# Test-Time Canonicalization by Foundation Models for Robust Perception

Utkarsh Singhal*   Ryan Feng*   Stella X. Yu   Atul Prakash

* denotes equal contribution

# Mobile Agents Face Difficult and Diverse Input Transformations

Viewpoint

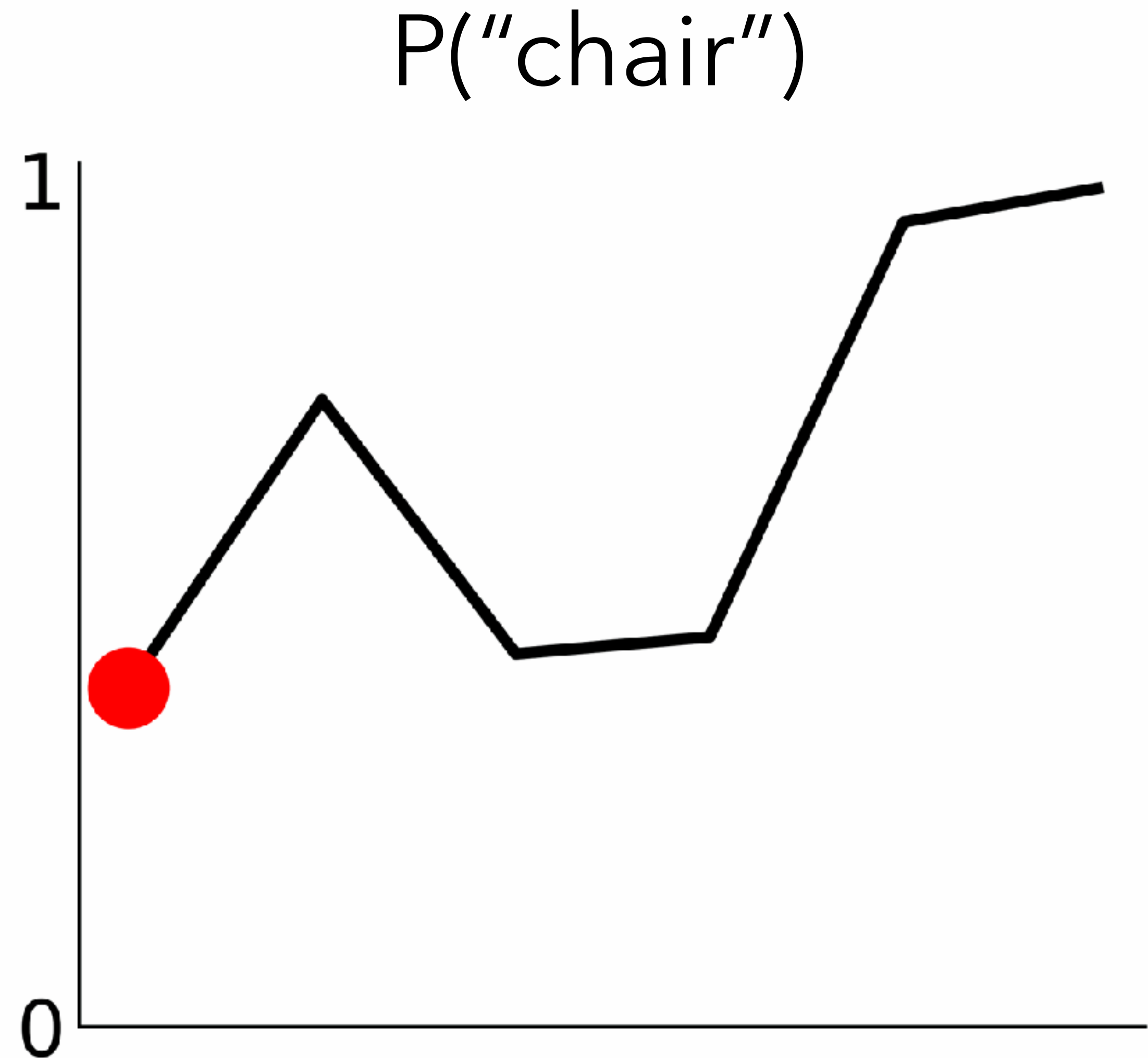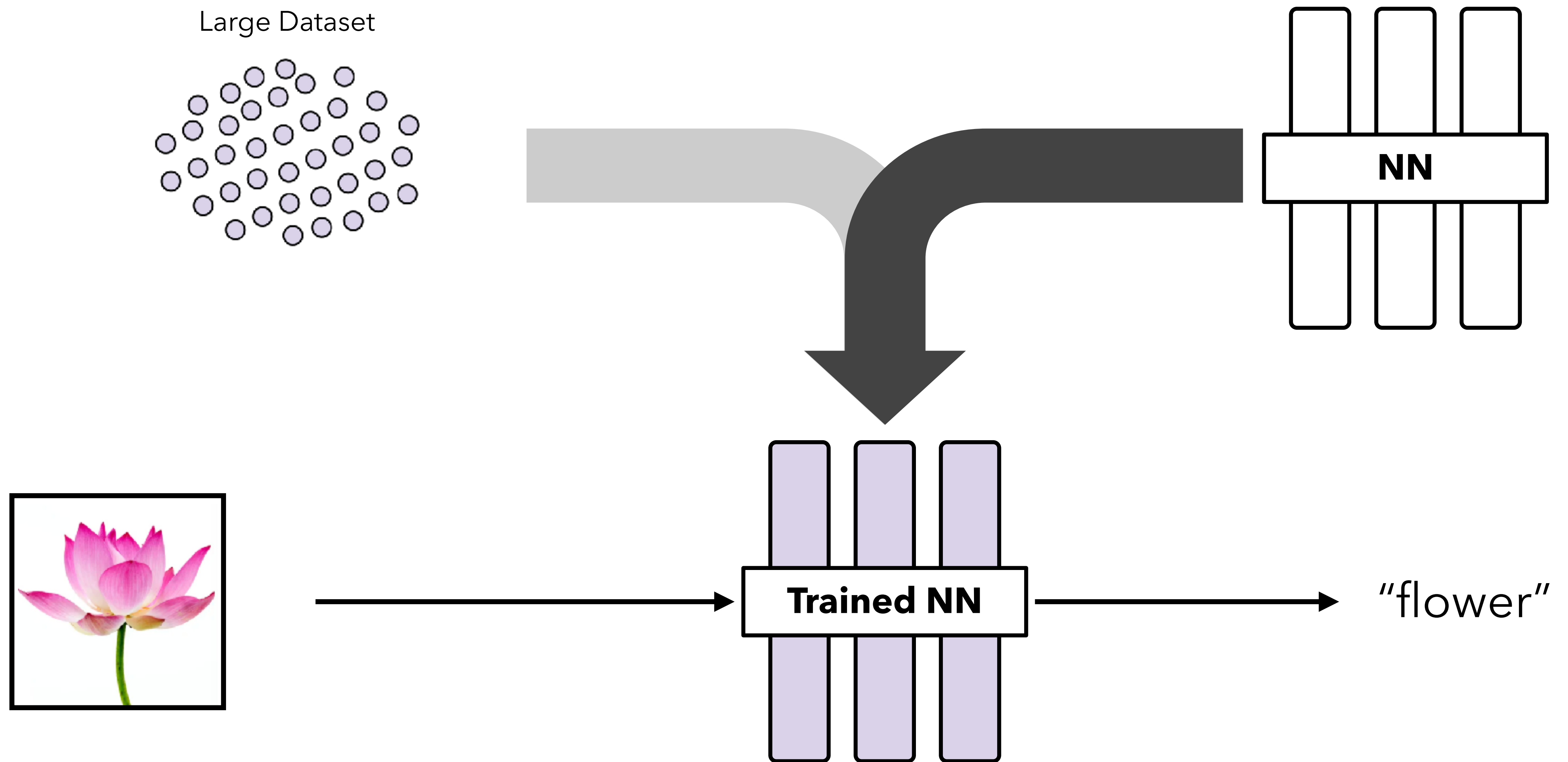Lighting

Environment

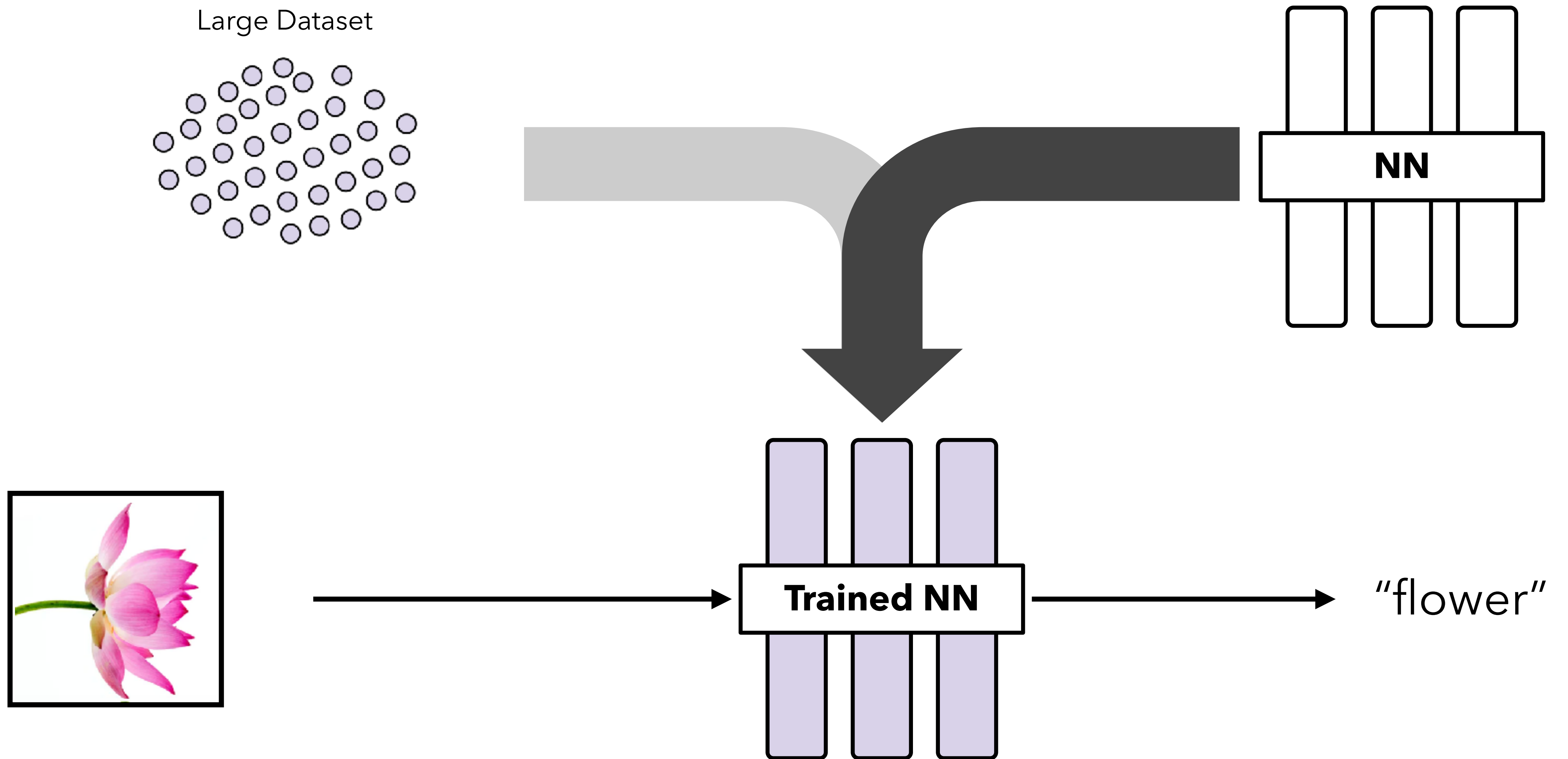# Foundation Models are Still not Robust Enough



