



Grounding Counterfactual Explanation of Image Classifiers to Textual Concept Space

Siwon Kim

Jinoh Oh

Sungjin Lee

Seunghak Yu

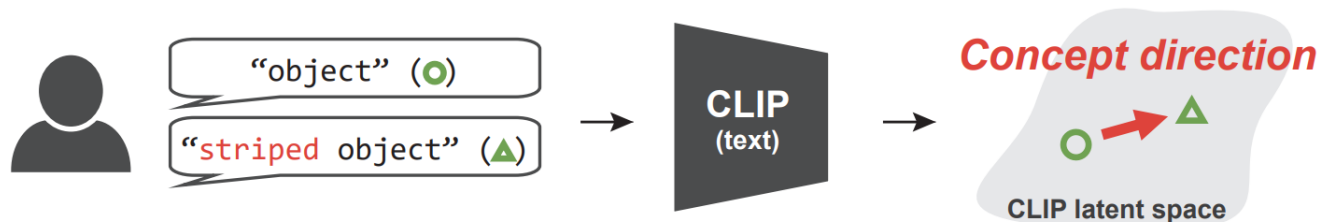
Jaeyoung Do

Tara Taghavi

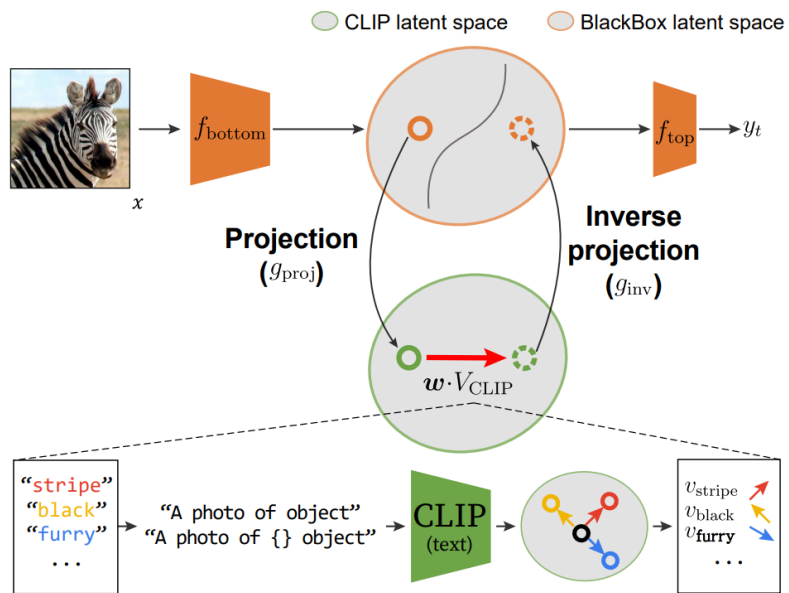


One-page Summary

Motivation

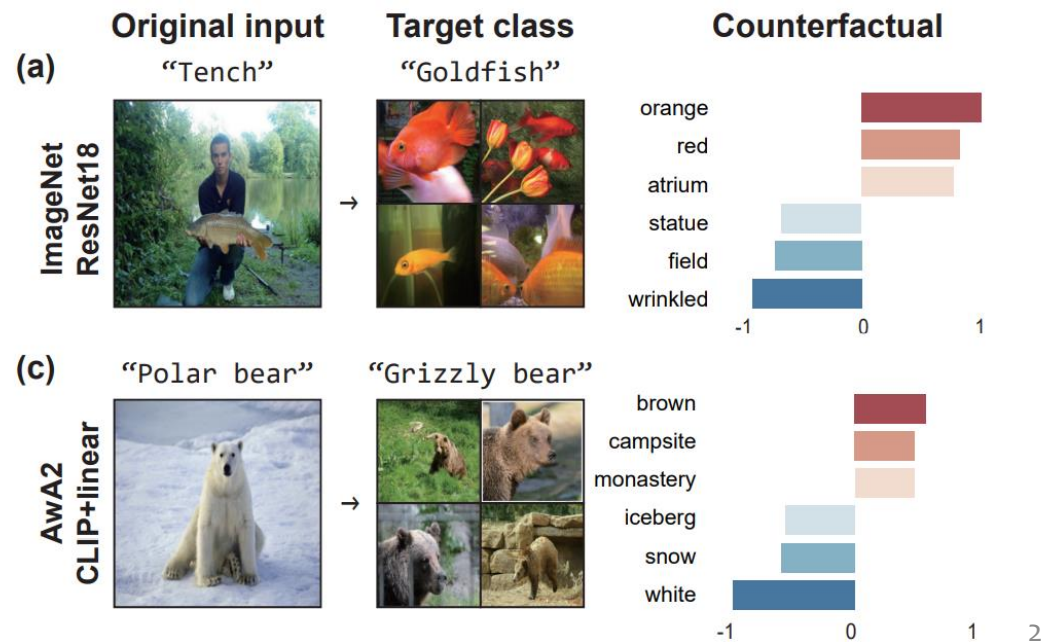


Method



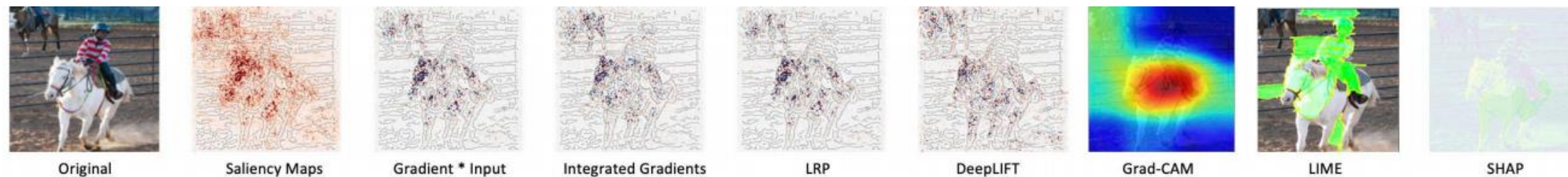
source:

Results

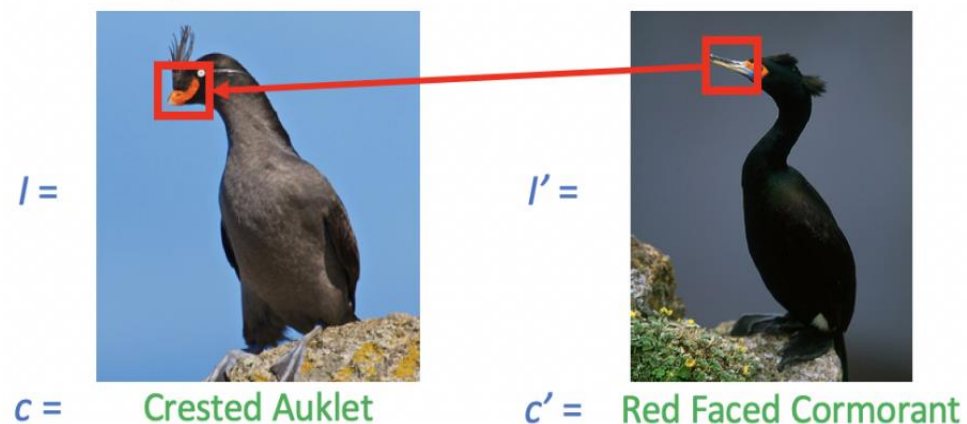


Are existing explanations human-understandable?

- Attribution-based explanation



- Visual counterfactual explanation



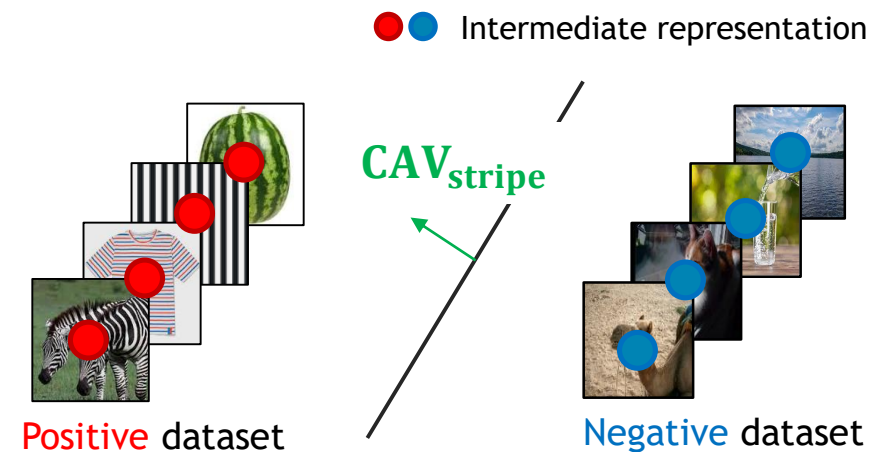
→ Conventional visual XAI methods

😊 Can: Provide important *regions*

🤔 Can't: Provide important *aspects*
(e.g. color? pattern?)

