# ConceptScope: Characterizing Dataset Bias via Disentangled Visual Concepts

NeurIPS 2025

Jinho Choi
KAIST AI
Kim Jaechul Graduate School

Hyesu Lim
KAIST AI
Kim Jaechul Graduate School

Steffen Schneider
HELMHOLTZ MUNICH

Jaegul Choo
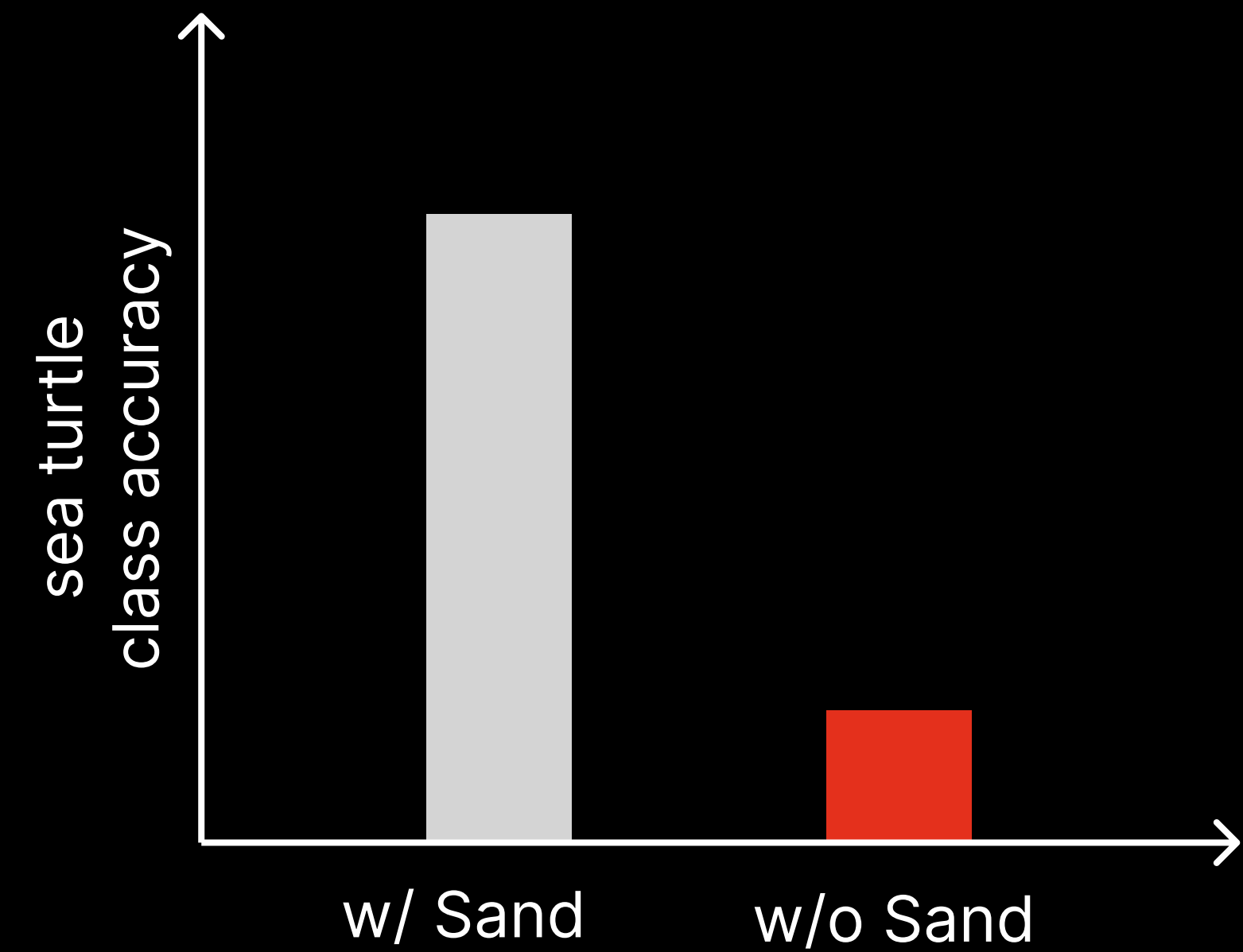KAIST AI
Kim Jaechul Graduate School

# Motivation

Training set

?

