# Does Leveraging the Human Ventral Visual Stream Improve Neural Network Robustness?

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

BIAI 03/21/2024

[1]Department of Psychology, UIUC

[2]Department of Computer Science, UIUC

[3]Department of Computer Science, UChicago

# Introduction
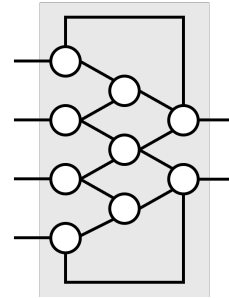## Robust human vision and vulnerable machine vision

Human visual perception achieves various invariance

*Biederman & Cooper, 1991;*
*Cave, Bost & Cobb, 1996;*
*Biederman & Gerhardstein, 1993;*
*...*

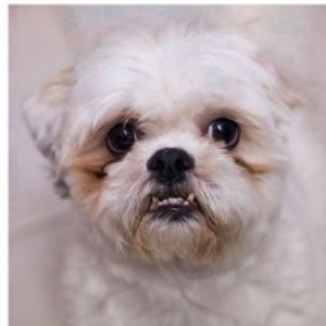Even imperceptible perturbations can lead to wrong prediction

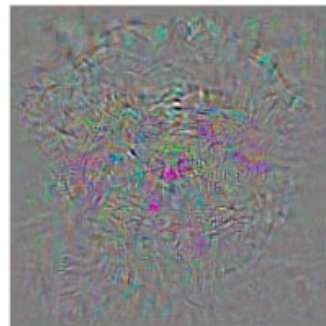*Szegedy et al., 2014;*
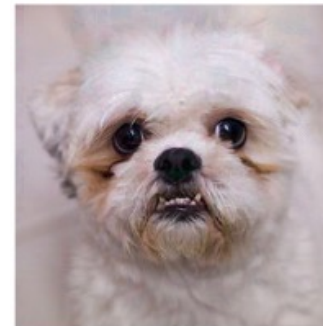*Carlini & Wagner, 2017;*
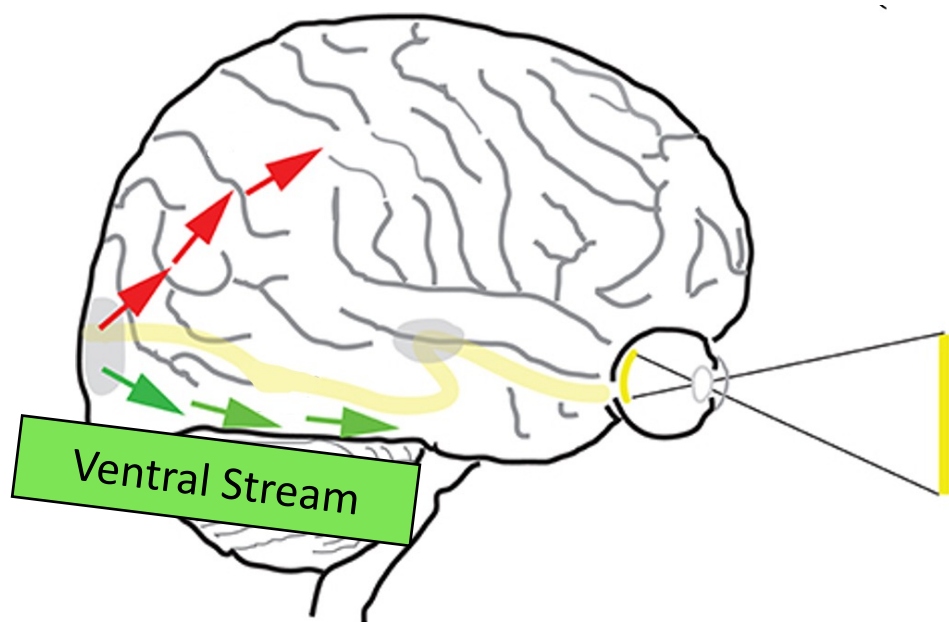*Kurakin et al., 2017;*
*...*

dog

Ostrich



(Szegedy et al., 2014)

# Introduction
## Achieving invariances along visual ventral stream

- The ventral visual stream forms a hierarchy, transitioning from basic visuals to more abstract and stable representations (Logothetis and Sheinberg, 1996, Zoccolan et al., 2007, Isik et al., 2014, Iordan et al., 2015…).
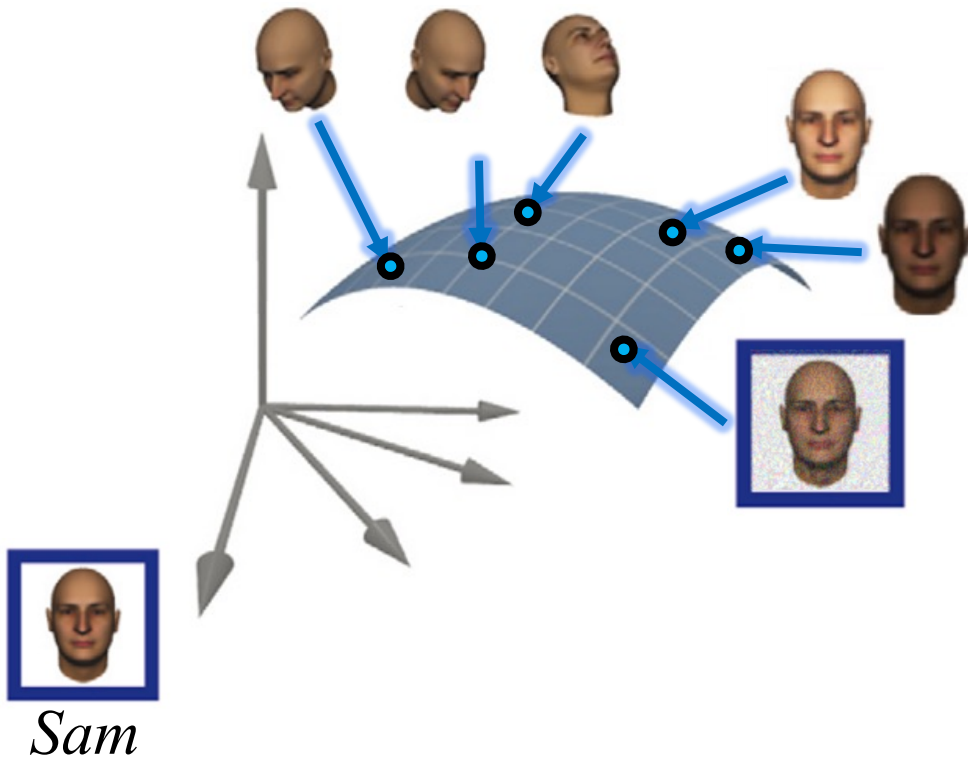- Evolving representations space achieved by separating object manifolds (Dicarlo & Cox, 2007).

Ventral Stream

*(Sheth & Young, 2016)*

# Introduction
## Achieving invariances along visual ventral stream

- The ventral visual stream forms a hierarchy, transitioning from basic visuals to more abstract and stable representations (Logothetis and Sheinberg, 1996, Zoccolan et al., 2007, Isik et al., 2014, Iordan et al., 2015…).
- Evolving representations space achieved by separating object manifolds (Dicarlo & Cox, 2007).



