



DeepREAL



**ICML**  
International Conference  
On Machine Learning

# Inside the Visual Mind: Neuroscience-Motivated Concept Circuits for Interpreting and Steering Vision Transformers

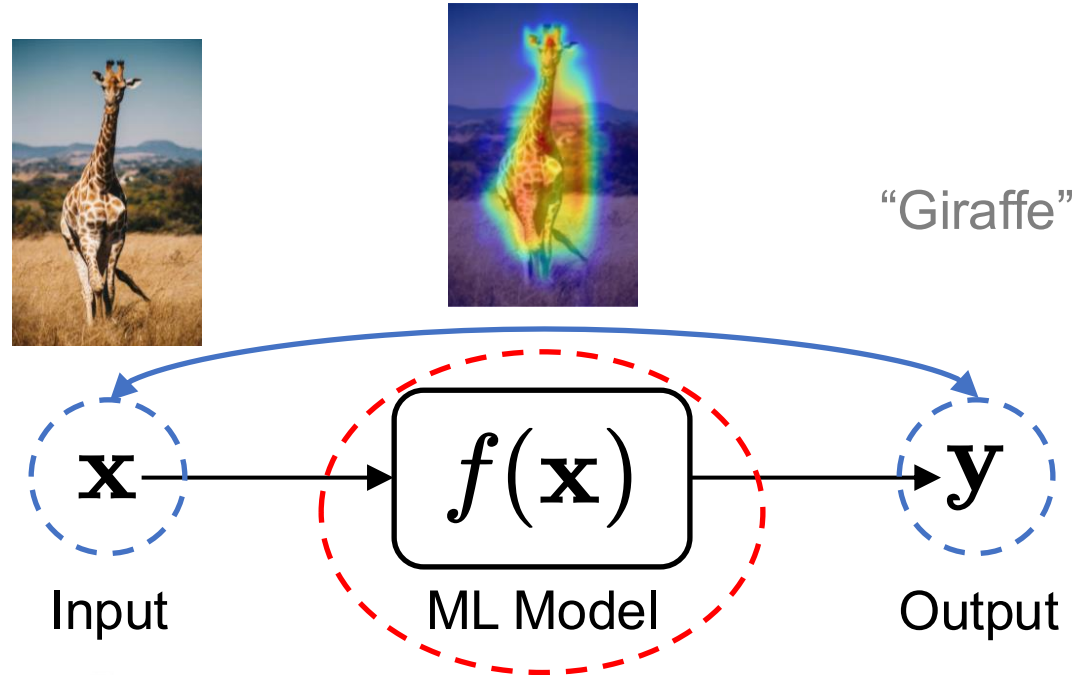
Tang Li   Yanlin Chen   Mengmeng Ma   Xi Peng

Deep Robust & Explainable AI Lab (DeepREAL)

Tue, Jul 7, 2026   10:30 AM – 12:15 PM KST   Coex: HALL A

# From “correlation” to “cognition”

“What” lead to the prediction



but “How” does it think inside?

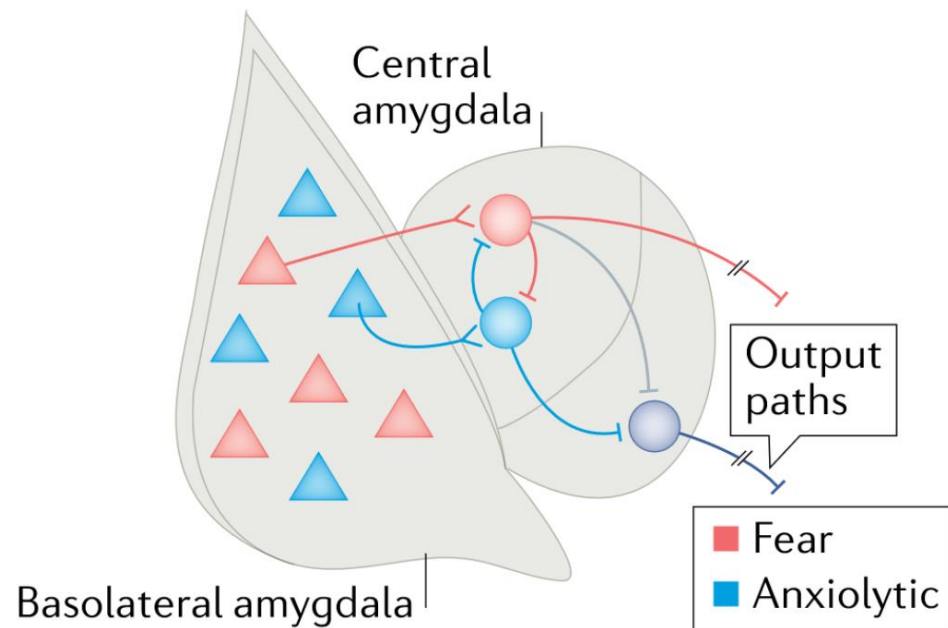
In other words, can we “read the mind” of AI?

# Two views to understand biological brains: Cognitive Science of Human

## Sherringtonian View

- Node-to-node connections
- Transfer between neuron computations

### Amygdala circuits for fear and anxiety

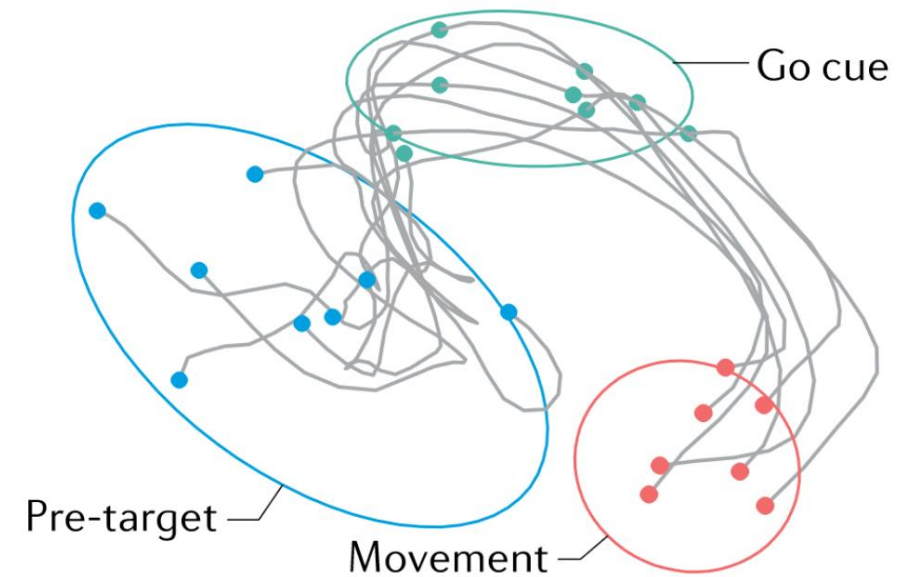


[Nature'15, Janak et al.]

## Hopfieldian View

- Representational spaces
- Transformations between spaces

### Trajectories in neural space in PMd during a reaching task in monkeys



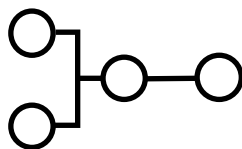
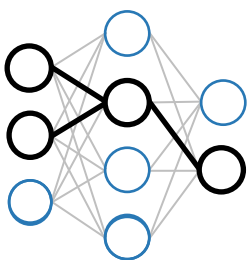
[Neurosci'10, Churchland et al.]