# Motivation



Training set

w/ Sand >> w/o Sand

**Collection bias exits in datasets**
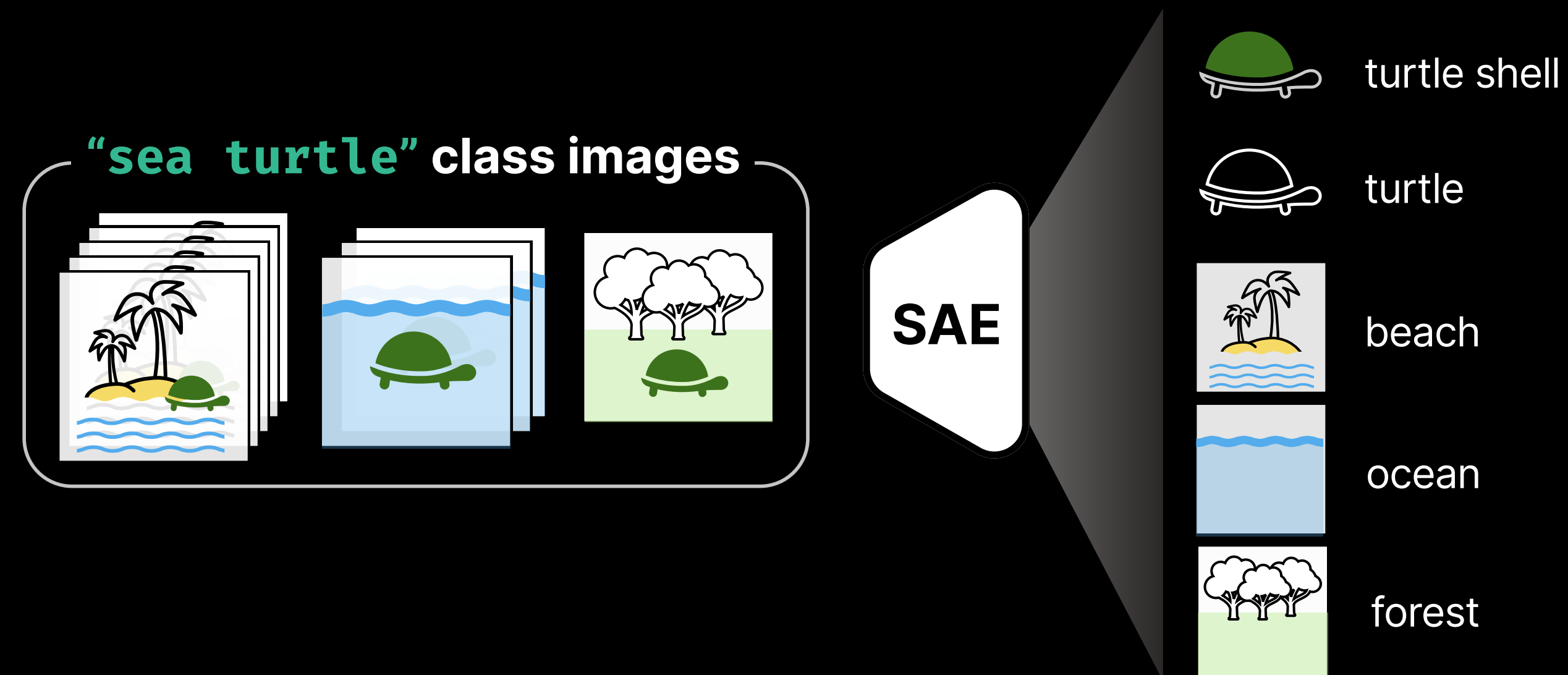


**Dataset bias lead to model bias**

# Existing Approach



👨‍💻