*Sam*

# Introduction
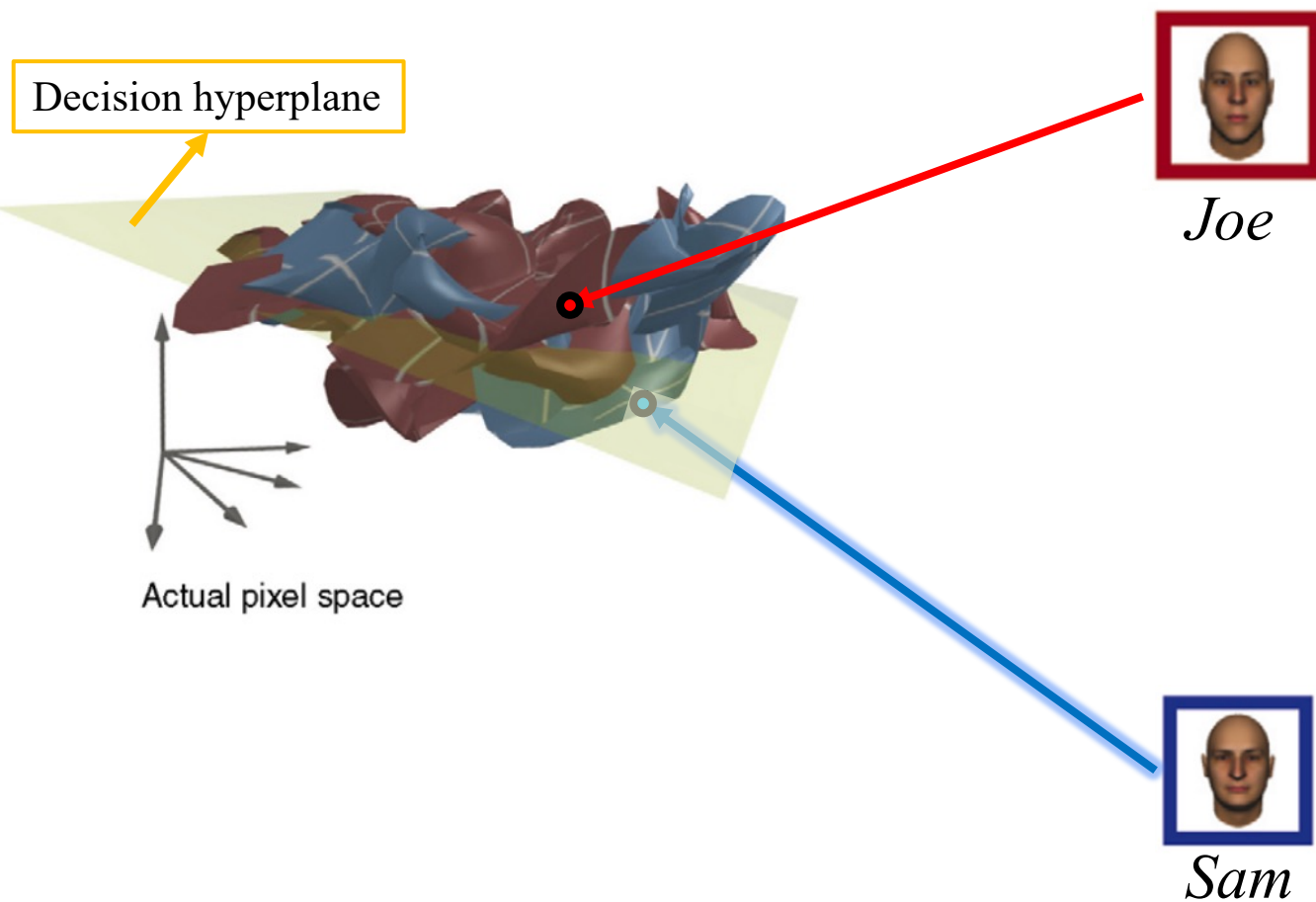## Achieving invariances along visual ventral stream

- The ventral visual stream forms a hierarchy, transitioning from basic visuals to more abstract and stable representations (Logothetis and Sheinberg, 1996, Zoccolan et al., 2007, Isik et al., 2014, Iordan et al., 2015…).
- Evolving representations space achieved by separating object manifolds (Dicarlo & Cox, 2007).



