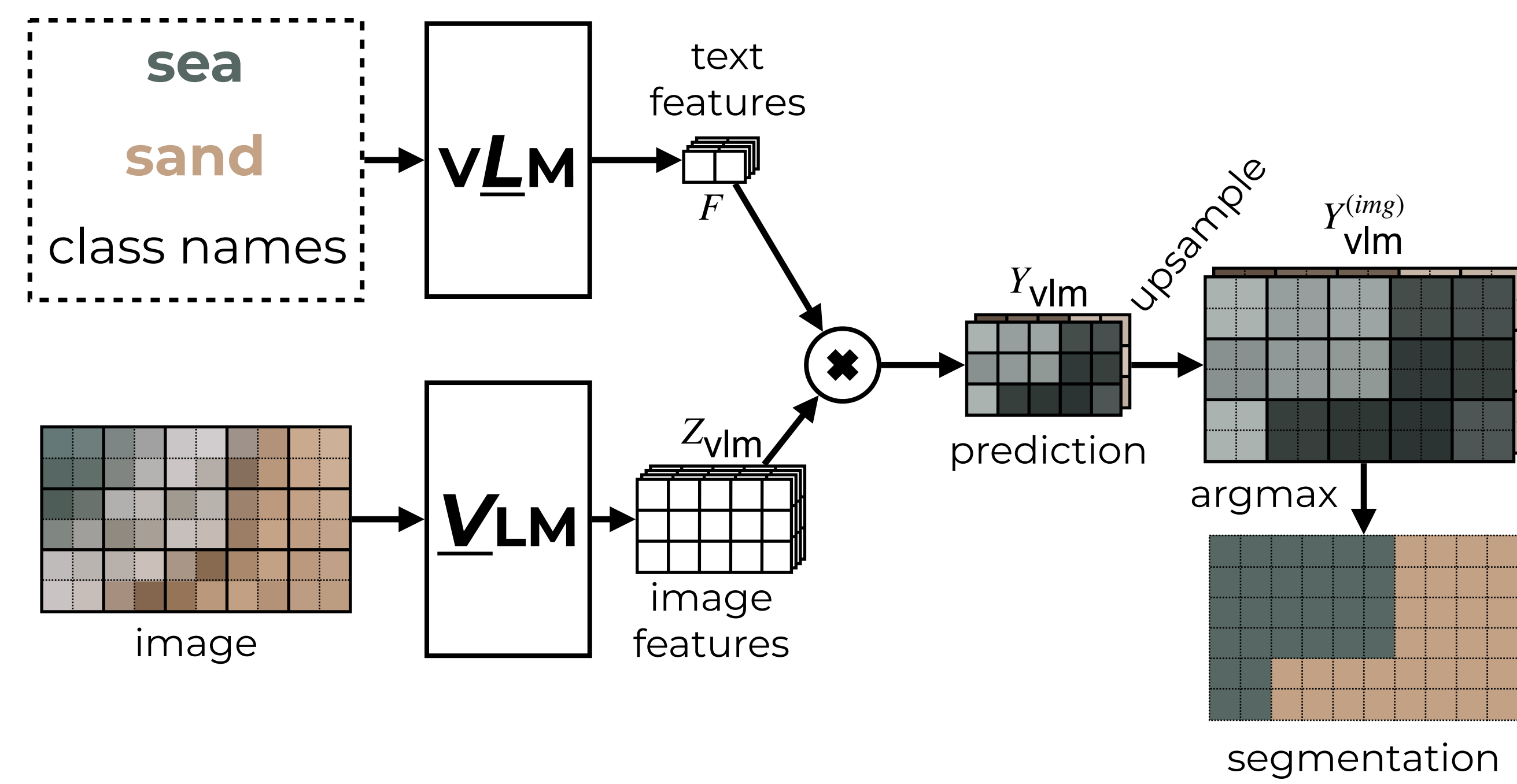


OVSS: Open-Vocabulary Semantic Segmentation

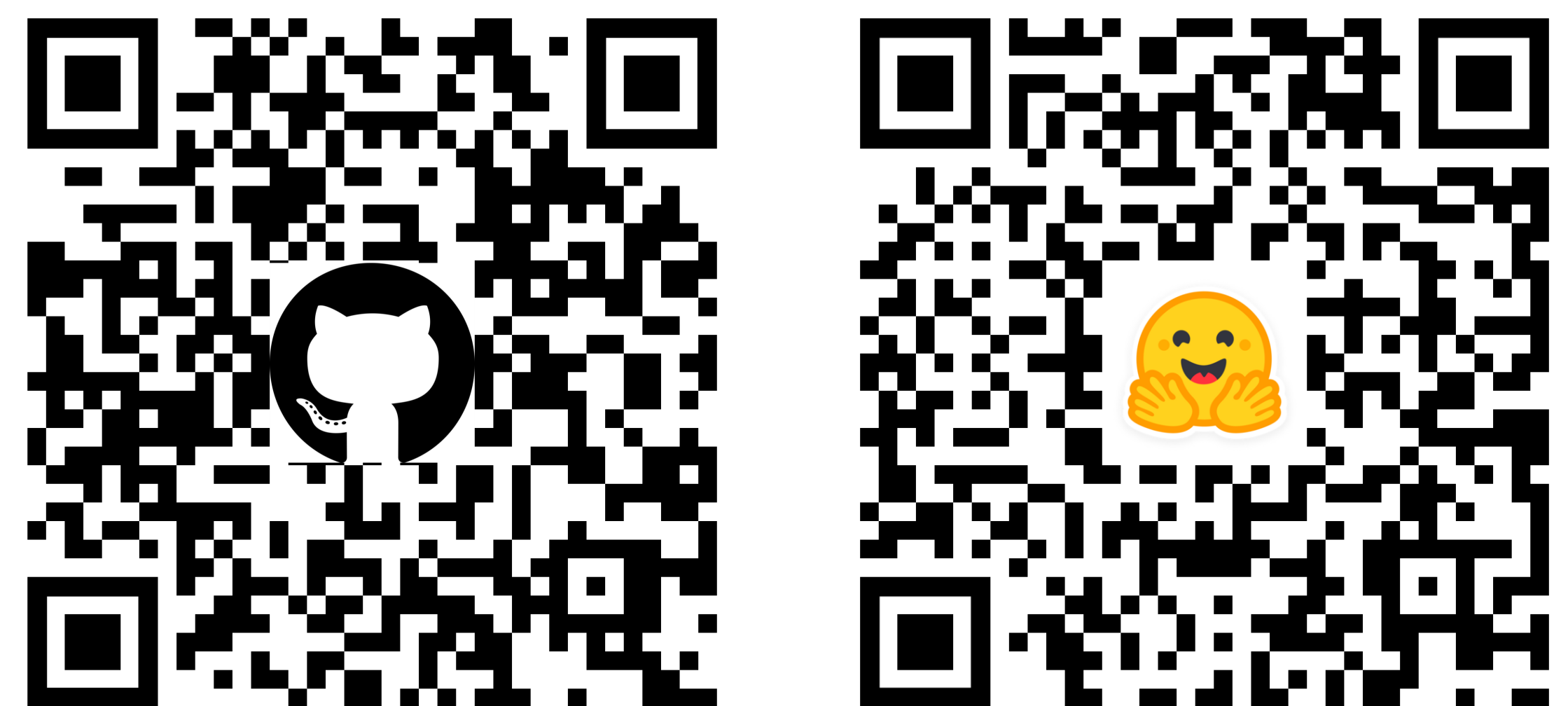
Task: Classify each pixel into one of the classes given only by their names



VLMs offer a training-free approach, but:

