

# 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, \*work conducted at an internship at NVIDIA

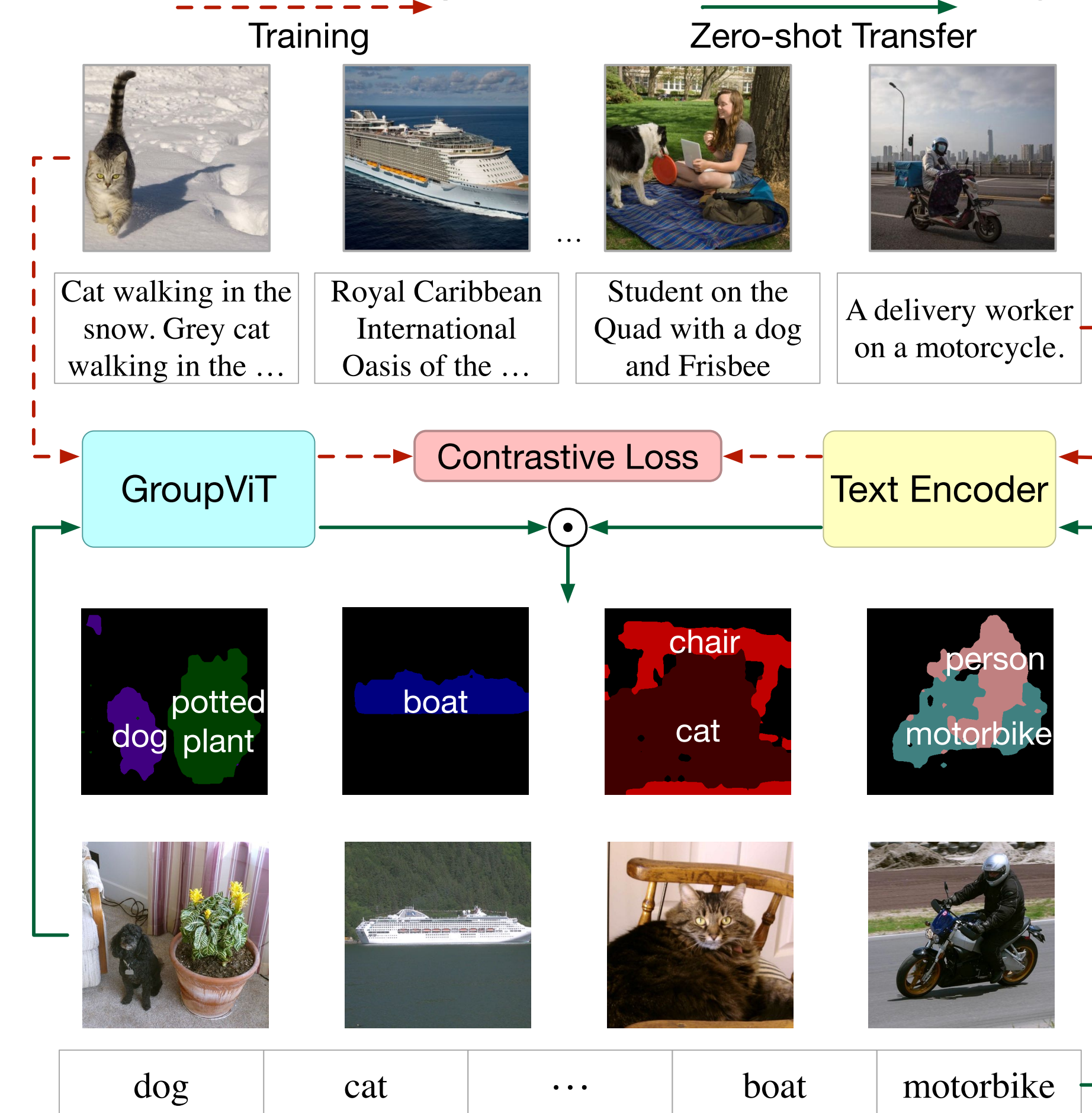
## Goal

To learn semantic segmentation with text supervision and **without** mask labels.