*Sam*

*Joe*

# Introduction

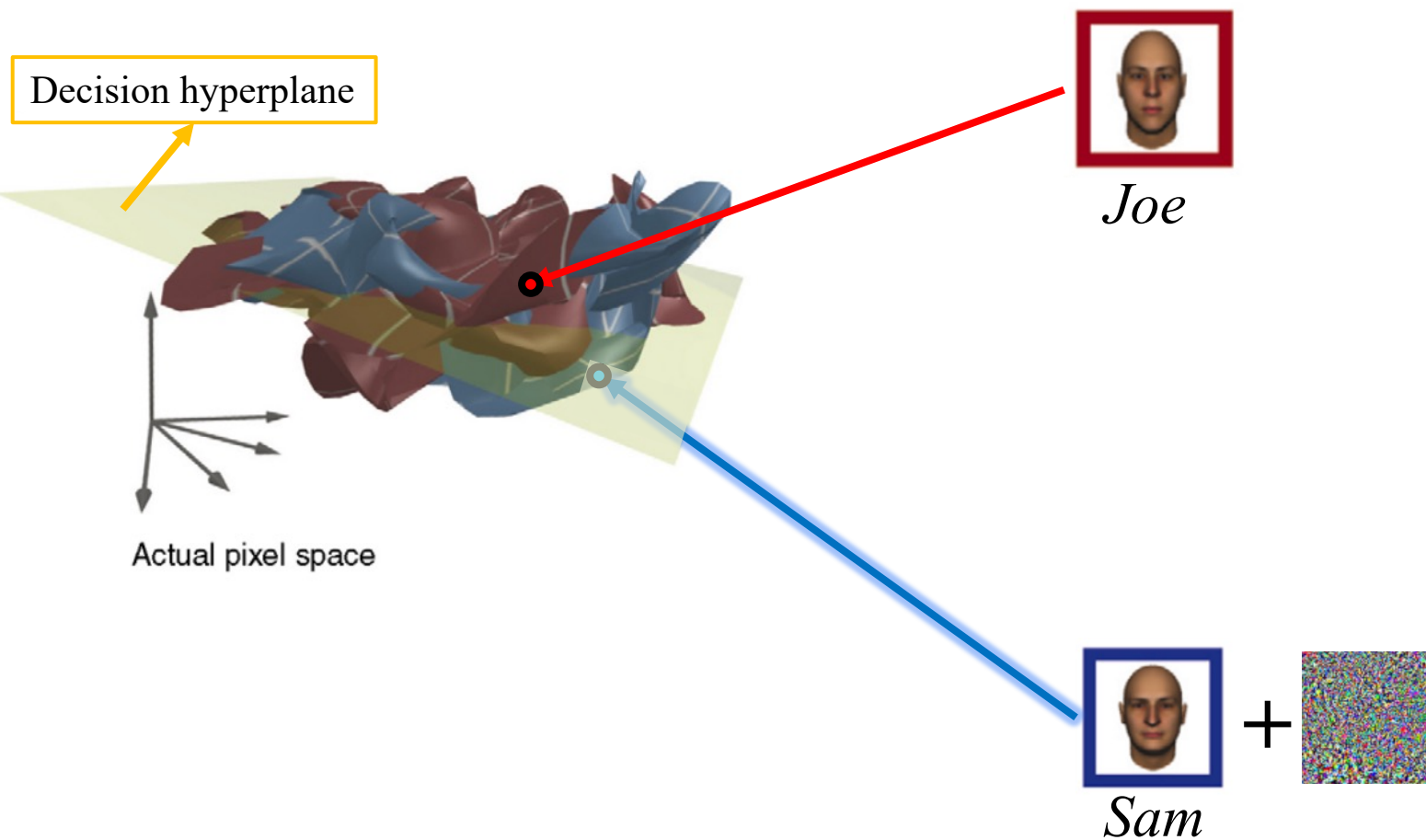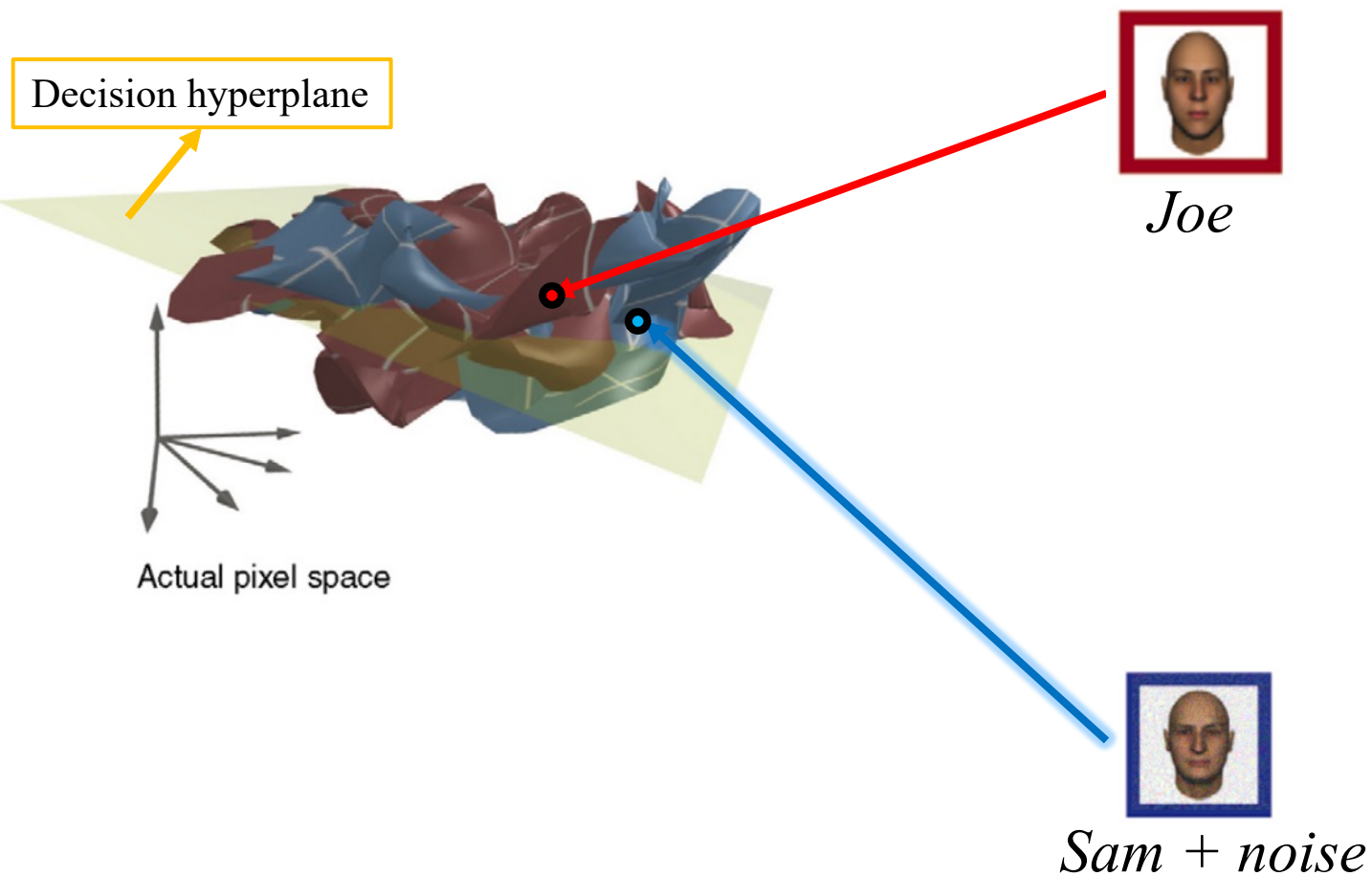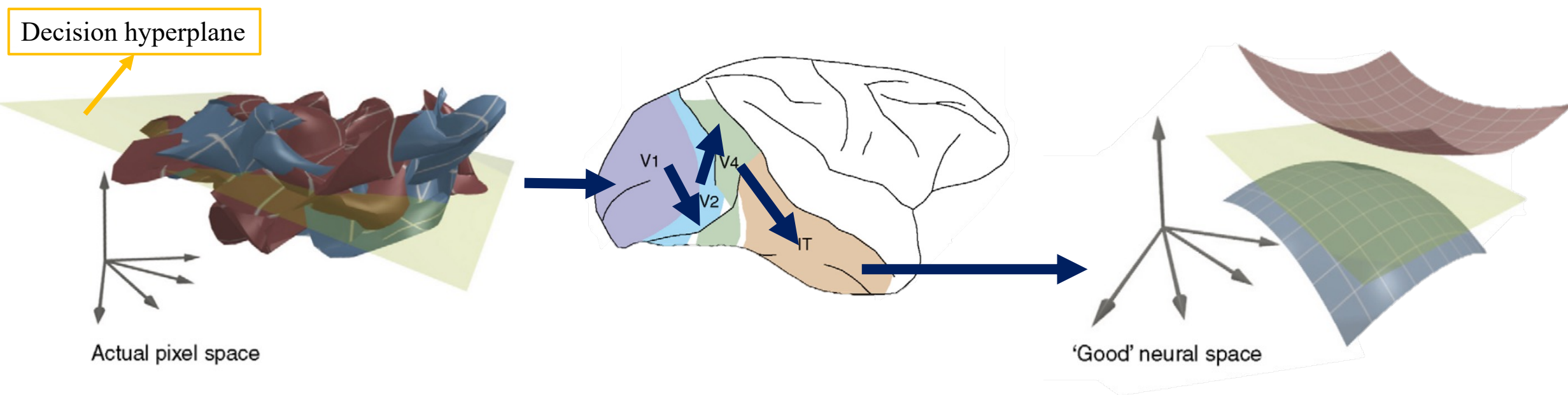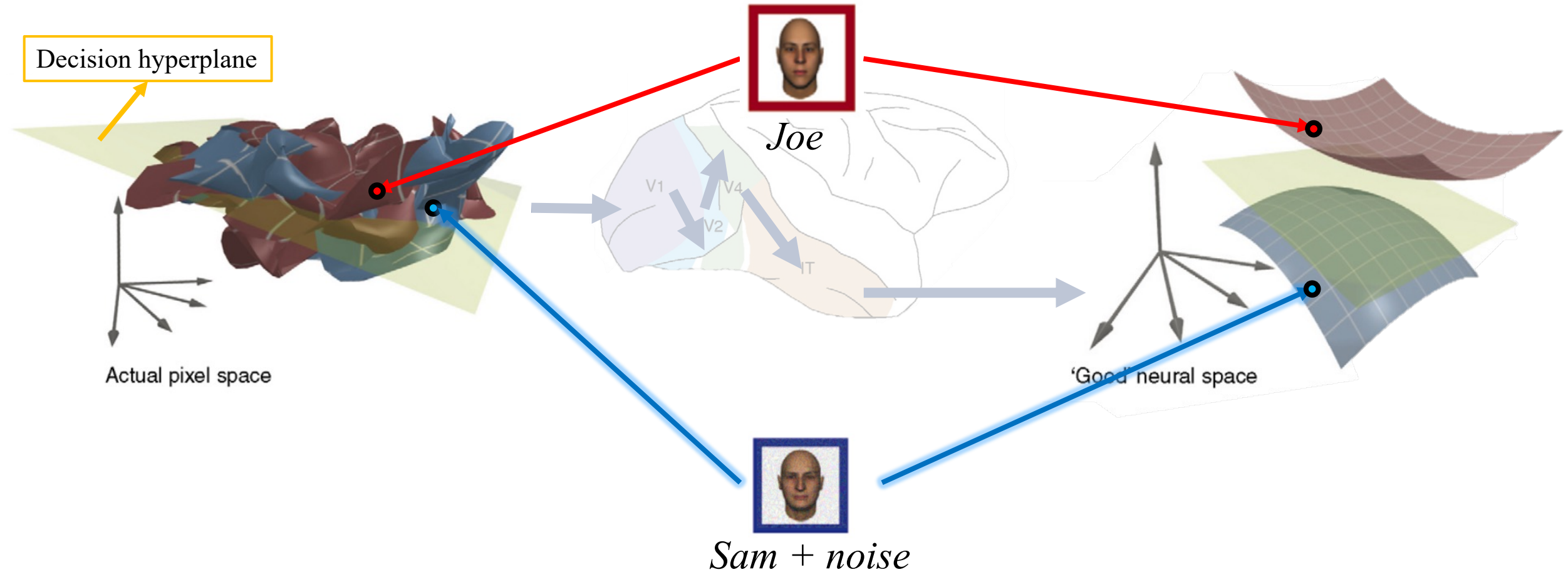## Achieving invariances along visual ventral stream

- Evolving representations achieved more separable object manifolds (Dicarlo & Cox, 2007).

# Introduction
Achieving invariances along visual ventral stream

- Evolving representations achieved more separable object manifolds (Dicarlo & Cox, 2007).

# Introduction
## Achieving invariances along visual ventral stream

- Evolving representations achieved more separable object manifolds (Dicarlo & Cox, 2007).



Decision hyperplane

Actual pixel space

*Joe*

*Sam + noise*

# Introduction
Achieving invariances along visual ventral stream

- Evolving representations achieved more separable object manifolds (Dicarlo & Cox, 2007).



