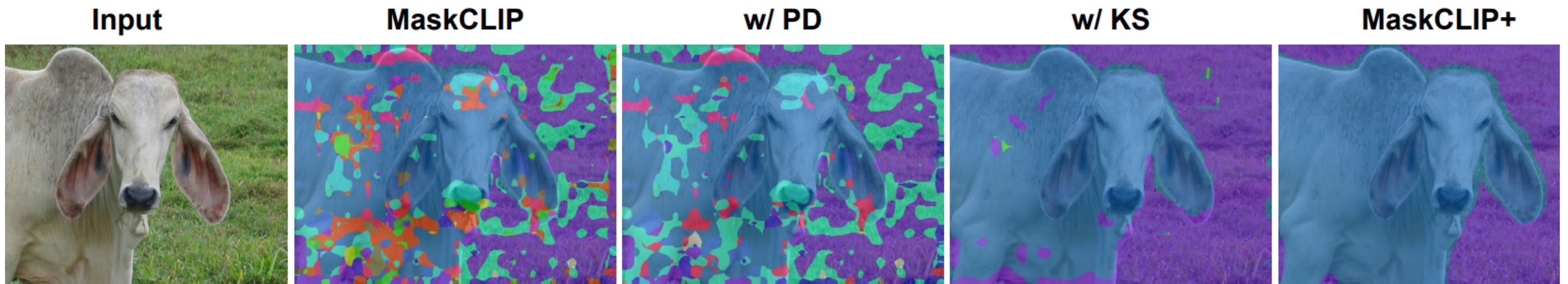


Extract Free Dense Labels from CLIP

Authors: Chong Zhou, Chen Change Loy , and Bo Dai



Sonia Castro Paniello
20.02.24

INTRODUCTION

Goal: Dense prediction with CLIP

CLIP

Learning Transferable Visual Models From Natural Language Supervision



Dense Prediction ?

Goal: Dense prediction with CLIP

CLIP

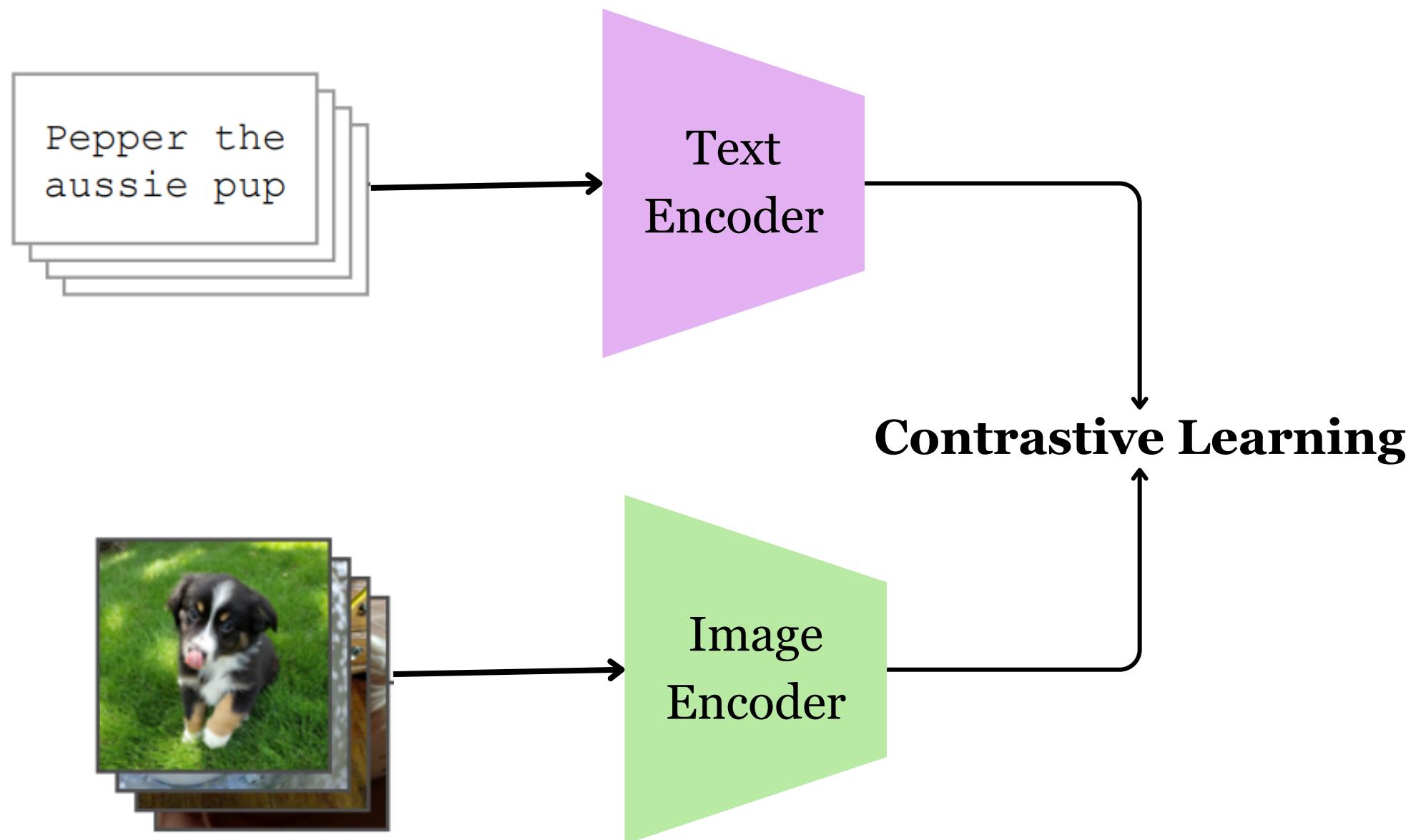
Learning Transferable Visual Models From Natural Language Supervision



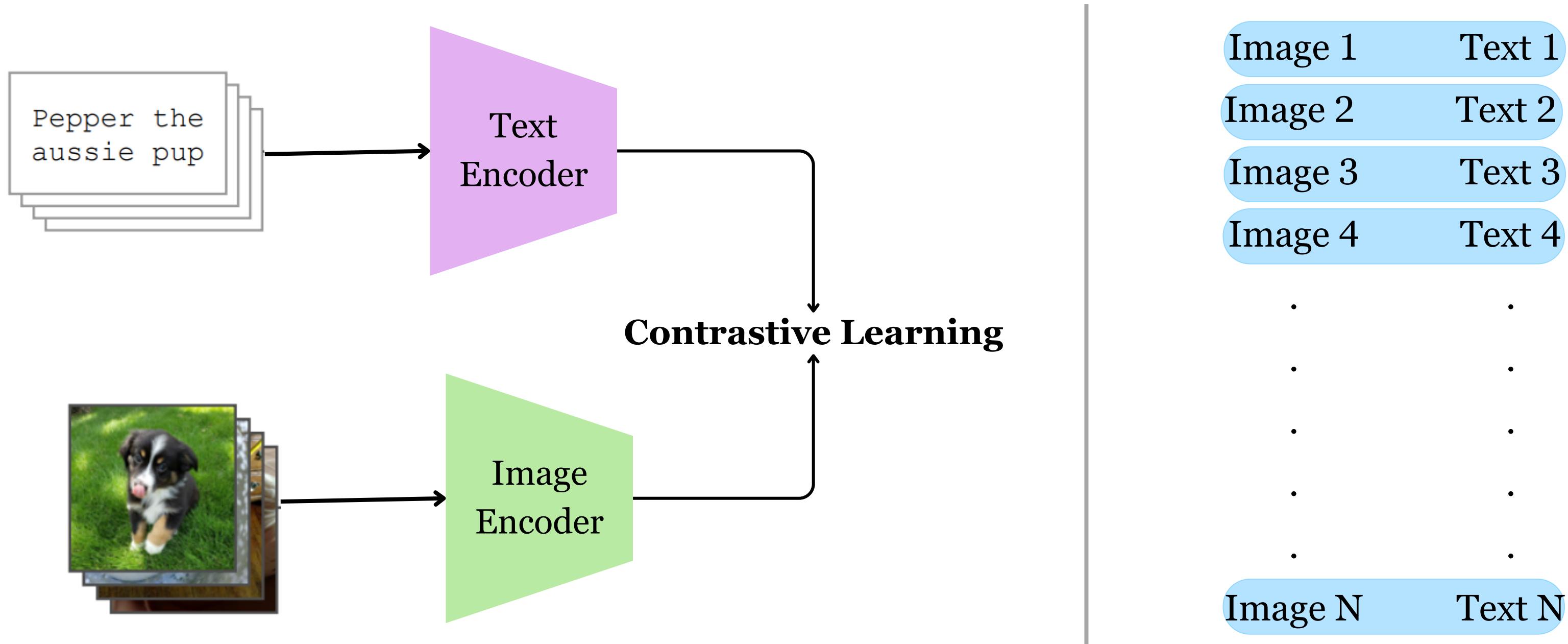
Dense Prediction ?

- Expensive Annotations
- Open Vocabulary

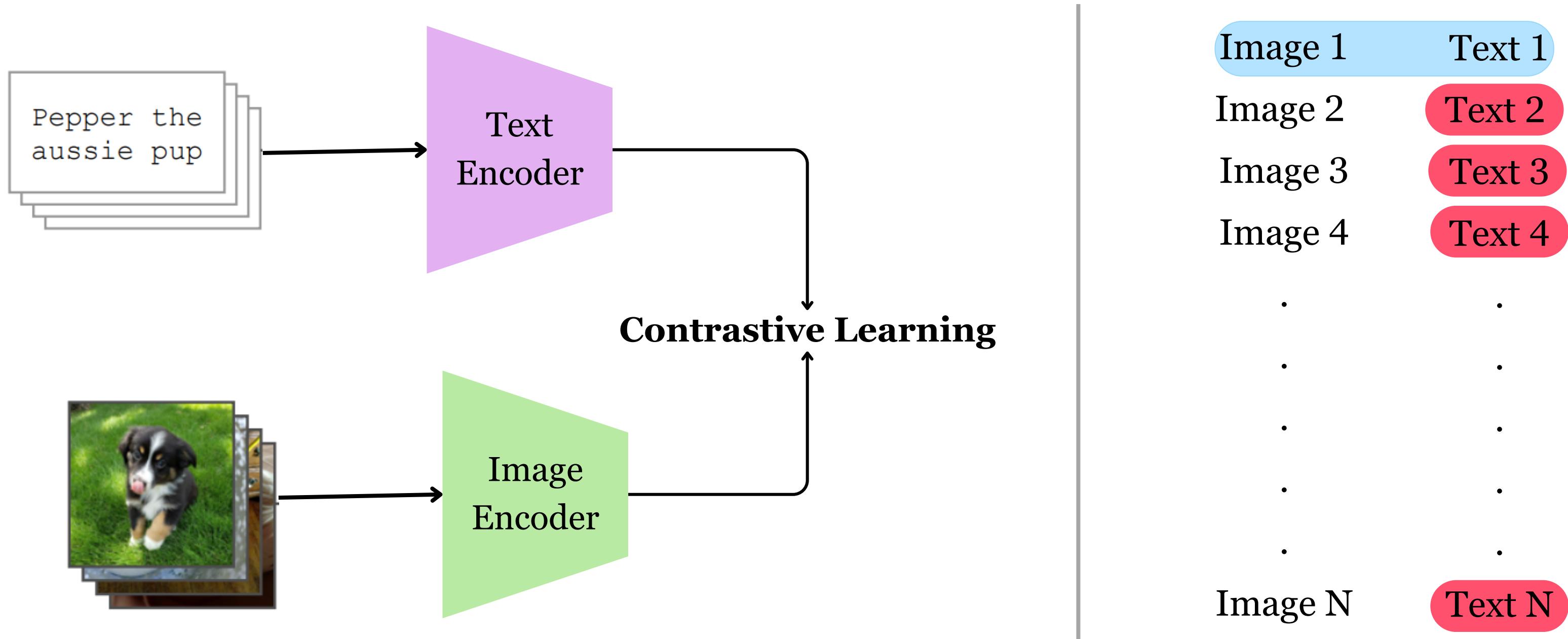
CLIP: Contrastative Image-Language Pretraining



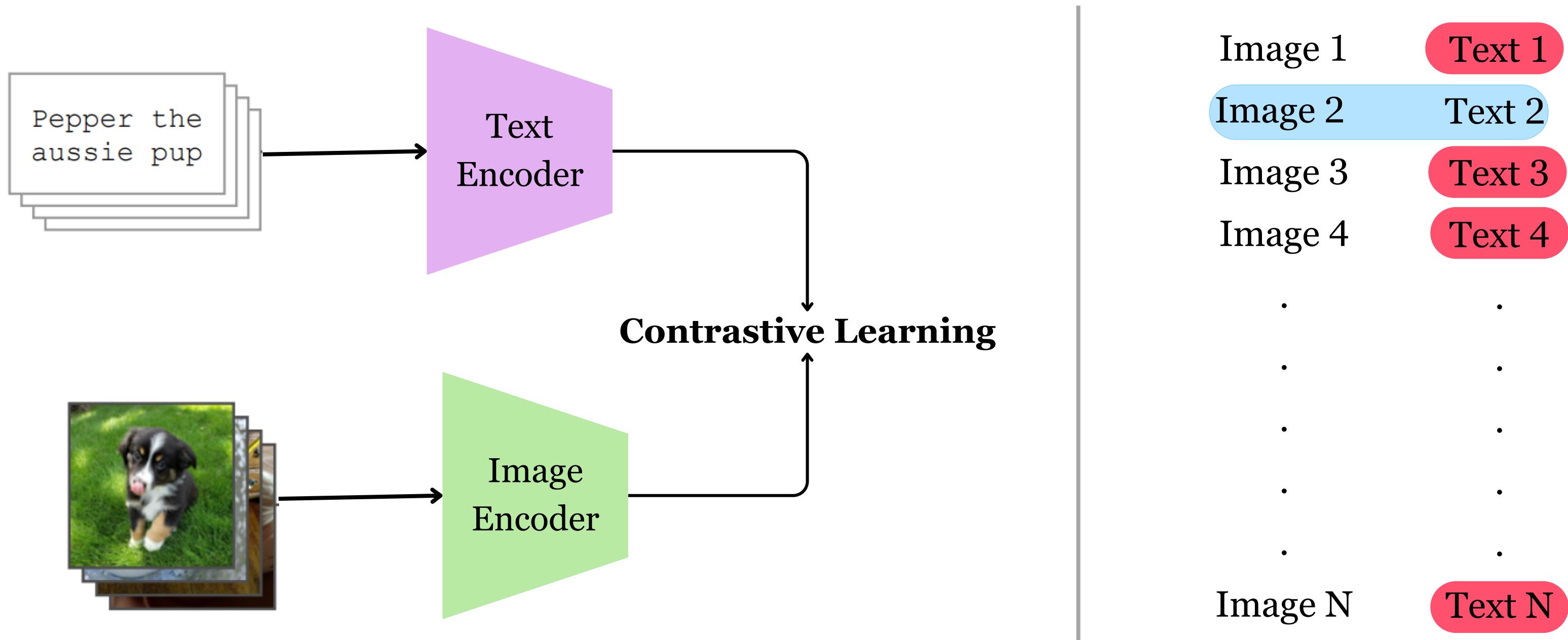
CLIP: Contrastative Image-Language Pretraining