Towards human understandable explanation: Concept-based explanation

- Concept?
 - The units of human-understandable high-level semantics
 - Typically defined by words such as “stripe”, “white”, ...
- Concept Activation Vector (CAV)
 - **Positive** dataset: samples that **exhibit** c
 - **Negative** dataset: samples that **exclude** c
 - $CAV = A$ vector normal to the linear hyperplane



Example: “Stripe” concept



1. Diverging CAVs & 2. Unintended entanglement

I. Diverging CAVs

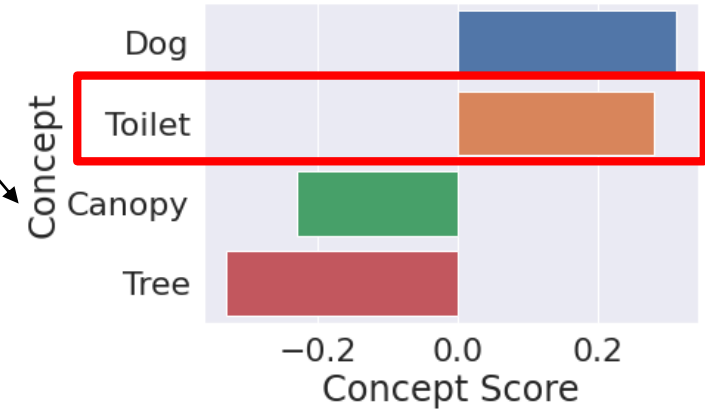
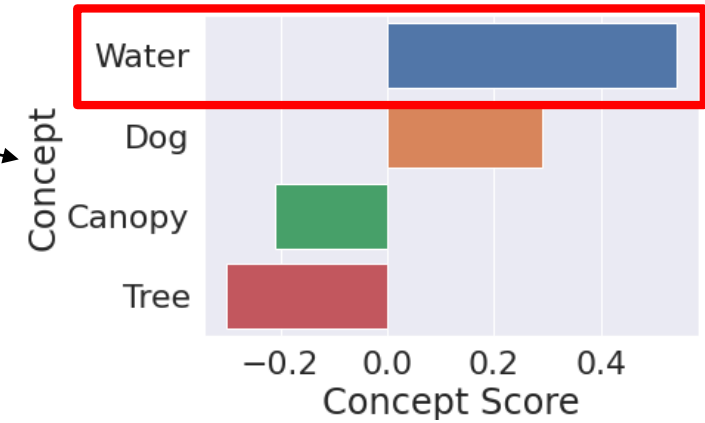
Positive dataset 1 for concept "water"



Positive dataset 2 for concept "water"



Negative dataset for concept "water"



2. Unintended entanglement



Positive images for “Green”

vs.

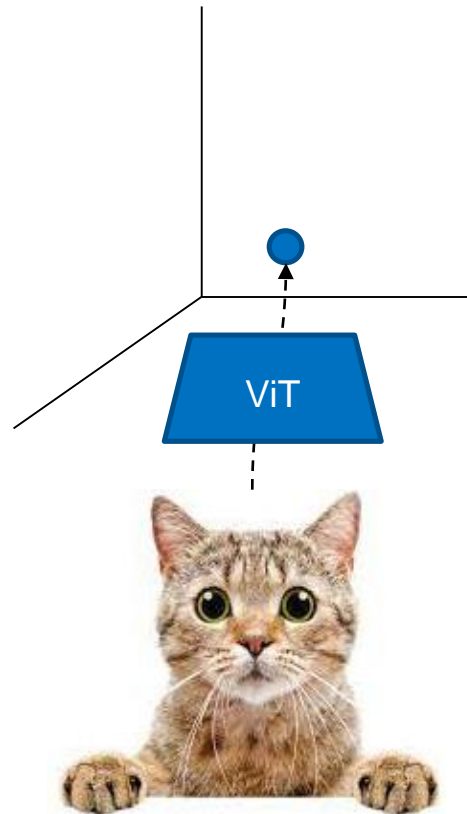


Positive images for “Grass”

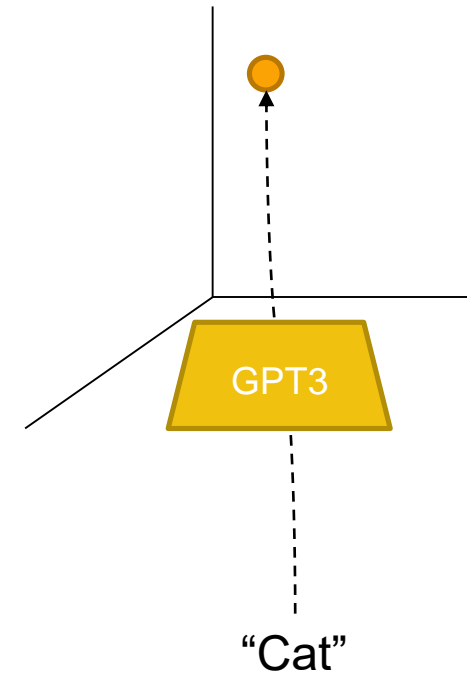
Counterfactual **TEXtual **EX**planation (Coun**TEX**)**

