# Neural-guidance by the Human Ventral Visual Stream Improves Neural Network Robustness

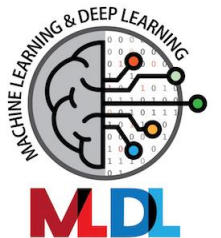**Zhenan Shao**[1,2], Linjian Ma[3], Bo Li[3,4], Diane M. Beck[1,2]

[1]Department of Psychology, University of Illinois Urbana-Champaign

[2]Beckman Institute, University of Illinois Urbana-Champaign

[3]Department of Computer Science, University of Illinois Urbana-Champaign
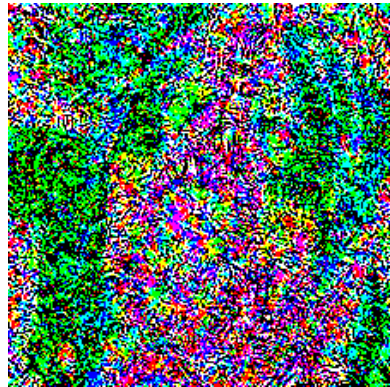
[4]Department of Computer Science, University of Chicago
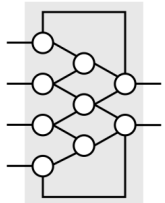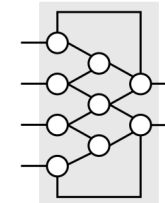
*MLDL Workshop at Sandia National Laboratories, 2024*

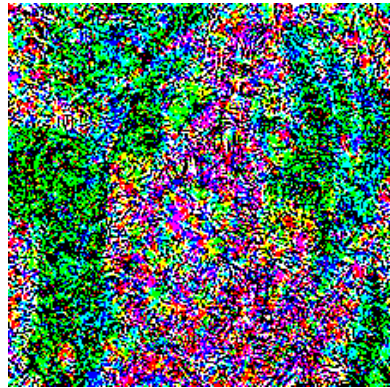# Vulnerable machine vision



Dog!

Dog!

Ventral Stream

Dog!
Conf:88%

Ostrich!
Conf:~100%

- Evolving representational space achieved by disentangling object manifolds along the ventral visual stream.

  *Dicarlo & Cox, Trends Cogn Sci. 2007*



Ventral Stream

*"Identity-preserving transformations"*

- Evolving representational space achieved by disentangling object manifolds along the ventral visual stream.

  *Dicarlo & Cox, Trends Cogn Sci. 2007*

Ventral Stream

- Evolving representational space achieved by disentangling object manifolds along the ventral visual stream.
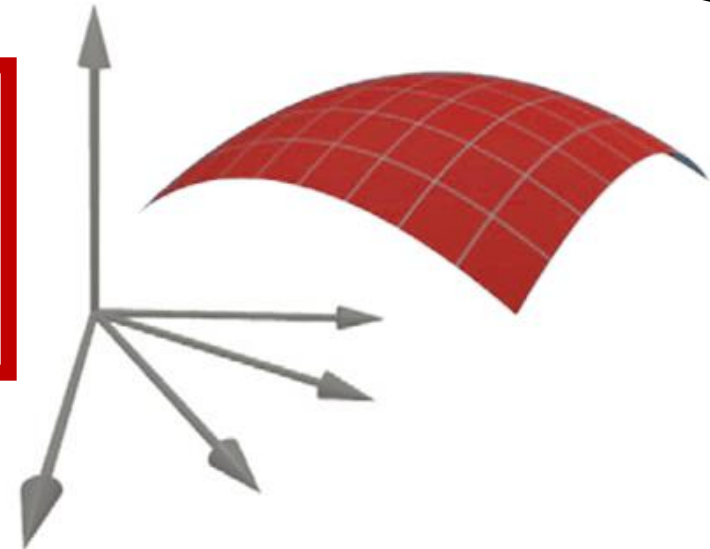
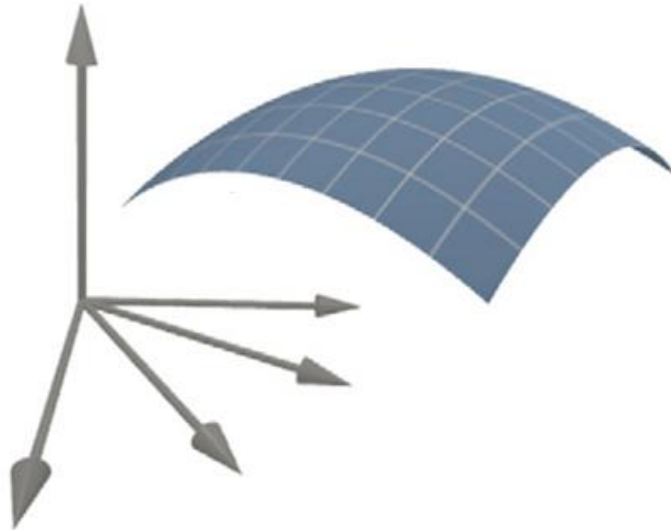*Dicarlo & Cox, Trends Cogn Sci. 2007*

Pixel space

"Good" neural space

Decision hyperplane

- Evolving representational space achieved by disentangling object manifolds along the ventral visual stream.
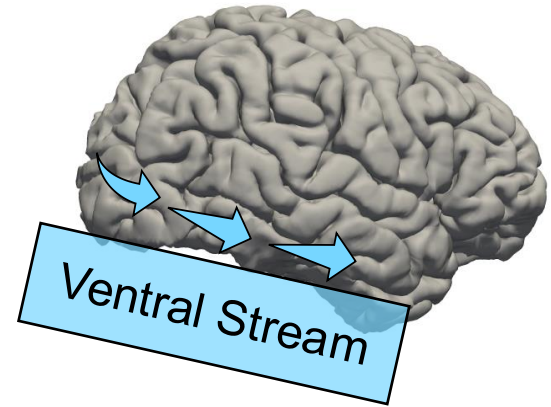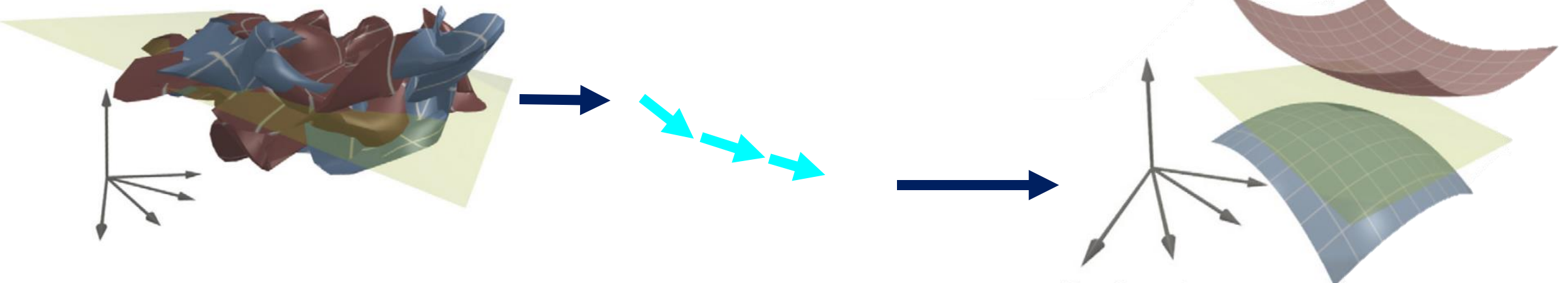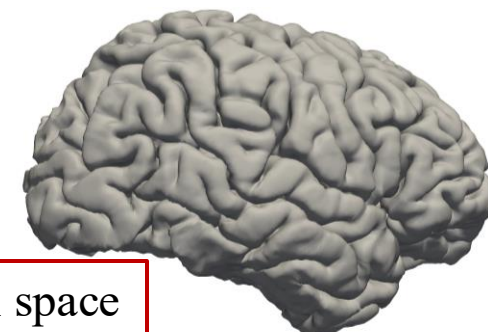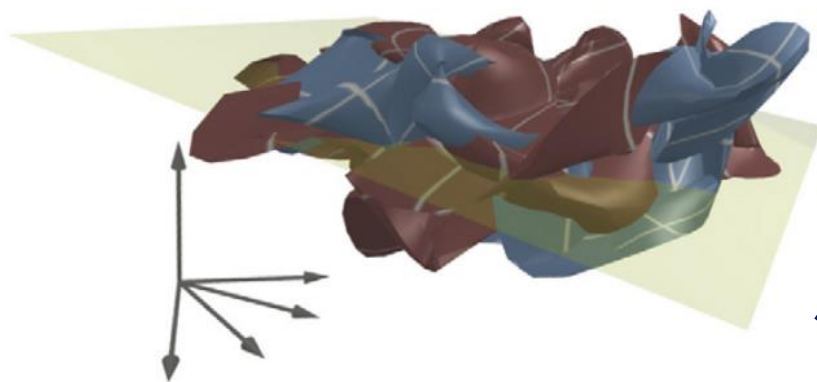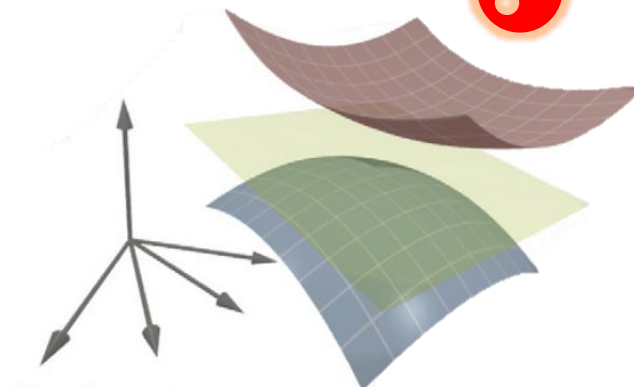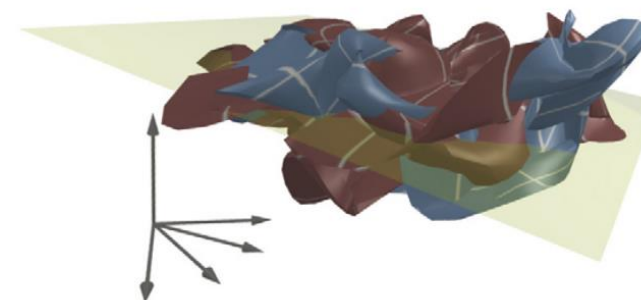
  *Dicarlo & Cox, Trends Cogn Sci. 2007*
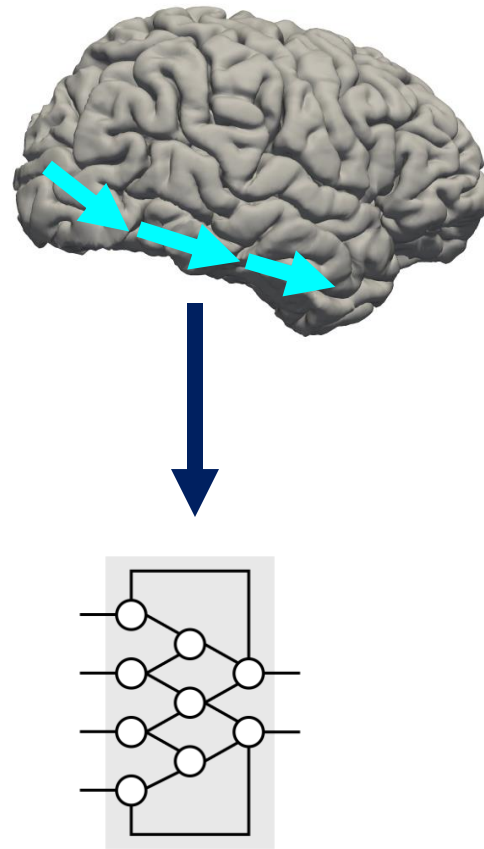
Pixel space