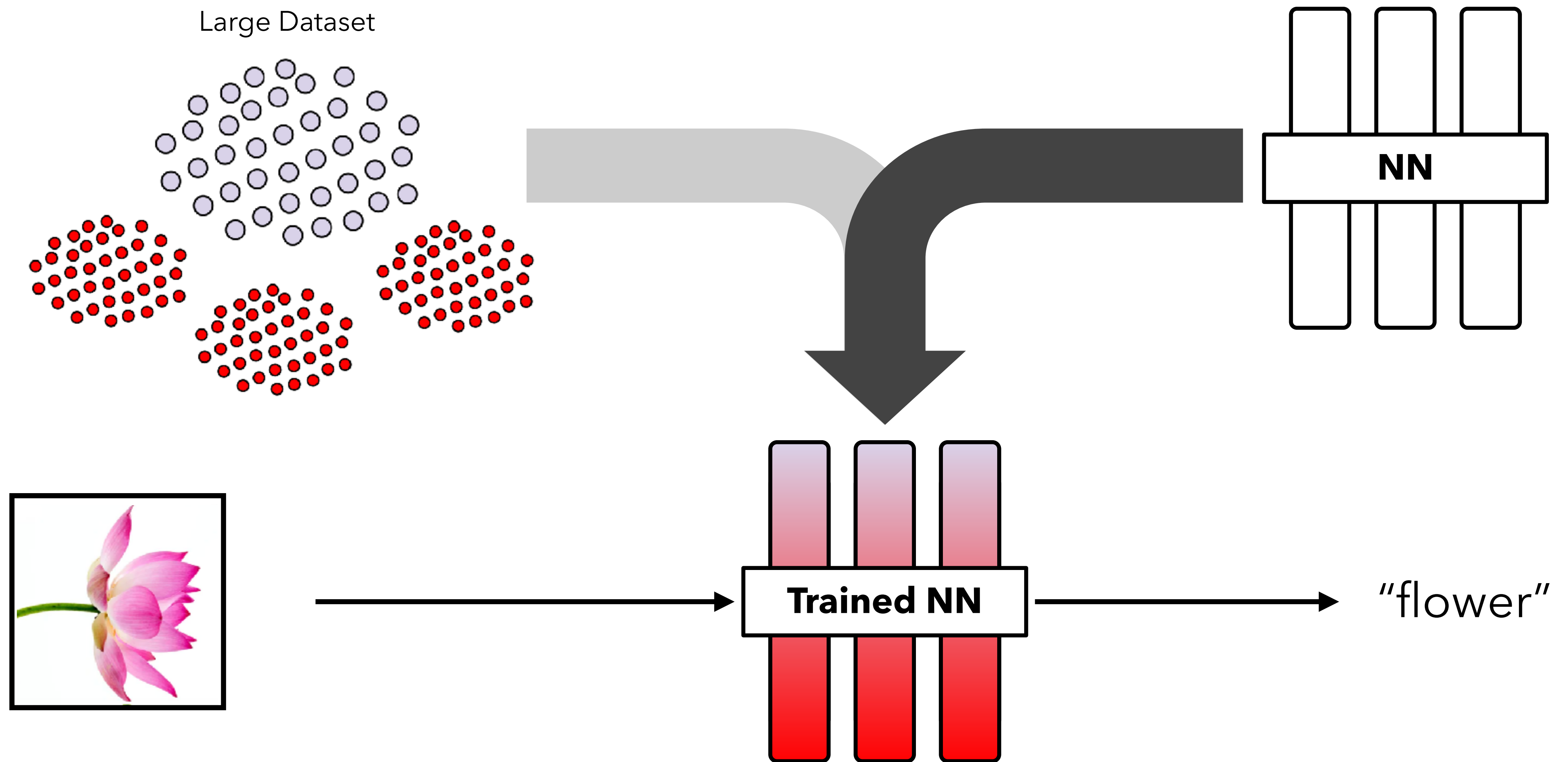P("chair")

# A Standard Pipeline Handles Upright Data

Large Dataset

**NN**

**Trained NN**

"flower"

# How to Handle Transformed Inputs?

Large Dataset
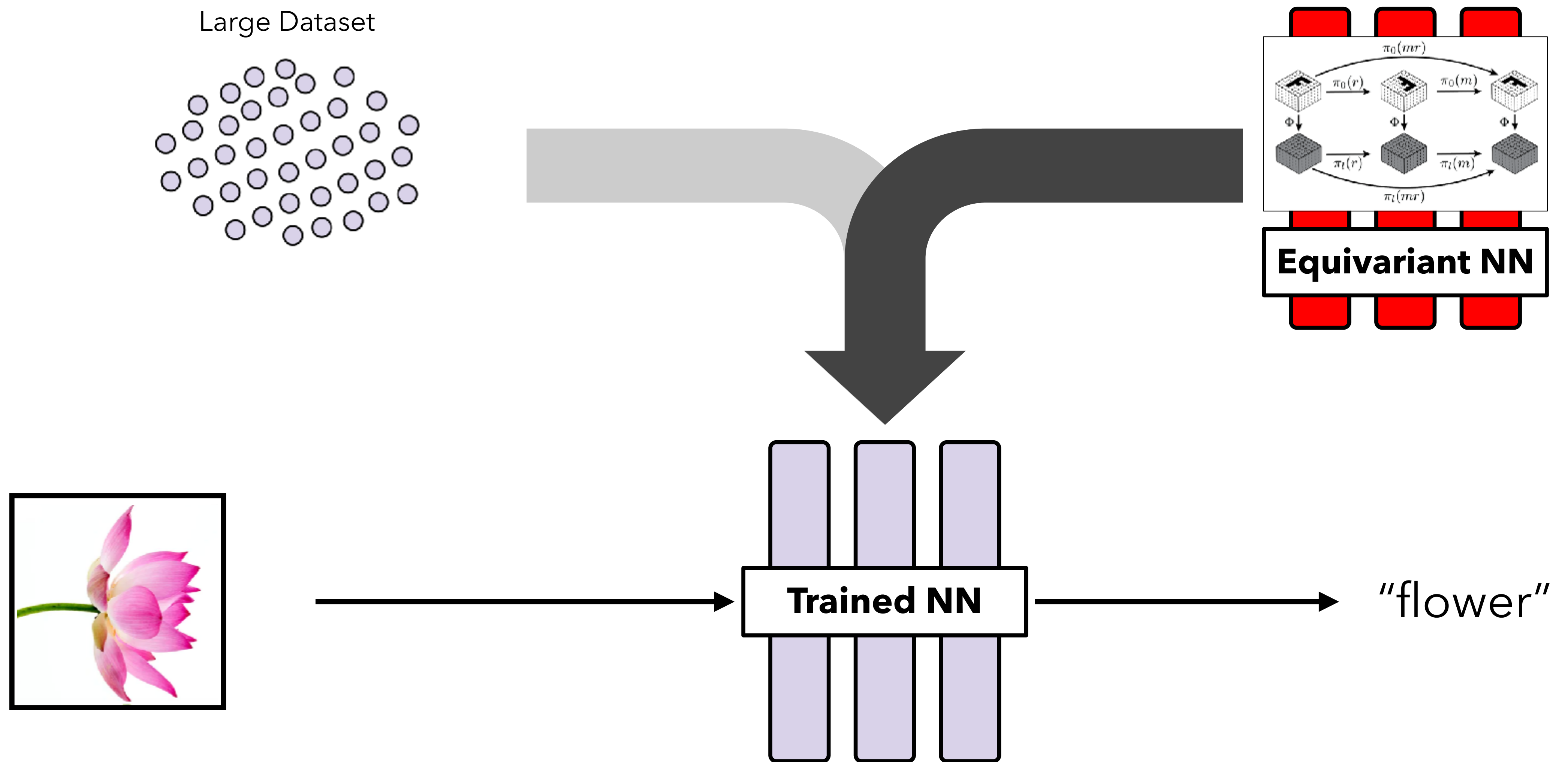
**NN**

**Trained NN**

"flower"

# Data Augmentation: Train on Transformed Data
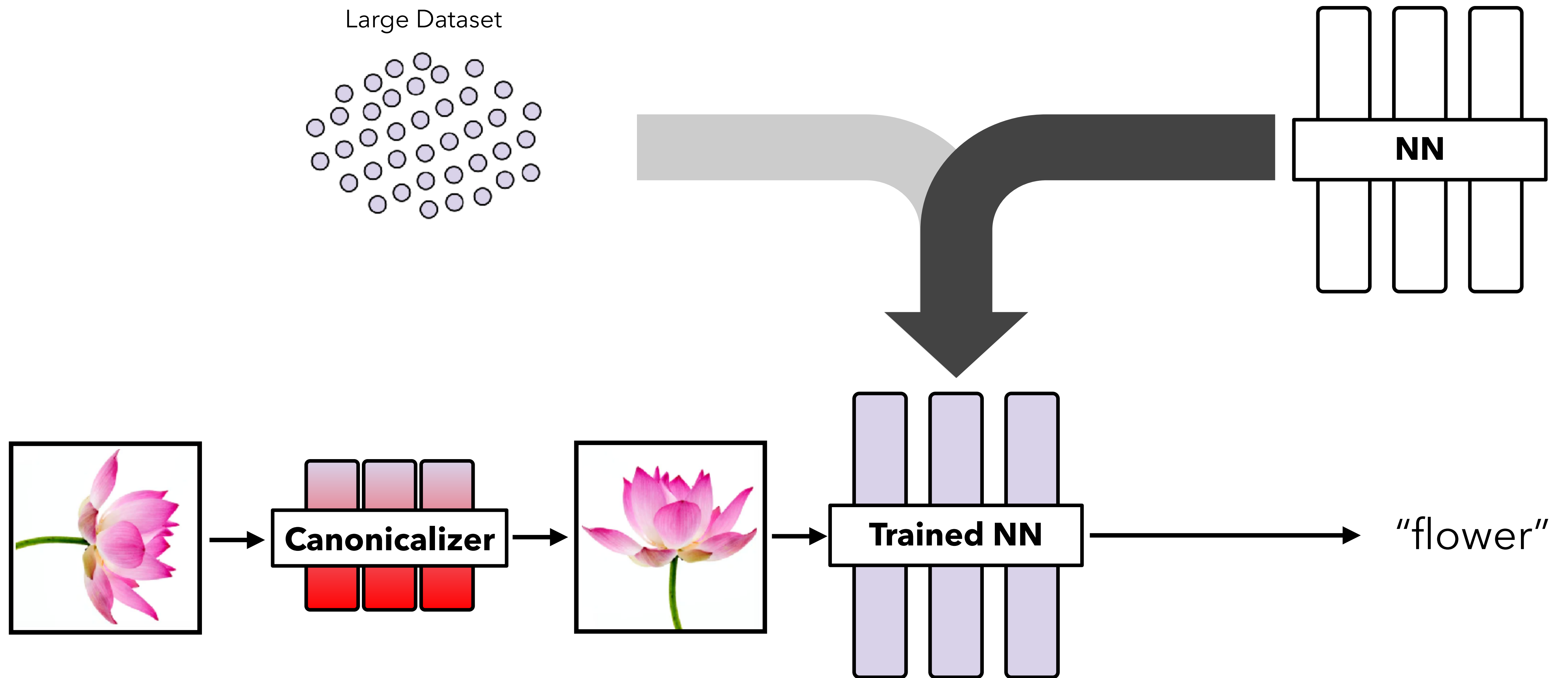
# Equivariant Networks: Transform-Specific Architectures

# Prior Canonicalization Work: Train a Network to Fix the Input



Large Dataset

NN

Canonicalizer

Trained NN

"flower"

# How to Help Foundation Models Face Diverse and Challenging Input Transformations?
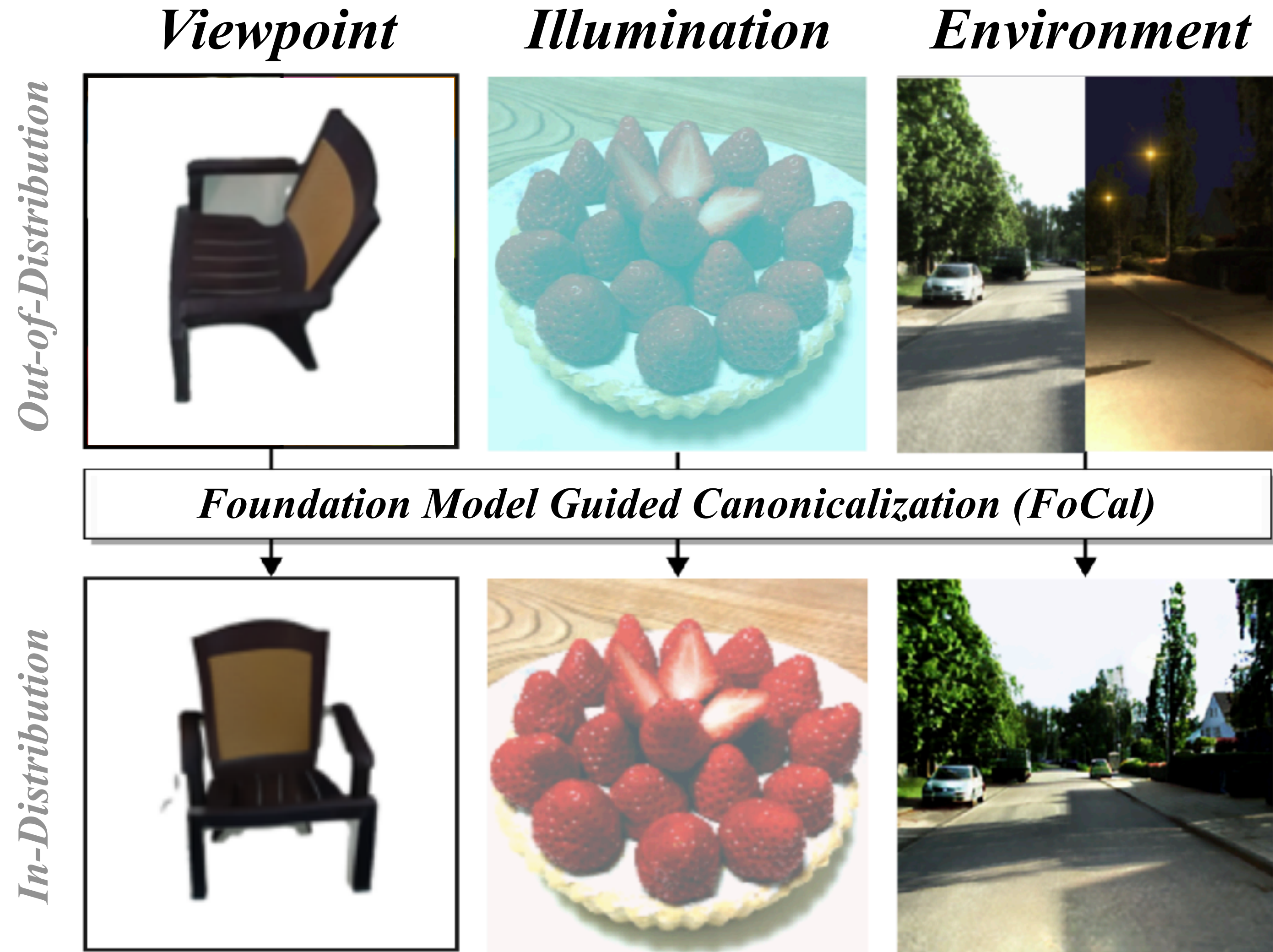
- **Data augmentation?**
  **Train-time transforms only**