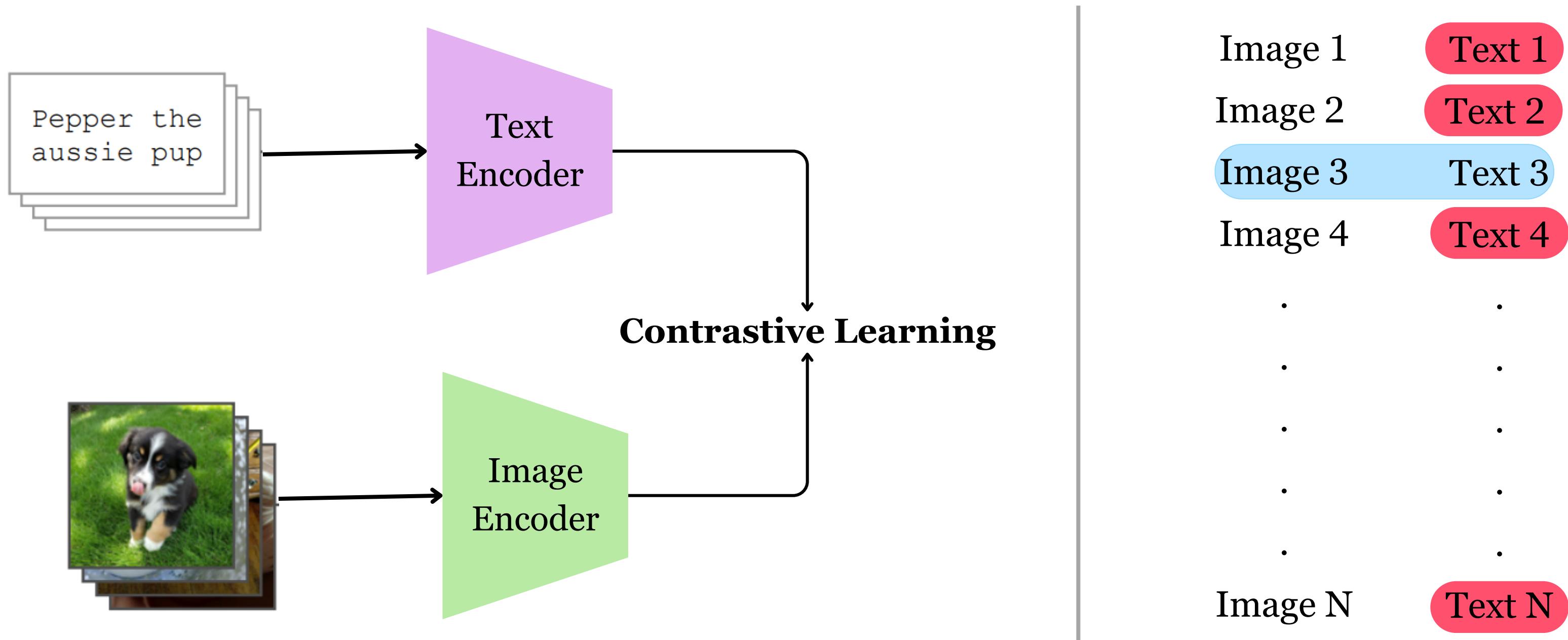
CLIP: Contrastative Image-Language Pretraining



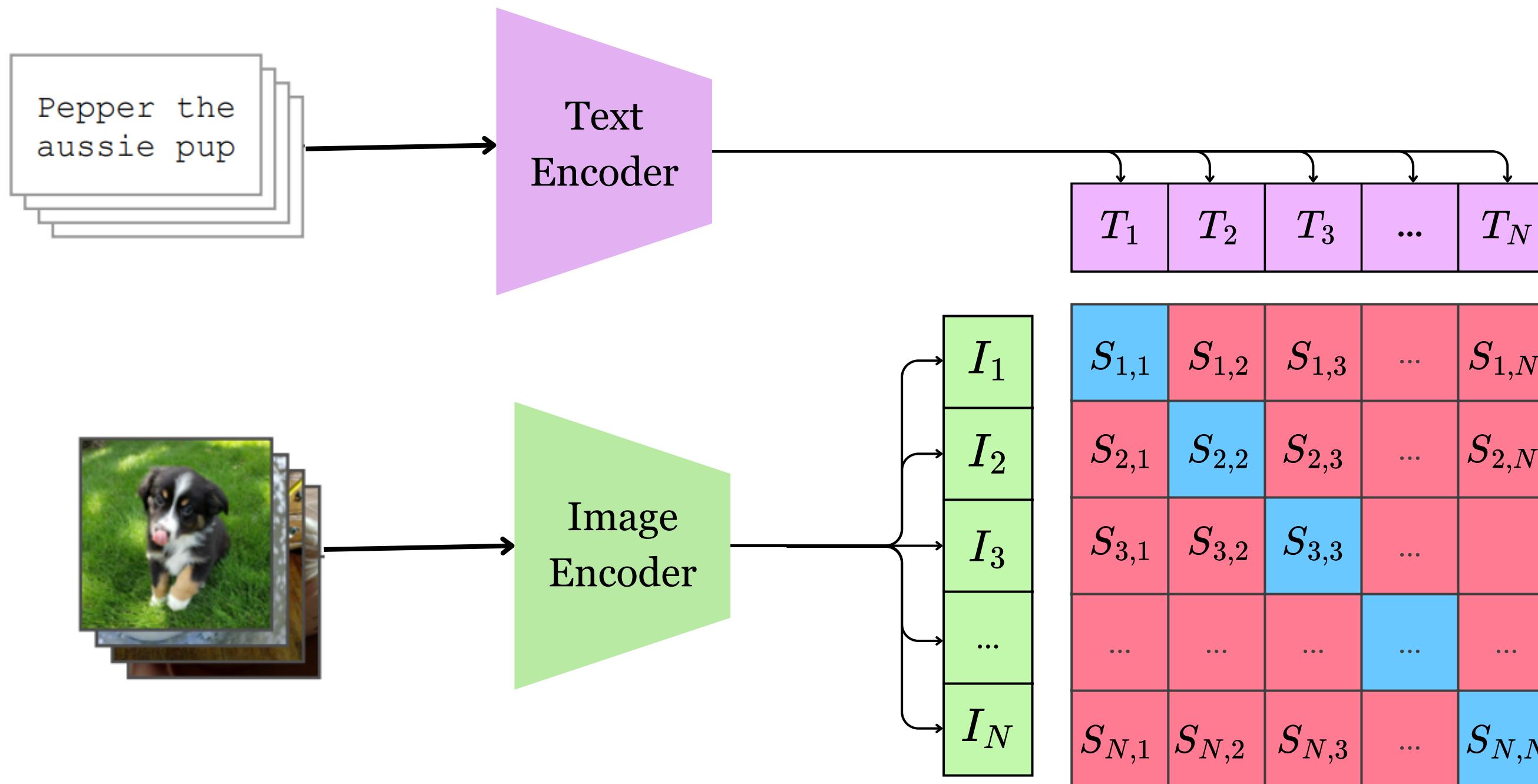
CLIP: Contrastative Image-Language Pretraining



CLIP: Contrastative Image-Language Pretraining



CLIP: Contrastative Image-Language Pretraining



Intuition: Descriptions contain Dense Semantic Guidance



The man at bat readies to swing at the pitch while the umpire looks on.

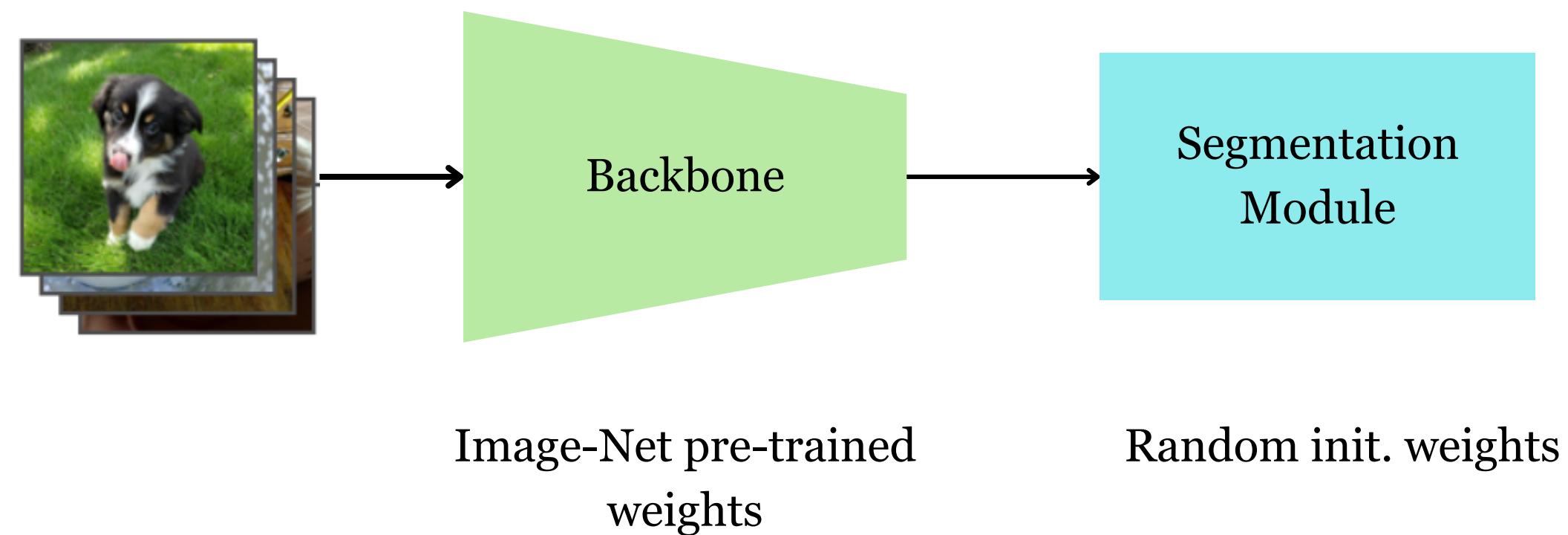
- Local image semantics
- Open vocabulary
- Contextual information



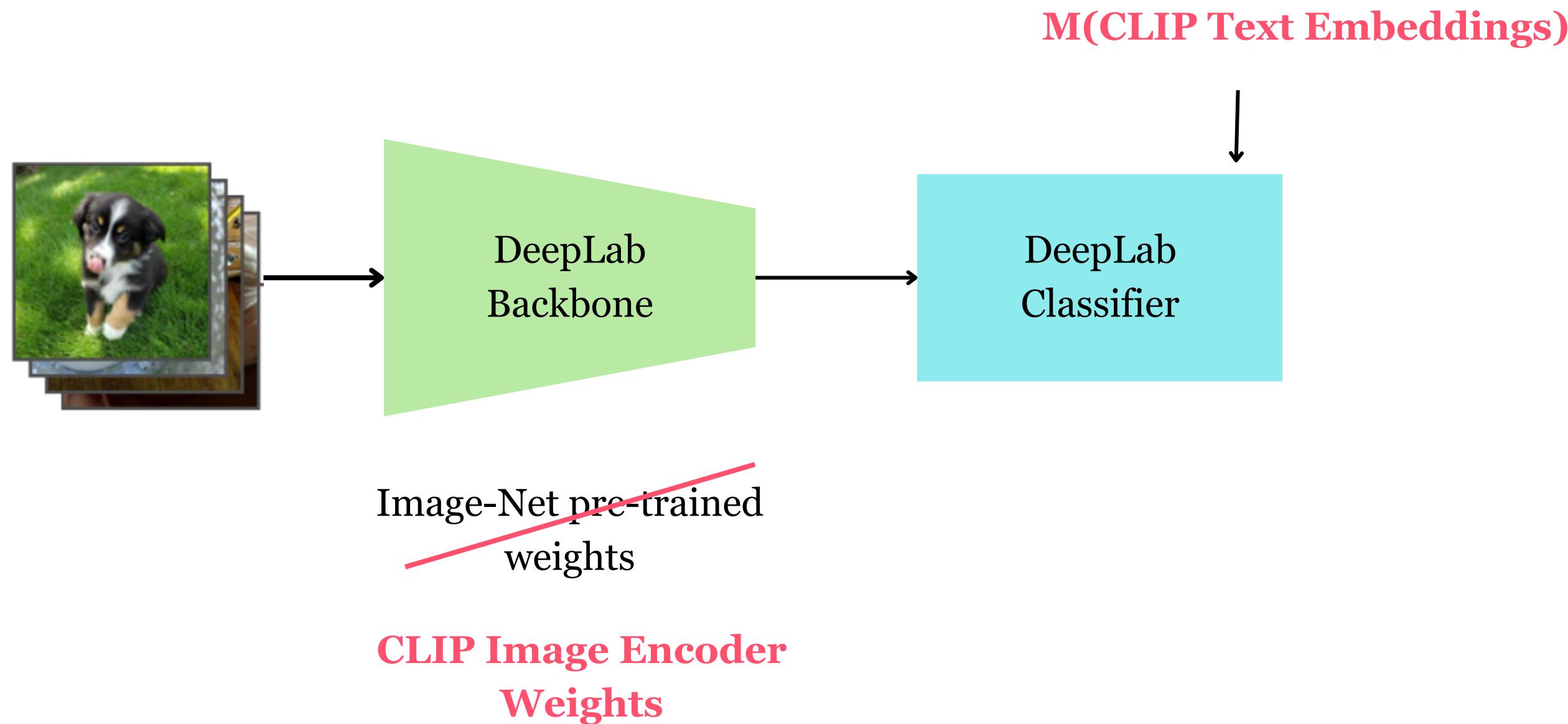
A horse carrying a large load of hay and two people sitting on it.

