

LEARNING TO GENERATE TEXT-GROUNDED MASK FOR OPEN-WORLD SEMANTIC SEGMENTATION FROM ONLY IMAGE-TEXT PAIRS

Junbum Cha, Jonghwan Mun, Byungseok Roh

kakaobrain

TARGET TASK

Open-world Semantic Segmentation

learning an universal model to segment **arbitrary concepts** beyond pre-defined categories

from Only Image-Text Pairs

using only image-text pairs **without any segmentation annotation**

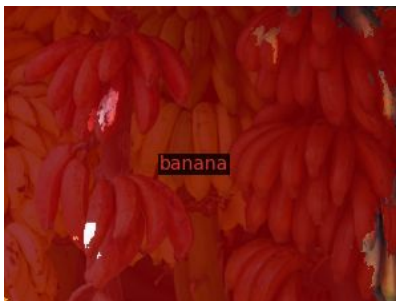
OPEN-WORLD SEMANTIC SEGMENTATION

- Open-world segmentation model can segment **different dog species**,



OPEN-WORLD SEMANTIC SEGMENTATION

- Open-world segmentation model can segment different dog species, **bananas by color**,



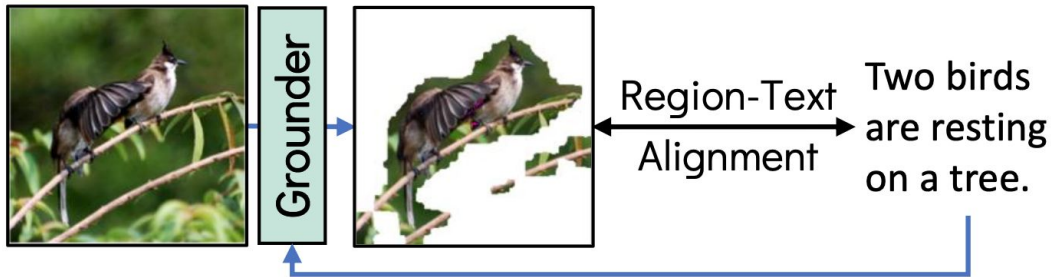
OPEN-WORLD SEMANTIC SEGMENTATION

- Open-world segmentation model can segment different dog species, bananas by color, and **even proper nouns** such as Frodo, Gollum, and Samwise



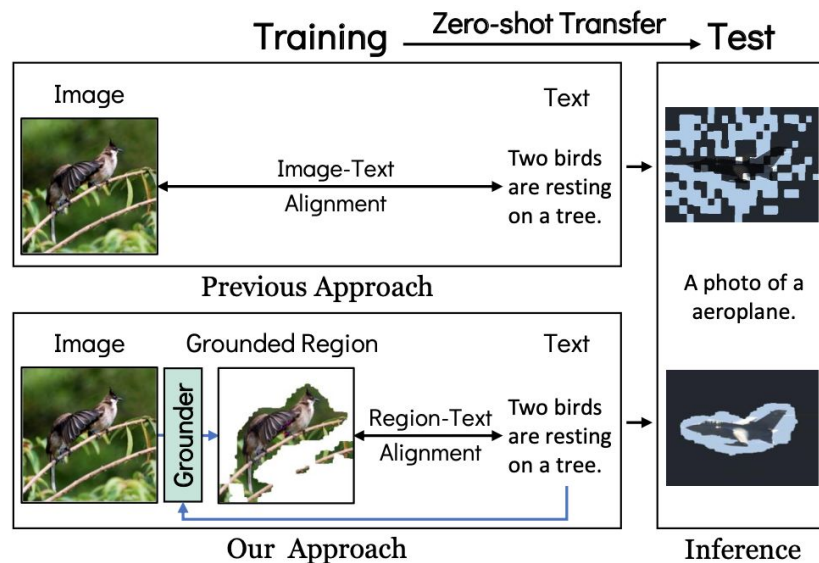
REGION-TEXT ALIGNMENT

- Open-world segmentation is conducted by region-text alignment
- In inference, the model identifies regions in the image that align with the given text queries