"Good" neural space

UNKNOWN

"Not so good" DNN space

Decision hyperplane
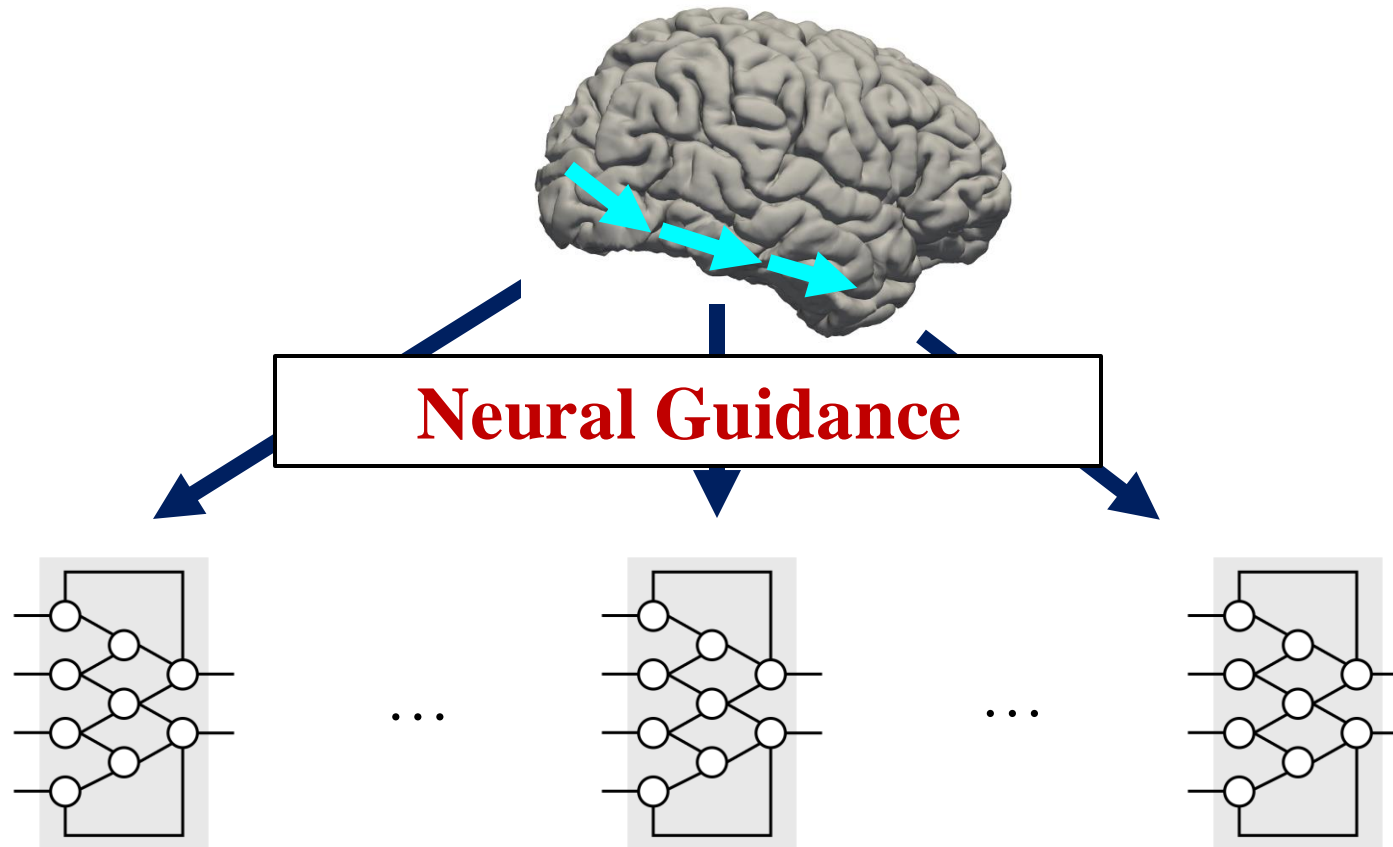
1. Does training guided by human ventral cortex activity improve DNN robustness?
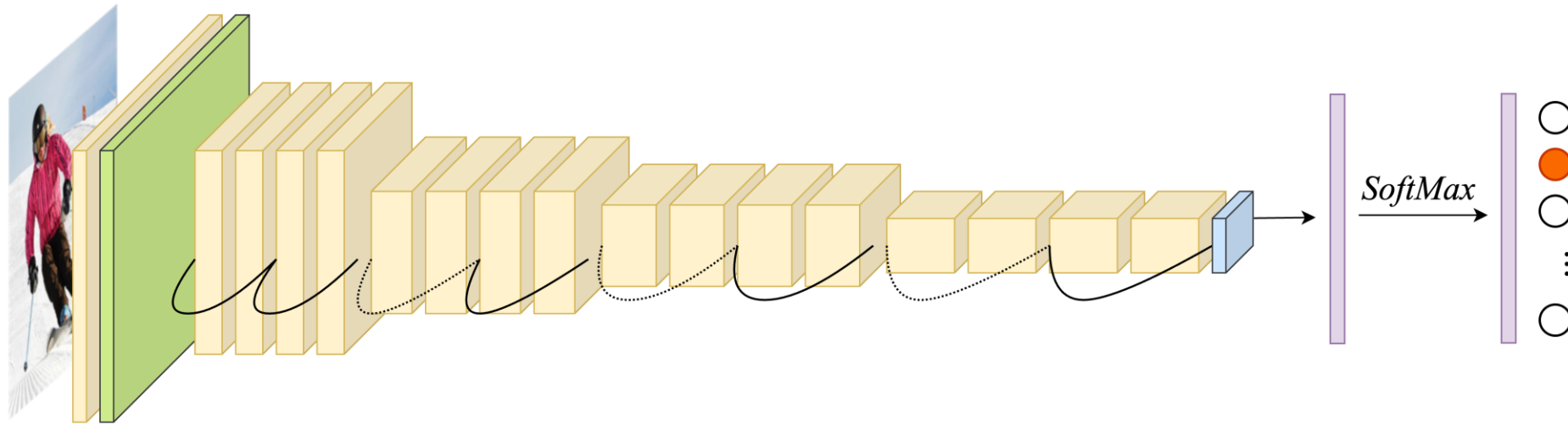
1. Does training guided by human ventral cortex activity improve DNN robustness?

2. Does such improvement increase as we ascend the ventral visual cortex?



**Neural Guidance**

$$Loss_{total} = L_{task}$$



Max Pooling

Avg Pooling

Conv

FC

$$Loss_{total} = L_{task}$$



Max Pooling

Avg Pooling

Conv

FC

*SoftMax*

*"Task Head"*

$$Loss_{total} = L_{task} + \left\| R_{DNN} - R_{neural} \right\|_2$$



*"Task Head"*

*"Neural Head"*

SoftMax

$l_2$

Max Pooling

Avg Pooling

Conv

FC

FC -- DNN neural representation

Human neural representation

# Training with **Neural Guidance**

$$Loss_{total} = \alpha L_{task} + (1 - \alpha)\big|\big|R_{DNN} - R_{neural}\big|\big|_2$$



*SoftMax*

*"Task Head"*

*"Neural Head"*

$l_2$

Max Pooling

Avg Pooling

Conv

FC

FC -- DNN neural representation

Human neural representation

**?**

- Brain activities were recorded with 7T fMRI while each human subject viewing 10,000 natural images. *(Natural Scene Dataset, Allen et al., Nat. Neurosci. 2022)*



Neural activity pattern

- Brain activities were recorded with 7T fMRI while each human subject viewing ~30,000 natural images.

- 7 bilateral Regions of Interest (ROIs) were used
  *(Wang et al., Cereb. Cortex, 2015)*
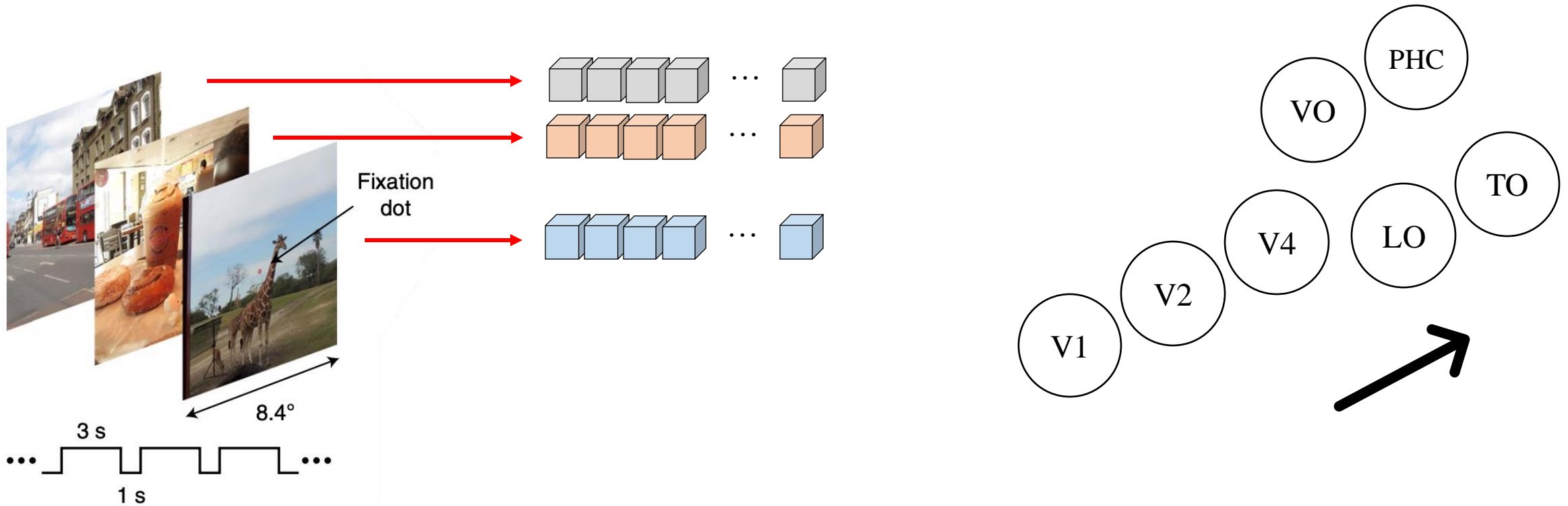


Ventral Visual Stream Hierarchy

- Brain activities were recorded with 7T fMRI while each human subject viewing ~30,000 natural images.

- 7 bilateral Regions of Interest (ROIs) were used



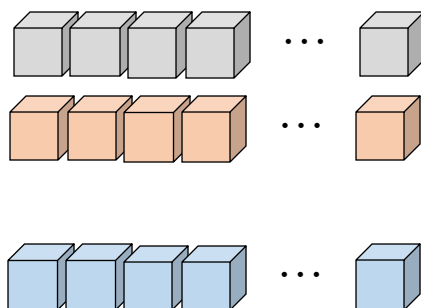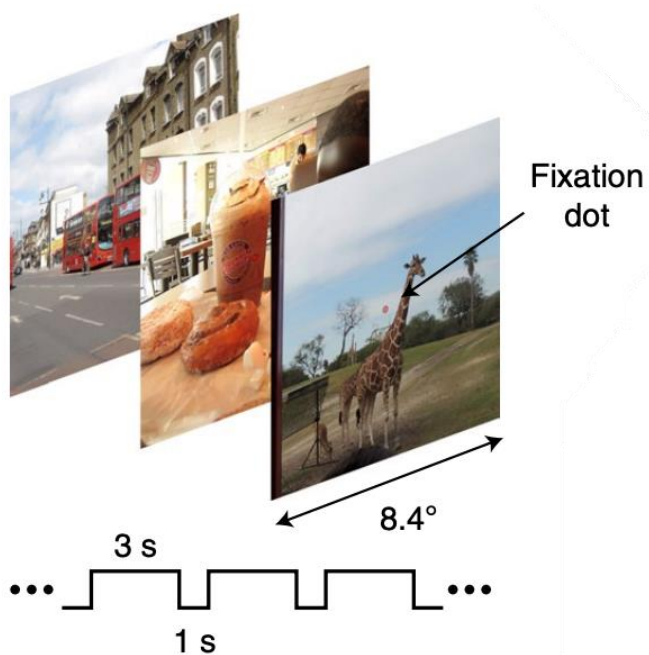*(NSD, Allen et al., Nat. Neurosci. 2022)*

- Brain activities were recorded with 7T fMRI while each human subject viewing ~30,000 natural images.

- 7 bilateral Regions of Interest (ROIs) were used



*(NSD, Allen et al., Nat. Neurosci. 2022)*

- Brain activities were recorded with 7T fMRI while each human subject viewing ~30,000 natural images.

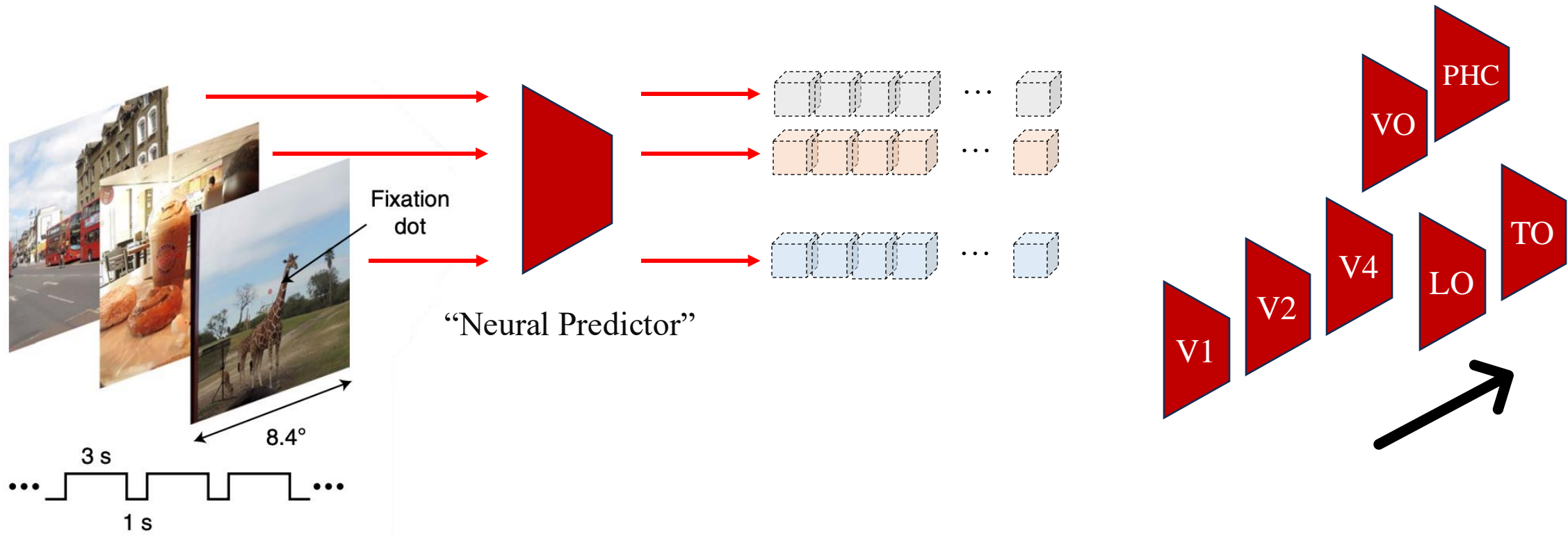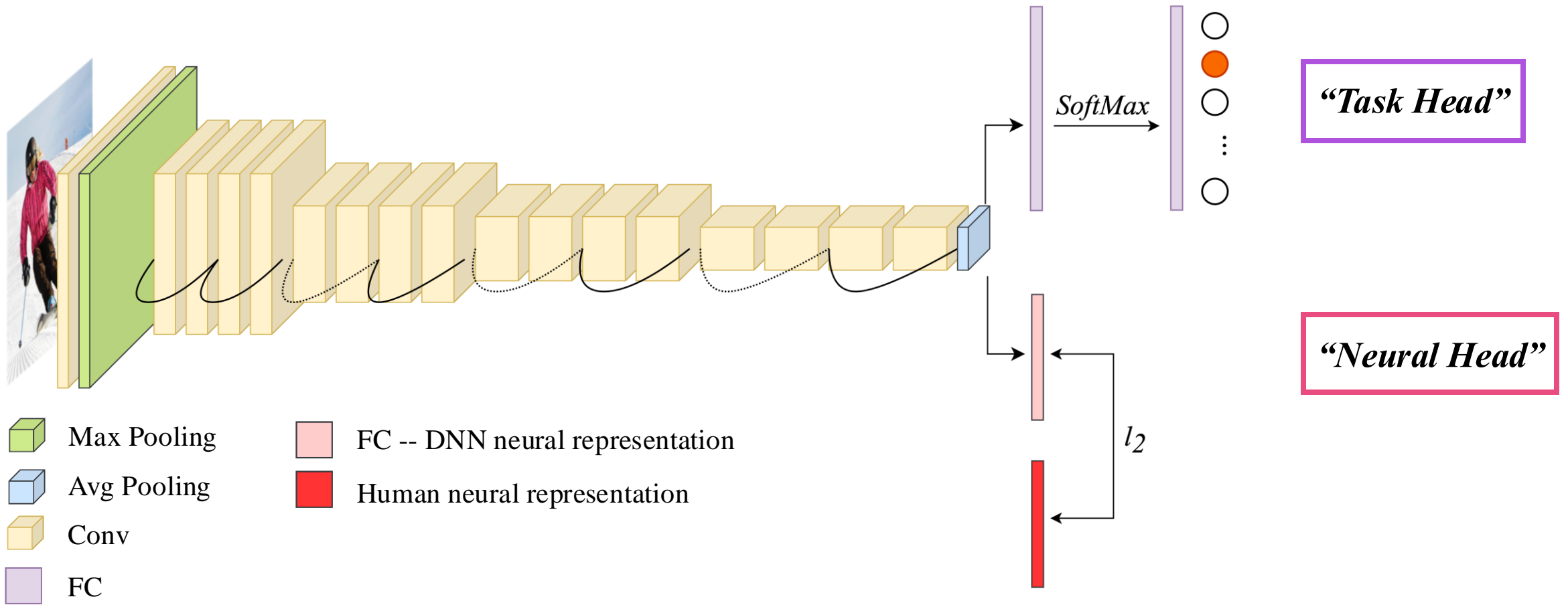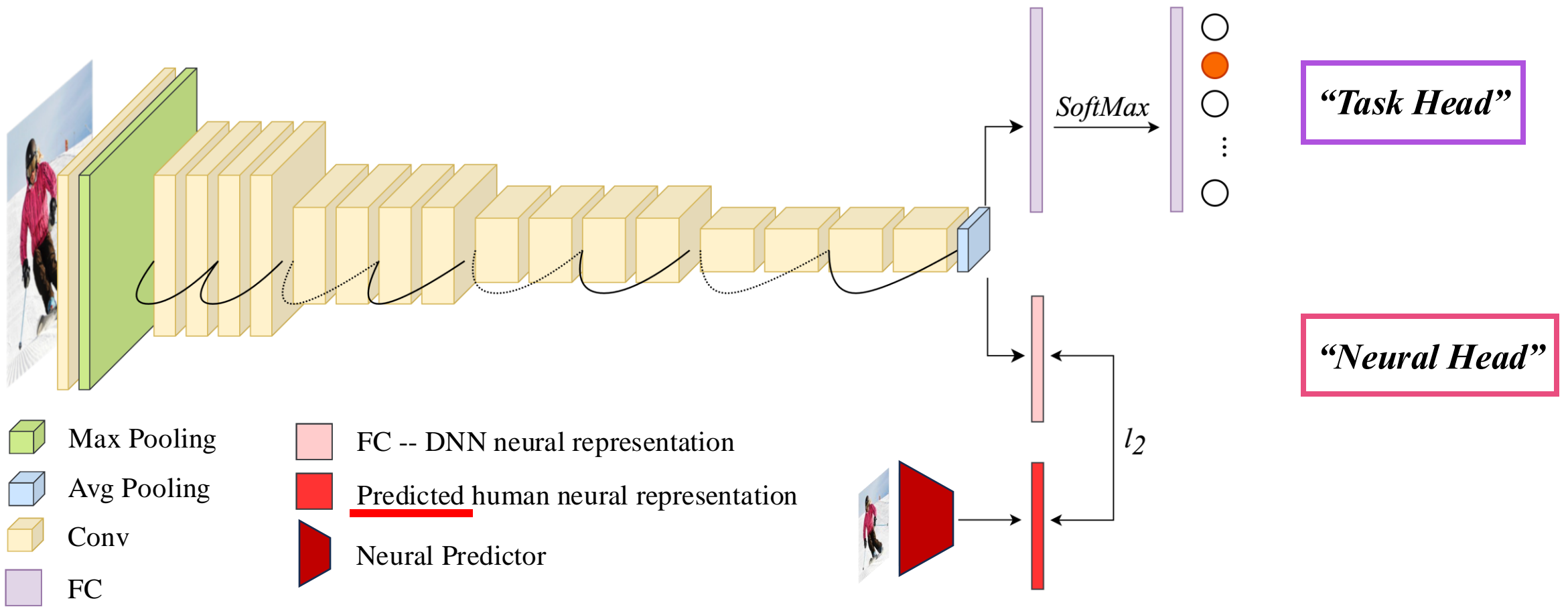- 7 bilateral Regions of Interest (ROIs) were used