*A small sea turtle crawls across the sandy beach toward the ocean waves*

👩‍💻

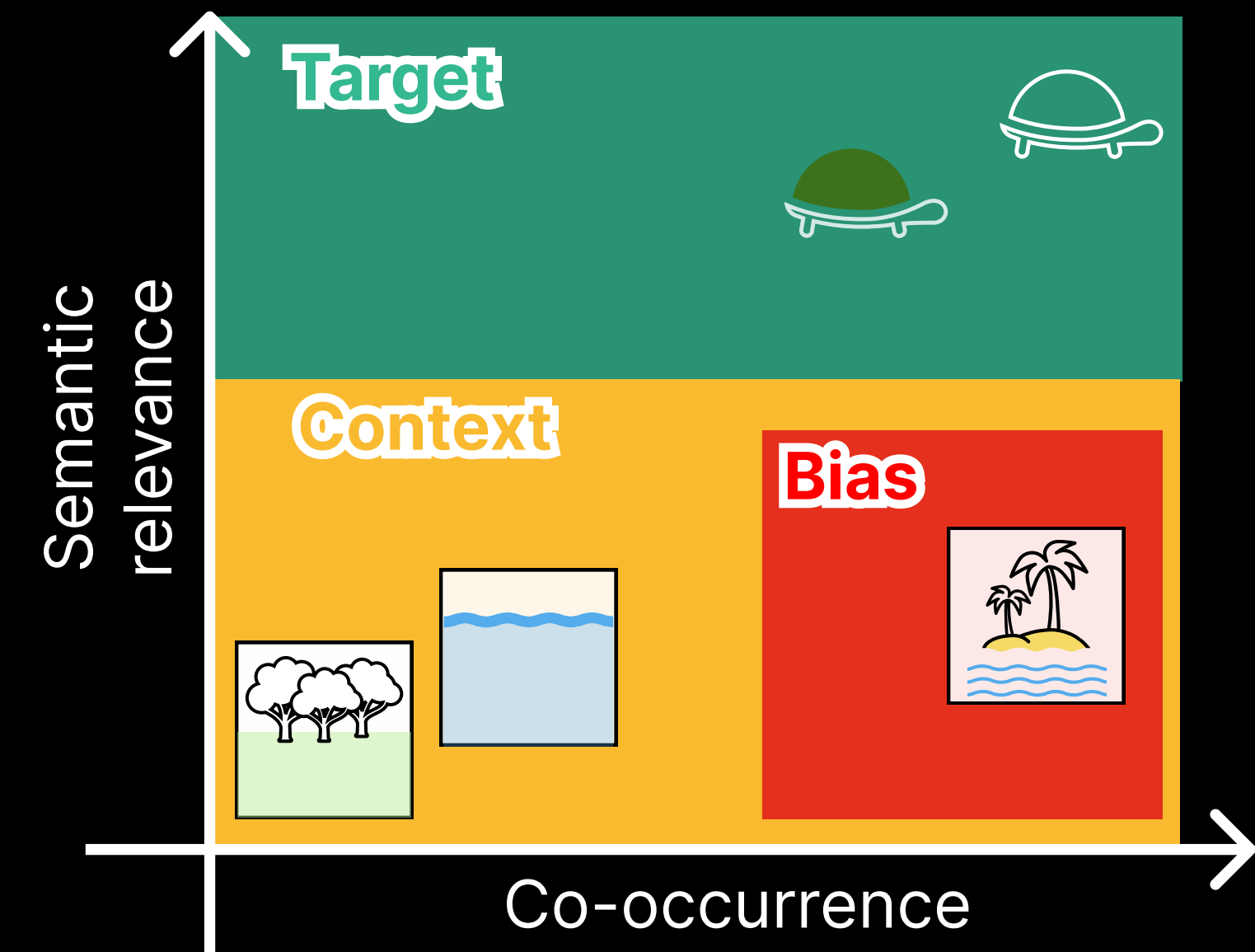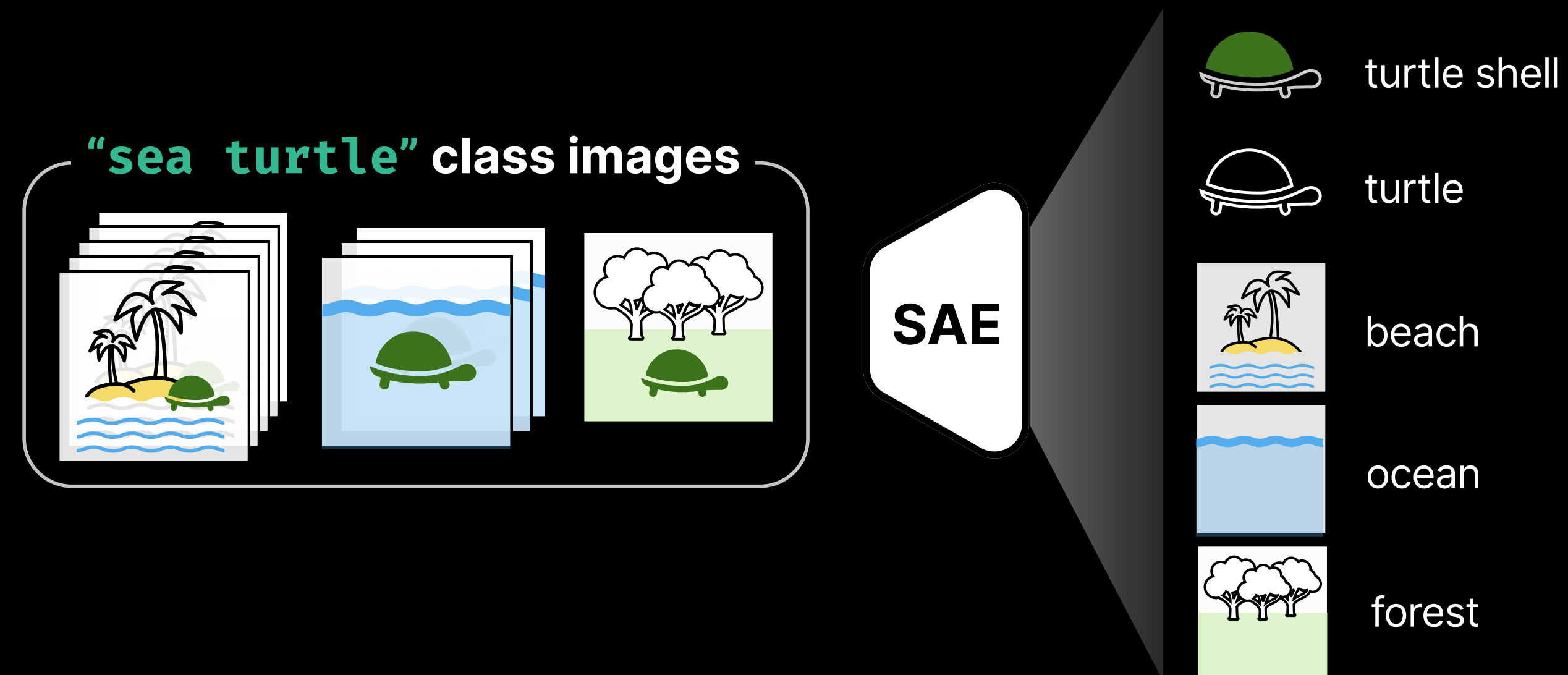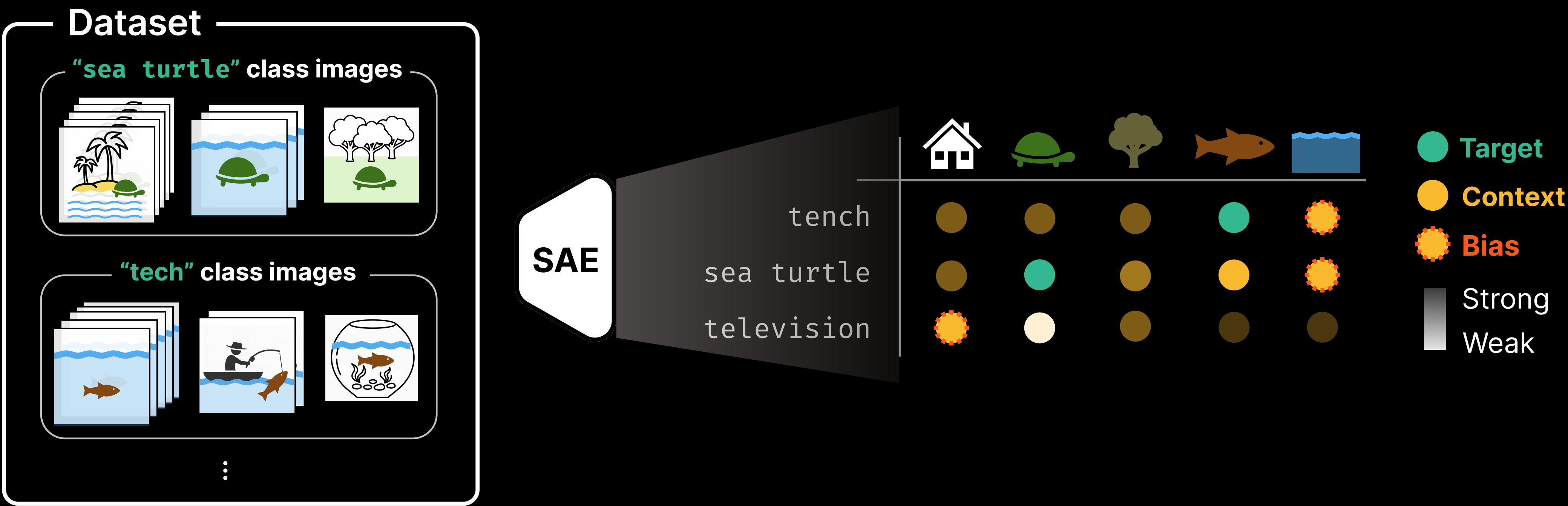*A newly hatched leatherback turtle makes a dash for the water*