- **Equivariant Neural Networks?**
  **Architecture-specific transforms only**


- **Prior Canonicalization work?**
  **Don't generalize to new datasets & transforms**

# FoCal: Foundation Model Guided Canonicalization

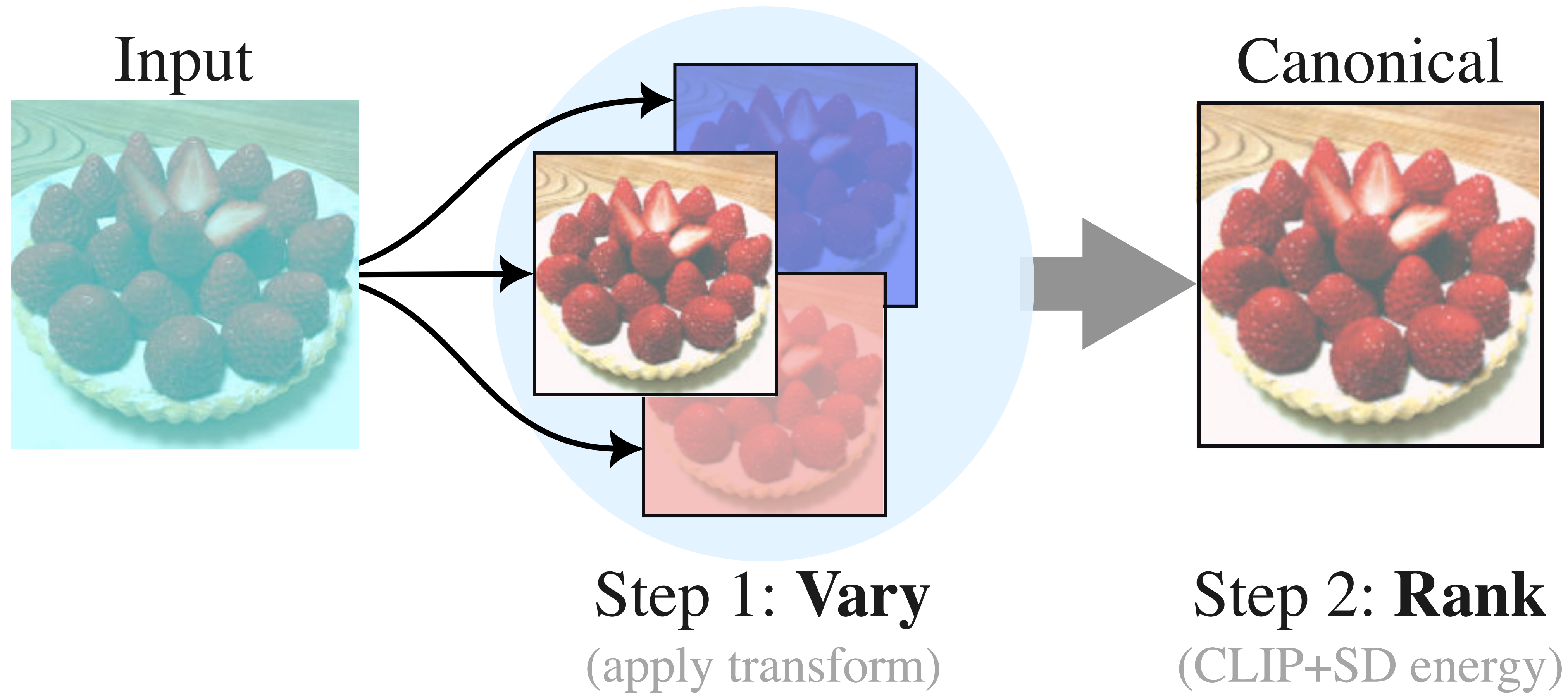**Insight**: Foundation models know what is typical; we use this to convert OOD inputs ⟶ typical

**Benefits:**

- Fully test-time approach
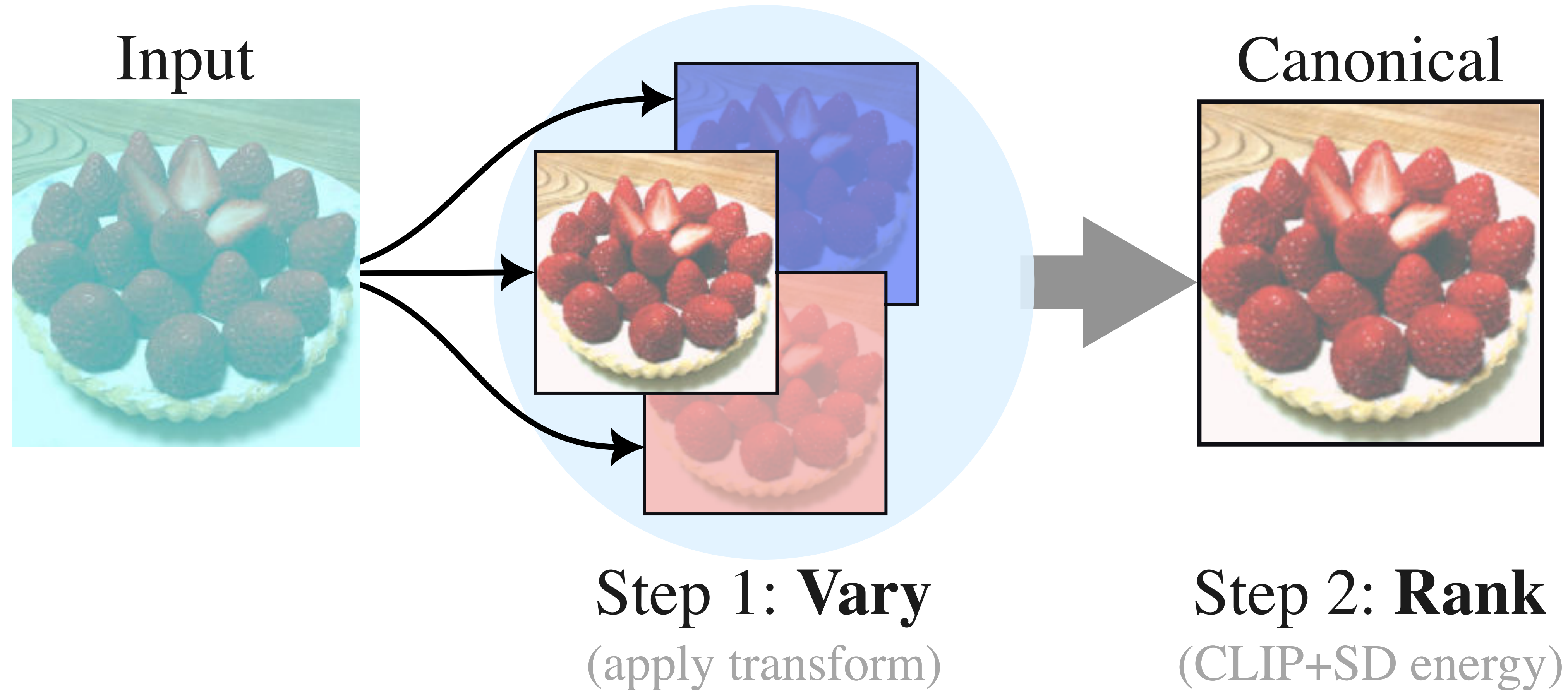
- Generalizes to diverse and complex transformations

# FoCal: Test-Time Canonicalization by Foundation Models



Input

Canonical

Step 1: **Vary**
(apply transform)

Step 2: **Rank**
(CLIP+SD energy)

# FoCal: Test-Time Canonicalization by Foundation Models



Input

Canonical

Step 1: **Vary**
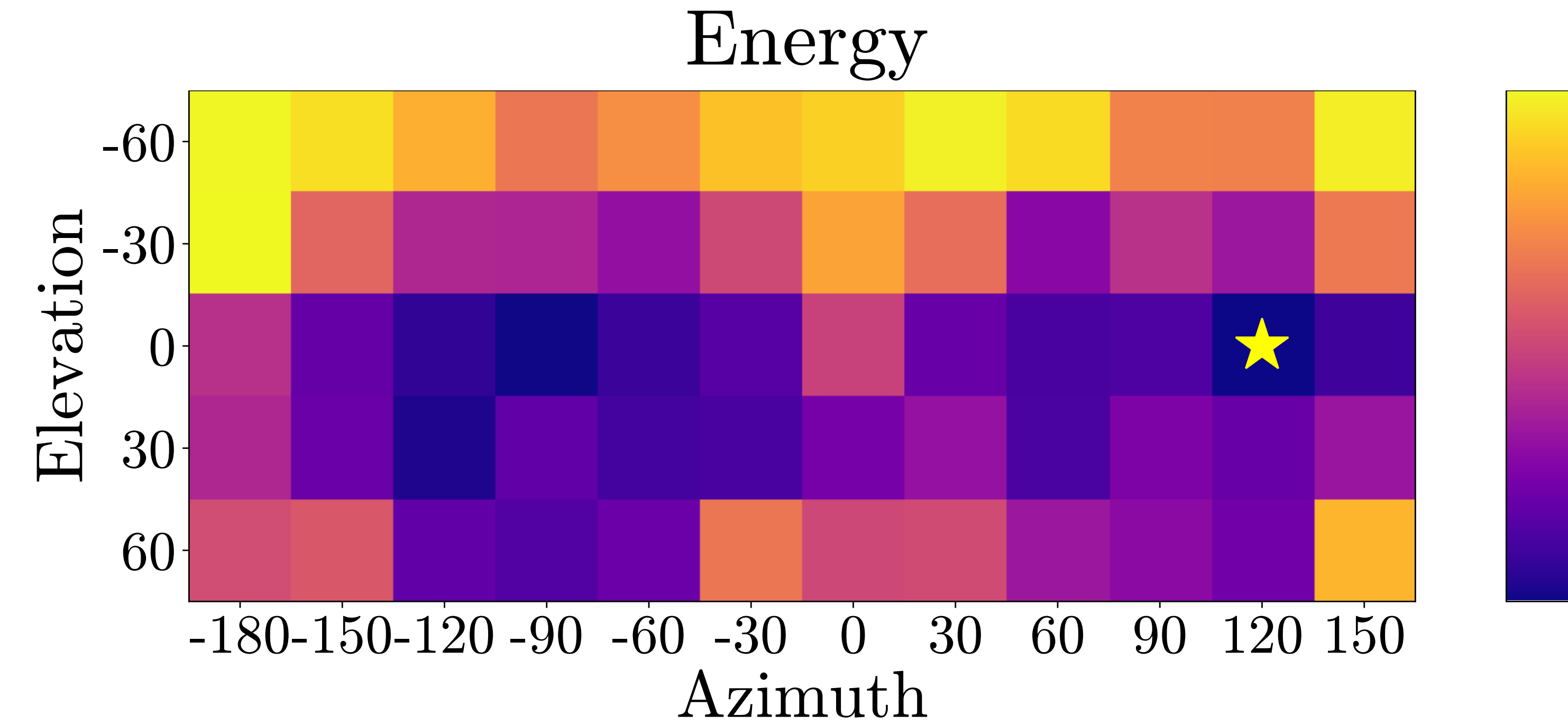(apply transform)

