



# Uncertainty Quantification & Adaptation of Vision-Language Models (VLMs)

FRQS Workshop in Digital Pathology and Vision-Language models 27/10/25



Nicolas Thome
CNRS, ILLS, Sorbonne Université, ISIR

#### Success of Al

**ChatGPT** 



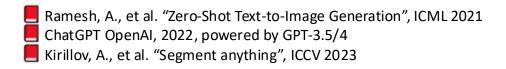
DALL-E



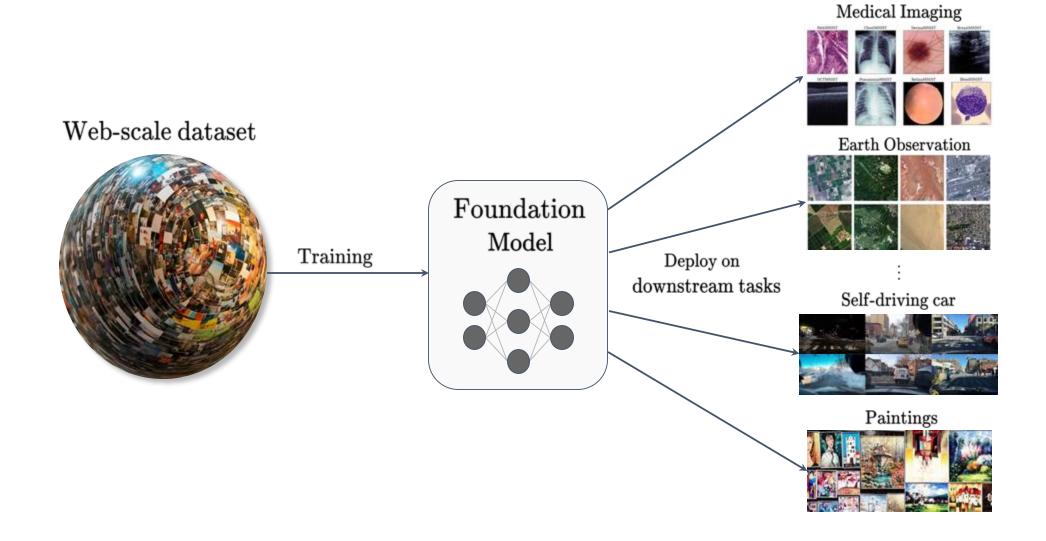
Segment Anything Model (SAM)



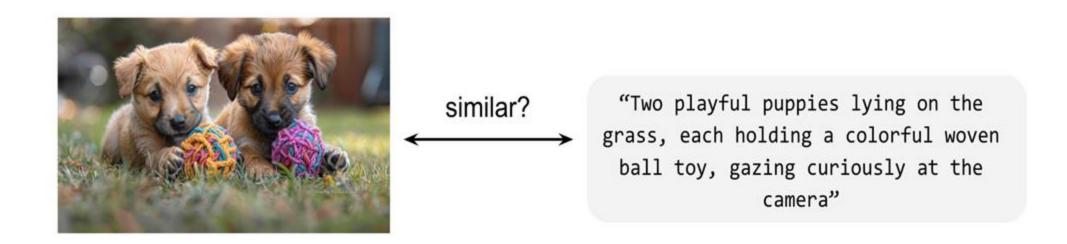
Deep Learning: underlying principle powering these breakthroughs



## Foundation models



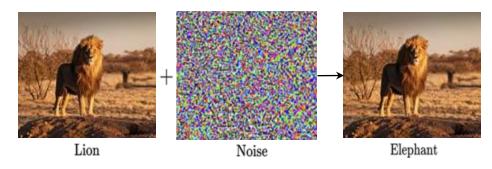
## Vision perception & Vision Language Models (VLMs)



- Shared representation space for image and text
- Contrastive loss, e.g., CLIP, SigLIP
- Zero-shot classification

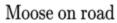
## Robustness of VLMs

#### Adversarial attacks



#### Uncertain inputs







Blurry image



Snowy road

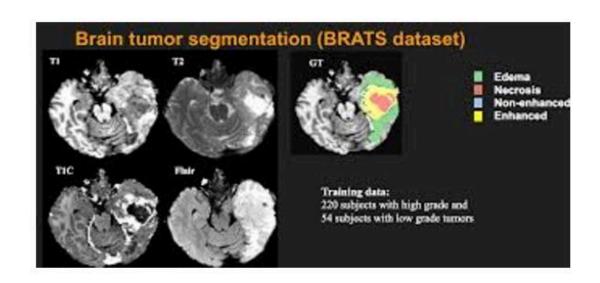
#### Spurrious correlations

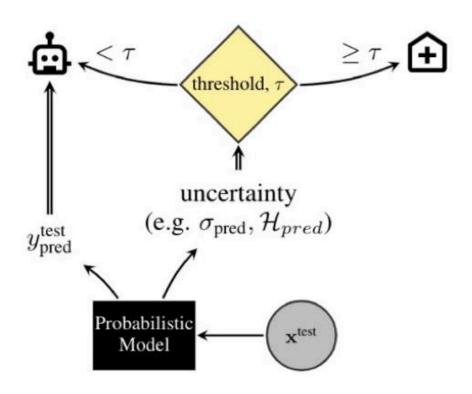




A nurse A firefighter

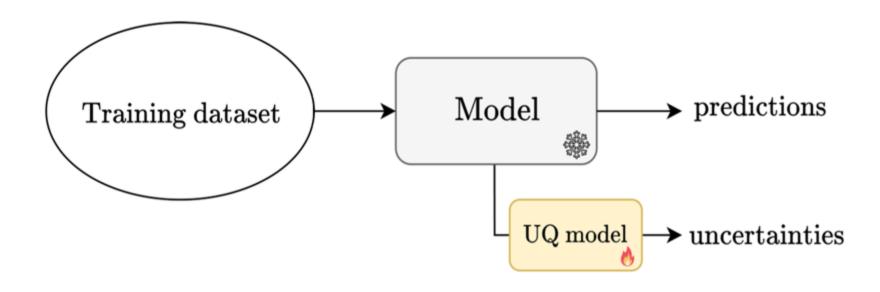
# Post-hoc uncertainty quantification





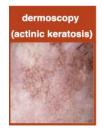
- Safety-critical applications, e.g., medical domain
- Selective classification: reject uncertain inputs

# Post-hoc uncertainty quantification



- Foundation model re-training: computationally expensive
- Post-hoc: predictive model frozen
- Auxiliary uncertainty module

# Post-hoc uncertainty quantification





novel class of target





different imaging view

#### Epistemic uncertainty



Moose on the road

- Out-of-distribution (OOD) (unknown)
- Reducible

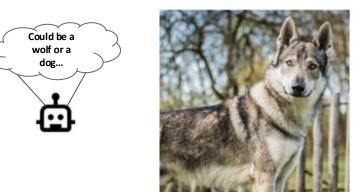
I have never

seen such a

thing...

