

# GroupViT: Semantic Segmentation Emerges from Text Supervision

Jiarui Xu<sup>1\*</sup> Shalini De Mello<sup>2</sup> Sifei Liu<sup>2</sup> Wonmin Byeon<sup>2</sup>  
 Thomas Breuel<sup>2</sup> Jan Kautz<sup>2</sup> Xiaolong Wang<sup>1</sup>  
<sup>1</sup>UC San Diego      <sup>2</sup>NVIDIA

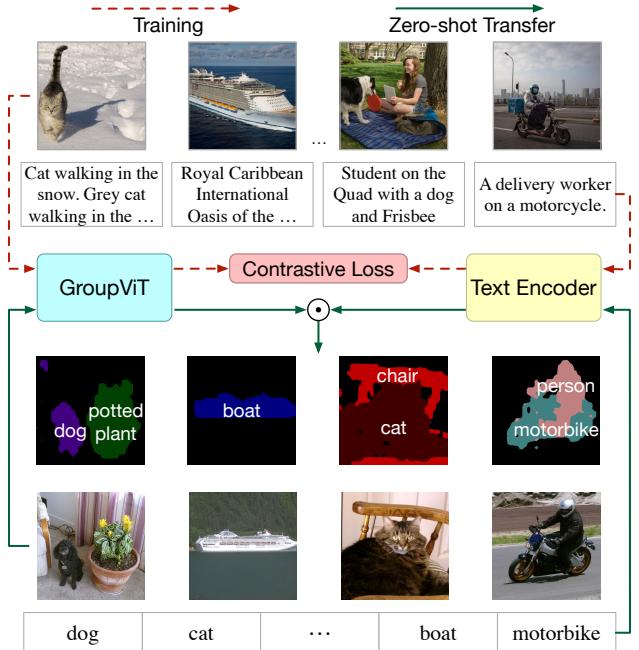
## Abstract

*Grouping and recognition are important components of visual scene understanding, e.g., for object detection and semantic segmentation. With end-to-end deep learning systems, grouping of image regions usually happens implicitly via top-down supervision from pixel-level recognition labels. Instead, in this paper, we propose to bring back the grouping mechanism into deep networks, which allows semantic segments to emerge automatically with only text supervision. We propose a hierarchical Grouping Vision Transformer (GroupViT), which goes beyond the regular grid structure representation and learns to group image regions into progressively larger arbitrary-shaped segments. We train GroupViT jointly with a text encoder on a large-scale image-text dataset via contrastive losses. With only text supervision and without any pixel-level annotations, GroupViT learns to group together semantic regions and successfully transfers to the task of semantic segmentation in a zero-shot manner, i.e., without any further fine-tuning. It achieves a zero-shot accuracy of 52.3% mIoU on the PASCAL VOC 2012 and 22.4% mIoU on PASCAL Context datasets, and performs competitively to state-of-the-art transfer-learning methods requiring greater levels of supervision. We open-source our code at <https://github.com/NVlabs/GroupViT>.*

## 1. Introduction

Visual scenes are naturally composed of semantically-related groups of pixels. The relationship between grouping and recognition has been studied extensively in visual understanding even before the deep learning era [56, 57]. In bottom-up grouping, the idea is to first re-organize pixels into candidate groups and then to process each group with a recognition module. This pipeline has been successfully applied in image segmentation from superpixels [64], constructing region proposals for object detection [80, 102] and semantic segmentation [3]. Beyond bottom-up inference, top-down feedback from recognition can also provide signals to perform better visual grouping [79, 101].

\*Jiarui Xu was an intern at NVIDIA during the project.



**Figure 1. Problem Overview.** First, we jointly train GroupViT and a text encoder using paired image-text data. With GroupViT, meaningful semantic grouping automatically emerges without any mask annotations. Then, we transfer the trained GroupViT model to the task of zero-shot semantic segmentation.

However, on moving to the deep learning era, the ideas of explicit grouping and recognition have been much less separated and more tightly coupled in end-to-end training systems. Semantic segmentation, e.g., is commonly achieved via a Fully Convolutional Network [51], where pixel grouping is only revealed at the output by recognizing each pixel’s label. This approach eliminates the need to perform explicit grouping. While this method is very powerful and still delivers state-of-the-art performance, there are two major limitations that come with it: (i) learning is limited by the high cost of per-pixel human labels; and (ii) the learned model is restricted only to a few labeled categories and cannot generalize to unseen ones.

Recent developments in learning visual representations from text supervision have shown tremendous success on

transferring to downstream tasks [63]. The learned model can not only be transferred to ImageNet classification in a zero-shot manner and achieve state-of-the-art performance, but can also perform recognition on object categories beyond ImageNet. Inspired by this line of research, we ask the question: Can we also learn a semantic segmentation model purely with text supervision, and without any per-pixel annotations, capable of generalizing to different sets of objects categories, or vocabularies, in a zero-shot manner?

To accomplish this, we propose to bring back the grouping mechanism into deep networks, which allows semantic segments to emerge automatically with only text supervision. An overview of our approach is illustrated in Fig. 1. By training on large-scale paired image-text data with contrastive losses, we enable the model to be zero-shot transferred to several semantic segmentation vocabularies, without requiring any further annotation or fine-tuning. Our key idea is to leverage the Vision Transformer (ViT) [24] and incorporate a new visual grouping module into it.

We call our model *GroupViT* (Grouping Vision Transformer). Compared to convolutional neural networks (ConvNets), which operate on regular grids, the global self-attention mechanism of Transformers naturally provides the flexibility to combine visual tokens into non-grid-like segments. Thus, instead of organizing visual tokens into grids, as recent ViT-based applications [17, 25, 48, 86] do, we propose to perform hierarchical grouping of visual tokens into irregular-shaped segments. Specifically, our GroupViT model is organized in different stages through a hierarchy of Transformer layers, where each stage contains multiple Transformers to perform information propagation among the group segments, and a grouping module that merges smaller segments into larger ones. With different input images, our model dynamically forms different visual segments, each intuitively representing a semantic concept.

We train GroupViT with text supervision only. To perform learning, we merge visual segment outputs in the final stage of GroupViT using average pooling. We then compare this image-level embedding to those derived from textual sentences via contrastive learning. We construct positive training pairs by using corresponding image and text pairs, and negative ones by using text from other images. We extract the text embedding with a Transformer model, trained jointly along with GroupViT from scratch. Interestingly, even though we only provide textual training supervision at the image level, we find that semantically meaningful segments automatically emerge using our grouping architecture.

During inference, for the task of semantic segmentation, given an input image, we extract its visual groups using GroupViT (Fig. 1). Each final group’s output represents a segment of the image. Given a vocabulary of label names

for segmentation, we use the text Transformer to extract each label’s textual embedding. To perform semantic segmentation, we then assign the category labels to image segments according to their mutual similarity in the embedding space. In our experiments, we show that GrouViT trained on the Conceptual Caption [11, 68] and Yahoo Flickr Creative Commons [74] datasets with text supervision alone, can transfer to semantic segmentation tasks on the PASCAL VOC [26] and PASCAL Context [58] datasets in a zero-shot manner. Without any fine-tuning, we achieve a mean intersection over union (mIoU) of 52.3% on PASCAL VOC 2012 and an mIoU of 22.4% on PASCAL Context, performing competitively to state-of-the-art transfer-learning methods requiring greater levels of supervision. To our knowledge, our work is the first to perform semantic segmentation on different label vocabularies in a zero-shot manner with text supervision alone, without requiring any pixel-wise labels.

Our contributions are the following:

- Moving beyond regular-shaped image grids in deep networks, we introduce a novel GroupViT architecture to perform hierarchical bottom-up grouping of visual concepts into irregular-shaped groups.
- Without any pixel-level labels and training and with only image-level text supervision using contrastive losses, GroupViT successfully learns to group image regions together and transfers to several semantic segmentation vocabularies in a zero-shot manner.
- To our knowledge, ours is the first work to explore zero-shot transfer from text supervision alone to several semantic segmentation tasks without using *any* pixel-wise labels and establishes a strong baseline for this new task.

## 2. Related Work

**Vision Transformer.** Inspired by the success of Transformers in NLP [22, 81], the Vision Transformer (ViT) [24] was recently proposed and has been successfully applied to multiple computer vision tasks, including image classification [48, 77, 78, 92], object detection [48, 84, 95], semantic segmentation [48, 87, 98] and action recognition [4, 6, 27, 49, 66]. However, much like ConvNets, most variants of ViT still operate on regular image grids. For example, Liu et al. [48] divide the image into regular shaped windows and apply a Transformer block to each one. The convolutional operations are also inserted back into the Transformer block in [17, 25, 86]. While these variants of ViT achieve remarkable performance, they don’t fully leverage the flexibility of the global self-attention mechanism in Transformers. That is, self-attention, by design, can be applied to any arbitrary image segments and is not limited to rectangular-shaped and scan-ordered ones only. Our GroupViT model, on the other hand, leverages this property of Transformers

to learn to group visual information into several arbitrary-shaped segments. With a hierarchical design, it further merges smaller segments into larger ones and yields different semantic groups for each image.

**Representation Learning with Text Supervision.** With large-scale image-text paired data available on the Internet, representation learning with text supervision [15, 20, 33, 35, 40, 42, 53, 63, 96] has been shown to be successful in transferring to various downstream tasks such as visual question answering [2, 100] and visual reasoning [94]. For example, Desai et al. [20] pre-train ConvNets with the image captioning task, and transfer the representation by fine-tuning with downstream task annotations, e.g., object detection labels. Recently, Radford et al. [63] propose to perform contrastive learning between image and text. They show that the learned model can be directly transferred to ImageNet classification [19] in a zero-shot manner without any fine-tuning. Going beyond image classification, our GroupViT model further explores zero-shot transfer to semantic segmentation tasks with only text supervision, which has not been shown in previous work to the best of our knowledge.

**Visual Grounding.** Visual grounding aims to learn image region-text correspondence. One line of research explores a fully supervised approach to detecting text-related bounding boxes within an image [15, 29, 36, 53, 61] using datasets such as Flickr30k Entities [62] and Visual Genome [38]. To scale up learning, weakly-supervised visual grounding has been introduced where the bounding box and text correspondence is not available during training [12, 31, 45, 46, 83, 91]. However, to localize object bounding boxes these approaches still rely on pre-trained object detectors [83, 91], which, in turn, utilize box annotations from other datasets. While related, we emphasize there are two main differences between our problem setting and that of visual grounding: (i) We train our model on millions of noisy image-text pairs from the web, while visual grounding requires human curated and annotated data at a relatively smaller scale; (ii) Our GroupViT provides a bottom-up mechanism for progressive visual grouping where object segments automatically emerge with text supervision, while visual grounding needs bounding box annotations borrowed from other datasets.