"Neural Predictor"

*(NSD, Allen et al., Nat. Neurosci. 2022)*

$$Loss_{total} = \alpha L_{task} + (1 - \alpha)\left|\left|R_{DNN} - R_{neural}\right|\right|_2$$



*"Task Head"*

*"Neural Head"*

SoftMax

$l_2$

Max Pooling

Avg Pooling

Conv

FC

FC -- DNN neural representation

Human neural representation

$$Loss_{total} = \alpha L_{task} + (1 - \alpha)\big|\big|R_{DNN} - R_{neural}\big|\big|_2$$



*SoftMax*

**"Task Head"**

**"Neural Head"**

$l_2$

Max Pooling

Avg Pooling

Conv

FC

FC -- DNN neural representation

Predicted human neural representation

Neural Predictor

- 7 DNNs trained with Neural Guidance



Ventral Visual Stream Hierarchy

- 7 DNNs trained with Neural Guidance

- 4 baseline models for comparison

*"None"*                    *"Random"*                         *"V1-shuffle"*
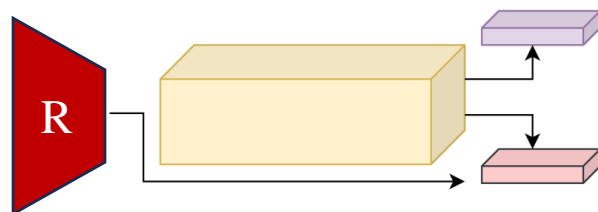
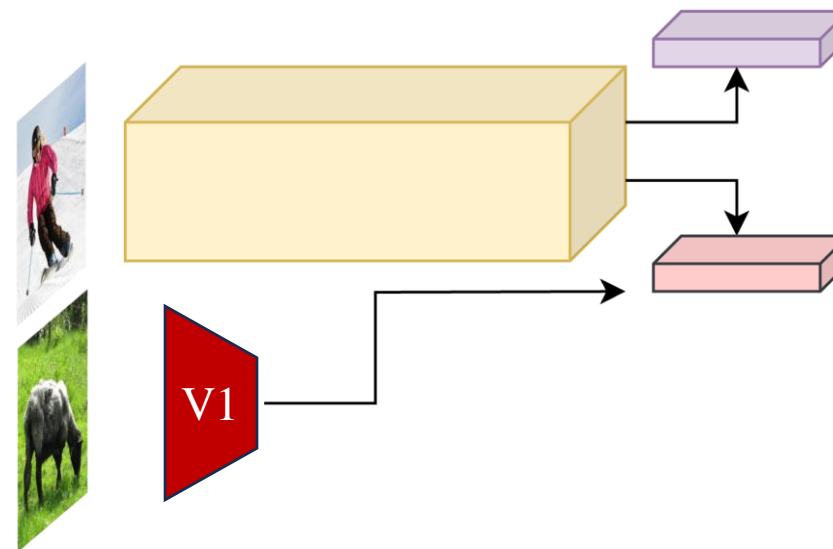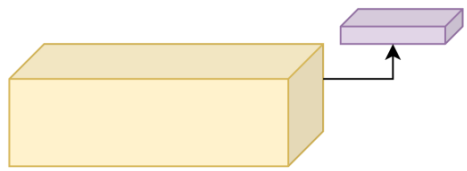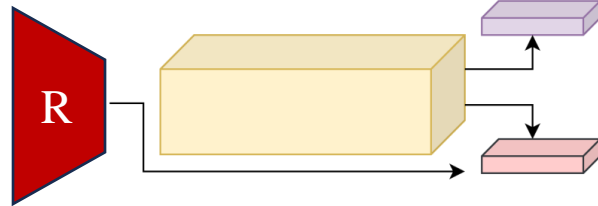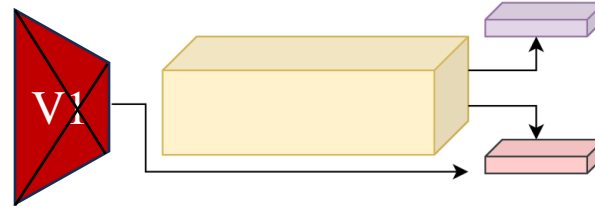- 7 DNNs trained with **Neural Guidance**



- 4 baseline models for comparison



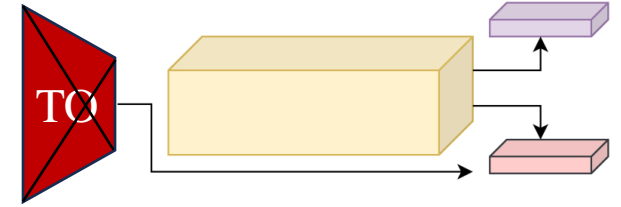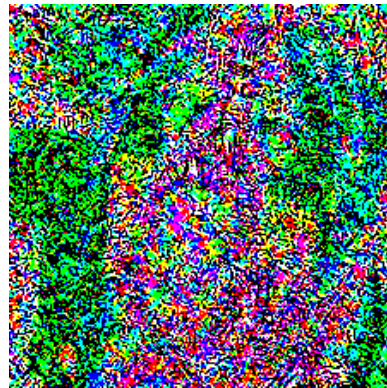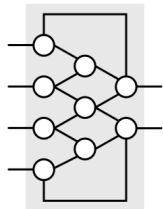*"None"*                     *"Random"*                     *"V1-shuffle"*                     *"TO-shuffle"*

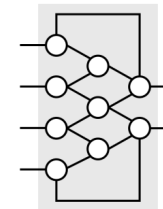- $l_p$-based adversarial attack: $\max_{||\tau||_p < \epsilon} l(f_\theta(x + \tau), y)$
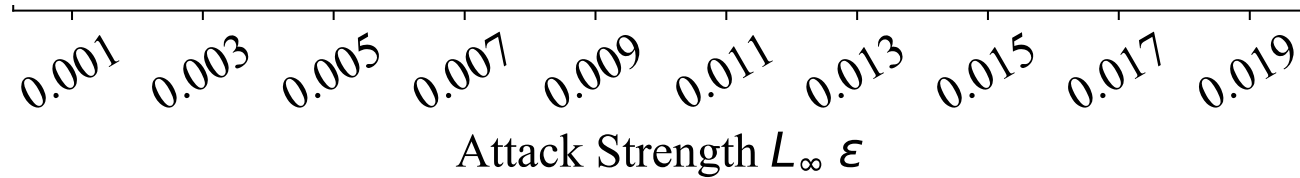


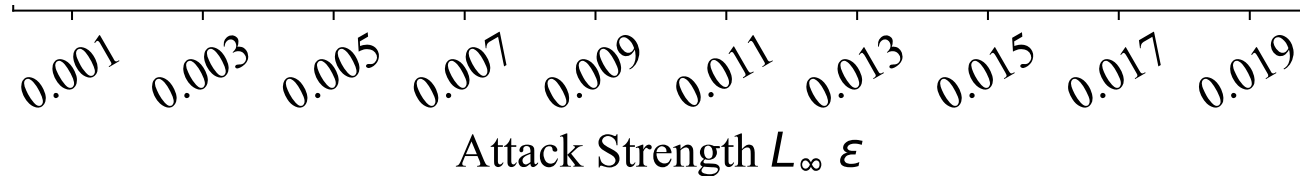*Dog!*
Conf:88%

*Ostrich!*
Conf:91%

- $l_p$-based adversarial attack: $\max_{||\tau||_p < \epsilon} l(f_\theta(x + \tau), y)$



Attack Strength $L_\infty$ $\varepsilon$

0.001  0.003  0.005  0.007  0.009  0.011  0.013  0.015  0.017  0.019

*Task: Image Classification*
*Dataset: ImageNet (Deng et al., 2009)*
*Attack: $l_\infty$-based PGD attack*



