Dongdong Chen, Weiming Zhang, Feifei Wang, Gang Hua, and Nenghai Yu. Diversity-aware meta visual prompting. In *CVPR*, 2023. 3
- [19] Menglin Jia, Luming Tang, Bor-Chun Chen, Claire Cardie, Serge Belongie, Bharath Hariharan, and Ser-Nam Lim. Visual prompt tuning. In *ECCV*, 2022. 3
- [20] Muhammad Uzair Khattak, Hanoona Rasheed, Muhammad Maaz, Salman Khan, and Fahad Shahbaz Khan. Maple: Multi-modal prompt learning. In *CVPR*, 2023. 3
- [21] Dongwon Kim, Namyup Kim, and Suha Kwak. Improving cross-modal retrieval with set of diverse embeddings. In *CVPR*, 2023. 3
- [22] Kuang-Huei Lee, Xi Chen, Gang Hua, Houdong Hu, and Xiaodong He. Stacked cross attention for image-text matching. In *ECCV*, 2018. 3
- [23] Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient prompt tuning. *arXiv preprint arXiv:2104.08691*, 2021. 3
- [24] Boyi Li, Kilian Q Weinberger, Serge Belongie, Vladlen Koltun, and Rene Ranftl. Language-driven semantic segmentation. In *ICLR*, 2022. 2
- [25] Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation. *arXiv preprint arXiv:2101.00190*, 2021. 3
- [26] Feng Liang, Bichen Wu, Xiaoliang Dai, Kunpeng Li, Yanan Zhao, Hang Zhang, Peizhao Zhang, Peter Vajda, and Diana Marculescu. Open-vocabulary semantic segmentation with mask-adapted clip. In *CVPR*, 2023. 2, 3
- [27] Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. *ACM Computing Surveys*, 2023. 3
- [28] Quande Liu, Youpeng Wen, Jianhua Han, Chunjing Xu, Hang Xu, and Xiaodan Liang. Open-world semantic segmentation via contrasting and clustering vision-language embedding. In *ECCV*, 2022. 1, 2, 6
- [29] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017. 5
- [30] Huaishao Luo, Junwei Bao, Youzheng Wu, Xiaodong He, and Tianrui Li. SegCLIP: Patch aggregation with learnable centers for open-vocabulary semantic segmentation. *ICML*, 2023. 6
- [31] Chaofan Ma, Yuhuan Yang, YanFeng Wang, Ya Zhang, and Weidi Xie. Open-vocabulary semantic segmentation with frozen vision-language models. In *BMVC*, 2022. 2
- [32] 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. *CVPR*, 2014. 6