# Our Approach: ConceptScpe

Sparse Autoencoder (SAE) as a concept extractor

# Our Approach: ConceptScpe

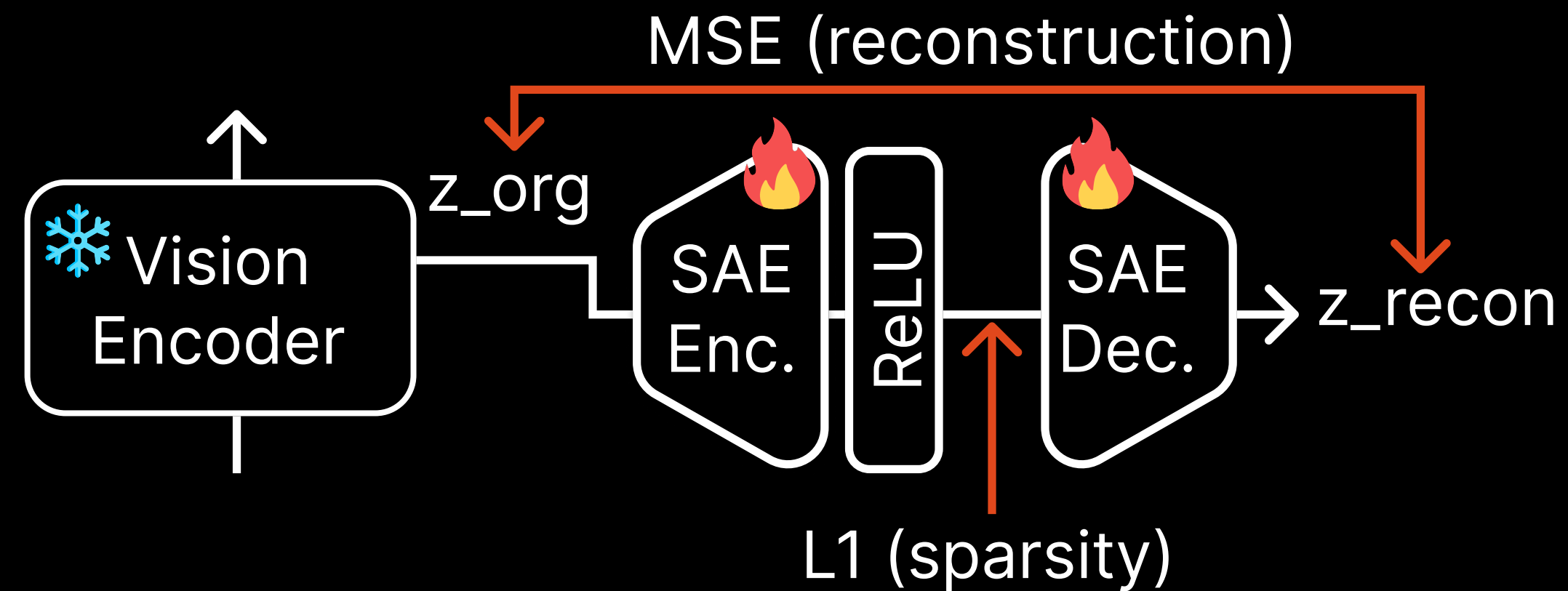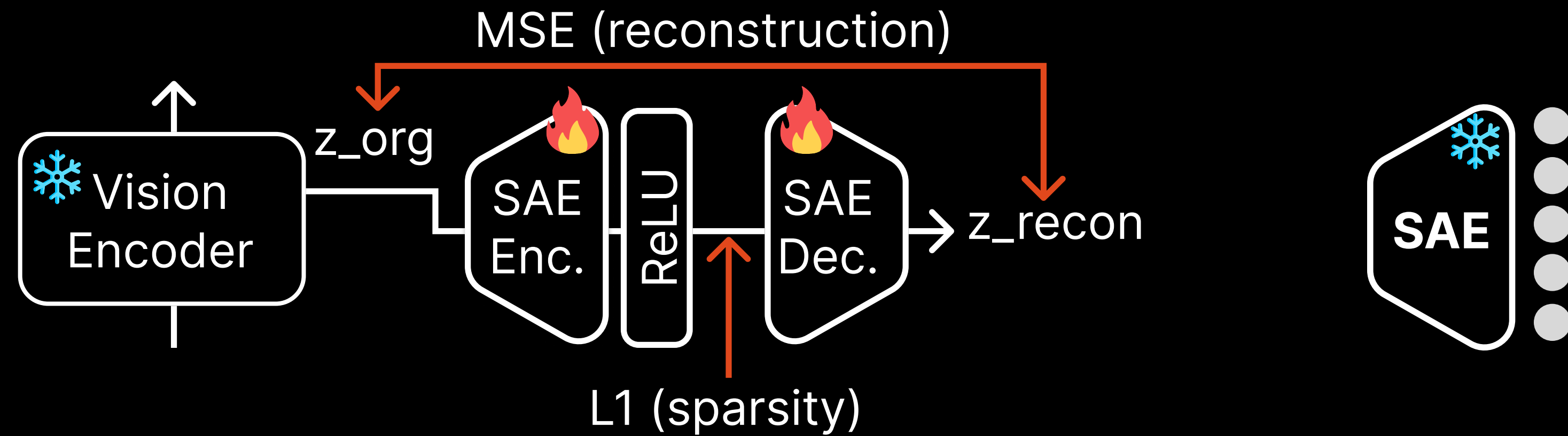Sparse Autoencoder (SAE) as a concept extractor