# Two views to understand artificial brains (ML models):

## Sherringtonian View

### Bottom-up View

- Node-to-node connections
- Identify functional circuits



Sub-computation  
Graph

Circuit

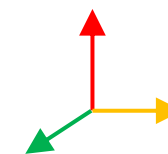
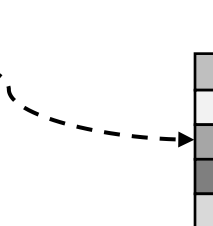
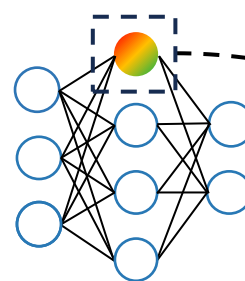
- *Activation patching* [Conmy'23]
- *Path patching* [Goldowsky-Dill'23]
- *Causal abstraction* [Geiger'23]

⊖ Local Semantics    ⊕ Global Structure

## Hopfieldian View

### Top-down View

- Representation space transforms
- Disentangle features from representations



Polysemantic  
Neuron

Representations

Monosemantic  
features

- *Probing* [Burns'23]
- *Representation Engineering* [Zou'23]
- *Sparse dictionary learning* [Cunningham'24]

⊕ Local Semantics    ⊖ Global Structure

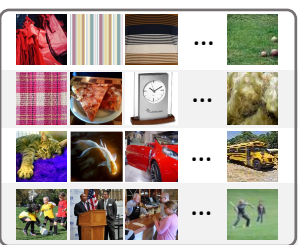


Can we bridge Bottom-up and Top-down Interpretations to develop the “Cognitive Science of AI”?

# Our Solution: ViSAE Toolbox for Interpreting ViT Inner Workings

## Neuroscience-motivated Probing Suite

### Probing Image Set

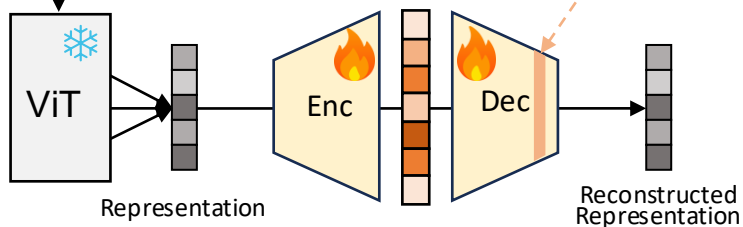


- Primitive-level
- Intermediate-level
- Object-level
- Scene-level

### Concept Set

"red"	"curve"
"edges"	"circle"
"wood"	"stripe"
"cat"	"wings"
"grass"	"jump"
"inside"	"look"

### Sparse Autoencoder (SAE)

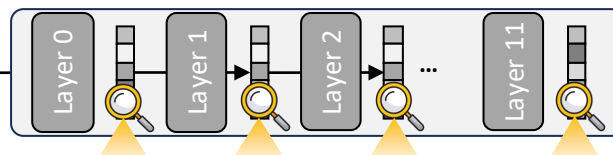


## Concept Circuits Tracing Algorithm

### Input Image

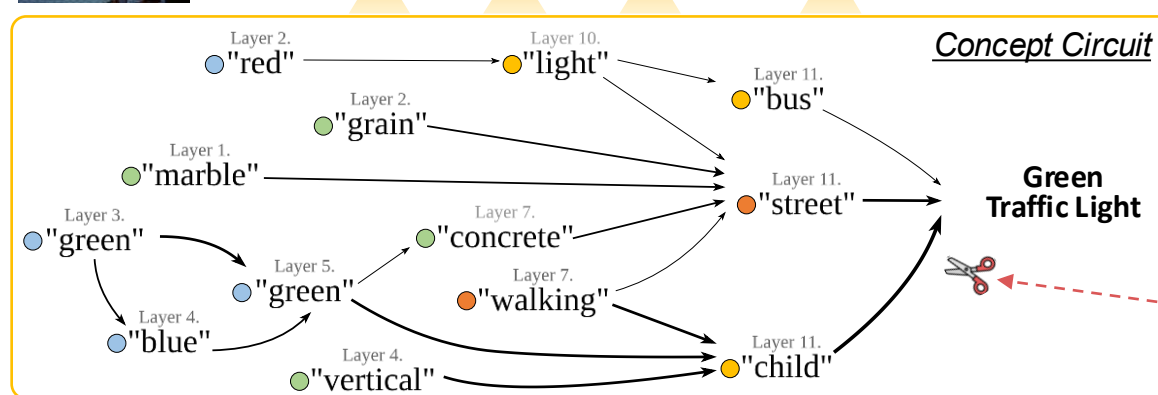


### ViT



### Prediction

**Green Traffic Light**



## Applications

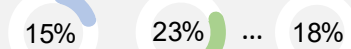
### Auditing

#### Concept Localization



#### Failure Mode Analysis

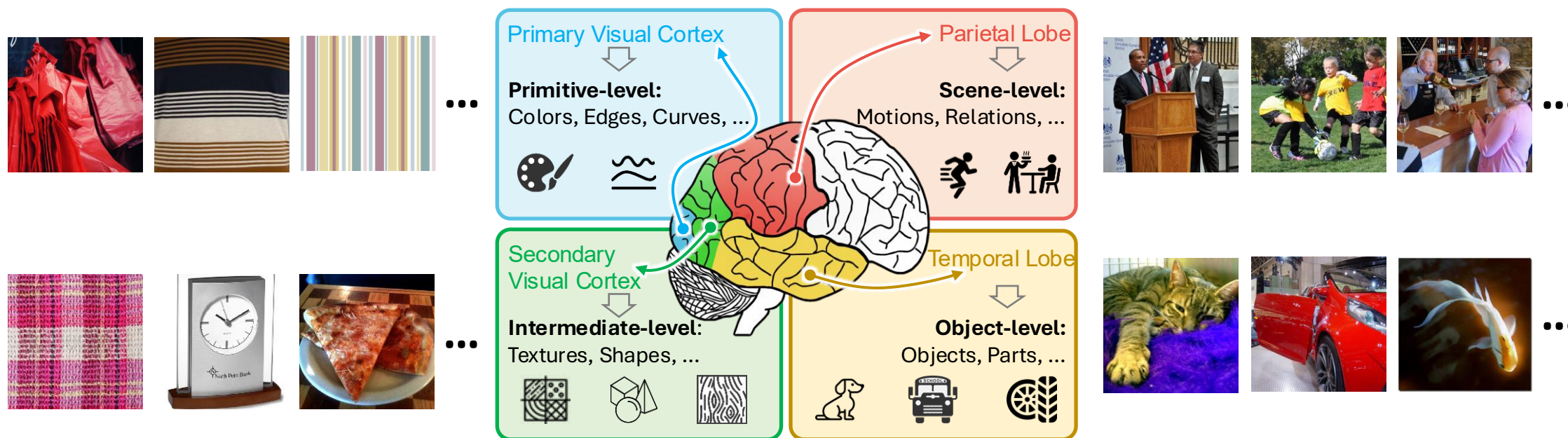
Part-whole confusion    Color bias    Background shortcut