Attack Strength $L_\infty$ $\varepsilon$

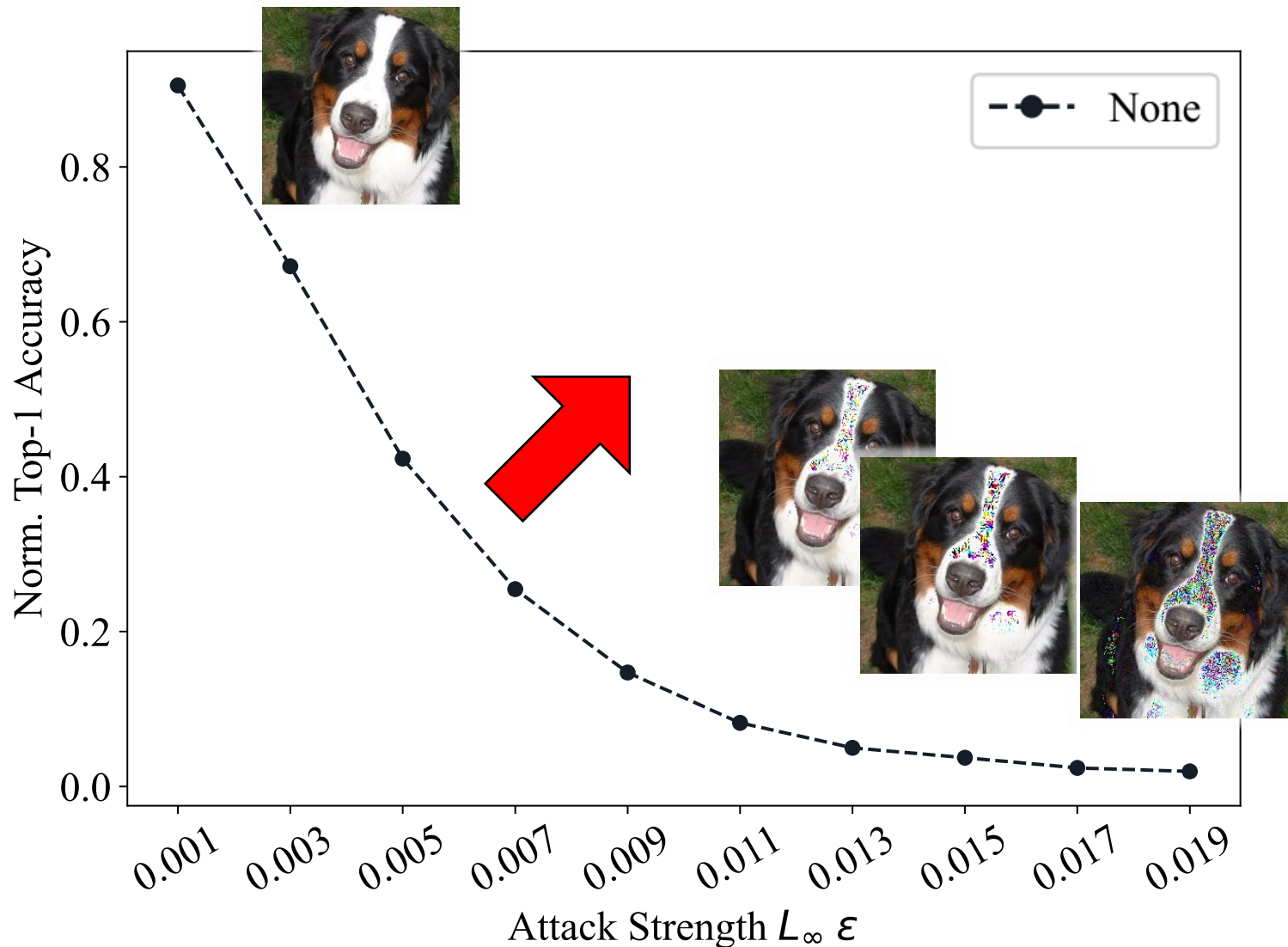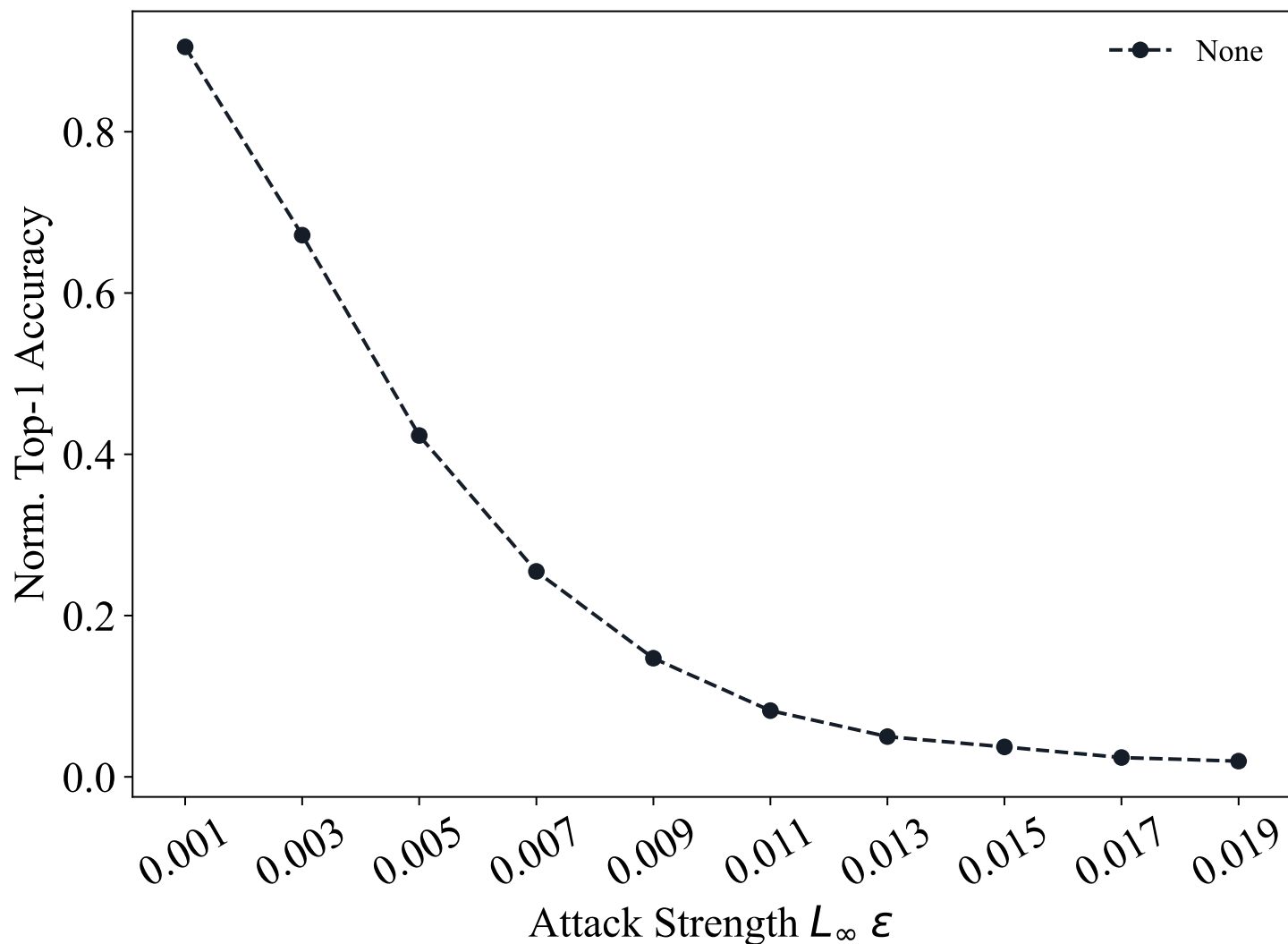0.001  0.003  0.005  0.007  0.009  0.011  0.013  0.015  0.017  0.019

Task: Image Classification
Dataset: ImageNet (Deng et al., 2009)
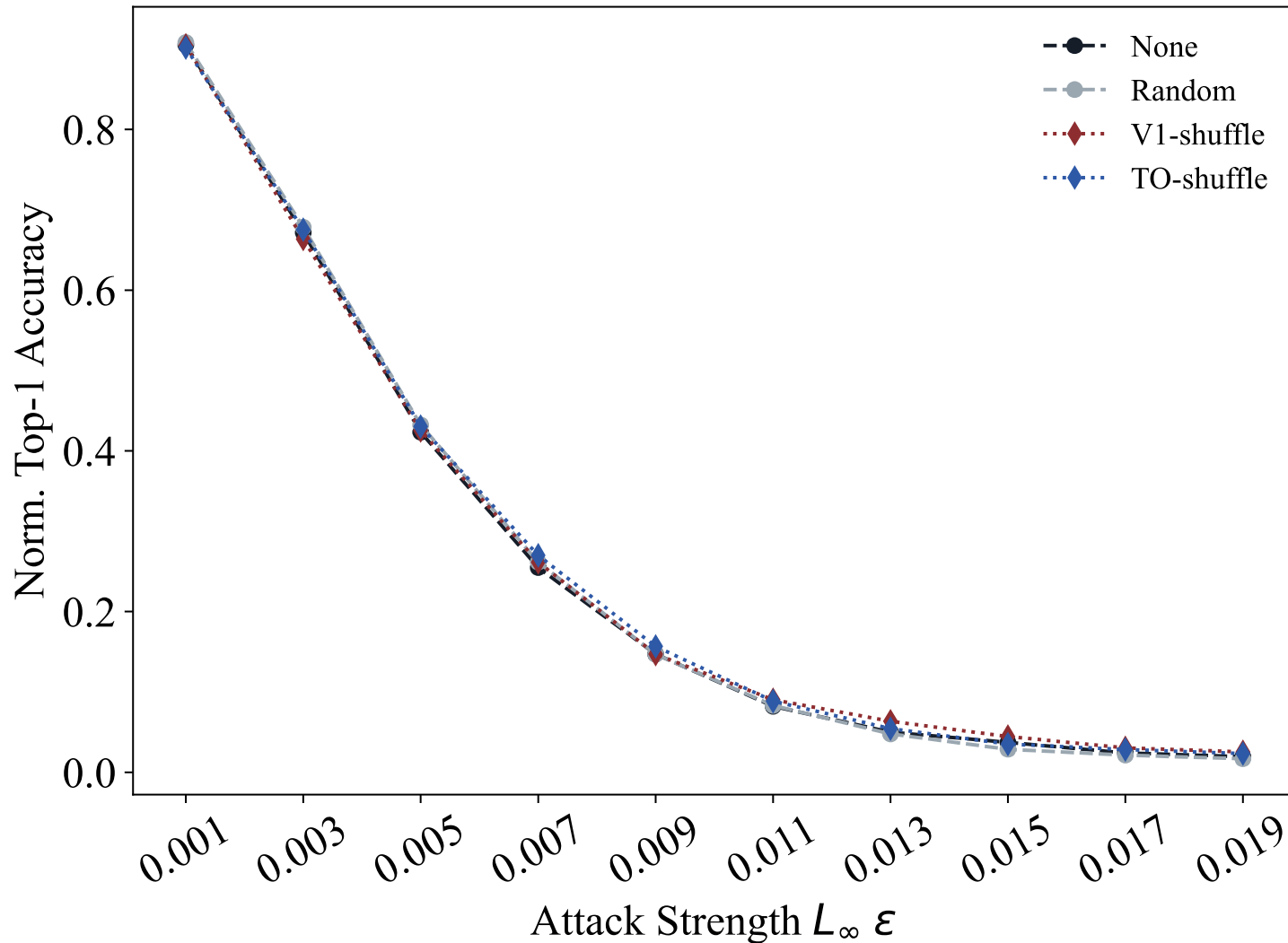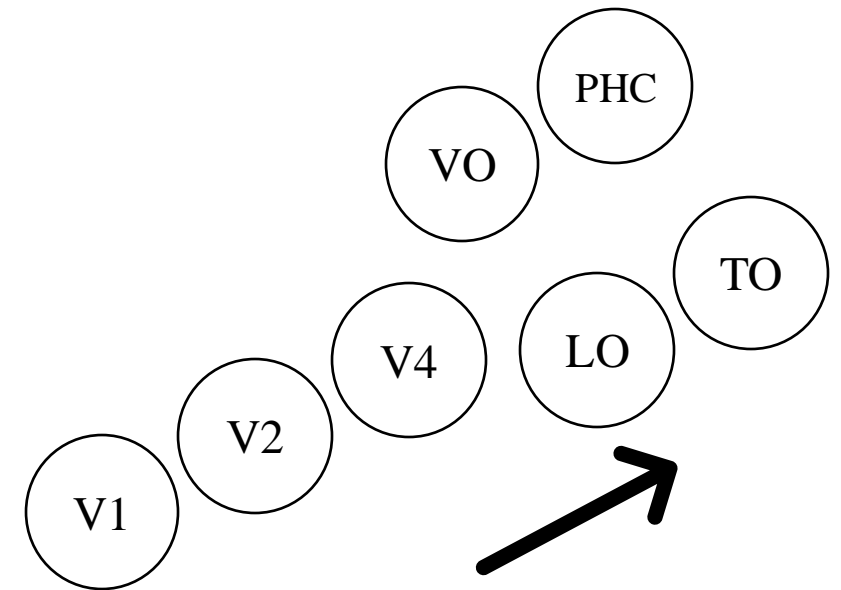Attack: $l_\infty$-based PGD attack

# Evaluating DNN robustness - Results



*Task: Image Classification*
*Dataset: ImageNet (Deng et al., 2009)*
*Attack: $l_\infty$-based PGD attack*

*None*   *Random*

*V1-shuffle*   *TO-shuffle*

Task: Image Classification
Dataset: ImageNet (Deng et al., 2009)
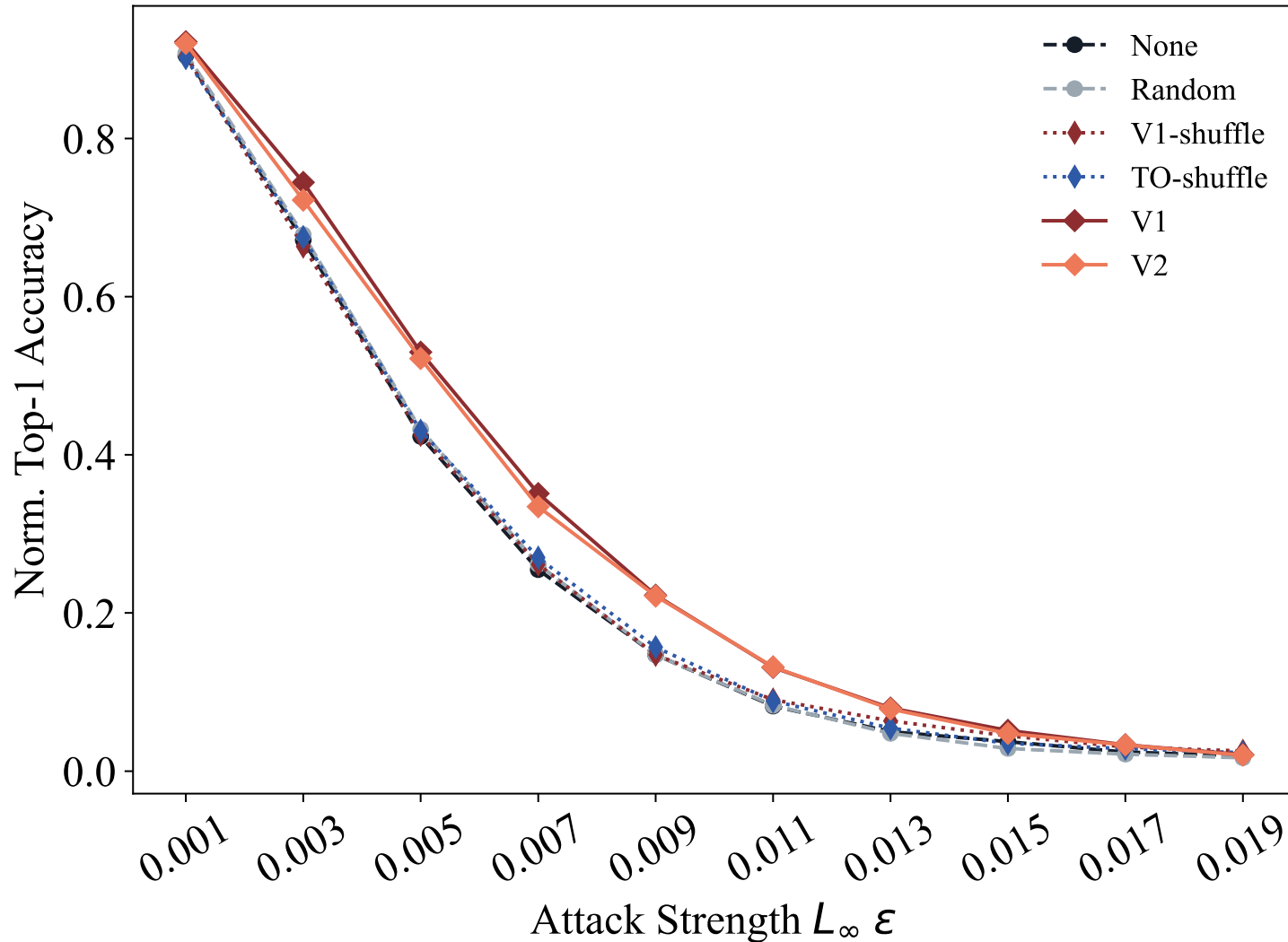Attack: $l_\infty$-based PGD attack

