# Exploring Blind Spots of Vision Models

Sriram Balasubramanian*   Gaurang Sriramanan*   Vinu Sankar Sadasivan   Soheil Feizi
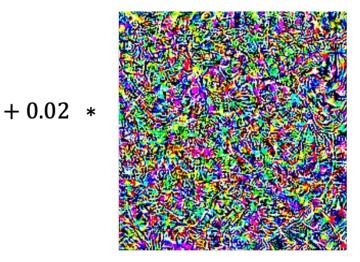
# Introduction

- Input **over-sensitivity** well studied in adversarial literature



Prediction: **Hamster**

Confidence = 99.99%

$+\ 0.02\ *$

50-step PGD targeted attack with $\varepsilon = \frac{8}{255}$ scaled by 50x

$=$

Prediction: **Banjo**

Confidence = 100%

# Introduction

- Input **_over-sensitivity_** well studied in adversarial literature

- We study input **_under-sensitivity_** for general models

- Uncover extent of excessive invariance in common vision models?

Over-Sensitivity

$$||f(x) - f(x')|| \quad \uparrow$$

$$x \approx x'$$

Under-Sensitivity

$$f(x) \approx f(x')$$

$$||x - x'|| \quad \uparrow$$

# Mathematical Preliminaries

- For $g : \mathbb{R}^d \to \mathbb{R}, L_g(c) = \{x \in \mathcal{X} : g(x) = c\}$ is called the Level Set

- Important Property: For any curve in the Level Set $\gamma(t) : [0,1] \to L_g(c)$

$$\frac{d}{dt}(g(\gamma(t))) = 0 = \langle \nabla g(\gamma(t)), \gamma'(t) \rangle$$

**Lemma 1.** *If $g : \mathbb{R}^d \to \mathbb{R}$ is a continuously differentiable function, then each of its regular level sets is an $(d-1)$ dimensional submanifold of $\mathbb{R}^d$.*

- How expansive are these submanifolds for common ML models?

# Can we Traverse along Level Sets?

goose

Source Image $x_s$

Confidence for class "goose" = 0.997
Confidence for class "Scottish Terrier" = 0

Scottish Terrier

Target Image $x_t$

Confidence for class "goose" = 0
Confidence for class "Scottish Terrier" = 1.0

# Level Set Traversal (LST) Algorithm

Repeat until Max Iterations

**Compute Input Gradient**

$$\Delta \boldsymbol{x} = \boldsymbol{x}_t - \boldsymbol{x}$$
$$\boldsymbol{g} = \nabla_{\boldsymbol{x}} CE(f(\boldsymbol{x}), y)$$

**Compute Orthogonal Projection**

$$c_{//} = (\boldsymbol{g} \cdot \Delta \boldsymbol{x}) / ||\boldsymbol{g}||^2$$
$$\Delta \boldsymbol{x}_\perp = \eta(\Delta \boldsymbol{x} - c_{//} \boldsymbol{g})$$

**Update Image**

$$\boldsymbol{x}_{||} = \Pi_\infty(\boldsymbol{x}_{||} - \epsilon \boldsymbol{g}, -\epsilon, \epsilon)$$
$$\boldsymbol{x}_{\text{new}} = \boldsymbol{x} + \Delta \boldsymbol{x}_\perp + \boldsymbol{x}_{||}$$

**Verify Model Confidence**

**if** $f(\boldsymbol{x}_s)[j] - f(\boldsymbol{x}_{\text{new}})[j] > \delta$ **then**
    **return** $\boldsymbol{x}$
$\boldsymbol{x} = \boldsymbol{x}_{\text{new}}$

