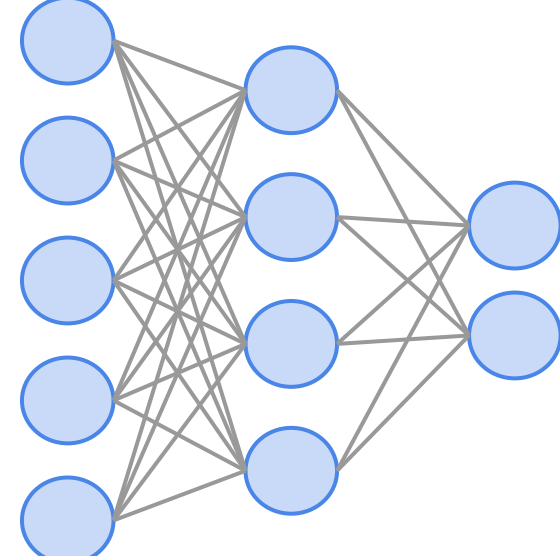
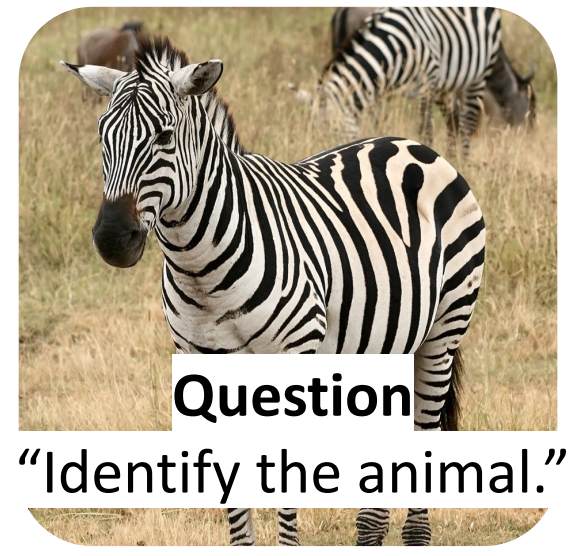




Introduction

Why did my Neural Network make that decision?



Response
"Zebra"

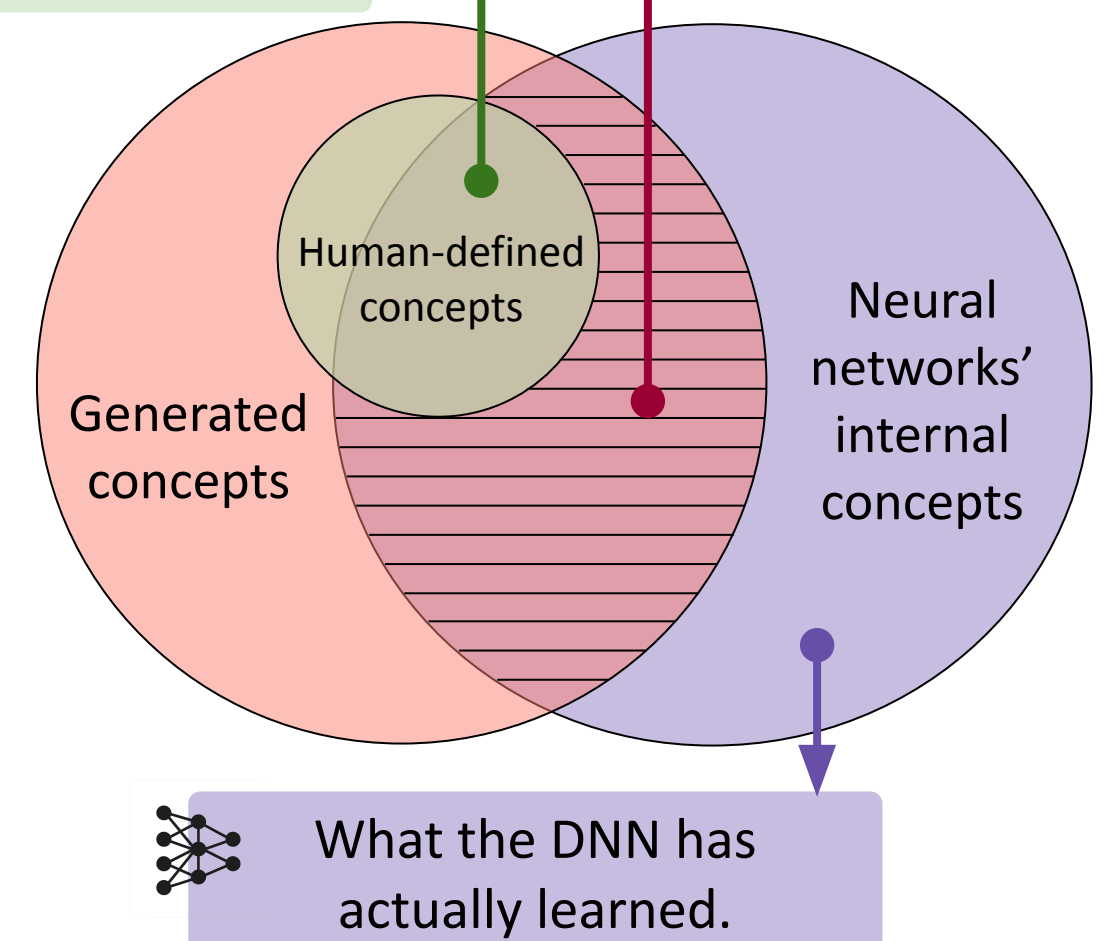
Why?