**Semantic Segmentation with Less Supervision.** Multiple research directions have been proposed to learn to segment with less supervision than dense per-pixel labels. For example, few-shot learning [23, 47, 54, 59, 75, 82, 90] and active learning [9, 67, 71, 72, 88] are proposed to perform segmentation with as few pixel-wise labels as possible. Going further, zero-shot approaches [7, 41] are proposed to learn segmentation models for unseen categories without using pixel-wise labels for them. However, it still requires learning with segmentation labels on seen categories as the initial step. Another line of re-

lated research is of weakly-supervised semantic segmentation [1, 10, 28, 34, 39, 43, 70, 73, 85], which aims to learn semantic segmentation with only image-level object category supervision. While it largely reduces supervision, it still requires manual labeling using a finite vocabulary on a carefully-curated image dataset. Different from all previous work, our approach completely gets rid of human annotations and GroupViT is trained with large-scale noisy text supervision. Instead of a fixed vocabulary, we show that GroupViT can be generalized to any set of categories in a zero-shot manner for semantic segmentation.

The concurrently developed unpublished text-supervised semantic segmentation methods [30, 89, 93, 99] also show promising results. One major difference between these methods and GroupViT is that, they exploit vision-language model [33, 63] pre-trained on well-prepared large-scale private dataset with 400M-1.8B image-text pairs, while our GroupViT is trained from scratch with much noisier public datasets (30M images in total) to learn grouping and segmentation and yet achieves competitive performance. Among these works, OpenSeg [30] also learns with class agnostic mask annotations to generate mask proposals, while our method does not require any mask annotations.

### 3. Method

We propose the GroupViT architecture for zero-shot transfer to semantic segmentation with text supervision only. GroupViT introduces a new hierarchical grouping Transformer architecture that exploits the global self-attention mechanism of Transformers to partition input images into progressively larger arbitrary-shaped groups. We first describe GroupViT’s architecture in detail in Sec. 3.1. To train it, we employ carefully-designed contrastive losses between image-text pairs, which we describe in Sec. 3.2. Lastly, we transfer the trained GroupViT model, without further fine-tuning, to the task of zero-shot semantic segmentation as described in Sec. 3.3.

#### 3.1. Grouping Vision Transformer

We introduce the GroupViT image encoder (Fig. 2), which performs hierarchical progressive grouping of visual concepts via a Transformer-based architecture. In GroupViT, we separate Transformer layers into multiple grouping stages. In each stage, we learn a number of group tokens (as learnable parameters) via self-attention that aggregate information globally from all image tokens (segments). We then use the learned group tokens to merge similar image tokens together via a *Grouping Block*. Through a hierarchy of grouping stages, we group smaller image segments into larger ones. We describe each component next.

**Architecture** Following the design of ViT [24], we first split an input image into  $N$  non-overlapping patches and

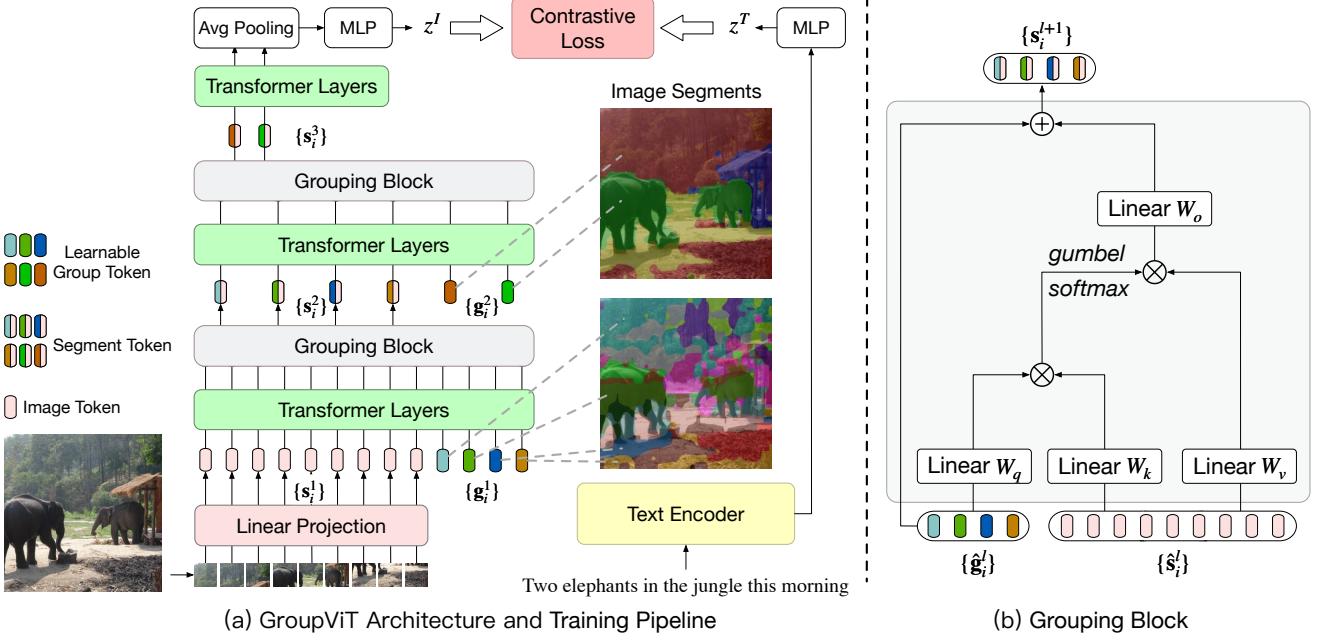


Figure 2. (a) **The Architecture and Training Pipeline of GroupViT.** GroupViT contains a hierarchy of Transformer layers grouped into stages, each operating on progressively larger visual segments. The images on the right show visual segments that emerge in the different grouping stages. The lower stage groups pixels into object parts, e.g., noses and legs of elephants; and the higher stage further merges them into entire objects, e.g., the whole elephant and the background forest. (b) **The Architecture of Grouping Block.** Each grouping stage ends with a grouping block that computes the similarity between the learned group tokens and segment (image) tokens. The assignment is computed via gumbel softmax over group tokens and converted into a one-hot hard assignment. The segment tokens assigned to the same group are merged together and represent new segment tokens that are input to the next grouping stage.

linearly project each into a latent space. We treat each projected patch as an input image token and denote the set of all of them as  $\{\mathbf{p}_i\}_{i=1}^N$ . In each grouping stage, besides the image tokens, we concatenate a set of learnable group tokens and input them into the Transformer for that stage.

**Multi-stage Grouping** As Fig. 2(a) shows, instead of forwarding all the  $N$  input image tokens through all the layers of the Transformer, we separate its layers into a hierarchy of grouping stages. Each stage incorporates a *Grouping Block* at its end to merge the smaller groups into larger ones.

Formally, suppose there are  $L$  grouping stages, each indexed by  $l$  and with a set of learnable group tokens  $\{\mathbf{g}_i\}_{i=1}^{M_l}$ . For simplicity, we treat the image patches  $\{\mathbf{p}_i\}_{i=1}^N$  input to the first grouping stage as the set of starting segments  $\{\mathbf{s}_i^1\}_{i=1}^{M_0}$ , where  $N = M_0$ . We simplify  $\{\mathbf{s}_i^l\}_{i=1}^{M_{l-1}}$  to  $\{\mathbf{s}_i^l\}$  and similarly  $\{\mathbf{g}_i^l\}_{i=1}^{M_l}$  to  $\{\mathbf{g}_i^l\}$ . Starting with  $l=1$ , for each grouping stage, we first concatenate  $\{\mathbf{s}_i^l\}$  and  $\{\mathbf{g}_i^l\}$  together and then input them into a number of Transformer layers, each of which performs information propagation between them via

$$\{\hat{\mathbf{g}}_i^l\}, \{\hat{\mathbf{s}}_i^l\} = \text{Transformer}([\{\mathbf{g}_i^l\}; \{\mathbf{s}_i^l\}]),$$

where  $[ ; ]$  denotes the concatenation operator. Then we group the updated  $M_{l-1}$  image segment tokens  $\{\hat{\mathbf{s}}_i^l\}$  into

$M_l$  new segment tokens  $\{\mathbf{s}_i^{l+1}\}_{i=1}^{l+1}$  via a Grouping Block as

$$\{\mathbf{s}_i^{l+1}\} = \text{GroupingBlock}(\{\hat{\mathbf{g}}_i^l\}, \{\hat{\mathbf{s}}_i^l\}).$$

In each grouping stage  $M_l < M_{l-1}$ , i.e., there are progressively fewer group tokens, resulting in progressively larger and fewer image segments. After the final grouping stage,  $L$ , we apply Transformer layers on all segment tokens and finally average their outputs to obtain the final global image representation  $z^I$  as

$$\{\hat{\mathbf{s}}_i^{L+1}\} = \text{Transformer}(\{\mathbf{s}_i^{L+1}\}), \quad (1)$$

$$z^I = \text{MLP}(\text{AvgPool}(\{\hat{\mathbf{s}}_i^{L+1}\})). \quad (2)$$

As shown in Fig. 2(a), GroupViT re-organizes visual information into arbitrary image segments after the first stage itself and thus is not confined to a regular-grid structure.

**Grouping Block** As shown in Fig. 2(b), the Grouping Block at the end of each grouping stage takes the learned group tokens  $\{\hat{\mathbf{g}}_i^l\}$  and segment tokens  $\{\hat{\mathbf{s}}_i^l\}$  as inputs. It merges all the segment tokens that are assigned to the same group token into a single new image segment, based on similarity in the embedding space.