### Steering

Suppress Spurious Concepts → Prediction: **Red Traffic Light** ✓

# Neuroscience-Motivated Probing Suite

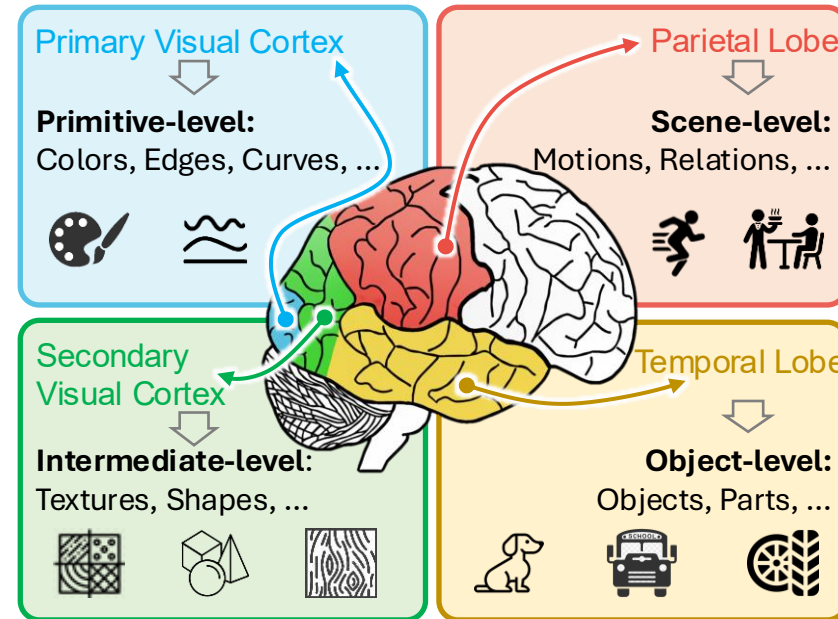


Probing Image Set	Data Source	# of Images	Concepts Covered by Images (%)					Coverage Efficiency ↑ (%/1K Images)
			Primitive	Intermediate	Object	Scene	Avg.	
ImageNet	ImageNet	1,281K	81.0	78.2	97.7	59.0	78.9	0.06
MSCOCO	MSCOCO	118K	69.6	65.4	80.4	63.1	69.6	0.59
Ours	Primitive Level: DTD, Broden; Intermediate Level: Broden, ShapeNet; Object Level: ImageNet, VisualGenome; Scene Level: Places365, MSCOCO;	64K	87.1	80.6	92.6	61.7	80.5	1.26

# Neuroscience Motivated Probing Suite

*“red”, “blue”, “lime”, “lines”,  
“edges”, “dots”, ...*

*“wooden”, “striped”,  
“triangle”, “steel”, “round”, ...*



*“playing basketball”, “leap”,  
“above”, ...*

*“dog”, “airplane”, “chair”,  
bird, ...*

### Concept Statistics

Level	# of Concepts
Primitive	1,073
Intermediate	1,723
Object	10,534
Scene	2,720
<b>Total</b>	<b>16,050</b>

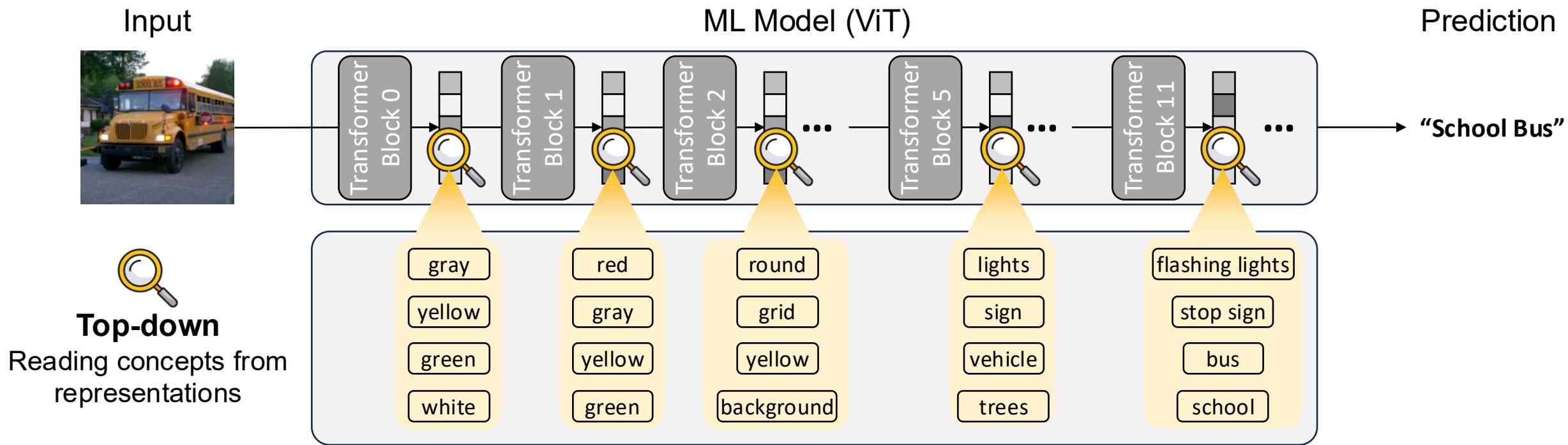
### Human Evaluations

Evaluator	Faithfulness	Completeness
Human_A	4.65 ± 0.06	4.72 ± 0.04
Human_B	4.83 ± 0.20	4.73 ± 0.06
Human_C	4.90 ± 0.05	4.78 ± 0.18
<b>Avg.</b>	<b>4.79</b>	<b>4.74</b>

### Comparison of Concept Quality

Concept Set	# of Concepts	Redundancy (↓)	Visually Grounded (↑)
LAION-freq	15K	0.419	0.602
Google-freq	20K	0.654	0.432
LaBo	10K	0.584	0.565
<b>Ours</b>	<b>16K</b>	<b>0.213</b>	<b>0.864</b>