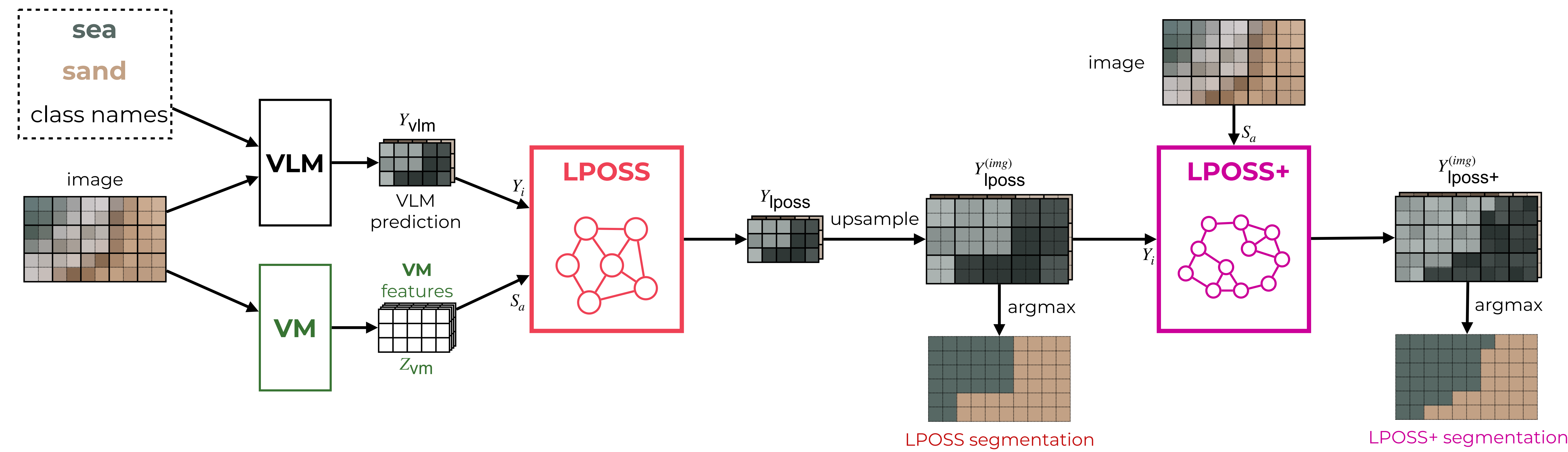
- predictions are performed independently at each patch (noisy)
- predictions are at patch level (low res.)

Code Demo



Code: <https://github.com/vladan-stojnic/LPOSS>
Demo: <https://hf.co/spaces/stojnvla/LPOSS>

Can **graph-based segmentation approaches** improve open-vocabulary semantic segmentation in a training-free way?



LPOSS: Label propagation on the patch level

- Refine VLM predictions Y_i
- Similar predictions for similar & nearby patches

$$Q(\hat{Y}) = (1 - \alpha) \sum_{i=1}^N \|\hat{Y}_i - Y_i\|^2 + \alpha \sum_{i,j=1}^N S_{ij} \left\| \frac{\hat{Y}_i}{\sqrt{d_i}} - \frac{\hat{Y}_j}{\sqrt{d_j}} \right\|^2 \quad (1)$$

Affinity $S = S_a \odot S_p$

S_a - appearance affinity: patch feature similarity

S_p - spatial affinity: from 2D distance of patches

Normalization factor $d_i = \sum_j S_{ij}$

- Use of VM (DINO) features to construct S_a

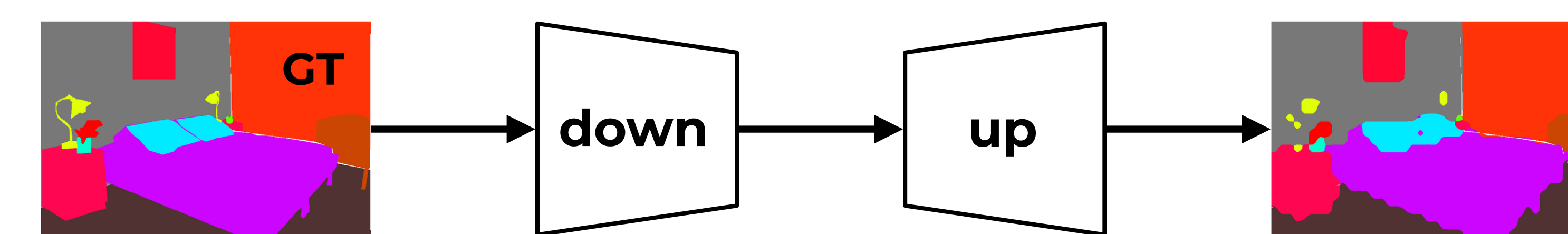
LPOSS+: Label propagation on the pixel level

