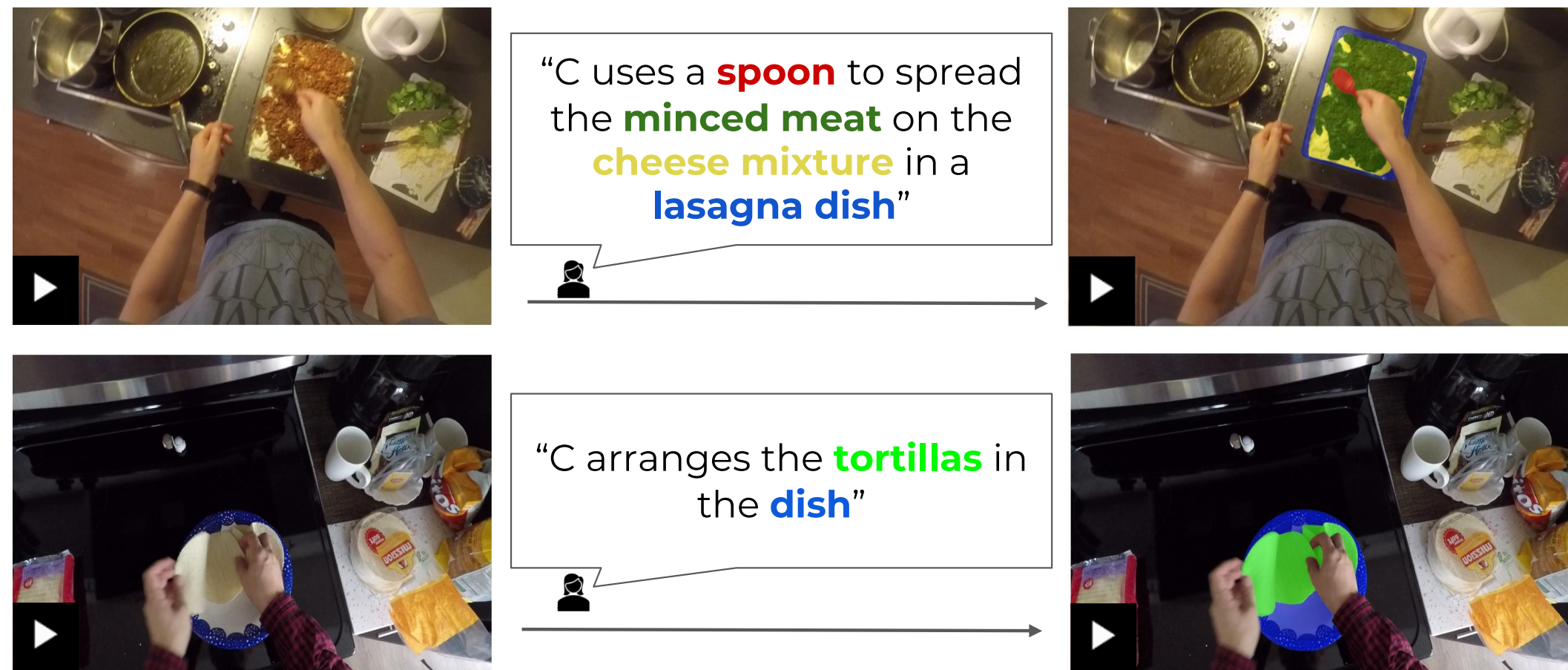


## Overview

- **Narration-based Video Object Segmentation (NVOS):** segment object instances mentioned in narrations for egocentric videos
- **Referred Object-Segment Aligner (ROSA):** a weakly-supervised framework for NVOS without spatial annotations
- **VISOR-NVOS:** an NVOS benchmark with newly-collected video clip narrations and associated segmentation masks



Narration-based Video Object Segmentation (Ours)