→ Out-of-distribution detection

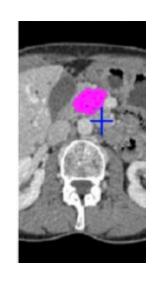
#### Aleatoric uncertainty



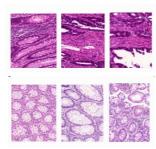
Czechoslovakian wolfdog

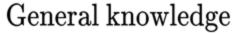
- Ambiguous example
- Irreducible
- → Failure prediction





## Adaptation of VLMs: domain shifts





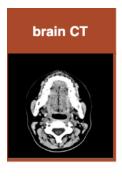








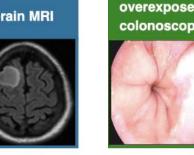


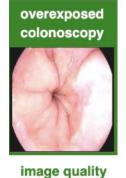




different

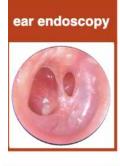
modality





issue

colonoscopy











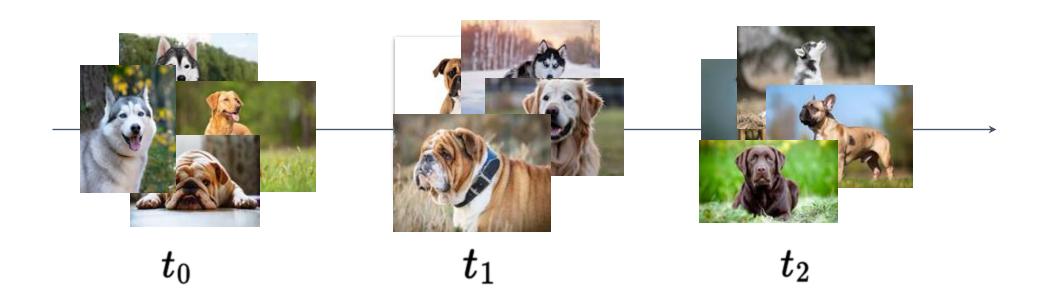
another cohort

**Domain shifts:** common in the medical domain: acquisition devices, centers, populations, etc

→ Adaptation needed for highly specialized applications

# Test-Time-Adaptation (TTA)

- No annotated samples
- Exploit incoming test-samples (online)



Additional requirement: Maintain or improve robustness after adaptation

# Today's talk

- ViLU [A]: UQ for VLMs, failure prediction
  - Fine text-image interactions for UQ
  - Large-scale training & generalization
- CLIP-TTA [B]: test-time adaptation for CLIP
  - Training loss adapted to CLIP
  - Robustness to pseudo-label errors & class collapse

#### [A] ViLU: Learning Vision-Language Uncertainties for Failure Prediction.

M. Lafon, Y. Karmim, J. Silva-Rodriguez, P. Couairon, C. Rambour, R. Fournier, I. Ben Ayed, J. Dolz, N. Thome. ICCV 2025.

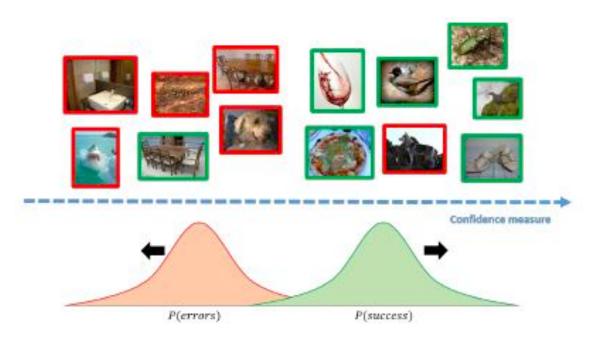
#### [B] CLIPTTA: Robust Contrastive Vision-Language Test-Time Adaptation.

M. Lafon, G. Vargas Hakim, C. Rambour, C. Desrosiers, N. Thome. NeurIPS 2025.

# Outline

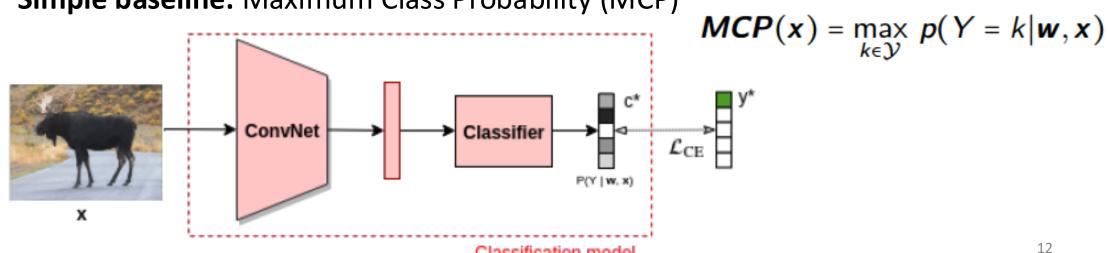
- 1. ViLU
- 2. CLIPTTA

# Failure Prediction with deep neural nets



 Confidence criterion C(x): separate correct from incorrect prediction

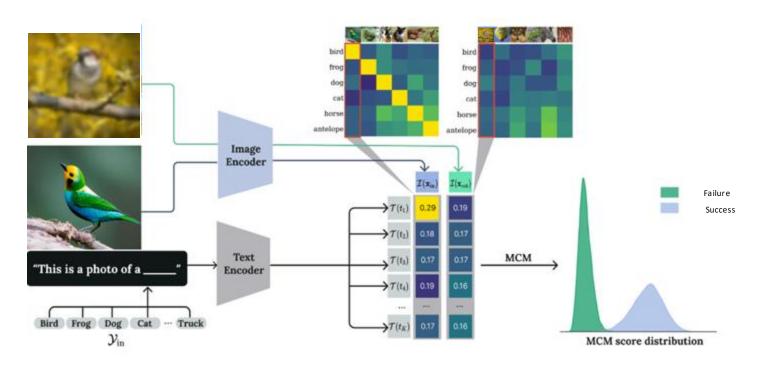
Simple baseline: Maximum Class Probability (MCP)