- [33] Prashant Pandey, Mustafa Chasmai, Monish Natarajan, and Brejesh Lall. A language-guided benchmark for weakly supervised open vocabulary semantic segmentation. *arXiv preprint arXiv:2302.14163*, 2023. 2
- [34] 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. 1, 3
- [35] Yongming Rao, Wenliang Zhao, Guangyi Chen, Yansong Tang, Zheng Zhu, Guan Huang, Jie Zhou, and Jiwen Lu. Denseclip: Language-guided dense prediction with context-aware prompting. In *CVPR*, 2022. 3
- [36] Pengzhen Ren, Changlin Li, Hang Xu, Yi Zhu, Guangrun Wang, Jianzhuang Liu, Xiaojun Chang, and Xiaodan Liang. Viewco: Discovering text-supervised segmentation masks via multi-view semantic consistency. *ICLR*, 2023. 2, 6
- [37] Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In *ACL*, 2018. 6
- [38] Gyungin Shin, Weidi Xie, and Samuel Albanie. Reco: Retrieve and co-segment for zero-shot transfer. *NIPS*, 2022. 2
- [39] Yale Song and Mohammad Soleymani. Polysemous visual-semantic embedding for cross-modal retrieval. In *CVPR*, 2019. 3
- [40] Robin Strudel, Ricardo Garcia, Ivan Laptev, and Cordelia Schmid. Segmenter: Transformer for semantic segmentation. In *ICCV*, 2021. 1
- [41] Jinglong Wang, Xiawei Li, Jing Zhang, Qingyuan Xu, Qin Zhou, Qian Yu, Lu Sheng, and Dong Xu. Diffusion model is secretly a training-free open vocabulary semantic segmenter. *arXiv preprint arXiv:2309.02773*, 2023. 2
- [42] Weijia Wu, Yuzhong Zhao, Mike Zheng Shou, Hong Zhou, and Chunhua Shen. Diffumask: Synthesizing images with pixel-level annotations for semantic segmentation using diffusion models. *ICCV*, 2023. 2
- [43] Enze Xie, Wenhai Wang, Zhiding Yu, Anima Anandkumar, Jose M Alvarez, and Ping Luo. Segformer: Simple and efficient design for semantic segmentation with transformers. *NIPS*, 2021. 1
- [44] Yun Xing, Jian Kang, Aoran Xiao, Jiahao Nie, Shao Ling, and Shijian Lu. Rewrite caption semantics: Bridging semantic gaps for language-supervised semantic segmentation. In *NIPS*, 2023. 1, 2, 6
- [45] Jiarui Xu, Shalini De Mello, Sifei Liu, Wonmin Byeon, Thomas Breuel, Jan Kautz, and Xiaolong Wang. Groupvit: Semantic segmentation emerges from text supervision. In *CVPR*, 2022. 1, 2, 6
- [46] Jilan Xu, Junlin Hou, Yuejie Zhang, Rui Feng, Yi Wang, Yu Qiao, and Weidi Xie. Learning open-vocabulary semantic segmentation models from natural language supervision. In *CVPR*, 2023. 1, 2, 6
- [47] Mengde Xu, Zheng Zhang, Fangyun Wei, Han Hu, and Xiang Bai. Side adapter network for open-vocabulary semantic segmentation. In *CVPR*, 2023. 2
- [48] Hantao Yao, Rui Zhang, and Changsheng Xu. Visual-language prompt tuning with knowledge-guided context optimization. In *CVPR*, 2023. 4
- [49] Muyang Yi, Quan Cui, Hao Wu, Cheng Yang, Osamu Yoshie, and Hongtao Lu. A simple framework for text-supervised semantic segmentation. In *CVPR*, 2023. 1, 2, 3, 6, 7
- [50] Xiaohui Yuan, Jianfang Shi, and Lichuan Gu. A review of deep learning methods for semantic segmentation of remote sensing imagery. *Expert Systems with Applications*, 2021. 1
- [51] Fei Zhang, Tianfei Zhou, Boyang Li, Hao He, Chaofan Ma, Tianjiao Zhang, Jiangchao Yao, Ya Zhang, and Yanfeng Wang. Uncovering prototypical knowledge for weakly open-vocabulary semantic segmentation. *NIPS*, 2023. 6
- [52] Sixiao Zheng, Jiachen Lu, Hengshuang Zhao, Xiatian Zhu, Zekun Luo, Yabiao Wang, Yanwei Fu, Jianfeng Feng, Tao Xiang, Philip H.S. Torr, and Li Zhang. Rethinking semantic segmentation from a sequence-to-sequence perspective with transformers. In *CVPR*, 2021. 1
- [53] Bolei Zhou, Hang Zhao, Xavier Puig, Tete Xiao, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Semantic understanding of scenes through the ade20k dataset. *IJCV*, 2019. 6
- [54] Chong Zhou, Chen Change Loy, and Bo Dai. Extract free dense labels from clip. In *ECCV*, 2022. 2
- [55] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Conditional prompt learning for vision-language models. In *CVPR*, 2022. 3
- [56] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Learning to prompt for vision-language models. *IJCV*, 2022. 3
- [57] Beier Zhu, Yulei Niu, Yucheng Han, Yue Wu, and Hanwang Zhang. Prompt-aligned gradient for prompt tuning. In *ICCV*, 2023. 3

Image-Text Co-Decomposition for Text-Supervised Semantic Segmentation

Supplementary Material

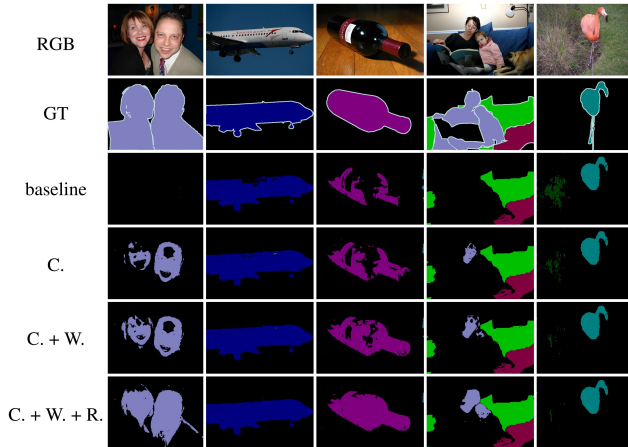


Figure 5. **Ablation studies.** We improve the baseline model by incrementally including (C.) the image-text co-decomposition module, (W.) the word highlighting prompt, and (R.) the region highlighting prompt. We present the segmentation results of the resulting models on the images of the PASCAL VOC [13] dataset.

7. Additional Qualitative Results

7.1. Ablation study visualization

In the following, we conduct ablation studies by visualizing the effects of the proposed components in our method, including the image-text co-decomposition method, the word highlighting prompt, and the region highlighting prompt. To this end, Fig. 5 offers the visual comparison of segmentation results produced by the variants of our method on five images of the PASCAL VOC [13] dataset.