Step 2: **Rank**
(CLIP+SD energy)

FoCal searches over the transformation space and selects the best version:

1. **Vary**: Generate candidate transformed versions
2. **Rank**: Find the best using foundation model priors (energy functions)
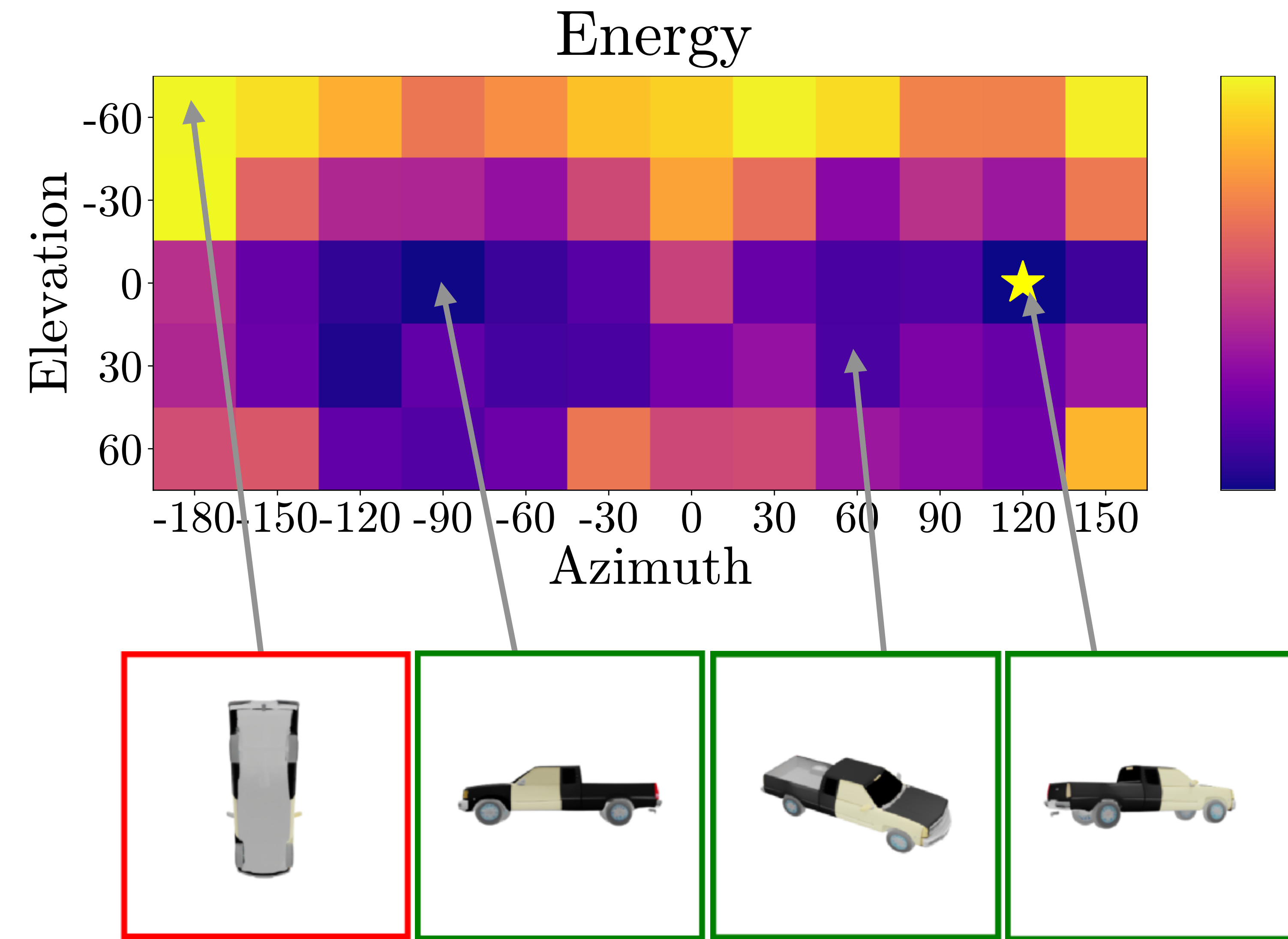
# FoCal Energy Function Picks Typical Viewpoints
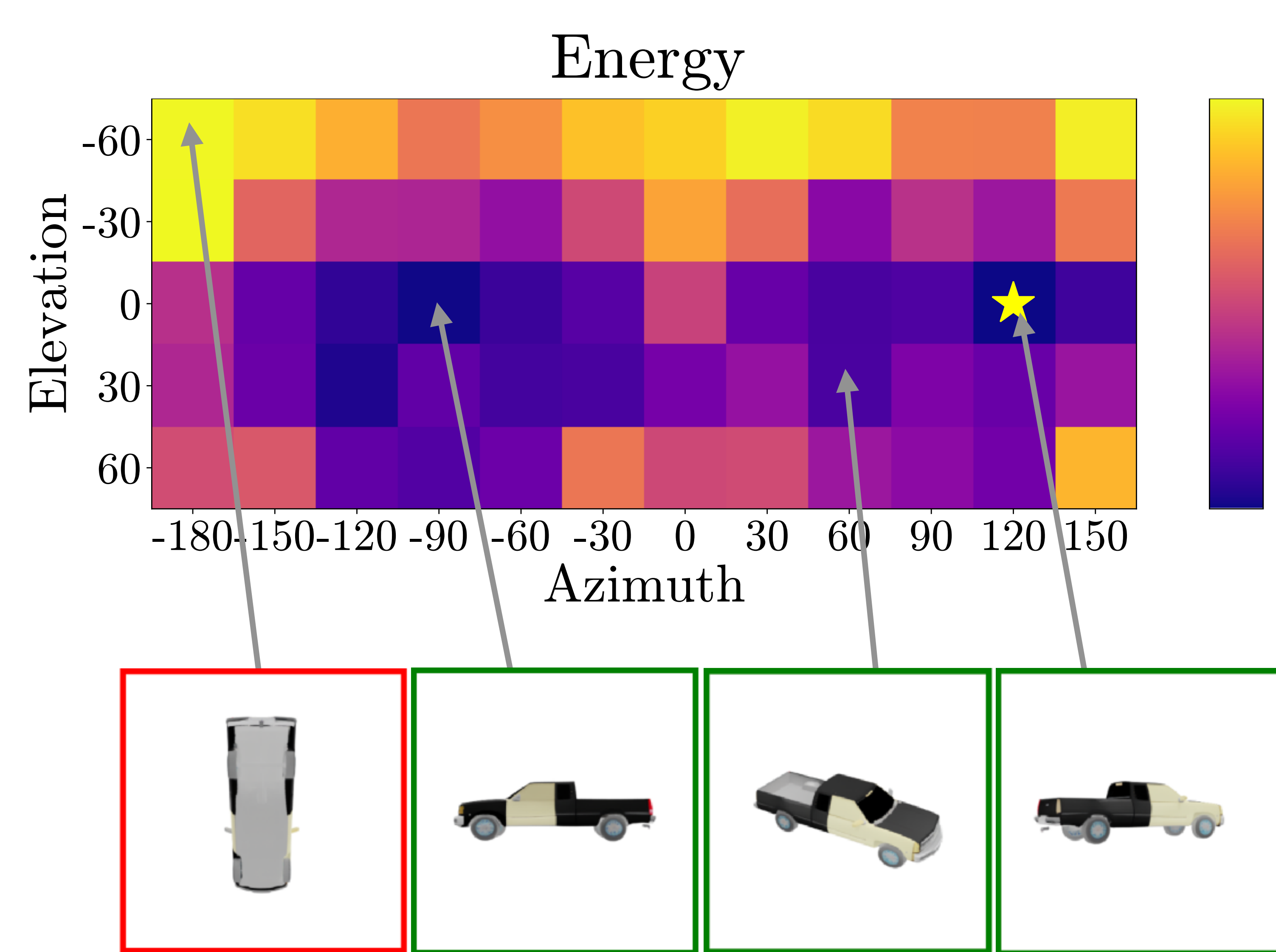
# FoCal Energy Function Picks Typical Viewpoints

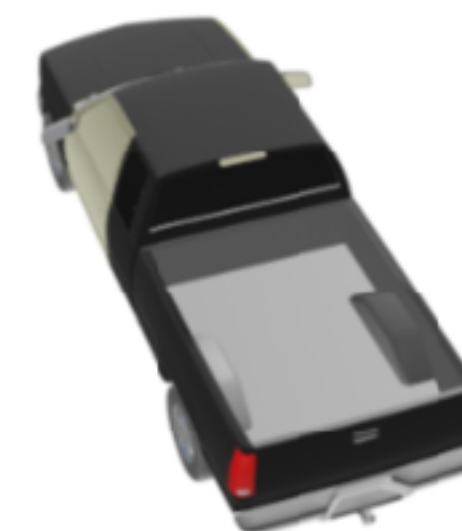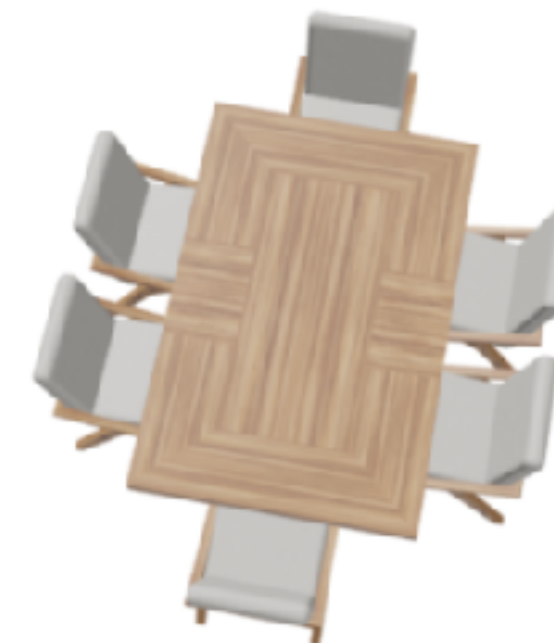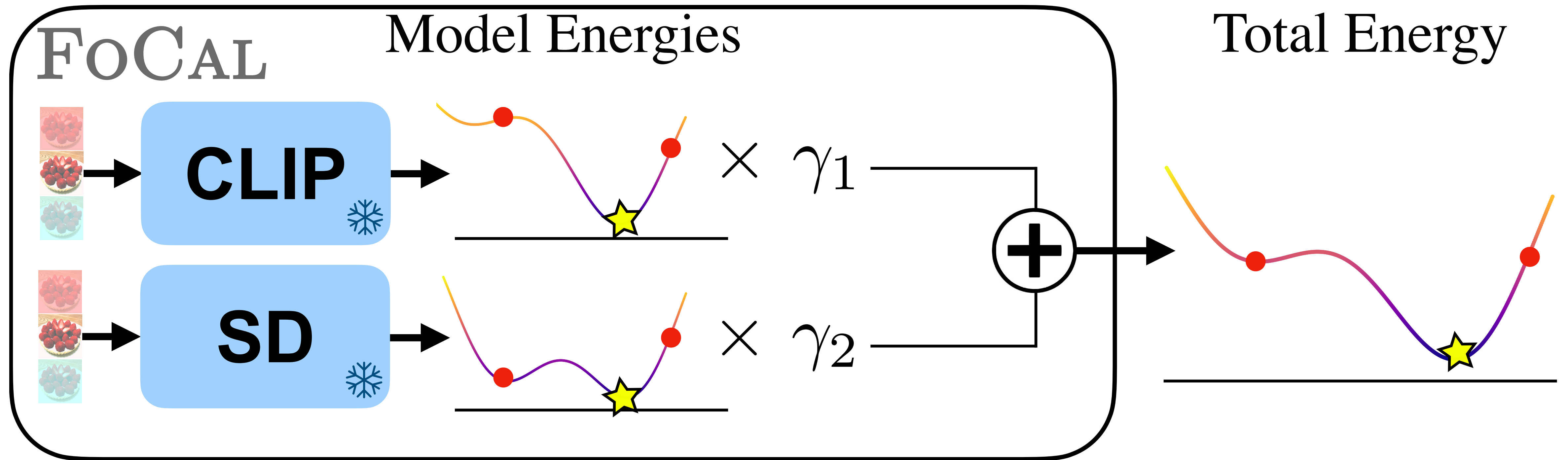# FoCal Energy Function Picks Typical Viewpoints