Vision (V) and Language (L) were Separated

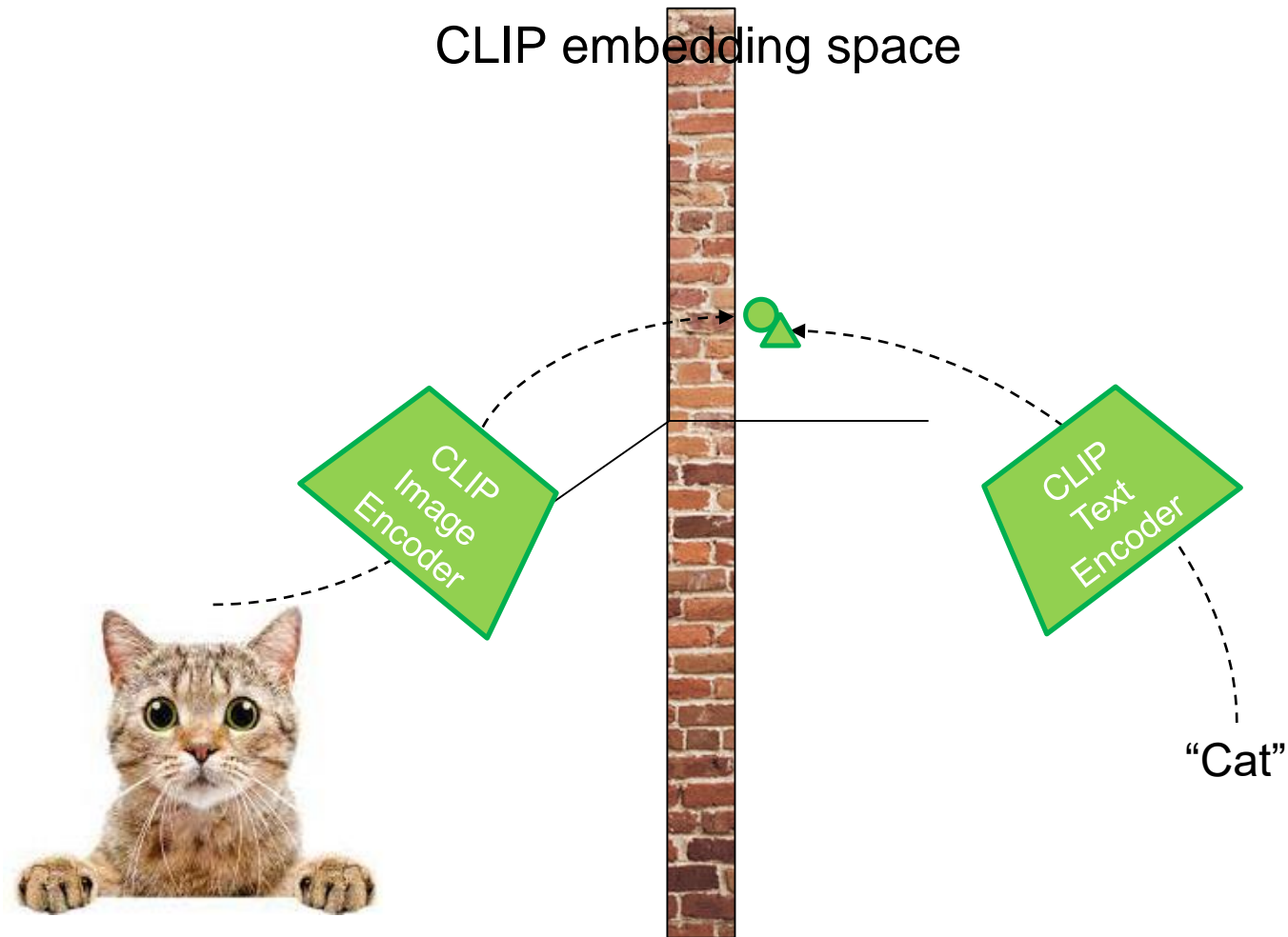
Image model embedding space



Language model embedding space

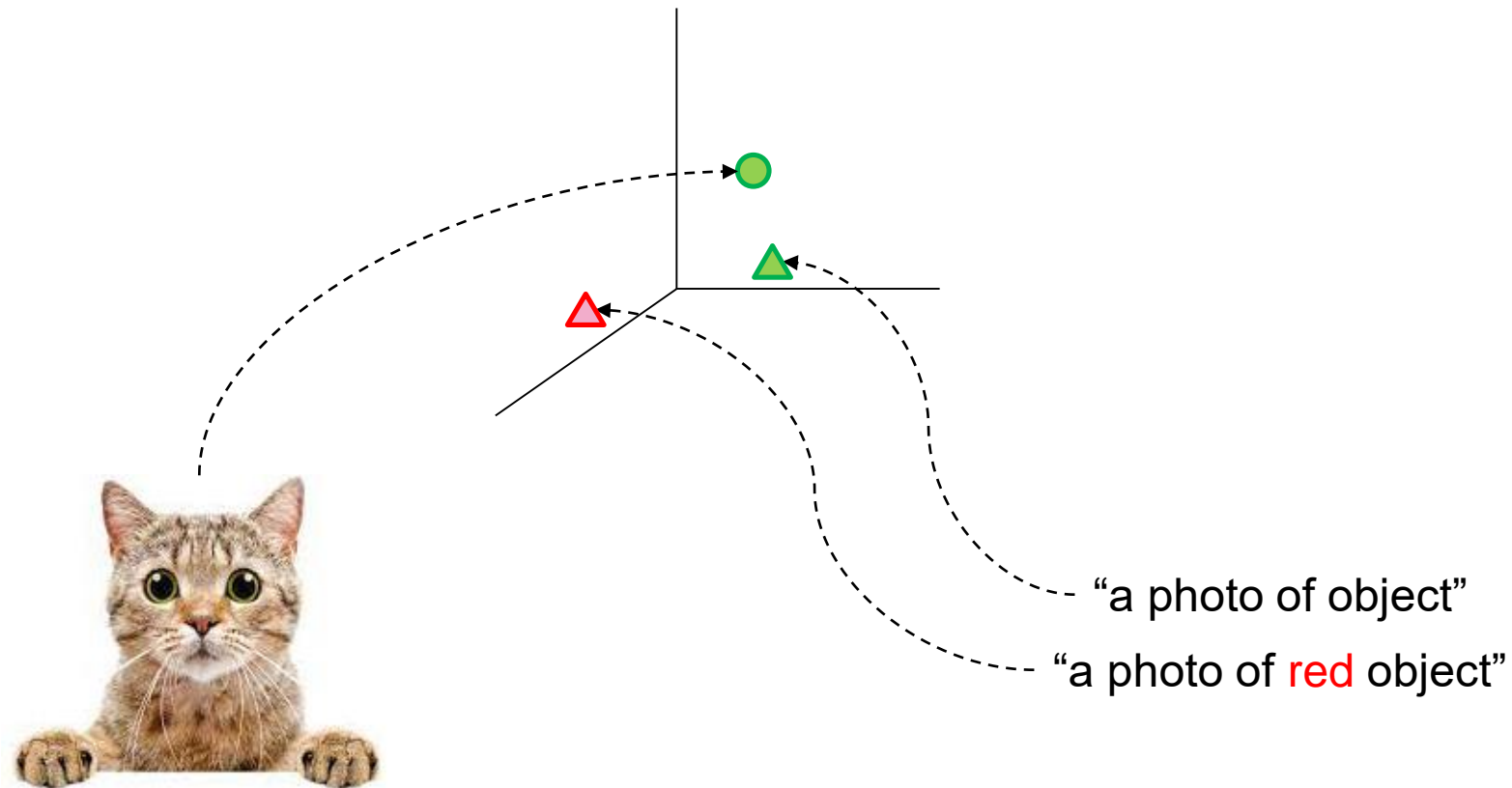


CLIP enables V&L Joint Embedding Space



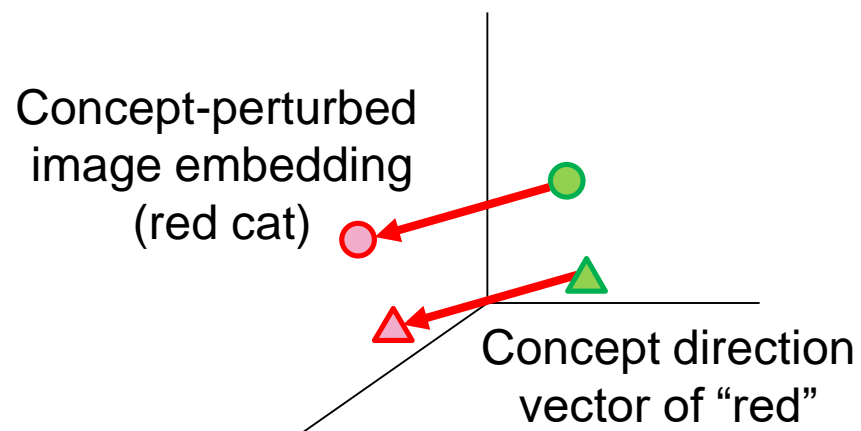
CLIP enables textual guidance on images

CLIP embedding space



CLIP enables textual guidance on images

CLIP embedding space

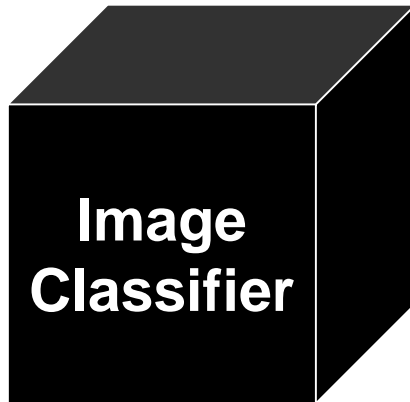


“a photo of cat”

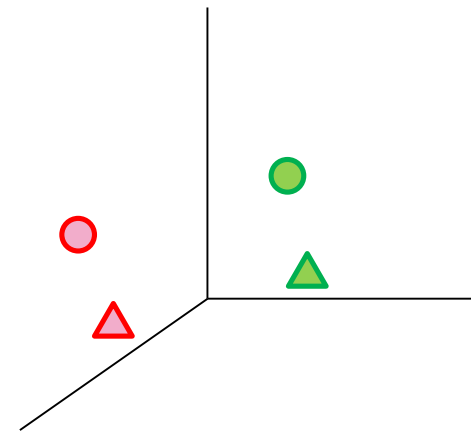
“a photo of red cat”

But, the target classifier is not CLIP

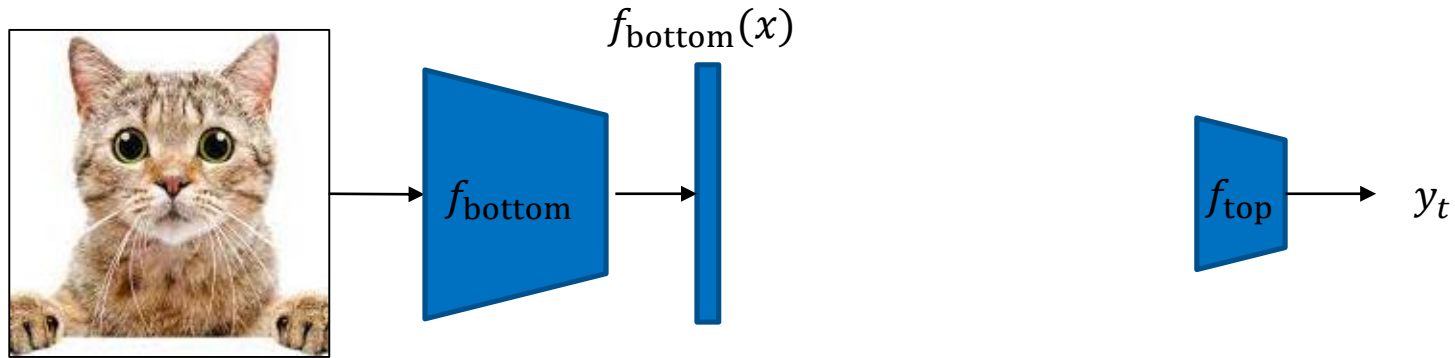
- How can we leverage the well-performing CLIP latent space?