- Refine LPOSS predictions **at the pixel level**
- Repeat optimization of Eq. (1), but:
 - Use upsampled LPOSS prediction as Y_i
 - Affinity S_{ij} over *pixels*

Motivation for LPOSS+

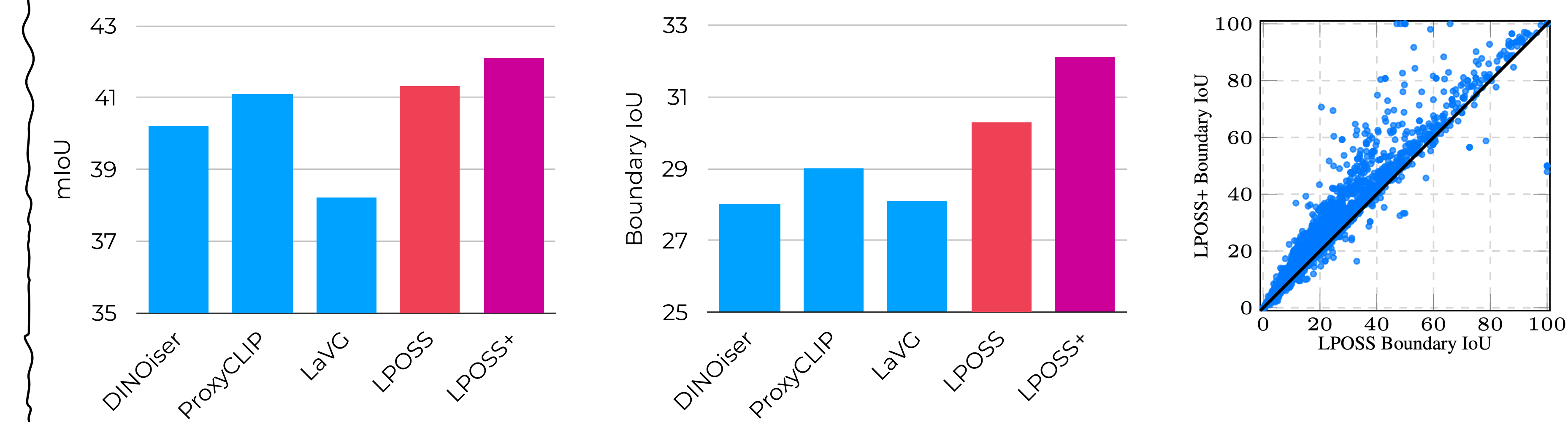
Oracle experiment: Downscale and upscale GT