Decision hyperplane

Actual pixel space

'Good' neural space

# Introduction
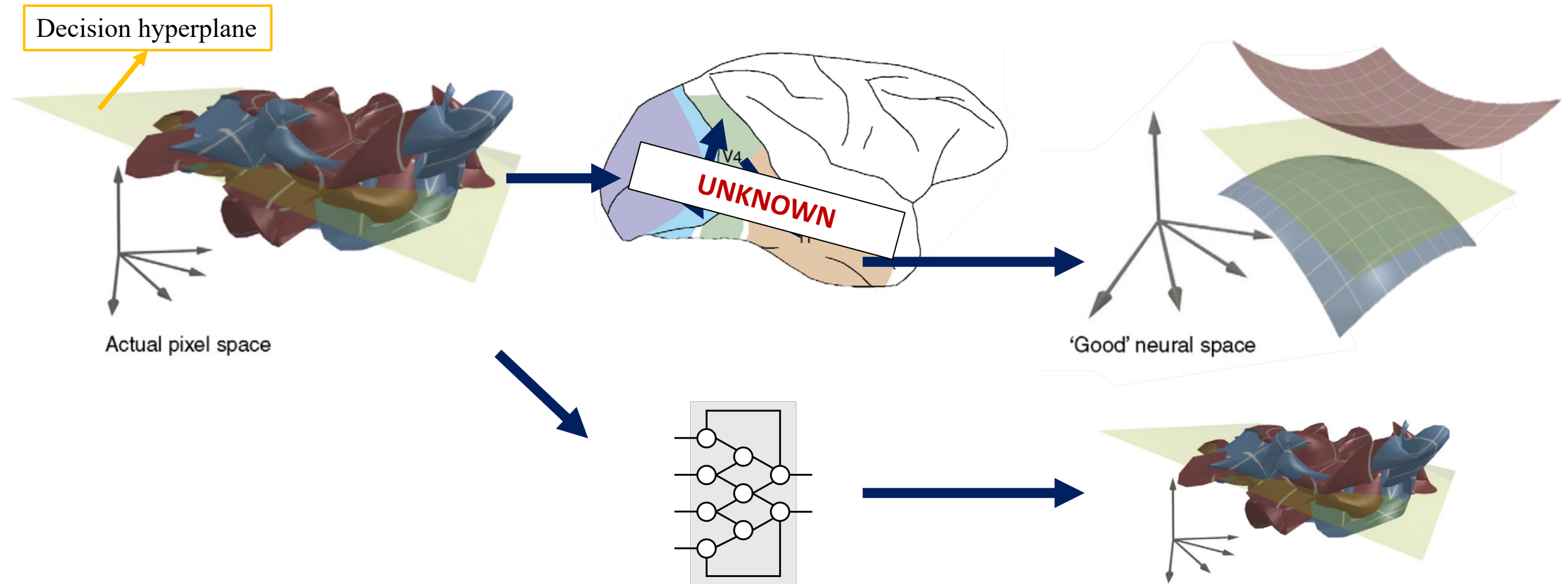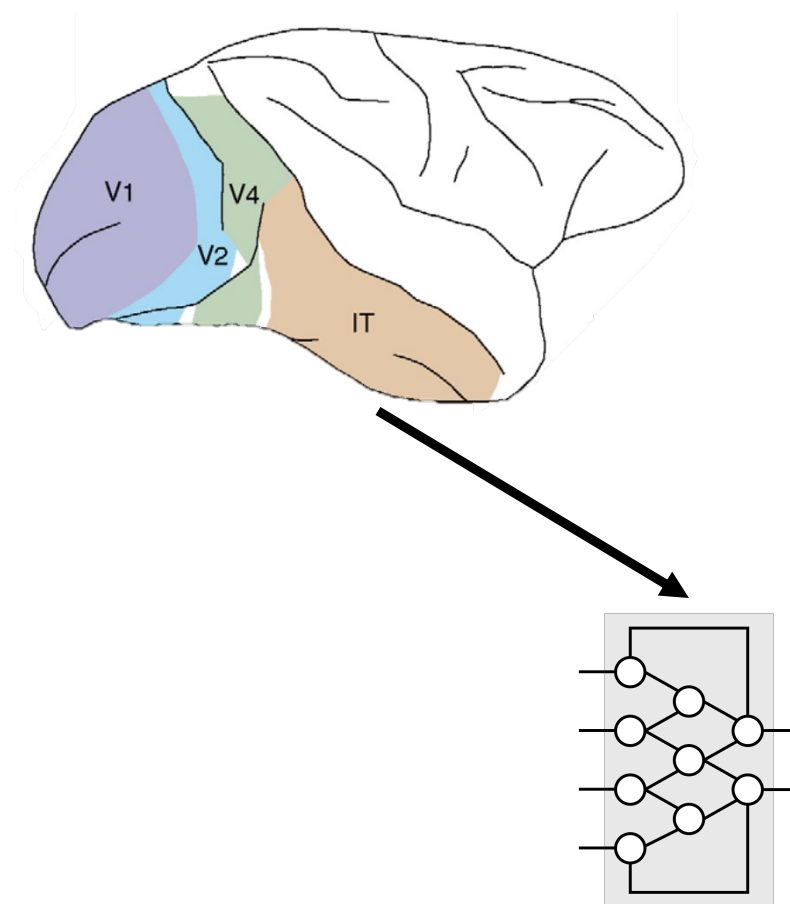## Achieving invariances along visual ventral stream

- Evolving representations achieved more separable object manifolds (Dicarlo & Cox, 2007).

# Introduction
## Achieving invariances along visual ventral stream

- Evolving representations achieved more separable object manifolds (Dicarlo & Cox, 2007).
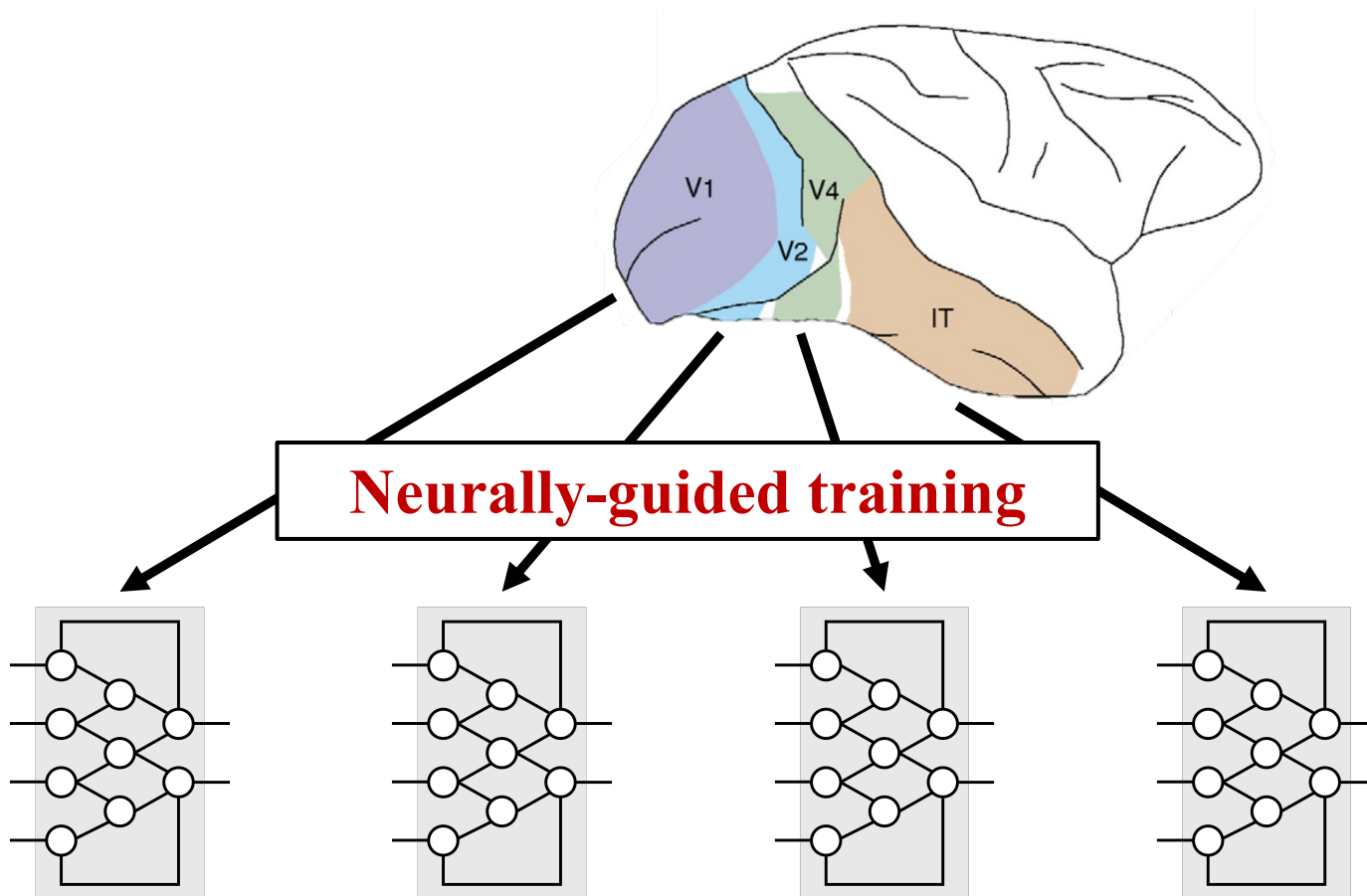
# Our question