TRAIN-TEST DISCREPANCY IN EXISTING WORKS

All of the existing methods train the model via **image-text alignment**, even though the target task requires **region-text alignment**



TCL: TEXT-GROUNDED CONTRASTIVE LEARNING

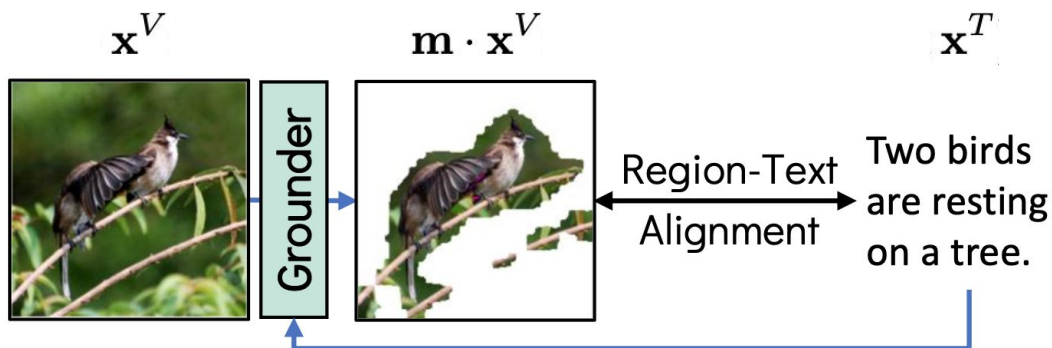
We propose Text-grounded Contrastive Learning (TCL) to let the model learn region-text alignment instead of image-text alignment in training time

CL $\arg \max_{\theta} I_{\theta}(\mathbf{x}^V; \mathbf{x}^T)$



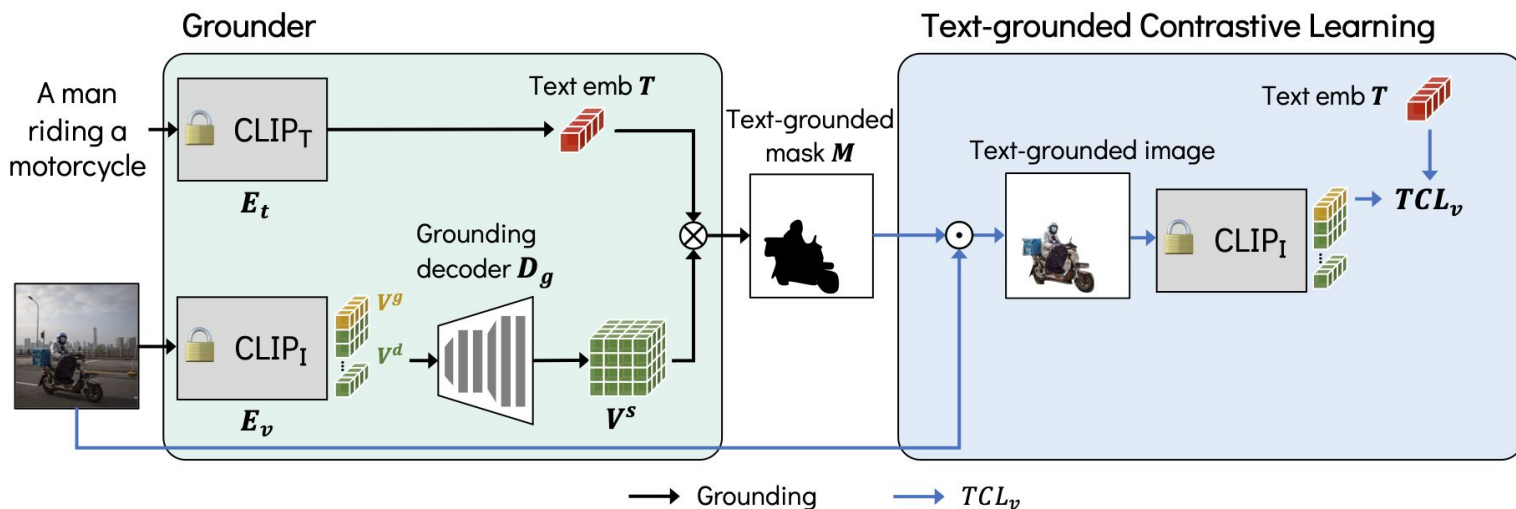
TCL $\arg \max_{\theta} I_{\theta}(\mathbf{m} \cdot \mathbf{x}^V; \mathbf{x}^T)$

