



**Adversarially trained neural representations
may already be as robust as corresponding
biological neural representations**

Zhiwei Bai

Shanghai Jiao Tong University

2022-12-04

About the author?

Chong Guo¹ Michael J. Lee^{1 2 3} Guillaume Leclerc⁴ Joel Dapello^{1 2 5} Yug Rao⁶ Aleksander Madry^{4 7}
James J. DiCarlo^{1 2 3}

James J. DiCarlo

- Born in 1968, neuroscientist
- Peter de Florez Professor, Brain and Cognitive Sciences
- Director, MIT Quest for Intelligence
- **"Aim to understand how a complex network of brain regions underlies our ability to recognize vast numbers of objects and faces rapidly."**



Outline

1. **Motivation:** why do we care about adversarial robustness?
2. **Result:** who is more robust between AT-DNNs and primate visual perception?
3. **Method:** how to measure adversarial sensitivity of IT neural sites?

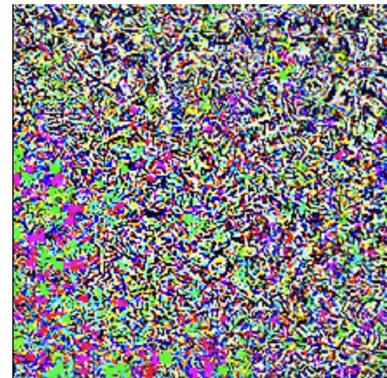
1. Motivation: why do we care about adversarial robustness?

Motivation: DNN is **VERY** brittle

- Pre-trained ResNet50



+



=



?

Hog (猪)

High confidence (0.996)

Wombat (袋熊)

High confidence (0.999)



Some specific perturbations can easily fool current deep neural networks.

How to create adversarial examples?

- Given data $S = \{(\mathbf{x}_j, \mathbf{y}_j)\}_{j=1}^n$, usual training goal:

$$\min_{\theta} R_S(\theta) = \mathbb{E}_S \ell(f_{\theta}(\mathbf{x}), \mathbf{y}) := \frac{1}{n} \sum_{j=1}^n \ell(f_{\theta}(\mathbf{x}_j), \mathbf{y}_j)$$

- Create adversarial examples: for a fixed $(\mathbf{x}, \mathbf{y}) \in S$,

$$\max_{\|\boldsymbol{\delta}\| < \epsilon} \ell(f_{\theta}(\mathbf{x} + \boldsymbol{\delta}), \mathbf{y})$$

- Targeted attack: $\max_{\|\boldsymbol{\delta}\| < \epsilon} (\ell(f_{\theta}(\mathbf{x} + \boldsymbol{\delta}), \mathbf{y}) - \ell(f_{\theta}(\mathbf{x} + \boldsymbol{\delta}), \mathbf{y}_{\text{target}}))$

What is adversarial training?

- Usual training goal:

$$\min_{\theta} R_S(\theta) = \mathbb{E}_S \ell(f_{\theta}(x), y) := \frac{1}{n} \sum_{j=1}^n \ell(f_{\theta}(x_j), y_j)$$

- Adversarial training goal:

$$\min_{\theta} \hat{R}_S(\theta) = \mathbb{E}_S \max_{\|\delta\| < \epsilon} \ell(f_{\theta}(x + \delta), y) := \frac{1}{n} \sum_{j=1}^n \max_{\|\delta_j\| < \epsilon} \ell(f_{\theta}(x_j + \delta_j), y_j)$$

- $R_S(\theta) \leq \hat{R}_S(\theta)$, intuitively $\hat{R}_S(\theta)$ is the **worst-case**.

Is adversarial training enough?

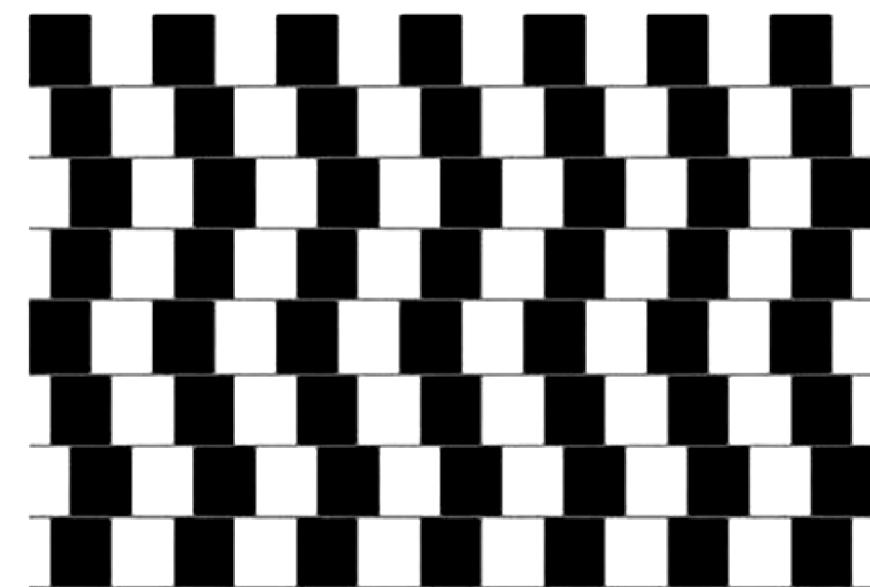
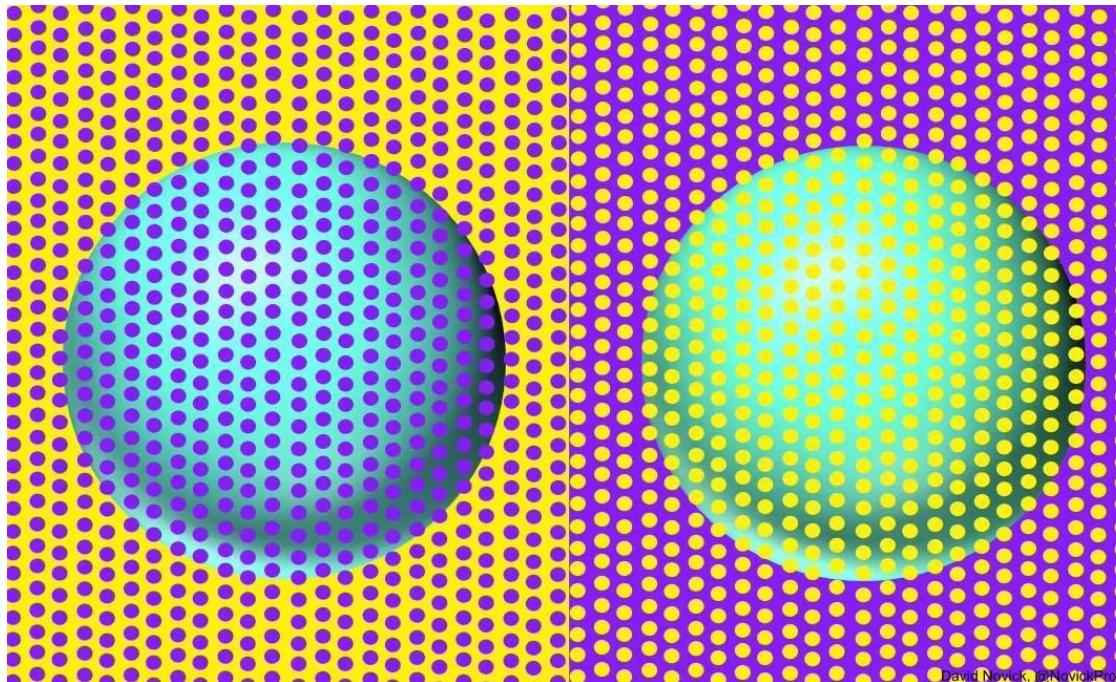
- Gold standard of robust perception: visual systems of primates