Energy

Elevation

Azimuth

Original    Canonicalized

# FoCal Energy = CLIP Prior + Stable Diffusion Prior



$$E_{\text{CLIP}}(\boldsymbol{x}; \alpha, \beta) = \left(\alpha \cdot \underset{c \in 1,2,\ldots,|C|}{\text{mean}} -\beta \cdot \max\right)\left(f_\theta(\boldsymbol{x})[c]\right)$$

Classifier Energy

$$E_{\text{diff}}(\boldsymbol{x}) = \frac{1}{T}\sum_{t=1}^{T} \mathbb{E}_{\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0},\mathbf{I})}\left[\|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_\theta(\mathbf{x_t}, t)\|^2\right]$$

Diffusion Energy

# "Distribution Slices": Why FoCal Generalizes Across Transforms

**Transform 1** — **Transform 2** — **Transform 3**

**Insight:** Transformed images form a "slice" of the natural image distribution, and foundation models have already learned a prior over this distribution.
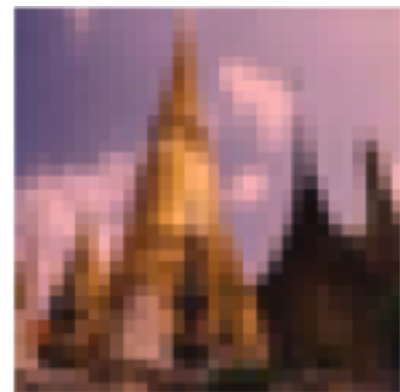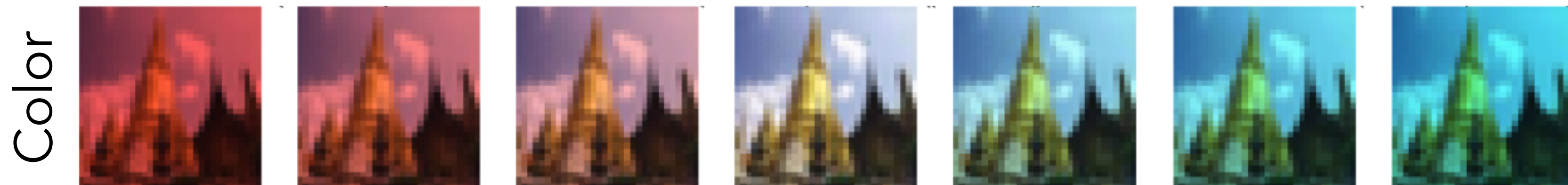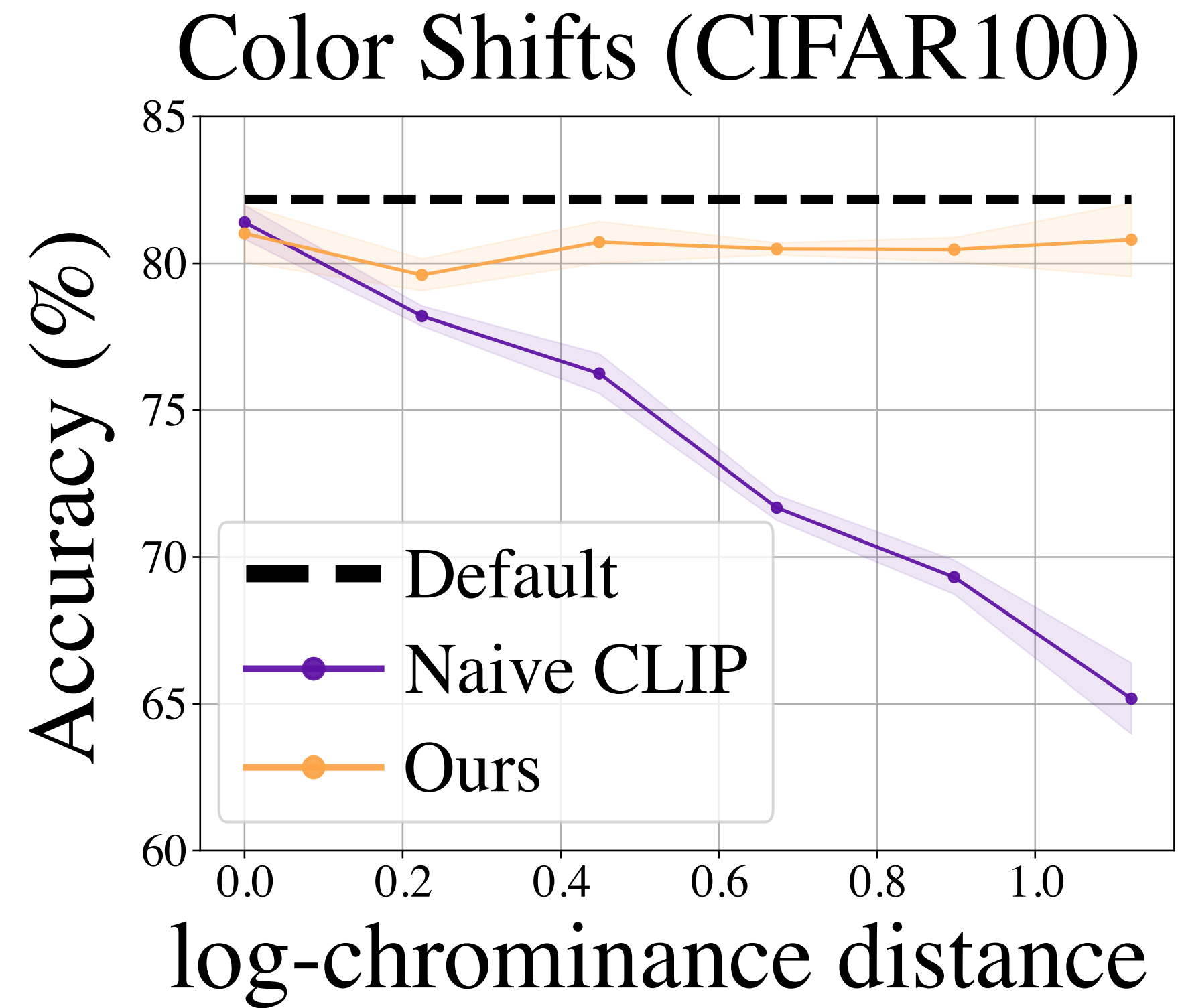
**Significant Improvement on Worst 3D Viewpoints**

Filtered Objaverse-LVIS

+50%

Accuracy

Viewpoint Difficulty Percentile

Finetuned CLIP — Ours

# Significant Robustness Boost for Color and Contrast



Color

# Significant Robustness Boost for Color and Contrast



Color Shifts (CIFAR100)

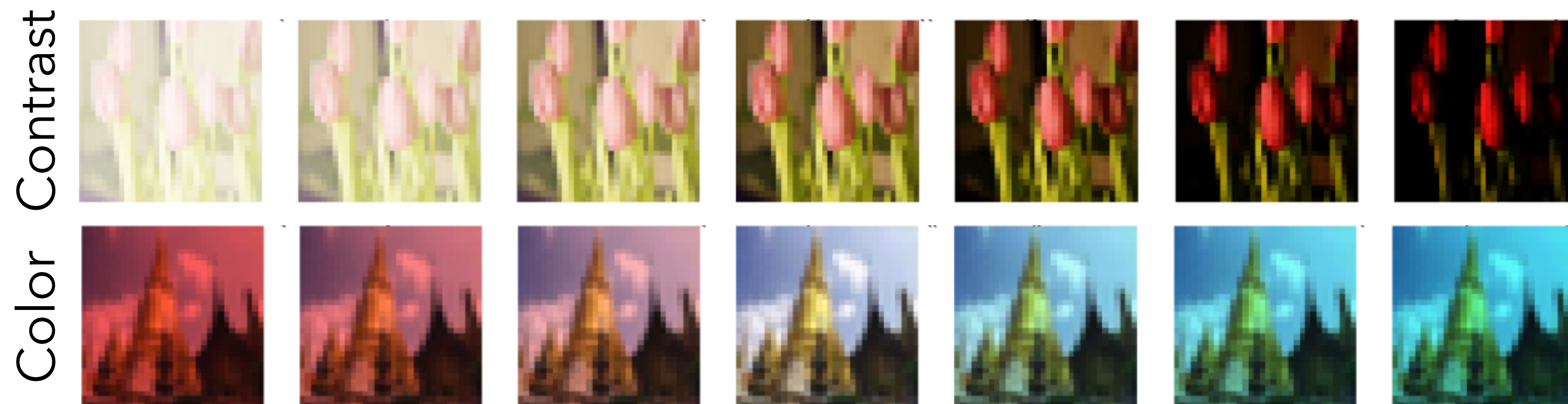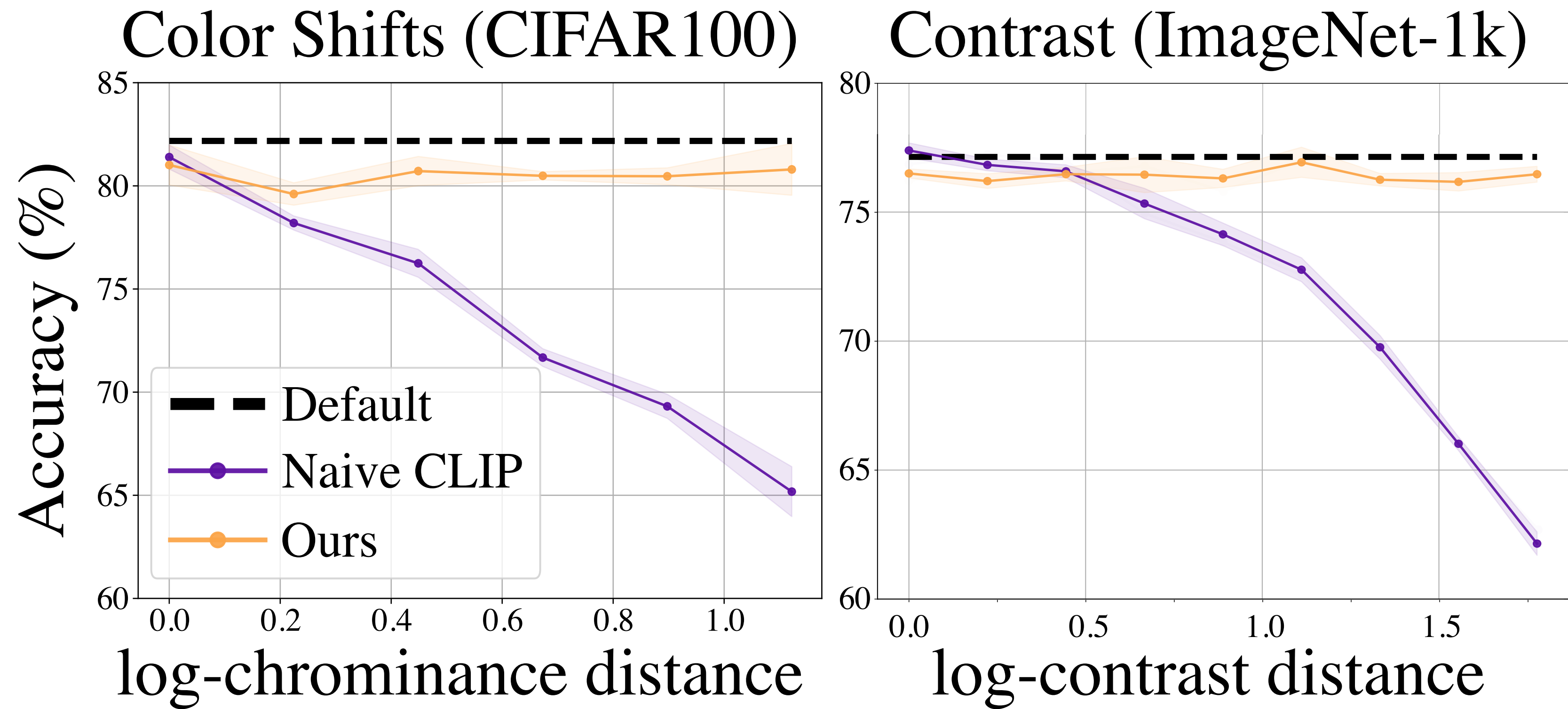Accuracy (%) vs log-chrominance distance

- Default
- Naive CLIP
- Ours

Color

# Significant Robustness Boost for Color and Contrast

Color Shifts (CIFAR100)