- Because there is grass.
- Because of black and white stripes.
- Random guess.



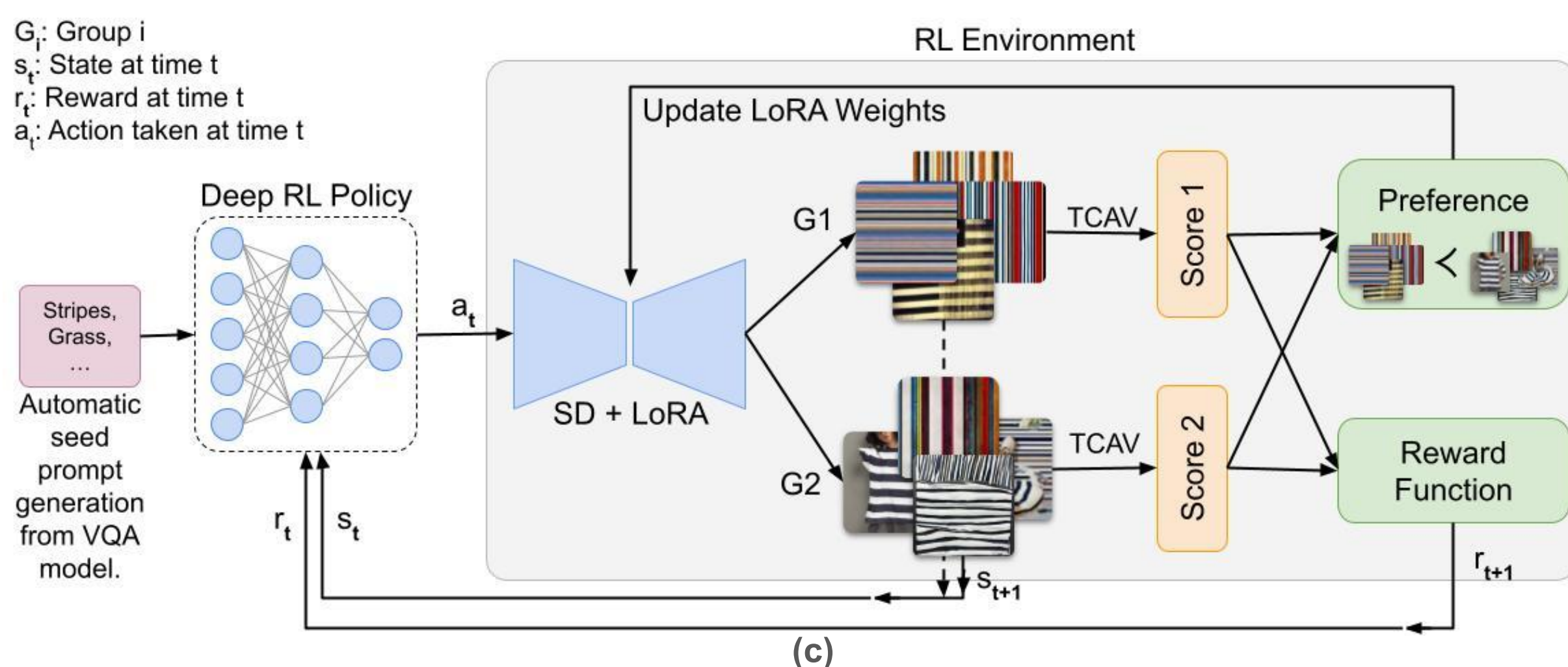