*Task: Image Classification*
*Dataset: ImageNet (Deng et al., 2009)*
*Attack: $l_\infty$-based PGD attack*

# Evaluating DNN robustness - Results



*Task: Image Classification*
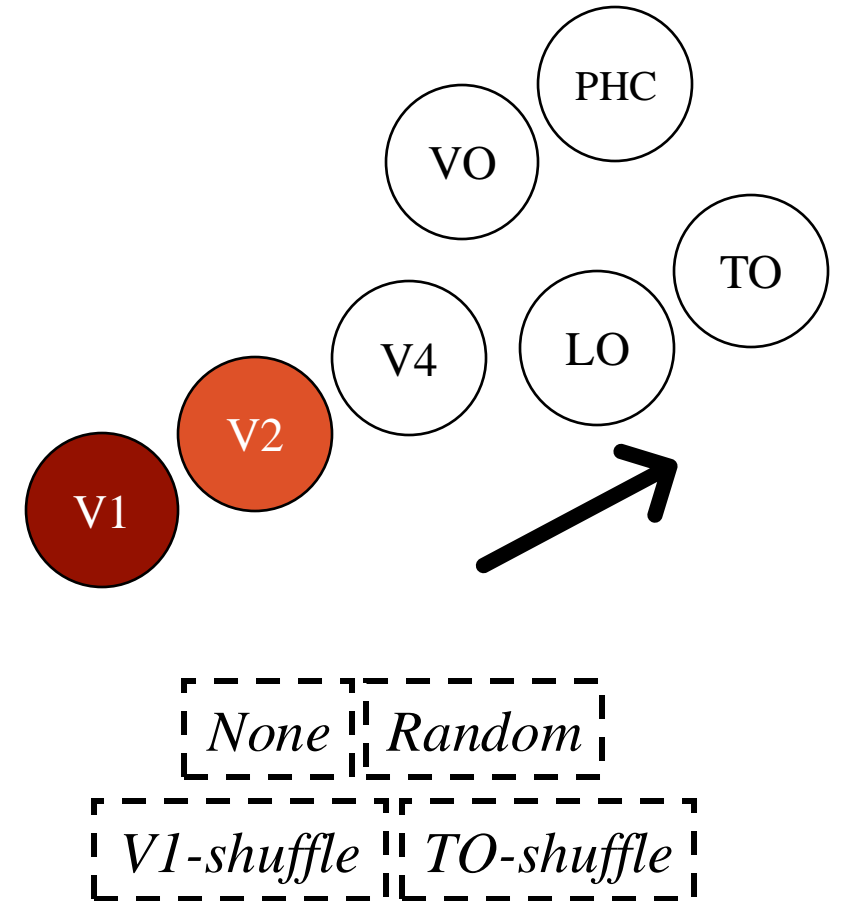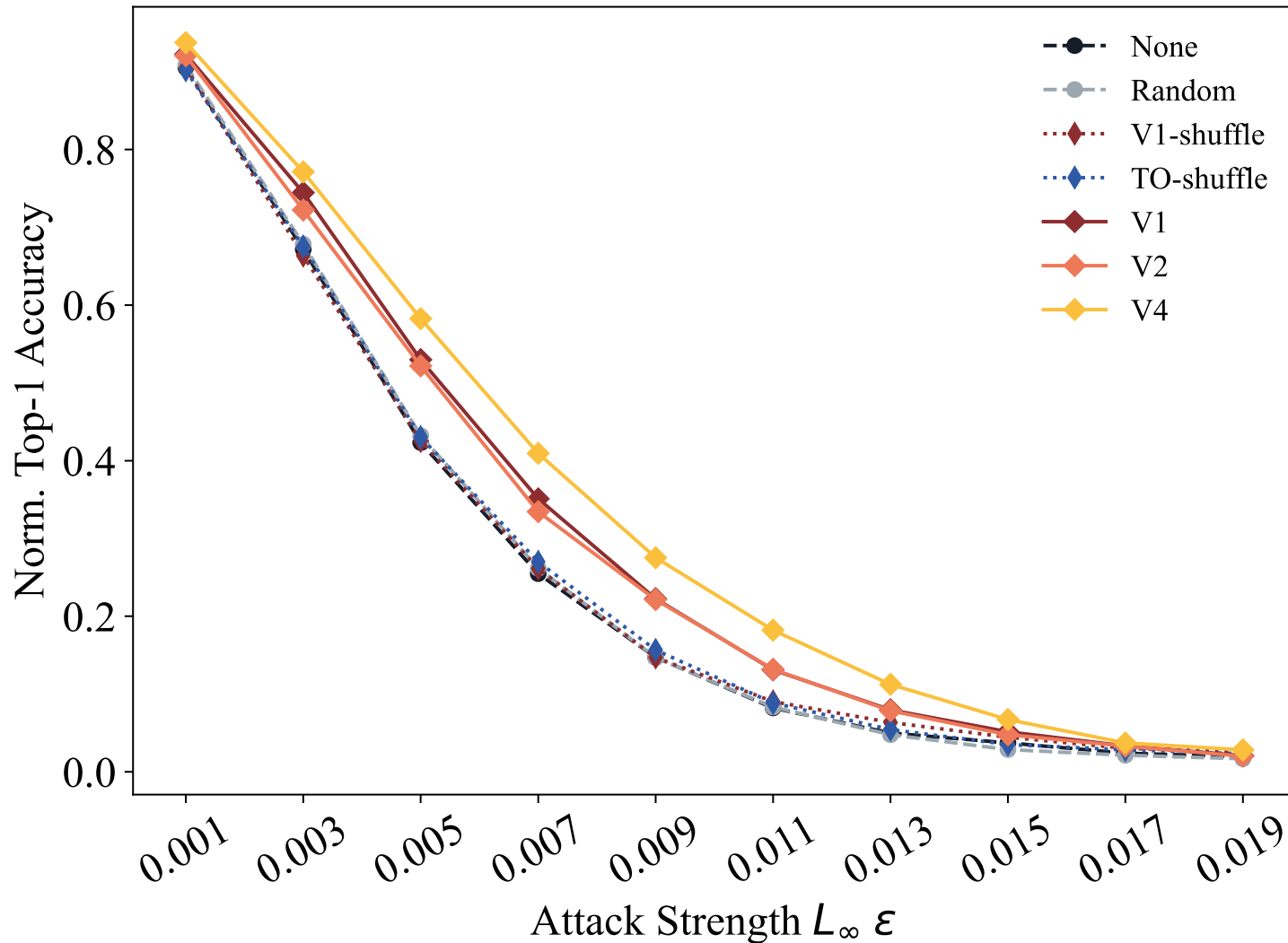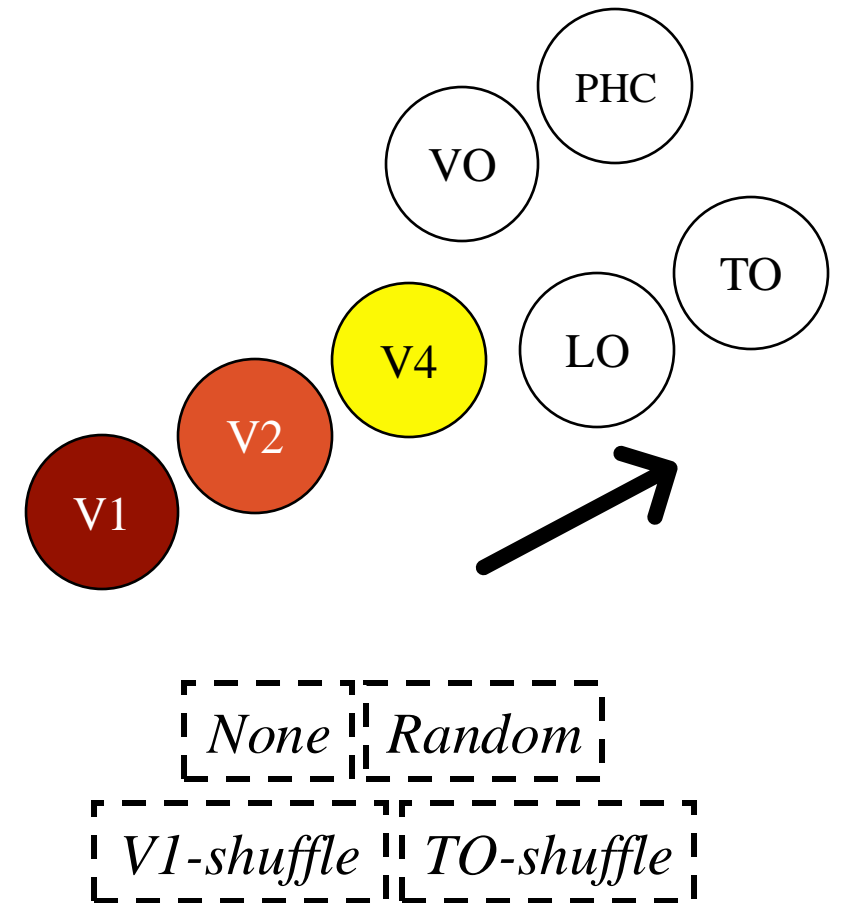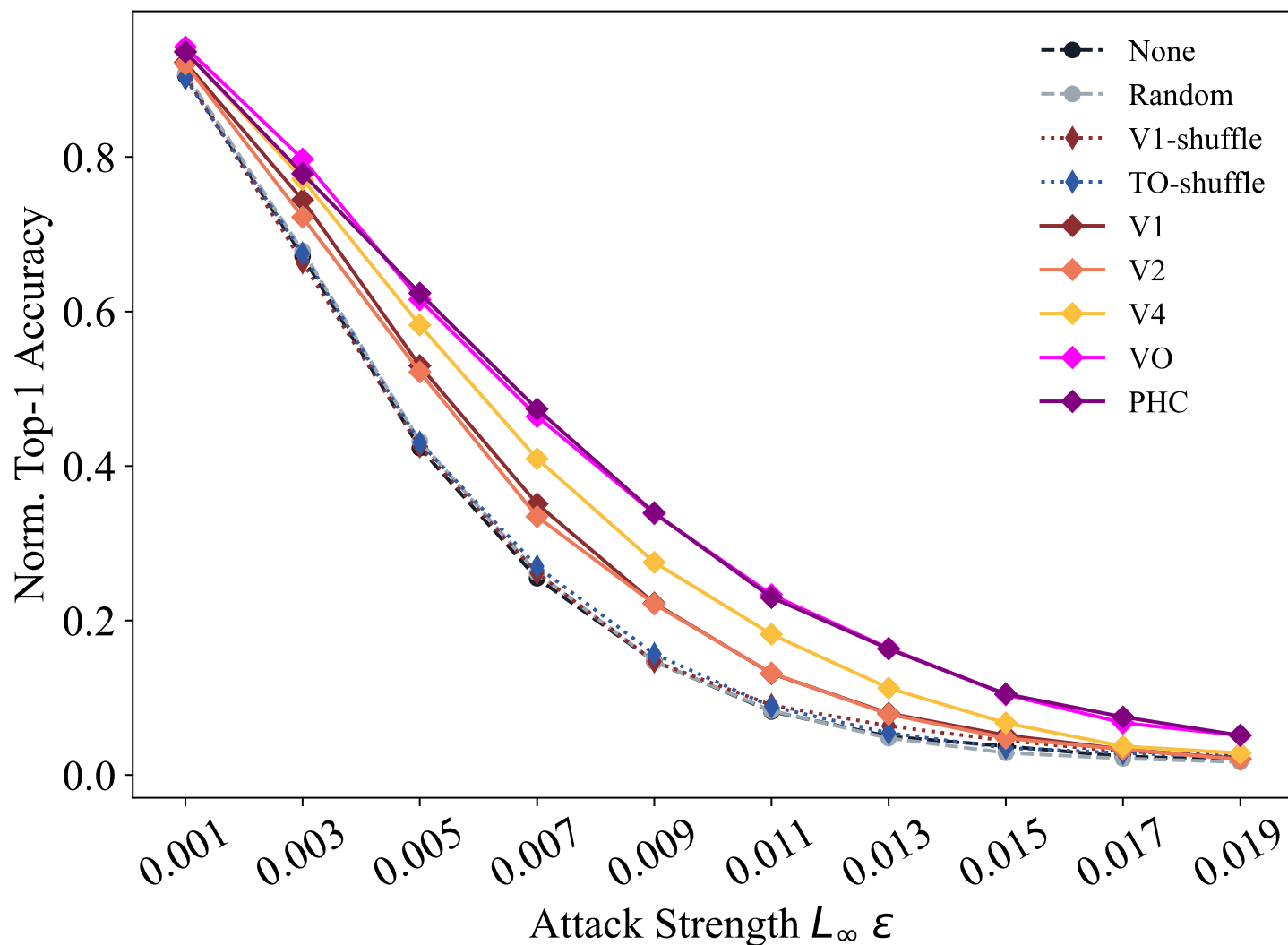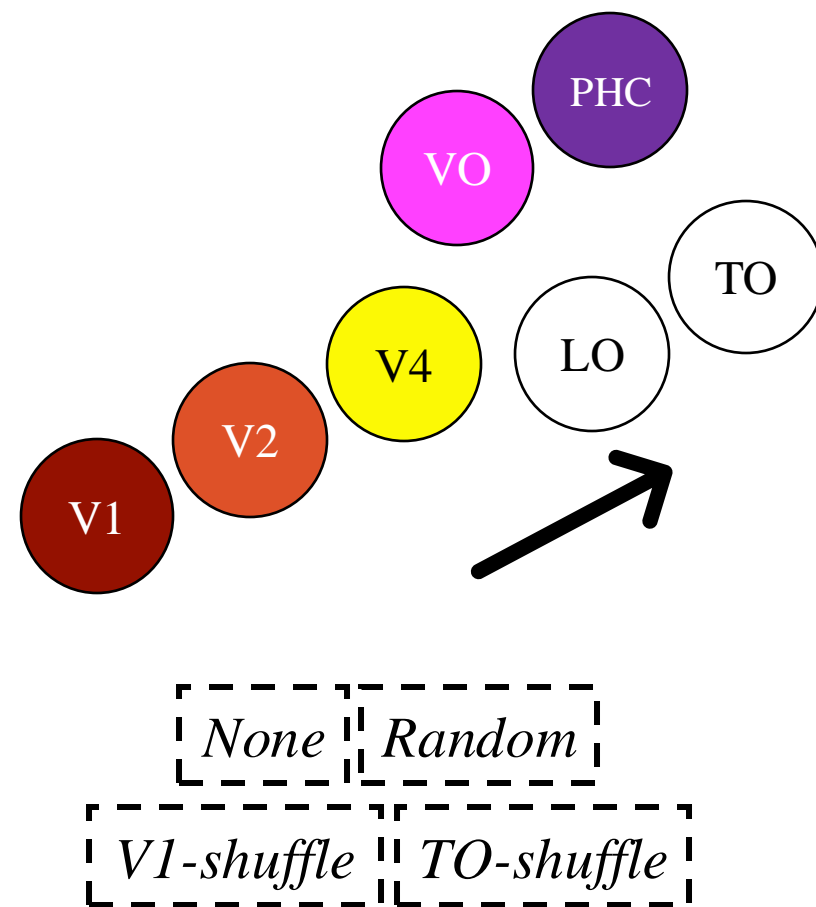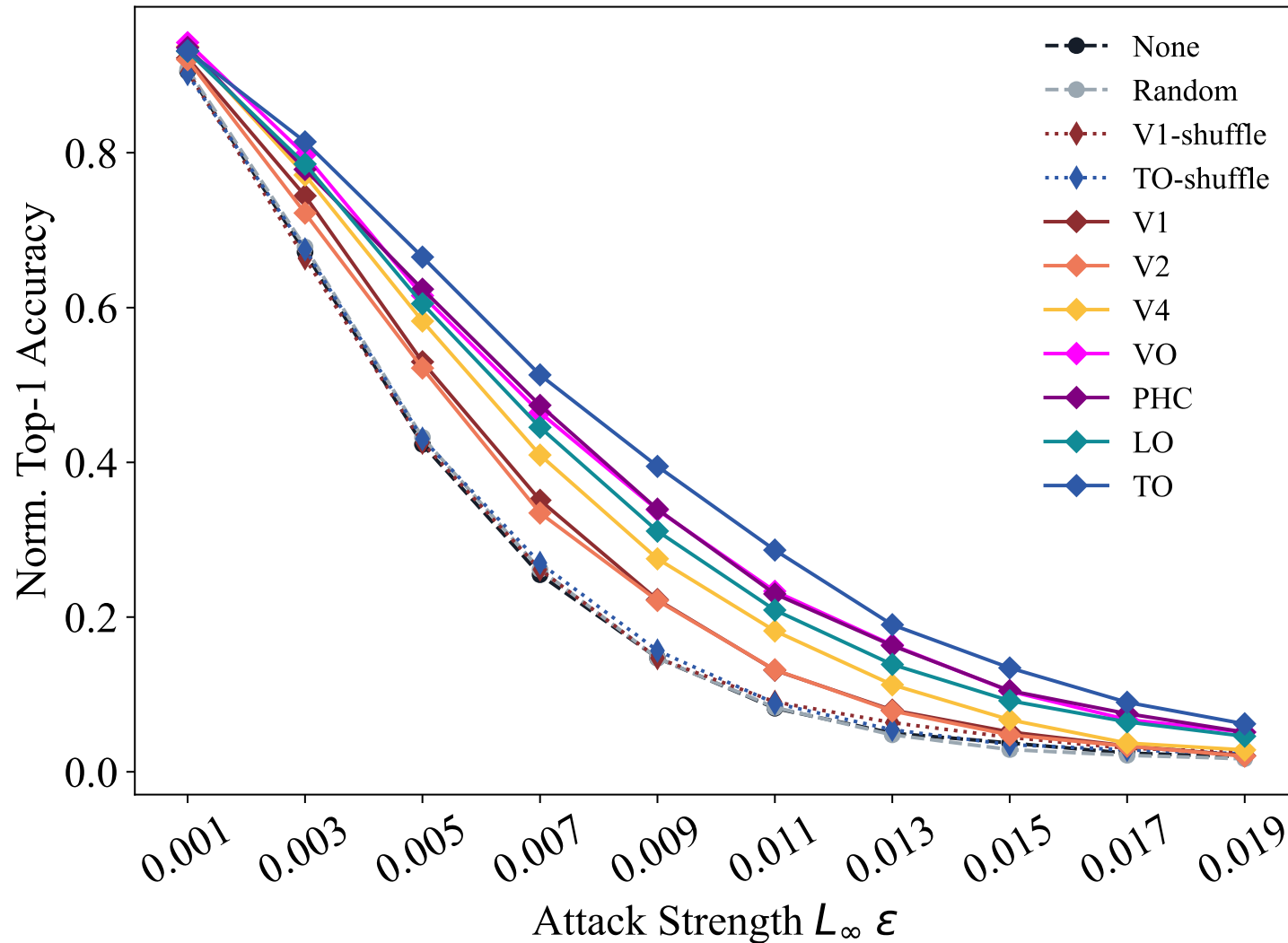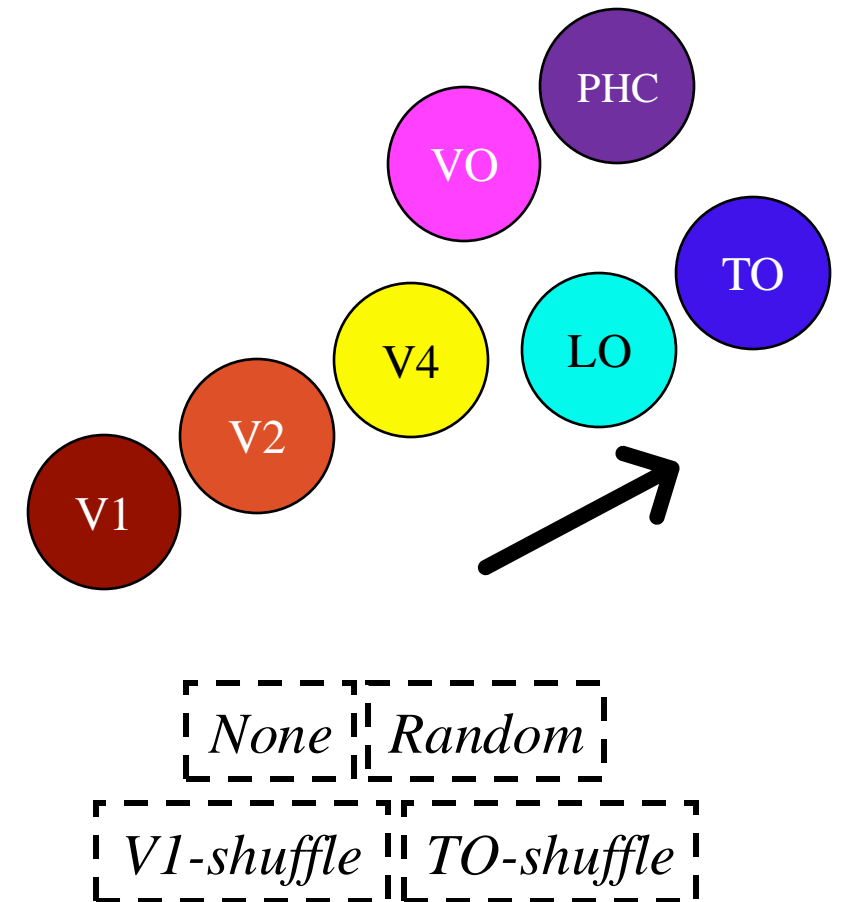*Dataset: ImageNet (Deng et al., 2009)*
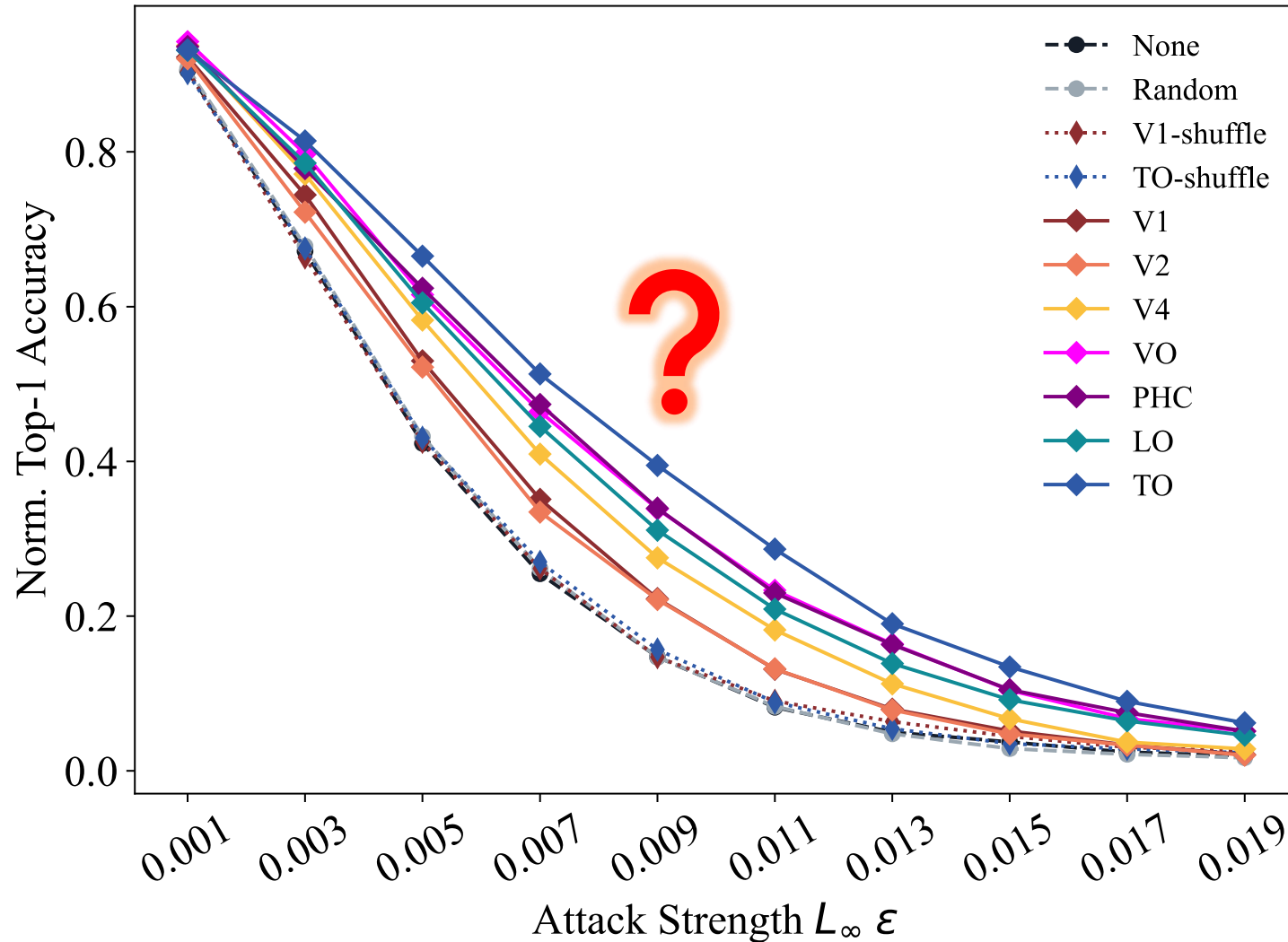*Attack: $l_\infty$-based PGD attack*