#### Failure Prediction with VLMs

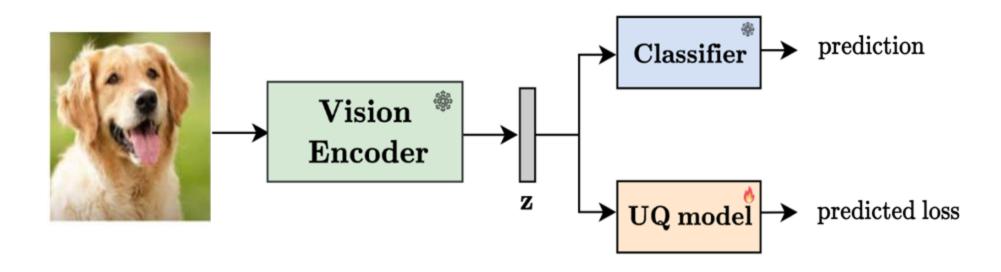
**Maximum Concept Matching (MCM) = Probability of predicted class with VLMs** 

MCP Extension to VLMs



- Pros:
  - Strong baseline
  - No training required
- <u>Cons:</u>
  - Overconfident by design for errors
  - Limited adaptability

## Failure Prediction for vision models: loss prediction



- Learning to predict the training loss
- High predicted loss = likely incorrect prediction
- Learning Visual Uncertainties (LVU): ConfidNet, PVU
  - → How to adapt loss prediction for VLMs?
- ConfidNet, Corbière, C., et al. "Addressing failure prediction by learning model confidence". NeurIPS 2019
- 📕 D. Yoo and In So Kweon. Learning loss for active learning. CVPR 2019
- Kirchhof, Michael, et al. "Pretrained visual uncertainties." arXiv preprint 2024

#### Taking task complexity into account



What is the uncertainty associated with this image?

→ It depends on the task

Task 1: cat vs dog classification





- Low class confusion
- Low uncertainty

#### Taking task complexity into account



What is the uncertainty associated with this image?

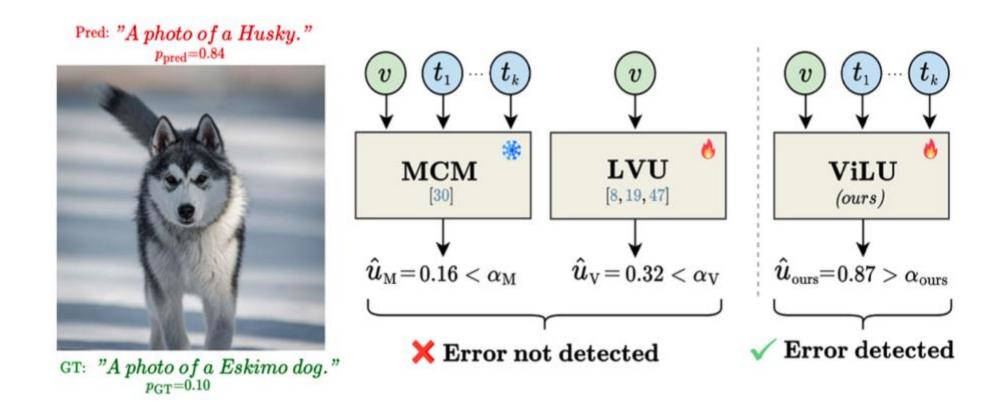
→ It depends on the task

Task 2: Dog breed classification



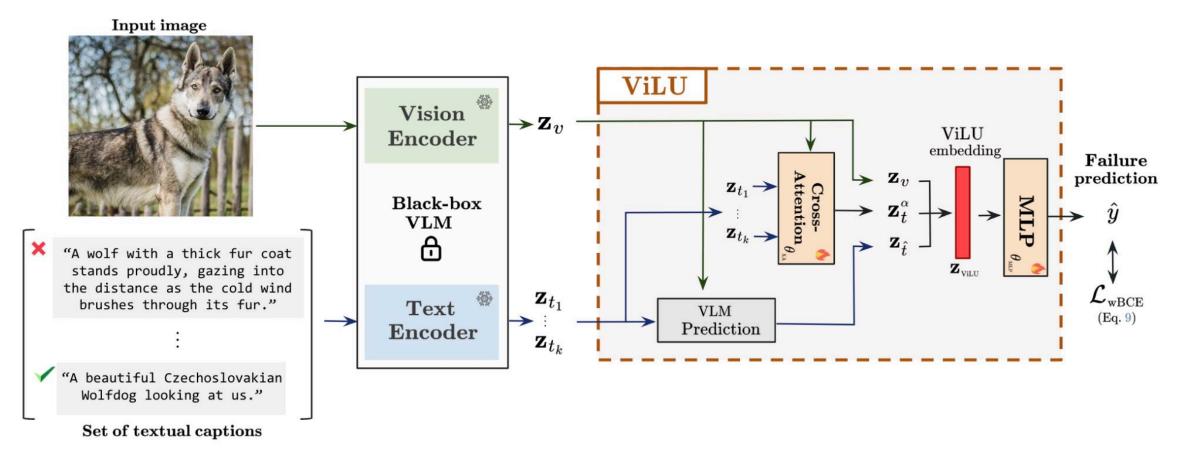
- Higher class confusion
- High uncertainty

## ViLU: Learning Vision-Language Uncertainty



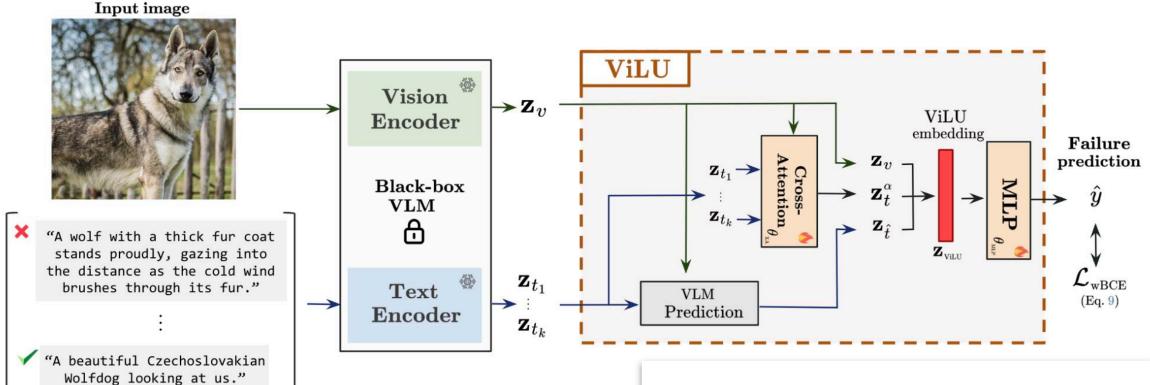
- LVU methods use visual input only
- ViLU → Learn uncertainty contextualized with the task information