CLIP embedding space

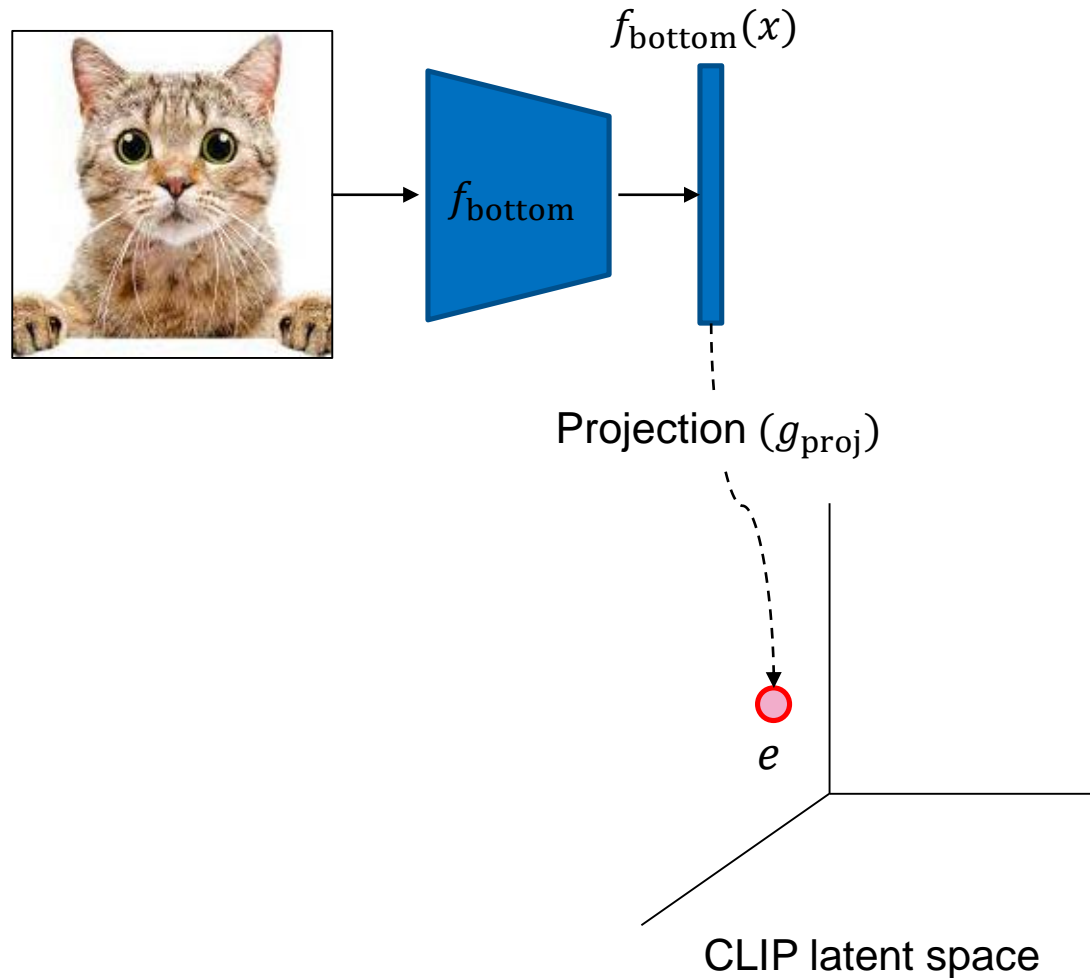


Counterfactual **TEX**tual **EX**planation (**CounTEX**)



$$1. f = [f_{\text{bottom}}, f_{\text{top}}]$$

Counterfactual **TEX**tual **EX**planation (**CounTEX**)

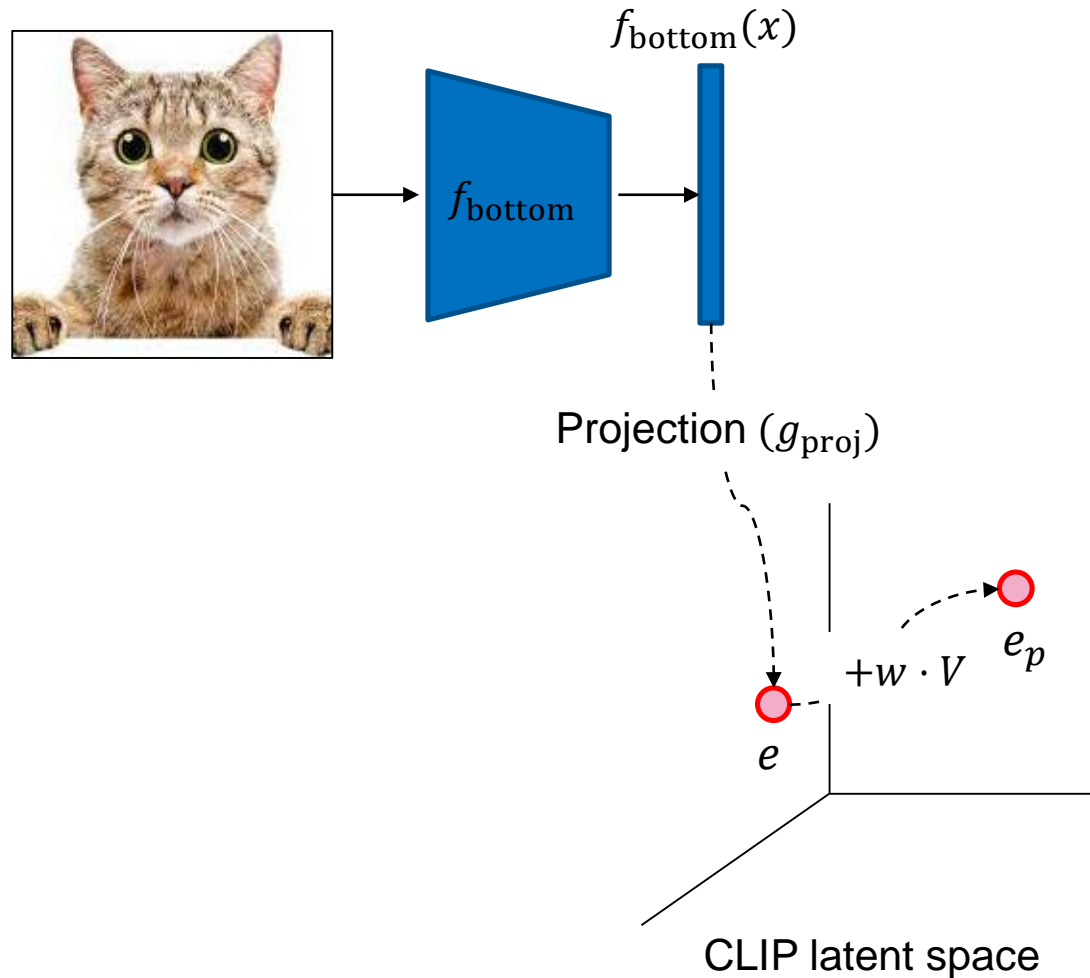


$$1. f = [f_{\text{bottom}}, f_{\text{top}}]$$

$$2. e = g_{\text{proj}}(f_{\text{bottom}}(x)),$$

where $g_{\text{proj}}(\cdot)$ is a projection function

Counterfactual **TEX**tual **EX**planation (**CounTEX**)