Robustness of the best of artificial NNs

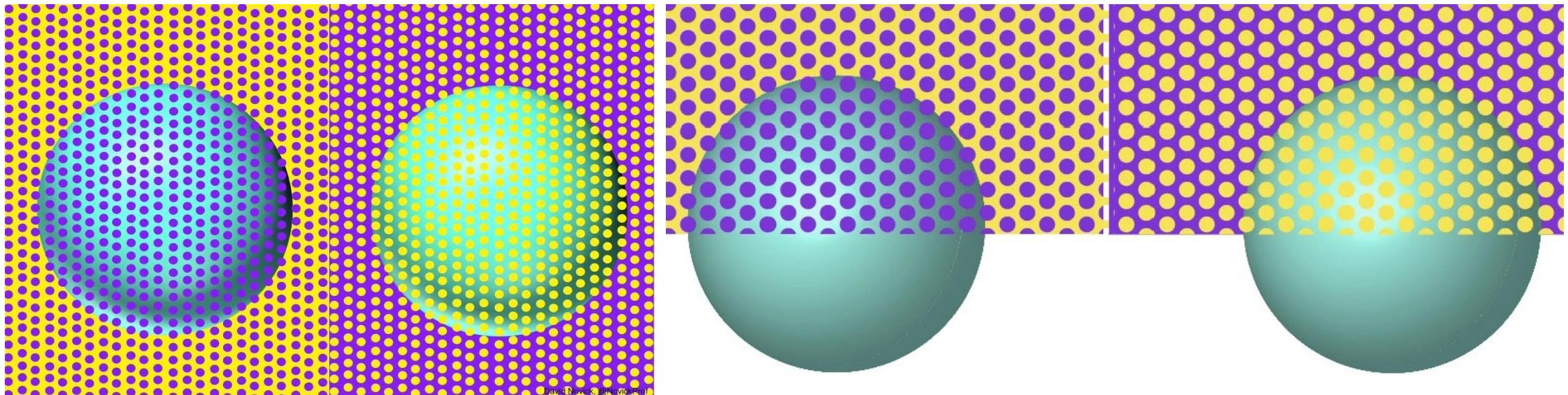


Robustness of visual systems of primates

Do visual systems of primates really robust?



Do visual systems of primates really robust?



Are these adversarial examples in primates widespread or only for a few specific tasks?

In particular, do adversarial examples exist for cat and dog classification tasks?

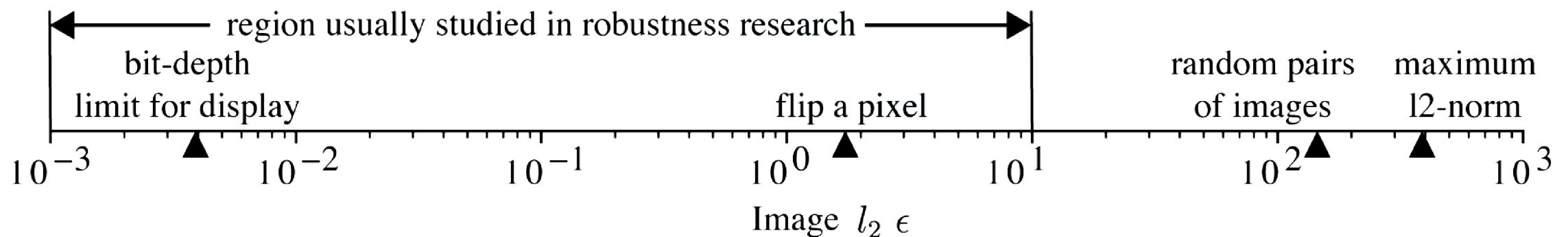
To answer the question, we need answer:

1. How to set ϵ in $\|\delta\| < \epsilon$?
2. "Worst-case" relies on detailed knowledge. How to create adversarial examples?

$$\max_{\|\delta\| < \epsilon} \ell(r(x + \delta), y)$$

How to set ϵ in $\|\delta\| < \epsilon$?

- Consider l_2 norm $\|\delta\|_2 < \epsilon$



- Restrict to a specific regime, such as $[1, 10]$

"Worst-case" relies on detailed knowledge. How to create adversarial examples?

$$\ell(\delta^*) := \max_{\|\delta\| < \epsilon} \ell(r(x + \delta), y)$$

1. Black-box attack: without using detailed knowledge.

- Rely on **random sampling** image perturbation directions,
"unlikely to yield good estimates of adversarial sensitivity".

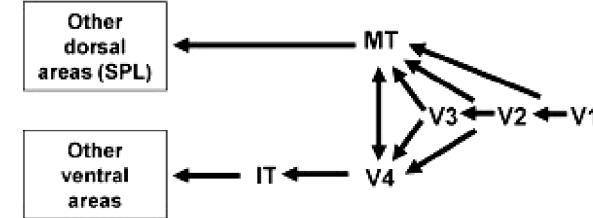
2. Build a "white-box" model to **estimate** the adversarial example.

- This paper develops an experimental method to do this.
- **Lower bound** is enough for this paper's claim.