The image-text co-decomposition module equips the model with the region-word alignment ability to localize objects in the images accurately. This module aligns words with corresponding regions in the image, leading to more precise segmentation results. Furthermore, both the word and region highlighting prompts contribute to feature extraction, improving the model’s ability to capture the details of the objects. Hence, the resultant model is more effective in segmenting the whole objects of interest.

7.2. Multi-noun queries

Fig. 6 shows predictions on wild web images with various text queries using the same images and queries selected from Fig. 5 of TCL [5]. Although our method is primarily designed and trained for single-noun queries, the figure demonstrates its effectiveness in processing more complex queries.

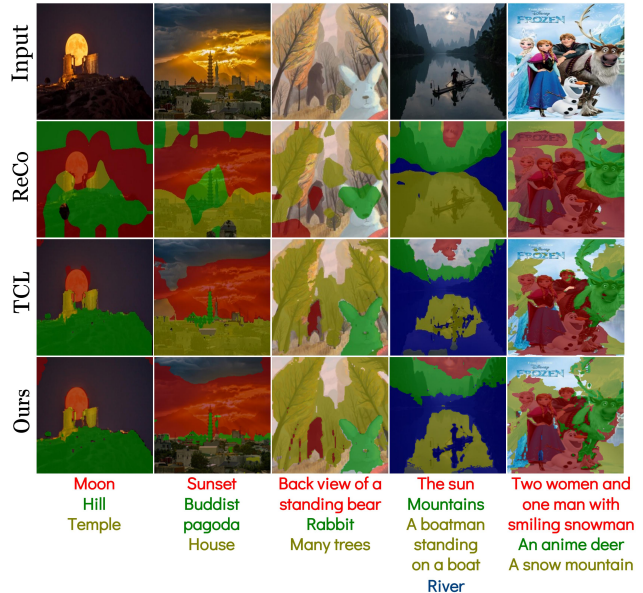


Figure 6. **Examples in the wild.** We show predictions on wild images with free-form text queries. Texts used as target classes are shown at the bottom of the images.

7.3. Failure case visualization

In Fig. 7, we show several failure cases of our method and two competing methods, TCL [5] and SimSeg [49], on the images of the PASCAL VOC [13] dataset.

The first example in Fig. 7a shows a common limitation of existing methods: When segmenting the “person” class, most methods focus on the most distinctive areas, namely the face in this example, and suffer from the variations in the clothes, resulting in the segment that does not cover the entire person. The second example in Fig. 7b depicts a scenario, where unexpected variations are present, *i.e.*, people showing in a television monitor. All three methods segment the outer borders of the monitor. Compared to TCL and SimSeg, our method can further segment the individuals within the monitor. Although the ground truth covers the entire TV monitor, this example validates the effectiveness of our model in localizing the individuals present on the screen.

Fig. 7c, Fig. 7d, and Fig. 7e showcase instances where co-occurrent objects, such as trains and tracks, airplanes and contrails, and boats and water, tend to be segmented together even though they are of different semantic categories. This is a challenge for our method and the two competing methods TCL [5] and SimSeg [49]. These visualization examples emphasize the difficulties of accurate

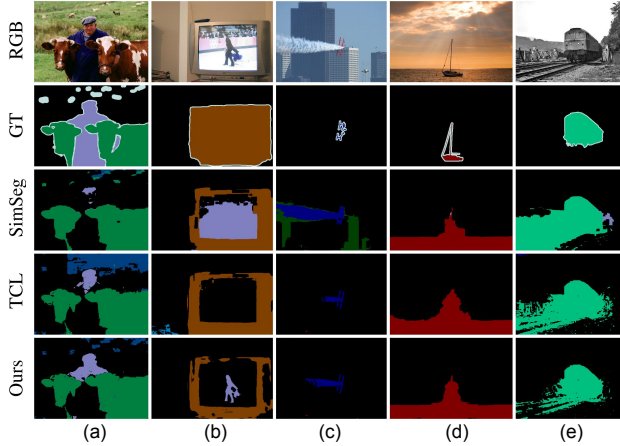


Figure 7. **Failure cases.** The proposed method is compared with the two most competitive methods, TCL [5] and SimSeg [49], on the images of the PASCAL VOC [13] dataset.

segmentation and the challenges in aligning model predictions with ground truth annotations. They provide insights into the limitations of current segmentation approaches and suggest future research directions.

8. More Implementation Details

Training time. On four NVIDIA 2080Ti GPUs, it takes eight hours to train the baseline model with only the image segmenter. On the same devices, it takes twelve hours to train our image-text co-decomposition method, which requires training an additional text segmenter. In light of the improved performance as described above, the longer training period can be justified.