METHOD

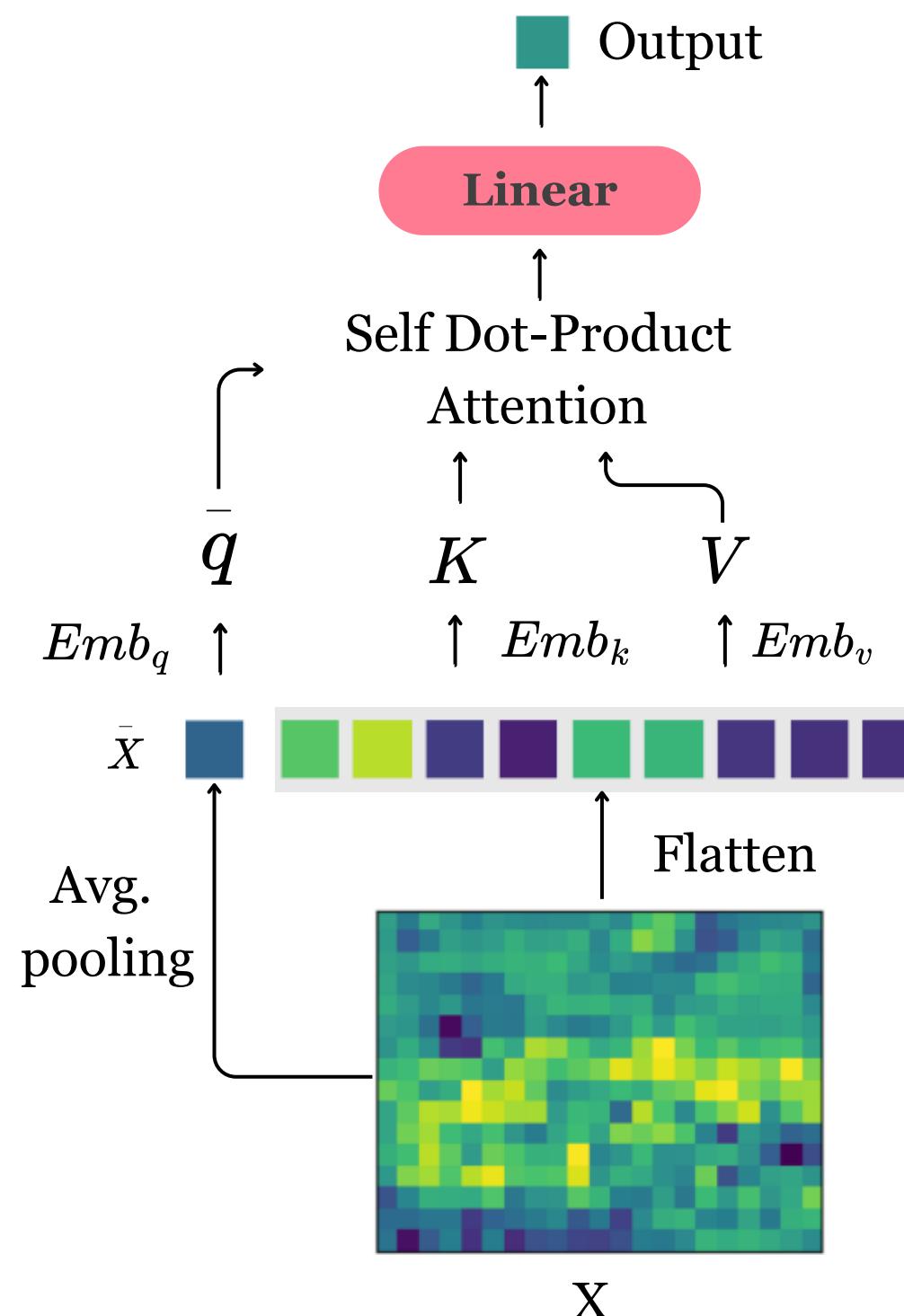
First Approach: Backbone + Segmentation Module



First Approach: Fails to segment unannotated classes

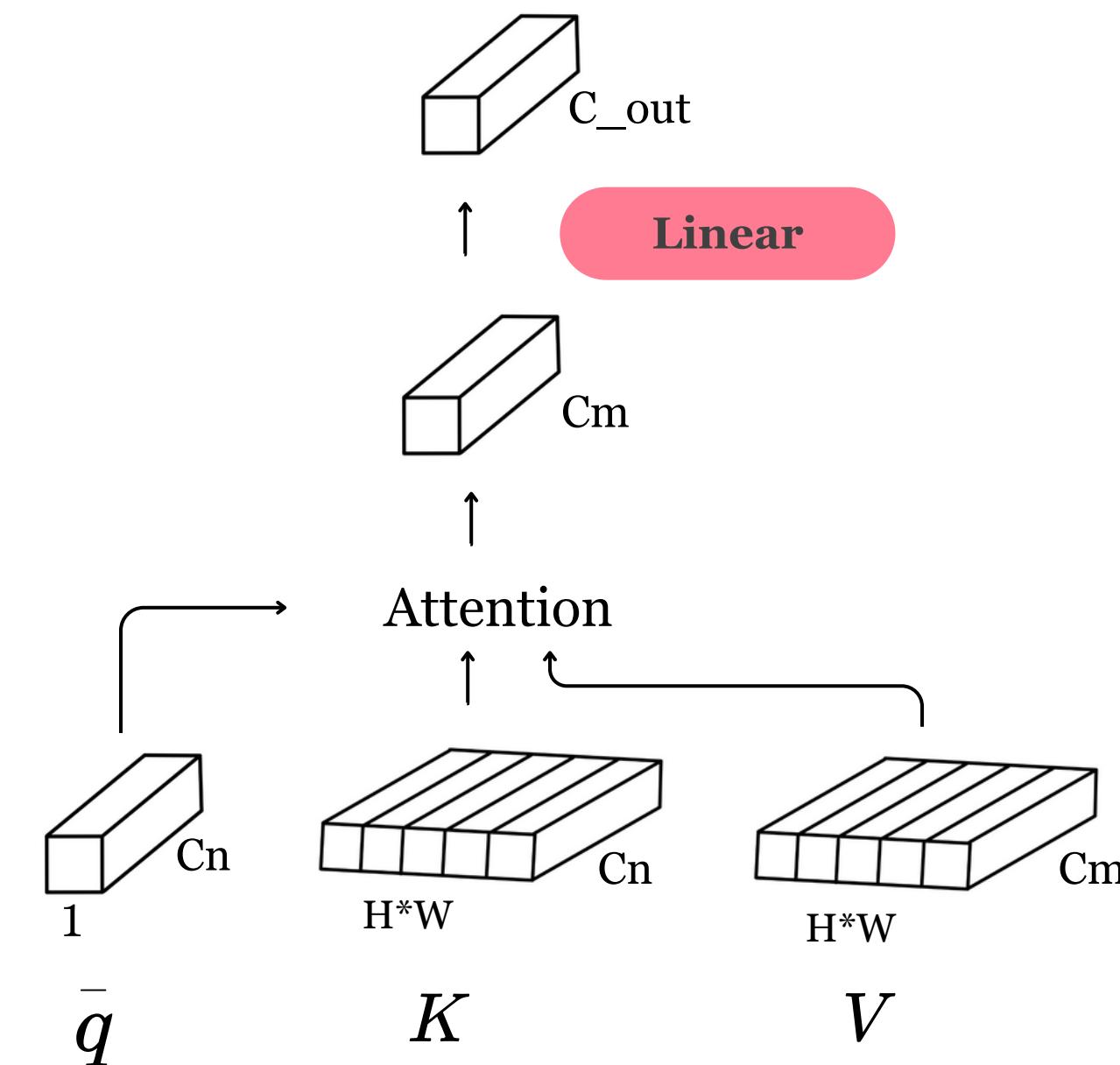


CLIP: Global Attention Pooling

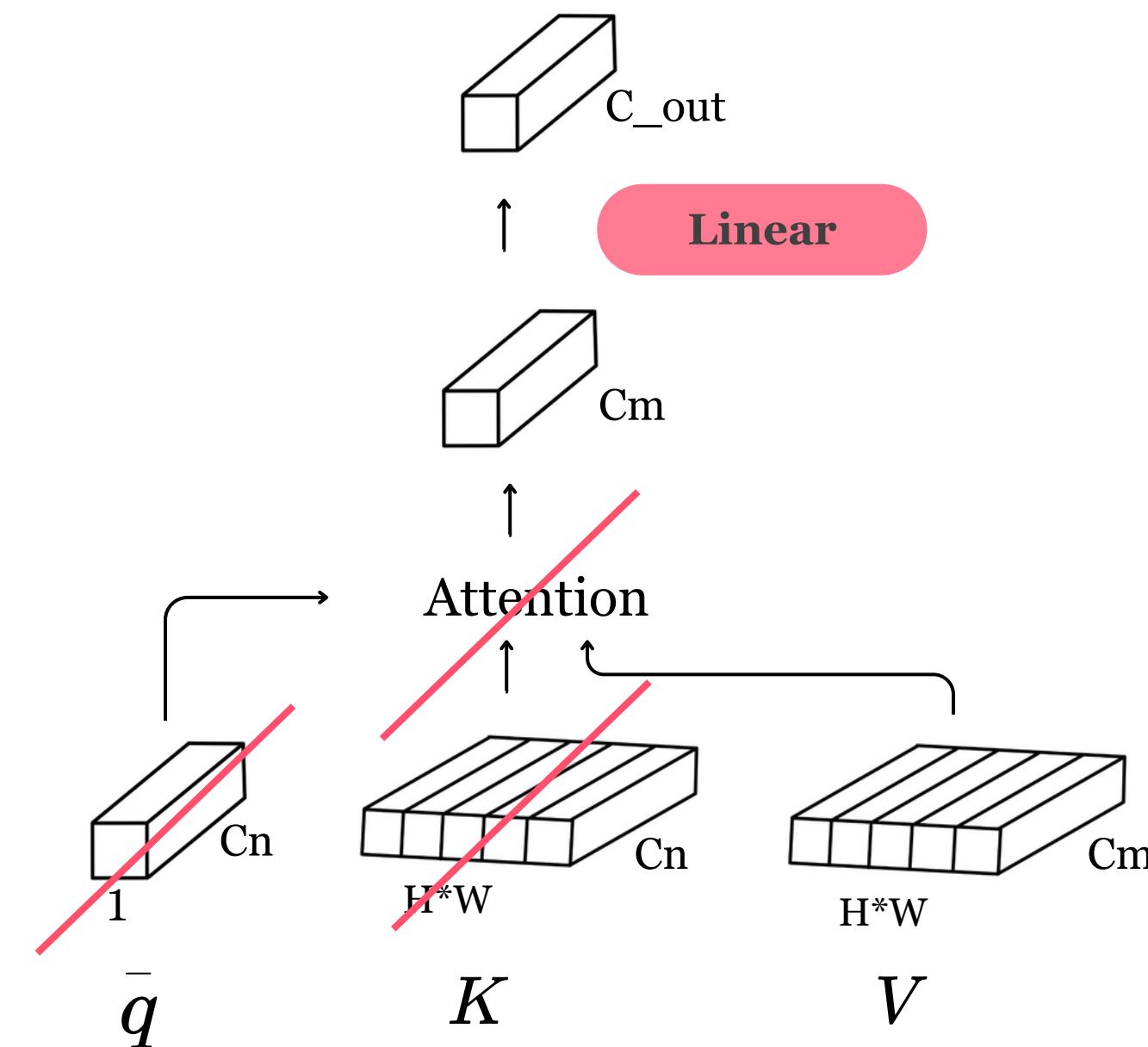


$$\begin{aligned} \text{AttnPool}(\bar{q}, k, v) &= \mathcal{F}\left(\sum_i \text{softmax}\left(\frac{\bar{q}k_i^\top}{C}\right)v_i\right) \\ &= \sum_i \text{softmax}\left(\frac{\bar{q}k_i^\top}{C}\right)\mathcal{F}(v_i), \end{aligned}$$