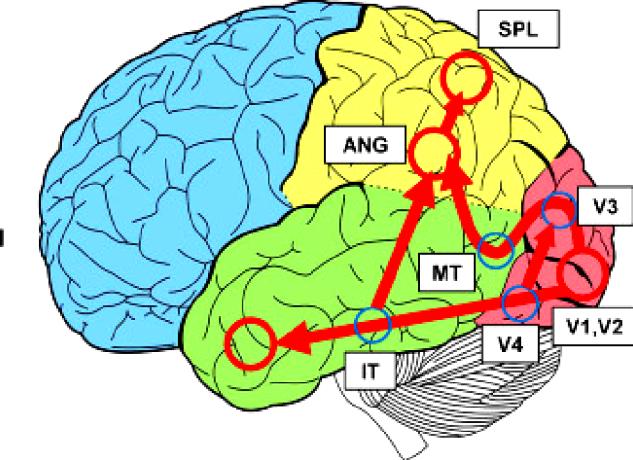
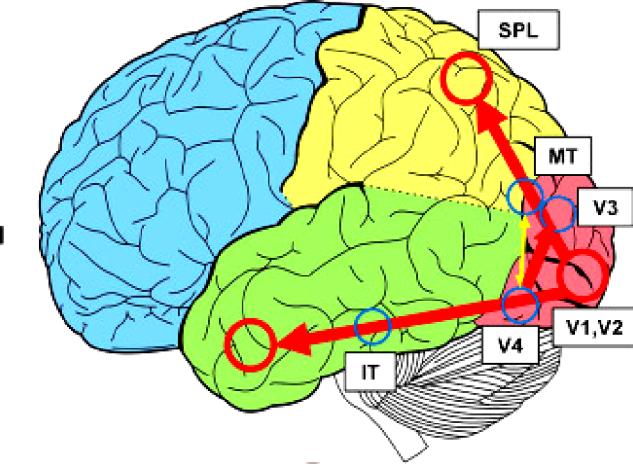
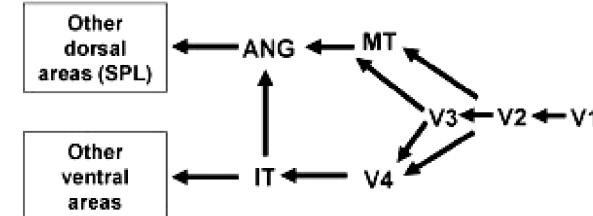
Which area of neurons to choose?

- V1 (Primary Visual cortex)
 - Support: Primates is more robust
- IT (Inferior Temporal cortex)
 - Support: Artificial NN is more robust
 - Compare with the penultimate layer of DNN

CLASSICAL MODEL



MODIFIED MODEL



How to quantify robustness?

- Create adversarial examples: for a fixed $(\mathbf{x}, \mathbf{y}) \in S$,

$$\ell(\boldsymbol{\delta}^*) := \max_{\|\boldsymbol{\delta}\| < \epsilon} \ell(r(\mathbf{x} + \boldsymbol{\delta}), \mathbf{y})$$

- Quantify sensitivity: for a fixed $\mathbf{x} \in S_x$ and a fixed i^{th} neural site,

$$s_i(\mathbf{x}, \epsilon) := \max_{\|\boldsymbol{\delta}\|_2 < \epsilon} |r_i(\mathbf{x} + \boldsymbol{\delta}) - r_i(\mathbf{x})|$$

Why not $\max_{\|\boldsymbol{\delta}\|_2 < \epsilon} \|r(\mathbf{x} + \boldsymbol{\delta}) - r(\mathbf{x})\|_1$?

Quantify robustness: individual unit level

- Quantify sensitivity: for a fixed $\mathbf{x} \in S_x$ and a fixed i^{th} neural site,

$$s_i(\mathbf{x}, \epsilon) := \max_{\|\boldsymbol{\delta}\|_2 < \epsilon} |r_i(\mathbf{x} + \boldsymbol{\delta}) - r_i(\mathbf{x})|$$

- Marginalizing the image distribution S_x :

$$s_i(\epsilon) := \mathbb{E}_{\mathbf{x} \sim S_x}[s_i(\mathbf{x}, \epsilon)]$$

- Normalized adversarial sensitivity:

$$\tilde{s}_i(\epsilon) = \frac{s_i(\epsilon)}{\sigma_i}, \text{ where } \sigma_i = (\text{Var}_{\mathbf{x} \sim S_x} r_i(\mathbf{x}))^{\frac{1}{2}}$$

2. Result: who is more robust between AT-DNNs and primate visual perception?

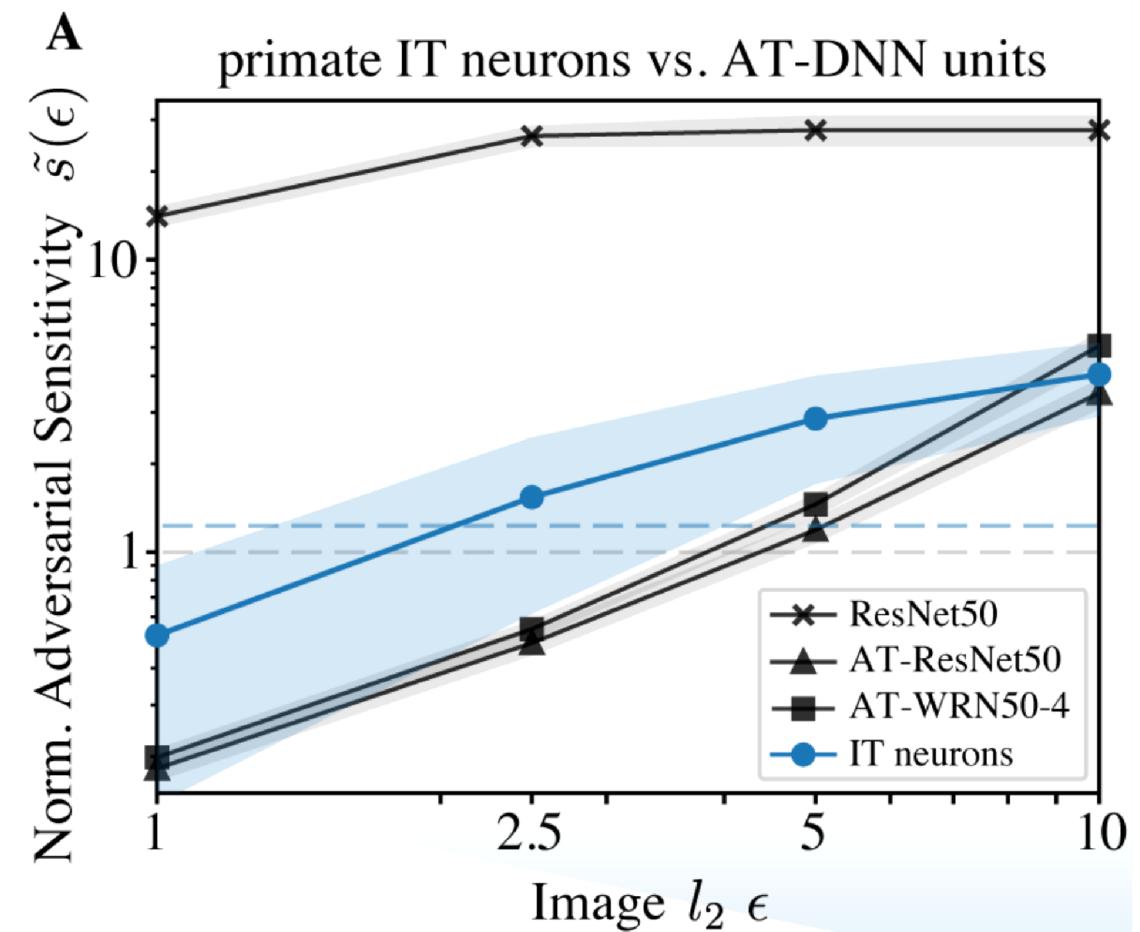
Who is more robust?

