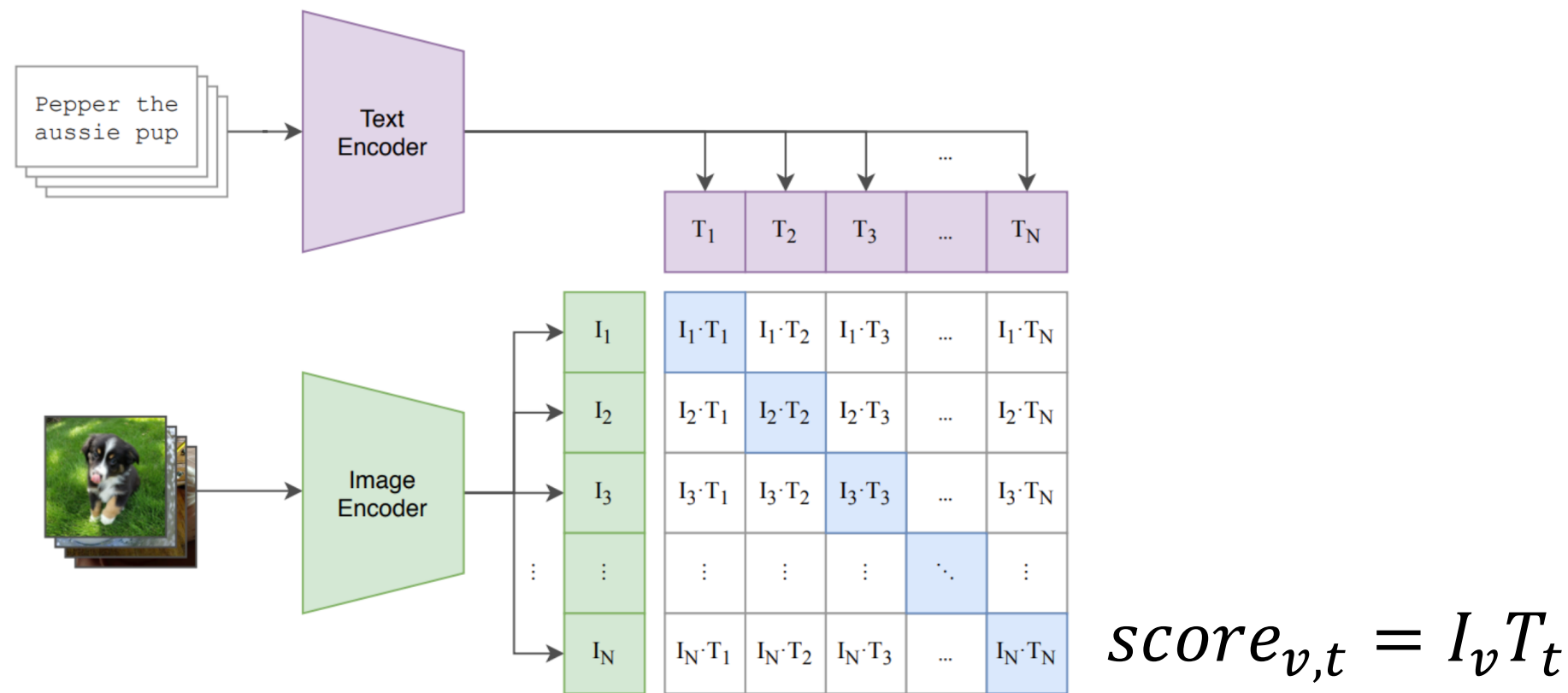


Motivation

- Text-to-image, image-to-text generative models still generate output **misaligned with their inputs**.
- While **CLIPScore** is widely used for measuring alignment between image and text, single scalar score **lacks interpretability**.
- Recent studies focus on detecting dense misalignments to enhance explainability and provide model feedback.
- However, most approaches heavily rely on large generative models, leading to **high computational costs**.

Preliminaries

CLIP



Gradient-based relevance map method

- Existing gradient-based relevance map extraction methods (e.g., GradCAM) primarily *focus on generating heatmaps* from images. They mostly **disregard negative gradients**, assuming these represent unrelated information.

- In transformer literature, GAE (Chefer et al., '21) derives output logits with respect to attention maps.

$$\nabla A_l^h = \frac{\partial \text{score}_{v,t}}{\partial A_l^h},$$

($A_l^h \in \mathbb{R}^{n \times n}$ denote attention map at l-th layer h-th head)

- Then *eliminates negative gradients* and applies an element-wise product with the attention map itself.

$$R_l^h = \text{ReLU}(\nabla A_l^h \odot A_l^h). \quad R_l = \frac{1}{H} \sum_{h=1}^H R_l^h.$$

- Across layers, relevance map is updated with $R_l \leftarrow R_l + R_l \odot R_{l-1}$. Final attribution is calculated by pooling [CLS] index row of relevance map.