Formally, we compute the similarity matrix  $\mathbf{A}^l$  between the group tokens  $\{\hat{\mathbf{g}}_i^l\}$  and segment tokens  $\{\hat{\mathbf{s}}_i^l\}$  via a Gumbel-Softmax [32, 55] operation computed over the group tokens as

$$\mathbf{A}_{i,j}^l = \frac{\exp(W_q \hat{\mathbf{g}}_i^l \cdot W_k \hat{\mathbf{s}}_j^l + \gamma_i)}{\sum_{k=1}^{M_l} \exp(W_q \hat{\mathbf{g}}_k^l \cdot W_k \hat{\mathbf{s}}_j^l + \gamma_k)}, \quad (3)$$

where  $W_q$  and  $W_k$  are the weights of the learned linear projections for the group and segment tokens, respectively, and  $\{\gamma_i\}$  are i.i.d random samples drawn from the Gumbel (0, 1) distribution. We compute the group to assign a segment token to by taking the one-hot operation of it  $\text{argmax}$  over all the groups. Since the one-hot assignment operation via  $\text{argmax}$  is not differentiable, we instead use the straight through trick in [60] to compute the assignment matrix as

$$\hat{\mathbf{A}}^l = \text{one-hot}(\mathbf{A}_{\text{argmax}}^l) + \mathbf{A}^l - \text{sg}(\mathbf{A}^l), \quad (4)$$

where  $\text{sg}$  is the stop gradient operator. With the straight through trick,  $\hat{\mathbf{A}}^l$  has the one-hot value of assignment to a single group, but its gradient is equal to the gradient of  $\mathbf{A}^l$ , which makes the Grouping Block differentiable and end-to-end trainable. We call this one-hot assignment strategy as *hard assignment*. After assigning the segment tokens to the different learned groups, we merge the embedding of all the tokens belonging to the same group to form a new segment token  $\mathbf{s}_i^{l+1}$ . For each group, the output of the Grouping Block is a weighted sum of the segment tokens assigned to that group and computed as

$$\mathbf{s}_i^{l+1} = \hat{\mathbf{g}}_i^l + W_o \frac{\sum_{j=1}^{M_{l-1}} \hat{\mathbf{A}}_{i,j}^l W_v \hat{\mathbf{s}}_j^l}{\sum_{j=1}^{M_{l-1}} \hat{\mathbf{A}}_{i,j}^l}, \quad (5)$$

where  $W_v$  and  $W_o$  are the learned weights to project the merged features. An alternative to hard assignment is soft assignment, which uses  $\mathbf{A}^l$  instead of  $\hat{\mathbf{A}}^l$  for computing Eqn. 5. Empirically, we found that hard assignment results in more effective grouping versus soft assignment (Table 1).

The Grouping Block works similarly to a single iteration of the previously proposed Slot Attention mechanism [50]. While Slot Attention learns instance-level grouping from self-supervision, our Grouping Block groups similar semantic regions with weak text supervision. For example, in the second row of Fig. 6, the two horses are grouped together.

### 3.2. Learning from Image-Text Pairs

To train GroupViT to perform hierarchical grouping, we employ carefully-designed contrastive losses between image-text pairs. We describe these next.

**Image-Text Contrastive Loss** To learn visual representations via text supervision, following [33, 63], we train a dual-encoder architecture via an image-text contrastive loss. In our case, GroupViT acts as the image encoder and a Transformer [81] as the text encoder. The final image embedding from GroupViT (Eqn. 2) is the average embedding of all its output segment tokens. The text embedding is the embedding of the last output token (end-of-sentence token) from the text Transformer. We forward the input image and

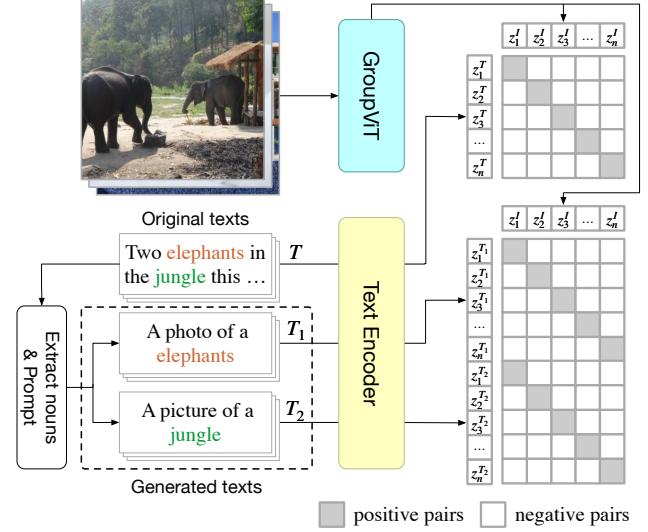


Figure 3. **Multi-label Image-Text Contrastive Loss.** Given an input image-text pair, we generate new text from the original text by extracting its nouns and by prompting them with several sentence templates. For contrastive learning, we treat only matched image and text pairs as positive ones. We train GroupViT and the text encoder to maximize the feature similarity between positive image-text pairs and to minimize it between negative pairs.

text in a pair through their respective encoders, project them into a common embedding space and compute a similarity measure between them. We consider all matched image-text pairs as positive pairs, and all other unmatched ones as negative ones. Our training objective is to pull the representations of the positive pairs closer to each other, while pushing those of the unmatched ones far away from each other.

Formally, assume a batch of  $B$  image-text pairs  $\{(x_i^I, x_i^T)\}_{i=1}^B$ , where  $x_i^I$  and  $x_i^T$  are the image and text inputs, respectively, of the  $i$ -th pair. We encode each of them, via their respective encoders, into embedding vectors  $z_i^I$  and  $z_i^T$  and  $l_2$ -normalize each. We then measure their similarity by computing their dot product. The total image-text contrastive loss is defined as

$$\mathcal{L}_{I \leftrightarrow T} = \mathcal{L}_{I \rightarrow T} + \mathcal{L}_{T \rightarrow I}, \quad (6)$$

which is composed of an image-to-text contrastive loss defined as

$$\mathcal{L}_{I \rightarrow T} = -\frac{1}{B} \sum_{i=1}^B \log \frac{\exp(z_i^I \cdot z_i^T / \tau)}{\sum_{j=1}^B \exp(z_i^I \cdot z_j^T / \tau)},$$

and a text-to-image contrastive loss defined as

$$\mathcal{L}_{T \rightarrow I} = -\frac{1}{B} \sum_{i=1}^B \log \frac{\exp(z_i^T \cdot z_i^I / \tau)}{\sum_{j=1}^B \exp(z_i^T \cdot z_j^I / \tau)},$$

where  $\tau$  is a learnable temperature parameter to scale the logits.

**Multi-Label Image-Text Contrastive Loss** To enable effective visual grouping, besides the image-text loss in

Eqn. 6, we propose a multi-label contrastive loss with text prompting. As illustrated in Fig. 3, we use the “prompting engineering” mechanism proposed in [63] to generate additional text labels for each image besides its originally provided sentence label. Specifically, we randomly select  $K$  noun words from a sentence  $x_i^T$ , and prompt each of them with a set of handcrafted sentence templates, e.g., “A photo of a {noun}”. The motivation to select nouns is that objects in images are more likely to be described by them. Besides training with the original image-text pairs  $\{(x_i^I, x_i^T)\}_{i=1}^B$ , we employ additional contrastive losses between the new sets of image-“prompted text” pairs  $\{(x_i^I, x_i^{T_1})\}_{i=1}^B, \{(x_i^I, x_i^{T_2})\}_{i=1}^B, \dots, \{(x_i^I, x_i^{T_K})\}_{i=1}^B$ , where  $\{x_i^{T_k}\}_{k=1}^K$  are all prompted sentences generated from the nouns sampled from  $x_i^T$ . As shown in Fig. 3, compared to the standard contrastive loss (Eqn. 6), which results in only one positive pair among the batch  $B$ , in our case, each image  $x_i^I$  has  $K$  positive text pairs and  $B(K - 1)$  negative ones.

Similarly to the standard image-text contrastive loss (Eqn. 6), our multi-label contrastive loss is defined as

$$\mathcal{L}_{I \leftrightarrow \{T_k\}_{k=1}^K} = \mathcal{L}_{I \rightarrow \{T_k\}_{k=1}^K} + \mathcal{L}_{\{T_k\}_{k=1}^K \rightarrow I}, \quad (7)$$

which is a sum of two two-way contrastive losses

$$\mathcal{L}_{I \rightarrow \{T_k\}_{k=1}^K} = -\frac{1}{B} \sum_{i=1}^B \log \frac{\sum_{k=1}^K \exp(z_i^I \cdot z_i^{T_k} / \tau)}{\sum_{k=1}^K \sum_{j=1}^B \exp(z_i^I \cdot z_j^{T_k} / \tau)}$$

and

$$\mathcal{L}_{\{T_k\}_{k=1}^K \rightarrow I} = -\frac{1}{KB} \sum_{k=1}^K \sum_{i=1}^B \log \frac{\exp(z_i^{T_k} \cdot z_i^I / \tau)}{\sum_{j=1}^B \exp(z_i^{T_k} \cdot z_j^I / \tau)}.$$

Finally, the total image-text contrastive loss for training GroupViT is defined as

$$\mathcal{L} = \mathcal{L}_{I \leftrightarrow T} + \mathcal{L}_{I \leftrightarrow \{T_k\}_{k=1}^K}. \quad (8)$$

### 3.3. Zero-Shot Transfer to Semantic Segmentation

Since GroupViT automatically groups images into semantically-similar segments, its output can be easily zero-shot transferred to semantic segmentation without any further fine-tuning. This zero-shot transfer pipeline is illustrated in Fig. 4. To infer the segments of an image belonging to a finite vocabulary of object classes, we forward a test image through GroupViT without applying AvgPool to its final  $L$  output segments, and obtain the embedding of each of them as  $\{z_i^I\}_{i=1}^{M_L}$ . Each segment token corresponds to an arbitrarily-shaped region of the input image. We then compute the similarity between the embedding of each segment token and the text embedding of all the semantic classes present in the dataset. We assign each image segment to the semantic class with the highest image-text embedding similarity.