*Task: Image Classification*
*Dataset: ImageNet (Deng et al., 2009)*
*Attack: $l_\infty$-based PGD attack*

Task: Image Classification
Dataset: ImageNet *(Deng et al., 2009)*
Attack: $l_\infty$-based PGD attack

*Task: Image Classification*
*Dataset: ImageNet (Deng et al., 2009)*
*Attack: $l_\infty$-based PGD attack*

☑ **Neural guidance improves robustness** (max: 22% accuracy increase)
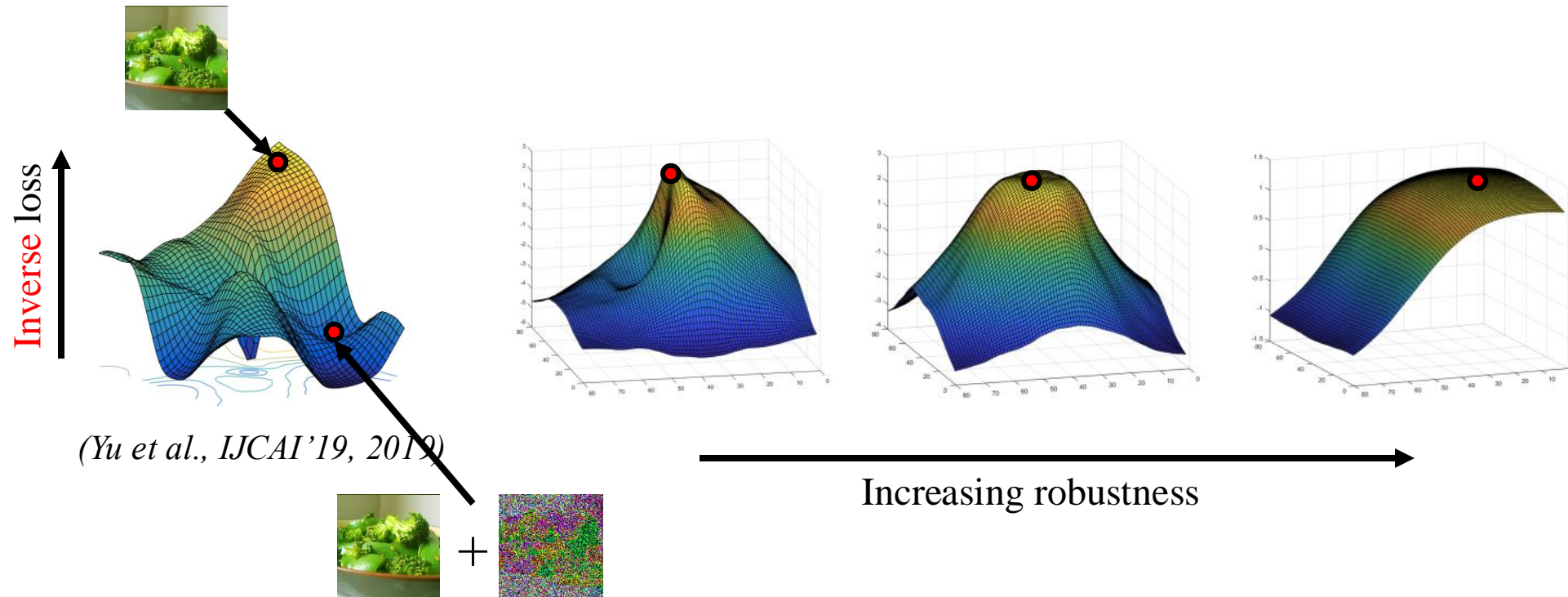
☑ **There exists a hierarchy of improvement's magnitude**