# Our Approach: ConceptScpe

Sparse Autoencoder (SAE) as a concept extractor

# ConceptScpe: Characterizing dataset bias via visual concepts
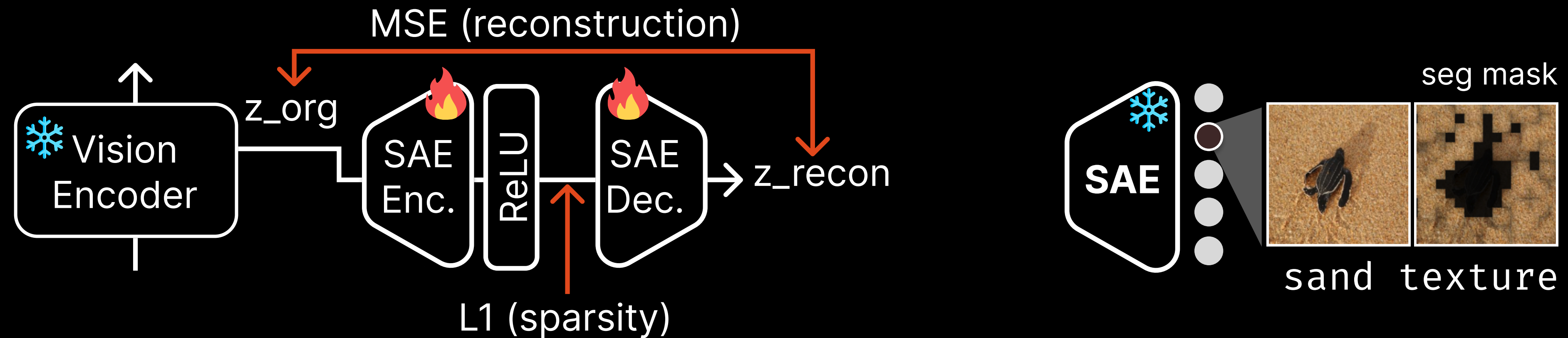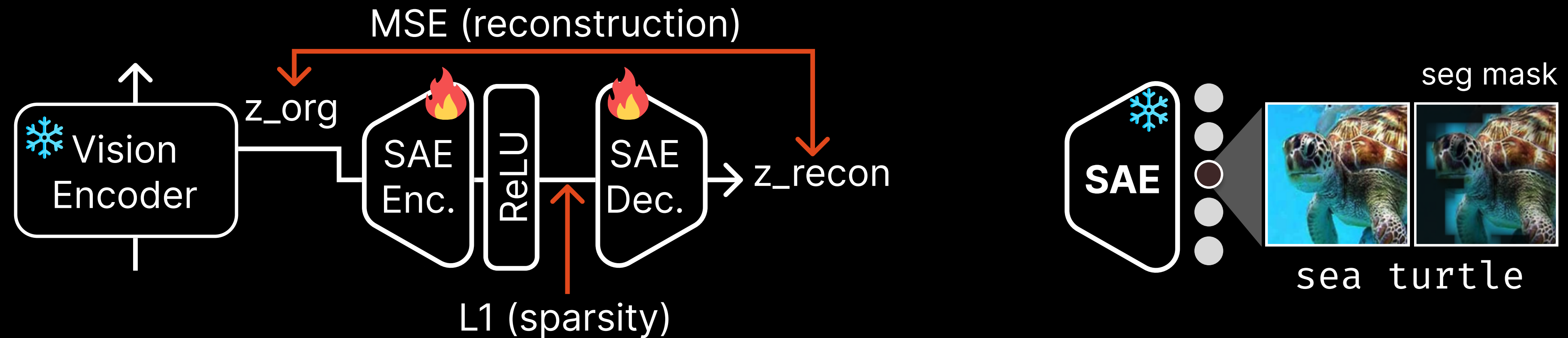
Training SAE and Interpreting Concepts



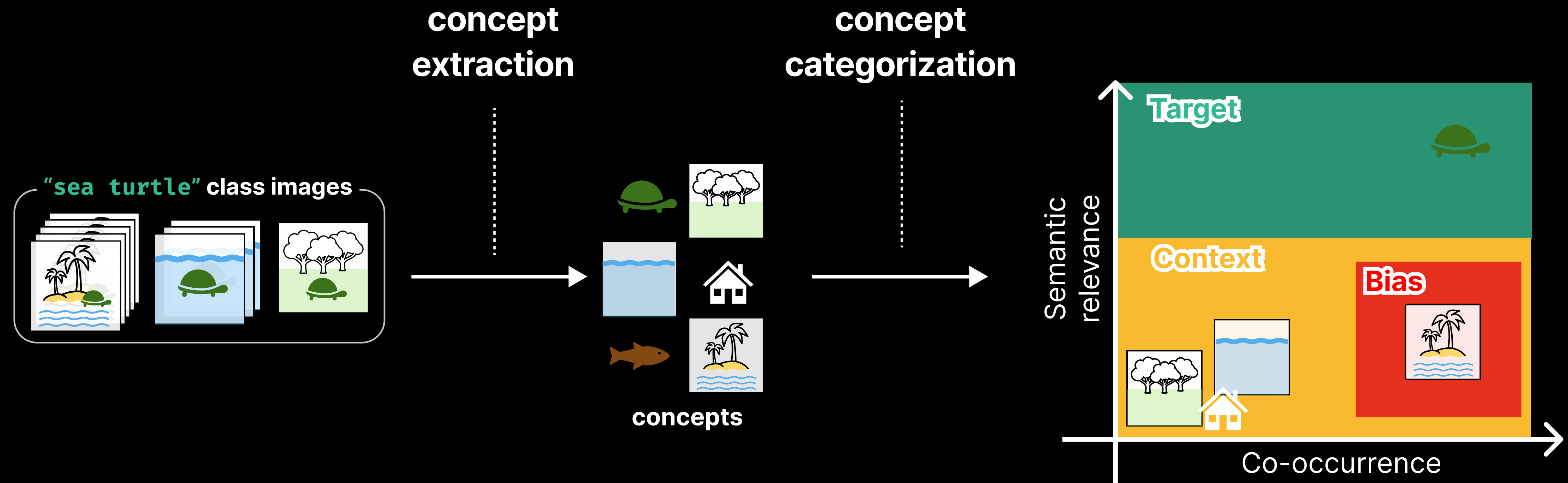Lim, Hyesu, et al. "Sparse autoencoders reveal selective remapping of visual concepts during adaptation." (ICLR 2025)
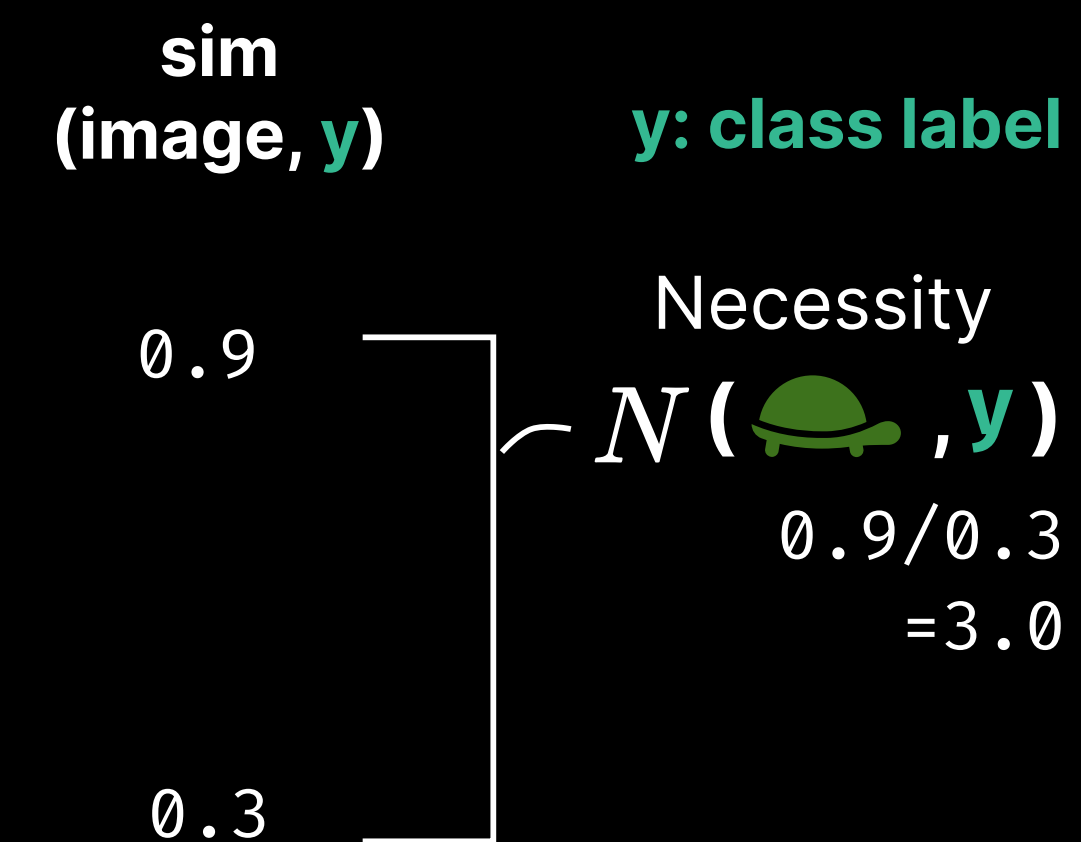
8 of 22