# Concept Circuit Tracing Algorithm



# Concept Circuit Tracing Algorithm

Input

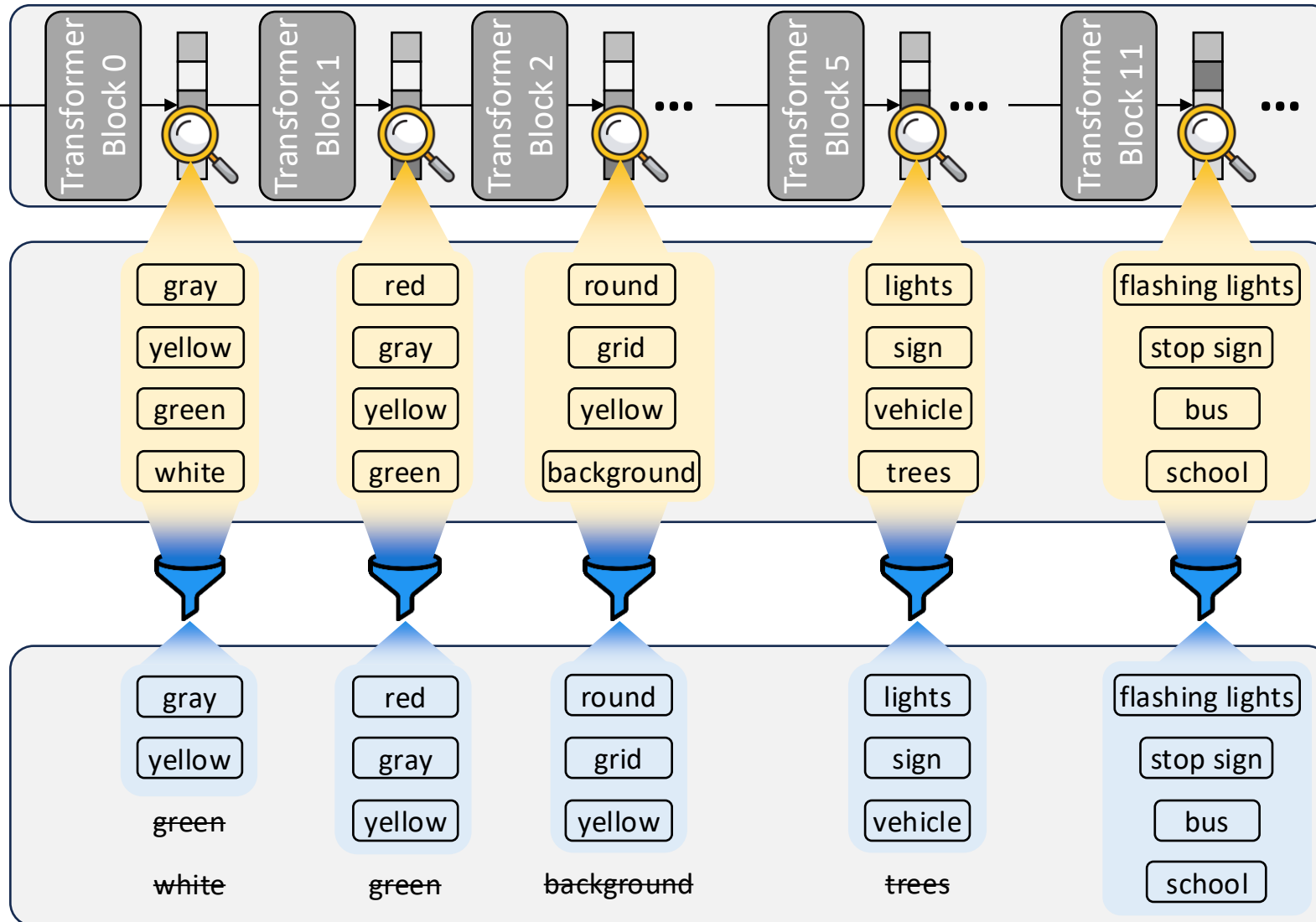


ML Model (ViT)

Prediction

"School Bus"

  
**Top-down**  
 Reading concepts from representations



# Concept Circuit Tracing Algorithm

Input

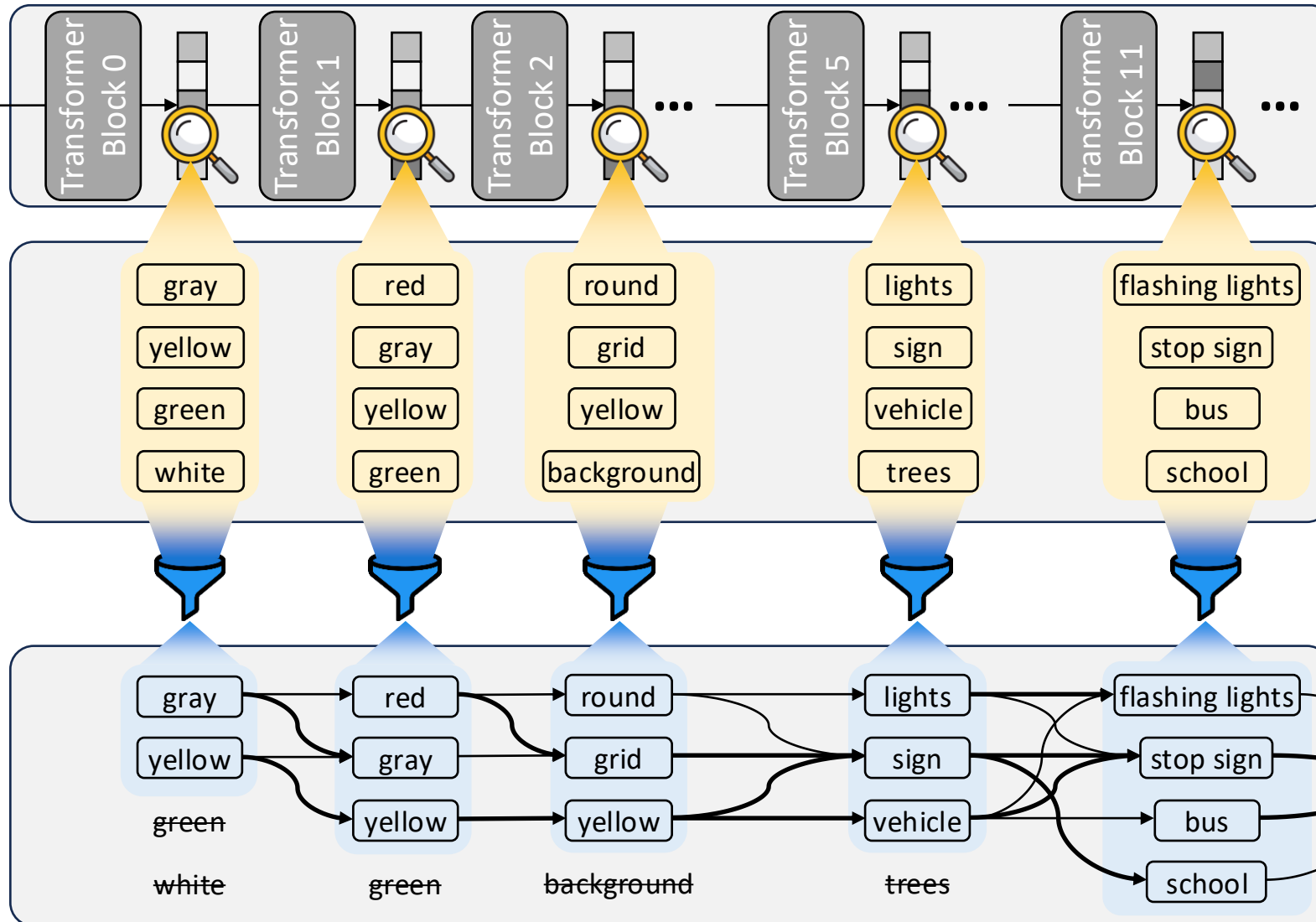


ML Model (ViT)

Prediction

"School Bus"

  
**Top-down**  
 Reading concepts from representations

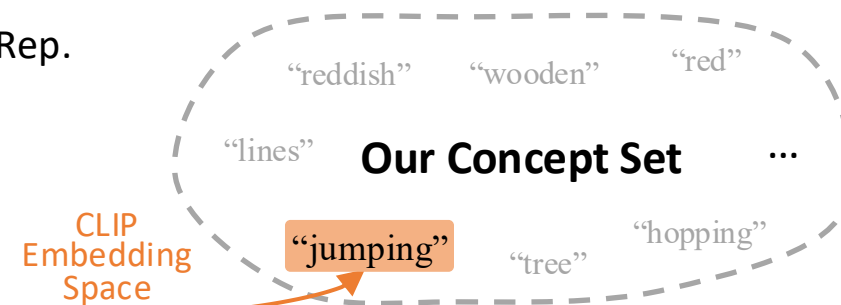
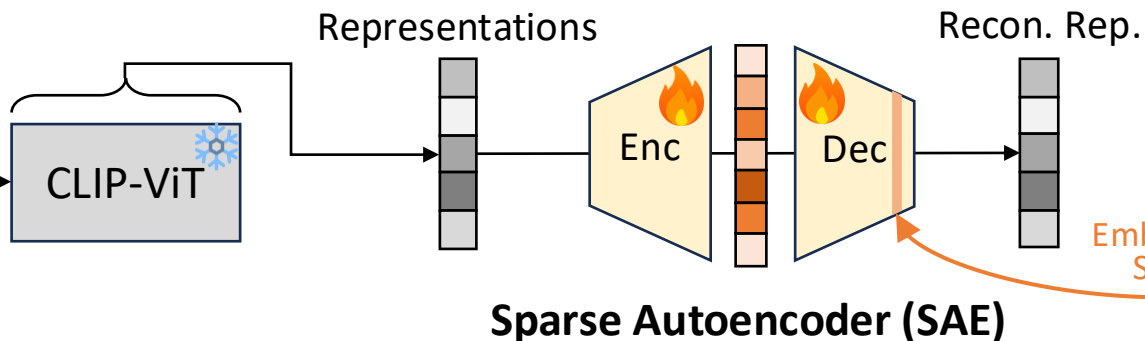


# Top-down Concept Reading

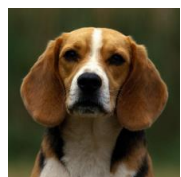
➤ **Step 1:** Train SAE with the probing image set.