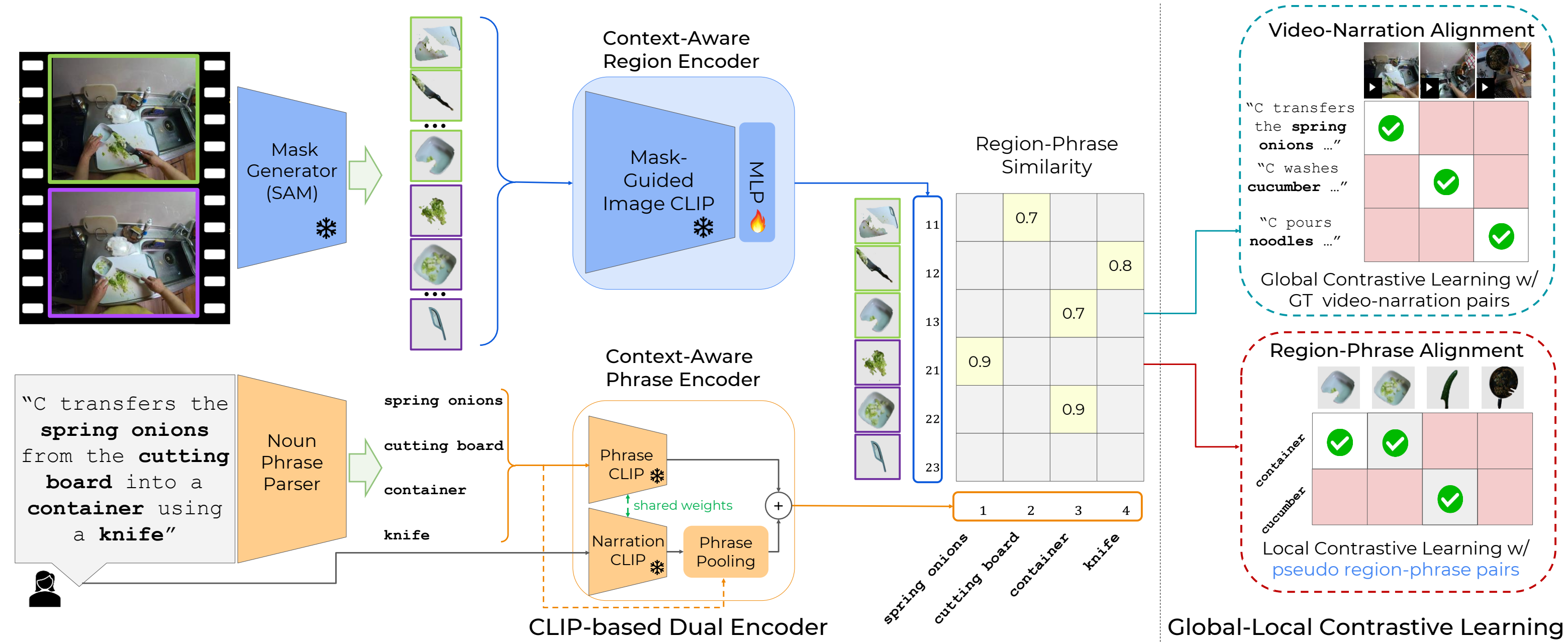
## Comparison with Related Tasks



## References

- [1] Kirillov et al. Segment anything. ICCV 2023.
- [2] Radford et al. Learning transferable visual models from natural language supervision. ICML 2021.
- [3] Darkhalil et al. Epic-kitchens visor benchmark: Video segmentations and object relations. NeurIPS 2022.
- [4] Tokmakov et al. Breaking the "Object" in Video Object Segmentation. CVPR 2023.

## Referred Object-Segment Aligner (ROSA)



- Generate *mask proposals* for video clips using SAM [1]; extract *object phrases* from narrations
- **CLIP-based Dual Encoder:** obtain *context-aware* representations for segmentation masks and object phrases via pretrained CLIP [2] models
- **Global-Local Contrastive Learning:** contrastive training via global *video-narration alignment* (VNA) and local *region-phrase alignment* (RPA)

