# What could go wrong?
## Discovering and describing failure modes in computer vision

Gabriela Csurka, Tyler L. Hayes, Diane Larlus, Riccardo Volpi

ECCV 2024

NAVER LABS Europe

## Summary

❖ We formalize **Language-Based Error Explainability** (**LBEE**)

❖ We propose a family of **task agnostic** methods **to tackle LBEE**

❖ We introduce a **set of metrics** to evaluate LBEE performance

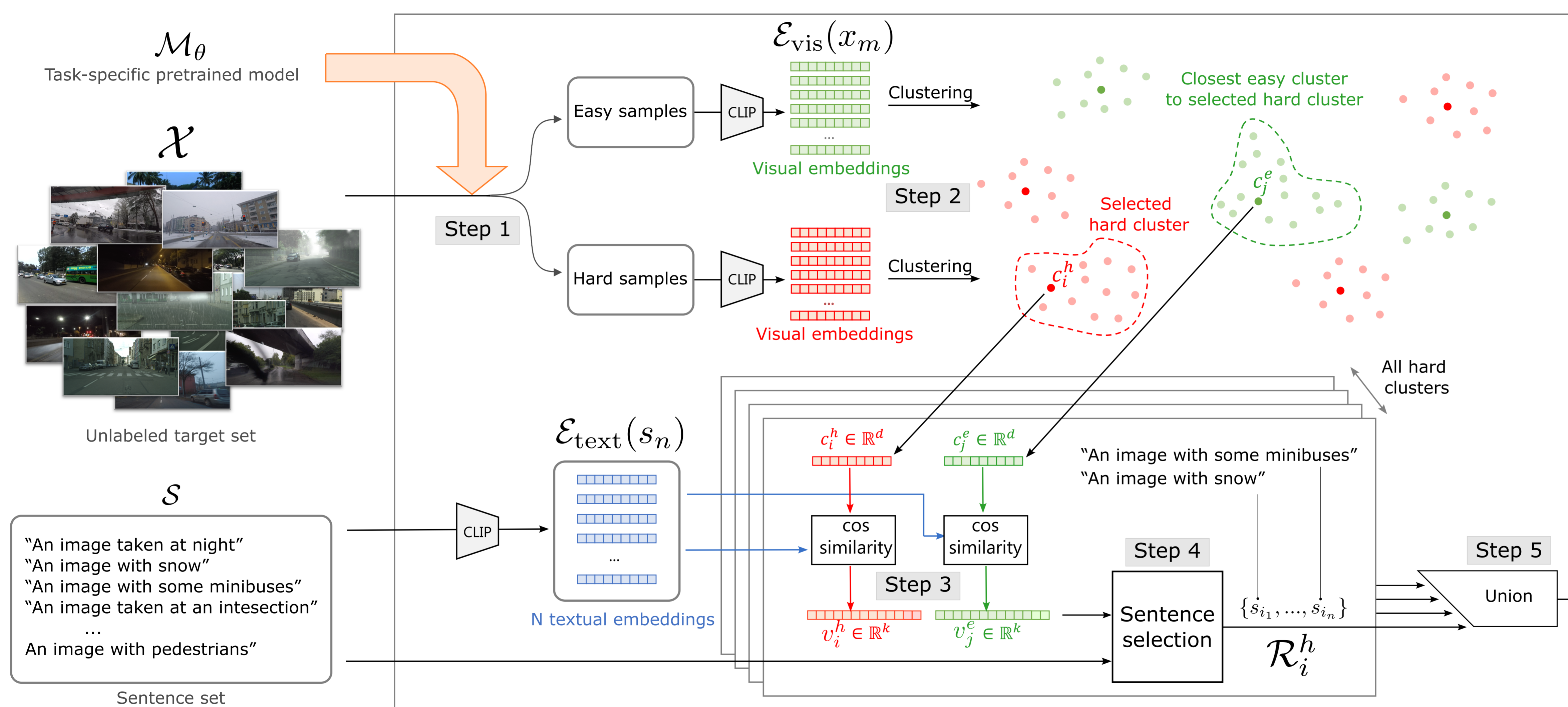❖ We show the effectiveness of the proposed methods on various tasks

## Contribution #1: Problem Formulation

Given a **target set** $X$ and a **model** $M_\theta$, our goal is to find sentences describing likely failure causes for the model

$$S_\beta^* = \left\{ s_n \in S \;\middle|\; \omega_\theta^{s_n} < \omega_\theta^{\mathrm{avg}} - \beta \right\}$$