Text-grounded mask



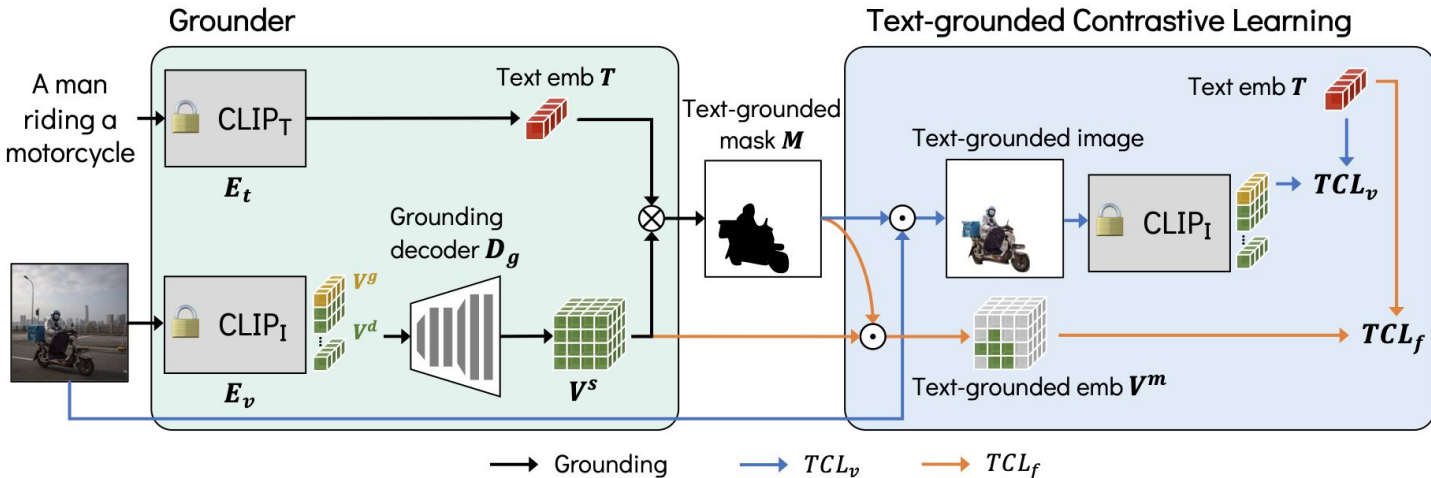
TCL FRAMEWORK

1. Grounder generates text-grounded mask M
2. TCL loss is computed by incorporating text-grounded image into contrastive learning



FEATURE-LEVEL TCL LOSS

- (Image-level) TCL loss only can consider “positive” text-grounded masks
- We introduce feature-level TCL loss to incorporate “negative” masks into our objective



PREVENTING TRIVIAL SOLUTION



- There is a trivial solution in this framework: generating full-mask independent of the text
- To prevent this trivial solution, we introduce area TCL loss:

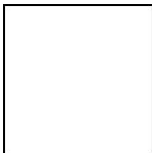


A man riding
a motorcycle

M^+



M^-



$$\mathcal{L}_{\text{area}} = \left\| p^+ - \underbrace{\mathbb{E} [\overline{M^+}]} \right\|_1 + \left\| p^- - \underbrace{\mathbb{E} [\overline{M^-}]} \right\|_1$$