☑ *Replicated across datasets, tasks, attacks…*

- CIFAR-100
- MSCOCO
- Image Captioning

- $L_\infty$ FGSM
- Auto-Attack (APGD-CE, APGD-T, FAB square)
- $L_2$ FGM
- $L_2$ Deepfool

- Robust DNNs have smoother output surfaces



(Yu et al., IJCAI'19, 2019)

Inverse loss

Increasing robustness

- Robust DNNs have smoother output surfaces



*None*

*V4-guided*

*TO-guided*

- Robust DNNs have smoother output surfaces

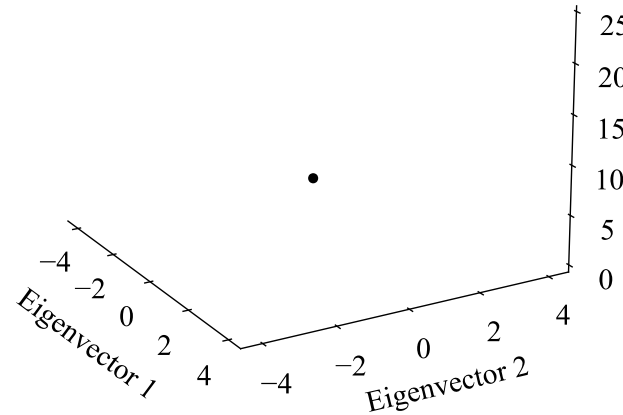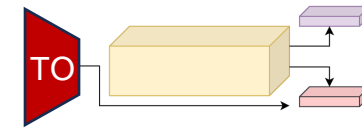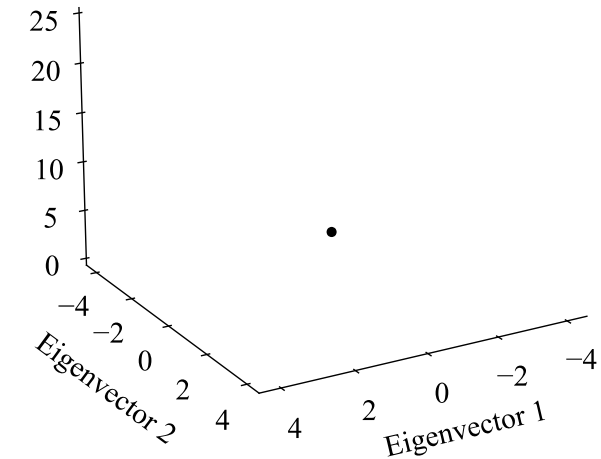- Robust DNNs have smoother output surfaces -- neurally-guided DNNs are indeed smoother!

- Conventional DNNs usually develop highly homogenous output surfaces



*"Transfer attack"*

**Adversarial examples are transferable across:**
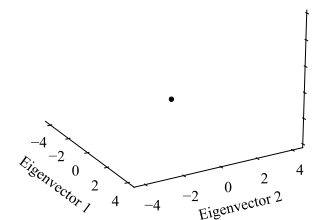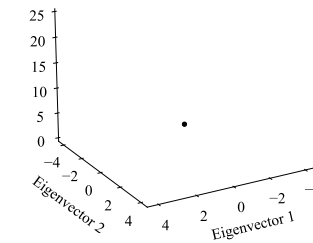
- Architectures *(Liu et al., 2017)*

- ML techniques *(Papernot et al., 2016)*

- Training datasets *(Lu et al., 2020)*

- Tasks *(Richards et al., 2021)*

- Conventional DNNs usually develop highly homogenous output surfaces

- Conventional DNNs usually develop highly homogenous output surfaces
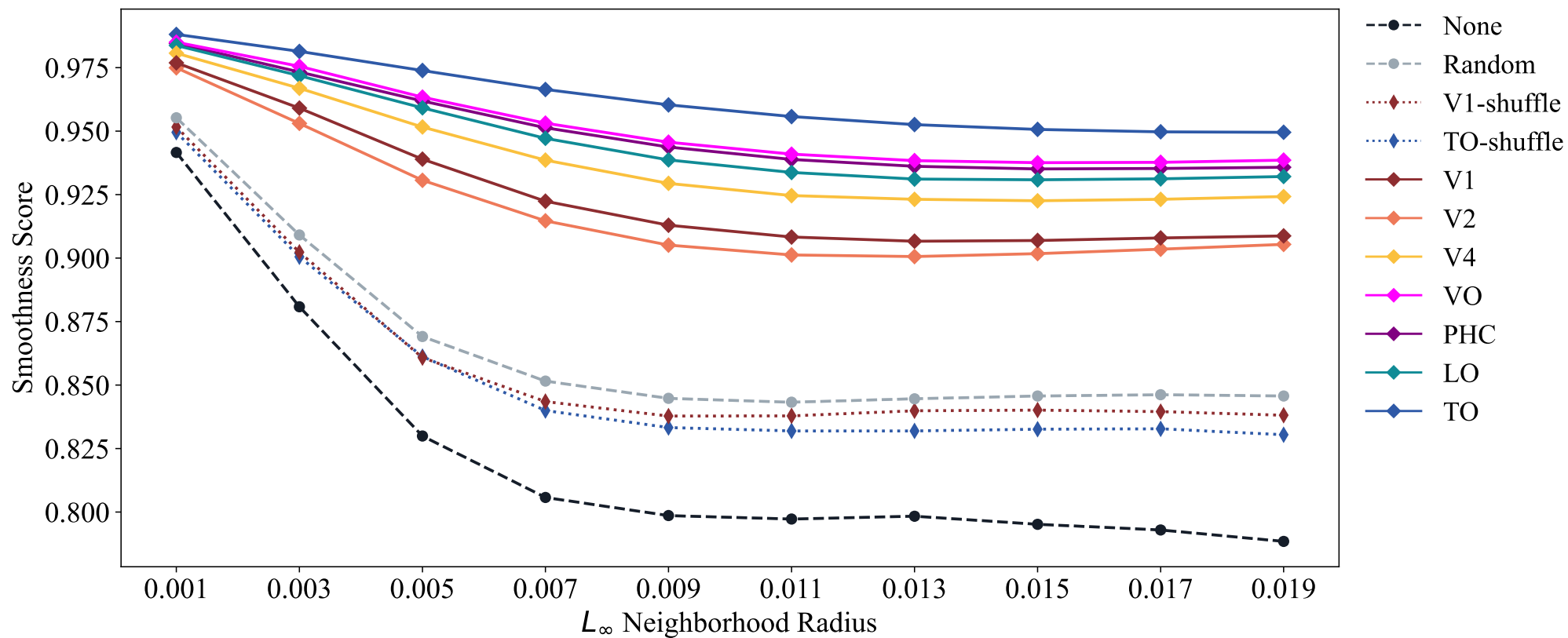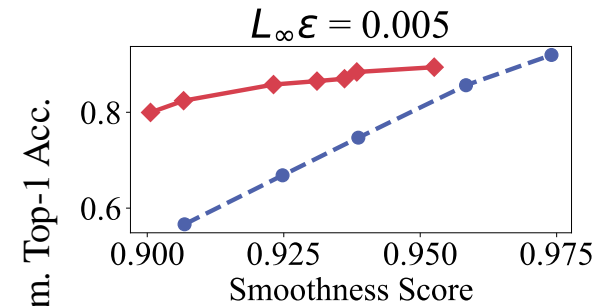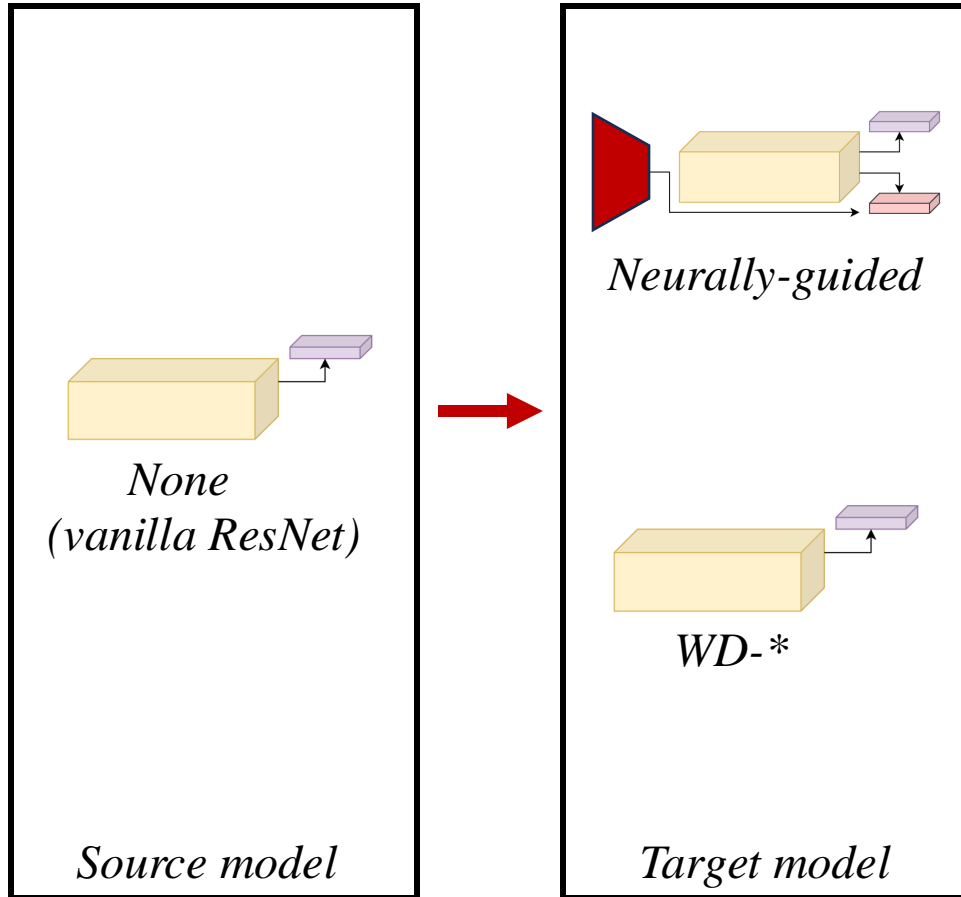
- Robust DNNs have smoother output surfaces -- neurally-guided DNNs are indeed smoother!

- Conventional DNNs usually develop highly homogenous output surfaces -- but neurally-guided DNNs have distinct surfaces!

- **Representational space**



Representational space                    Output surface

- Representational spaces of neurally guided DNNs (Representational Similarity Analysis) *(Kriegeskorte et al., Front. Syst. Neurosci., 2008)*



*(Nilli et al., Plos. Comp. Bio., 2014)*

# Neural guidance leads to distinct representational space

- Representational spaces of neurally guided DNNs (Representational Similarity Analysis) *(Kriegeskorte et al., Front. Syst. Neurosci., 2008)*

- **Representational spaces** of neurally guided DNNs are **distinct** from conventionally trained ones

- Robust DNNs have smoother output surfaces -- <span style="color:red">neurally-guided DNNs are indeed smoother!</span>