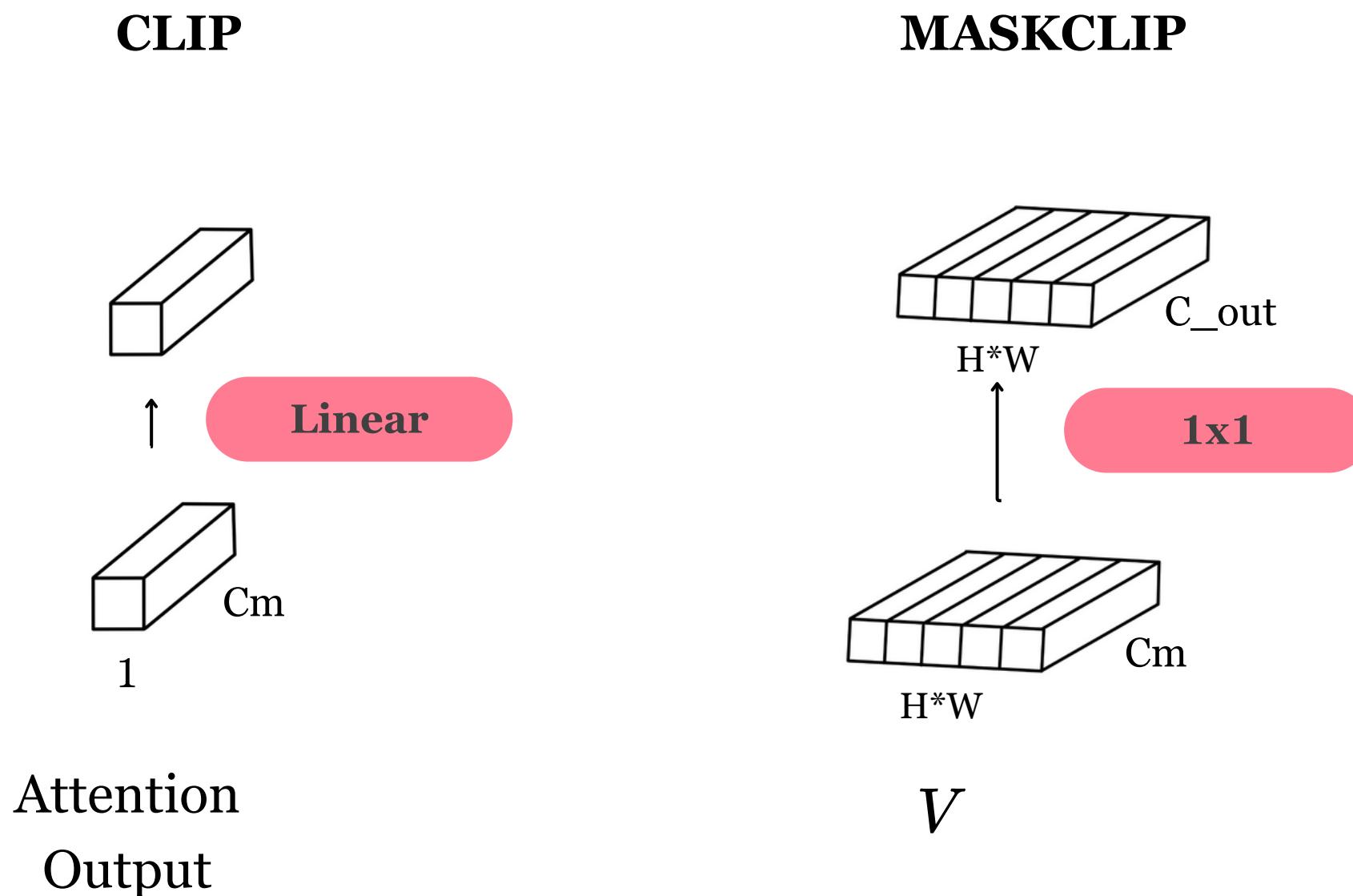
CLIP: Global Attention Pooling



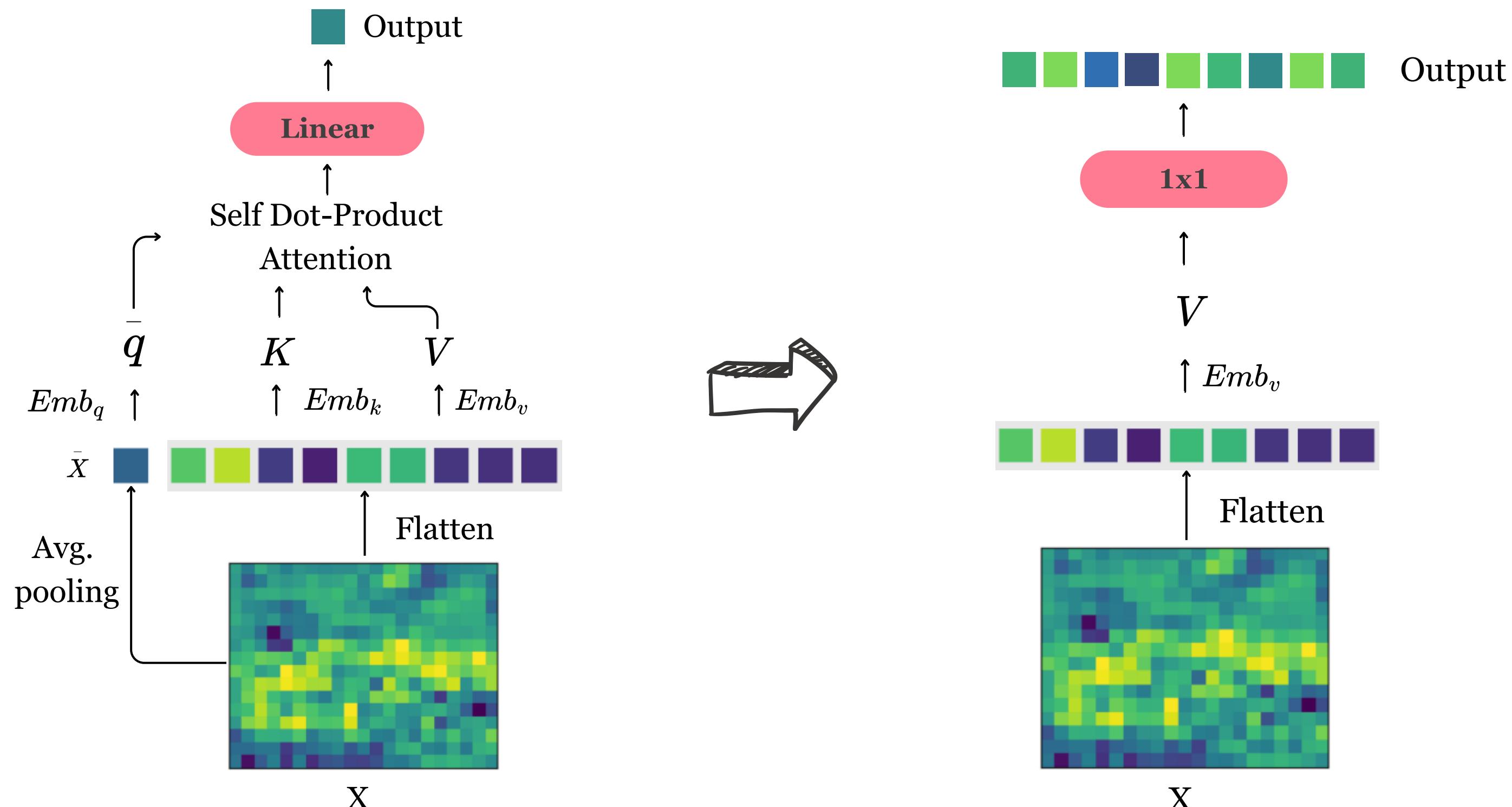
Remove Query and Key Embedding Layers



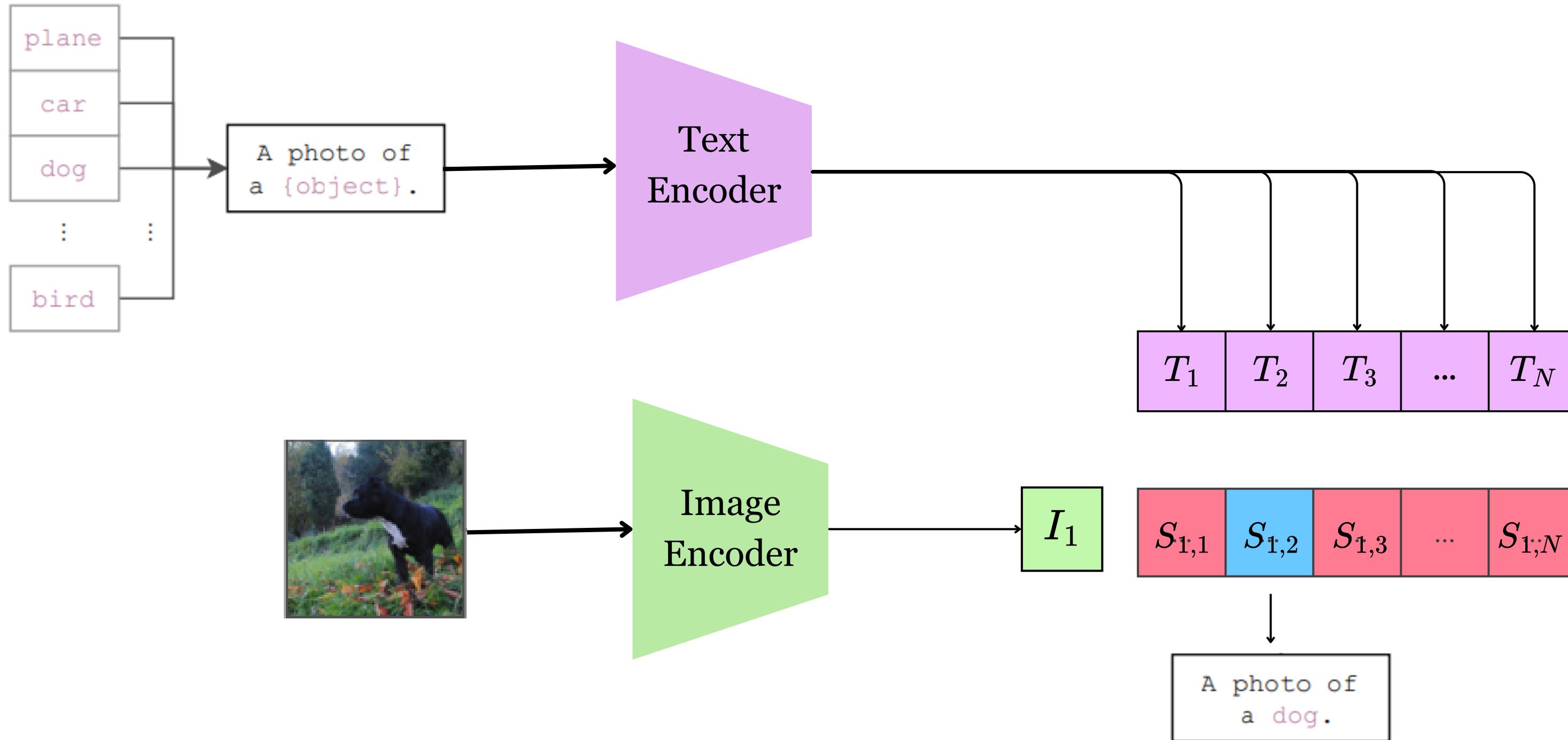
Reformulate Linear Layer to 1x1 Convolution



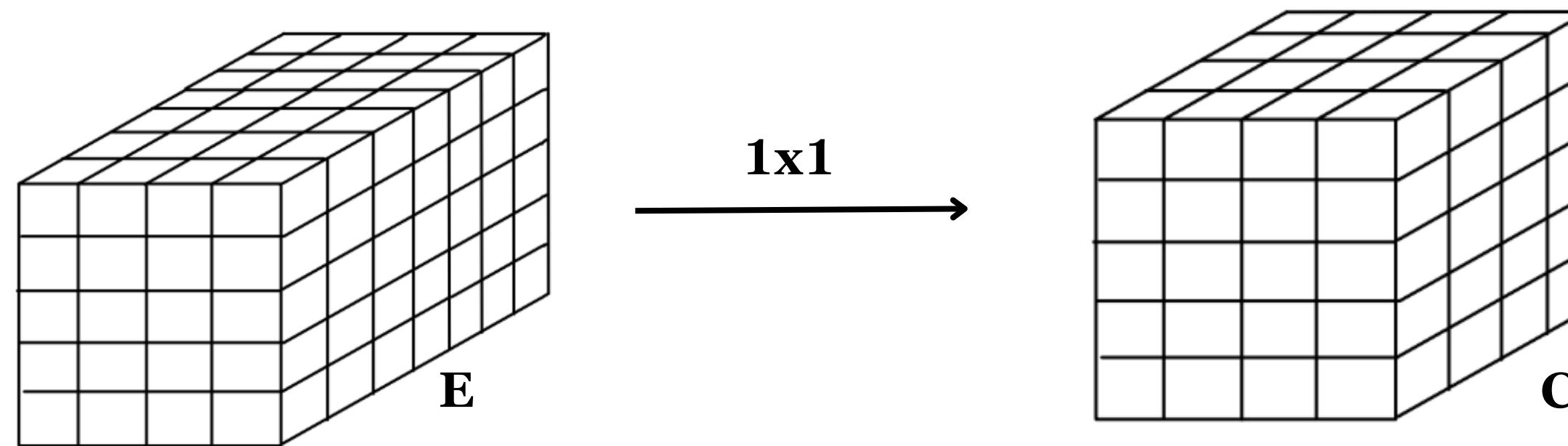
MaskCLIP: Image Encoder



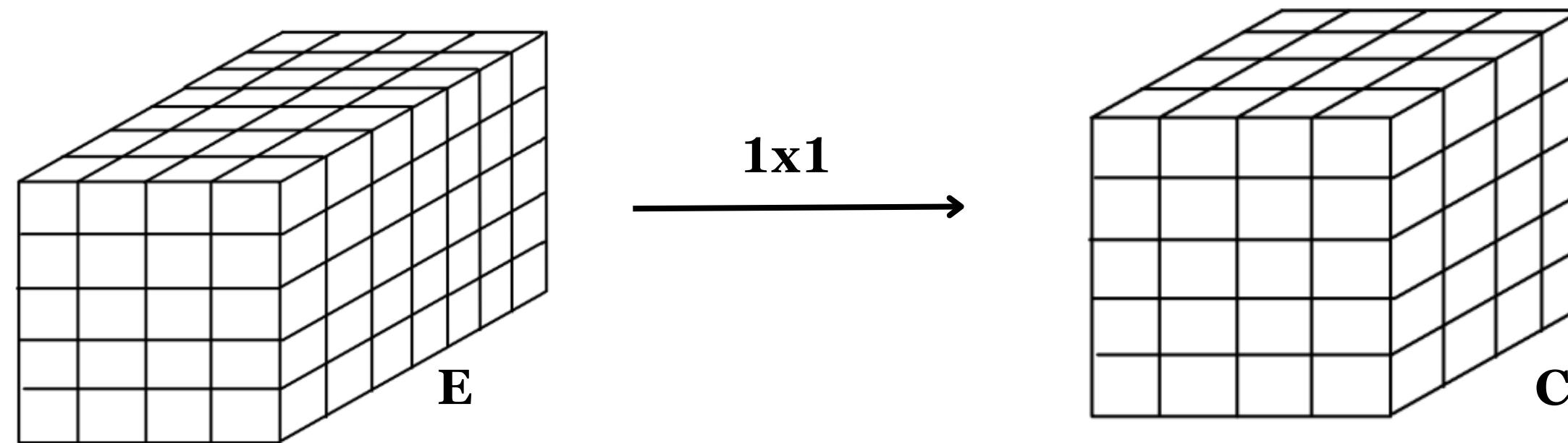
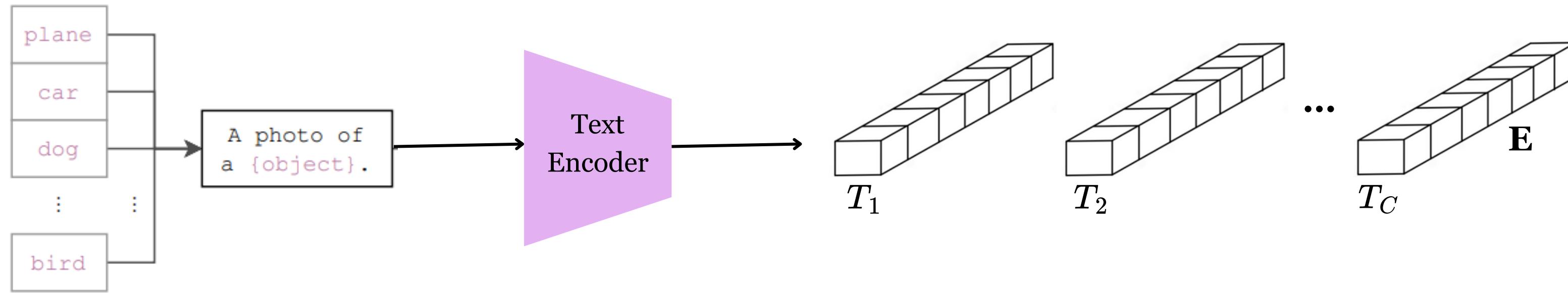
CLIP: Inference