# Significant Robustness Boost for Color and Contrast

# Day-Night Results

### Night Image



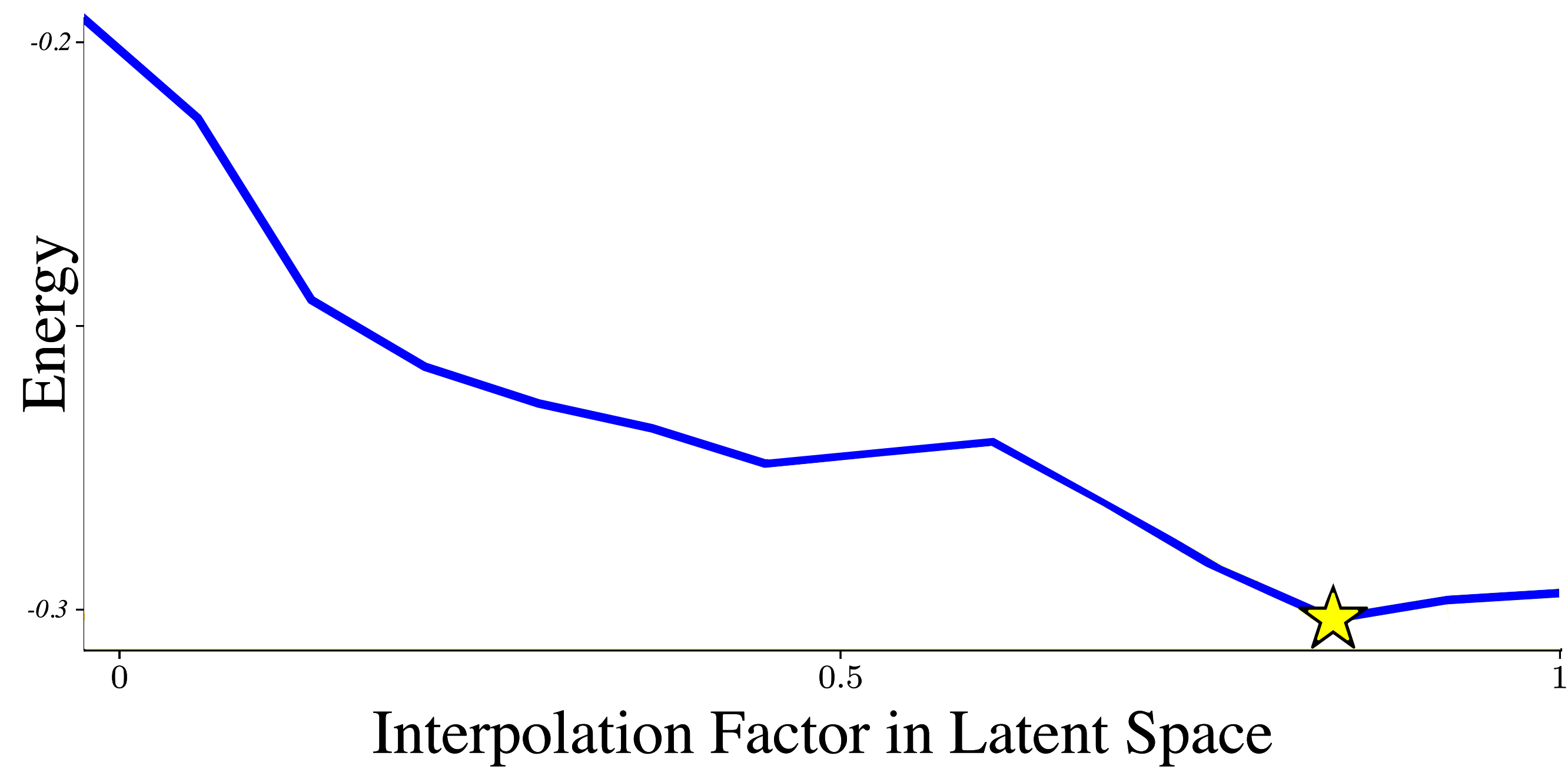### Day Image

# Day-Night Results



Night Image

Day Image

Night                                                                    Day
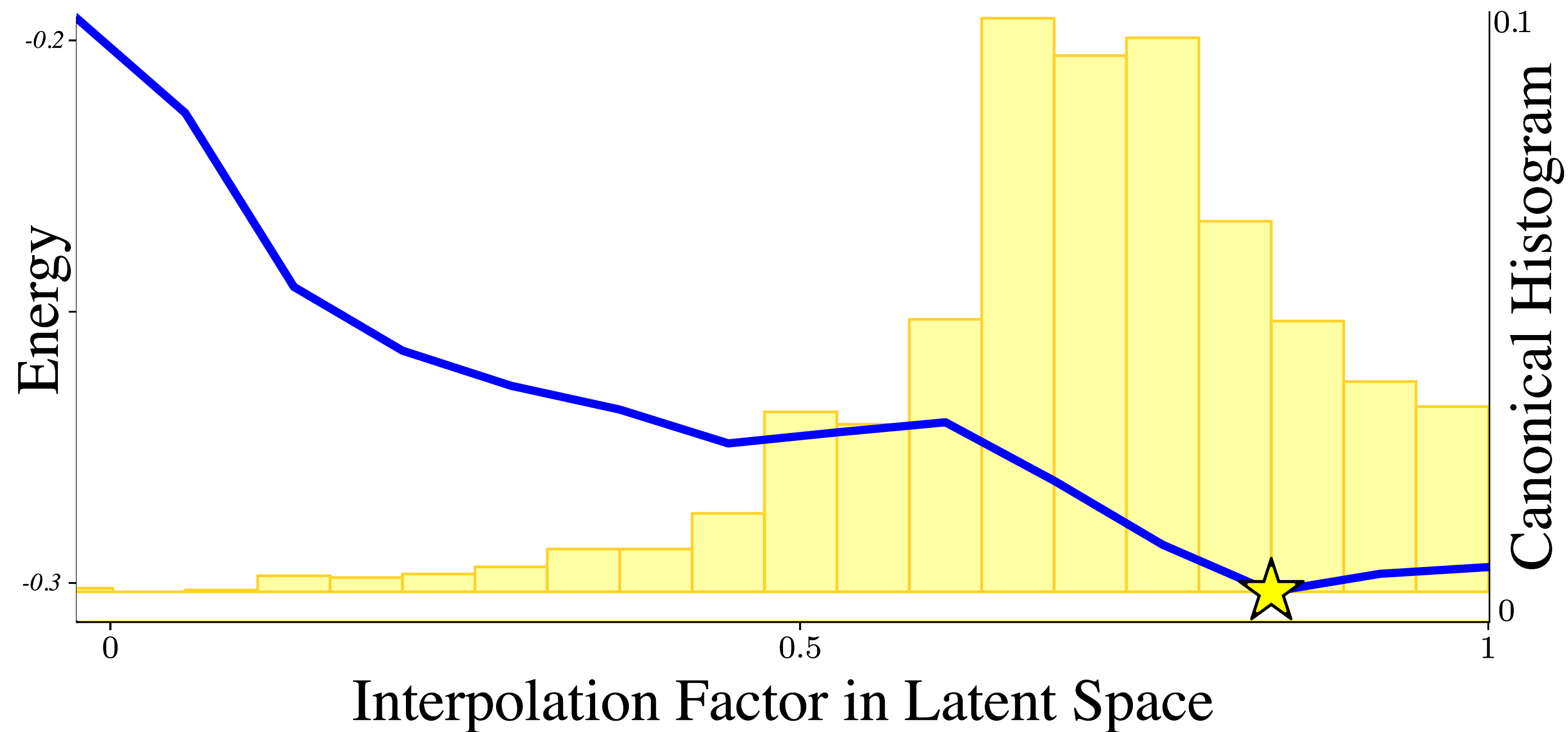
# Day-Night Results



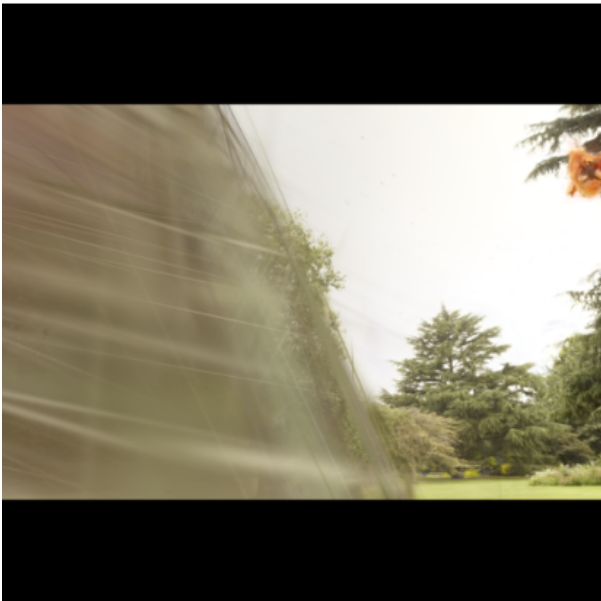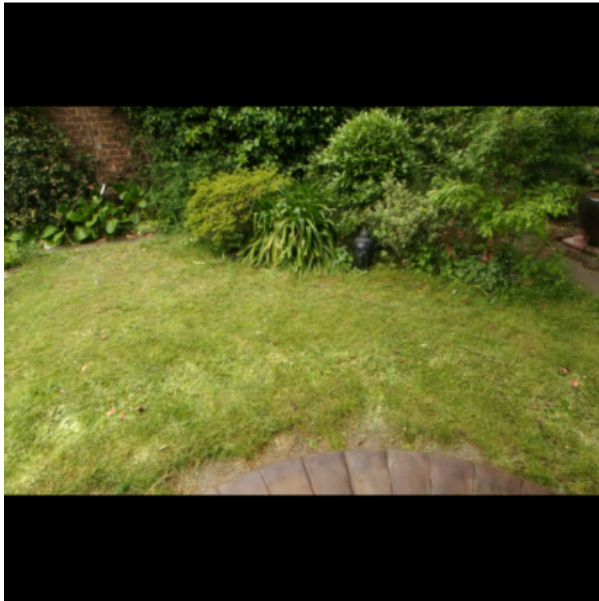Night Image

Day Image

Energy

-0.2

-0.3

0                                          0.5                                          1

Interpolation Factor in Latent Space

Night                                                                                    Day

# Day-Night Results

Night Image



Day Image



Energy

Canonical Histogram

-0.2

0.1

-0.3

0

0

0.5

1

Interpolation Factor in Latent Space

Night

Day

# Day-Night Results



Night Image

Optimized Image

Day Image



Energy

Canonical Histogram

Interpolation Factor in Latent Space

Night

Day
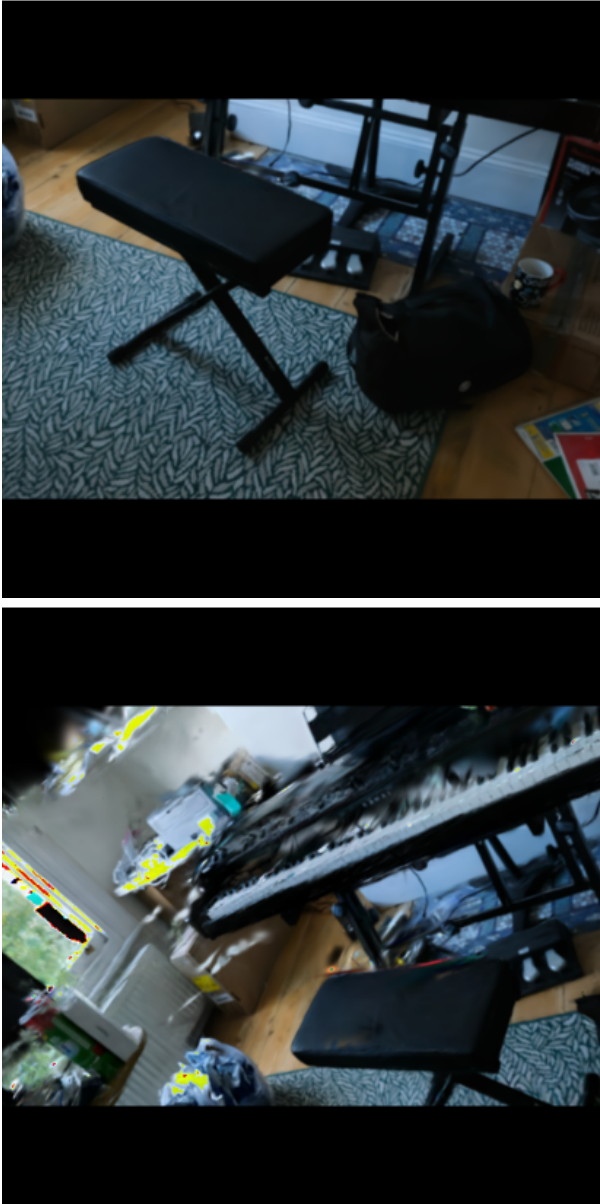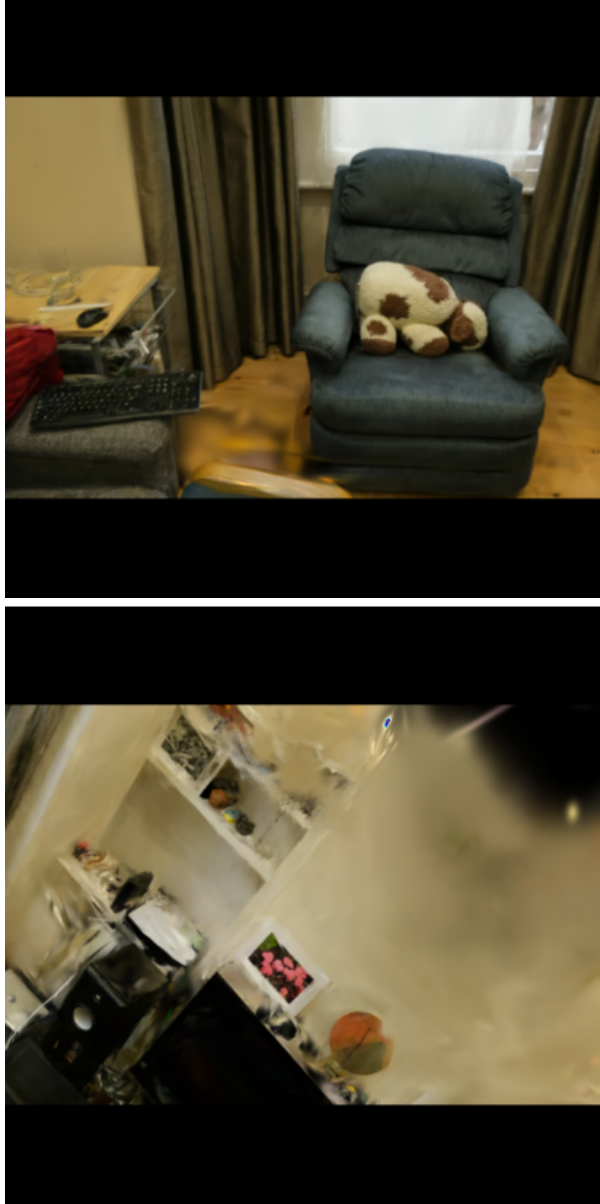
# Active Vision Results