# LST Path in Input Space for ResNet-50



goose

Source Image $\boldsymbol{x}_s$

Confidence for class "goose" = 0.997
Confidence for class "Scottish Terrier" = 0

Scottish Terrier

Target Image $\boldsymbol{x}_t$

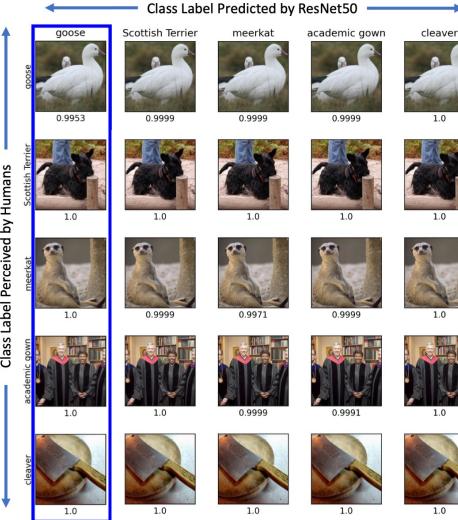Confidence for class "goose" = 0
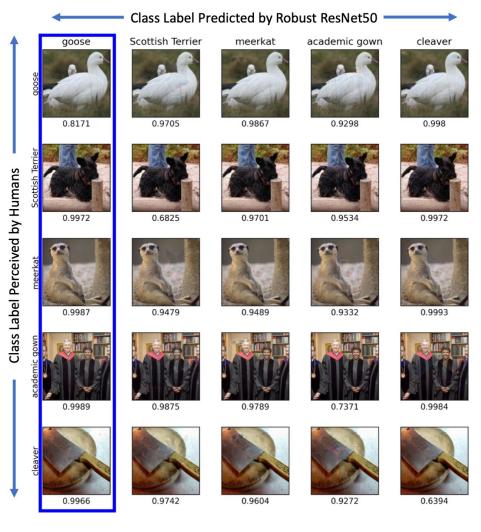Confidence for class "Scottish Terrier" = 1.0

| Step: 0 | Step: 10 | Step: 20 | Step: 40 | Step: 60 | Step: 80 | Step: 120 | Step: 160 | Step: 200 | Step: 300 | Step: 400 |
|---|---|---|---|---|---|---|---|---|---|---|
| Confidence: 0.997 | Confidence: 0.998 | Confidence: 0.999 | Confidence: 0.999 | Confidence: 1.0 | Confidence: 1.0 | Confidence: 1.0 | Confidence: 1.0 | Confidence: 1.0 | Confidence: 1.0 | Confidence: 1.0 |

LST Blind Spots

# LST over arbitrary Source-Target pairs



Normally Trained ResNet-50

Adversarially Trained ResNet-50

# Star-like Substructure of Level Sets



$$\{\boldsymbol{x} \in \mathcal{X} : f^{\mathrm{goose}}(\boldsymbol{x}) \geq f^{\mathrm{goose}}(\boldsymbol{x_s}) - \delta\}$$

Scottish Terrier

goose

$\boldsymbol{x_s}$

academic gown

meerkat

# Star-like Substructure of Level Sets



Normally Trained ResNet-50

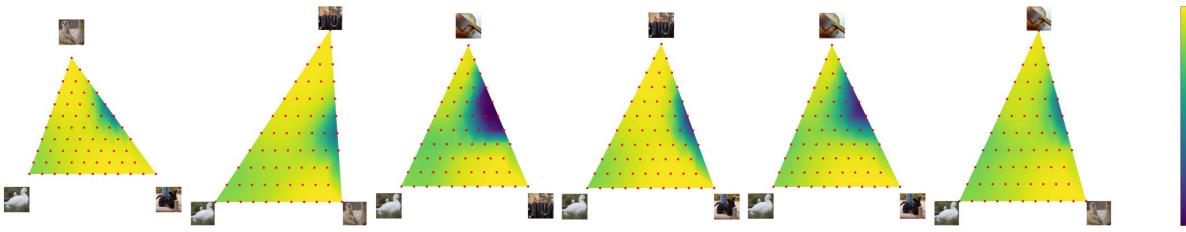Adversarially Trained ResNet-50

# Star-like Substructure of Level Sets

Normally Trained ResNet-50:



Adversarially Trained ResNet-50:

# Quantitative Analysis of Blind Spot Invariance

Distance metrics:

1. RMSE:

   Root mean squared error

2. Max norm ($\ell_\infty$):

   Maximum absolute difference

3. SSIM:

   Structural Similarity Index

4. LPIPS:

   Perceptual Image Similarity

Confidence metrics:

1. Source confidence ($p_{src}$):

   Confidence of the model for the source image

2. Average path confidence

   Mean confidence over the linear paths connecting the source image to LST outputs

3. Average Δ confidence:

   Mean confidence over the enclosed triangle

4. Average Δ fraction for a given $\delta$:

   Fraction of triangle over which confidence is at least $p_{src} - \delta$

# Quantitative Analysis of Blind Spot Invariance

Table 1: Quantitative image distance metrics between output of Level Set Traversal and target images.

| Models | RMSE : $\mu \pm \sigma$ | $\ell_\infty$ dist: $\mu \pm \sigma$ | SSIM: $\mu \pm \sigma$ | LPIPS dist: $\mu \pm \sigma$ |
|---|---|---|---|---|
| ResNet-50 (Normal) | $0.008 \pm 0.001$ | $0.046 \pm 0.020$ | $0.990 \pm 0.021$ | $0.002 \pm 0.004$ |
| ResNet-50 (AT) | $0.029 \pm 0.008$ | $0.746 \pm 0.124$ | $0.915 \pm 0.041$ | $0.057 \pm 0.037$ |
| DeiT-S (Normal) | $0.011 \pm 0.002$ | $0.116 \pm 0.030$ | $0.973 \pm 0.024$ | $0.024 \pm 0.017$ |
| DeiT-S (AT) | $0.046 \pm 0.010$ | $0.821 \pm 0.117$ | $0.898 \pm 0.041$ | $0.219 \pm 0.068$ |

Table 2: Quantitative confidence metrics over the triangular convex hull ($\Delta$) of a given source image and two target LST blindspot image-pairs and over linear interpolant paths between source and blindspot images. (For reference, a random classifier would have confidence of 0.001)

| Models | $p_{\text{src}}$ $(\mu \pm \sigma)$ | Avg $\Delta$ Conf. $(\mu \pm \sigma)$ | Avg $\Delta$ Frac. $(\mu \pm \sigma)$ | | | | Avg Path Conf. $(\mu \pm \sigma)$ |
|---|---|---|---|---|---|---|---|
| | | | $\delta = 0.0$ | $\delta = 0.1$ | $\delta = 0.2$ | $\delta = 0.3$ | |
| ResNet-50 (Normal) | $0.99 \pm 0.02$ | $0.56 \pm 0.10$ | $0.13 \pm 0.15$ | $0.51 \pm 0.11$ | $0.53 \pm 0.1$ | $0.54 \pm 0.10$ | $0.96 \pm 0.05$ |
| ResNet-50 (AT) | $0.88 \pm 0.11$ | $0.83 \pm 0.09$ | $0.49 \pm 0.29$ | $0.79 \pm 0.13$ | $0.85 \pm 0.1$ | $0.88 \pm 0.09$ | $0.93 \pm 0.06$ |
| DeiT-S (Normal) | $0.85 \pm 0.06$ | $0.68 \pm 0.05$ | $0.54 \pm 0.11$ | $0.67 \pm 0.06$ | $0.71 \pm 0.06$ | $0.73 \pm 0.06$ | $0.94 \pm 0.02$ |
| DeiT-S (AT) | $0.76 \pm 0.08$ | $0.59 \pm 0.07$ | $0.20 \pm 0.09$ | $0.43 \pm 0.14$ | $0.63 \pm 0.15$ | $0.76 \pm 0.12$ | $0.73 \pm 0.06$ |

# Conclusions

- Using LST, we find that the level sets of common vision models is **remarkably expansive**

- The **linear** path from any given source image to LST blind spot outputs retain **high model confidence** throughout for arbitrary targets

- This unveils a **star-like substructure** for the equi-confidence level sets of common models

- Adversarially trained models are significantly more **under-sensitive,** over inputs **well beyond** the original threat model