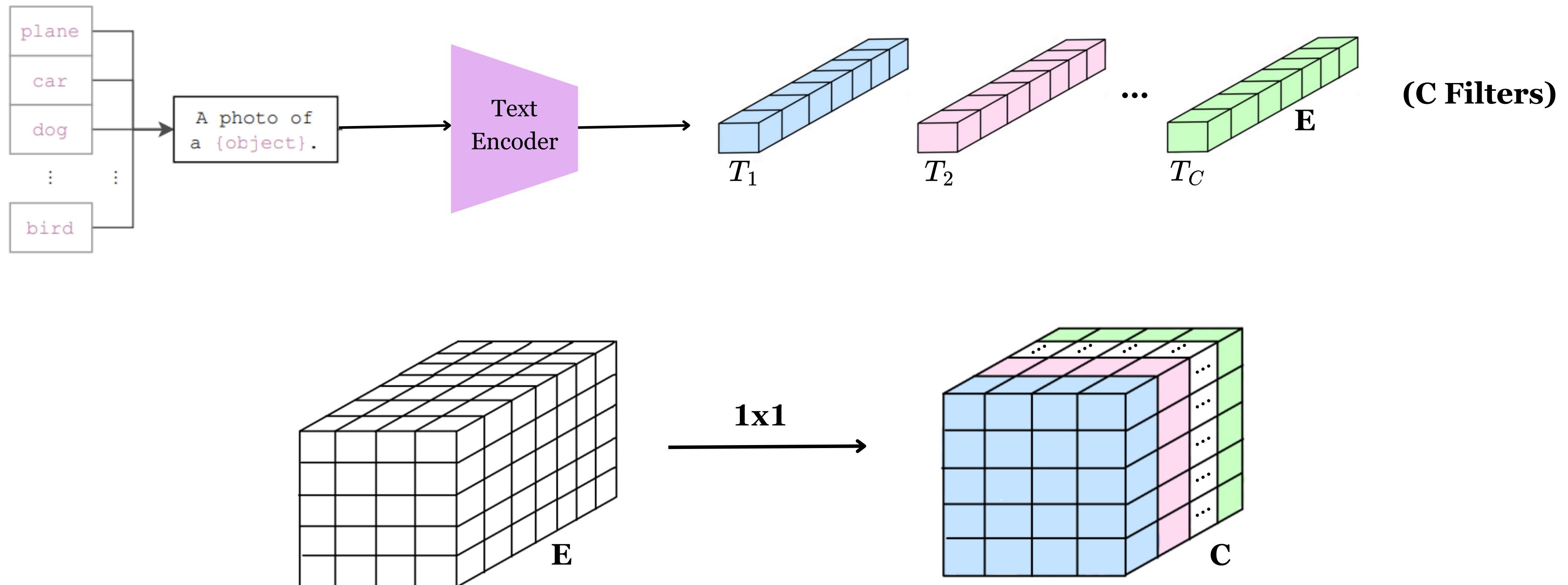
MaskCLIP: CLIP Text Embeddings as Classifier



MaskCLIP: CLIP Text Embeddings as Classifier



MaskCLIP: CLIP Text Embeddings as Classifier



MaskCLIP Refinement Techniques

- KS: Key smoothing

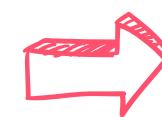
$$\text{pred}_i = \sum_j \cos\left(\frac{k_i}{\|k_i\|_2}, \frac{k_j}{\|k_j\|_2}\right) \text{pred}_j,$$

- PD: Prompt denoising

MaskCLIP+: Use Segmentation Architectures

Image encoder of CLIP → Segmentation Architectures

How??



Backbone + Segmentation Head 

MaskCLIP+: Use Segmentation Architectures

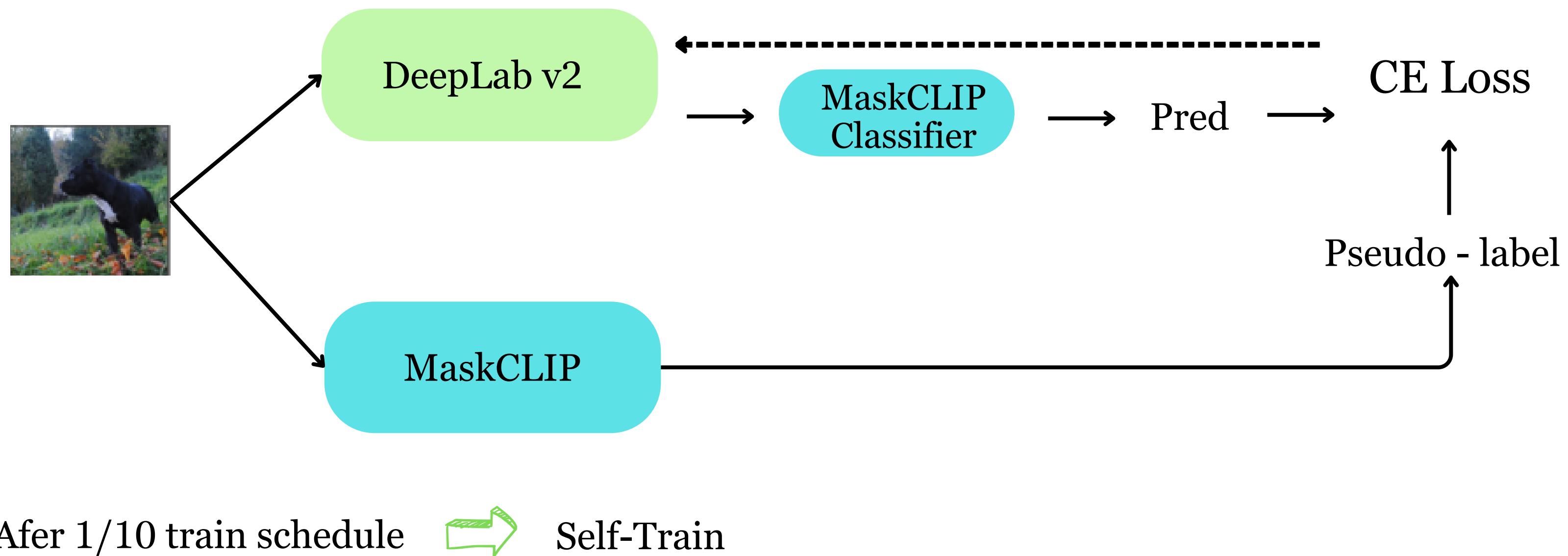
Image encoder of CLIP → Segmentation Architectures

How??

➡ Backbone + Segmentation Head 

➡ Use MaskCLIP predictions as a **train-time pseudo GT labels**

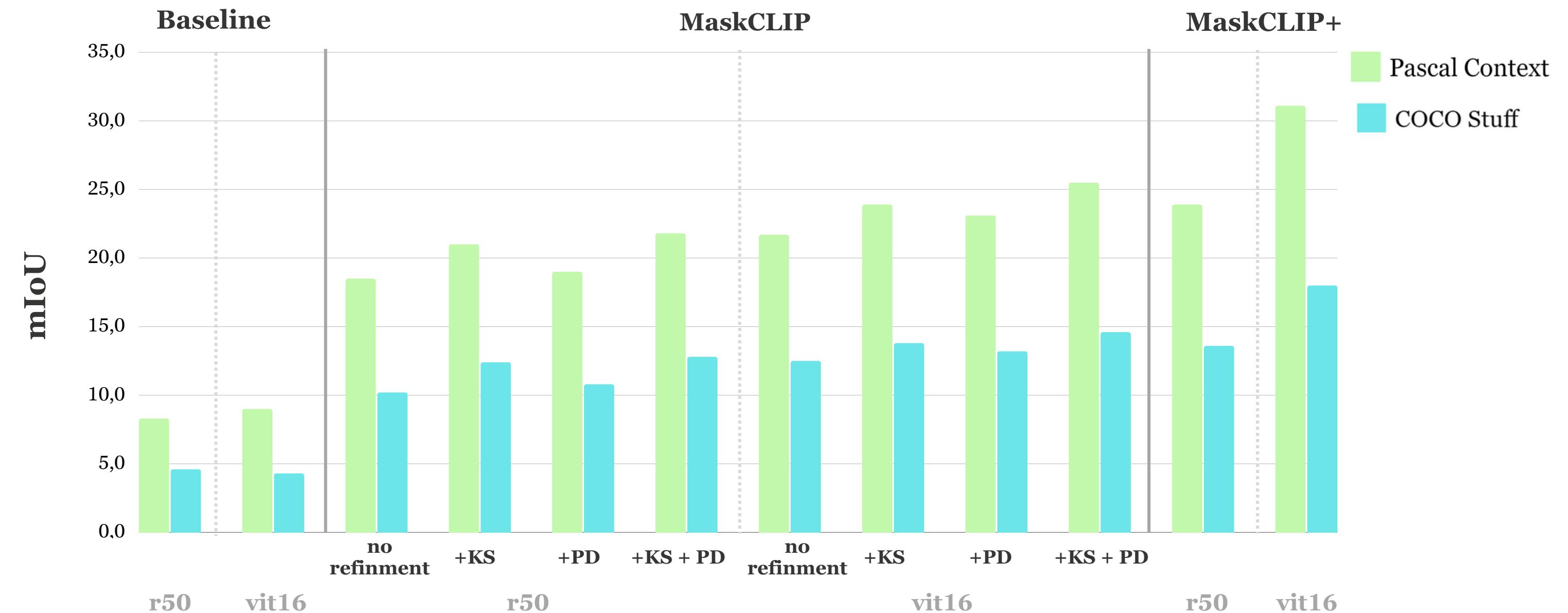
MaskCLIP+: Pseudo-labels and Self-trainning