CLIP for Dense Misalignment (CLIP4DM)

CLIP4DM

1. We **remove ReLU** so that negative gradients can represent misalignments.
2. We aggregate relevance map per layer with **averaging** instead of a product.
3. Then, we merge tokens into a word by **averaging their attributions** and predict a word whose attribution is **lower than epsilon** as misaligned.



Caption: A man riding snowboard down a snow covered slope.

CLIPScore: 61.3

Ours: A man riding snowboard down a snow covered slope.

Misaligned word: snowboard

F-CLIPScore

To classify global misalignments, we propose F-CLIPScore which aggregate sum of negative attributions and CLIPScore.

$$\text{mis}(w_j) = \begin{cases} 1, & \text{if } w_j < \epsilon \\ 0, & \text{otherwise.} \end{cases}$$

$$\text{F-CLIPScore}(v, t) = (1 - \text{score}_{v,t}) \cdot \sum_j \text{mis}(w_j) \cdot w_j.$$

Experiments

We experiment five dense misalignment detection benchmarks.

- **① FOIL** : natural image & natural caption.
- **② nocaps-FOIL** : natural image & natural caption.

Results

- 1) Our results demonstrate **state-of-the-art performance** in *detecting dense misalignments* among zero-shot models.
- 2) F-CLIPScore **outperforms CLIPScore** in *classifying misaligned pairs*.
- 3) Our method achieves **significantly higher efficiency**, with **27–75 × better FPS**, as it relies solely on CLIP, unlike baselines that depend on large generative models.

Method	FPS	LA	AP	LA	AP
CLIPScore (ViT-B/32)	13.4	-	0.71	-	0.69
CLIPScore (ViT-H/14)	8.7	-	0.76	-	0.72
RefCLIPScore (ViT-B/32)	8.7	-	0.75	-	0.74
ALoHa	0.2	0.40	0.61	0.45	0.70
Ours (ViT-B/32)	12.0	0.73	0.71	0.60	0.69
Ours (ViT-H/14)	7.1	0.84	0.81	0.72	0.80

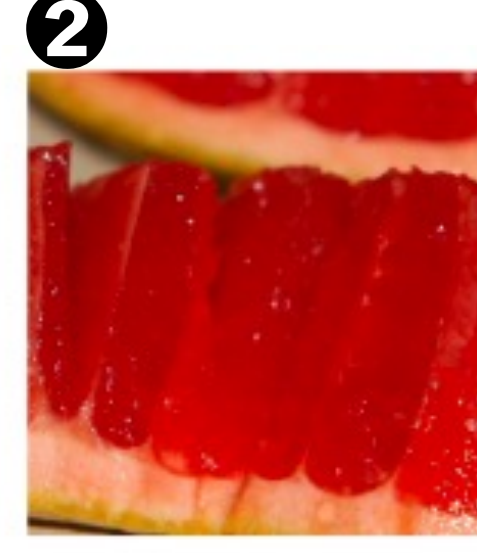
method	ref. captions	FPS	LA	AP	Model	ft.	FPS	NLI score
CLIPScore		18.8	-	0.39	PaLI 5B		-	0.226
RefCLIPScore	✓	9.0		0.43	MiniGPT-v2 (LLaMa2-7B)		0.28	0.560
ALoHa	✓	0.2	0.20	0.49	Ours (ViT-B/32)		7.90	0.605
Ours (ViT-B/32)		9.6	0.19	0.36	Ours (ViT-H/14)		5.81	0.660
Ours (ViT-H/14)		6.6	0.35	0.36	PaLI 5B	✓	-	0.765
					PaLI 17B	✓	-	0.785

Model	ft.	F1	precision	recall	Model	ft.	pearson	spearman
ALoHa		0.34	0.31	0.39	CLIPScore (ViT-B/32)		0.19	0.13
Ours (ViT-B/32) $\epsilon=-0.00001$		0.40	0.33	0.50	PickScore (ViT-H/14)		0.35	0.34
Ours (ViT-H/14) $\epsilon=-0.00001$		0.43	0.37	0.52	Ours (ViT-B/32) $\epsilon=-0.00001$		0.28	0.33
Ours (ViT-H/14) $\epsilon=-0.00005$		0.31	0.49	0.23	Ours (ViT-H/14) $\epsilon=-0.00001$		0.37	0.43
Rich-HF (multi-head)	✓	0.43	0.63	0.33	CLIPScore (ViT-B/32)	✓	0.40	0.39
Rich-HF (augmented prompt)	✓	0.44	0.61	0.341	Rich-HF (multi-head)	✓	0.49	0.50
					Rich-HF (augmented prompt)	✓	0.47	0.50