- The number of IT sites $m = 21$

$$\tilde{s}(\epsilon) = \frac{1}{m} \sum_{i=1}^m \tilde{s}_i(\epsilon)$$

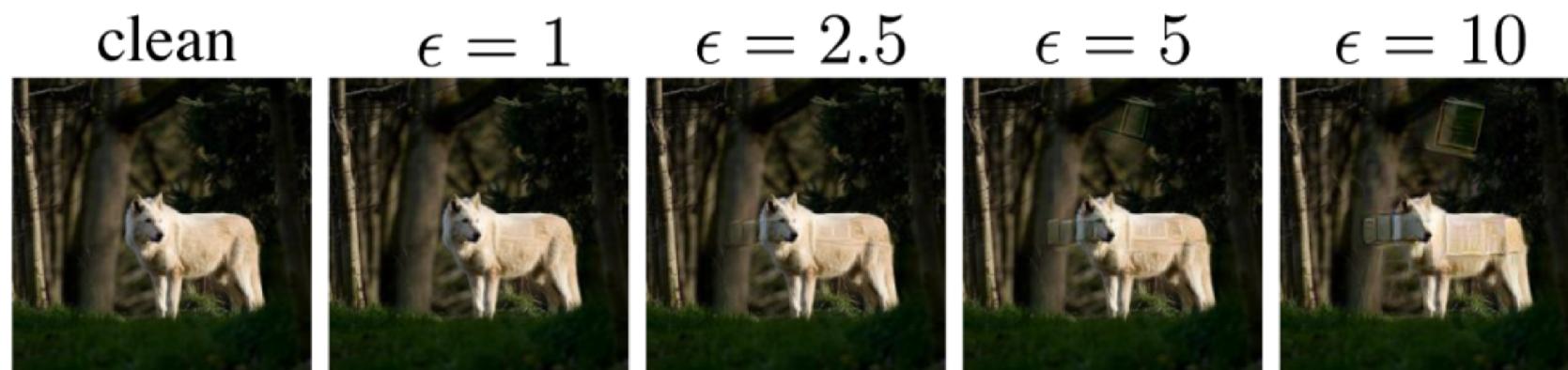
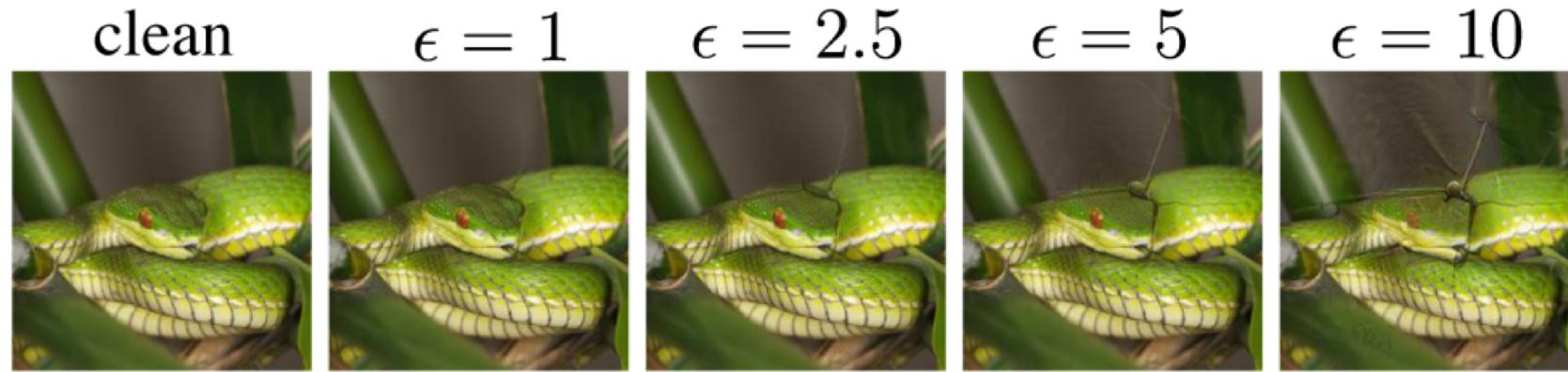
- Adversarially trained NN ($l_2\epsilon = 3$), 10-fold smaller
- Grey dashed: standard deviation
- Blue dashed: random pairs of images

$$\frac{1}{C_n^2} \sum_{j \neq k} |r_i(\mathbf{x}_j) - r_i(\mathbf{x}_k)| / \sigma_i$$