## ViLU: Learning Vision-Language Uncertainty



- ullet Inputs: visual representation  $oldsymbol{z}_v$  + K concept embeddings  $Z_t = \left\{oldsymbol{z}_{t_j}
  ight\}_{1 \leq j \leq K}$
- ViLU embedding:  $m{z}_{ ext{ViLU}}=(m{z}_v,m{z}_{\hat{t}},m{z}_{\hat{t}}^lpha)$ ,  $m{z}_{\hat{t}}$  pred. Class embedding
  - Image-text cross-attention module => $z_t^{\alpha}$
  - Query  $z_v$  , keys/values  $Z_t$

## ViLU: Learning Vision-Language Uncertainty



- Failure prediction from ViLU embedding: binary classification
- Consistent generalization of MCM

Set of textual captions

$$egin{aligned} oldsymbol{z}_{ ext{ViLU}} &= (oldsymbol{z}_v, oldsymbol{z}_{\hat{t}}, oldsymbol{z}_t^lpha) \ g_{ heta_{ ext{MLP}}}(oldsymbol{z}_{ ext{ViLU}}) &pprox rac{1}{2} oldsymbol{z}_{ ext{ViLU}}^T A oldsymbol{z}_{ ext{ViLU}} &= oldsymbol{z}_v^T oldsymbol{z}_{\hat{t}} \ \end{aligned}$$
 with  $A = \left(egin{array}{ccc} oldsymbol{0} & I_d & oldsymbol{0} & oldsymbol{$ 

## Experiments

- ViLU trained and evaluated on each downstream dataset
- 13 classification datasets (e.g. ImageNet)
- 3 Image-caption datasets (e.g. CC12M)
- Metrics: FPR95 and AUC
- Baselines:
  - MCM, entropy, DOCTOR
  - Data-driven predictors: LVU, Rel-U, BayesVLM

## Results

	CIFAR-10		Image	ImageNet-1k		C12M
	AUC↑	FPR95↓	AUC↑	FPR95↓	AUC↑	FPR95↓
MCM	89.9	52.1	80.8	71.3	88.8	58.8
Entropy	88.7	59.9	78.3	76.8	80.2	74.0
DOCTOR	89.5	56.5	80.3	72.9	88.6	59.9
BayesVLM	92.6	44.9	81.5	70.3	90.9	53.3
LVU - ConfidNet	96.4	21.8	78.7	77.0	74.0	76.5
LVU + XA	97.9	10.8	88.8	50.1	88.9	48.9
LVU + Pred.	97.7	11.4	86.1	63.5	93.6	37.0
ViLU (ours)	98.3	8.2	89.5	47.4	$\boldsymbol{95.2}$	25.2

- LVU competitive on small datasets (e.g. CIFAR-10)
- ViLU achieves best performance

## Results

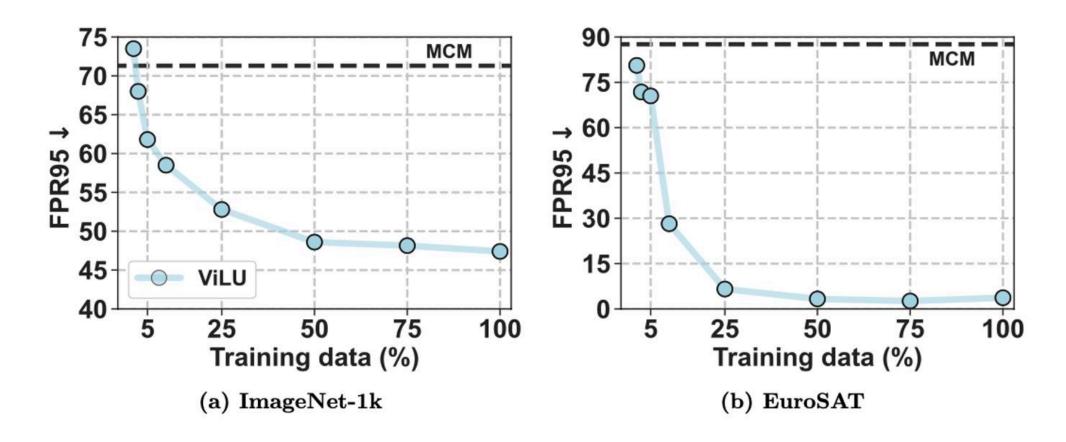
	CIFAR-10		Image	ImageNet-1k		C12M
	AUC↑	FPR95↓	AUC↑	FPR95↓	AUC↑	FPR95↓
MCM	89.9	52.1	80.8	71.3	88.8	58.8
Entropy	88.7	59.9	78.3	76.8	80.2	74.0
DOCTOR	89.5	56.5	80.3	72.9	88.6	59.9
BayesVLM	92.6	44.9	81.5	70.3	90.9	53.3
LVU - ConfidNet	96.4	21.8	78.7	77.0	74.0	76.5
LVU + XA	97.9	10.8	88.8	50.1	88.9	48.9
LVU + Pred.	97.7	11.4	86.1	63.5	93.6	37.0
ViLU (ours)	98.3	8.2	89.5	47.4	$\boldsymbol{95.2}$	25.2

• LVU struggles on challenging datasets (e.g. ImageNet)

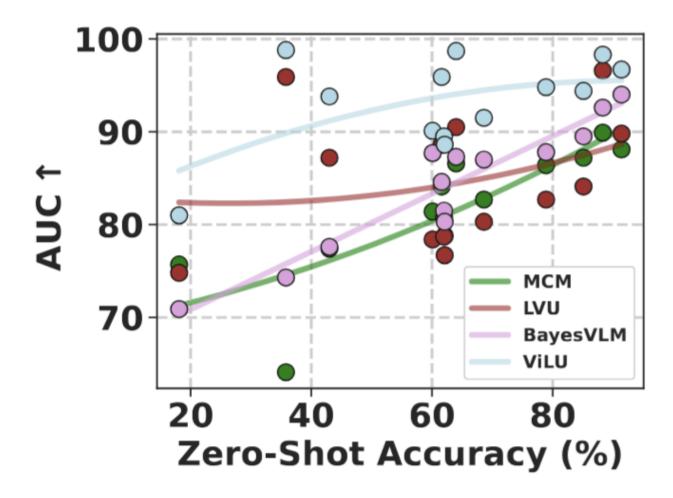
## Results

	CIFAR-10		Image	ImageNet-1k		C12M
	AUC↑	FPR95↓	AUC↑	FPR95↓	AUC↑	FPR95↓
MCM	89.9	52.1	80.8	71.3	88.8	58.8
Entropy	88.7	59.9	78.3	76.8	80.2	74.0
DOCTOR	89.5	56.5	80.3	72.9	88.6	59.9
BayesVLM	92.6	44.9	81.5	70.3	90.9	53.3
LVU - ConfidNet	96.4	21.8	78.7	77.0	74.0	76.5
LVU + XA	97.9	10.8	88.8	50.1	88.9	48.9
LVU + Pred.	97.7	11.4	86.1	63.5	93.6	37.0
ViLU (ours)	98.3	8.2	89.5	47.4	$\boldsymbol{95.2}$	25.2