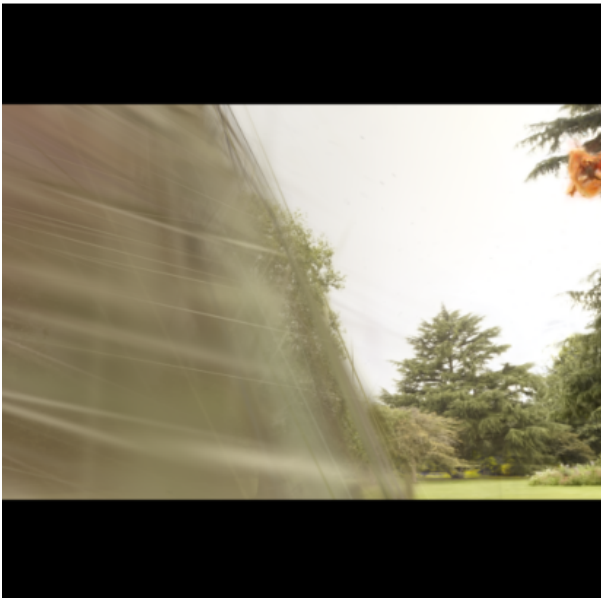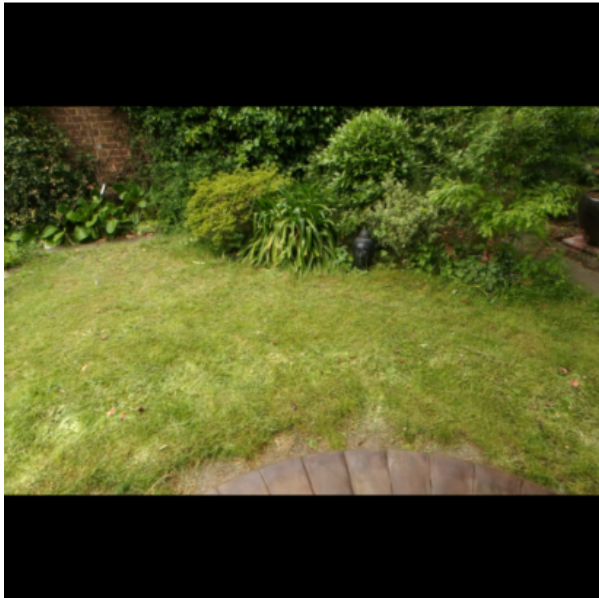Initialization: *flowers*, *garden*, *treehill*, *bicycle*, *bonsai*, *room*, *stump*, *kitchen*

# Active Vision Results

# Active Vision Results

# FoCal: a *Scalable, Data-Driven, Test-Time* Approach to Robust Perception

# Test-Time Canonicalization by Foundation Models for Robust Perception

Utkarsh Singhal*    Ryan Feng*    Stella X. Yu    Atul Prakash

***Summary:*** *Test-time search makes models more robust to natural input variations by converting the varied versions of the input into a 'typical' version.*

## Motivation

### Foundation models (FM) are *brittle*



"chair" ✓    "bench" ✗

✓    ✗

### Prior approaches rely on *transform-specific training*

1. **Data Augmentation**
Train on transformed data

2. **Equivariant NN**
Train with modified architecture

NN → Equivariant NN

3. **Previous Canonicalizers**
Train NN to 'upright' the image

Canonicalizer

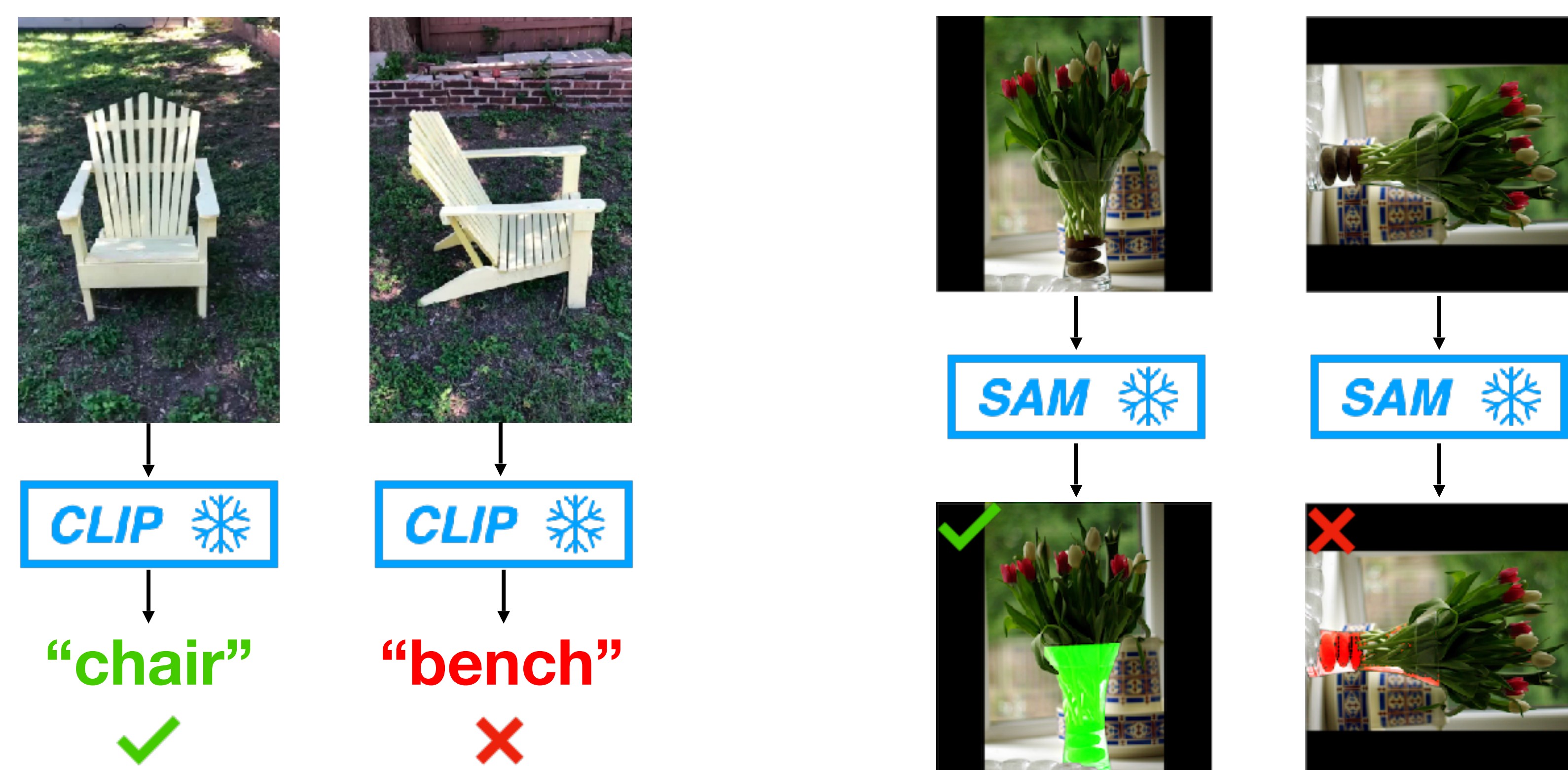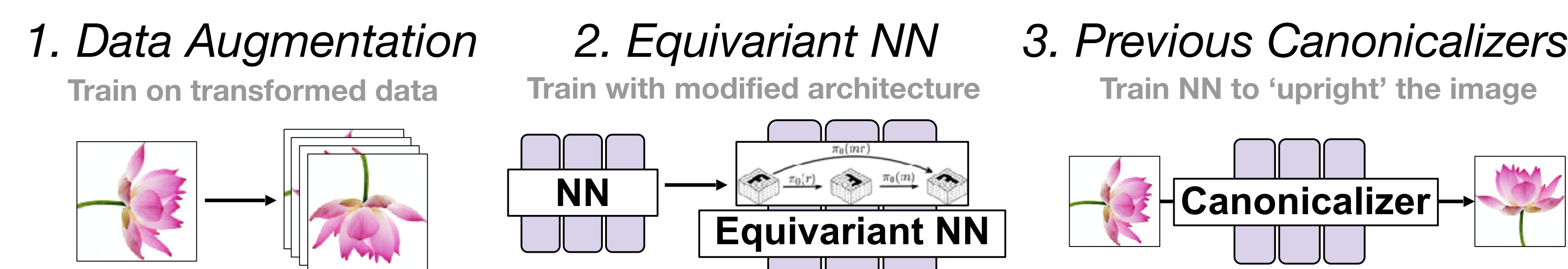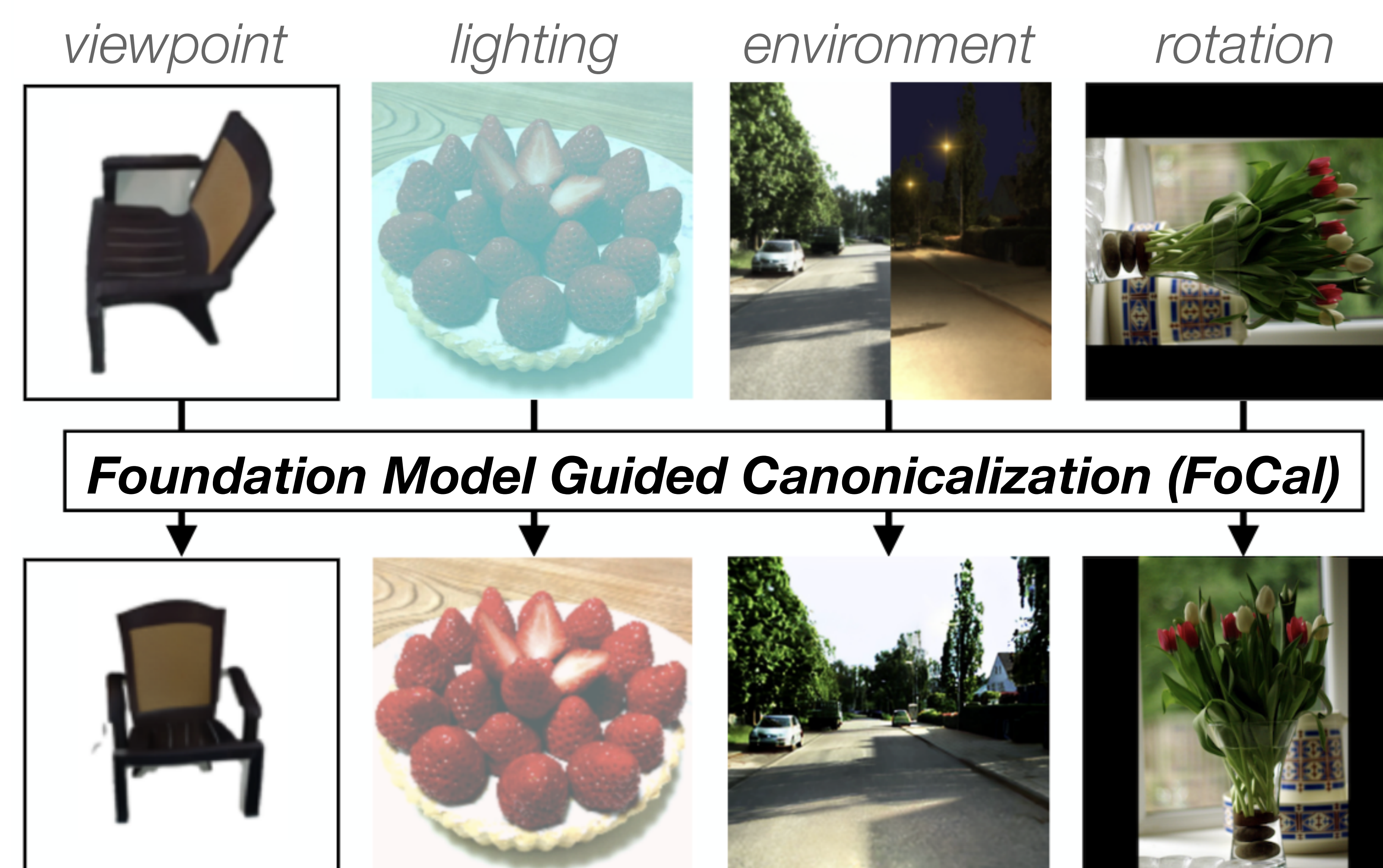***Problem:*** *These methods struggle with OOD input variations at test-time*