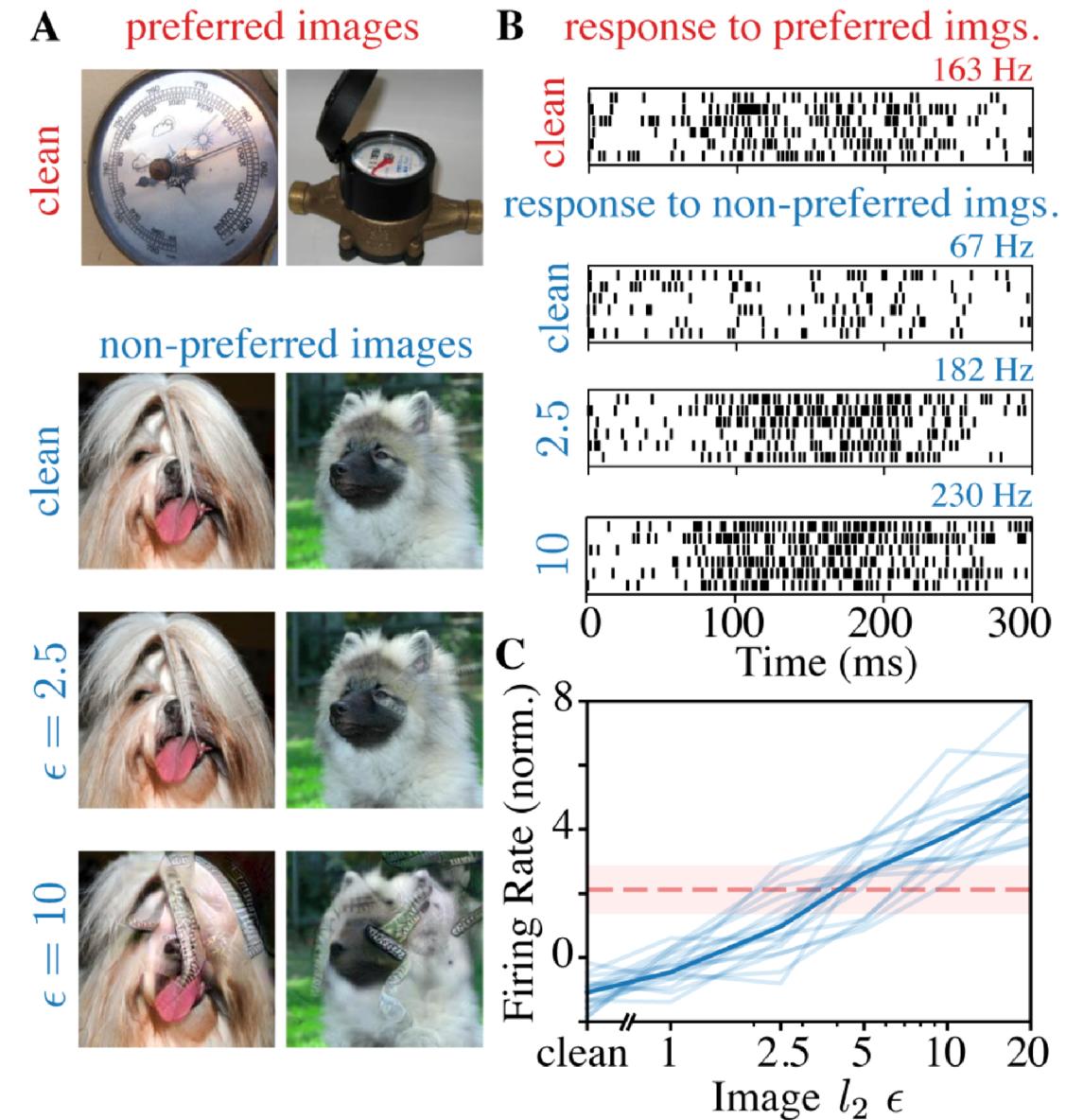
What do the adversarial examples of primate IT neurons look like?

B Sample adversarial images for primate IT neurons



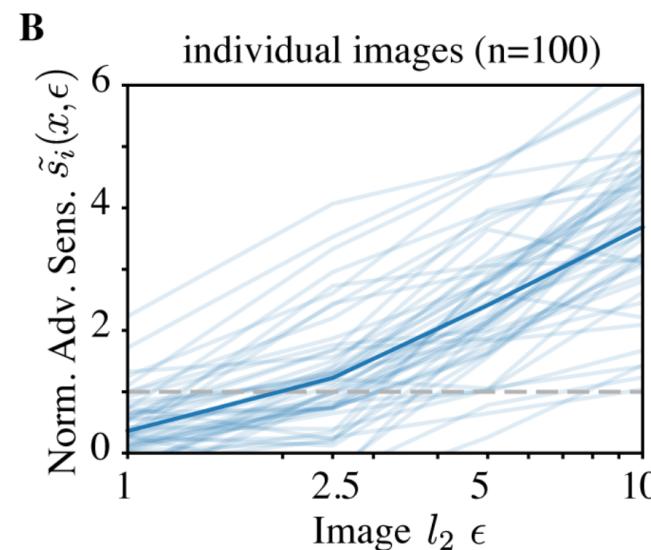
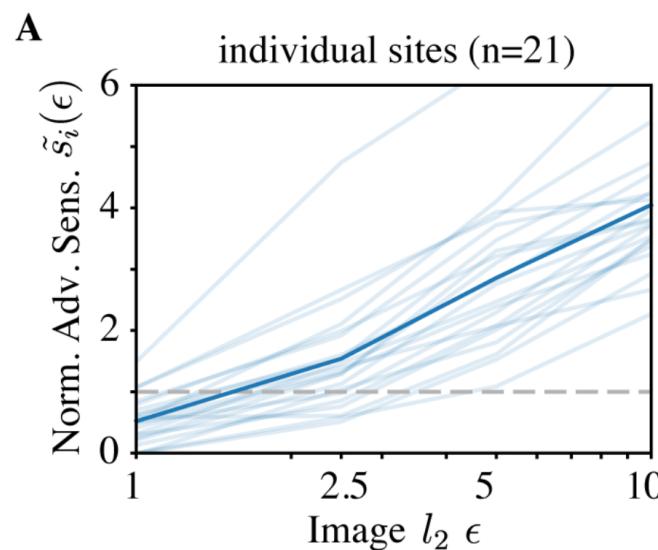
How stable is “category preference” of each IT neural site?

- Fix a site, identify the most and least preferred categories
- **Q: Is "category preference" well-defined?**
- Perform targeted adversarial perturbation
- $\epsilon = 2.5$ highly-preferred
- $\epsilon = 10$ "super-stimuli"
- Red dashed: preferred images



Are all IT neurons susceptible, or could the average results be due to just a few strongly modulated neurons?