# ConceptScpe: Characterizing dataset bias via visual concepts
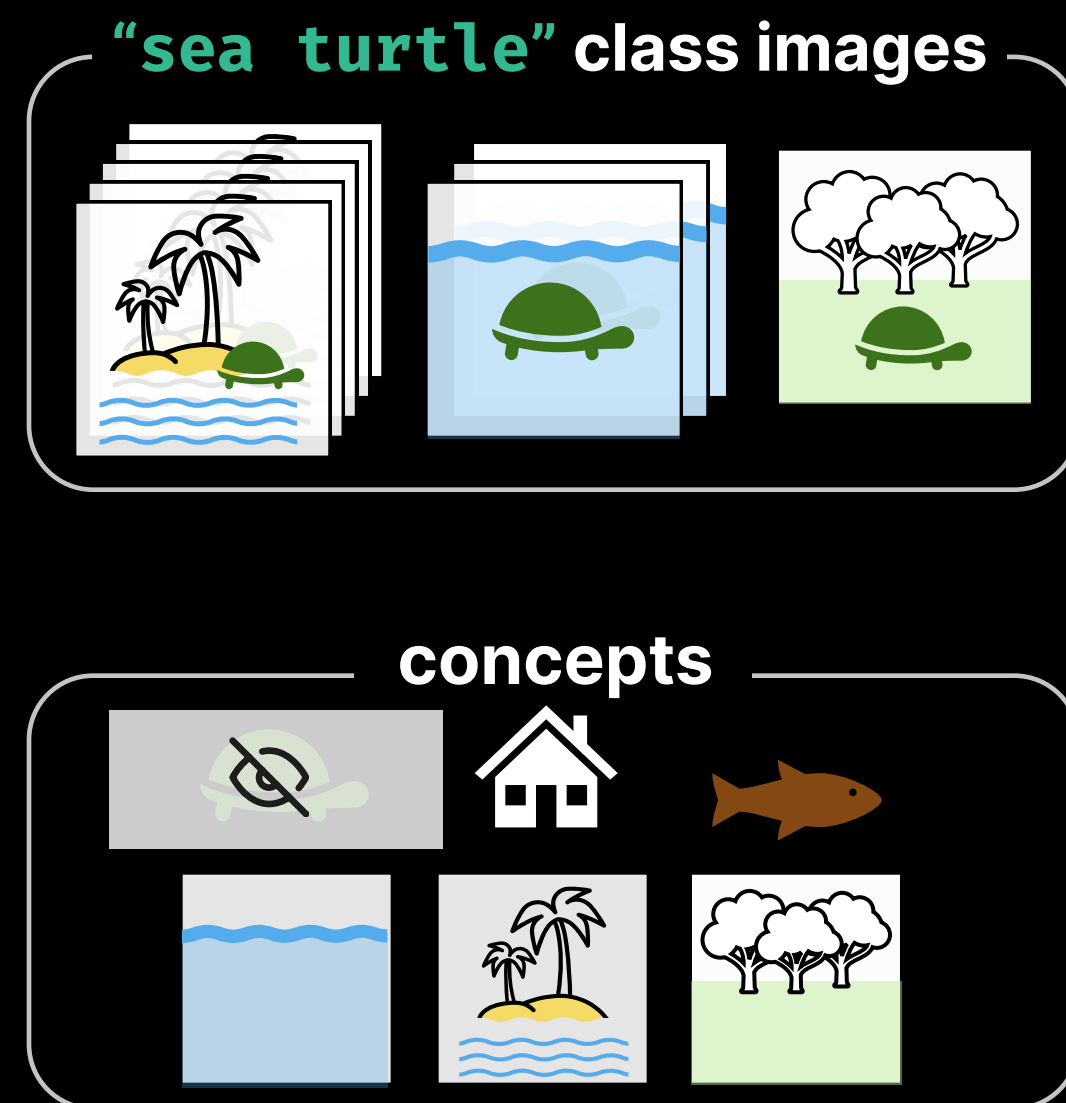
Training SAE and Interpreting Concepts

Lim, Hyesu, et al. "Sparse autoencoders reveal selective remapping of visual concepts during adaptation." (ICLR 2025)