Predefined sentence set

Model average performance on images relevant to $s_n$

Model average performance on $X$

Predefined margin

## Contribution #2: A Family of Methods



$\mathcal{M}_\theta$ Task-specific pretrained model

$\mathcal{X}$ Unlabeled target set

$\mathcal{S}$ Sentence set
"An image taken at night"
"An image with snow"
"An image with some minibuses"
"An image taken at an intesection"
...
An image with pedestrians"

Easy samples → CLIP → $\mathcal{E}_{\mathrm{vis}}(x_m)$ Visual embeddings → Clustering

Hard samples → CLIP → Visual embeddings → Clustering

Step 1

Step 2

Closest easy cluster to selected hard cluster — $c_j^e$

Selected hard cluster — $c_i^h$

All hard clusters

$\mathcal{E}_{\mathrm{text}}(s_n)$ CLIP → N textual embeddings

$c_i^h \in \mathbb{R}^d$ → cos similarity → $v_i^h \in \mathbb{R}^k$

$c_j^e \in \mathbb{R}^d$ → cos similarity → $v_j^e \in \mathbb{R}^k$

Step 3

"An image with some minibuses"
"An image with snow"

Step 4 — Sentence selection → $\{s_{i_1}, ..., s_{i_n}\}$ → $\mathcal{R}_i^h$

Step 5 — Union → $\mathcal{R}_S$

**Step 1**: split the images into easy and hard sets based on the model's confidence

**Step 2**: embed images in the CLIP space and cluster the hard and easy sets independently

**Step 3**: assign to each hard prototype the closest easy prototype in this space

**Step 4**: select sentences for hard clusters in the CLIP space based on cosine similarities with the cluster prototypes that are not relevant for the closest easy clusters

**Step 5**: aggregate cluster-specific sentence sets to produce the global output ($\mathcal{R}_S$)

## Contribution #3: Evaluation Metrics

Given a hard cluster $c_i^h$ and set of selected sentences $\mathcal{R}_i^h$
❖ **Hardness ratio (HR):** ratio of sentences pointing to reasons for model failure
❖ **Correctness Ratio (CR):** average ratio of images that are relevant to individual sentences

$$HR_i = \frac{\left|\{\forall s_k \in \mathcal{R}_i^h \mid \omega_\theta^{avg} - \omega_\theta^k > \beta\}\right|}{|\mathcal{R}_i^h|} \qquad CR_i = \frac{1}{|\mathcal{R}_i^h|} \sum_{s_k \in \mathcal{R}_i^h} \frac{1}{|c_i^h|} \sum_{x \in c_i^h} \Gamma(x, s_k)$$

Given $S_\beta^*$ and the overall output ($\mathcal{R}_S = \cup \mathcal{R}_i^h$)
❖ **True positive rate (TPR):** evaluates how well $S_\beta^*$ is covered
❖ **Jaccard Index (JI):** measures coverage while penalizing false positives

$$TPR = \frac{|S_\beta^* \cap \mathcal{R}_S|}{|S_\beta^*|} \qquad JI = \frac{|S_\beta^* \cap \mathcal{R}_S|}{|S_\beta^* \cup \mathcal{R}_S|}$$