- Adversarial images can be found on all recorded IT sites
- Adversarial images can be found very close to any clean images



- Adversarial samples for biological neurons are dense in the image space similar to that of artificial neural networks

3. Method: how to measure adversarial sensitivity of IT neural sites?

Measure adversarial sensitivity of IT neural sites

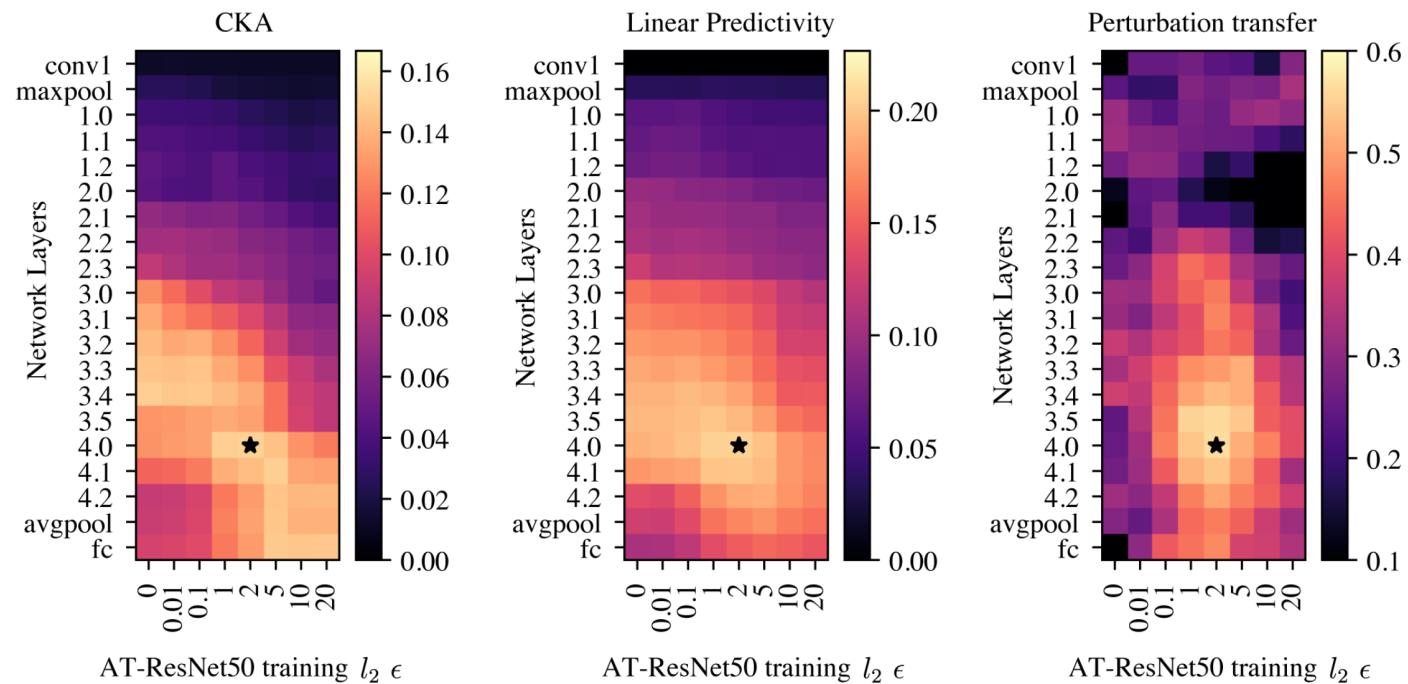
- Recall goal:

$$s_i(\mathbf{x}, \epsilon) := \max_{\|\boldsymbol{\delta}\|_2 < \epsilon} |r_i(\mathbf{x} + \boldsymbol{\delta}) - r_i(\mathbf{x})|$$

- Method: Build a “white-box” model to estimate the adversarial example
 - Iteratively generate better lower bound of adversarial example

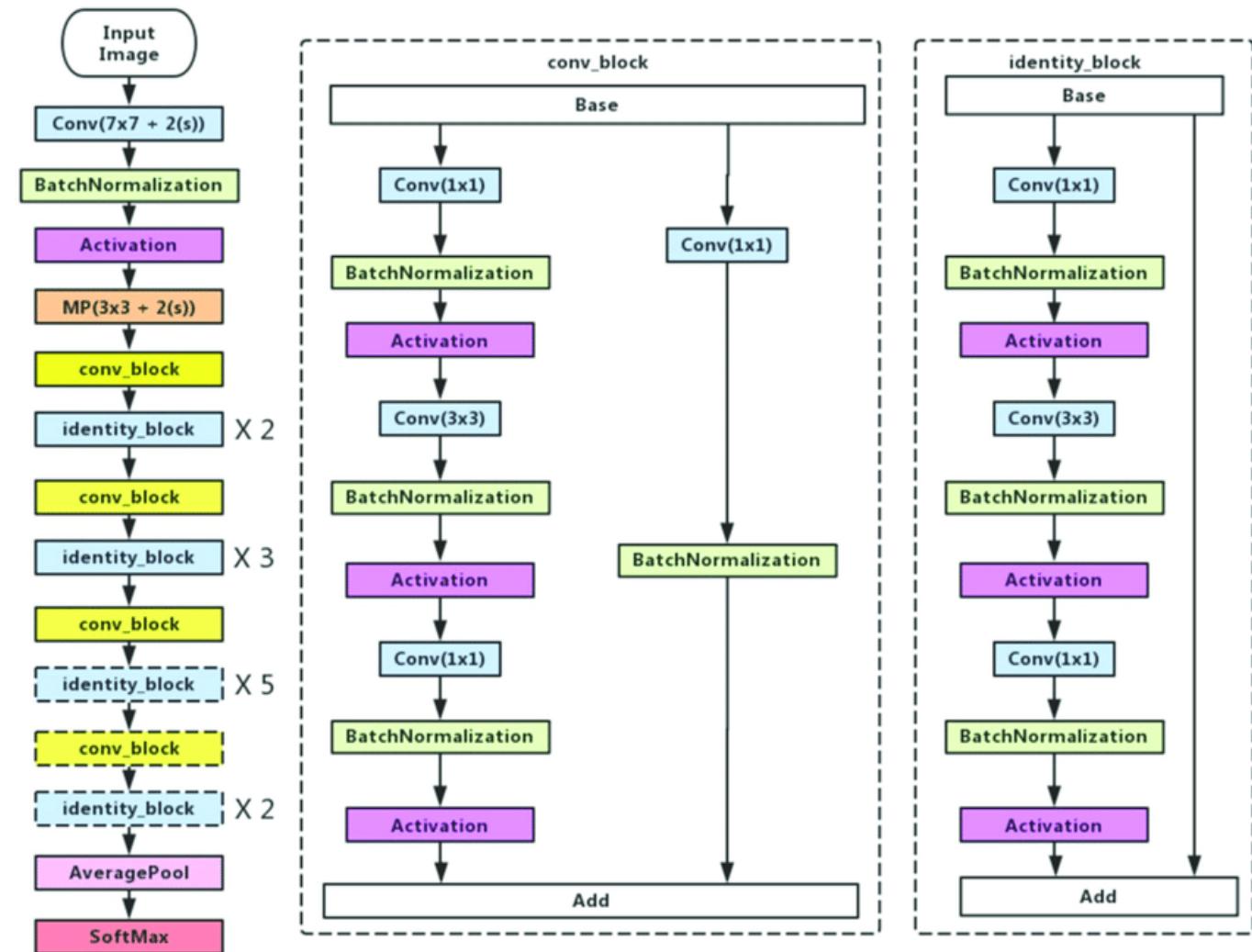
Screen a baseline model, criteria

1. Global representational similarity to IT as measured by CKA
2. Cross-validated linear predictivity for IT responses
3. How well does perturbations targeted toward a model layer transfers to IT neurons without any explicit mapping between the two systems