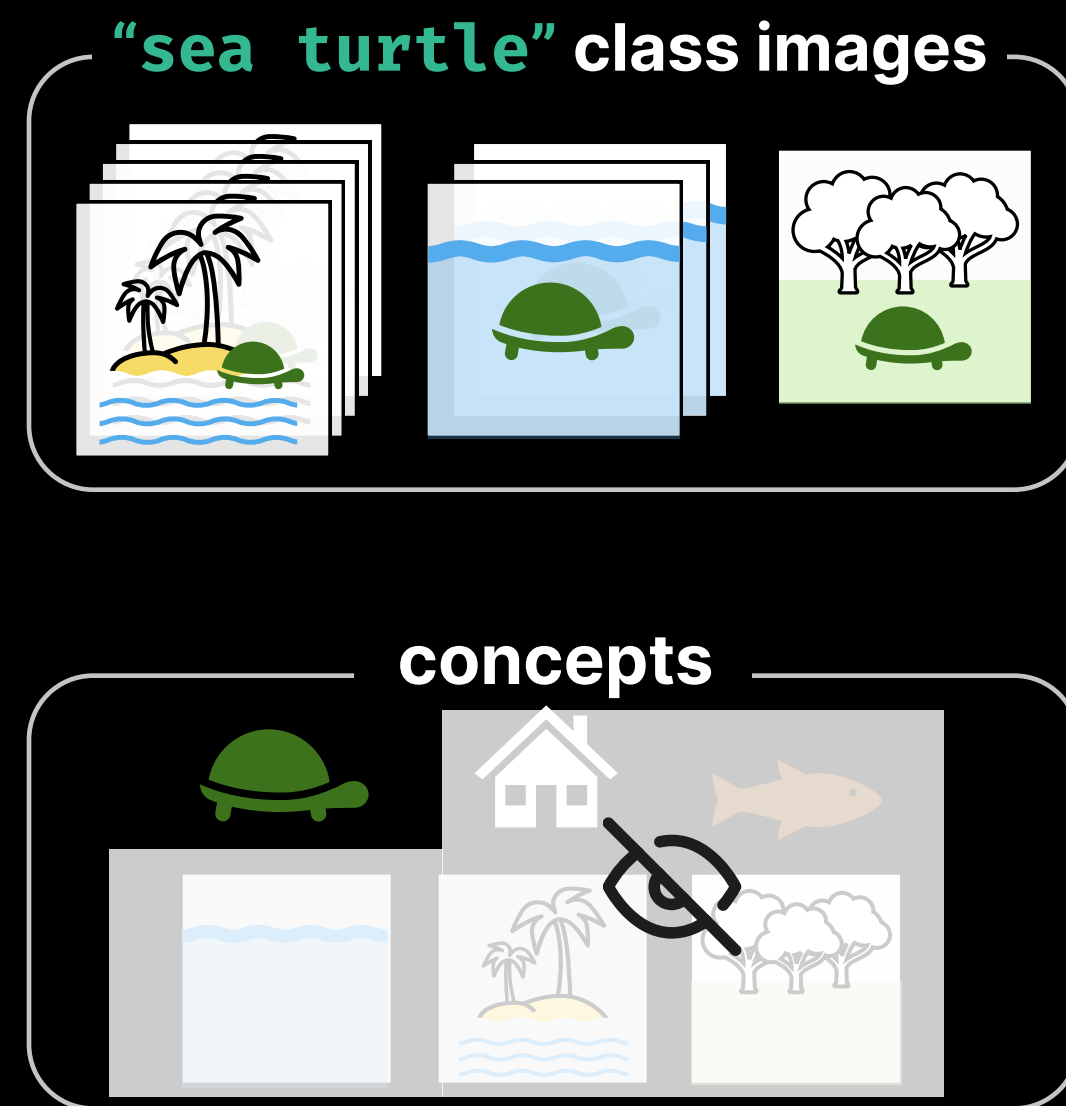# ConceptScpe: Characterizing dataset bias via visual concepts

Training SAE and Interpreting Concepts



MSE (reconstruction)

z_org

Vision Encoder

SAE Enc.

ReLU

SAE Dec.

z_recon

L1 (sparsity)

SAE

seg mask

sand texture

Lim, Hyesu, et al. "Sparse autoencoders reveal selective remapping of visual concepts during adaptation." (ICLR 2025)

# ConceptScpe: Characterizing dataset bias via visual concepts

Training SAE and Interpreting Concepts

Lim, Hyesu, et al. "Sparse autoencoders reveal selective remapping of visual concepts during adaptation." (ICLR 2025)

# ConceptScpe: Characterizing dataset bias via visual concepts

Training SAE and Interpreting Concepts



Lim, Hyesu, et al. "Sparse autoencoders reveal selective remapping of visual concepts during adaptation." (ICLR 2025)

# ConceptScpe: Characterizing dataset bias via visual concepts

Categorizing Concepts

# ConceptScpe: Characterizing dataset bias via visual concepts

Computing alignment score

# ConceptScpe: Characterizing dataset bias via visual concepts

Computing alignment score

# ConceptScpe: Characterizing dataset bias via visual concepts

Computing alignment score

# ConceptScpe: Characterizing dataset bias via visual concepts

Categorizing Concepts

# Results: SAEs can discover a wide range of visual concepts


knitted fabric


yellow objects


group of dogs


seat belts


lace patterns


orange objects


puppies


car gauges

| Method | Metric | Caltech101 (Objects) | DTD (Textures) | Waterbird (Backgrounds) | CelebA (Facial Attr.) | RAF-DB (Emotions) | Stanford40 (Actions) | Average |
|---|---|---|---|---|---|---|---|---|
| BLIP-2 | $F_1$ | $0.64 \pm 0.35$ | $0.38 \pm 0.25$ | $0.37 \pm 0.10$ | $0.27 \pm 0.24$ | $0.24 \pm 0.17$ | $0.66 \pm 0.18$ | 0.43 |
| LLaVA-NeXT | $F_1$ | $0.61 \pm 0.35$ | $0.40 \pm 0.21$ | $0.57 \pm 0.12$ | $0.62 \pm 0.24$ | $0.45 \pm 0.18$ | $0.80 \pm 0.16$ | 0.58 |
| **ConceptScope** | $F_1$ | $0.83 \pm 0.21$ | $0.57 \pm 0.20$ | $0.78 \pm 0.07$ | $0.81 \pm 0.11$ | $0.55 \pm 0.18$ | $0.78 \pm 0.13$ | 0.72 |
| | AUPRC | $0.89 \pm 0.19$ | $0.57 \pm 0.23$ | $0.83 \pm 0.09$ | $0.85 \pm 0.13$ | $0.59 \pm 0.21$ | $0.82 \pm 0.15$ | 0.76 |