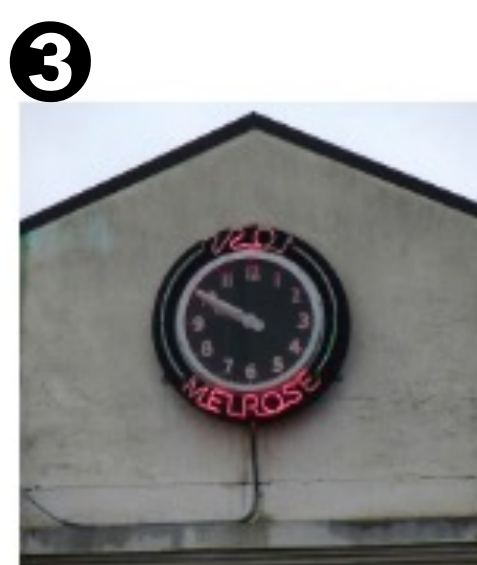
a woman wearing a net on her head cutting a pizza.

Ours: pizza



Slices are cut into a wedge of watermelon.

Ours: watermelon



A large clock mounted to the side of a building in a city with red neon lettering spelling out the word "welcome".

Ours: welcome



an image of the character batman inside a clock

Ours: batman



Aliens invading McDonald's

Ours: McDonald's



Santa Monica pierr in summer

Ours: Santa Monica

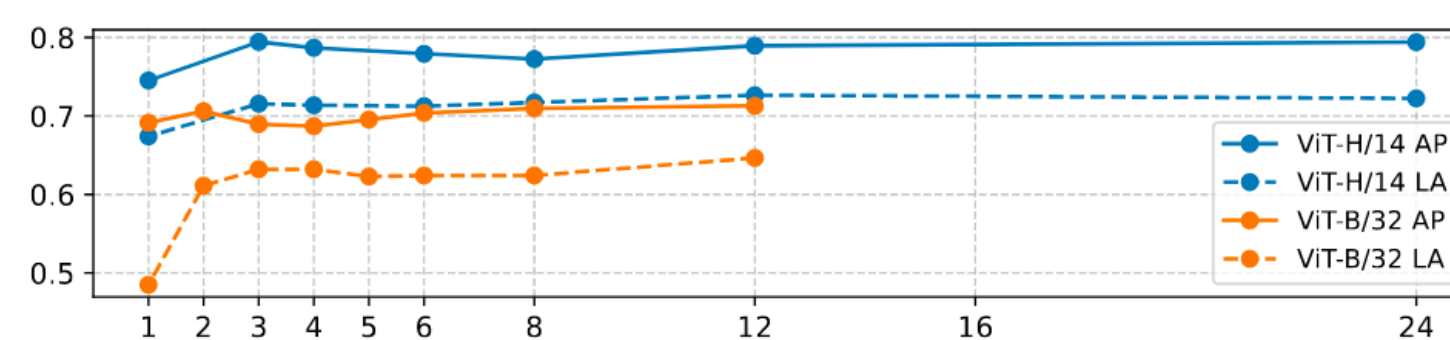
Analysis

Ablation on full-gradients

$\text{ReLU}(-\nabla A_l^h)$	$\text{ReLU}(-\nabla A_l)$	LA	AP
✓		0.698	0.779
	✓	0.700	0.776
		0.716	0.794

Using both negative and positive attributions is helpful.

Number of layers



Part-of-speech

Metric	NOUN	PROPN	VERB	ADV	ADJ	NUM	ADP
F1	0.393	0.312	0.301	0.258	0.258	0.132	0.177
Precision	0.470	0.602	0.567	0.444	0.417	0.500	0.278
Recall	0.337	0.211	0.205	0.182	0.187	0.076	0.130

While ours can extract all sorts of POS, it shows that it follows CLIP's strength and weakness.

Comparison with baselines



A small kitchen area with an oven and white fridge.

ALOHa: kitchen area
Ours: white



Two elephants with one calf standing in the dirt and ...

ALOHa: dirt
Ours: calf



A cat laying in a pan on a wooden floor.

ALOHa: pan
Ours: wooden floor



A street sign that reads 'broadway' on a traffic light.

MiniGPT-v2: The word 'broadway' on the sign does not align with the image. The image shows a ... but the word 'broadway' is not visible in the image.
Ours: **broadway**



A collection of stuffed animals and a doll that is holding a cell phone.

MiniGPT-v2: The doll holding a cell phone is not a stuffed animal.
Ours: **cell phone**



Two zebras standing in a dry grass field with trees

ALOHa: tree
Ours: two



A bird walking on a beach at either sunset or night, with waves ...

ALOHa: wave
Ours: sunset



A man is sitting on a bench outside, while the birds ...

ALOHa: bench
Ours: birds

Our method captures a *wider range of misalignments* (e.g., attributes, entity-level objects, intangible objects) compared to a combination of foundation models. (③)

Also our method output *concise and precise feedback* with much enhanced FPS compared to generative models. (④)