## Method Overview

First, we jointly train GroupViT and a text encoder using paired image-text data. With GroupViT, meaningful semantic image groups automatically emerge. Then, we transfer the trained GroupViT to zero-shot semantic segmentation.



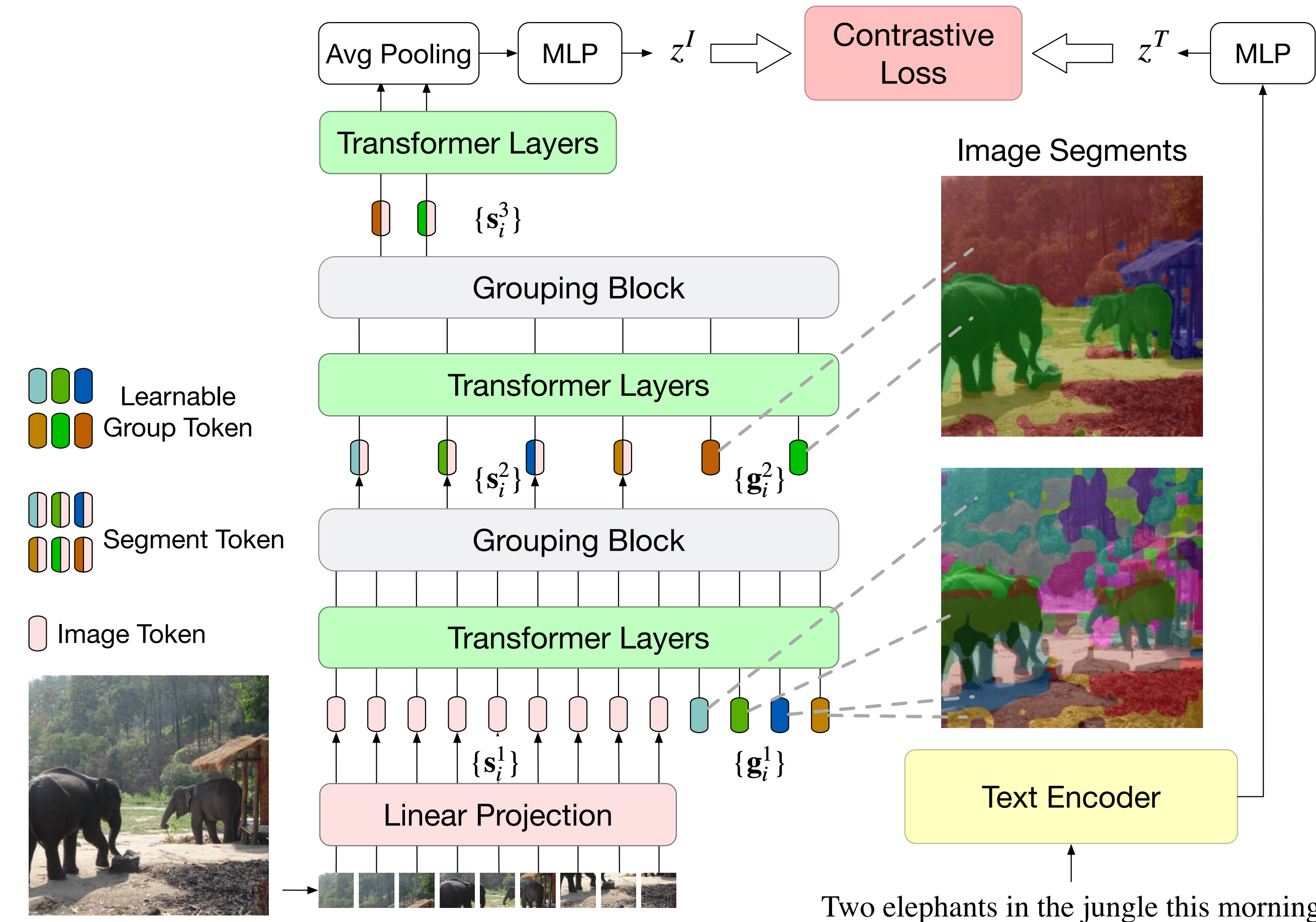
Code and demo!