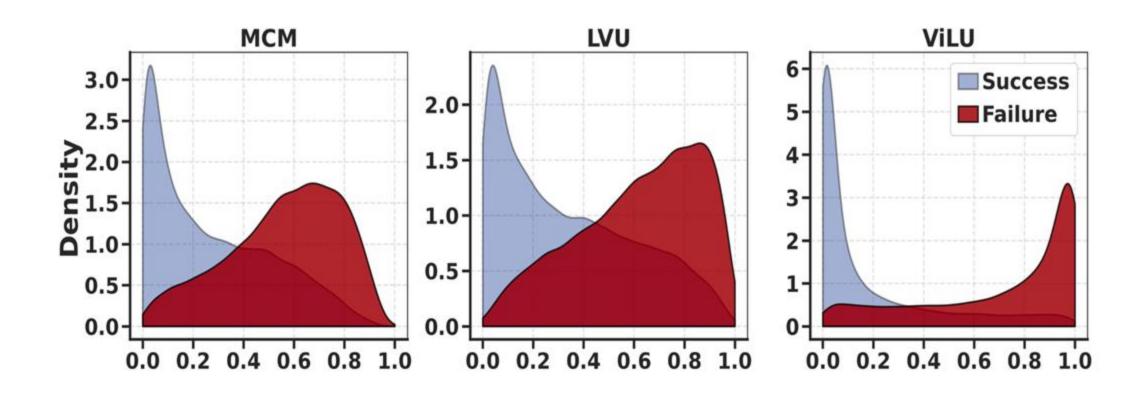
• Importance of ViLU embeddings: XA + predicted class



Data efficiency: ViLU effective even with a small fraction of train set



Effective with low ZS accuracy ≠ MCM, BayesVLM



• ViLU better separation



GT: "container ship"

Pred: "ocean liner"



GT: "mailbox"

Pred: "birdhouse"

Misclassification detection with ViLU ≠ MCM or LVU

# Towards zero-shot uncertainty quantification

		FPR95↓	
Dataset	MCM	$\mathbf{LVU}$	m ViLU
CIFAR-10	52.1	77.2	54.2
CIFAR-100	67.3	83.8	59.9
Caltech101	68.7	82.5	48.8
Flowers102	68.0	96.8	$\boldsymbol{67.4}$
OxfordPets	59.9	93.1	58.1
Food 101	63.3	87.2	67.4
${\bf FGVCAircraft}$	82.9	94.5	82.3
${ m EuroSAT}$	87.6	88.2	85.7
DTD	77.9	93.1	78.2
SUN397	75.9	90.1	72.7
$\operatorname{StanfordCars}$	<b>73.4</b>	92.6	84.1
UCF101	68.9	90.4	63.8
Average	70.5	89.1	68.6

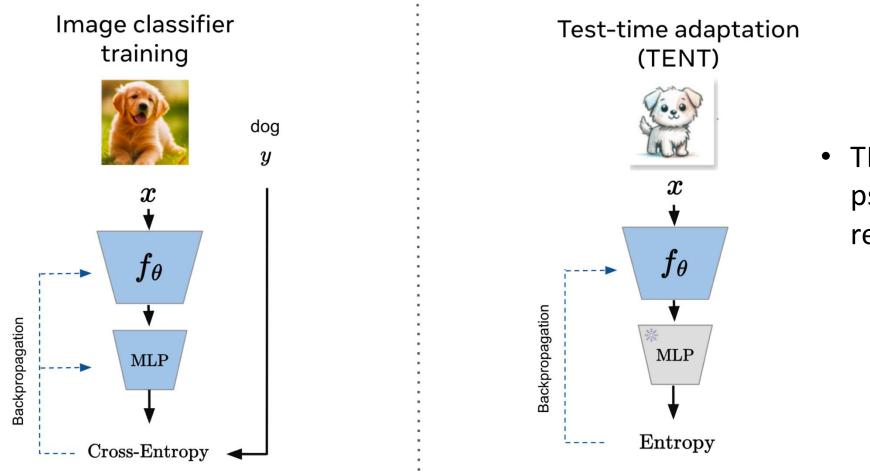
• ViLU pre-trained on CC12M transfers well on downstream datasets

# Outline

## 1. ViLU

# 2. CLIPTTA

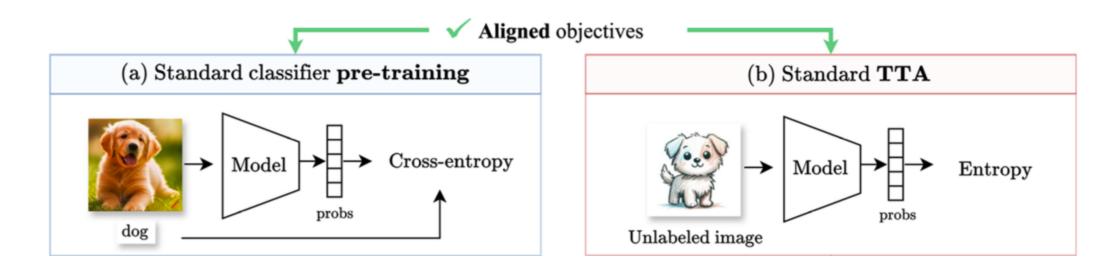
## Test-Time-Adapation (TTA) methods of vision models



 TENT-like methods: pseudo-labels, reinforce predictions

<sup>■</sup> Wang, Dequan, et al. "Tent: Fully test-time adaptation by entropy minimization." ICLR 2021.