## FoCal: Test-Time Canonicalization

**Idea:** transform input to the most 'typical' version

*viewpoint*    *lighting*    *environment*    *rotation*


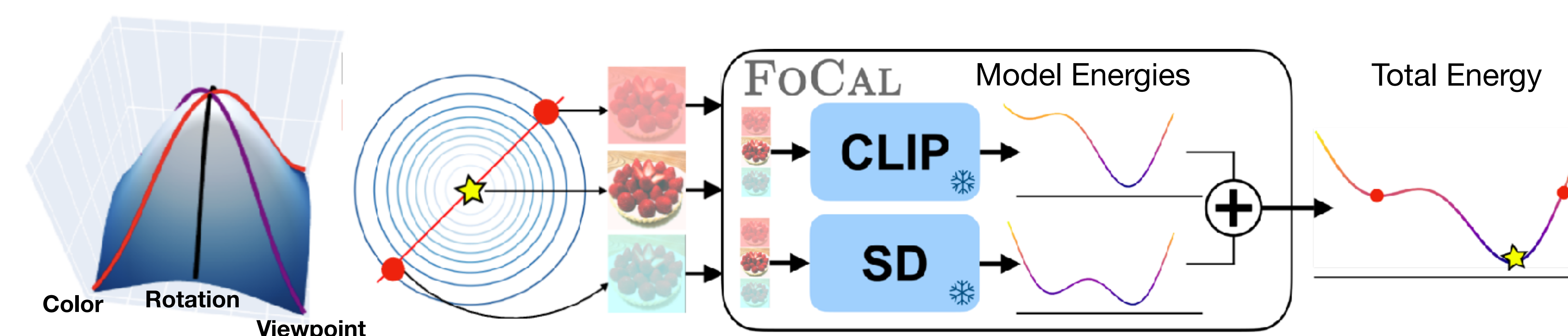
***Foundation Model Guided Canonicalization (FoCal)***

+ *Robustness for many natural variations*
+ *Any downstream task, any model, no training*
+ *Guaranteed invariance for invertible transforms*
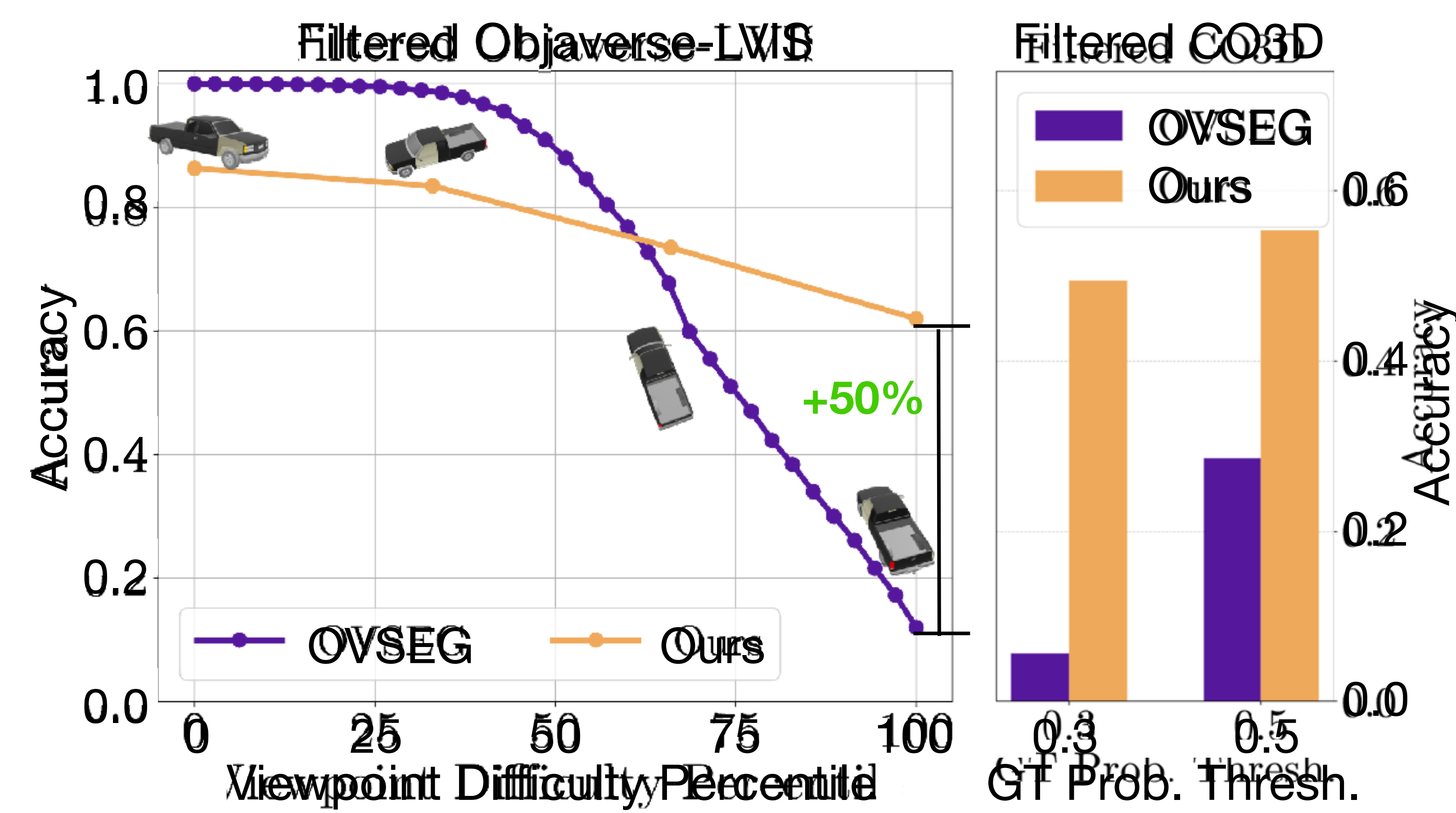
## Approach: Vary & Rank



**Input**    **Canonical**    **FoCal Energy**

Step 1: **Vary**    Step 2: **Rank**
(apply transform)    (CLIP+SD energy)

**Insight:** FM energy estimates the input 'typicality' for many natural variations. Minimizing FM energy over a transform yields robustness.

Color    Rotation    Viewpoint

**Model Energies**    **Total Energy**

FoCal    CLIP    SD

## Results: FoCal Improves Robustness Across Many Domains

### 3D Viewpoint Shifts
*Better on difficult views*



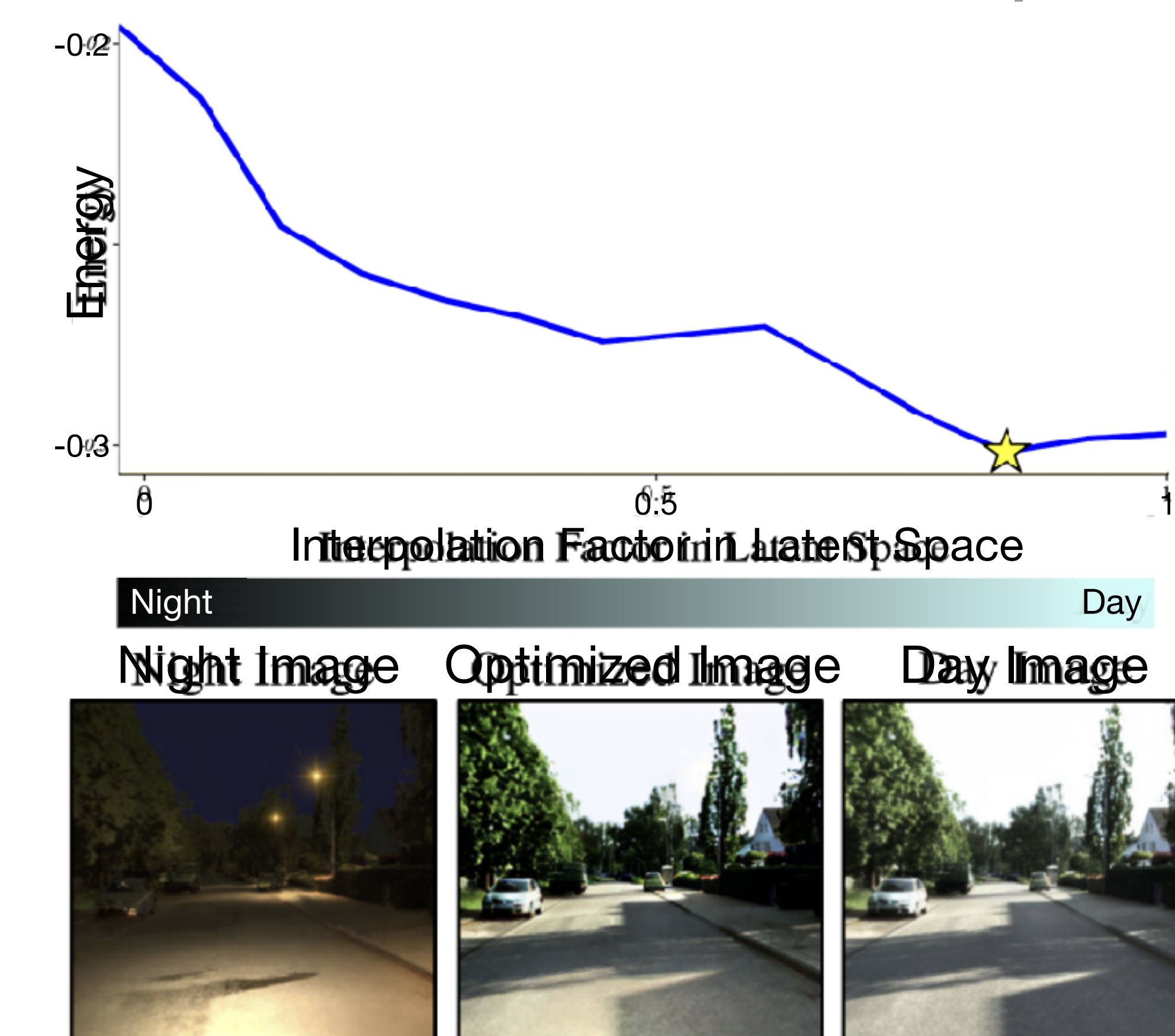Filtered Objaverse-LVIS    Filtered CO3D

OVSEG    Ours

+50%

### Color and Contrast
*Better robustness across color / contrast shifts*



Color Shifts (CIFAR100)    Contrast (ImageNet-1k)

Default    Naive CLIP    Ours

log-chrominance distance    log-contrast distance

Color    Contrast

### Day-Night
*Canonicalization in SD latent space*



Interpolation Factor in Latent Space

Night    Day

Night Image    Optimized Image    Day Image

### Active Vision (Exploring Virtual Environment)
*FoCal looks at salient objects in upright poses*



*garden*    *bicycle*

Initialization

Random Pose

Optimized