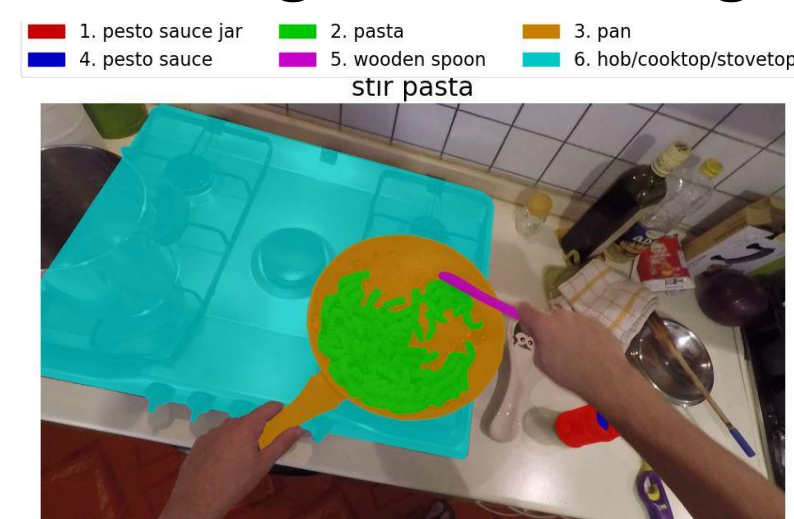
$$\mathcal{L}_{VNA} = \frac{1}{B} \sum_{i=1}^B -\log \frac{e^{\phi(v_i, s_i)/\tau}}{\sum_j e^{\phi(v_i, s_j)/\tau}} - \log \frac{e^{\phi(v_i, s_i)/\tau}}{\sum_j e^{\phi(v_j, s_i)/\tau}}$$

$$\mathcal{L}_{RPA} = -\sum_t \log \frac{e^{g(\tilde{r}_{tn}, w_n)/\tau}}{\sum_{n'} e^{g(\tilde{r}_{tn}, w_{n'})/\tau}} - \log \frac{\sum_t e^{g(\tilde{r}_{tn}, w_n)/\tau}}{\sum_{n'} \sum_t e^{g(\tilde{r}_{tn'}, w_n)/\tau}}$$

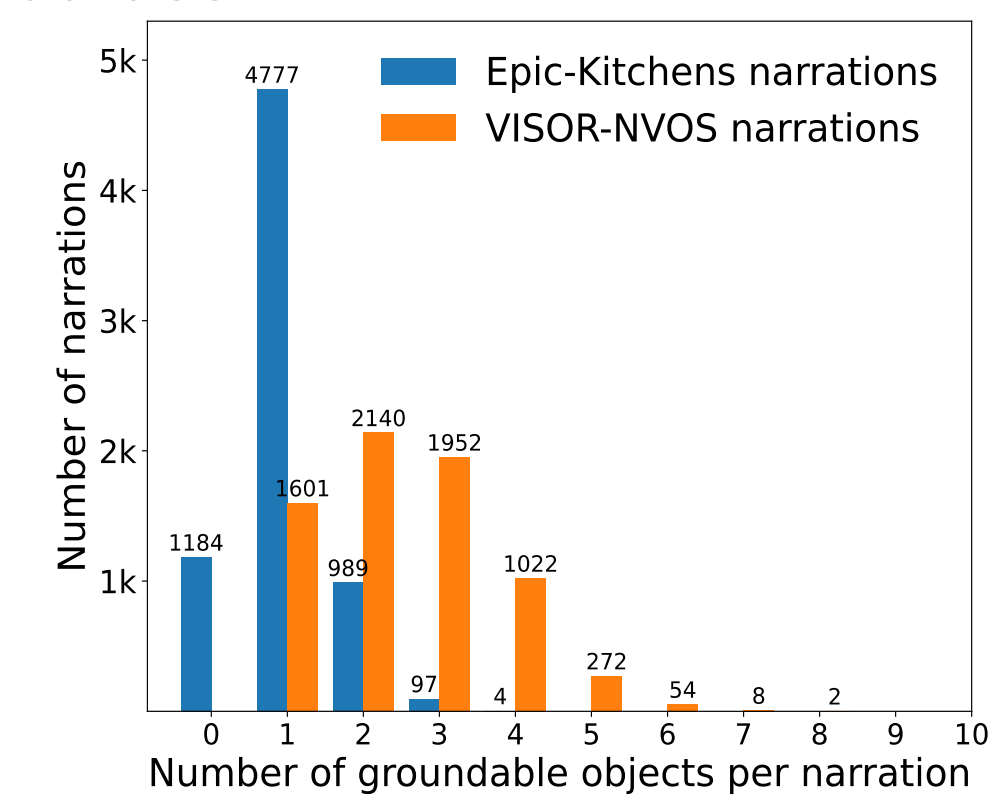
region to phrase                      phrase to region