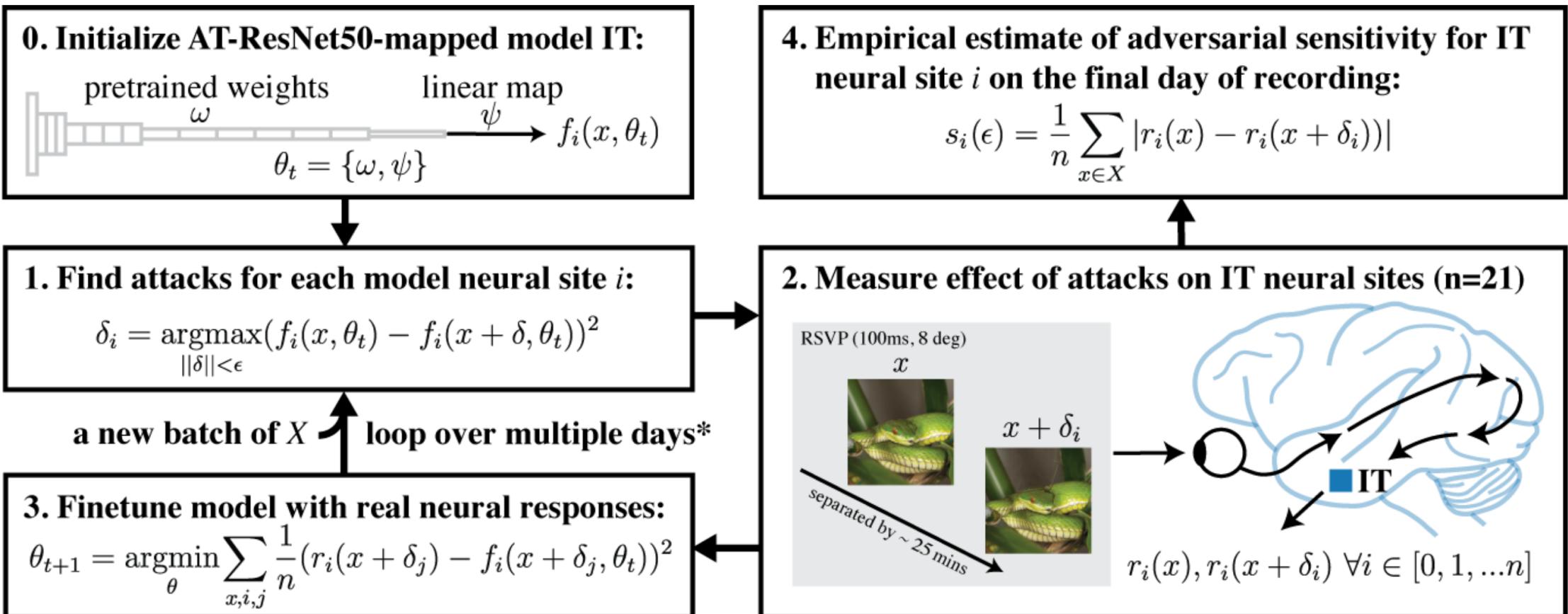


ResNet50

1 (conv)
+ 3×4 (conv_block)
+ $3 \times (2 + 3 + 5 + 2)$ (identity_block)
+1 (Fully-connected)
=50



Create adversarial examples iteratively

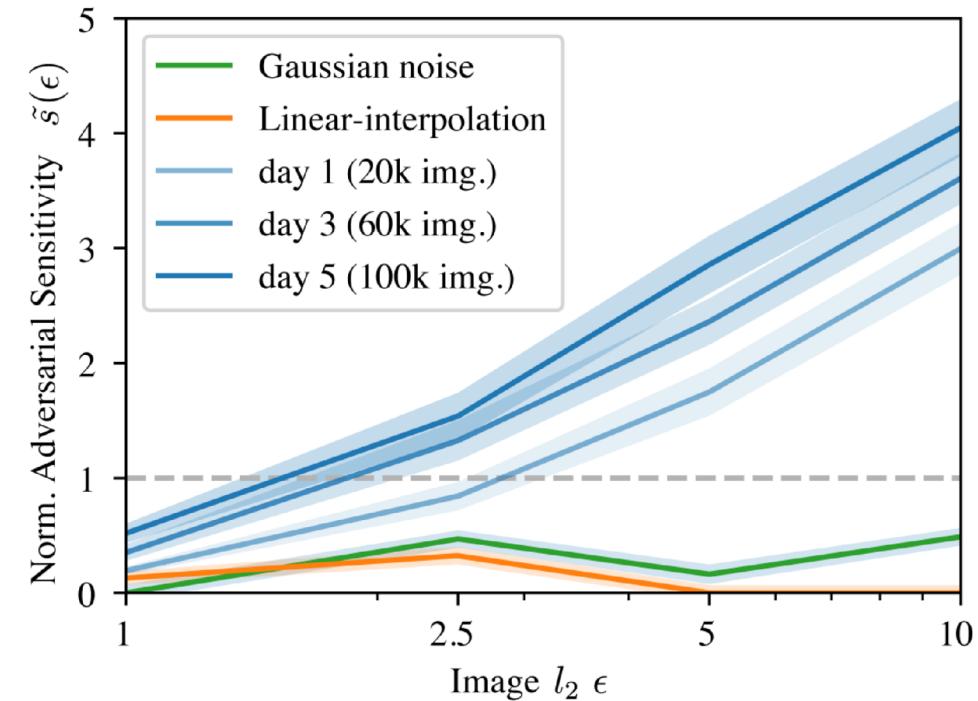


Details of measure effect of attacks on IT neural sites

- Show both clean and attack images to a fixating monkey
- The visual stimuli are presented 8 degrees over the visual field for 100ms followed by a 100ms grey mask as in a standard rapid serial visual presentation (RSVP) task.
- The average temporal separation between a clean image and its perturbed pair is 25 minutes
- Total of 6 days.
- For Figure 1A, we report IT sensitivity from the last day of experiment which sampled 882 unique images per perturbation ϵ (i.e. 42 images per neural site).
- Measure the total number spikes between 70ms-170ms after image presentation.