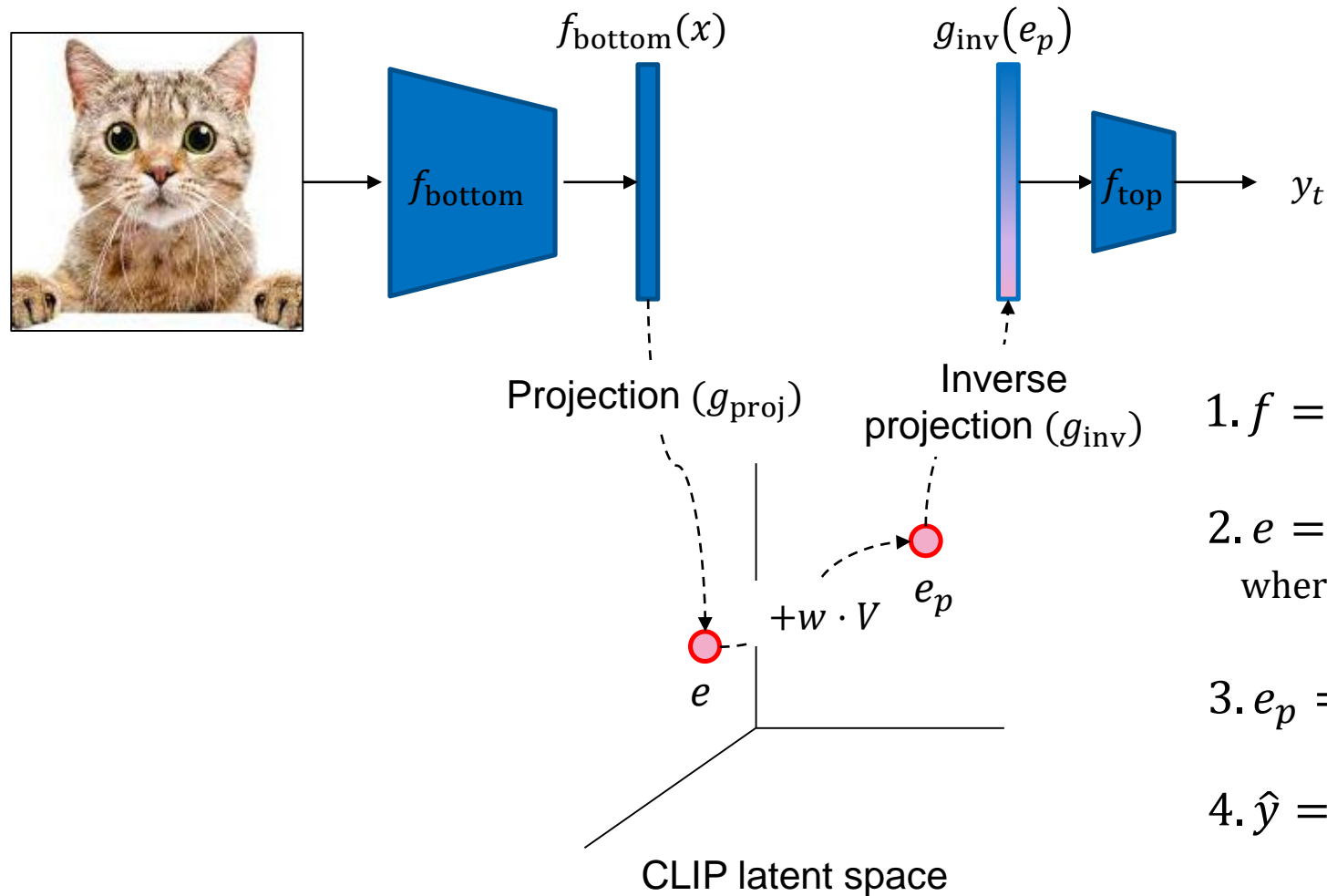
$$1. f = [f_{\text{bottom}}, f_{\text{top}}]$$

$$2. e = g_{\text{proj}}(f_{\text{bottom}}(x)),$$

where $g_{\text{proj}}(\cdot)$ is a projection function

$$3. e_p = e + w \cdot V$$

Counterfactual **TEX**tual **EX**planation (**CounTEX**)



$$1. f = [f_{\text{bottom}}, f_{\text{top}}]$$

$$2. e = g_{\text{proj}}(f_{\text{bottom}}(x)),$$

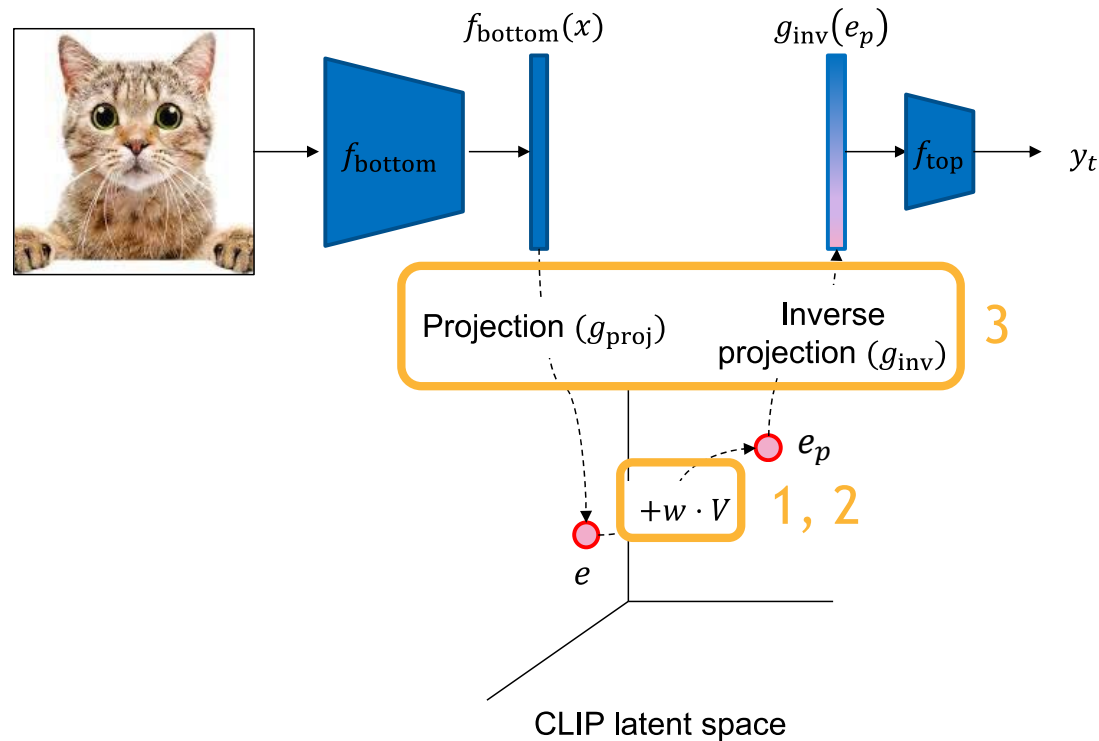
where $g_{\text{proj}}(\cdot)$ is a projection function

$$3. e_p = e + w \cdot V$$

$$4. \hat{y} = f_{\text{top}}(g_{\text{inv}}(e_p))$$

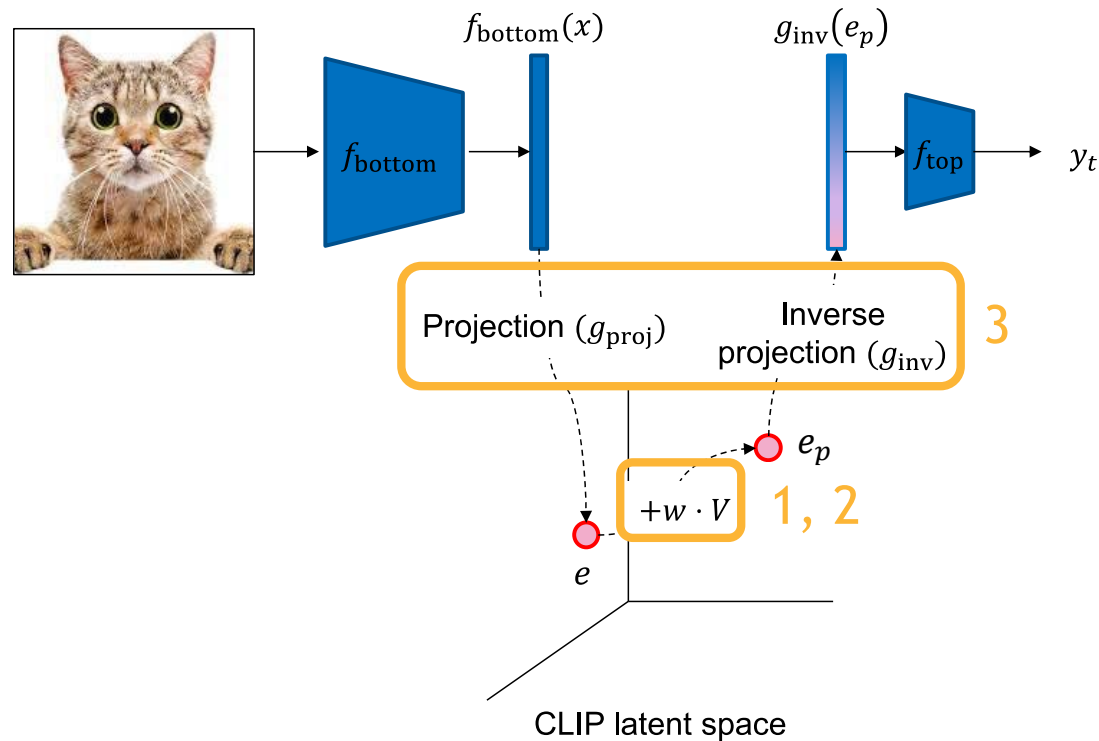
Research Questions

1. How can we obtain concept direction bank V ?
2. How can we obtain weight w ?
3. How can we implement projection and inverse projection?