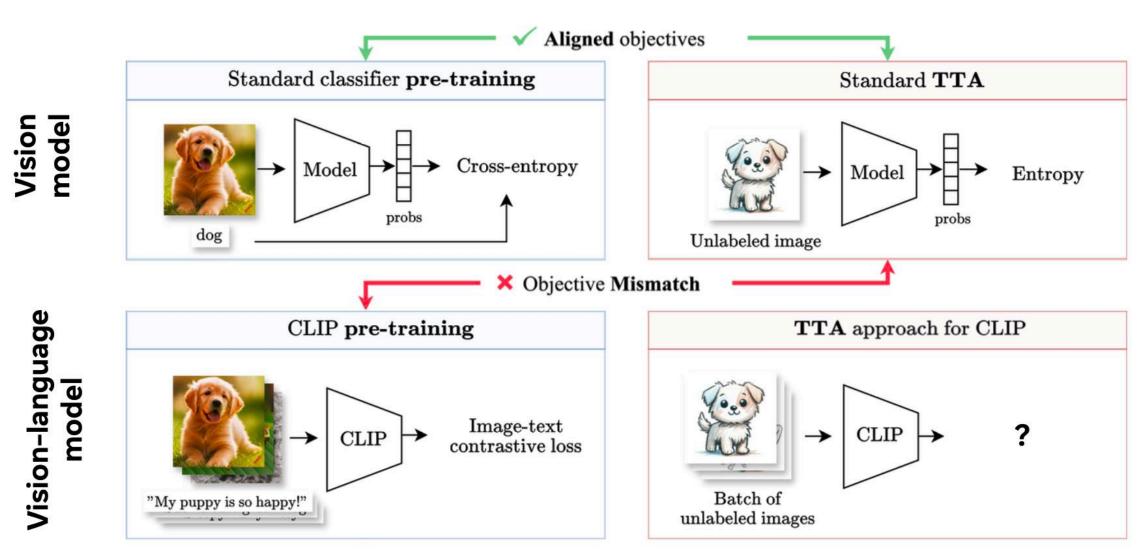
#### TENT-like TTA methods for vision models



Entropy minimization: soft version of cross-entropy

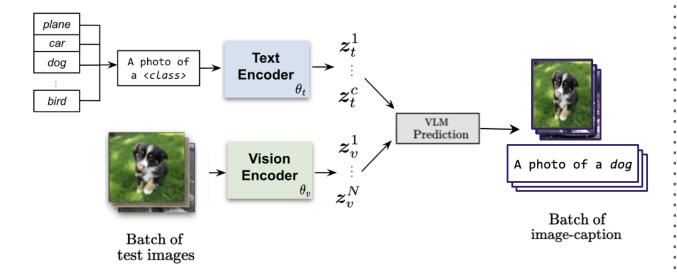
$$\mathcal{L}_{\mathrm{XE}} = -\log p(y|\boldsymbol{x}) \qquad \Longleftrightarrow \qquad \mathcal{L}_{\mathrm{Entropy}} = -\sum_{c} p(y_{c}|\boldsymbol{x}) \log p(y_{c}|\boldsymbol{x})$$
when  $p(y_{c}|\boldsymbol{x}) = \begin{cases} 1 & \text{if } y_{c} = y \\ 0 & \text{otherwise} \end{cases}$ 

# TTA on VLMS, e.g., CLIP: objective mismatch



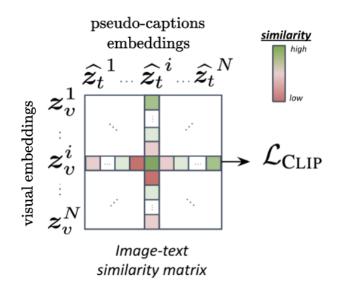
#### CLIP-TTA: how to use CLIP objective at test-time?

#### 1) Pseudo-captioning



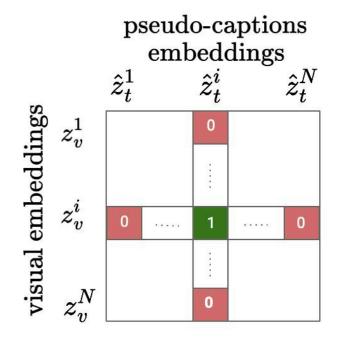
 Associate a pseudo-caption to each test image

#### 2) Contrastive loss



 Compute CLIP loss using batch of paired image-caption

#### CLIP-TTA: soft-contrastive loss

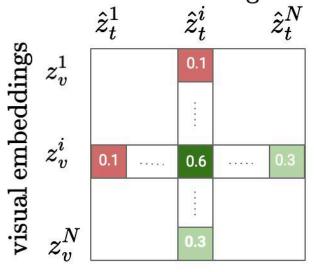


$$\mathcal{L}_{ ext{CLIP}}( heta) \coloneqq -\sum_{i=1}^{N} \left[ \qquad \underbrace{\log p(\hat{m{t}}_i | m{x}_i)}_{ ext{image} o ext{text}} + \qquad \underbrace{\log p(m{x}_i | \hat{m{t}}_i)}_{ ext{text} o ext{image}} 
ight]$$

- Use CLIP loss for the predicted text / caption
- Does not account for the uncertainty among the pseudo-captions

#### CLIP-TTA: soft-contrastive loss

pseudo-captions embeddings



$$egin{aligned} \mathcal{L}_{ ext{CLIP}}( heta) \coloneqq -\sum_{i=1}^{N} \left[ & \underbrace{\log p(\hat{oldsymbol{t}}_i | oldsymbol{x}_i)}_{ ext{image} o ext{text}} & + & \underbrace{\log p(oldsymbol{x}_i | \hat{oldsymbol{t}}_i)}_{ ext{text} o ext{image}} 
ight] \ \mathcal{L}_{ ext{s-cont}}( heta) \coloneqq \sum_{i=1}^{N} \left[ -\sum_{j=1}^{N} p(\hat{oldsymbol{t}}_j | oldsymbol{x}_i) \log p(\hat{oldsymbol{t}}_j | oldsymbol{x}_i) & -\sum_{j=1}^{N} p(oldsymbol{x}_j | \hat{oldsymbol{t}}_i) \log p(oldsymbol{x}_j | \hat{oldsymbol{t}}_i) 
ight] \end{aligned}$$

 $image \rightarrow text$ 

- Still aligned with CLIP loss, exploit uncertainty among the pseudo-captions
- => Entropy version of CLIP loss
- Loss computed on batch + on a memory batch  ${\cal M}$  + regularization loss  ${\cal L}_{
  m reg} = -\sum_{c=1}^C ar q_c \log ar q_c$

$$\mathcal{L}_{ ext{CLIPTTA}}( heta) = rac{1}{2} \Big[ \mathcal{L}_{ ext{s-cont}}( heta) + \mathcal{L}_{ ext{s-cont}}^{\mathcal{M}}( heta) \Big] + \lambda_{ ext{reg}} \mathcal{L}_{ ext{reg}}( heta),$$

text→image

# CLIP-TTA soft-contrastive loss: properties Gradient analysis

#### **Entropy loss**

$$abla_{oldsymbol{z}_i} \mathcal{L}_{ ext{TENT}} = -\sum_{k=1}^C \; \left[ \sum_{c=1}^C \log rac{q_{ik}}{q_{ic}} \; q_{ic} \; 
ight] \, q_{ik} \; oldsymbol{z}_t^k$$