Expectation of **positive**
text-grounded mask area

Expectation of **negative**
text-grounded mask area

SMOOTH REGULARIZATION

Finally, we introduce a smooth regularization loss via total variation (TV) loss:

$$\mathcal{L}_{\text{TV}} = \underbrace{\|\mathbf{M}\|_{\text{TV}}}_{\text{text-grounded mask}} + \underbrace{\|\mathbf{V}^s\|_{\text{TV}}}_{\text{dense image embedding}}$$

where

$$\|\mathbf{y}\|_{\text{TV}} = \sum_{i,j} \sqrt{|y_{i+1,j} - y_{i,j}|^2} + \sqrt{|y_{i,j+1} - y_{i,j}|^2} = \sum_{i,j} |y_{i+1,j} - y_{i,j}| + |y_{i,j+1} - y_{i,j}|$$

SUM UP

We introduce TCL loss to let the model learn **region-text alignment**, instead of image-text alignment

$$\mathcal{L}_{\text{final}} = \underbrace{\lambda_{\text{TCL}}(\mathcal{L}_{\text{TCL}_v} + \mathcal{L}_{\text{TCL}_f}) + \lambda_{\text{area}}\mathcal{L}_{\text{area}}}_{\text{TCL losses}} + \underbrace{\lambda_{\text{tv}}\mathcal{L}_{\text{tv}}}_{\text{regularization}}$$

image-level TCL feature-level TCL area TCL smooth regularization

ZERO-SHOT EVALUATION PROTOCOL

- Since the open-world semantic segmentation task is introduced recently, evaluation protocols vary across studies
- For a fair comparison, we present a unified evaluation protocol

Prohibited examples

floor-other
floor-tile
floor-wood
floor-stone

} floor

Class integration



Crop

vegetation -> tree

Rephrasing

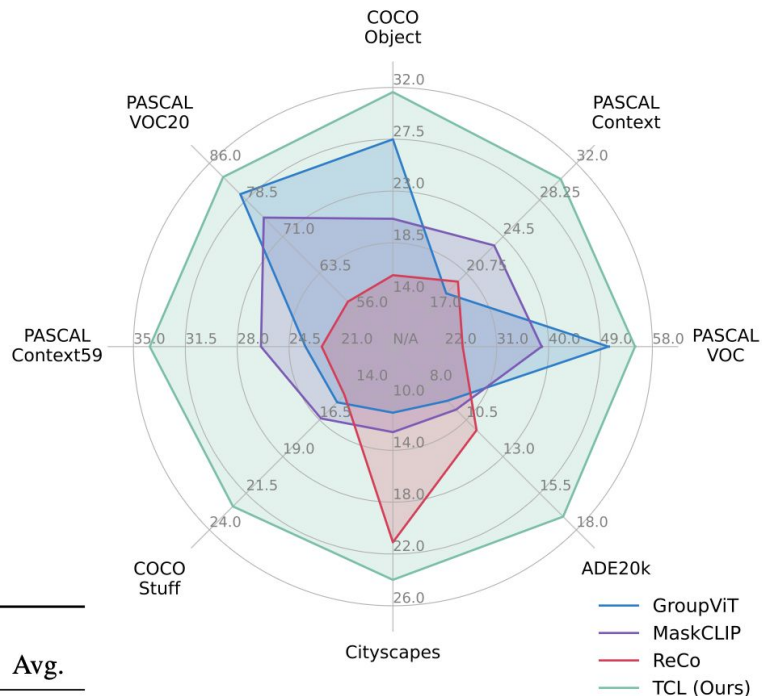
VOC: threshold=0.95
Context: threshold=0.35
...