➤ **Step 2:** Map SAE features to concepts.

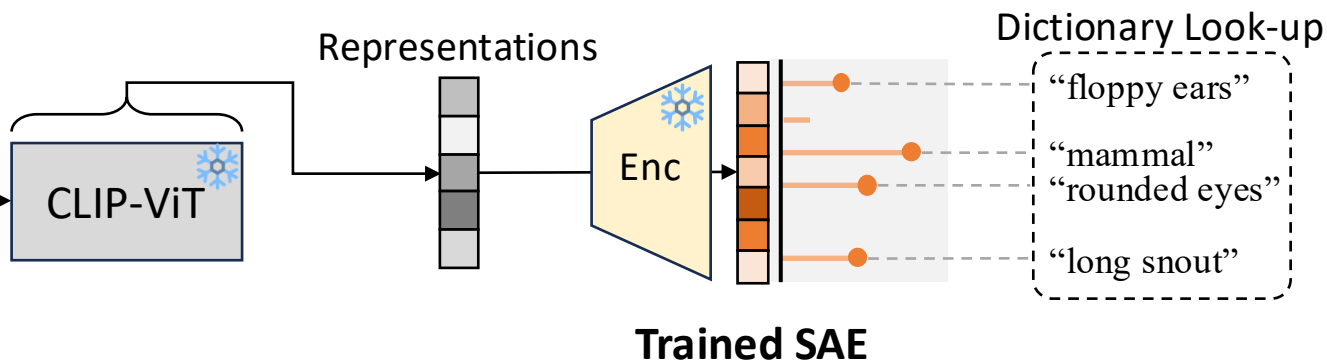
Our Probing Images



➤ **Step 3:** Read concepts from representations.