## VISOR-NVOS Benchmark

- Annotate object-based narrations for video clips from VISOR [3] dataset
- 7,561 validation videos and 7,051 test videos
- 37,170 referred objects with associated segmentation masks
- Average number of groundable objects per narration: 2.54



**Narration:**  
The person uses a [wooden spoon]<5> to stir [pasta]<2> in a [pan]<3> on the [stovetop]<6>.



## Evaluation

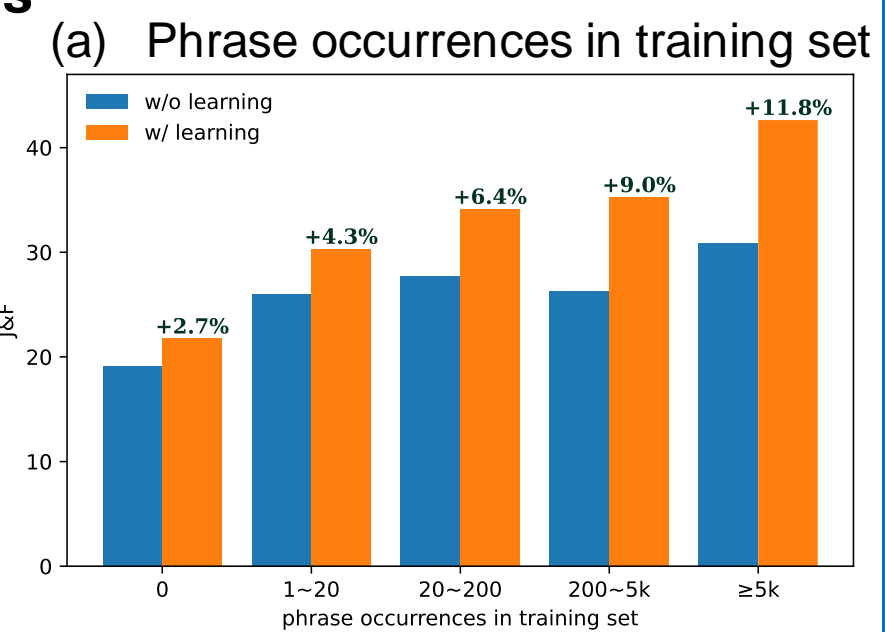
### Quantitative Results

Method	Supervision		VISOR-NVOS		VOST		
	Cross-Modal	Ego4D	$\mathcal{J}$	$\mathcal{F}$	$\mathcal{J}\&\mathcal{F}$	$\mathcal{J}_{union}$	
SAM upper bound			70.3	75.6	73.0	62.2	65.8
<i>Trained w/ labeled regions</i>							
ODISE	mask-text		29.0	32.8	30.9	17.6	21.1
GroundedSAM	bbox-text		37.3	41.8	39.5	21.8	25.3
<i>Trained w/o labeled regions</i>							
SAM + CLIP (ViT-B/16)	image-text		22.2	25.8	24.0	16.5	19.2
CoMMa + SAM	video-narration	✓	15.3	25.3	20.3	9.4	11.5
<b>ROSA (ViT-B/16)</b>	video-narration	✓	<b>34.9</b>	<b>41.2</b>	<b>38.1</b>	<b>22.2</b>	<b>25.4</b>
<b>ROSA (ViT-L/14)</b>	video-narration	✓	<b>38.7</b>	<b>46.0</b>	<b>42.4</b>	<b>23.2</b>	<b>26.7</b>