Dataset-specific HPs

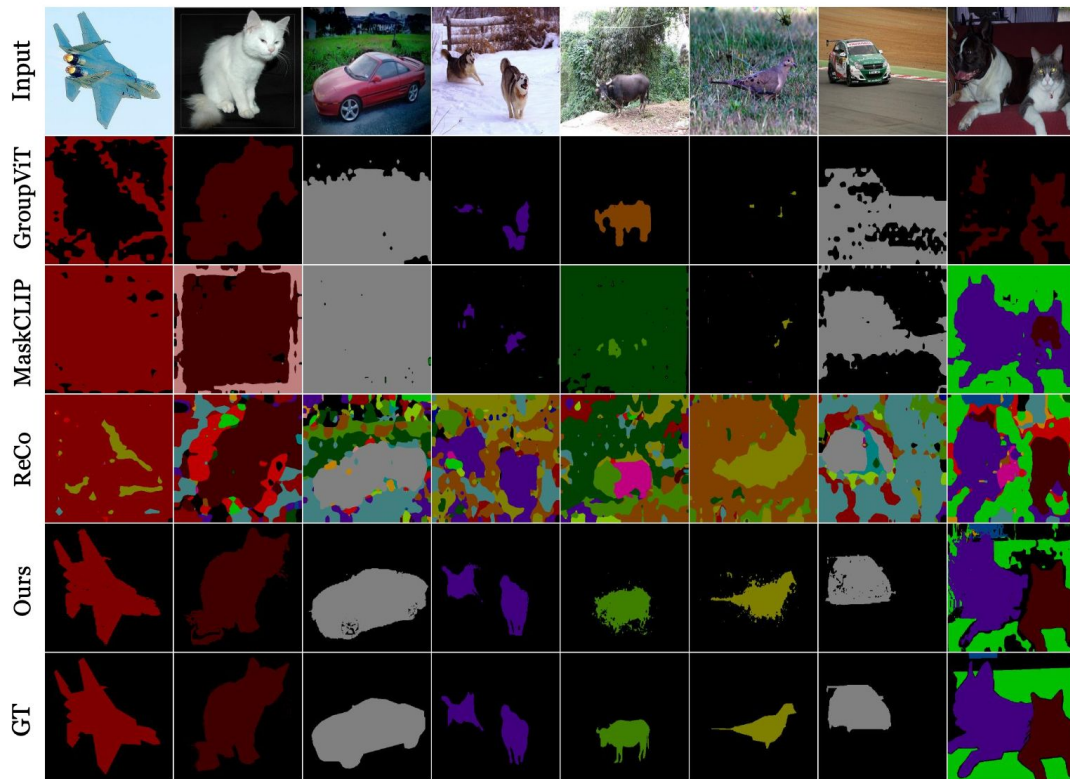
QUANTITATIVE RESULTS

TCL remarkably outperforms previous methods in every dataset

Methods	with background class			without background class					Avg.
	VOC	Context	Object	VOC20	Context59	Stuff	City	ADE	
GroupViT (YFCC)	49.5	19.0	24.3	74.1	20.8	12.6	6.9	8.7	27.0
GroupViT (RedCaps)	<u>50.4</u>	18.7	<u>27.5</u>	<u>79.7</u>	23.4	15.3	11.1	9.2	<u>29.4</u>
MaskCLIP [†]	29.3	21.1	15.5	53.7	23.3	14.7	<u>21.6</u>	10.8	23.7
MaskCLIP	38.8	<u>23.6</u>	20.6	74.9	<u>26.4</u>	<u>16.4</u>	12.6	9.8	27.9
ReCo	25.1	19.9	15.7	57.7	22.3	14.8	21.1	<u>11.2</u>	23.5
TCL (Ours)	55.0 (+4.6)	30.4 (+6.8)	31.6 (+4.1)	83.2 (+3.5)	33.9 (+7.5)	22.4 (+6.0)	24.0 (+2.4)	17.1 (+5.9)	37.2 (+7.8)