Specifically, let  $\hat{\mathbf{A}}^l$  be the assignment matrix of the  $l$ -th grouping stage described in Sec. 3.1, which indicates

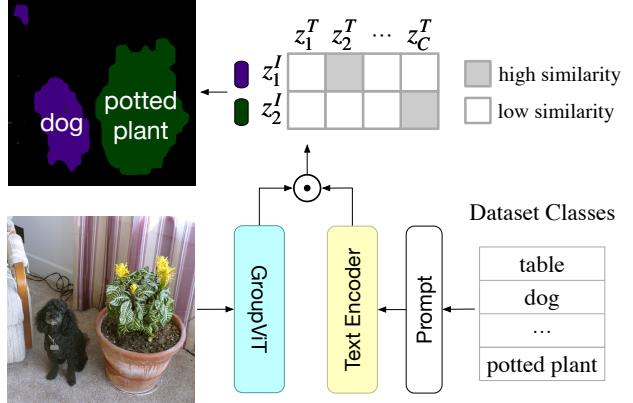


Figure 4. **Zero-Shot Transfer of GroupViT to Semantic Segmentation.** Each output segment’s embedding from GroupViT corresponds to a region of the image. We assign each output segment to the object class with the highest image-text similarity in the embedding space.

the mapping between the input and output segments of  $l$ -th stage. Multiplying all the stage-level assignment matrices  $\prod_{l=L}^1 \hat{\mathbf{A}}^l$  yields the final assignment between the input patches  $\{\mathbf{p}_i\}_{i=1}^N$  and the final-stage output tokens  $\{z_i^I\}_{i=1}^{M_L}$ . We use the same “prompting engineering” as described in Sec. 3.2 to transform all the semantic segmentation label names into sentences. The embedding of label names in the dataset is  $\{z_i^T\}_{i=1}^C$ , where  $C$  is the number of classes. As shown in Fig. 4, to classify an image segment  $z_i^I$  to one of  $C$  classes, we compute the dot product between  $l_2$ -normalized class name embedding vectors  $\{z_i^T\}_{i=1}^C$  and  $z_i^I$ , and assign it to the class with the highest similarity.

## 4. Experiments

### 4.1. Implementation Details

**Architecture** The architecture of GroupViT is based on ViT-S [24, 77] with 12 Transformer layers, each with a hidden dimension of 384. We use input images of size  $224 \times 224$  and a patch size of  $16 \times 16$ . We add a learnable positional embedding to each patch after linearly projecting it. We experiment with 1-stage and 2-stage architectures for GroupViT. Both architectures output 8 tokens after the final grouping stage. In 1-stage GroupViT, we learn 64 group tokens and insert the grouping block after the sixth Transformer layer. Before the grouping block, we project the 64 group tokens into 8 tokens using an MLP-Mixer layer [76] and output 8 segment tokens. In 2-stage GroupViT, there are 64 and 8 group tokens in the first and second grouping stages, respectively. We insert grouping blocks after the sixth and ninth Transformer layers. Our text-encoder is the same as [63]. We use a 2-layer MLP to project the visual and text embedding vectors into the same latent space.

**Training** We use the CC12M [11] and the filtered YFCC [74] datasets for training, containing 12M and 14M image-text pairs, respectively. Our batch size is 4096 with a learning rate initialized to 0.0016 and decayed via the cosine schedule [52]. We use the Adam [37] optimizer with a weight decay of 0.05. We train GroupViT for 30 epochs with the 5 initial epochs containing linear warm-up. For the multi-label contrastive loss, we set  $K = 3$ . We use the same text templates as in [63] for generating text prompts.

**Zero-shot Transfer to Semantic Segmentation** We evaluate GroupViT for the task of zero-shot transfer to semantic segmentation on the validation splits of the PASCAL VOC 2012 [26] and PASCAL Context [58] datasets. They each contain 20 and 59 foreground classes, respectively, with an additional background class. During inference, GroupViT predicts only the foreground classes by thresholding the softmax-normalized-similarity between the embedding of the output image segments and the text segmentation labels, where we set the threshold to 0.9 and 0.5 for PASCAL VOC 2012 and PASCAL Context, respectively. We resize each input image to have a shorter side length of 448.

## 4.2. Ablation Study

To discern the contribution of each component of GroupViT, we conduct an ablation study. For all experiments, we train a 1-stage GoupViT with the CC12M dataset, unless otherwise specified. We report mIoU (mean intersection over union) of the predicted and ground truth segmentation masks on the PASCAL VOC 2012 validation set.

**Hard vs. Soft Assignment** In each Grouping Block, we assign image segment tokens to group tokens via hard or soft assignment (Sec. 3.1). For soft assignment, we use the original  $\mathbf{A}^l$  matrix instead of  $\hat{\mathbf{A}}^l$  used for hard assignment to compute Eqn. 5. The impact of this is shown in the first column of Table 1. We find that hard assignment improves over soft assignment by a large margin, >10% mIoU. We conjecture that with soft assignment, the features of new segment tokens  $\{\mathbf{s}_i^{l+1}\}$  are likely to be more correlated with each other due to absence of zero values in  $\mathbf{A}^l$ . Hence, each group may contain information from the same image patches increasing ambiguity while assigning text labels to image segments. With hard assignment, however, the affinity matrix  $\hat{\mathbf{A}}^l$  assigns image segments to groups in a mutually exclusive manner, making groups more differentiated and their assignment to text labels less ambiguous.

**Multi Label Contrastive Loss** We investigate the effect of adding the multi-label contrastive loss in the second column of Table 1. Adding the multi-label contrastive loss to the standard one (Eqn. 8) improves performance both with hard and soft assignment, by 13.1% and 2.6%, respectively. With the multi-label contrastive loss, the input text during

arch	hard assignment	multi-label loss	mask mIoU
GroupViT			12.0
GroupViT	✓		36.7
GroupViT		✓	25.1
GroupViT	✓	✓	<b>39.3</b>

Table 1. **Ablation results of hard vs. soft assignment and multi-label contrastive loss.**

arch	# group tokens	# output tokens	mask mIoU
GroupViT	16	4	28.6
GroupViT	16	8	37.1
GroupViT	32	8	38.3
GroupViT	64	8	<b>39.3</b>
GroupViT	64	16	38.0

Table 2. **Ablation results of different numbers of group and output tokens.**

arch	dataset	# stages	mask mIoU	boundary mIoU
GroupViT	CC12M	1	39.3	31.6
GroupViT	CC12M	2	41.1	33.5
GroupViT	CC12M+YFCC	1	37.2	32.3
GroupViT	CC12M+YFCC	2	<b>52.3</b>	<b>40.3</b>

Table 3. **Ablation results of single-stage and multi-stage grouping.**

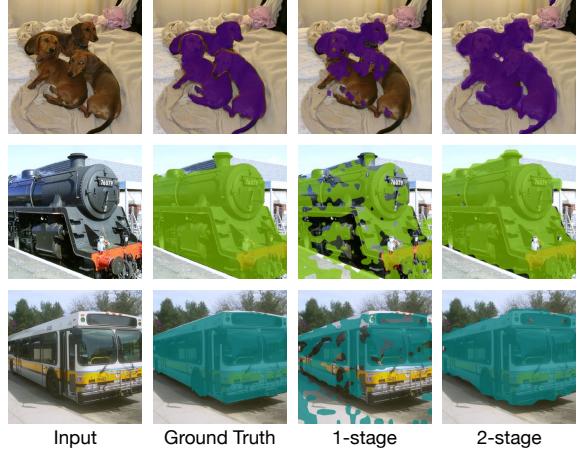


Figure 5. **Visual results of 1-stage and 2-stage GroupViT.** The segmentation maps generated by 2-stage GroupViT are smoother and more accurate than those of 1-stage GroupViT.

training and inference is in a similar prompted format. We conjecture that this consistency helps GroupViT better classify the learned image segments into label categories.

**Group Tokens** In Table 2, we compare different group and output tokens. We observe that increasing group tokens consistently improves performance. Conceptually, each group token represents a distinct semantic concept. So more group tokens presumably help GroupViT learn to group more semantic concepts. Note that although the number of

arch	method	mask mIoU
ViT	pixel-wise	20.1
ViT	K-means	25.0
ViT	Mean-shift [18]	20.7
ViT	Spectral clustering [69]	19.7
GroupViT	Ours	<b>52.3</b>

Table 4. Comparisons with zero-shot baselines.

group tokens is much less than the number of classes in the real world, each group token is a feature vector in a 384-D embedding space, which can represent many more concepts than just 1. We also experiment with different output tokens and find 8 to be optimal, similar to findings in [66].

**Multi Stage Grouping** In Table 3, we compare the 1-stage and 2-stage GroupViT architectures. We also compare their visual zero-shot semantic segmentation results in Fig. 5. We find that the 2-stage GroupViT generates smoother segmentation maps compared to its 1-stage counterpart. To quantify the smoothness of the segmentation maps, we also report the boundary mIoU [16] in Table 3, which computes the IoU of boundaries only. The 2-stage GroupViT improves the mask mIoU of the 1-stage variant by 1.8% and the boundary mIoU by 1.9%. We also train both models on a combination of the CC [11] and YFCC [74] datasets. While the 1-stage model does not benefit as much from the expanded dataset, the 2-stage model improves significantly both in terms of the mask and boundary mIoU values by  $\sim 7\%$ . These results demonstrate that our hierarchical grouping mechanism is effective, especially when training with larger datasets. We adopt the 2-stage GroupViT in the following experiments.

### 4.3. Visualization

**Qualitative Results on PASCAL VOC 2012** We show selected qualitatively segmentation results of GroupViT in Fig. 6. We select examples with a single object (row 1), multiple object of the same class (row 2), and multiple objects from different classes (row 3). GroupViT could generate plausible segmentation. We provide more qualitative results in the supplement Sec. B.

**Concepts Learnt by Group Tokens** We visualize what the group tokens learn in Fig. 7. We select some group tokens and highlight the attention regions across images in the PASCAL VOC 2012. We found different group tokens are learning different semantic concepts. In the first stage, group tokens usually focus on mid-level concepts such as “eyes” (row 1) and “limbs”(row 2). Interestingly, the group token 36 attends to “hands” if people are in the image, while focusing on “feet” if animals like bird and dog are present. Group tokens in the second stage are more associated with high-level concepts, e.g., “grass”, “body” and “face”. The figure also shows that the learnt concepts in the first stage

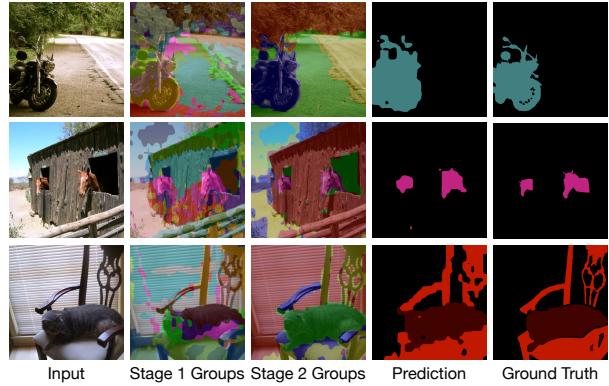


Figure 6. Qualitative results on PASCAL VOC 2012. Stage 1/2 Groups are grouping results prior to assigning labels.