1. Does training guided by human ventral cortex activity improve neural network robustness?

# Our question

1. Does training guided by human ventral cortex activity improve neural network robustness?

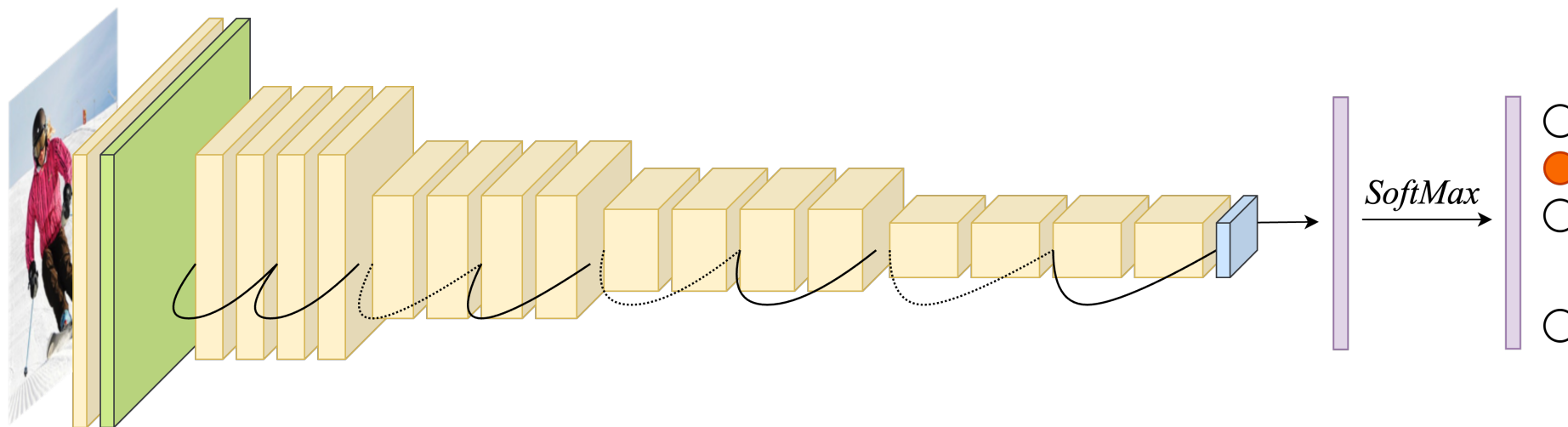2. Does such improvement increase as we ascend the ventral visual cortex?

# Method

Neurally-guided training

- DNN visual task training:

$$L = L_{task}$$
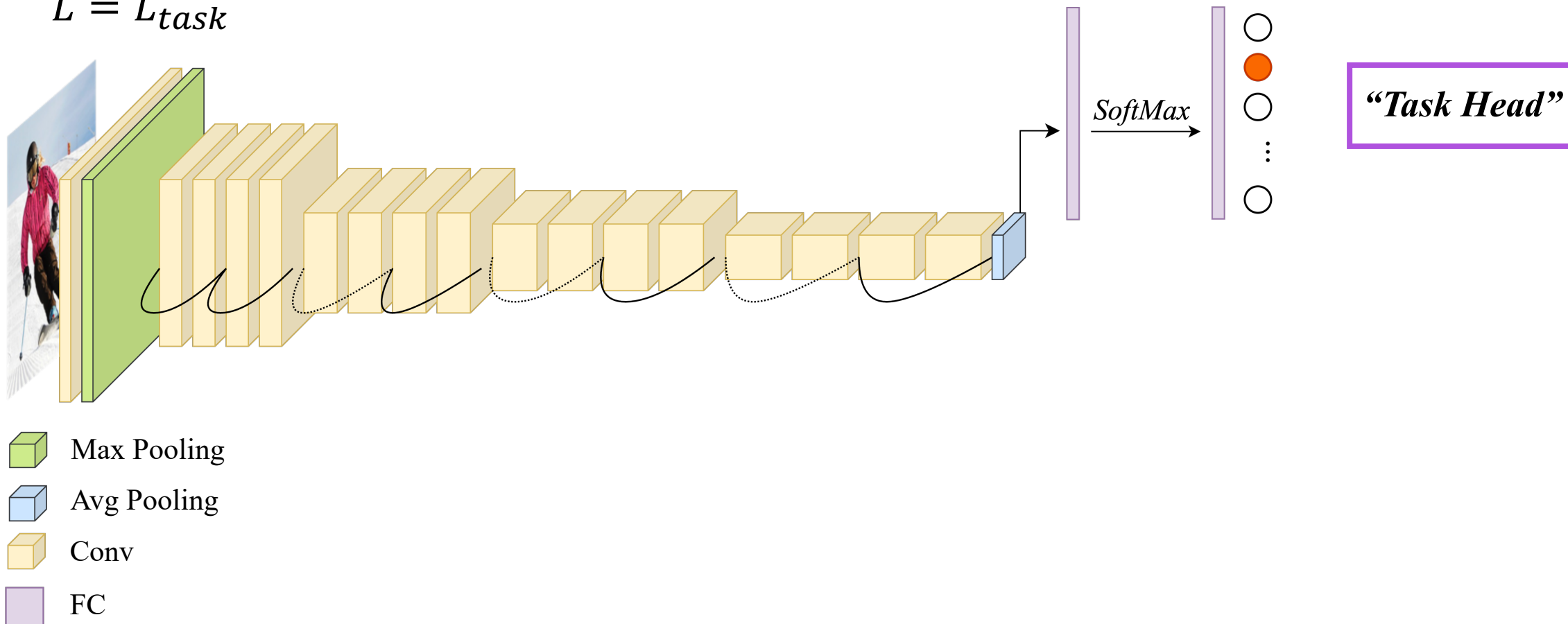


Max Pooling

Avg Pooling

Conv

FC

# Method
## Neurally-guided training

- DNN visual task training:

$$L = L_{task}$$



*SoftMax*

*"Task Head"*
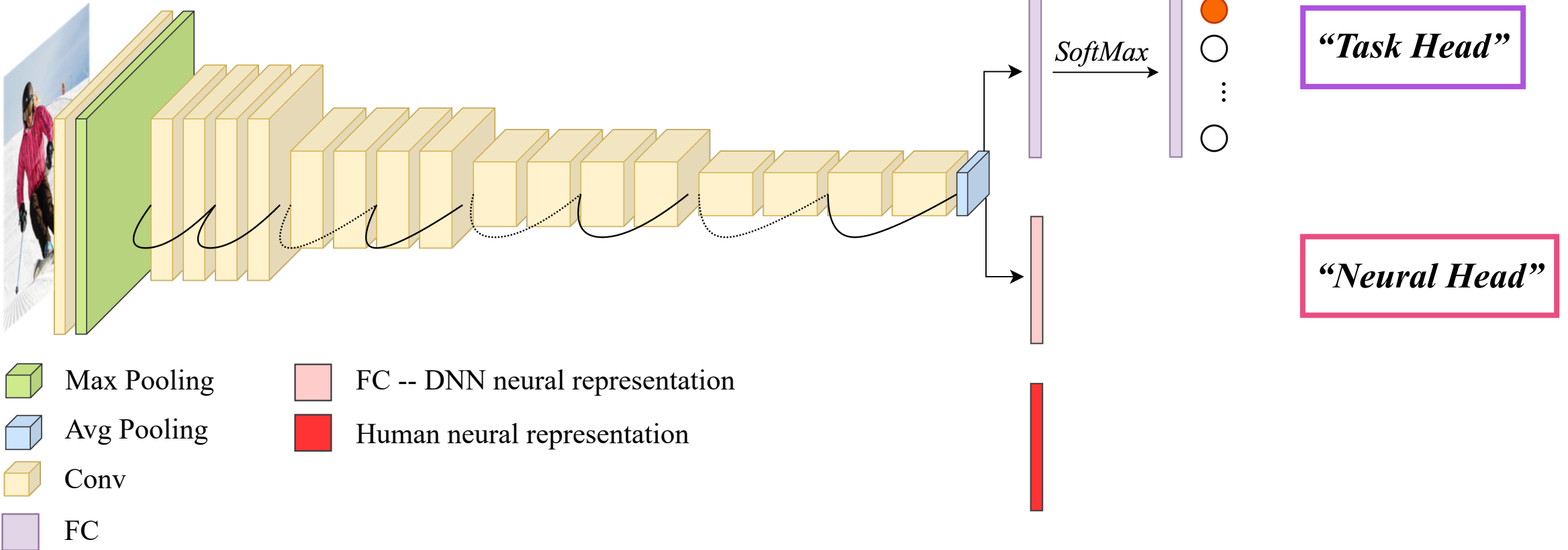
Max Pooling

Avg Pooling

Conv

FC

# Method
Neurally-guided training
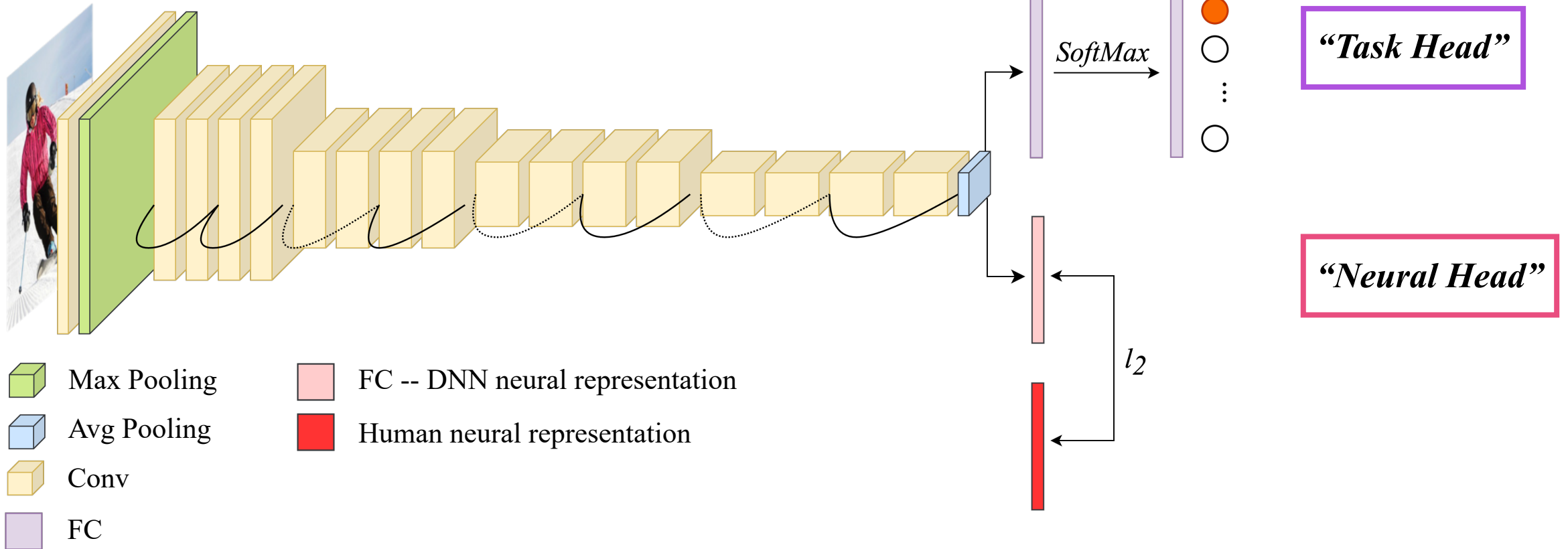
- DNN visual task training with Neural Guidance:

$$L = L_{task}$$

# Method
## Neurally-guided training

- DNN visual task training with Neural Guidance:

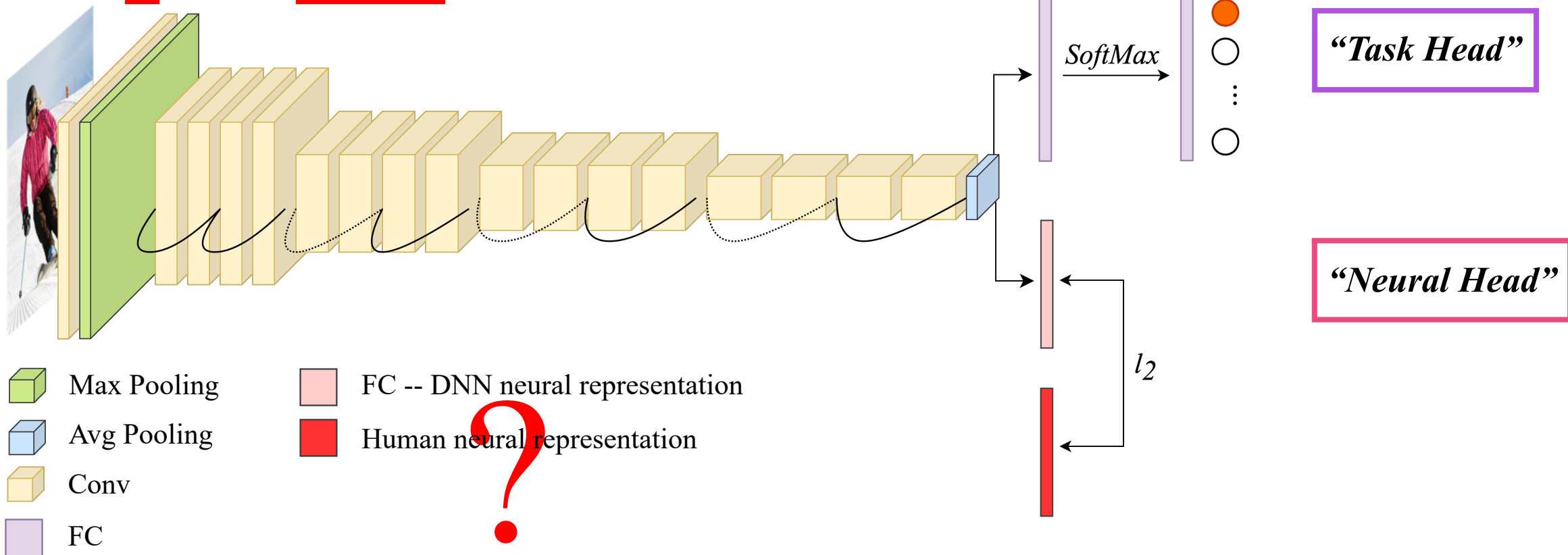$$L = L_{task} + ||R_{DNN} - R_{neural}||_2$$



*SoftMax*

*"Task Head"*

*"Neural Head"*

$l_2$

Max Pooling

Avg Pooling

Conv

FC

FC -- DNN neural representation

Human neural representation

# Method
## Neurally-guided training

- DNN visual task training with Neural Guidance:

