

Brian Pulfer<sup>\*1</sup>, Yury Belousov<sup>\*1</sup>, Vitaliy Kinakh<sup>1</sup>, Teddy Furon<sup>2</sup>, Slava Voloshynovskiy<sup>1</sup>

<sup>1</sup>University of Geneva

<sup>2</sup>University of Rennes, Inria, CNRS, IRISA

## Introduction

VFM have become increasingly popular and most tasks (classification, segmentation, captioning, VQA, ...) today are tackled through pre-trained VFMs.

We introduce Task Agnostic Attacks (TAAs), which degrade performances across tasks by maximally perturbing feature representations of VFMs independently of the task. We find that feature space is easily manipulated and that our attacks are competitive with PGD task-specific attacks (TSAs).

## Method

We adversarially modify images as to minimize the cosine similarity of the features extracted with the VFM

$$\mathcal{L}(\mathbf{x}_{\text{adv}}) = \frac{f(\mathbf{x})f(\mathbf{x}_{\text{adv}})}{\|f(\mathbf{x})\| \|f(\mathbf{x}_{\text{adv}})\|}$$

Where  $f$  is the feature extraction backbone. For ViTs models, we compute the average cosine similarity across patch tokens.

## Experimental Results

Backbone	Attack	Type	Classification abs↓ (rel↑)	Segmentation abs↓ (rel↑)
ViT-S	No attack		96.3 (0%)	81.4 (0%)
	Class token	TAA	7.9 (92%)	19.2 (86%)
	Patch tokens	TAA	<b>0.1 (100%)</b>	<b>11.6 (97%)</b>
	Class+patch tokens	TAA	2.0 (98%)	13.3 (94%)
	Classification	TSA	<b>0.0 (100%)</b>	19.8 (85%)
	Segmentation	TSA	40.6 (58%)	<b>9.2 (100%)</b>
ViT-B	No attack		97.0 (0%)	80.8 (0%)
	Class token	TAA	11.8 (88%)	23.5 (80%)
	Patch tokens	TAA	<b>0.0 (100%)</b>	<b>7.5 (102%)</b>
	Class + patch tokens	TAA	2.1 (98%)	11.8 (96%)
	Classification	TSA	<b>0.0 (100%)</b>	14.1 (93%)
	Segmentation	TSA	43.9 (55%)	<b>8.9 (100%)</b>

Attack	Captioning COCO				Question answering VQAv2		
PSNR	BLEU-4 ↓	METEOR ↓	ROUGE-L ↓	CIDEr ↓	number ↓	yes/no ↓	other ↓
No attack	29.6	30.3	59.0	131.4	72.5	95.9	76.9
45 dB	5.8	17.3	32.2	33.0	50.6	83.4	53.4
40 dB	3.8	13.6	27.2	16.4	38.7	76.0	41.6
35 dB	1.9	9.7	22.3	3.6	25.6	67.5	28.5