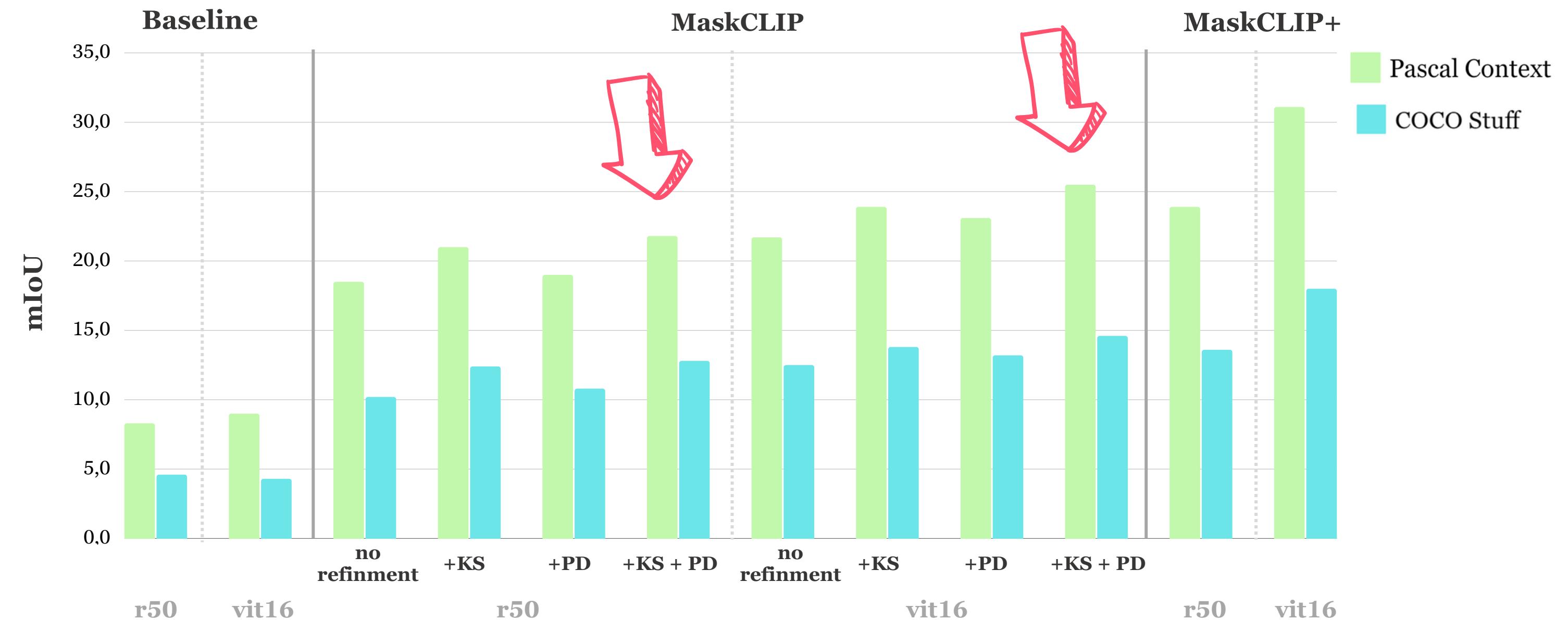
EXPERIMENTS AND RESULTS

Annotation Free Segmentation

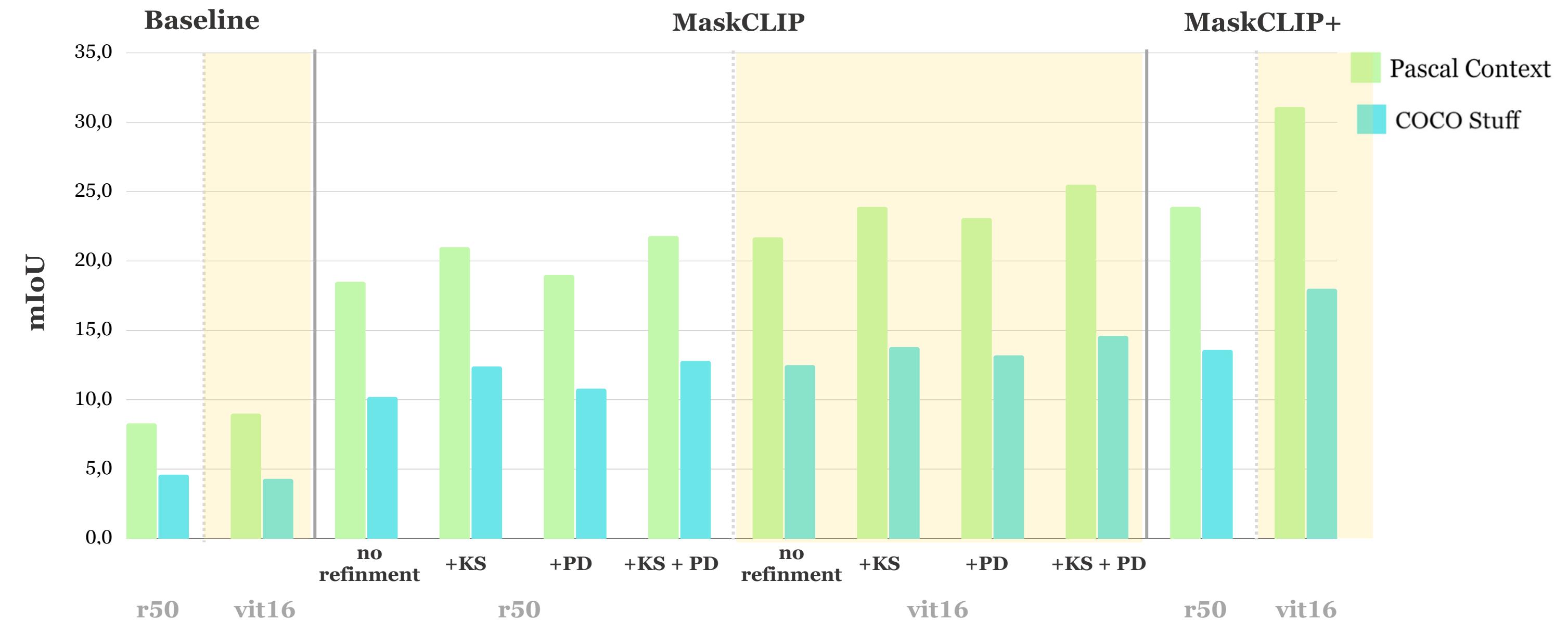
Baseline << MaskCLIP << MaskCLIP+



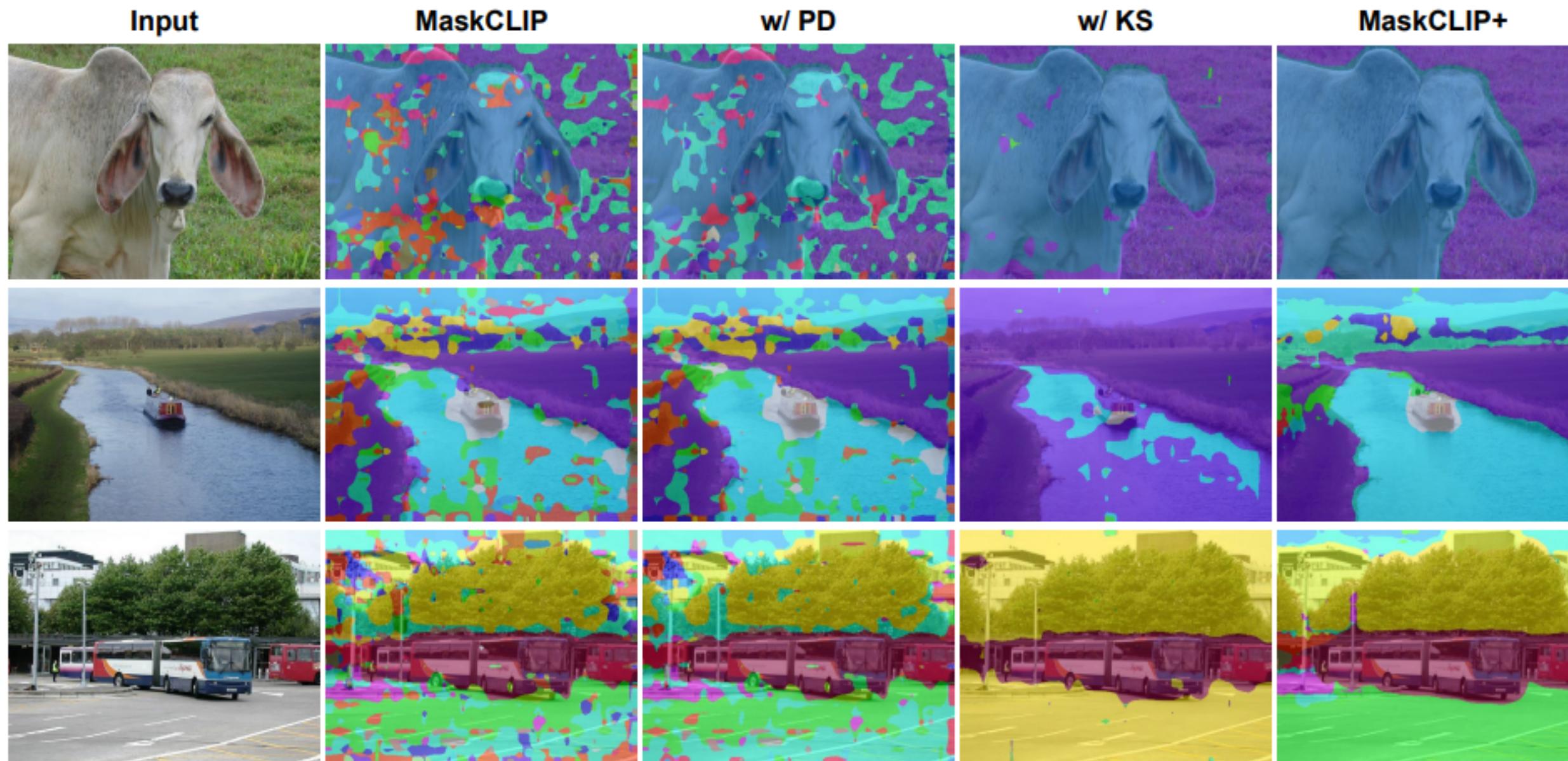
KS and PD are effective and orthogonal



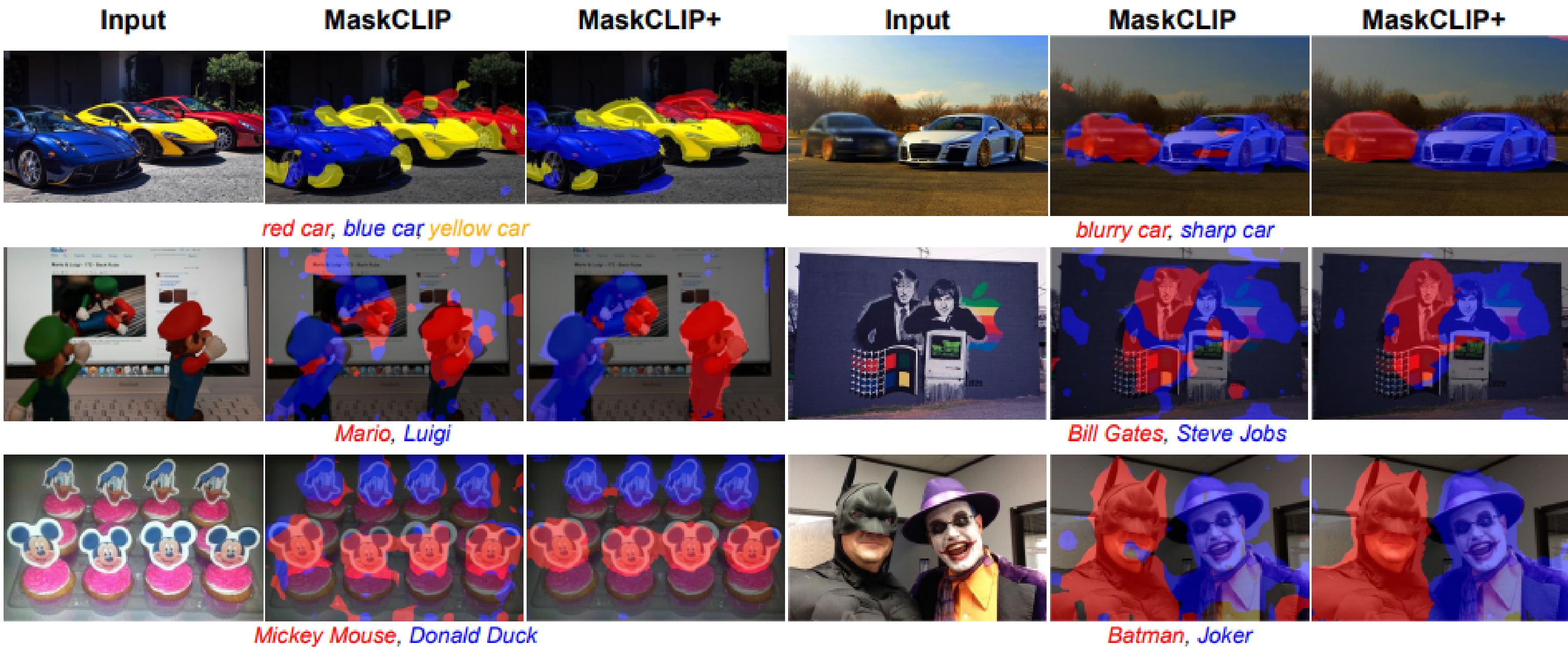
ViT models surpass ResNet models



MaskCLIP+ achieves the best performance

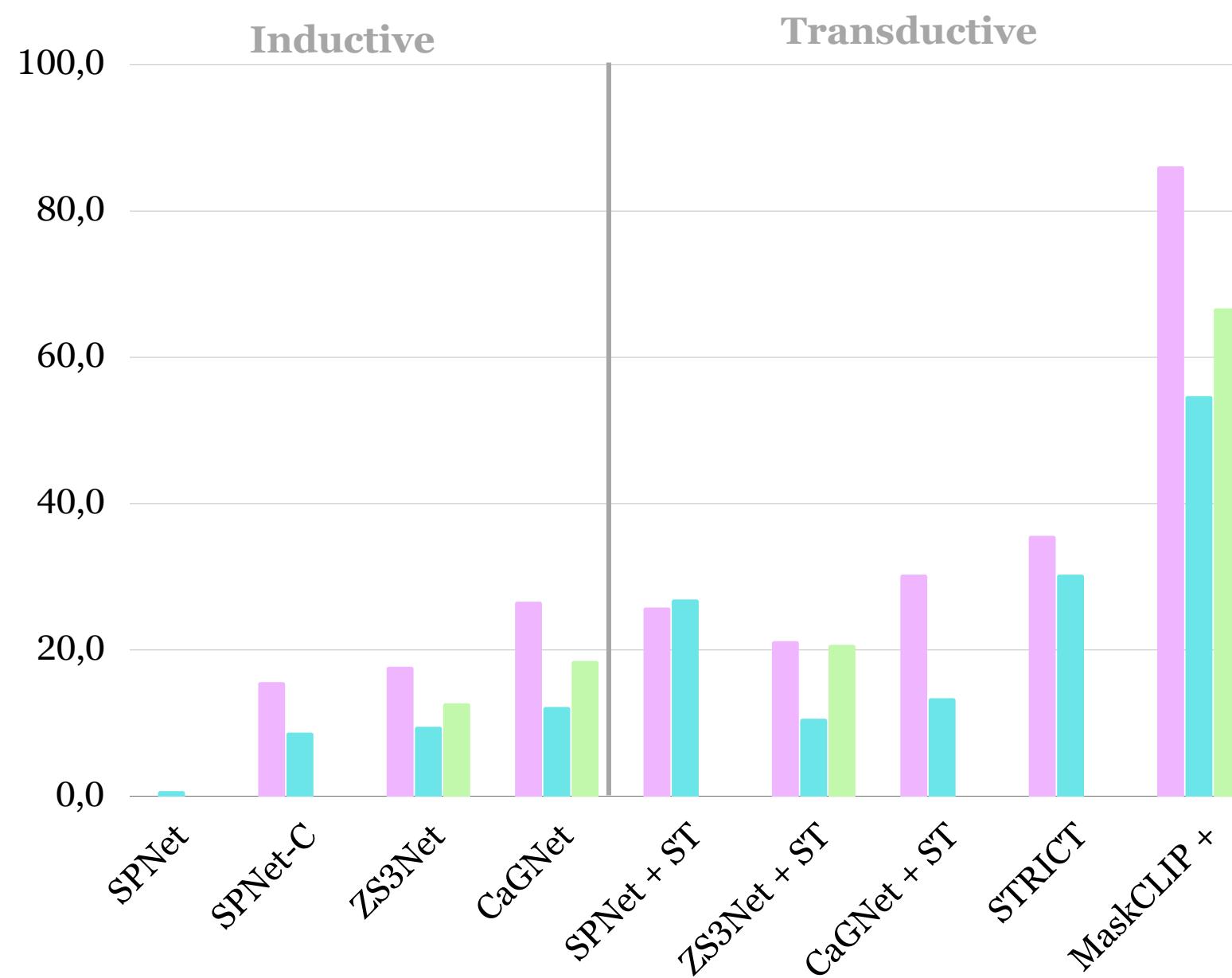


Open-vocabulary Ability

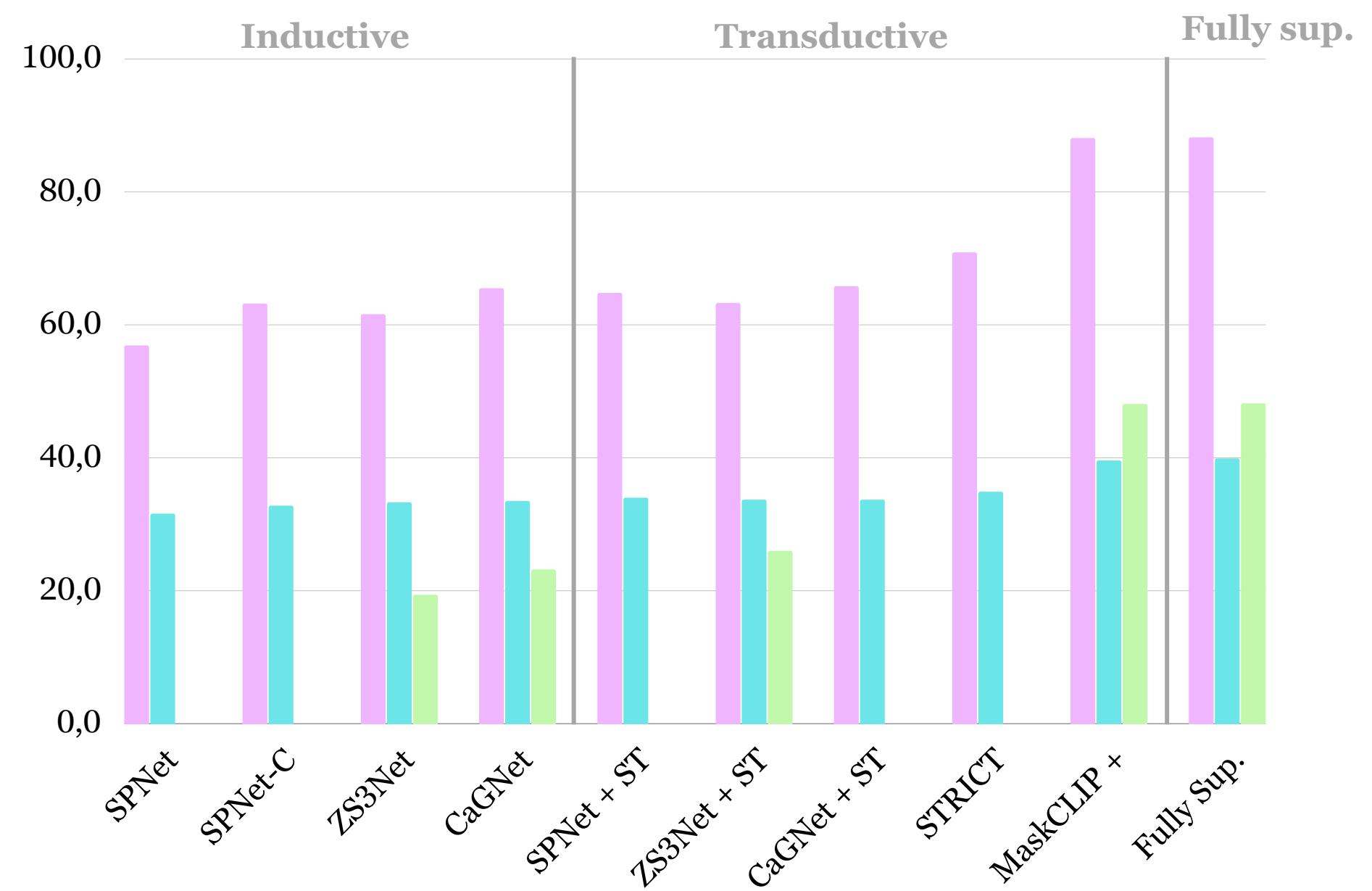


ZeroShot Segmentation

mIoU (unseen)