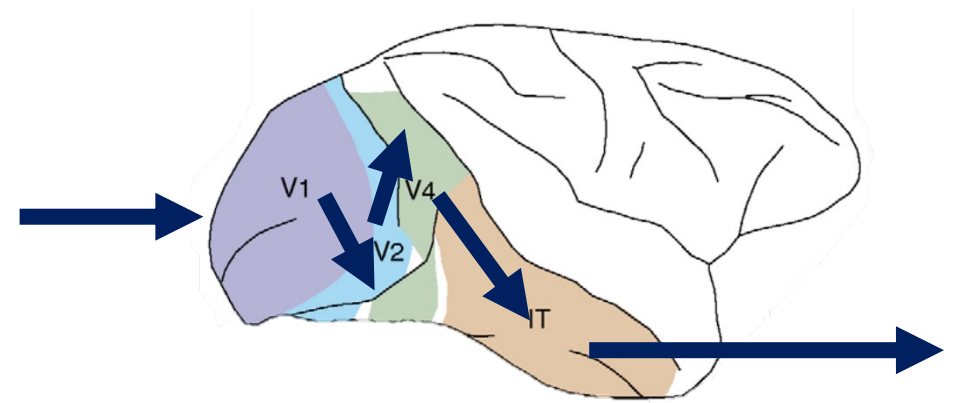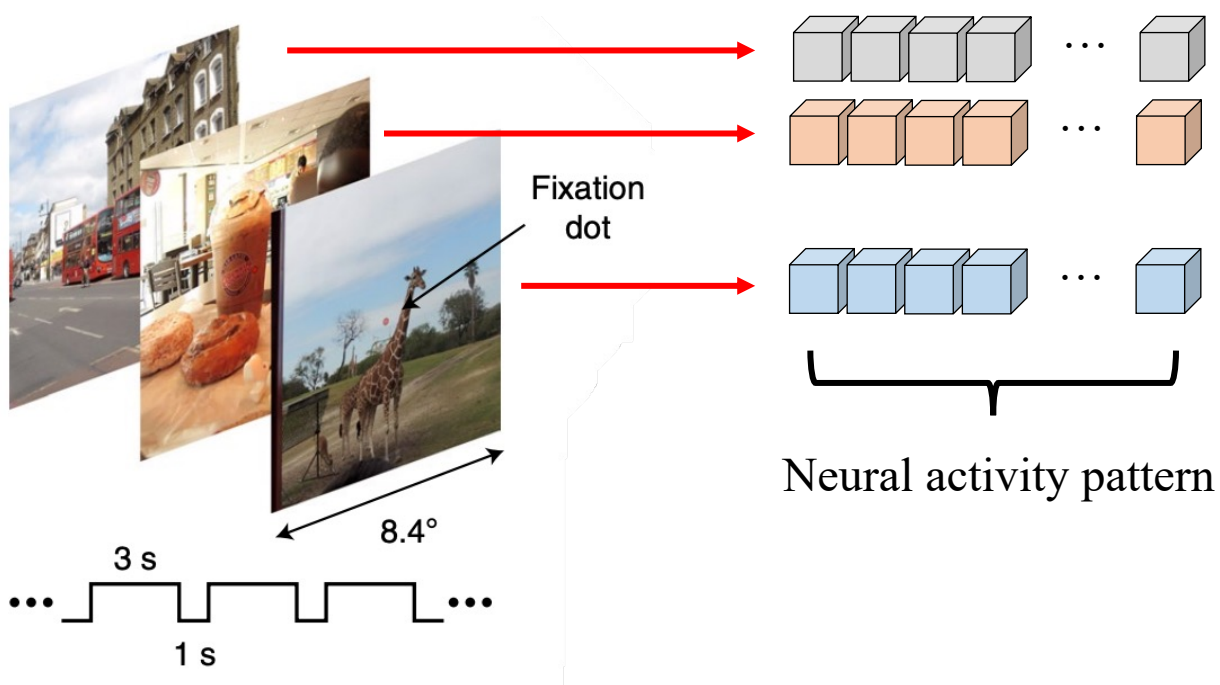$$L = \alpha L_{task} + (1 - \alpha) \, ||R_{DNN} - R_{neural} \, ||_2$$

*SoftMax*

*"Task Head"*

*"Neural Head"*

$l_2$

Max Pooling

Avg Pooling

Conv

FC

FC -- DNN neural representation

Human neural representation

# Method
## Neural data

- Each human subject viewed ~30,000 images (~9,000 unique).

- Brain activities were recorded with 7T fMRI.



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

# Method
## Neural data

- Each human subject viewed ~30,000 images (~9,000 unique).

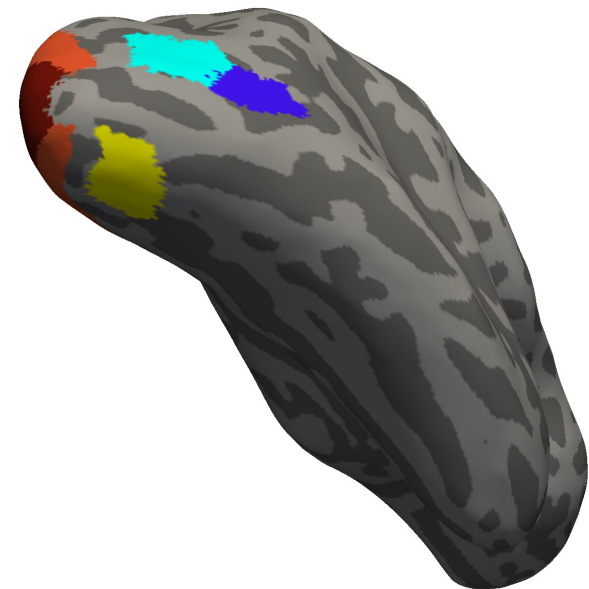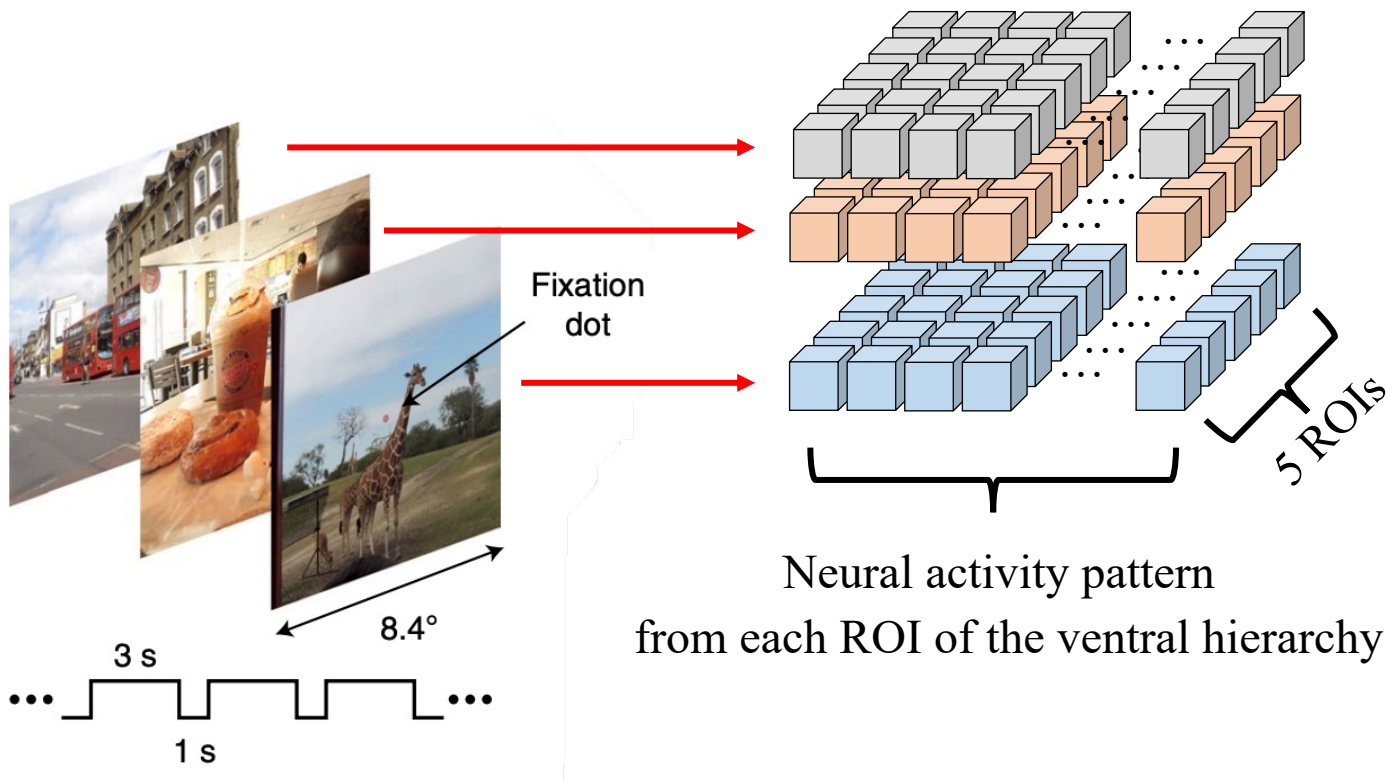- Brain activities were recorded with 7T fMRI.



Neural activity pattern

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

# Method
## Neural data

- Each human subject viewed ~30,000 images (~9,000 unique).