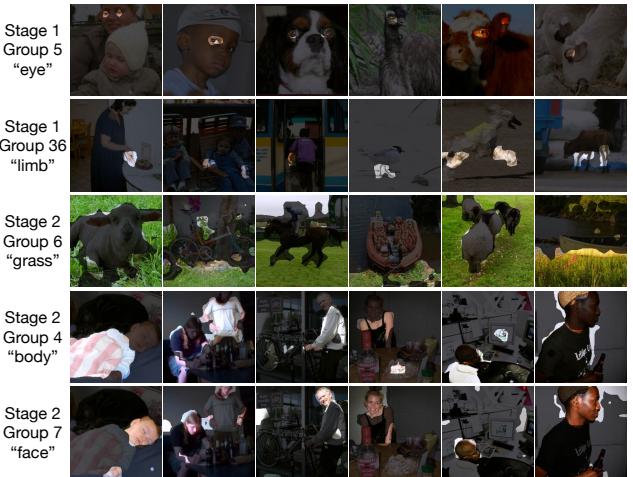


Figure 7. Concepts Learnt by Group Tokens.. We highlight the regions that group tokens attend to in different stages.

could be aggregated into higher level concepts in the second stage.

### 4.4. Comparisons with Existing Methods

We compare the zero-shot semantic segmentation performance of GroupViT with other zero-shot baselines and with methods for fully supervised transfer, based on ViT-S.

**Comparison with Zero-Shot Baselines** We train ViT and a text encoder with the image-text contrastive loss defined in CLIP [63], for comparison. To zero-shot transfer CLIP to semantic segmentation, during inference, we first apply non-parametric grouping on its output features. We then compute the similarity between the average feature of each group and the text embedding of the segmentation labels of the dataset. In this way, any non-parametric grouping method for ViT combined with CLIP can be considered as a zero-shot semantic segmentation baseline. We also report a “pixel-wise” baseline, which treats each pixel as a group and performs classification independently. As Table 4 shows that GroupViT outperforms other grouping methods

arch	model	pre-training		zero-shot	transfer	
		dataset	supervision		PASCAL VOC	PASCAL Context
ViT	DeiT [77]	ImageNet	class	✗	53.0	35.9
ViT	DINO [8]	ImageNet	self	✗	39.1	20.4
ViT	DINO [8]	CC12M+YFCC	self	✗	37.6	22.8
ViT	MoCo [14]	ImageNet	self	✗	34.3	21.3
ViT	MoCo [14]	CC12M+YFCC	self	✗	36.1	23.0
GroupViT	Ours	CC12M+YFCC	text	✓	<b>52.3</b>	<b>22.4</b>

Table 5. **Comparisons with fully supervised transfer.** Zero-shot ✓ means transfer to semantic segmentation without any fine-tuning. We report mIoU on the validation split of the PASCAL VOC 2012 and PASCAL Context datasets.

by a large margin. This demonstrates that, compared to ViT trained with CLIP, our GroupViT is more effective at zero-shot transfer to semantic segmentation. In the Table C.1, we also show that GroupViT’s ImageNet classification performance is comparable to that of ViT.

**Comparison with Fully-Supervised Transfer** We compare the performance of GroupViT with fully-supervised transfer to semantic segmentation. For fully-supervised transfer, we fine-tune a semantic segmentation head jointly with a pre-trained representation [13,97] on the training sets of the PASCAL VOC 2012 and PASCAL Context datasets separately and report their performances in Table 5. For a fair comparison, we employ a ViT architecture comparable to GroupViT’s for all baselines. Specifically, we append a  $1 \times 1$  convolution layer to a pre-trained ViT, trained with images of size  $224 \times 224$  and fine-tune the whole network with ground truth masks for 4k iterations. During inference, we resize the input images to have a shorter side length of 448 pixels. For fully-supervised transfer, we compare both supervised and self-supervised pre-training methods against GroupViT (Table 5). GroupViT (without fine-tuning) outperforms all variants of ViT pre-trained with self-supervision (with supervised fine-tuning) by a large margin on PASCAL VOC 2012 and is comparable to them on PASCAL Context. This implies that GroupViT, without any pixel-level annotations is able to transfer to several semantic segmentation datasets and can outperform existing state-of-the-art transfer-learning methods requiring more supervision (i.e., pixel-level labels for supervised transfer). Interestingly, on PASCAL VOC 2012, the zero-shot performance of GroupViT (mIoU of 52.3%) approaches that of fully-supervised ViT (mIoU of 53%) trained with both image classification and segmentation labels, which is significant.

## 5. Discussion

**Conclusion** We take the first step towards learning semantic segmentation with text alone and without any explicit human supervision. We show that, with GroupViT, the representation learned from large-scale noisy image-text pairs can be transferred to semantic segmentation in a zero-shot manner. This work also demonstrates that besides image

classification, text supervision could also be transferred to finer-grained vision tasks, which hasn’t yet been explored previously and opens up an exciting research direction.

**Limitations and Future Work** There are two potential improvements of GroupViT to explore in the future. Firstly, GroupViT’s performance is lower on PASCAL Context versus PASCAL VOC. This happens due to the presence of background classes, e.g., ground and road in PASCAL Context, which are less likely to be labeled in text; and misclassification of correctly grouped segments into incorrect textual classes (details in supplement Sec. C.3). Secondly, GroupViT’s architecture currently doesn’t integrate segmentation-specific enhancements, e.g., dilated convolutions [13], pyramid pooling [97] or a U-Net [65].

**Acknowledgements.** Prof. Wang’s lab is supported, in part, by grants from NSF CCF-2112665 (TILOS) and DARPA LwLL.

## References

- [1] Jiwoon Ahn and Suha Kwak. Learning pixel-level semantic affinity with image-level supervision for weakly supervised semantic segmentation. In *CVPR*, 2018. 3
- [2] Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh. Vqa: Visual question answering. In *ICCV*, 2015. 3
- [3] Pablo Arbeláez, Bharath Hariharan, Chunhui Gu, Saurabh Gupta, Lubomir Bourdev, and Jitendra Malik. Semantic segmentation using regions and parts. In *CVPR*, 2012. 1
- [4] Anurag Arnab, Mostafa Dehghani, Georg Heigold, Chen Sun, Mario Lučić, and Cordelia Schmid. Vivit: A video vision transformer. In *ICCV*, 2021. 2
- [5] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016. 13
- [6] Gedas Bertasius, Heng Wang, and Lorenzo Torresani. Is space-time attention all you need for video understanding? In *ICML*, 2021. 2
- [7] Maxime Bucher, Tuan-Hung Vu, Matthieu Cord, and Patrick Pérez. Zero-shot semantic segmentation. In *NeurIPS*, 2019. 3
- [8] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *ICCV*, 2021. 9, 13

- [9] Arantxa Casanova, Pedro O Pinheiro, Negar Rostamzadeh, and Christopher J Pal. Reinforced active learning for image segmentation. In *ICLR*, 2020. 3
- [10] Yu-Ting Chang, Qiaosong Wang, Wei-Chih Hung, Robinson Piramuthu, Yi-Hsuan Tsai, and Ming-Hsuan Yang. Weakly-supervised semantic segmentation via sub-category exploration. In *CVPR*, 2020. 3
- [11] 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. 2, 7, 8, 18
- [12] Kan Chen, Jiyang Gao, and Ram Nevatia. Knowledge aided consistency for weakly supervised phrase grounding. In *CVPR*, 2018. 3
- [13] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *IEEE TPAMI*, 2017. 9
- [14] Xinlei Chen, Saining Xie, and Kaiming He. An empirical study of training self-supervised vision transformers. In *ICCV*, 2021. 9, 13
- [15] Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and Jingjing Liu. Uniter: Learning universal image-text representations. In *ECCV*, 2020. 3
- [16] Bowen Cheng, Ross Girshick, Piotr Dollár, Alexander C. Berg, and Alexander Kirillov. Boundary IoU: Improving object-centric image segmentation evaluation. In *CVPR*, 2021. 8
- [17] Xiangxiang Chu, Zhi Tian, Yuqing Wang, Bo Zhang, Haibing Ren, Xiaolin Wei, Huaxia Xia, and Chunhua Shen. Twins: Revisiting the design of spatial attention in vision transformers. In *NeurIPS*, 2021. 2
- [18] Dorin Comaniciu and Peter Meer. Mean shift: A robust approach toward feature space analysis. *IEEE TPAMI*, 2002. 8
- [19] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, 2009. 3, 13
- [20] Karan Desai and Justin Johnson. Virtex: Learning visual representations from textual annotations. In *CVPR*, 2021. 3
- [21] Karan Desai, Gaurav Kaul, Zubin Aysola, and Justin Johnson. RedCaps: Web-curated image-text data created by the people, for the people. In *NeurIPS*, 2021. 18
- [22] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *NAACL*, 2018. 2
- [23] Nanqing Dong and Eric P Xing. Few-shot semantic segmentation with prototype learning. In *BMVC*, 2018. 3
- [24] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2021. 2, 3, 6, 13
- [25] Alaaeldin El-Nouby, Hugo Touvron, Mathilde Caron, Piotr Bojanowski, Matthijs Douze, Armand Joulin, Ivan Laptev, Natalia Neverova, Gabriel Synnaeve, Jakob Verbeek, et al. Xcit: Cross-covariance image transformers. *arXiv preprint arXiv:2106.09681*, 2021. 2
- [26] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *IJCV*, 2010. 2, 7
- [27] Haoqi Fan, Bo Xiong, Karttikeya Mangalam, Yanghao Li, Zhicheng Yan, Jitendra Malik, and Christoph Feichtenhofer. Multiscale vision transformers. In *ICCV*, 2021. 2
- [28] Junsong Fan, Zhaoxiang Zhang, Chunfeng Song, and Tie-niu Tan. Learning integral objects with intra-class discriminator for weakly-supervised semantic segmentation. In *CVPR*, 2020. 3
- [29] Zhe Gan, Yen-Chun Chen, Linjie Li, Chen Zhu, Yu Cheng, and Jingjing Liu. Large-scale adversarial training for vision-and-language representation learning. In *NeurIPS*, 2020. 3
- [30] Golnaz Ghiasi, Xiuye Gu, Yin Cui, and Tsung-Yi Lin. Open-vocabulary image segmentation. *arXiv preprint arXiv:2112.12143*, 2021. 3
- [31] Tanmay Gupta, Arash Vahdat, Gal Chechik, Xiaodong Yang, Jan Kautz, and Derek Hoiem. Contrastive learning for weakly supervised phrase grounding. In *ECCV*, 2020. 3
- [32] Eric Jang, Shixiang Gu, and Ben Poole. Categorical reparameterization with gumbel-softmax. In *ICLR*, 2017. 4
- [33] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc V Le, Yunhsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. In *ICML*, 2021. 3, 5
- [34] Sanghyun Jo and In-Jae Yu. Puzzle-cam: Improved localization via matching partial and full features. In *ICIP*, 2021. 3
- [35] Armand Joulin, Laurens Van Der Maaten, Allan Jabri, and Nicolas Vasilache. Learning visual features from large weakly supervised data. In *ECCV*, 2016. 3
- [36] Aishwarya Kamath, Mannat Singh, Yann LeCun, Gabriel Synnaeve, Ishan Misra, and Nicolas Carion. Mdetr-modulated detection for end-to-end multi-modal understanding. In *ICCV*, 2021. 3
- [37] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015. 7, 13
- [38] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, et al. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *IJCV*, 2017. 3
- [39] Jungbeom Lee, Eunji Kim, Sungmin Lee, Jangho Lee, and Sungroh Yoon. Ficklenet: Weakly and semi-supervised semantic image segmentation using stochastic inference. In *CVPR*, 2019. 3
- [40] Liunian Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. Visualbert: A simple and performant baseline for vision and language. *arXiv preprint arXiv:1908.03557*, 2019. 3