Research Questions

1. How can we obtain concept direction bank V ?
2. How can we obtain weight w ?
3. How can we implement projection and inverse projection?

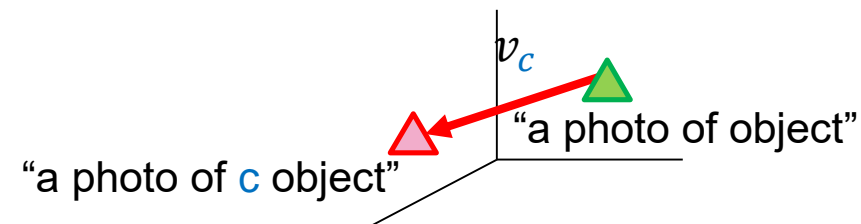


I. Prepare concept direction bank

I. Build concept direction vector bank V

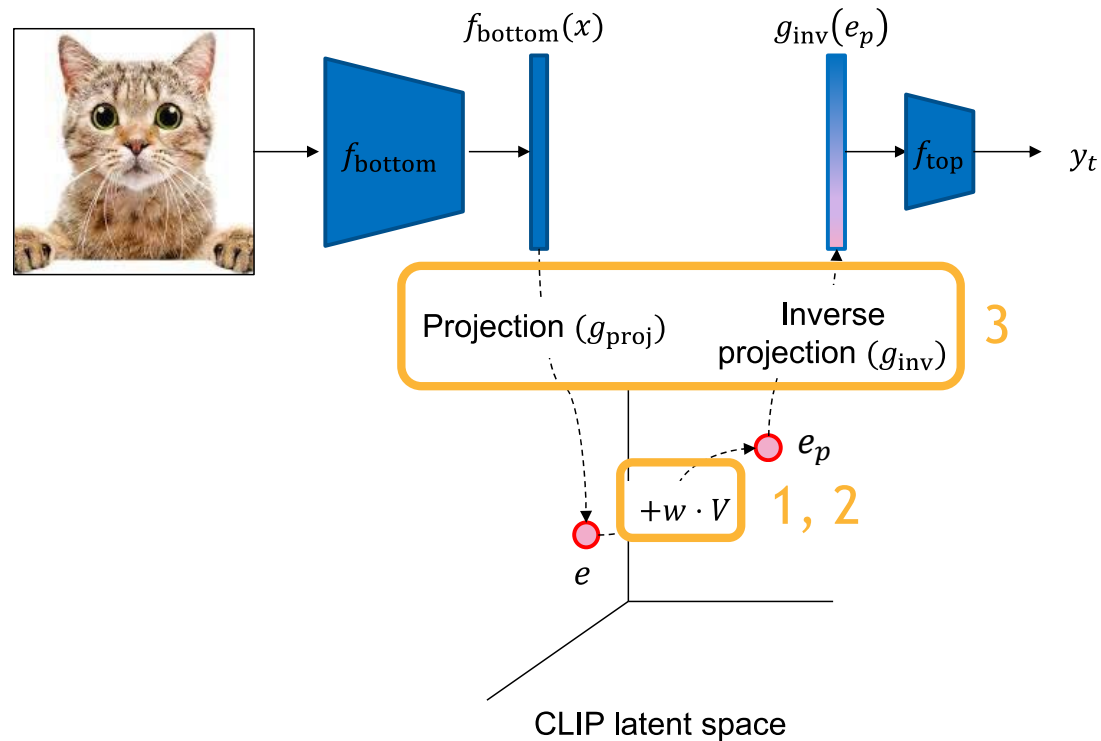
- Given a predefined concept library C ,
- Generate prompt pair for concept c
 - $[t_{src}^c, t_{trg}^c] = [\text{"a photo of object"}, \text{"a photo of } c \text{ object"}]$
- $v_c = \text{MinMaxNormalize}(\text{CLIP}_{\text{text}}(t_{src}^c) - \text{CLIP}_{\text{text}}(t_{trg}^c))$
- $V = \{v_c | c \in C\}$, where C is a predefined concept set

Category	Prompt template
Color	"A photo of {} object "
Texture	"A photo of {} object "
Scene	"A photo of object on {}"
Material	"A photo of object made of {}"
Part	"A photo of object containing {}"
Object	"A photo of object along with {}"



Research Questions

1. How can we obtain concept direction bank V ?
2. How can we obtain weight w ?
3. How can we implement projection and inverse projection?



2. Optimize w until the prediction changes to y_t

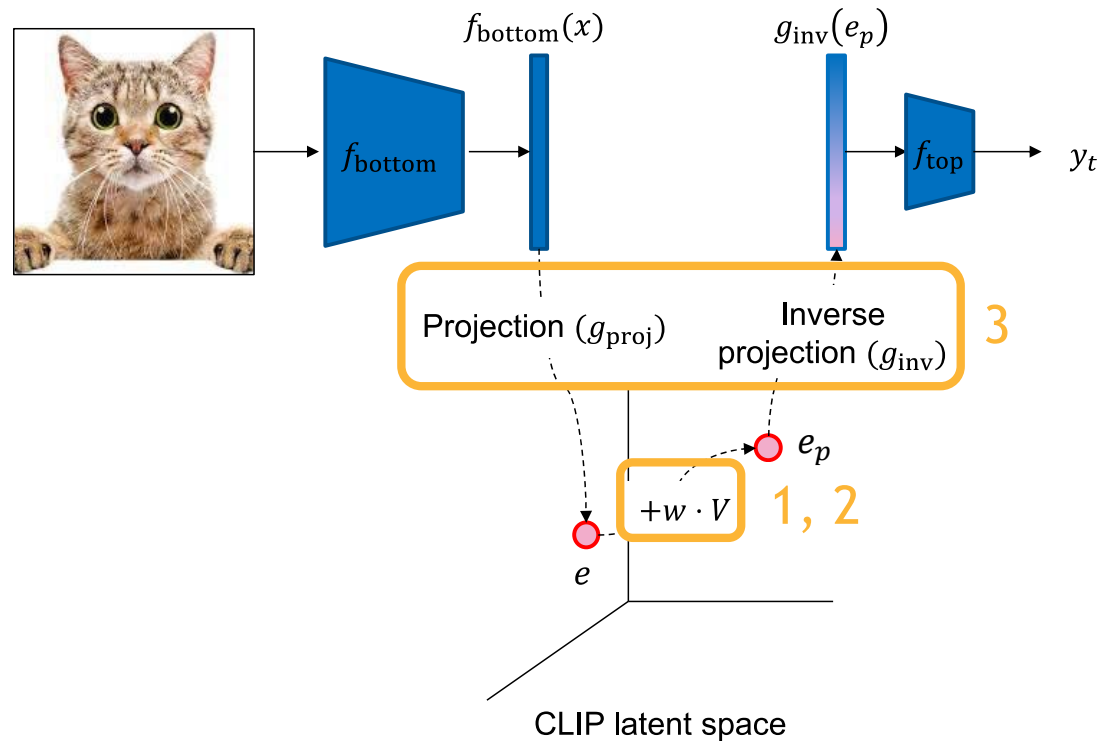
2. Optimize w with the objective function below

$$\min_w \mathcal{L} = \mathcal{L}_{\text{CE}} + \alpha \cdot \mathcal{L}_{\text{reg}} + \beta \cdot \mathcal{L}_{\text{id}}$$

- $\mathcal{L}_{\text{CE}} = \text{CrossEntropy}(f_{\text{top}}(g_{\text{inv}}(e_p)), y_t)$
 - Change the prediction to the target class by minimizing the cross-entropy loss
- \mathcal{L}_{reg}
 - Elastic net regularization: regularize concept importance to be 1) sparse and 2) unique minimum
- $\mathcal{L}_{\text{id}} = \|e_p - e\|^2$
 - Counterfactual approach requires the minimal modification
 - Identity loss to constrain the minimal perturbation