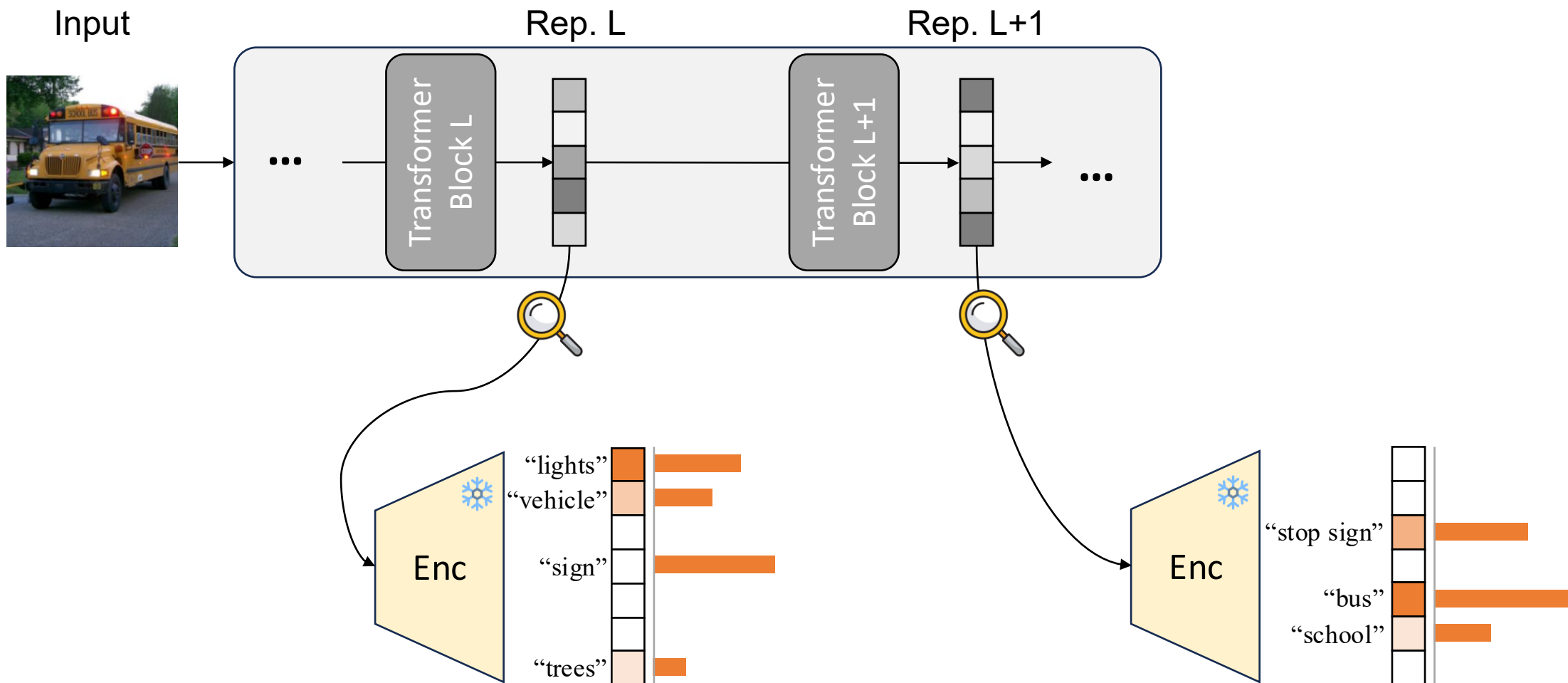
Testing Image





# Bottom-up Circuit Tracing

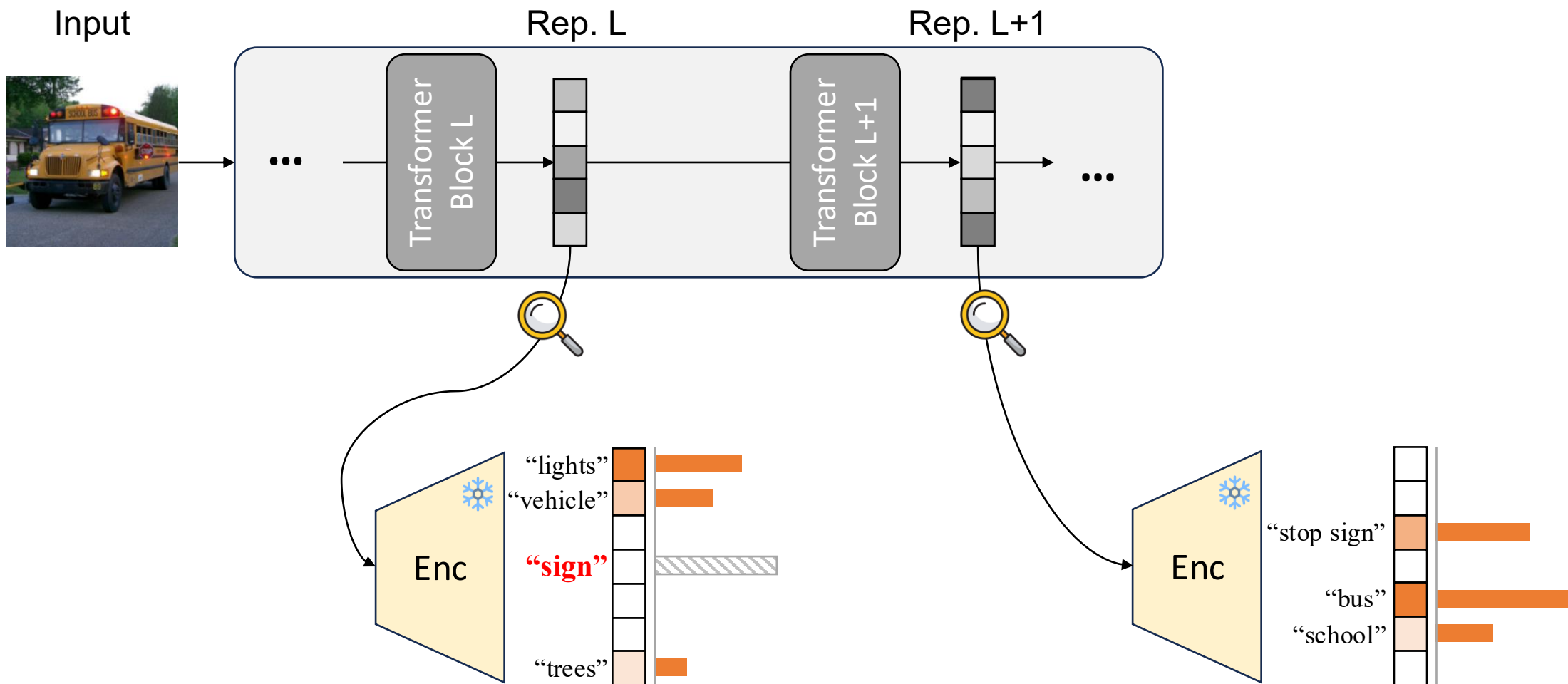
$$\begin{aligned} \text{IE}_{i \rightarrow j}^{s \rightarrow t} &= \text{IE}(\alpha_j^t; c_i^s; r_{\text{clean}}, r_{\text{patch}}) \\ &= \alpha_j^t(r_{\text{clean}}) - \alpha_j^t(r_{\text{clean}} \mid \text{do}(\alpha_i^s = \alpha_i^s(r_{\text{patch}}))). \end{aligned}$$





# Bottom-up Circuit Tracing

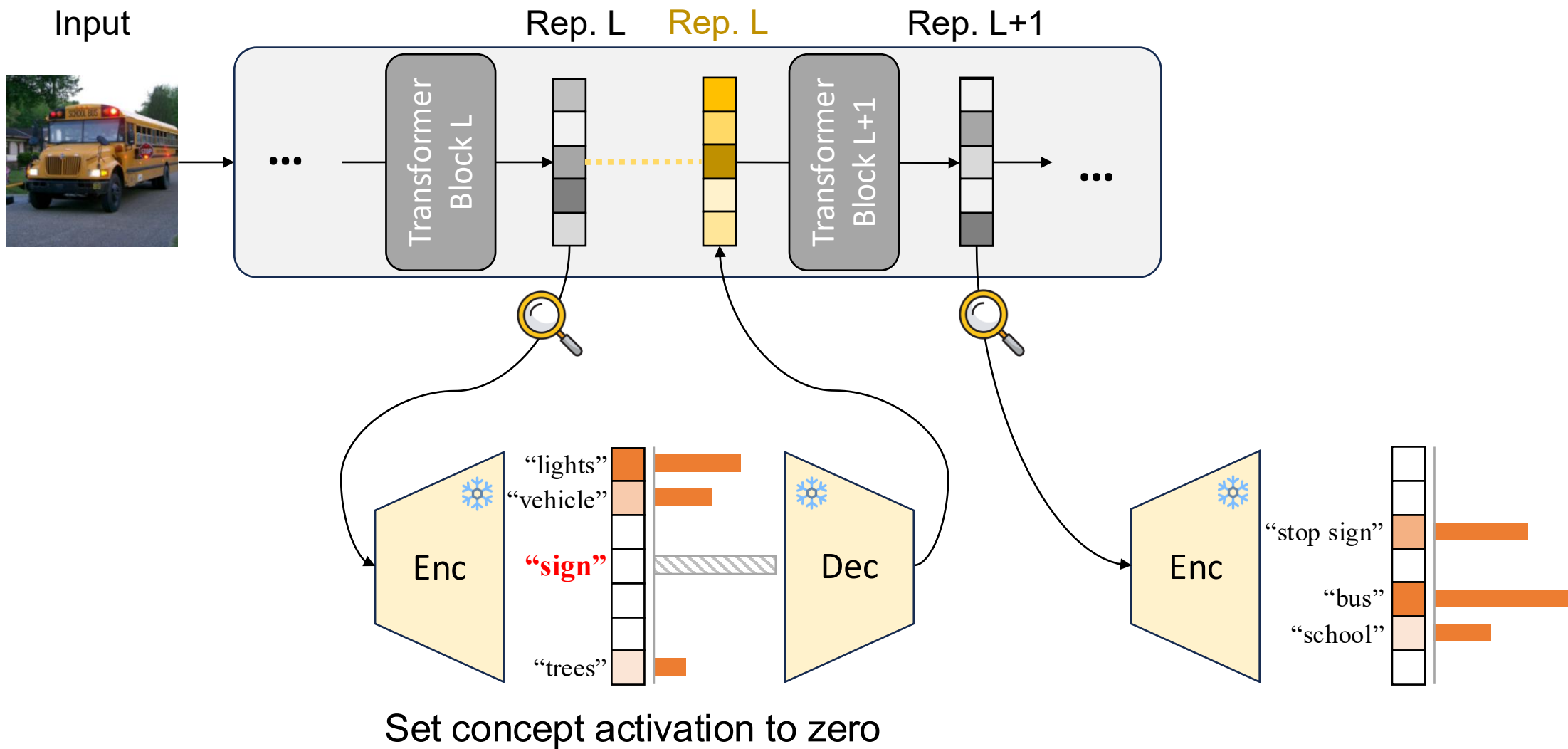
$$\begin{aligned} \text{IE}_{i \rightarrow j}^{s \rightarrow t} &= \text{IE}(\alpha_j^t; c_i^s; r_{\text{clean}}, r_{\text{patch}}) \\ &= \alpha_j^t(r_{\text{clean}}) - \alpha_j^t(r_{\text{clean}} \mid \text{do}(\alpha_i^s = \alpha_i^s(r_{\text{patch}}))). \end{aligned}$$





# Bottom-up Circuit Tracing

$$\begin{aligned} \text{IE}_{i \rightarrow j}^{s \rightarrow t} &= \text{IE}(\alpha_j^t; c_i^s; r_{\text{clean}}, r_{\text{patch}}) \\ &= \alpha_j^t(r_{\text{clean}}) - \alpha_j^t(r_{\text{clean}} \mid \text{do}(\alpha_i^s = \alpha_i^s(r_{\text{patch}}))). \end{aligned}$$