## Main Contributions

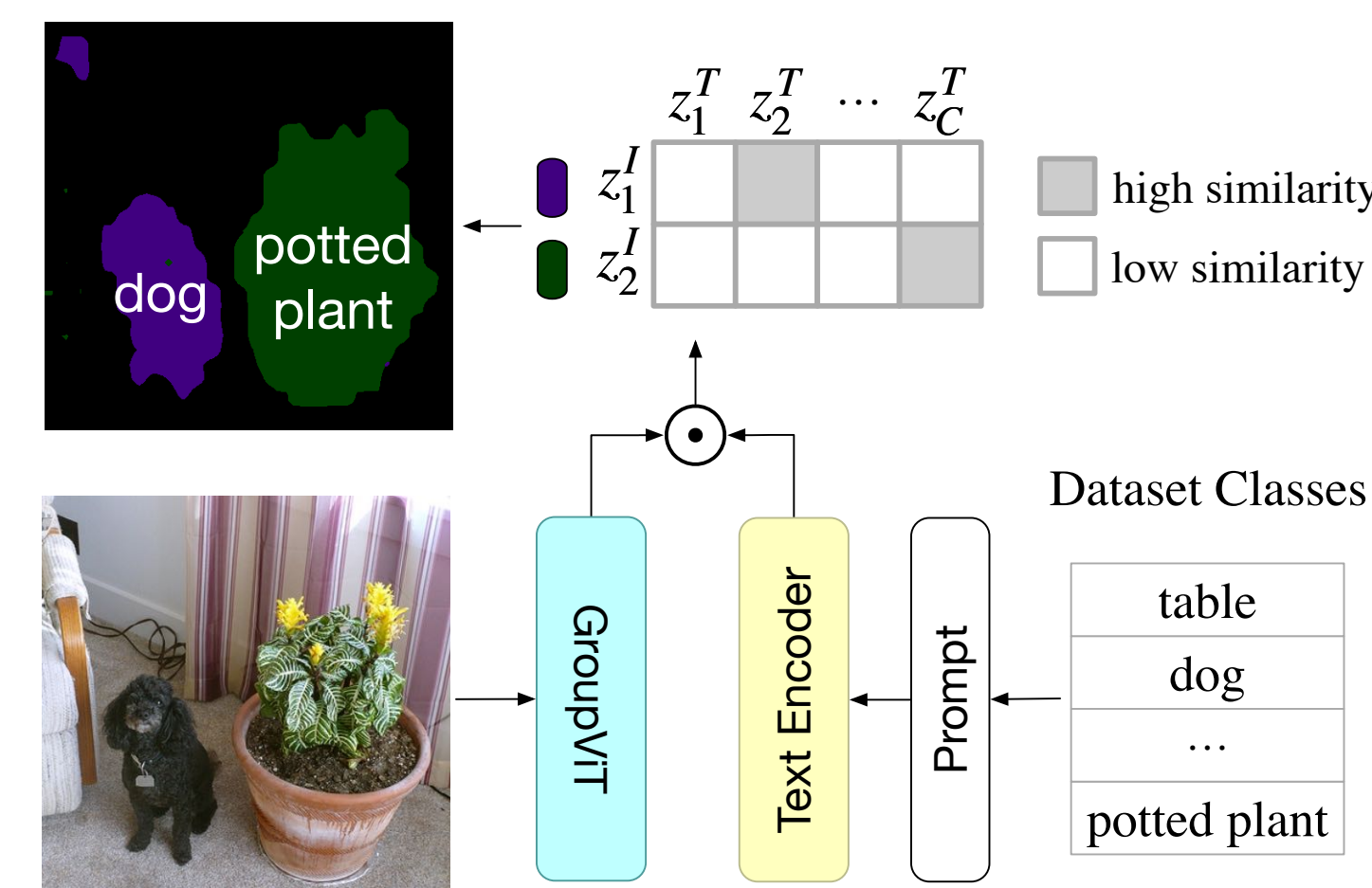
- Moving beyond regular-shaped image grids in deep networks and performing bottom-up grouping of visual concepts
- Zero-shot transfer to semantic segmentation without any pixel-level labels