- Always points towards predicted class
- Always Reinforces errors

#### Soft contrastive loss

$$abla_{oldsymbol{z}_i} \mathcal{L}_{ ext{CLIPTTA}} = \sum_{j=1}^N eta_{i,j} [-\widehat{oldsymbol{z}}_t^j + \sum_{k=1}^C w_{k,i} \ oldsymbol{z}_t^k]$$

- Can points towards a class ≠ prediction
- Batch-aware gradient: leverage correct predictions to correct errors
- Mitigate pseudo-label drift

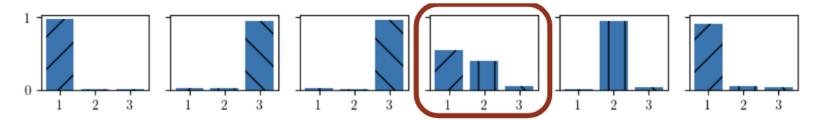
# CLIP-TTA soft-contrastive loss: properties

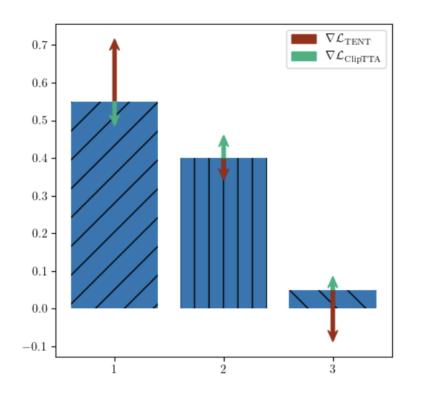
#### Toy example:

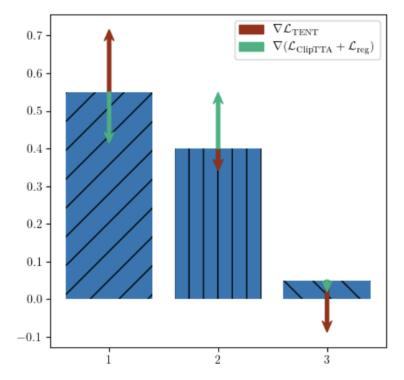
3-class classification problem, batch of size 6, sample 4 wrongly classified

- CLIPTTA gradient points towards correct class
- Can be combined with a regularization loss

$$\mathcal{L}_{ ext{reg}} = -\sum_{c=1}^C ar{q}_c \log ar{q}_c$$

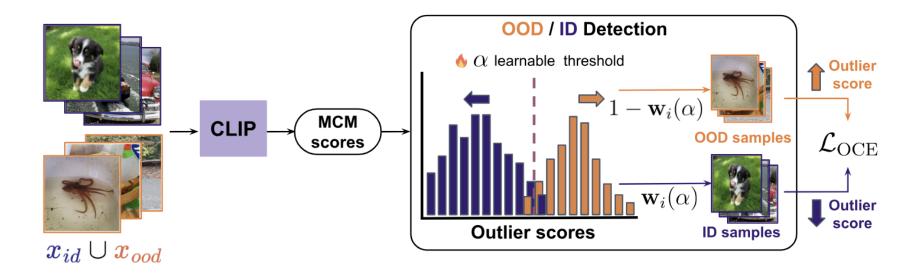






## Extension to open-set TTA

• Deal with In-Distribution (ID) & Out-Of-Distribution (OOD) samples in a batch



- Use MCM confidence, learn a separation threshold  $\alpha => w_i = \operatorname{sigmoid}(s_i \alpha)$
- Outlier Contrastive Exposure (OCE) loss to increase ID/OOD separability

$$\mathcal{L}_{ ext{OCE}} = - \left[ \underbrace{ \sum_{i=1}^{N} w_i s_i}_{\mu_{ ext{id}}} - \underbrace{ \sum_{i=1}^{N} (1 - w_i) s_i}_{\mu_{ ext{ood}}} \right]^2.$$

## Experiments

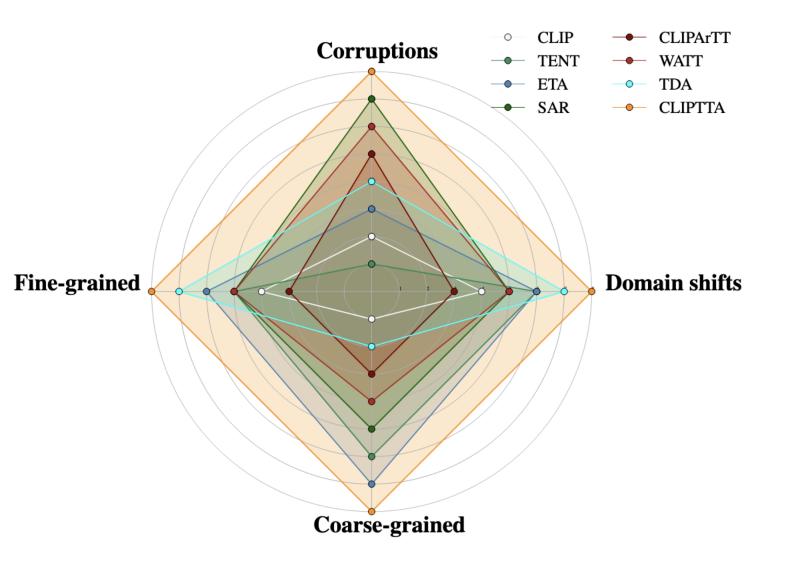
#### Experiments on 75 datasets

- Corruptions: 15 corruptions applied to CIFAR-10, CIFAR-100, ImageNet
- Domain shifts: VisDA-C, PACS, OfficeHome, Imagenet-Domains
- Coarse grained classification: CIFAR-10, CIFAR-100
- Fine-grained classification: Imagenet, + 10 datasets from the CLIP zeroshot suite

#### Baselines

- TENT-like: ETA, SAR, RoTTA, + CLIP TENT-like: CLIPArTT, WATT
- More generic methods:
  - Test-Prompt-Tuning (TPT)
  - Parameter-free: Training-free Dynamic Adapter (TDA)
- Results in non-episodic setting (no reset)