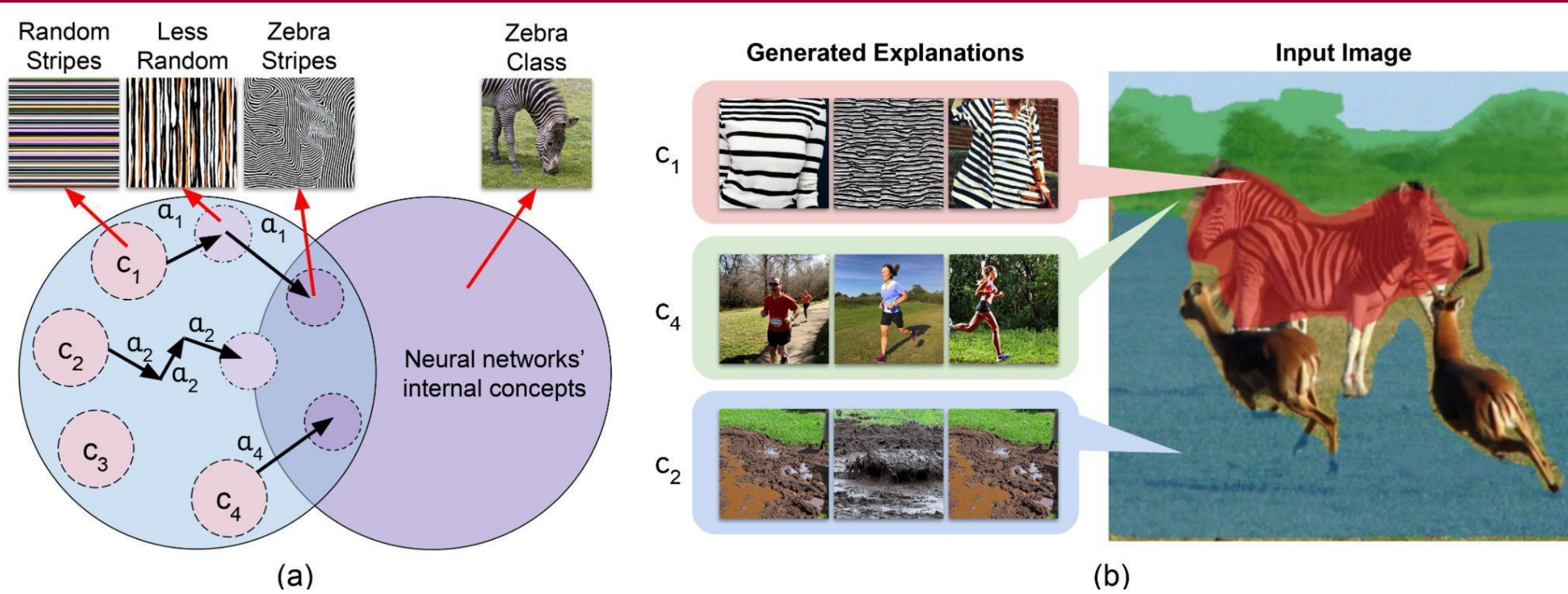
What human can guess it has learned.

What a generative model can explain.