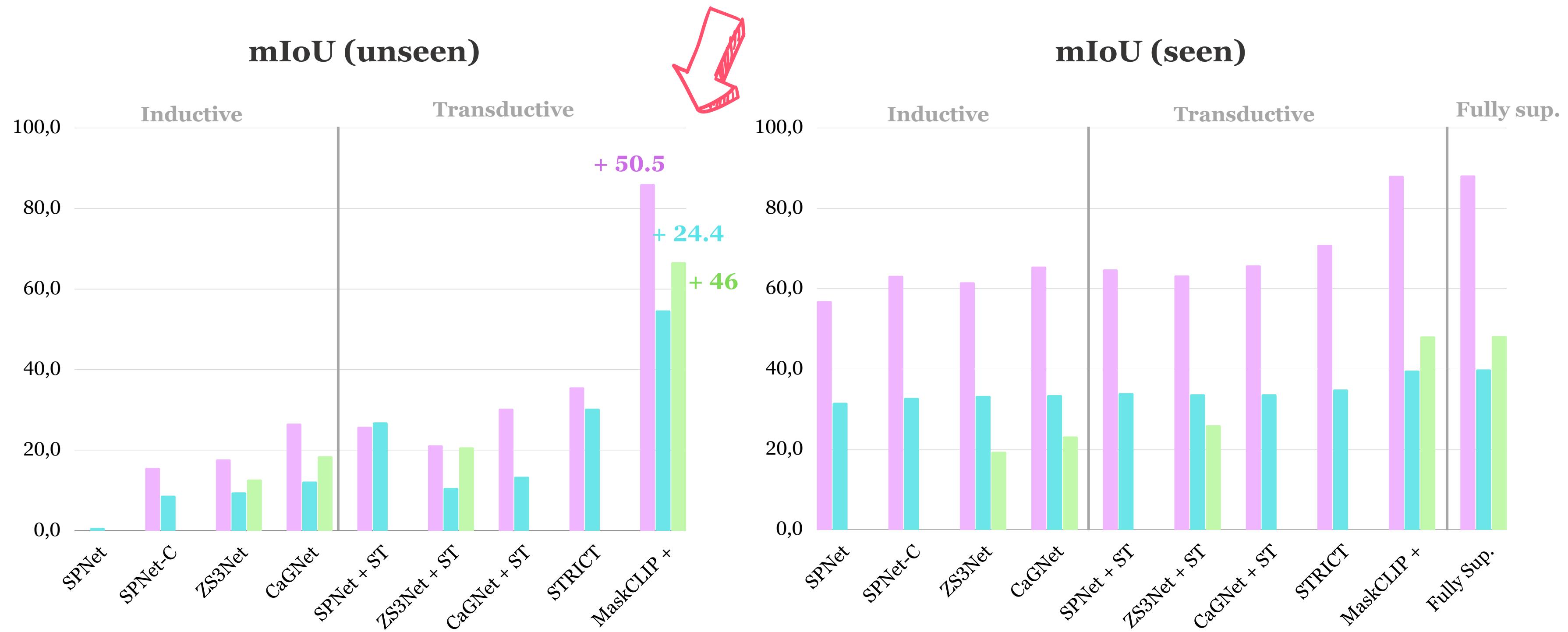


mIoU (seen)



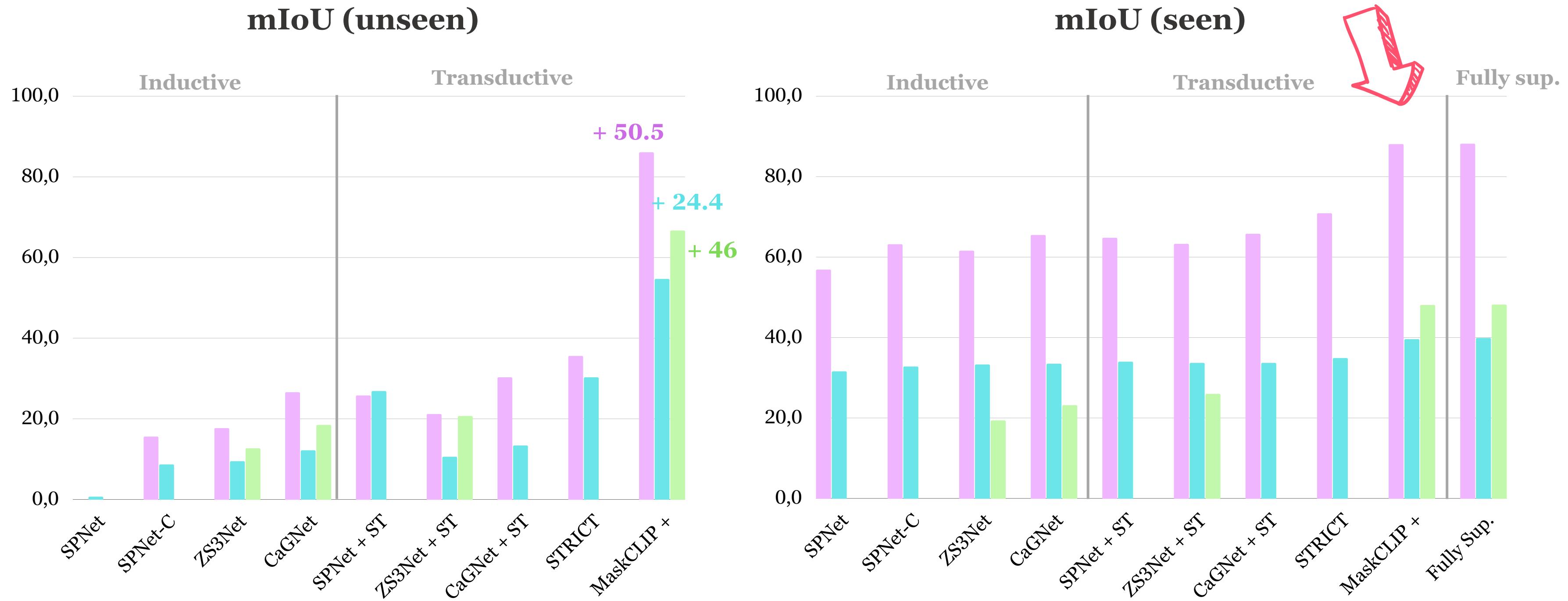
SOTA in unseen classes by large margins

█ Pascal VOC
█ COCO Stuff
█ Pascal Context



On par to fully supervised BL in seen classes

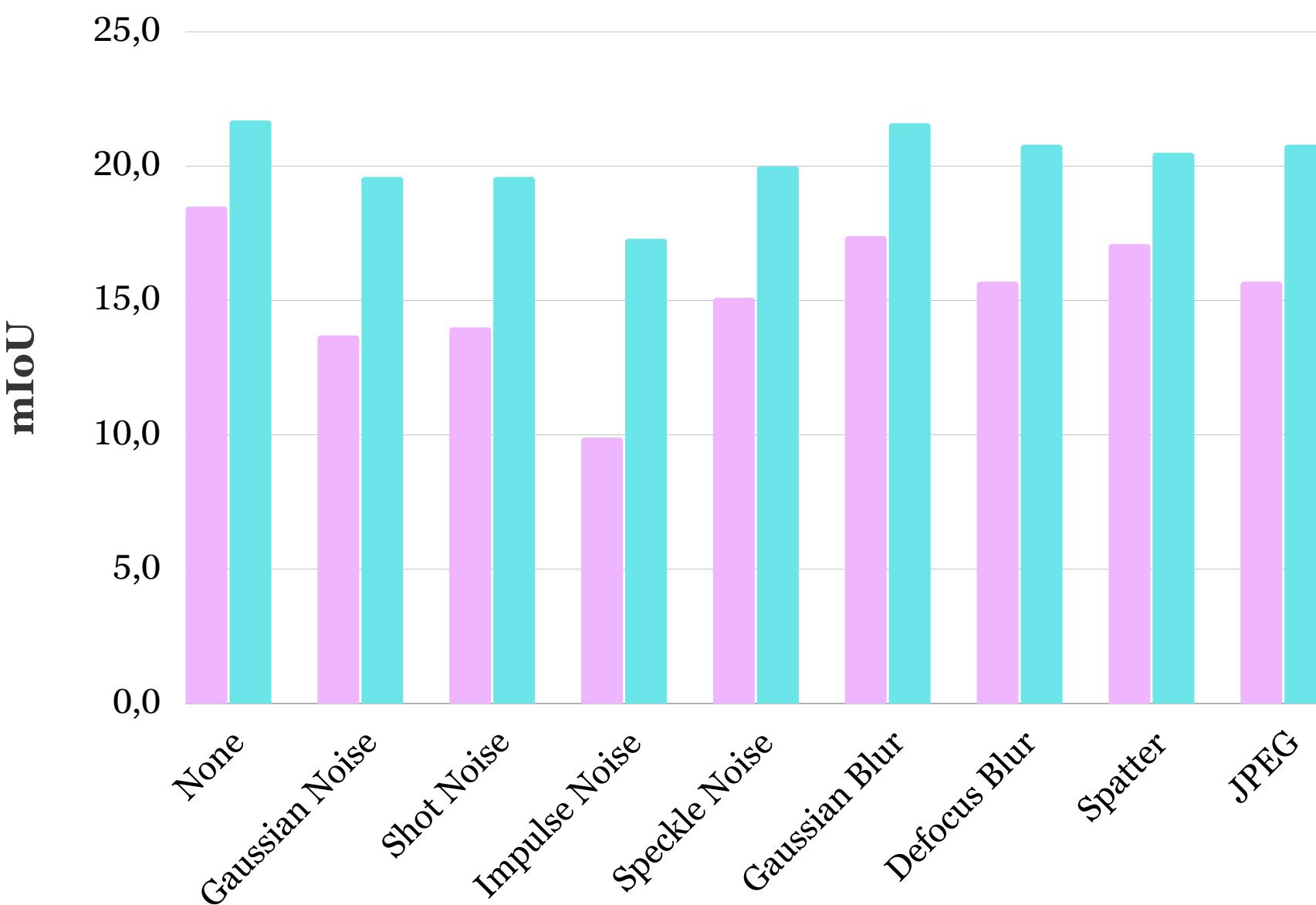
█ Pascal VOC
█ COCO Stuff
█ Pascal Context



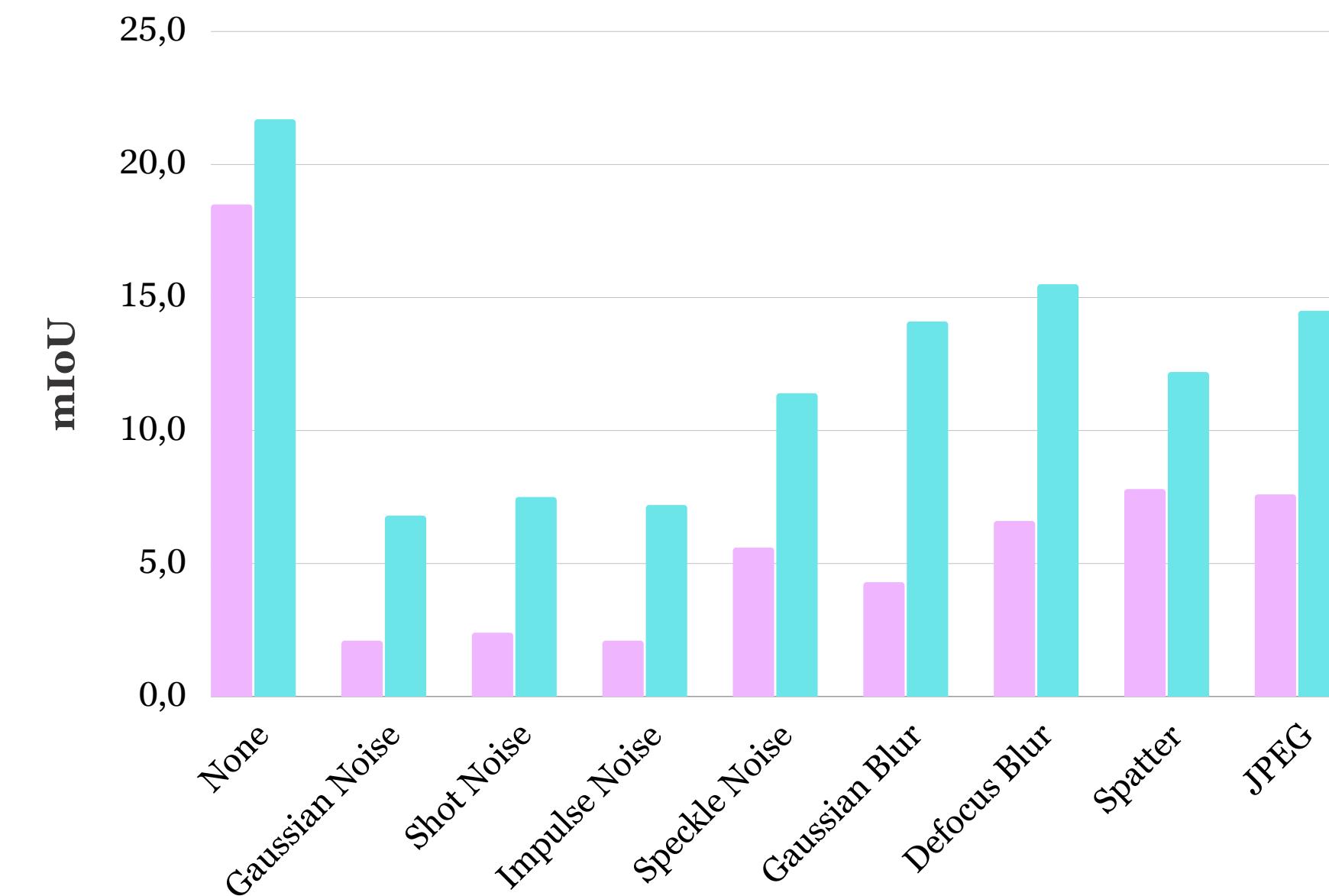
ViT16 models are more robust than ResNet50