#### Overall results



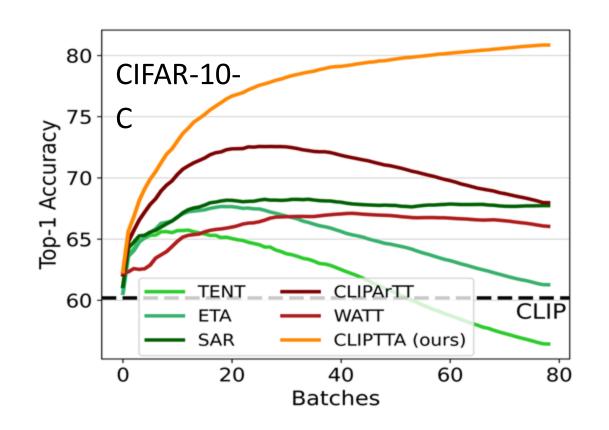
- Outperforms SOTA results on each class of dasets
- ≠ Specific methods effective only on some datasets, e.g., TDA on ImageNet

#### CLIPTTA Results: detailed results

	Corruptions			Domain shifts					
	CIFAR-10-C	CIFAR-100-C	Imagenet-C	Average	VisDA-C	PACS	OfficeHome	ImageNet-D	Average
CLIP	60.2	35.2	25.5	40.3	87.1	96.1	82.5	59.4	81.3
TPT (NeurIPS '22)	58.0	33.6	24.6	38.7	85.0	94.0	81.7	62.4	80.8
TDA (CVPR '24)	63.4	37.4	26.8	42.5	86.6	96.1	83.0	65.0	<u>82.8</u>
TENT (ICLR '21)	56.4	31.4	17.6	35.1	89.3	96.6	83.4	60.2	82.3
ETA (ICML '22)	61.3	38.9	26.8	42.3	88.3	96.7	84.1	59.9	82.3
SAR (ICLR '23)	67.8	43.2	33.6	48.2	87.8	96.2	83.8	60.6	82.1
RoTTA (CVPR '23)	58.0	33.6	24.6	38.7	83.7	95.8	82.5	61.6	80.9
CLIPArTT (WACV '25)	<u>68.1</u>	<u>48.0</u>	33.3	<u>49.8</u>	84.1	96.3	82.0	60.7	80.8
WATT (NeurIPS '24)	66.0	38.5	26.0	43.5	87.7	96.2	83.4	61.8	82.1
CLIPTTA (ours)	80.7	52.5	41.1	58.1	89.6	97.5	84.2	<u>63.4</u>	83.7

Huge performance boost on corrupted datasets, where CLIP's ZS performance low

# CLIP-TTA: robustness to pseudo-label drift

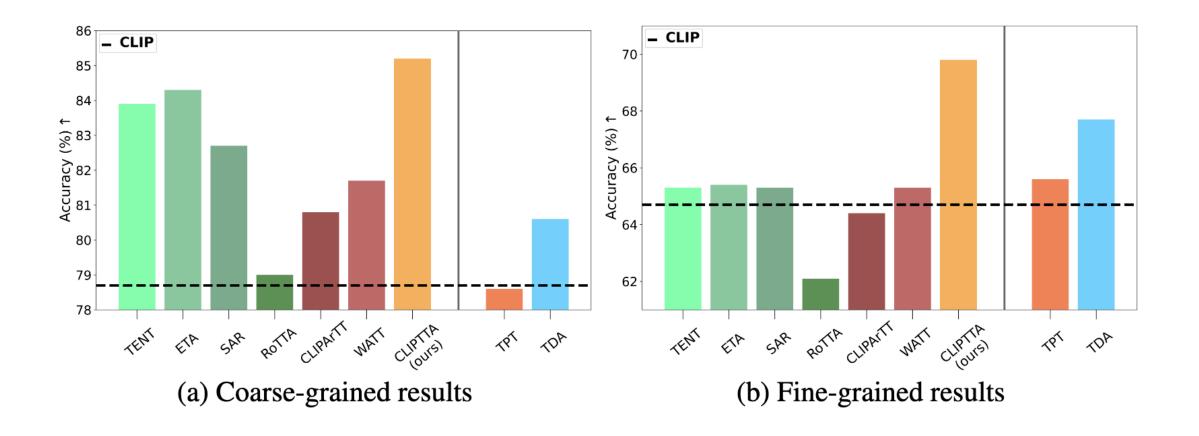


	C-100	C-100-C	IN	IN-C	Avg.
CLIP	68.1	35.2	66.7	25.5	48.9
TENT	72.9	31.4	66.5	17.6	47.1
$\mathcal{L}_{ ext{s-cont}}$ $\mathcal{L}_{ ext{s-cont}} + \mathcal{L}_{ ext{reg}}$ $\mathcal{L}_{ ext{s-cont}} + \mathcal{L}_{ ext{reg}} + \mathcal{M}$	74.2	50.8	68.8	40.3	58.5
	74.9	52.4	69.1	38.6	58.8
	<b>75.3</b>	<b>52.6</b>	<b>69.6</b>	<b>41.1</b>	<b>59.6</b>

- CLIPTTA improves during adaptation
- TENT-based methods collapse

- Big improvements due to  $\mathcal{L}_{s\text{-cont}}$
- Reg and memory can further boost results

# CLIP-TTA: Coarse/fine-grained classification



## CLIP-TTA: open-set TTA

	ACC↑	AUC↑	FPR95↓
CLIP [1]	66.7	90.1	43.8
TENT [11]	12.4	49.9	89.4
ETA [12]	<u>67.1</u>	89.6	46.1
SAR [13]	58.8	62.0	75.7
CLIPArTT [6]	31.2	61.1	87.5
WATT [7]	<u>67.1</u>	87.4	53.4
TDA [4]	66.8	82.1	59.8
CLIPTTA (ours)	67.6	93.5	25.7
OSTTA [17] †	66.9	84.9	59.2
SoTTA[26] †	66.7	89.3	47.1
STAMP [27] †	29.7	63.0	80.2
UniEnt [28] †	65.2	<u>95.4</u>	<u>17.1</u>
CLIPTTA + OCE (ours) †	67.6	97.7	9.7

- ImageNet ID, Places OOD
- OOD detection: baselines <</li>
   CLIP ≠ CLIPTTA (ours)
- CLIPTTA + OCE boost OOD detection with maintained accuracy

# Conclusion & perspectives

- ViLU: failure prediction for VLMs
- CLIPTTA: open-set TTA for CLIP, OOD detection

#### **Perspectives**

- ViLU: generalization & combining failure, OOD + calibration
- Extension to semantic segmentation & theoretical guarantees (conformal)
- CLIPTTA: detecting ID failures before adaptation, UQ to overcome low ZS accuracy & pseudo-label drift

# Thank you for your attention!

# •Questions?





