Does the method work?

1. Neural perturbation magnitude has an **consistent improvement** over days.
2. The perturbations achieved with our method is **significantly larger** than that achieved with a model-free method.
3. This explains why the field has systematically **underestimated** the sensitivity.



Introduce multiple methods to drastically improve convergence beyond the basic PGD

1. 100 **independent runs** for solving the adversarial images
2. Optimizing $\max_{\|\delta\|_2 < \epsilon} f_i(\mathbf{x} + \delta) - f_i(\mathbf{x})$ and $\max_{\|\delta\|_2 < \epsilon} f_i(\mathbf{x}) - f_i(\mathbf{x} + \delta)$
separately – reduce the chances to be stuck at saddle point.
3. **Larger** ϵ converges faster. First with a ball of radius 2ϵ and finally with one of radius ϵ . **Why?**
4. **Simulated annealing** with restarts: we begin with steps of size ϵ and reduce them by 10% every time no progress is made.

PGD

- Goal:

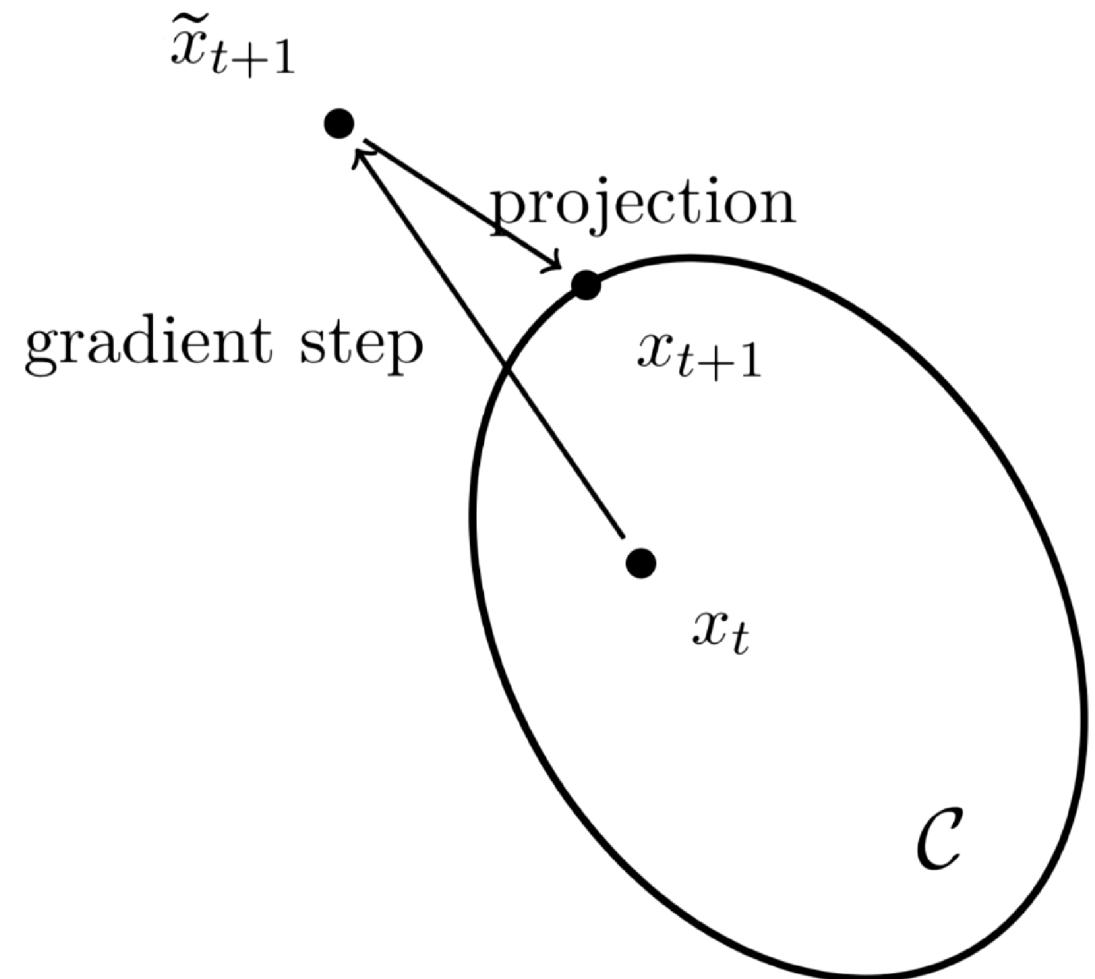
$$s_i(\mathbf{x}, \epsilon) := \max_{\|\boldsymbol{\delta}\|_2 < \epsilon} |f_i(\mathbf{x} + \boldsymbol{\delta}) - f_i(\mathbf{x})|$$

- Projected Gradient Descent

$$\min_{\|\boldsymbol{\delta}\|_2 < \epsilon} g(\boldsymbol{\delta}) := -|f_i(\mathbf{x} + \boldsymbol{\delta}) - f_i(\mathbf{x})|$$

$$\tilde{\boldsymbol{\delta}}_{t+1} = \boldsymbol{\delta}_t - \eta \nabla g(\boldsymbol{\delta}_t)$$

$$\boldsymbol{\delta}_{t+1} = \operatorname{argmin}_{\|\boldsymbol{\delta}\|_2 < \epsilon} \left\| \tilde{\boldsymbol{\delta}}_{t+1} - \boldsymbol{\delta} \right\|_2$$



Conclusion

The representations learned by adversarially trained artificial neural networks have already exceeded that of the corresponding biological neural representation in terms of their individual unit level adversarial robustness.

Paradox: how is it that primate visual perception seems so robust yet its fundamental units of computation are far more sensitive than expected?

1. Visual object recognition behavior in primate is actually
NOT adversarial robust
2. There is an unknown error-correction mechanism at the
population level in IT or in a down-stream area that
decodes object identity

Future work & new start

- Population level robustness
- Provides us with a set of standardized procedure

Thank you!