- [41] Peike Li, Yunchao Wei, and Yi Yang. Consistent structural relation learning for zero-shot segmentation. In *NeurIPS*, 2020. 3
- [42] Xiujun Li, Xi Yin, Chunyuan Li, Pengchuan Zhang, Xiaowei Hu, Lei Zhang, Lijuan Wang, Houdong Hu, Li Dong, Furu Wei, et al. Oscar: Object-semantics aligned pre-training for vision-language tasks. In *ECCV*, 2020. 3
- [43] Yi Li, Zhanghui Kuang, Liyang Liu, Yimin Chen, and Wayne Zhang. Pseudo-mask matters in weakly-supervised semantic segmentation. In *ICCV*, 2021. 3
- [44] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *ECCV*, 2014. 16
- [45] Xuejing Liu, Liang Li, Shuhui Wang, Zheng-Jun Zha, Li Su, and Qingming Huang. Knowledge-guided pairwise reconstruction network for weakly supervised referring expression grounding. In *ACM MM*, 2019. 3
- [46] Yongfei Liu, Bo Wan, Lin Ma, and Xuming He. Relation-aware instance refinement for weakly supervised visual grounding. In *CVPR*, 2021. 3
- [47] Yongfei Liu, Xiangyi Zhang, Songyang Zhang, and Xuming He. Part-aware prototype network for few-shot semantic segmentation. In *ECCV*, 2020. 3
- [48] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *ICCV*, 2021. 2
- [49] Ze Liu, Jia Ning, Yue Cao, Yixuan Wei, Zheng Zhang, Stephen Lin, and Han Hu. Video swin transformer. *arXiv preprint arXiv:2106.13230*, 2021. 2
- [50] Francesco Locatello, Dirk Weissenborn, Thomas Unterthiner, Aravindh Mahendran, Georg Heigold, Jakob Uszkoreit, Alexey Dosovitskiy, and Thomas Kipf. Object-centric learning with slot attention. In *NeurIPS*, 2020. 5
- [51] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *CVPR*, 2015. 1
- [52] Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts. In *ICLR*, 2016. 7
- [53] Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. In *NeurIPS*, 2019. 3
- [54] Zhihe Lu, Sen He, Xiatian Zhu, Li Zhang, Yi-Zhe Song, and Tao Xiang. Simpler is better: Few-shot semantic segmentation with classifier weight transformer. In *ICCV*, 2021. 3
- [55] Chris J Maddison, Andriy Mnih, and Yee Whye Teh. The concrete distribution: A continuous relaxation of discrete random variables. In *ICLR*, 2017. 4
- [56] Jitendra Malik. Visual grouping and object recognition. In *ICIP*, 2001. 1
- [57] Jitendra Malik, Pablo Arbeláez, João Carreira, Katerina Fragkiadaki, Ross Girshick, Georgia Gkioxari, Saurabh Gupta, Bharath Hariharan, Abhishek Kar, and Shubham Tulsiani. The three r's of computer vision: Recognition, reconstruction and reorganization. *Pattern Recognition Letters*, 72:4–14, 2016. 1
- [58] 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 *CVPR*, 2014. 2, 7
- [59] Khoi Nguyen and Sinisa Todorovic. Feature weighting and boosting for few-shot segmentation. In *ICCV*, 2019. 3
- [60] Aaron van den Oord, Oriol Vinyals, and Koray Kavukcuoglu. Neural discrete representation learning. In *NeurIPS*, 2017. 5
- [61] Bryan A Plummer, Paige Kordas, M Hadi Kiapour, Shuai Zheng, Robinson Piramuthu, and Svetlana Lazebnik. Conditional image-text embedding networks. In *ECCV*, 2018. 3
- [62] Bryan A Plummer, Liwei Wang, Chris M Cervantes, Juan C Caicedo, Julia Hockenmaier, and Svetlana Lazebnik. Flickr30k entities: Collecting region-to-phrase correspondences for richer image-to-sentence models. In *ICCV*, 2015. 3
- [63] 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. 2, 3, 5, 6, 7, 8, 13
- [64] Xiaofeng Ren and Jitendra Malik. Learning a classification model for segmentation. In *ICCV*, 2003. 1
- [65] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *MICCAI*, 2015. 9
- [66] Michael S Ryoo, AJ Piergiovanni, Anurag Arnab, Mostafa Dehghani, and Anelia Angelova. Tokenlearner: What can 8 learned tokens do for images and videos? In *NeurIPS*, 2021. 2, 8
- [67] Ozan Sener and Silvio Savarese. Active learning for convolutional neural networks: A core-set approach. In *ICLR*, 2017. 3
- [68] Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In *Proceedings of ACL*, 2018. 2, 18
- [69] Jianbo Shi and Jitendra Malik. Normalized cuts and image segmentation. *IEEE TPAMI*, 2000. 8
- [70] Wataru Shimoda and Keiji Yanai. Self-supervised difference detection for weakly-supervised semantic segmentation. In *ICCV*, 2019. 3
- [71] Gyungin Shin, Weidi Xie, and Samuel Albanie. All you need are a few pixels: semantic segmentation with pixelpick. In *ICCVW*, 2021. 3
- [72] Samarth Sinha, Sayna Ebrahimi, and Trevor Darrell. Variational adversarial active learning. In *ICCV*, 2019. 3
- [73] Guolei Sun, Wenguan Wang, Jifeng Dai, and Luc Van Gool. Mining cross-image semantics for weakly supervised semantic segmentation. In *ECCV*, pages 347–365. Springer, 2020. 3
- [74] Bart Thomee, David A Shamma, Gerald Friedland, Benjamin Elizalde, Karl Ni, Douglas Poland, Damian Borth, and Li-Jia Li. Yfcc100m: The new data in multimedia research. *Communications of the ACM*, 59(2):64–73, 2016. 2, 7, 8, 18

- [75] Zhuotao Tian, Hengshuang Zhao, Michelle Shu, Zhicheng Yang, Ruiyu Li, and Jiaya Jia. Prior guided feature enrichment network for few-shot segmentation. *IEEE TPAMI*, 2020. 3
- [76] Ilya Tolstikhin, Neil Houlsby, Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Thomas Unterthiner, Jessica Yung, Andreas Steiner, Daniel Keysers, Jakob Uszkoreit, et al. Mlp-mixer: An all-mlp architecture for vision. In *NeurIPS*, 2021. 6, 13
- [77] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. In *ICML*, 2021. 2, 6, 9, 13
- [78] Hugo Touvron, Matthieu Cord, Alexandre Sablayrolles, Gabriel Synnaeve, and Hervé Jégou. Going deeper with image transformers. In *ICCV*, 2021. 2
- [79] Zhuowen Tu, Xiangrong Chen, Alan L Yuille, and Song-Chun Zhu. Image parsing: Unifying segmentation, detection, and recognition. *IJCV*, 2005. 1
- [80] Jasper RR Uijlings, Koen EA Van De Sande, Theo Gevers, and Arnold WM Smeulders. Selective search for object recognition. *IJCV*, 2013. 1
- [81] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, 2017. 2, 5
- [82] Haochen Wang, Xudong Zhang, Yutao Hu, Yandan Yang, Xianbin Cao, and Xiantong Zhen. Few-shot semantic segmentation with democratic attention networks. In *ECCV*, 2020. 3
- [83] Liwei Wang, Jing Huang, Yin Li, Kun Xu, Zhengyuan Yang, and Dong Yu. Improving weakly supervised visual grounding by contrastive knowledge distillation. In *CVPR*, 2021. 3
- [84] Wenhai Wang, Enze Xie, Xiang Li, Deng-Ping Fan, Kaitao Song, Ding Liang, Tong Lu, Ping Luo, and Ling Shao. Pyramid vision transformer: A versatile backbone for dense prediction without convolutions. In *ICCV*, 2021. 2
- [85] Yude Wang, Jie Zhang, Meina Kan, Shiguang Shan, and Xilin Chen. Self-supervised equivariant attention mechanism for weakly supervised semantic segmentation. In *CVPR*, 2020. 3
- [86] Haiping Wu, Bin Xiao, Noel Codella, Mengchen Liu, Xiyang Dai, Lu Yuan, and Lei Zhang. Cvt: Introducing convolutions to vision transformers. *arXiv preprint arXiv:2103.15808*, 2021. 2
- [87] Enze Xie, Wenhai Wang, Zhiding Yu, Anima Anandkumar, Jose M Alvarez, and Ping Luo. Segformer: Simple and efficient design for semantic segmentation with transformers. In *NeurIPS*, 2021. 2
- [88] Shuai Xie, Zunlei Feng, Ying Chen, Songtao Sun, Chao Ma, and Mingli Song. Deal: Difficulty-aware active learning for semantic segmentation. In *ACCV*, 2020. 3
- [89] Mengde Xu, Zheng Zhang, Fangyun Wei, Yutong Lin, Yue Cao, Han Hu, and Xiang Bai. A simple baseline for zero-shot semantic segmentation with pre-trained vision-language model. *arXiv preprint arXiv:2112.14757*, 2021. 3
- [90] Lihe Yang, Wei Zhuo, Lei Qi, Yinghuan Shi, and Yang Gao. Mining latent classes for few-shot segmentation. In *ICCV*, 2021. 3
- [91] Raymond A Yeh, Minh N Do, and Alexander G Schwing. Unsupervised textual grounding: Linking words to image concepts. In *CVPR*, 2018. 3
- [92] Li Yuan, Yunpeng Chen, Tao Wang, Weihao Yu, Yujun Shi, Zihang Jiang, Francis EH Tay, Jiashi Feng, and Shuicheng Yan. Tokens-to-token vit: Training vision transformers from scratch on imagenet. In *ICCV*, 2021. 2
- [93] Nir Zabari and Yedid Hoshen. Semantic segmentation in-the-wild without seeing any segmentation examples. *arXiv preprint arXiv:2112.03185*, 2021. 3
- [94] Rowan Zellers, Yonatan Bisk, Ali Farhadi, and Yejin Choi. From recognition to cognition: Visual commonsense reasoning. In *CVPR*, 2019. 3
- [95] Pengchuan Zhang, Xiyang Dai, Jianwei Yang, Bin Xiao, Lu Yuan, Lei Zhang, and Jianfeng Gao. Multi-scale vision longformer: A new vision transformer for high-resolution image encoding. In *ICCV*, 2021. 2
- [96] Yuhao Zhang, Hang Jiang, Yasuhide Miura, Christopher D Manning, and Curtis P Langlotz. Contrastive learning of medical visual representations from paired images and text. *arXiv preprint arXiv:2010.00747*, 2020. 3
- [97] Hengshuang Zhao, Jianping Shi, Xiaojuan Qi, Xiaogang Wang, and Jiaya Jia. Pyramid scene parsing network. In *CVPR*, 2017. 9
- [98] 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 *CVPR*, 2021. 2
- [99] Chong Zhou, Chen Change Loy, and Bo Dai. Denseclip: Extract free dense labels from clip. *arXiv preprint arXiv:2112.01071*, 2021. 3
- [100] Luowei Zhou, Hamid Palangi, Lei Zhang, Houdong Hu, Jason Corso, and Jianfeng Gao. Unified vision-language pre-training for image captioning and vqa. In *AAAI*, 2020. 3
- [101] Song-Chun Zhu and David Mumford. *A stochastic grammar of images*. Now Publishers Inc, 2007. 1
- [102] C Lawrence Zitnick and Piotr Dollár. Edge boxes: Locating object proposals from edges. In *ECCV*, 2014. 1