## Architecture and Training Pipeline of GroupViT

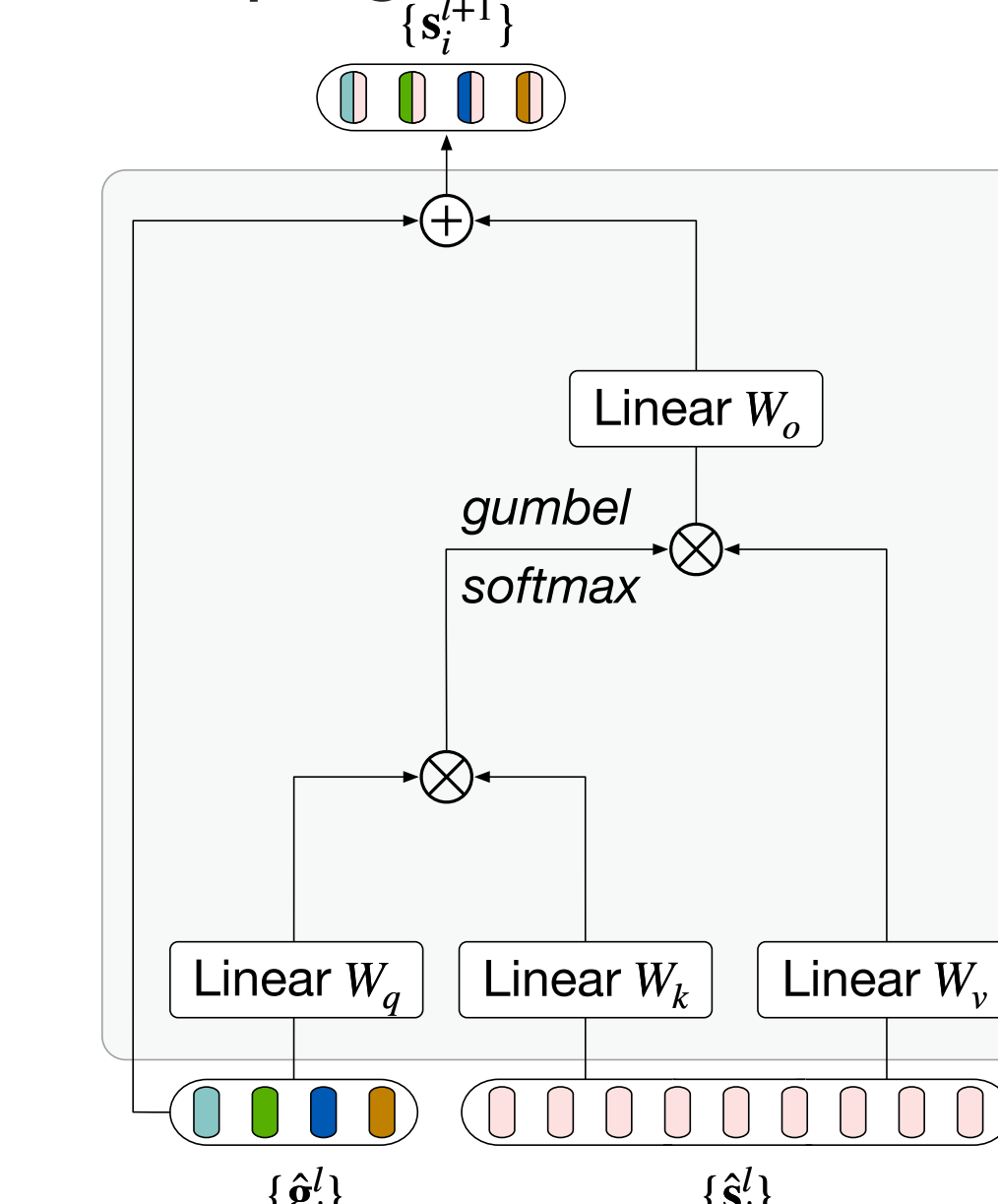
GroupViT contains Transformer layers with a hierarchy of grouping stages. Each stage groups image pixels into progressively larger visual segments. The images on the right show visual segments that emerge in the different grouping stages.



## Zero-shot Transfer to Semantic Segmentation



## Grouping Block



## Visualization Results without Training on any Mask annotations