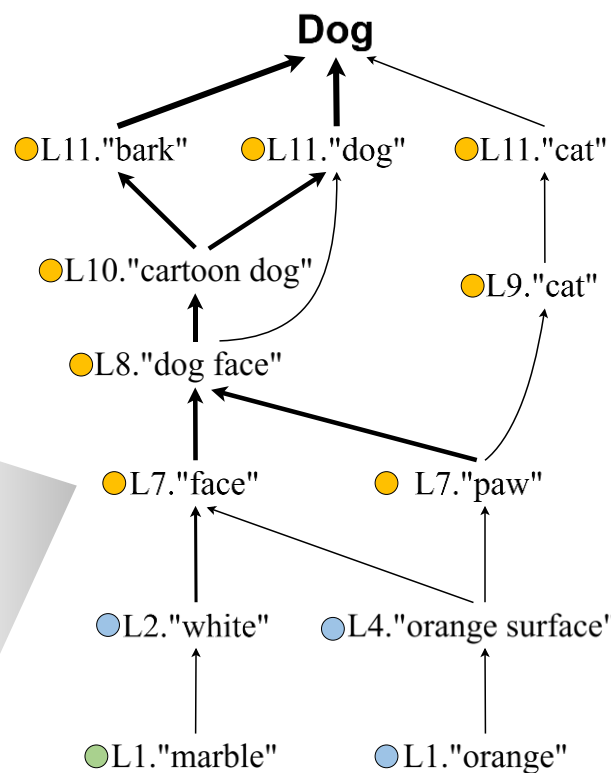
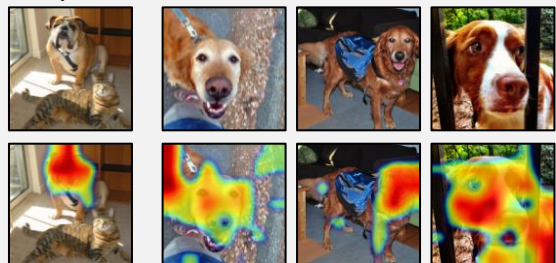




# Audit inner workings

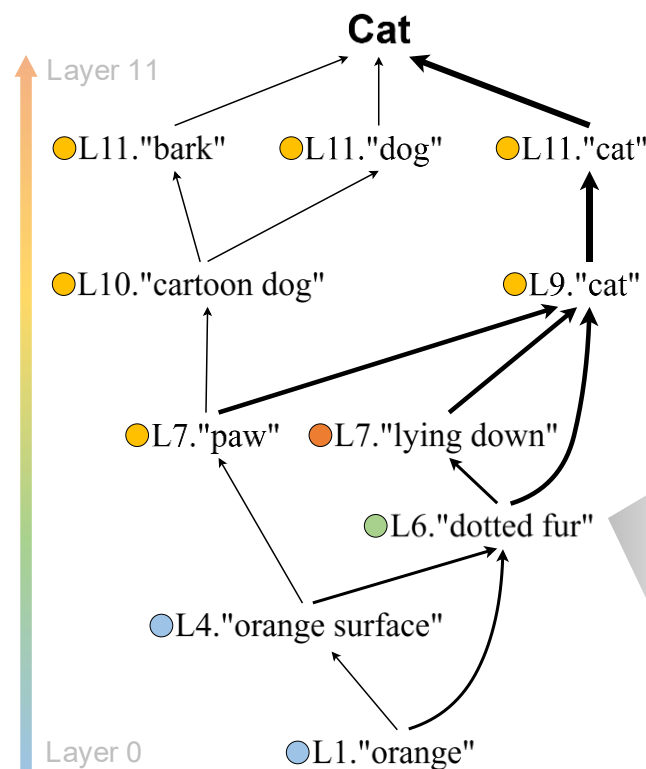
Layer 11, CLS token, #3167: **“bark”**

Input Top-3 Activated Images



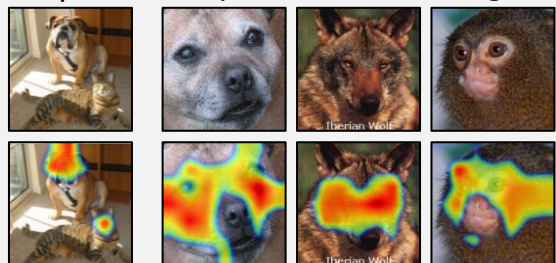
Layer 11, CLS token, #1354: **“cat”**

Input Top-3 Activated Images



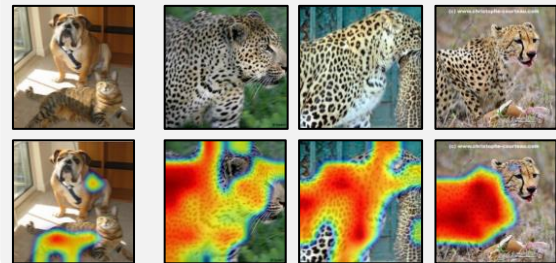
Layer 7, IMG token, #4114: **“face”**

Input Top-3 Activated Images



Layer 6, IMG token, #4375: **“dotted fur”**

Input Top-3 Activated Images



Our method captures the nuanced difference of model in making predictions.

# Locate concepts on pixels

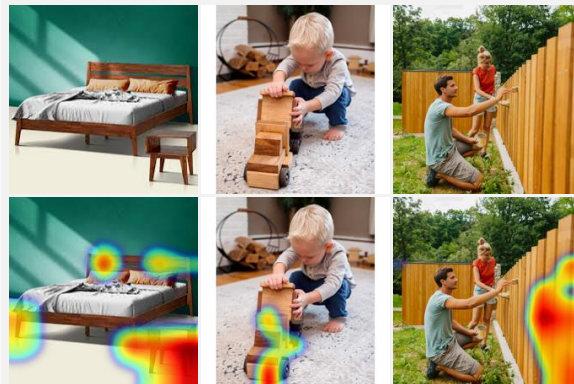
## Primitive-Level

Concept 2196 – “yellow”



## Intermediate-Level

Concept 1067 – “wooden”



## Object-Level

Concept 2924 – “wings”



## Scene-Level

Concept 1274 – “looking at”



Concept 1747 – “lines”



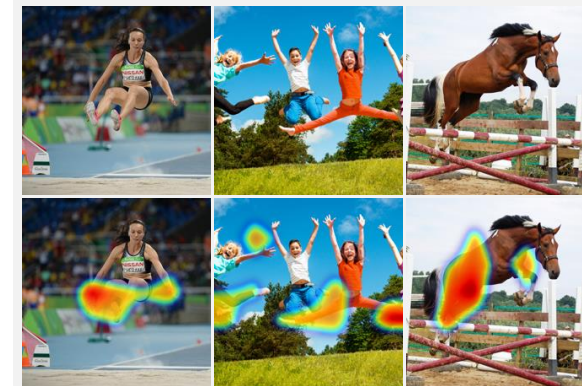
Concept 1875 – “circle”



Concept 139 – “cat”



Concept 2554 – “jumping”



Our method localize concepts on pixels, even for highly abstract ones.

# Analyze failure modes

## Failure Images

38 classes, 1.9K images

Label: canoe  
Pred: paddle



## Right vs. Wrong

Difference in SAE Concepts

["holding",  
"people",  
"mountain",  
"water",  
"shade", ...]

Label: parallel bar  
Pred: high bar (single)



["gymnast",  
"flying",  
"bar",  
"downward",  
"indoor", ...]