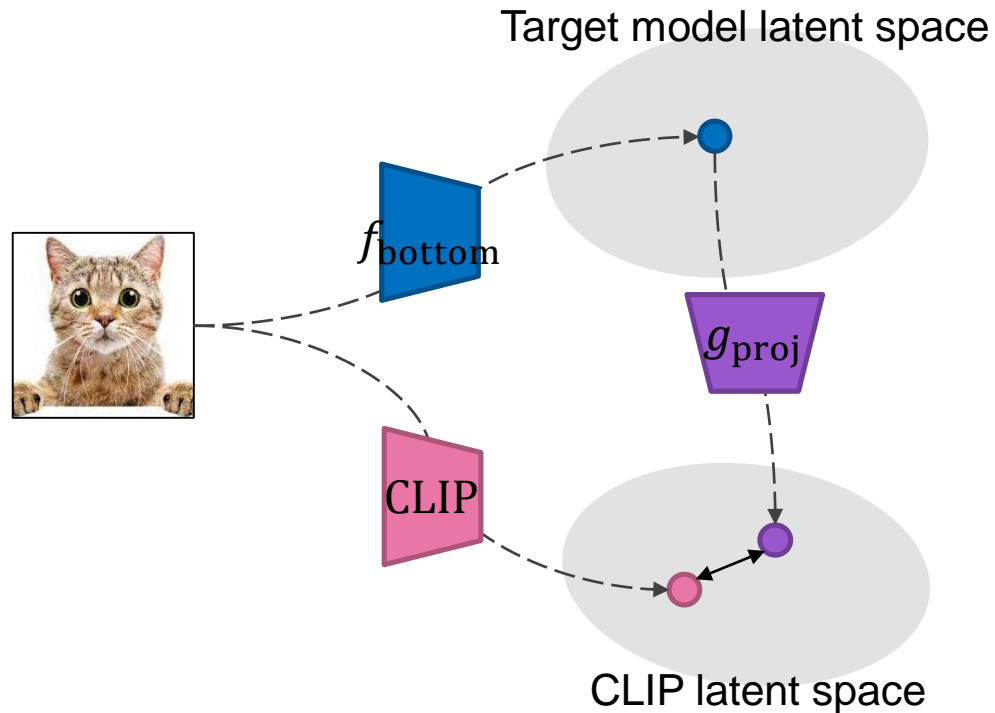
Research Questions

1. How can we obtain concept direction bank V ?
2. How can we obtain weight w ?
3. How can we implement projection and inverse projection?



3. Prepare projector and inverse projector

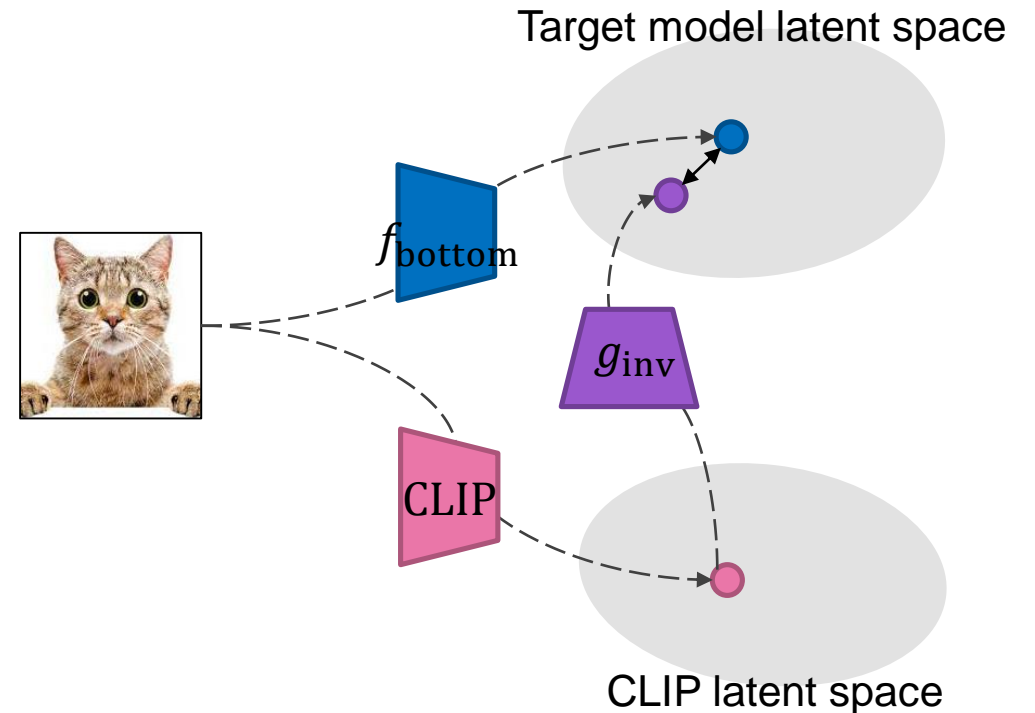
Projection



$$\mathcal{L}_{\text{proj}} = \left\| g_{\text{proj}}(f_{\text{bottom}}(x)) - \text{CLIP}(x) \right\|^2$$

source:

Inverse projection



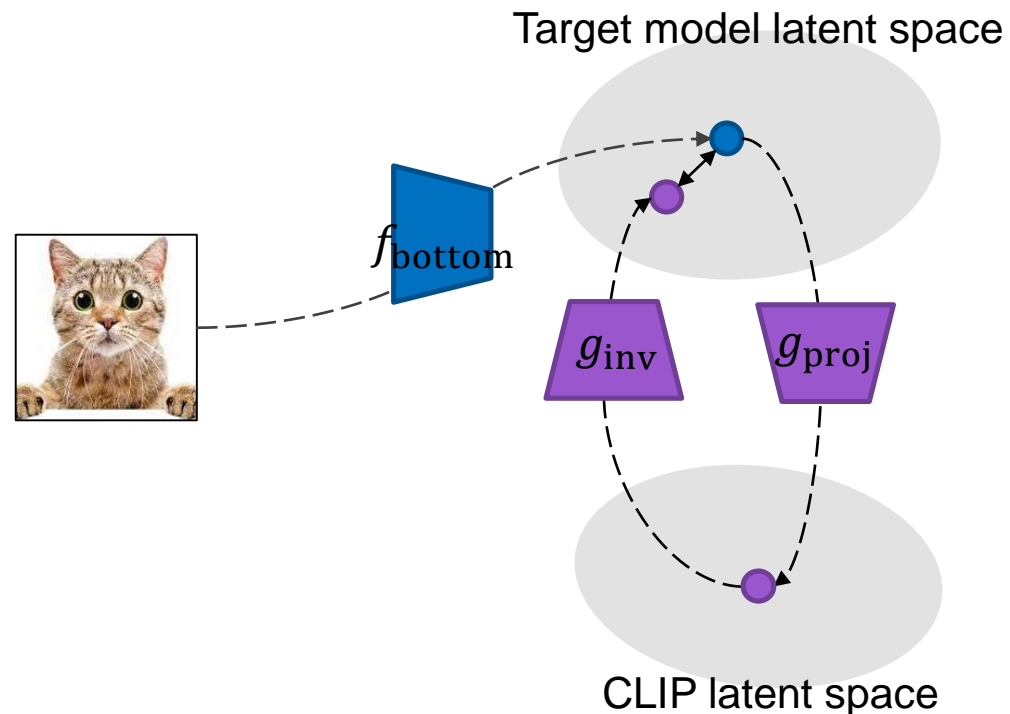
$$\mathcal{L}_{\text{inv}} = \left\| g_{\text{inv}}(\text{CLIP}(x)) - f_{\text{bottom}}(x) \right\|^2$$

3. Prepare projector and inverse projector

- Additionally finetune with cycle consistency loss

- $\mathcal{L}_{\text{finetune}} = \mathcal{L}_{\text{proj}} + \mathcal{L}_{\text{inv}} + \mathcal{L}_{\text{cycle}}$