- Conventional DNNs usually develop highly homogenous output surfaces -- <span style="color:red">but neurally-guided DNNs have distinct surfaces!</span>
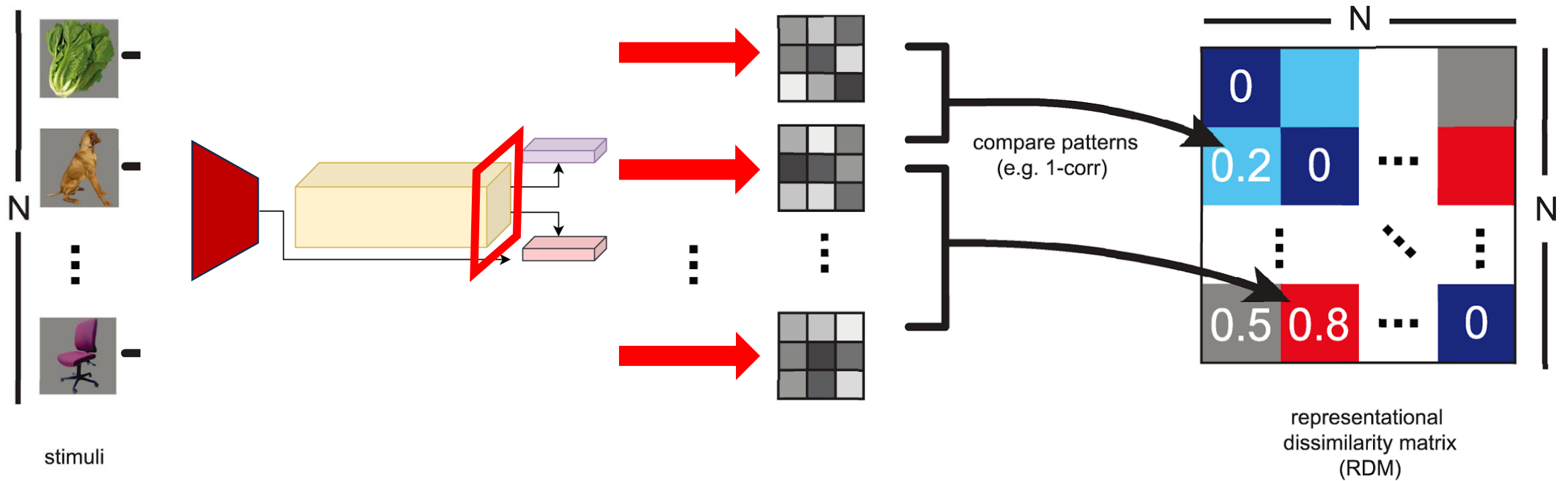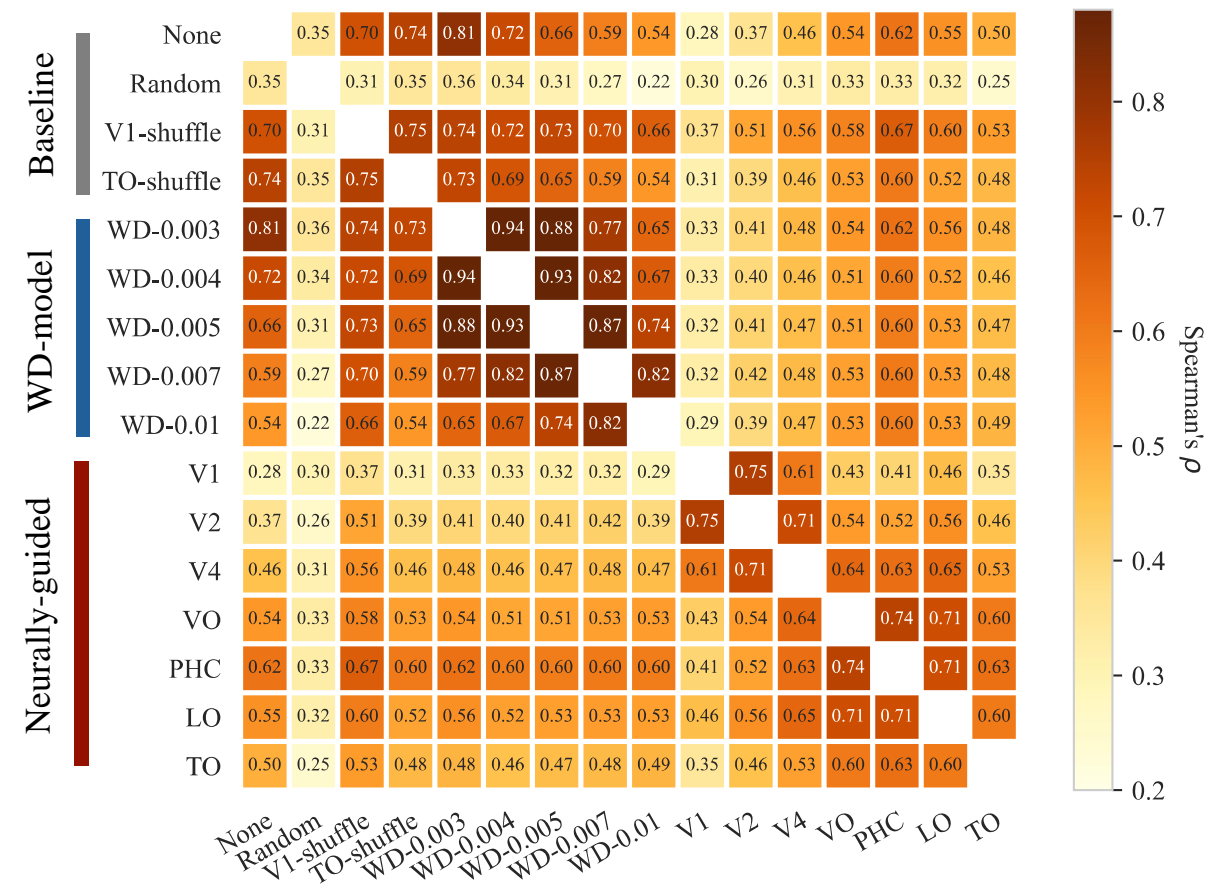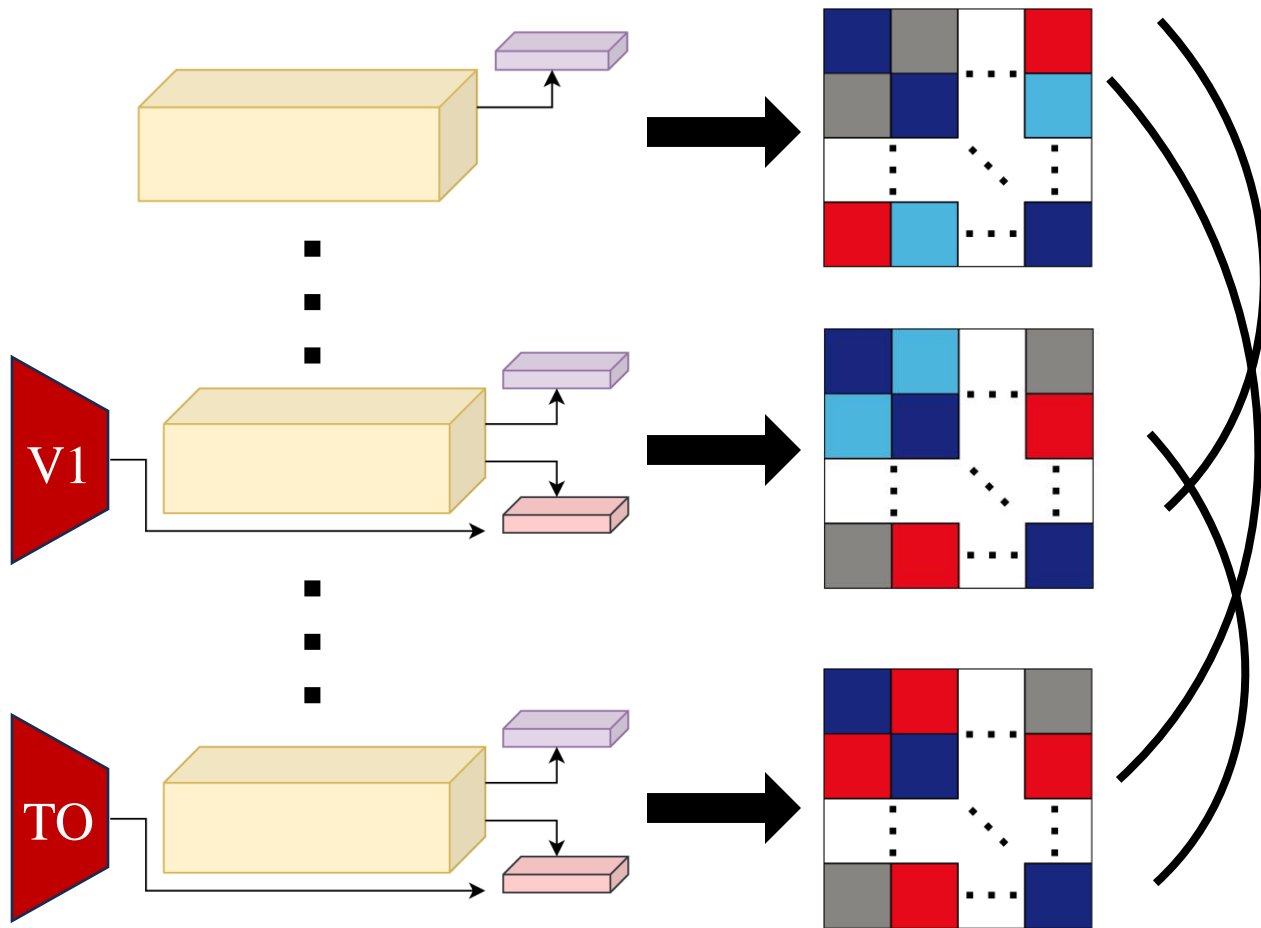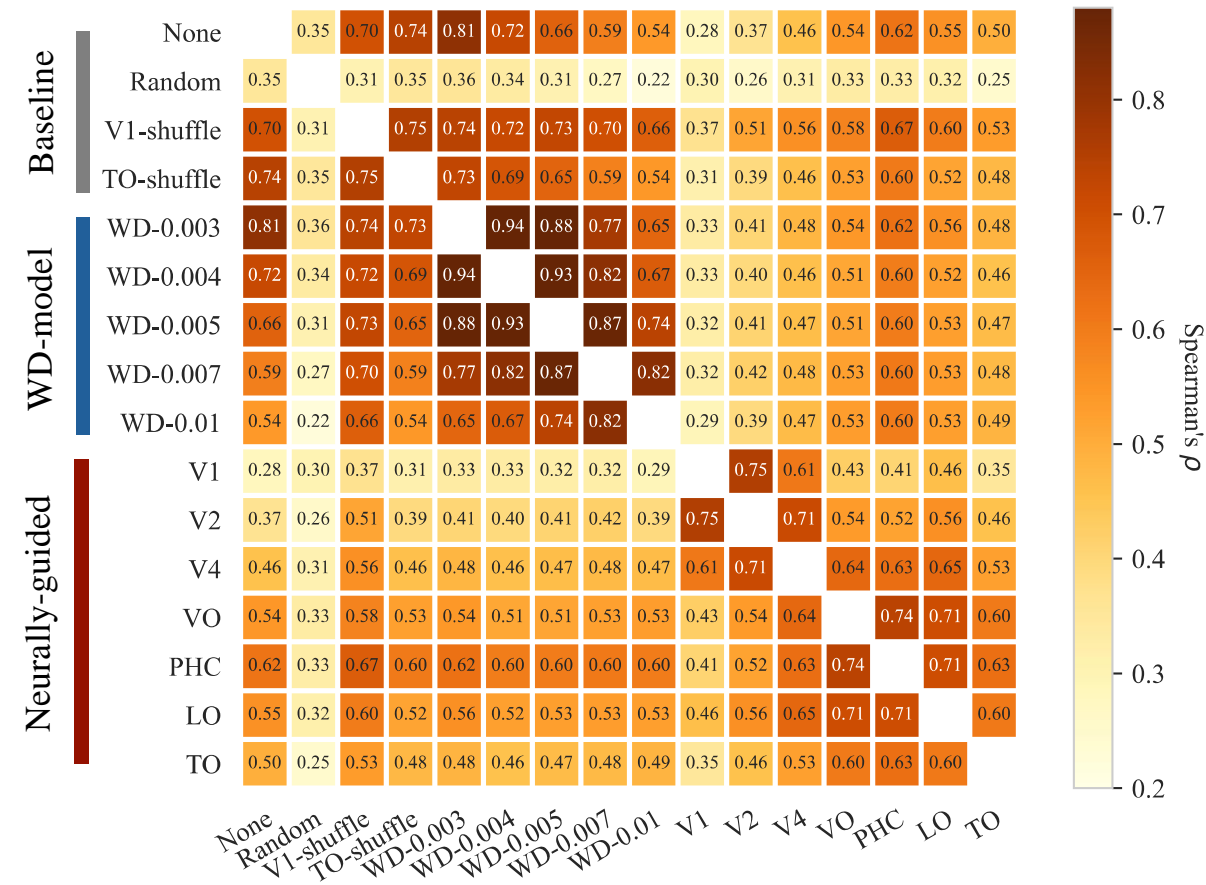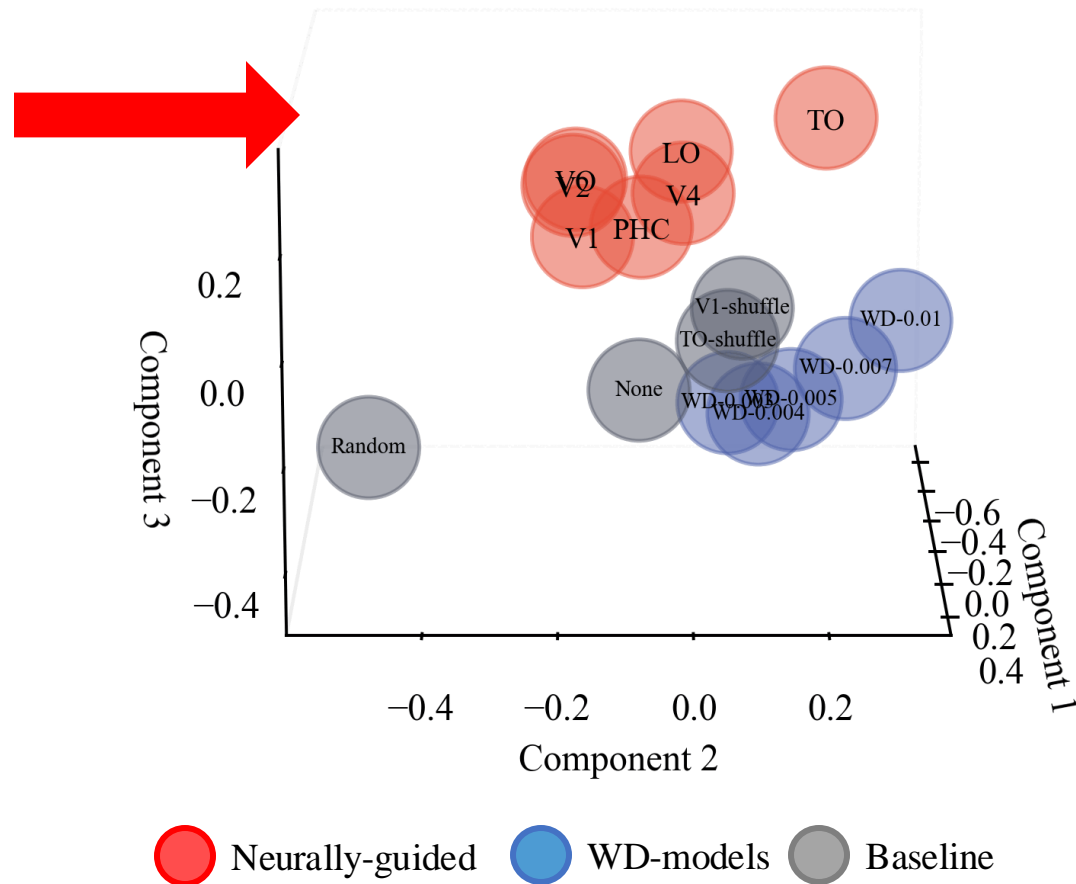
- Representational space -- <span style="color:red">neurally-guided DNNs developed distinct and better representational geometry!</span>

# Conclusion & Discussion

- We found:
  - **Hierarchical improvements** in DNN adversarial robustness with **neural guidance**
  - Neurally-guided DNNs developed **distinct** and **hierarchically smoother output surfaces** directly contributing to robustness
  - Neurally-guided DNNs developed **distinct representational spaces**
  - Neurally-guided DNNs are **progressively more shape biased**.

- Implications:
  - Robustness emerges from the evolving representational space along the ventral visual stream
  - Potential for understanding human representational space and advancing DNN architectural developments

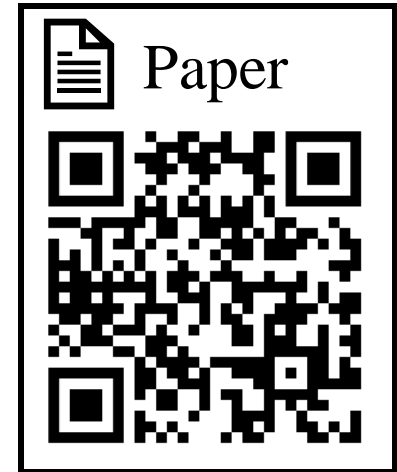📄 Paper

# Thank you! Questions?

Zhenan Shao

*zhenans2@illinois.edu*

Linjian Ma

Bo Li

Diane M. Beck

Paper