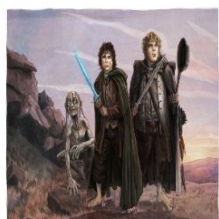


QUALITATIVE COMPARISON ON PASCAL VOC



QUALITATIVE EXAMPLES IN THE WILD

Input



TCL



Moon
Hill
Temple

Sunset
Buddist pagoda
House

Back view of
a standing bear
Rabbit
Many trees

Red banana
Green banana
Yellow banana

Eagle mark
MMU
Turkish
Fighter
Satellite

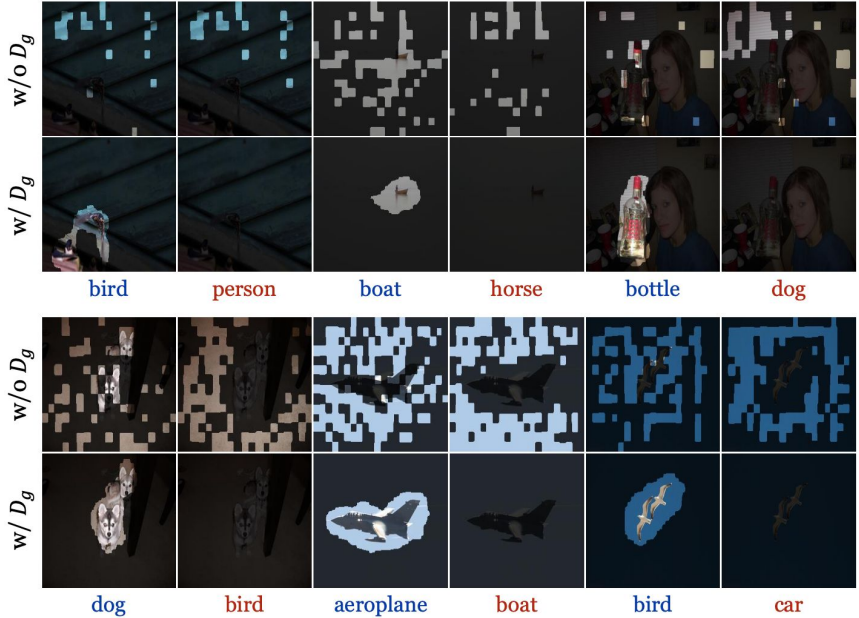
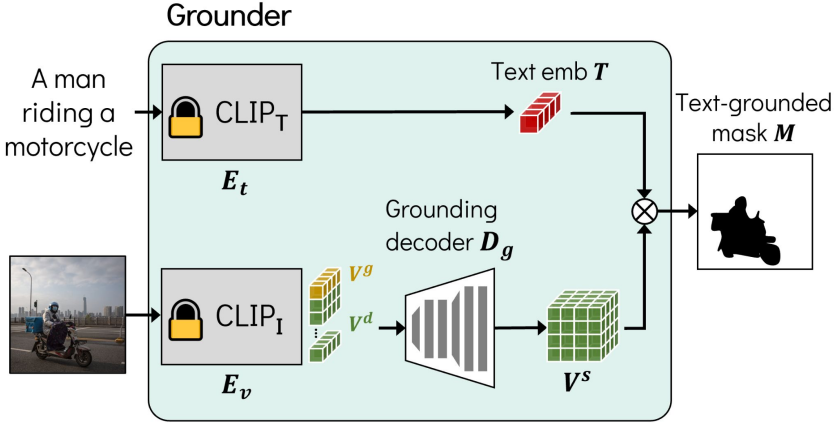
Pyramid
Sphinx

Frodo
Samwise
Gollum

Corgi
Shepherd

GROUNDING VISUALIZATION

Grounding decoder lets the model learn region-text alignment



ABLATION STUDIES

Table (a) indicates that TCL loss significantly improves segmentation performance by learning region-text alignment

[B] Training the grounding decoder without TCL loss does not improve performance

Method	VOC20	TCL _v	TCL _f	$\mathcal{L}_{\text{area}}$	CL	VOC20
A Baseline	53.2				✓	61.1
B + Decoder	52.3	✓		✓		74.6
C + TCL	77.4		✓	✓		76.0
		✓	✓	✓		77.4
		✓	✓	✓	✓	75.6
		✓	✓			67.1

(a) **Baseline to TCL.**

(b) **TCL losses.**

THANK YOU!

Contact: junbum.cha@kakaobrain.com

Code: <https://github.com/kakaobrain/tcl>

Demo: <https://huggingface.co/spaces/khanrc/tcl>