# Results: ConceptScope captures diverse visual states within each class



class name | **target** | **bias** | **context**

**cauliflower**
cauliflower | green cauliflower | green leaves | bowl | market | cooked

**bubble**
soap bubbles | glass art | child | urban | festival | park

**mountain**
landscapes | mountain | snow | blue sky | summits | trees

**hair wig**
hair wig | wearing a wig | mannequin | toy models | models | colorful hairs

# Results: ConceptScpe discovers real-world dataset bias



High bias ▭ Low bias ▭

**ImageNet** - "**balance beam**"
biased to "**competition**"

**ImageNet** - "**afghan hound**"
biased to "**dog show**"

**SUN397** - "**ice skating rink**" class
biased to "**New York**"

**Food101** - "**hotdog**"
biased to " **food wrappers**"

**Food101** - "**bibmbap**"
biased to " **fried eggs**"

**ImageNet** - "**bridgeroom**" class
biased to "**east asian culture**"

# Results: ConceptScpe discovers real-world dataset bias
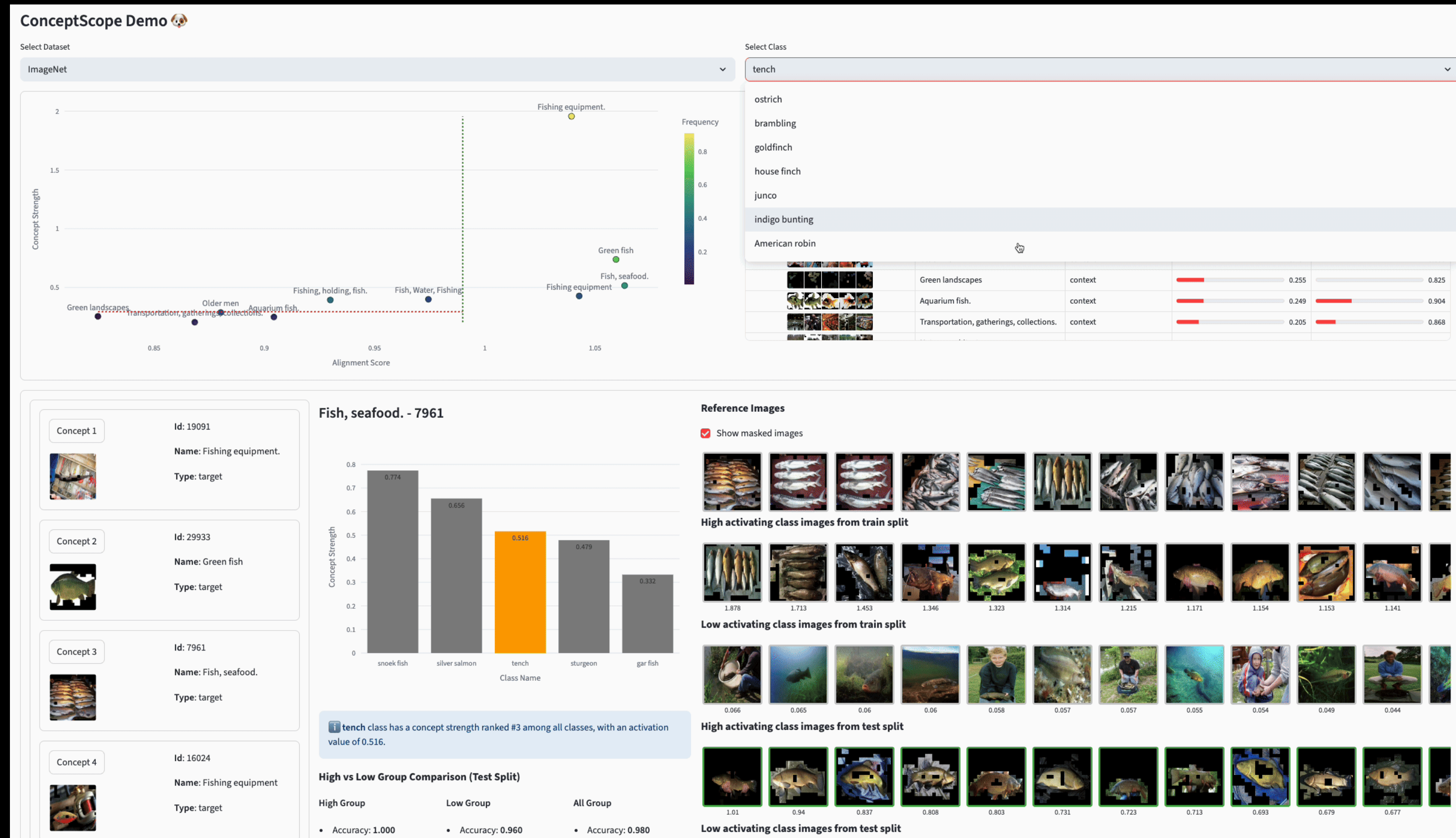
**CelebA** - "**blond hair**"
biased to "**female**"



**Waterbirds** - "waterbirds"
biased to "**ocean**"



**Nico++** - "**mammal**" class
biased to "**rock**"

| Method | Waterbirds | CelebA | Nico++(75) | Nico++ (90) | Nico++ (95) |
|---|---|---|---|---|---|
| DOMINO | 90.0% | 87.0% | 24.0% | 24.0% | 24.0% |
| FACTS | 100.0 % | 100.0% | 55.0% | 60.8% | 61.0% |
| ViG-Bias | 100.0% | 100.0% | 60.0% | 66.7% | 65.0% |
| **ConceptScope (Ours)** | **100.0%** | **100.0%** | **72.9%** | **73.1%** | **74.0%** |

# ConceptScpe: Characterizing dataset bias via visual concepts



## Project page

https://jjho-choi.github.io/ConcepScope-projectpage/

## Code & Demo

https://github.com/jjho-choi/ConceptScope