- mIoU drops to 85.2%
- Boundary IoU drops to 69.5%

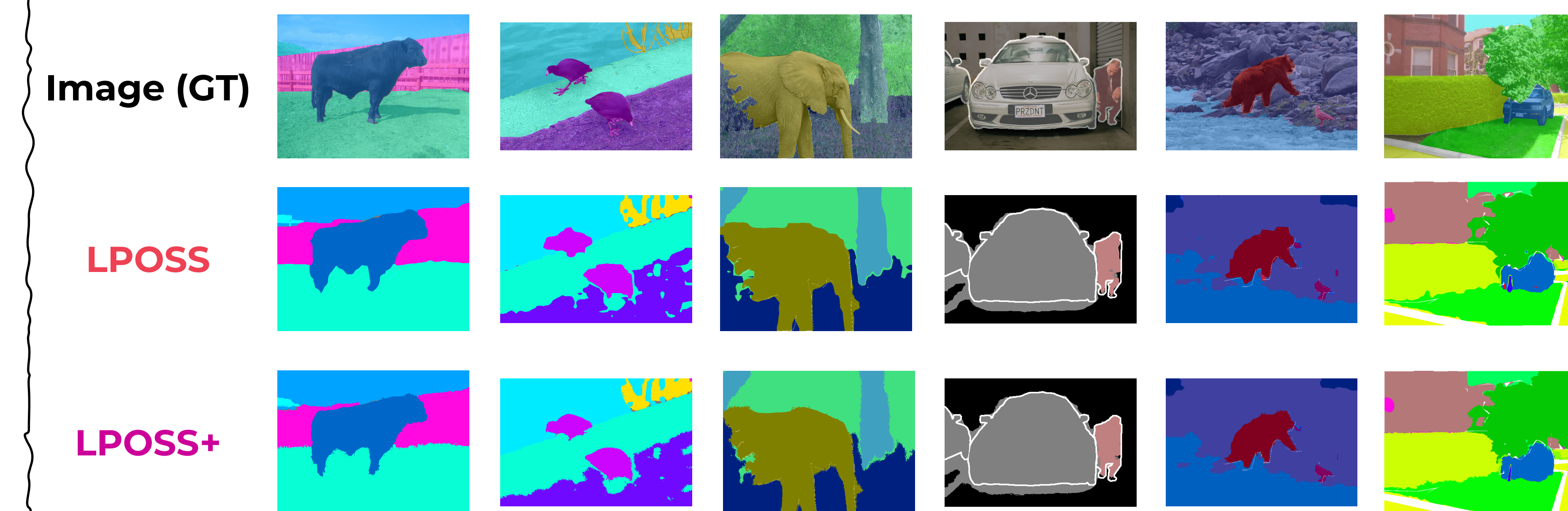


Quantitative results

Performance reported across 8 datasets

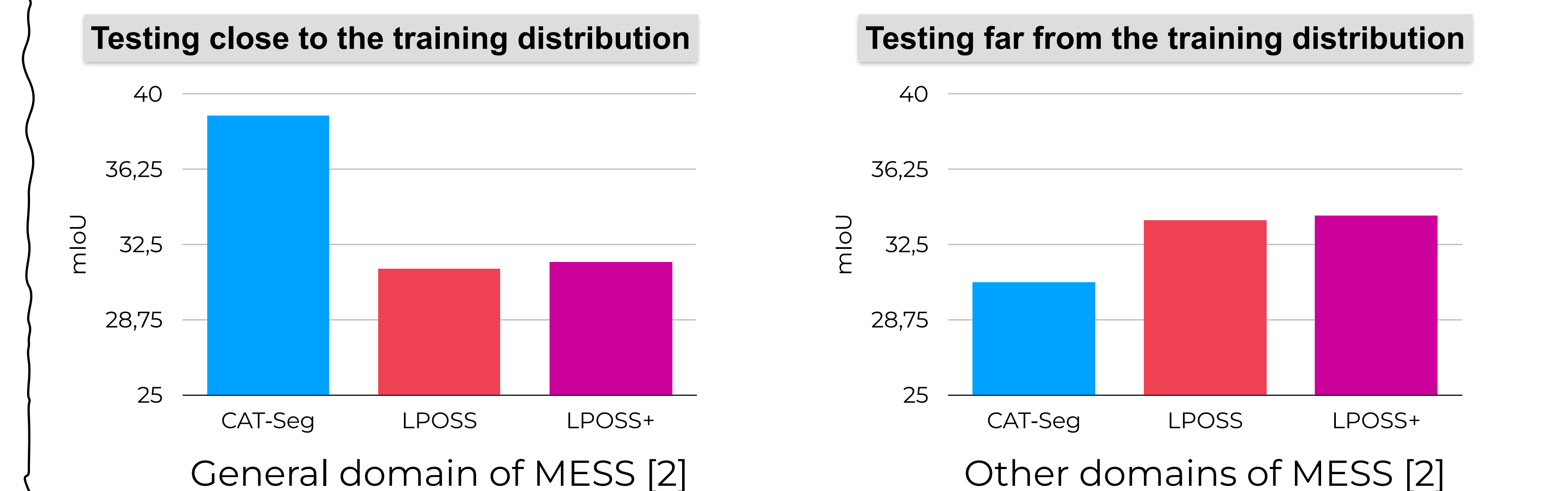


Qualitative results



Training-free vs finetuning on COCO-Stuff

- CAT-Seg [1] finetunes the VLM for OVSS → poor generalization



[1] Cho et al. CVPR'24, CAT-Seg: Cost Aggregation for Open-Vocabulary Semantic Segmentation

[2] Blumenstiel et al. NeurIPS'24, What a MESS: Multi-Domain Evaluation of Zero-Shot Semantic Segmentation