1. There could be infinite explanations when we ask why a DNN behave in a particular way – testing each individually [1] is challenging.
2. We propose a method, named RLPO, to **generate explanations**—represented as visual or natural language concepts—that truly matter to the neural network's decision by analyzing activation vectors.

Methodology



Seed Prompts: We start by questioning class images (using a VQA model) about their colors, textures, etc. After filtering, a set of top-scoring keywords is chosen as seed prompts.

Concept Generation: Using these prompts, we produce candidate concept sets.

RL-based Preference Optimization (RLPO): The RL agent selects seed prompts, guided by TCAV scores [1] as rewards.

$$TS_{c,m} = \frac{1}{|X_m|} \sum_{x_m} \mathbb{I} \left(\frac{\partial \text{output}}{\partial \text{activations}} \cdot (c \text{ direction}) > 0 \right)$$

$$= \frac{1}{|X_m|} \sum_{x_i \in X_m} \mathbb{I} \left(\frac{\partial f(x_i)}{\partial f_1(x_i)} \cdot v > 0 \right)$$

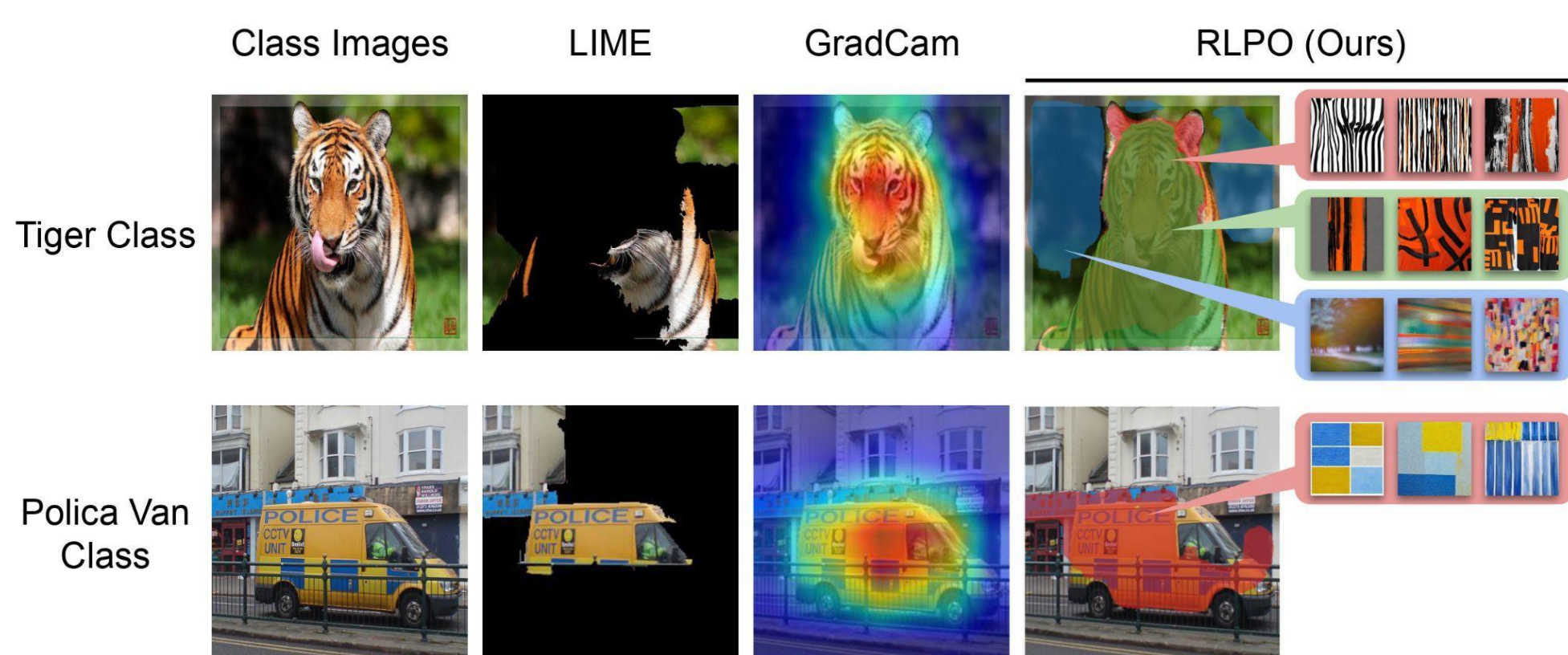
$$Q^*(s, a) = \mathbb{E}_{s' \sim P(\cdot | s, a)} [\xi_t r(s, a) + \gamma \max_{a' \in A} Q_{\text{target}}(s', a')]$$