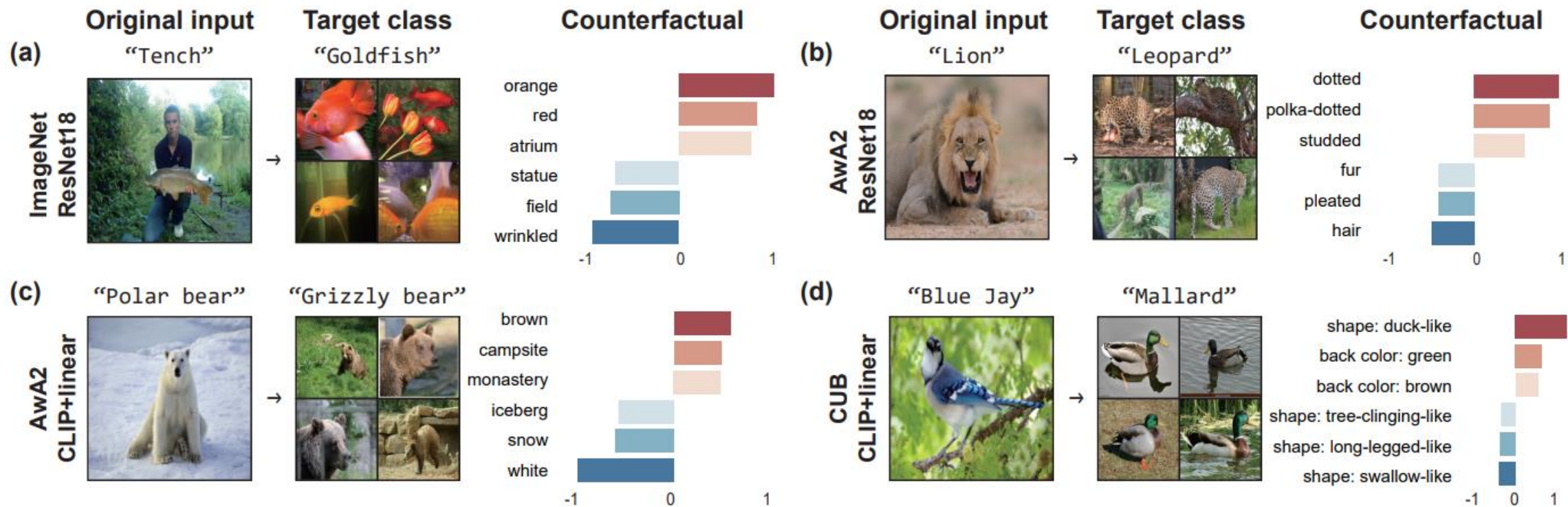
- $\mathcal{L}_{\text{cycle}} = \left\| f_{\text{bottom}}(x) - g_{\text{inv}}(g_{\text{proj}}(f_{\text{bottom}}(x))) \right\|^2$



Experimental setup

- Target model
 - CLIP+linear models
 - Trained on ImageNet/Animal with Attributes2 (AWA2)/Caltech-UCSD Birds-200-2011 (CUB)
 - Shares the embedding space with CLIP → Projection is not needed
 - ResNet18 models
 - Trained on ImageNet/AWA2/CUB
 - Does not share the embedding space with CLIP → Projection is needed
- Concept library \mathcal{C}
 - Reduce BRODEN for ImageNet-trained models
 - AWA2, CUB for AWA2-trained and CUB-trained models, respectively

Qualitative evaluation



Qualitative Comparison between Ours and CCE

- Spurious correlation in dataset collection of CCE led to inaccurate interpretation

Original Input Target class

“Granny smith”



“Strawberry”

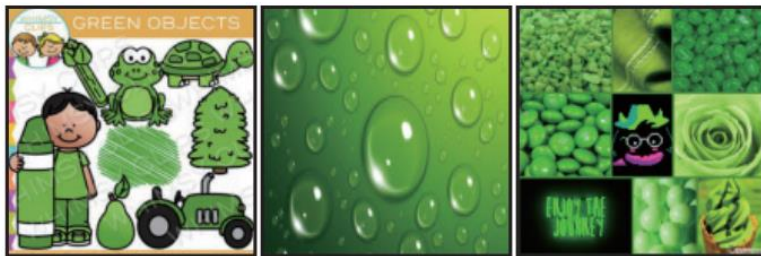
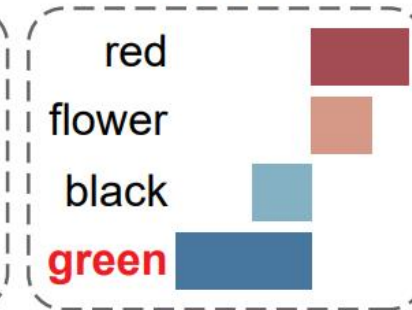


Counterfactual

CCE



Ours



Positive images for “Green”

vs.



Positive images for “Grass”

Debugging misclassification cases

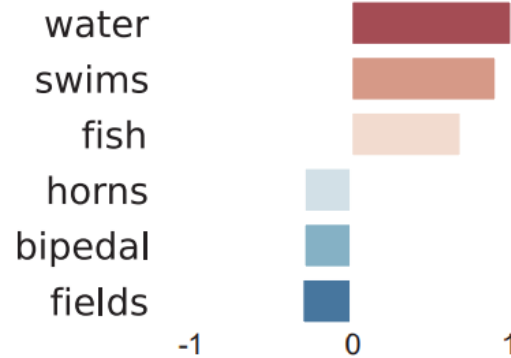
(a) Misclassified as
"Rhinosceros"



Original image

→

Correct answer
"Hippopotamus"



Counterfactual

(b) Background-
edited image



Prediction
"Hippopotamus"

(c) Examples of "Rhinosceros"



Examples of "Hippopotamus"



vs.

Quantitative Evaluation

- Lack of ground truth explanation
 - Especially for conceptual explanation
 - Some previous works often skip quantitative evaluation
- Repurpose existing datasets with class-wise attributes
 1. AWA2 dataset
 - 85 binary attributes are provided for 50 animal classes
 2. CUB dataset
 - 312 continuous attributes are provided for 200 bird classes

Examples of Ground Truth Attributes

AWA2 attribute labels



Class-wise binary attributes

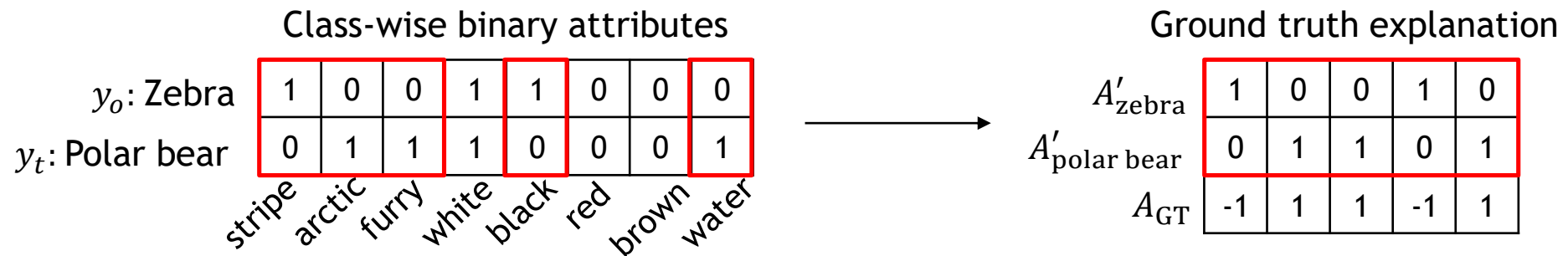
Zebra	1	0	0	1	1	0	0	0
Polar bear	0	1	1	1	0	0	0	1
	stripe	arctic	furry	white	black	red	brown	water

CUB attribute labels



Quantitative Evaluation Protocol

- Since our method adopts *counterfactual* approach,
 - A_{y_o} : Attributes of the original class, A_{y_t} : Attributes of the target class
 - Only the difference between two attributes is our interest
 - Consider only if the two elements of the attributes are different
 - $A_{GT} = [A_{y_t}[i] - A_{y_o}[i]$ for i in range $(\text{len}(A_{y_o}))$ if $A_{y_o}[i] \neq A_{y_t}[i]$



- Report $\text{AUROC}(A_{GT}, I)$

Quantitative Comparison w/ Baseline-CCE

- 26% improved AUROC compared to Baseline-CCE

Target model	Dataset	Library	CCE	Ours
CLIP+Linear	AwA2	C_{AwA2}	0.6436	0.8132
	CUB	C_{CUB}	0.7066	0.7891
ResNet18	AwA2	C_{AwA2}	0.6113	0.7314
	CUB	C_{CUB}	0.6979	0.7750
ResNet50	AwA2	C_{AwA2}	0.5811	0.7316
	CUB	C_{CUB}	0.6811	0.7336

Thank You!