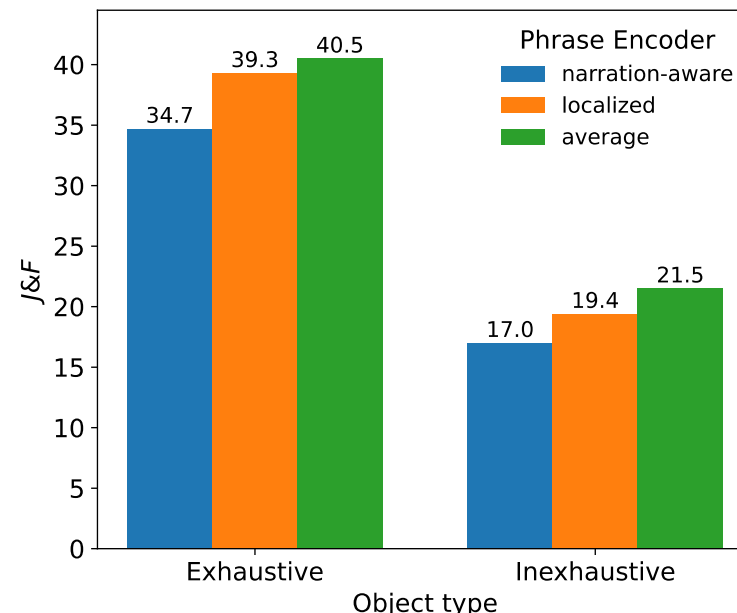
$\mathcal{J}$ : pixel-wise IoU;  $\mathcal{F}$ : F measure of mask boundaries;  $\mathcal{J}\&\mathcal{F}$ : average of  $\mathcal{J}$  and  $\mathcal{F}$   
VOST [4]: an auxiliary benchmark of grounding objects under complex transforms  
 $\mathcal{J}_{union}$ :  $\mathcal{J}$  with the union of all instances;  $\mathcal{J}_{ins}$ :  $\mathcal{J}$  with the best matched instance

### Ablation Studies

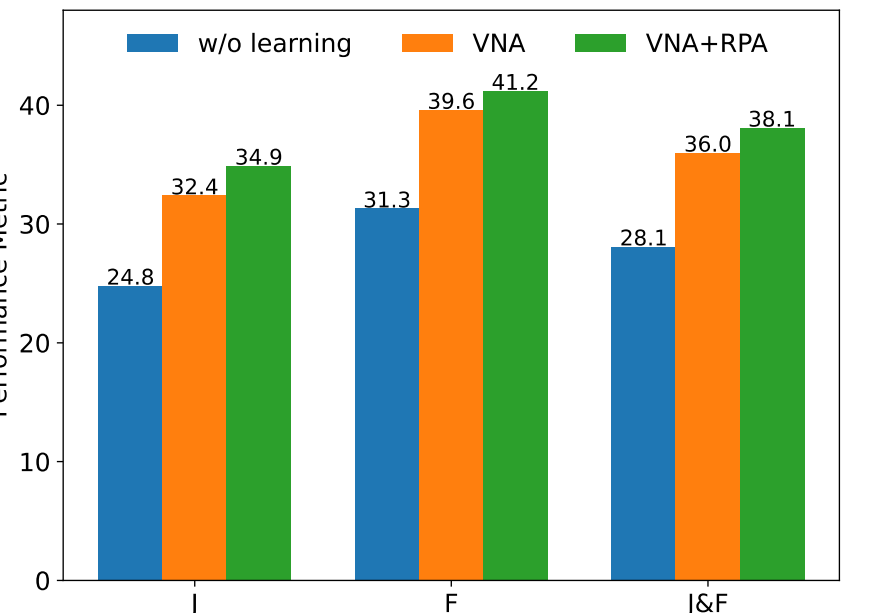
- More phrase occurrences in training set lead to more gains; ROSA generalizes well to unseen objects (+2.7%)
- Adding contexts improves for both exhaustive (+1.2%) and inexhaustive (+2.1%) objects
- Both VNA and RPA are effective at learning better region-phrase similarities



### (b) Context in Phrase Encoder



### (c) Global-Local Contrastive Learning



### Qualitative Results