If results are poor, preference optimization nudges the generative model toward the more promising concepts. Iteratively, the agent and preference optimization together refine the generation process, producing increasingly meaningful concepts.

Why RL? The search space is extremely large and fine-tuning is sequential and dependent on TCAV.

Experiments & Analysis

1. Our method produces **diverse, novel, and unique concepts** which triggers the DNN.



	Laymen (n=260)	Expert (n=240)
EG (Retrieval)	6.54%	10.45%
EG (Ours)	91.54%	65.45%
Odds (Retrieval)	14.29	8.57
Odds (Ours)	0.09	0.53

2. RLPO generates concepts that were **not thought by humans**.

Methods	Concepts	$TS_{c,m}(\uparrow)$	CS (\downarrow)	ED (\uparrow)	RCS (\downarrow)	RED (\uparrow)
EAC	C	1.0	0.76 \pm 0.03	7.21 \pm 0.63	0.67 \pm 0.14	6.34 \pm 2.16
Lens	C1	1.0	0.77 \pm 0.02	7.17 \pm 0.34	0.50 \pm 0.18	9.70 \pm 3.20
	C2	1.0	0.72 \pm 0.04	8.02 \pm 0.87	0.42 \pm 0.10	10.90 \pm 2.80
	C3	1.0	0.69 \pm 0.05	8.45 \pm 0.96	0.45 \pm 0.05	11.03 \pm 2.17
CRAFT	C1	1.0	0.76 \pm 0.04	7.37 \pm 0.62	0.57 \pm 0.16	8.80 \pm 3.20
	C2	1.0	0.72 \pm 0.02	8.25 \pm 0.39	0.50 \pm 1.90	9.90 \pm 3.40
	C3	1.0	0.73 \pm 0.04	7.98 \pm 0.79	0.44 \pm 0.07	10.80 \pm 1.90
RLPO (Ours)	C1	1.0	0.52 \pm 0.04	10.48 \pm 0.50	0.04 \pm 0.01	16.80 \pm 1.40
	C2	1.0	0.49 \pm 0.02	10.65 \pm 0.20	0.02 \pm 0.02	17.20 \pm 0.80
	C3	1.0	0.49 \pm 0.02	10.74 \pm 0.30	0.03 \pm 0.01	17.60 \pm 4.40

Positive Prompt

The customer service team was very helpful and responsive when I reached out for support. They were patient and provided clear instructions on how to address some of the issues, which improved the situation slightly.

Generated Concepts

Customer: client, purchaser, consumer, user, shopper
Team: group, crew, unit, squad, alliance, partnership
Helpful: supportive, useful, valuable, beneficial, productive
Clear: transparent, unclouded, open, lucid, distinct
Address: speak, contact, communicate, interact, approach
Issues: problems, concerns, matters, challenges, disputes

3. RLPO can reveal **influential word level concepts** in NLP tasks like sentiment analysis.

1. Kim, Been, et al. "Interpretability beyond feature attribution: Quantitative testing with concept activation vectors (tcav)." International conference on machine learning. PMLR, 2018.
2. Lisa Schut, Nenad Tomasev, Tom McGrath, Demis Hassabis, Ulrich Paquet, and Been Kim. Bridging the human-ai knowledge gap: Concept discovery and transfer in alphazero. arXiv preprint arXiv:2310.16410, 2023.
3. Som Sagar*, Aditya Taparia*, Harsh Mankodiya, Pranav Bidare, Yifan Zhou, and Ransalu Senanayake, "Trustworthy Conceptual Explanations for Neural Networks in Robot Decision-Making," IROS 2025.