## Experimental Setup

Tasks and datasets:
❖ **Urban scene segmentation:** ACDC, IDD, WD2
❖ **Classification with spurious correlations**: $\mathrm{NICO}_{++}^{75/85/95}$
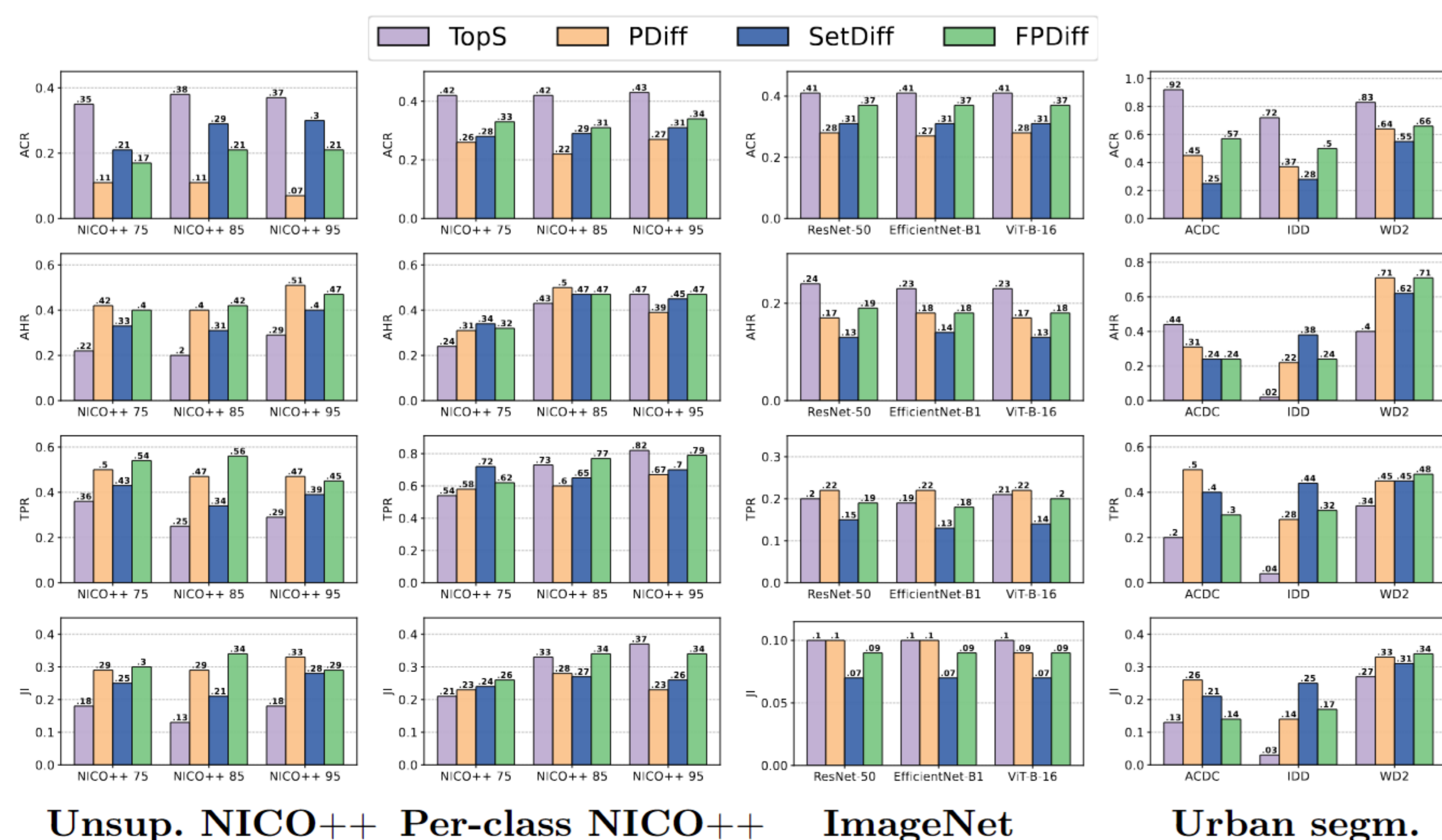❖ **ImageNet-1K classification**

Methods:
❖ **TopS**: top ranked sentences based on cosine similarity
❖ **PDiff:** rank based on prototype difference
❖ **FPDiff:** Pdiff filtered with TopS
❖ **SetDiff:** Sentence set differences.

Default design choices:
❖ Open-CLIP, 15 clusters, 3 sentences, $\beta = .2 * \omega_\theta^{\mathrm{std}}$

## Quantitative Results



TopS  PDiff  SetDiff  FPDiff

Unsup. NICO++     Per-class NICO++     ImageNet     Urban segm.

## Qualitative Results ($\mathcal{R}_i^h$)



**ACDC**
TopS: "taken in a rainy weather" "taken in a dull weather" "taken in a cloudy weather"
SetDiff: "shadows on the road" "people on the road" "of a sidewalk"
PDiff: "water on the road" "taken in a rainy weather" "rickshaw on the road"
FPDiff: "water on the road" "taken in a rainy weather" "taken in a stormy weather"

**WD2**
TopS: "vehicle on the road" "obstacle on the road" "jeep on the road"
SetDiff: "mud on the road" "taken in a windy weather" "round-about scene"
PDiff: "mountain" "terrain" "forest"
FPDiff: "mud on the road" "rocks on the road" "jeep on the road"

**NICO$_{++}^{85}$**
TopS: "train" "train in the rocks" "bus in the rocks"
SetDiff: "train" "train in the rocks" "bus in the rocks"
PDiff: "train in the rocks" "bus in the rocks" "train in the grass"
FPDiff: "train in the rocks" "bus in the rocks" "train in the grass"

**ImageNet**
TopS: "stage indoor" "arena performance" "taken in a basement"
SetDiff: "basement" "stage indoor" "thumbnail image"
PDiff: "with motion blur" "blurry image" "stage outdoor"
FPDiff: "blurry image" "stage outdoor" "stage indoor"