## Failure Modes

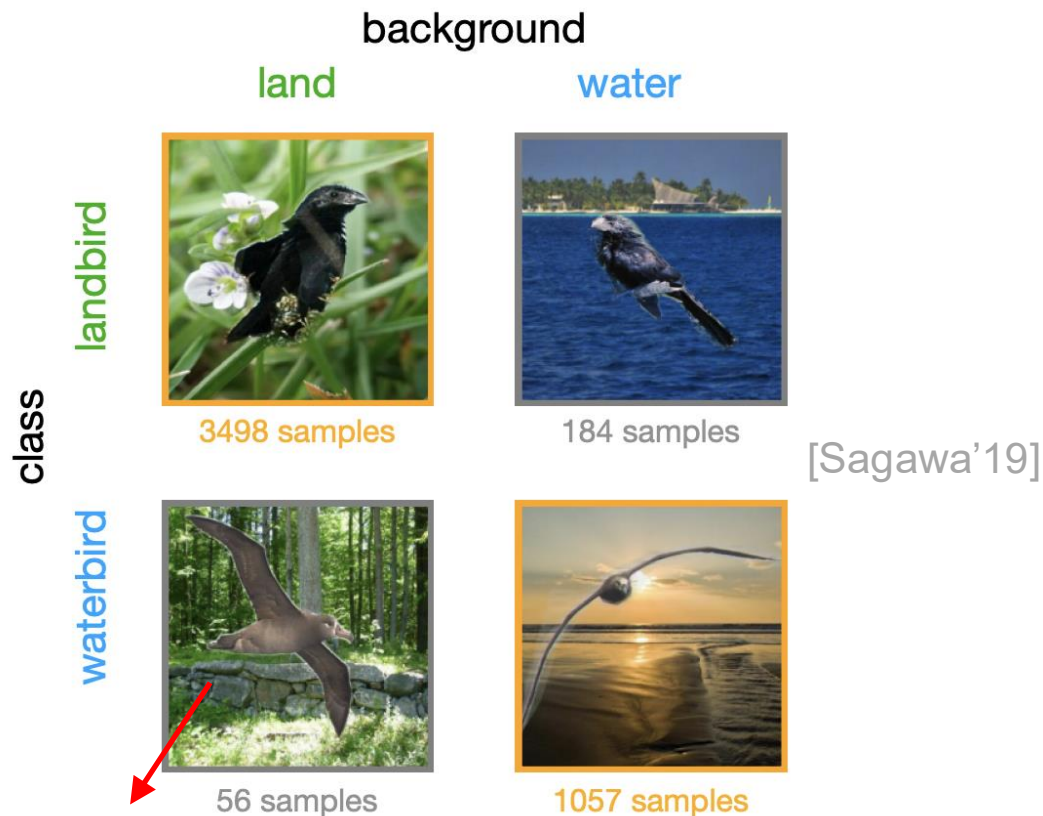
Statistics (%)

Our method can efficiently analyze the reasons behind failures.

# Steer model prediction by editing concepts

## WaterBird Dataset

*Bird categories spuriously correlated to backgrounds.*



**Worst Group:** waterbirds on land backgrounds.

Label: **waterbird**



Pred: **landbird**  
Confidence: **54%**

SAE Top-4 Concepts

“bird”: 7.09  
“water”: 3.90  
“tree”: 2.53  
“rock”: 2.46

**Remove**

Land Concepts



“bird”: 7.09  
“water”: 3.90  
“tree”: **0.00**  
“rock”: **0.00**

Pred: **waterbird**  
Confidence: **89%**

**Enhance**

Land Concepts



“bird”: 7.09  
“water”: 3.90  
“tree”: **5.00**  
“rock”: **5.00**

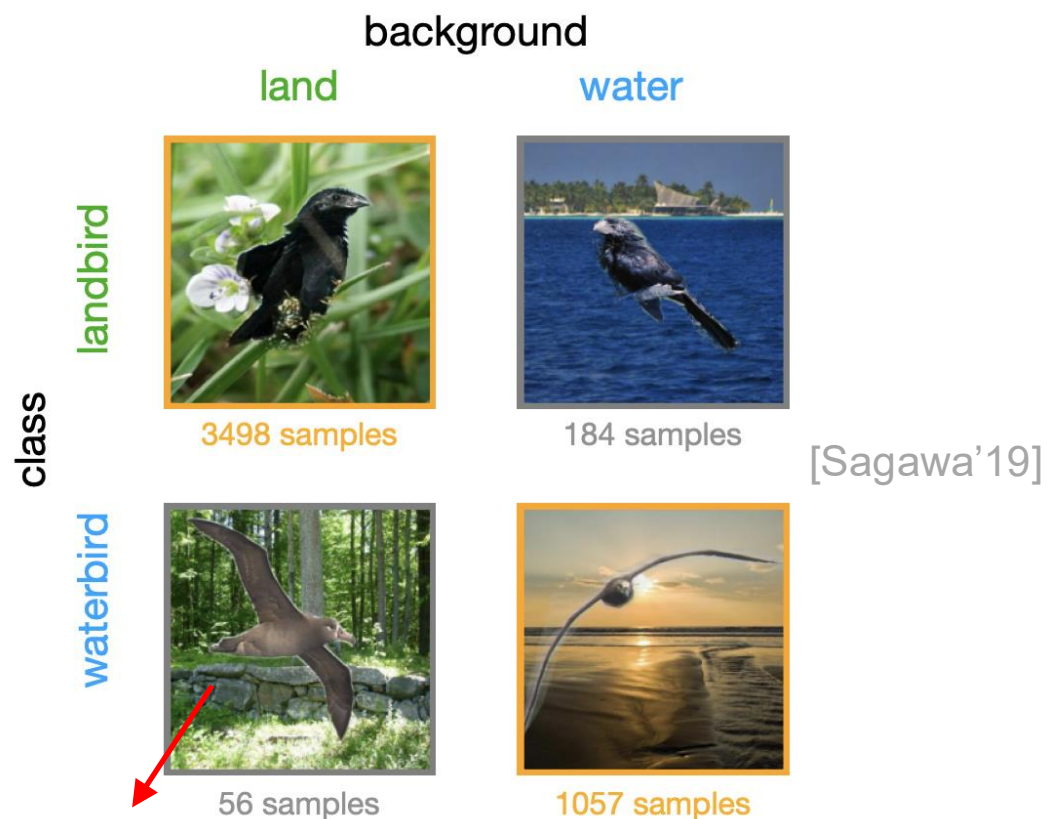
Pred: **landbird**  
Confidence: **94%**

Our SAE works as precise control “knobs” of concepts.

# Steer model prediction by editing concepts

## WaterBird Dataset

*Bird categories spuriously correlated to backgrounds.*



**Worst Group:** waterbirds on land backgrounds.

## Comparison of Steering Power

Method	Steer Spuri. Corr.	Overall Acc. (%)	Worst Group Acc. (%)	$\Delta$
CBM	None	-	37.3	-
	Remove	-	51.8	+ 14.5
SpLiCE	None	-	48.0	-
	Remove	-	60.0	+ 12.0
DN-CBM	None	-	57.5	-
	Remove	-	71.3	+ 13.8
PCBM	None	-	50.3	-
	Remove	-	74.7	+ 24.4
Ours	None	79.7	50.3	-
	Enhance	74.5	5.3	- 45.0
	Remove	85.2	98.5	+ 48.2

Our SAE works as precise control “knobs” of concepts.



DeepREAL



**ICML**  
International Conference  
On Machine Learning

# Thank You!

**Inside the Visual Mind: Neuroscience-Motivated Concept Circuits for Interpreting and Steering Vision Transformers**

Tang Li   Yanlin Chen   Mengmeng Ma   Xi Peng

Deep Robust & Explainable AI Lab (DeepREAL)

Tue, Jul 7, 2026   10:30 AM – 12:15 PM KST   Coex: HALL A