## A. Implementation Details

### A.1. Architecture

The architecture of GroupViT is based on ViT-S [24, 77] with 12 Transformer layers. Each layer consists of a multi-head self-attention and an MLP block. The input to each block is normalized by layer normalization [5]. We connect the group tokens in the different grouping stages via MLP-Mixer layers [76]. Our text-encoder consists of 12 Transformer layers, each with a hidden dimension of 256. Following [63], the Transformer operates on a lower-cased byte pair encoding (BPE) representation of the text with a vocabulary of 49,152 words.

### A.2. Fully-Supervised Transfer to Semantic Segmentation

To implement the baselines for fully-supervised transfer to semantic segmentation, we fine-tune the pre-trained ViT model jointly with a  $1 \times 1$  convolutional layer appended to it for pixel-wise classification. We scale each input image by a randomly selected factor in the range of  $[0.5, 2]$  and then crop random  $224 \times 224$  patches from each image during training. We use the Adam [37] optimizer with a weight decay of 0.05 and a learning rate 0.001. We train all models for 4k iterations with a batch size of 16. During inference, we resize each input image to have a shorter side of size 448 pixels. We open-source our code at <https://github.com/NVlabs/GroupViT>.

## B. Qualitative Results

**PASCAL VOC 2012** We show additional qualitative results of GroupViT on the PASCAL VOC 2012 dataset, i.e. examples with single object in Fig. B.1; multiple objects from the same category in Fig. B.2; and multiple objects from different categories in Fig. B.3. Observe that GroupViT successfully groups and correctly classifies the objects in these various challenging scenarios.

**PASCAL Context** We show more qualitative results of GroupViT on the PASCAL Context dataset in Fig. B.4. The PASCAL Context dataset annotates not only *object* classes from PASCAL VOC 2012, e.g. car and dog, but also *stuff* classes related to the context, e.g. sky and water. Observe that GroupViT successfully segments *object* and *stuff* classes in the PASCAL Context dataset, e.g., cat and window in the second row, and dog and water in the sixth row.

## C. Additional Experiments and Analysis

### C.1. Image Classification

We compare the performance of the GroupViT and ViT architectures for the task of object classification on Im-

geNet. Following CLIP [63], here we train both architectures using supervision from text only via an image-text contrastive loss. In Table C.1, we report both the zero-shot and the linear probing accuracy on the ImageNet [19] validation split. The zero-shot and linear probing evaluation follow the same setting as CLIP [63]. GroupViT’s ImageNet classification performance is comparable to (if not better than) that of ViT, thus demonstrating that our proposed grouping mechanism enhances the baseline ViT architecture with the capability to perform semantic pixel grouping and zero-shot transfer to semantic segmentation, without affecting its object classification performance.

model	zero-shot Acc@1	linear Acc@1
ViT	42.4	69.2
GroupViT	42.9	69.8

Table C.1. **ImageNet Accuracy.**

### C.2. Mask Probing

We follow the procedure outlined in DINO [8] to evaluate the quality of the masks generated by GroupViT and by the baseline ViT model pre-trained using prior methods in a fully supervised [77], self-supervised [8, 14] or text-supervised [63] manner. For the ViT models, similar to DINO [8] for each final attention head, we compute its similarity to the [CLS] token and derive an attention mask for the pixels with the highest attention values. We then compute the Jaccard similarity of each head’s attention mask to the ground truth mask and retain the attention mask with the highest similarity. As for GroupViT, it does not have a multi-head design in the Grouping Block. Thus, we directly select the group most similar, as measured by the Jaccard index, to the ground truth mask for each image. As Table C.2 shows, the mask probing result of GroupViT is significantly better than that of all variants of the baseline ViT architecture. Hence, compared to ViT, our GroupViT more effectively groups semantically-related visual inputs together.

### C.3. Limitations

We found that the mIoU of GroupViT on PASCAL Context is significantly lower than that on PASCAL VOC

arch	model	dataset	supervision	Jaccard Similarity
ViT	Random	-	-	23.6
ViT	DeiT [77]	ImageNet	class	24.6
ViT	MoCo [14]	ImageNet	self	28.2
ViT	DINO [8]	ImageNet	self	45.9
ViT	DINO [8]	CC12M+YFCC	self	41.8
ViT	CLIP [63]	CC12M+YFCC	text	28.6
GroupViT	Ours	CC12M+YFCC	text	51.8

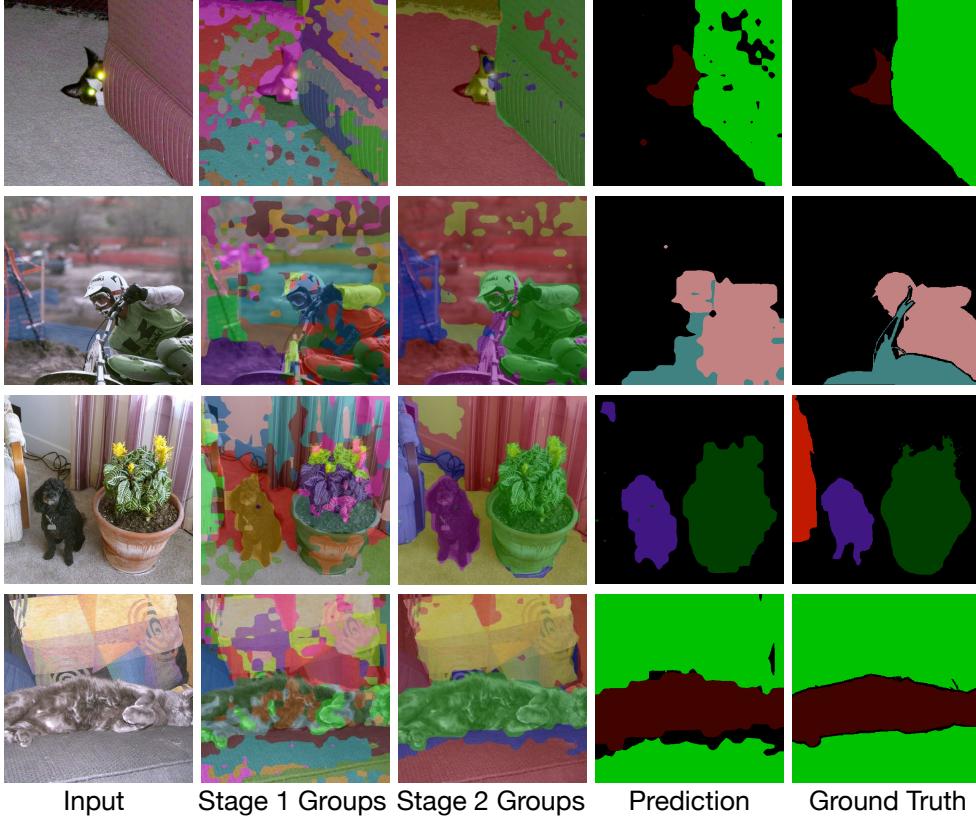
Table C.2. **Comparison of mask probing performance**  
GroupViT outperforms all other variants of the baseline ViT architecture at effectively grouping image regions on semantic groups.



**Figure B.1. Qualitative Results of GroupViT on PASCAL VOC 2012.** The results in columns labeled “Stage 1/2” show grouping results prior to assigning labels, where the regions belong to the same group are in the same color. All these examples contain a single object from a category.



**Figure B.2. Qualitative Results of GroupViT on PASCAL VOC 2012.** The results in columns labeled “Stage 1/2” show grouping results prior to assigning labels. The regions belong to the same group are in the same color. These examples contain multiple objects from the same category.



**Figure B.3. Qualitative Results of GroupViT on PASCAL VOC 2012.** The results in columns labeled “Stage 1/2” show grouping results prior to assigning labels, where the regions belong to the same group are in the same color. These examples contain multiple objects from multiple different categories.