r50 vit16

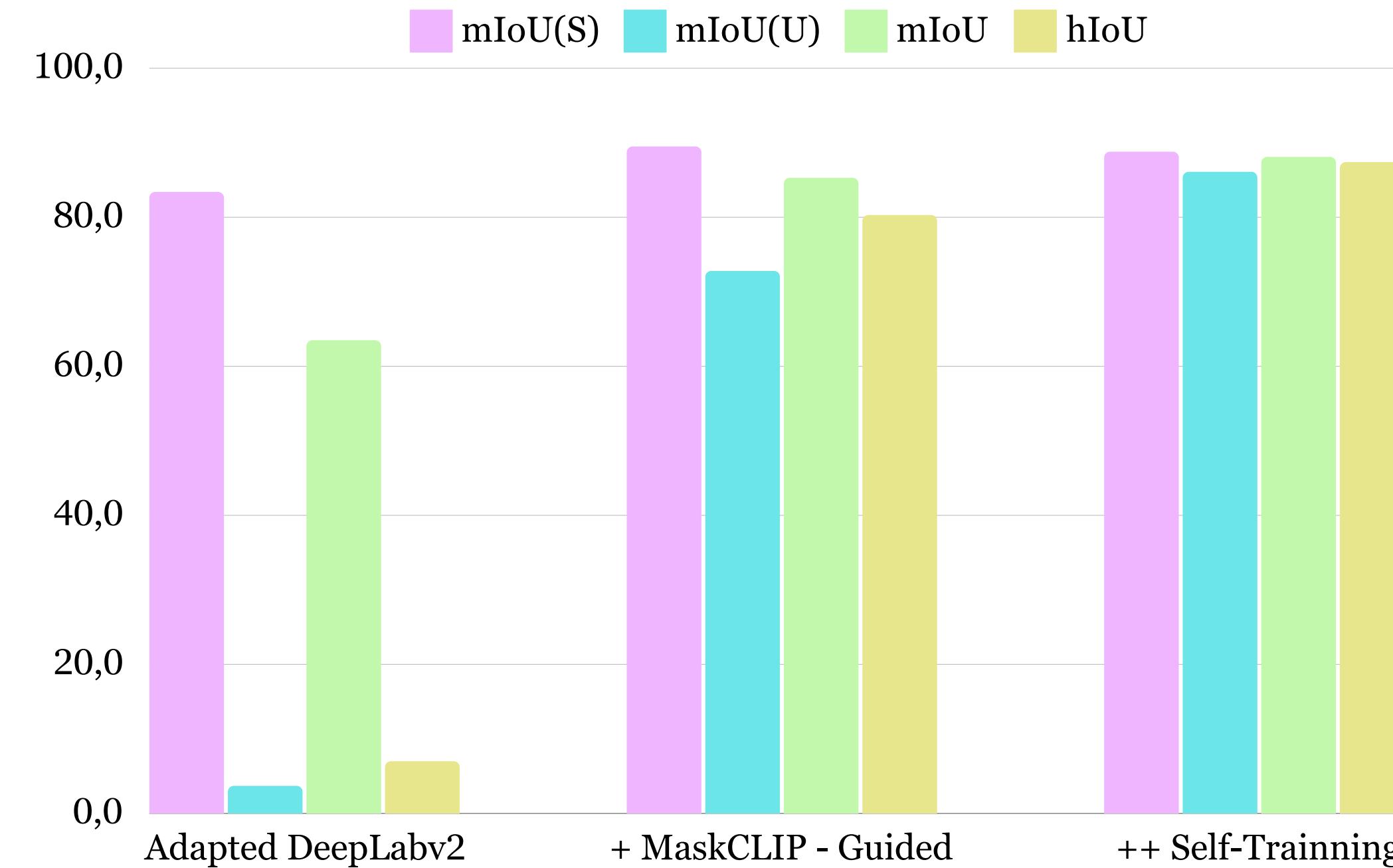
Level 1



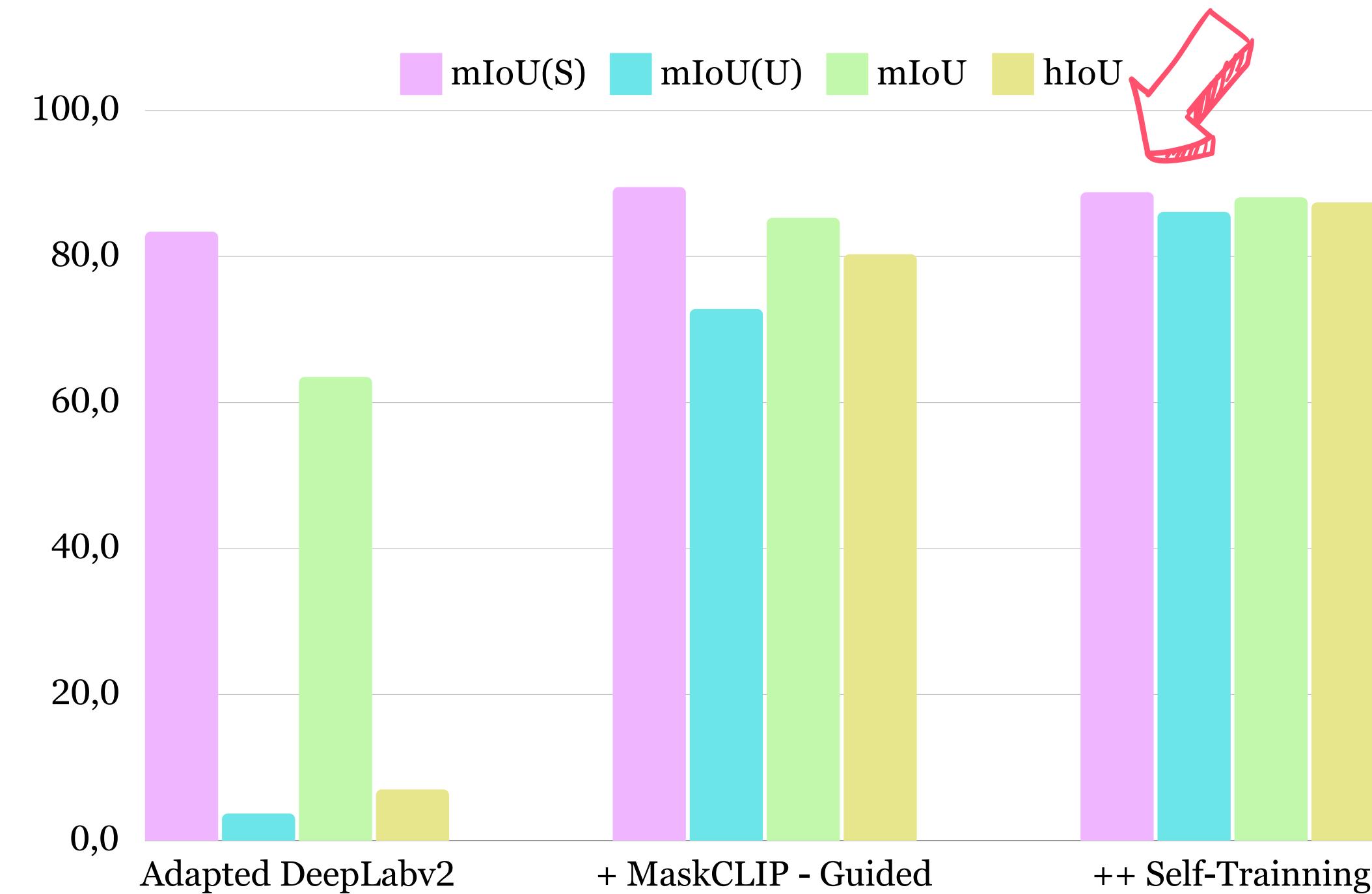
Level 5



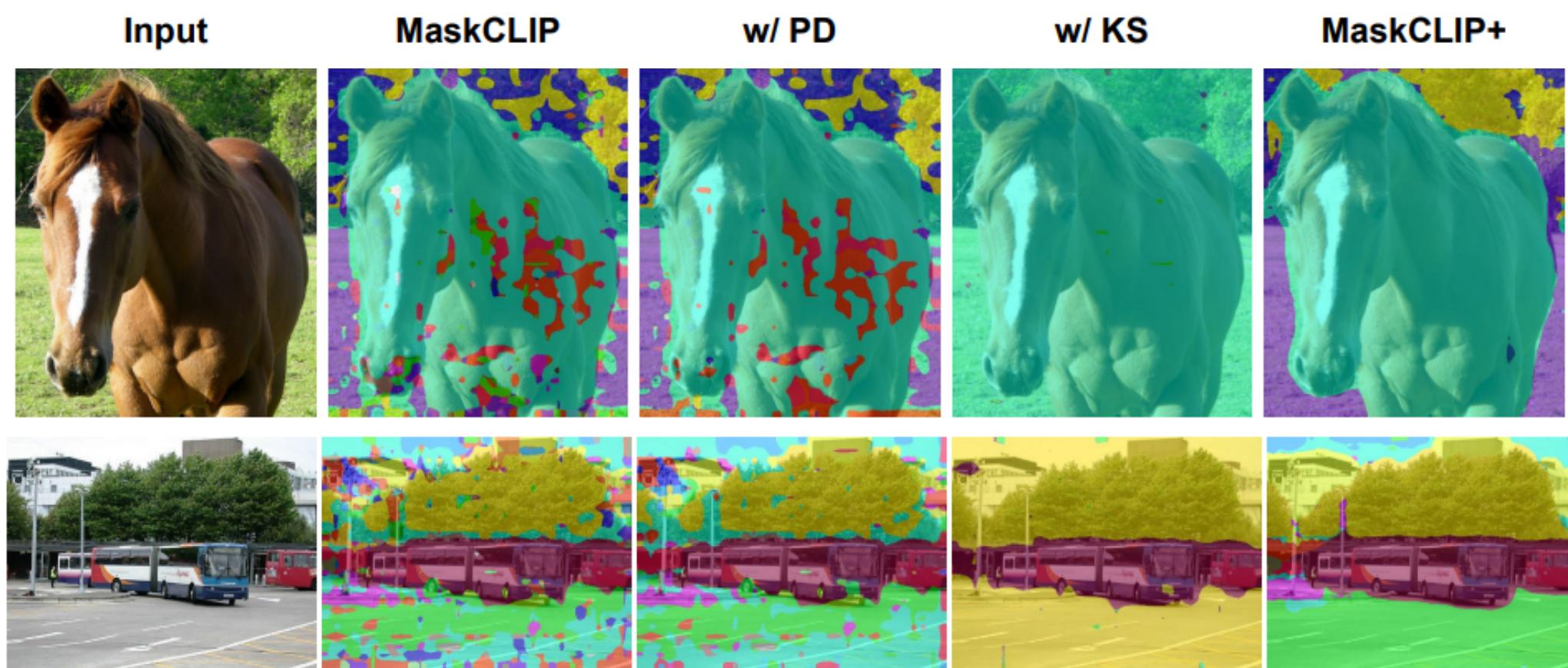
MaskCLIP-guiding and Self-trainning are Effective



Slight degradation on seen classes when Self-training



Failure Cases



MaskCLIP Open Vocabulary



CONCLUSIONS

Architectures Summary

- **MaskCLIP**

- Image encoder: V embedings + 1x1

- Classifier - CLIP Text Embeddings

- KS + PD

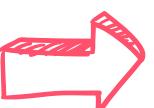
- **MaskCLIP+**

- Segmentation model

- MaskCLIP pseduo-labels

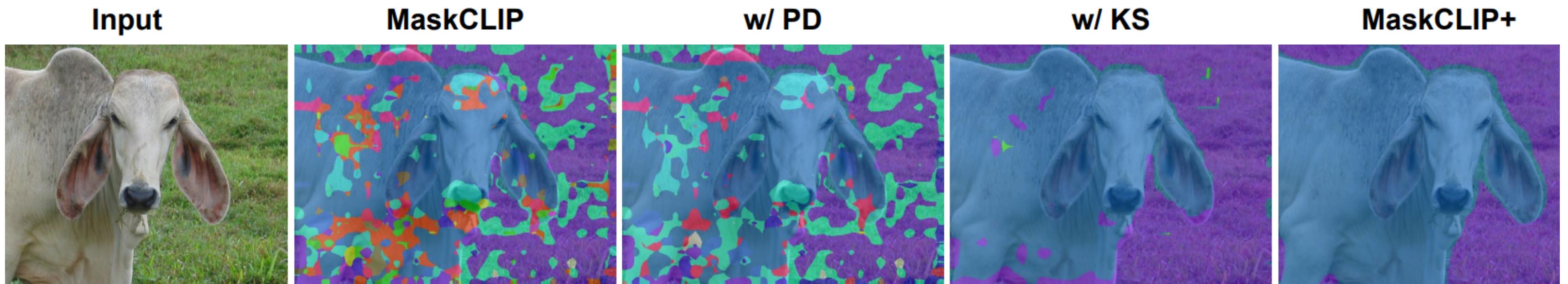
- Self-Train

Towards Cost-Effective Semantic Segmentation

- Expensive Annotations  Text-image Pairs
- Expensive Annotations  MaskCLIP: Free Segmentation Annotator
- Open Vocabulary

Extract Free Dense Labels from CLIP

Authors: Chong Zhou, Chen Change Loy , and Bo Dai



Sonia Castro Paniello
20.02.24