2012. This could be attributed to the presence of background classes in PASCAL Context, e.g., ground, road and wall that result in low IoU ( $\sim 1.5$ ) on zero-shot transferring GroupViT to semantic segmentation on PASCAL Context. Through visual inspection, we found that while the pixels belonging to these background classes are typically correctly grouped into a single group by GroupViT, the group as a whole may be miss-classified into the wrong class on being compared to the text embedding of the various class labels. We hypothesize that this, in turn, happens due to the low probability of the background classes being described in textual sentences used during training. We show examples of such failure case in Fig. C.5. We further conduct an oracle experiment to verify this finding. In the oracle experiment, for each output group from GroupViT, we compute its IoU with all ground truth masks and assign to each group the class label that results in the maximum IoU. This represents the upper bound of GroupViT’s performance since here we leverage ground truth masks to predict each group’s class label. We use our 2-stage GroupViT trained on CC12M and YFCC datasets for this oracle experiment, which is the same model labeled ”Ours” in Table 5 of the main paper. We report the oracle experiment’s

results on PASCAL Context in Table C.3. The large gap between the performance of the original and oracle mIoU values on the PASCAL Context dataset, shows that while GroupViT’s grouping results are reasonably good, there is room to further improve the groups’ classification to segmentation class labels via image-text embedding similarity.

arch	dataset	mask mIoU	oracle mask mIoU
GroupViT	PASCAL VOC	52.3	73.7
GroupViT	PASCAL Context	22.4	54.6
GroupViT	COCO	24.3	54.0

Table C.3. Original versus oracle results.

#### C.4. COCO Dataset

We evaluate the performance of GroupViT on the COCO dataset [44], which contains 80 object classes. We combine the instance masks of the same category to get the semantic segmentation mask for each image. We report semantic segmentation mIoU on COCO in Table C.3. It demonstrates that GroupViT is able to transfer to complex datasets with various number of classes.

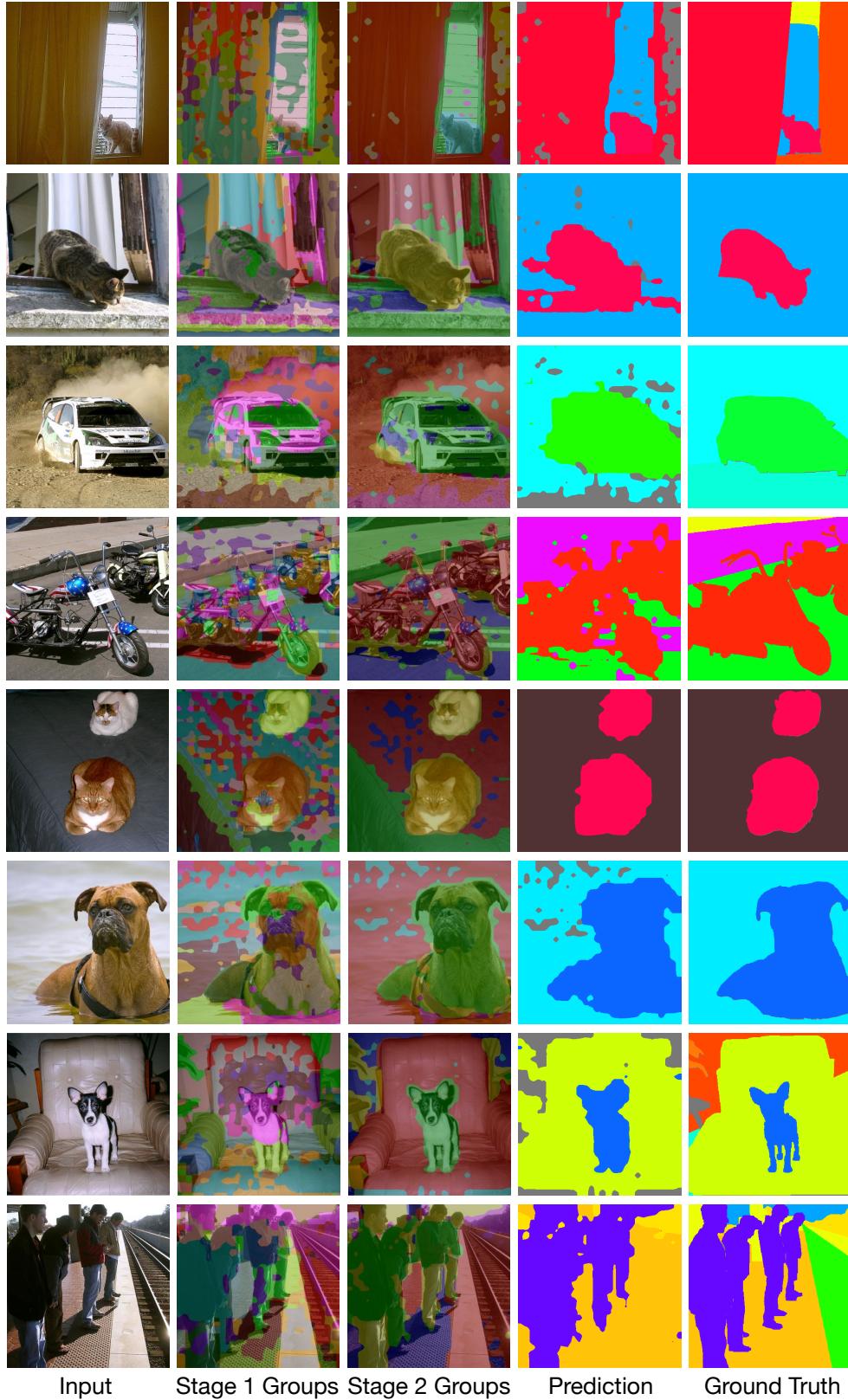


Figure B.4. **Qualitative Results of GroupViT on PASCAL Context.** Columns labeled “Stage 1/2” show grouping results prior to assigning labels, where the regions belong to the same group are in the same color. GroupViT can successfully segment *object* and *stuff* classes, e.g. cat and window in row 2, dog and water in row 6.

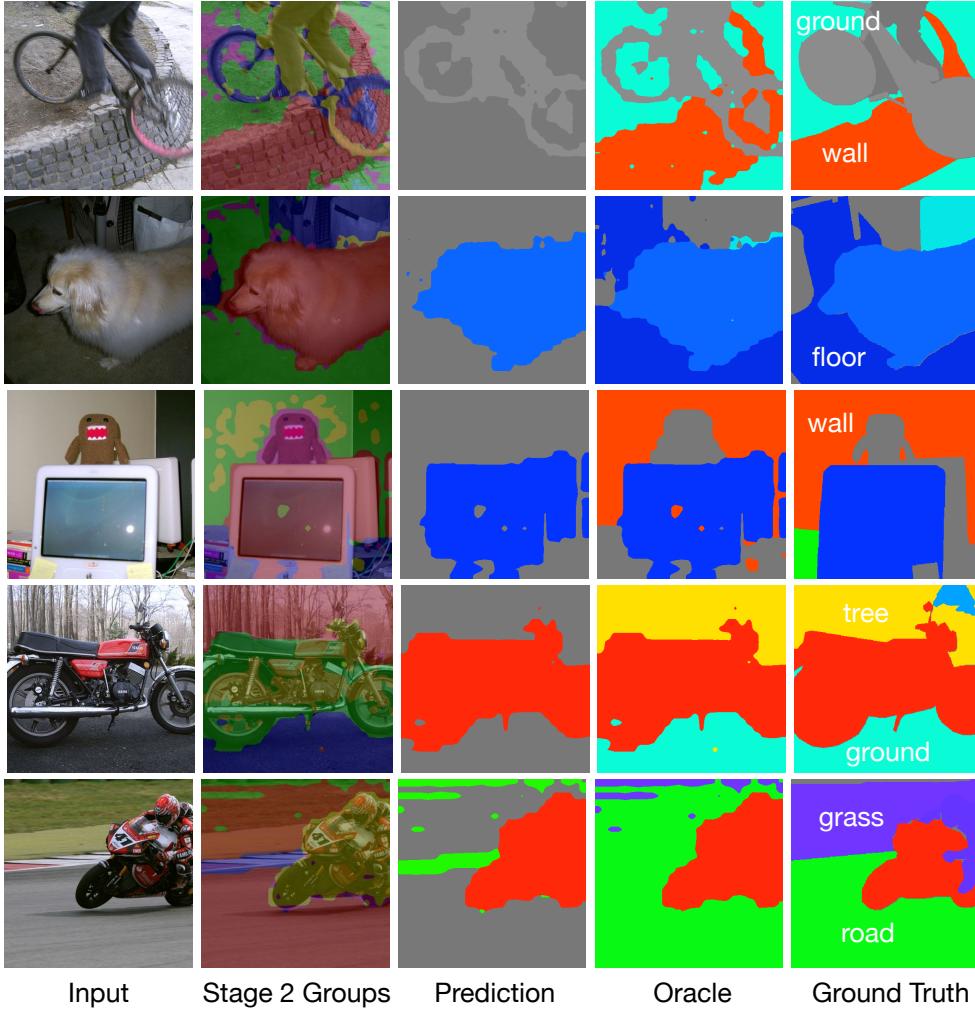


Figure C.5. **Failure cases on PASCAL Context.** ‘‘Oracle’’ shows the results of assigning groups to segmentation classes based on their IoU with the ground truth masks. Although GroupViT successfully groups *stuff* classes, e.g. ground, road and wall, it is not able to classify them correctly using the similarity between the visual and text embedding.

### C.5. Training on RedCaps

To show that our approach is generalizable to other training datasets, besides CC [11, 68] and filtered YFCC [74], we also train GroupViT on the recently released RedCaps dataset [21], which contains 12 millions image-text pairs from Reddit, of similar size as filtered YFCC. We report mIoU for zero-shot transfer to various image segmentation benchmarks datasets in Table C.4. Replacing YFCC with RedCaps yields similar accuracy on Pascal VOC, PASCAL Context and COCO datasets. It demonstrates that GroupViT is able to learn grouping with properly filtered image text pairs.

arch	Training Dataset	PASCAL VOC	PASCAL Context	COCO
GroupViT	CC+YFCC	52.3	22.4	24.3
GroupViT	CC+RedCaps	50.8	23.7	27.5

Table C.4. **Results trained with CC+RedCaps.**