- Brain activities were recorded with 7T fMRI.
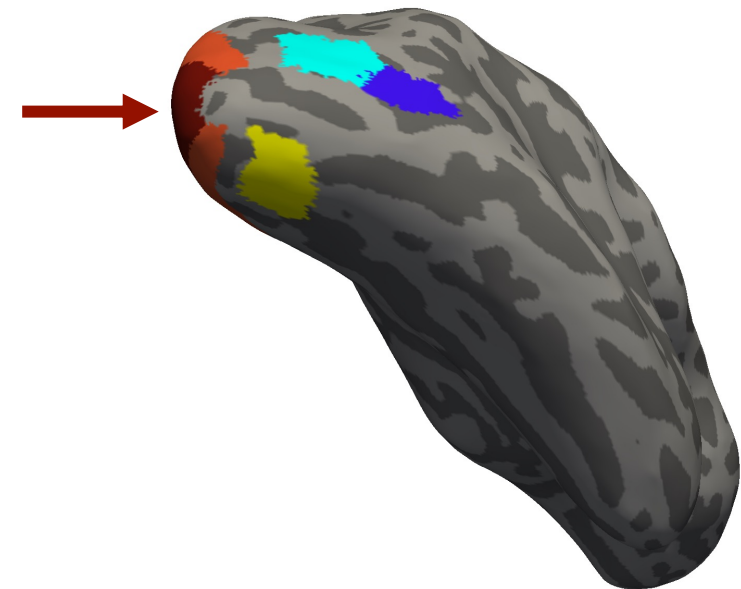
- 5 bilateral Regions of Interest (**ROIs**) were used



Fixation dot
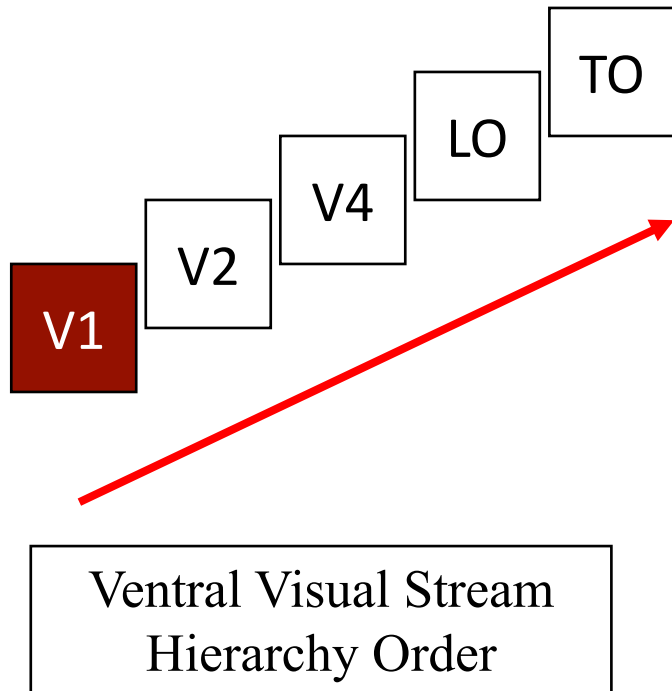
8.4°

3 s

1 s

Neural activity pattern
from each ROI of the ventral hierarchy

5 ROIs

Subject 1 ROIs
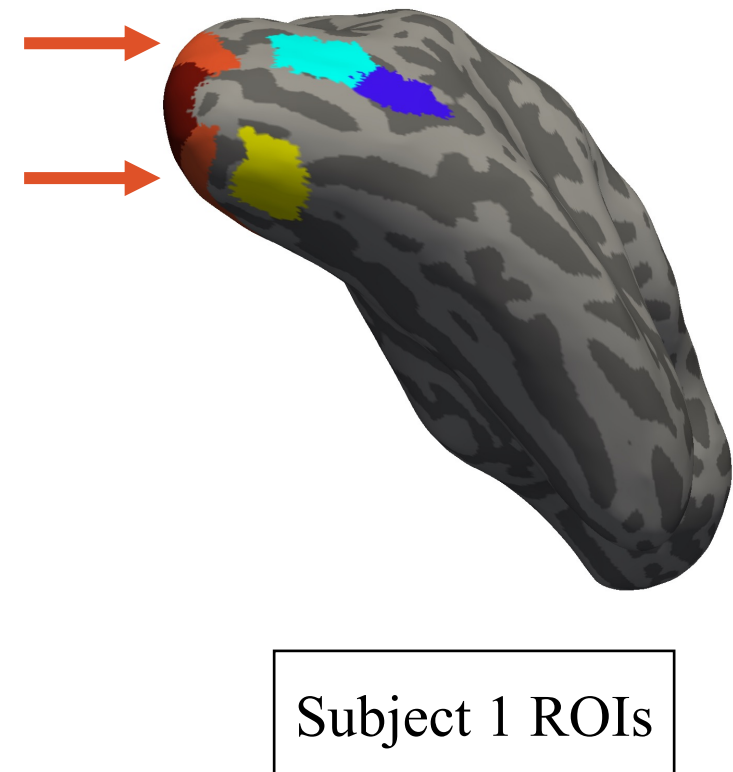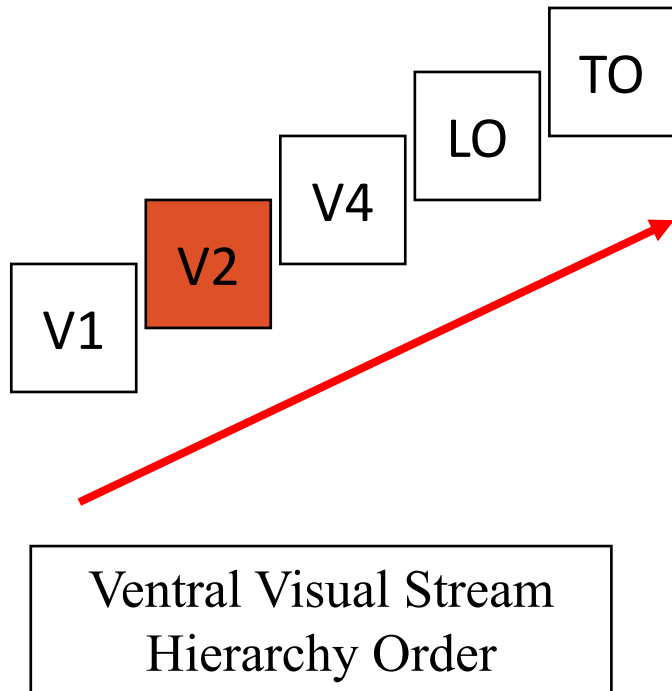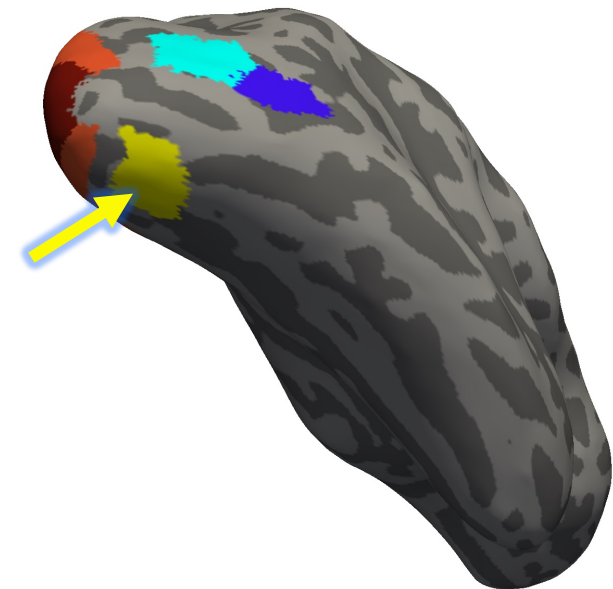
*(NSD, Allen et al., 2022)*

# Method
## Neural data

- Each human subject viewed ~30,000 images (~9,000 unique).

- Brain activities were recorded with 7T fMRI.

- 5 bilateral ROIs were used:



TO

LO

V4

V2

V1

Ventral Visual Stream
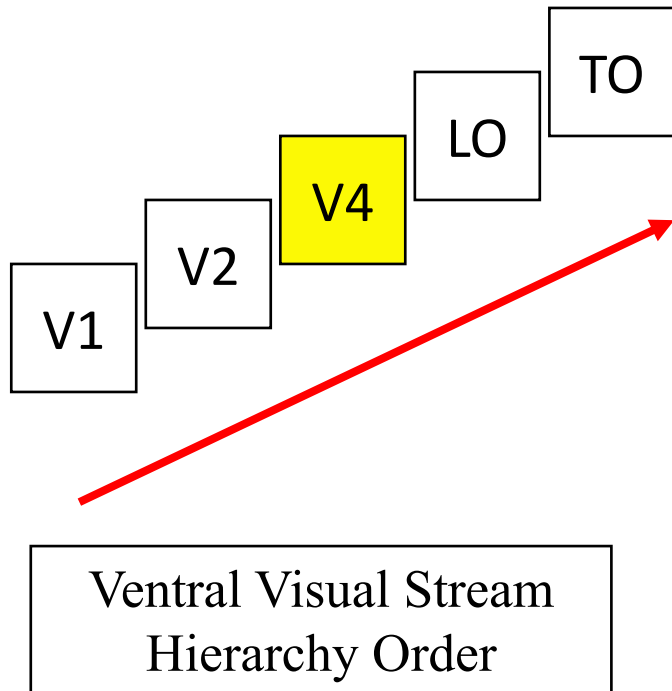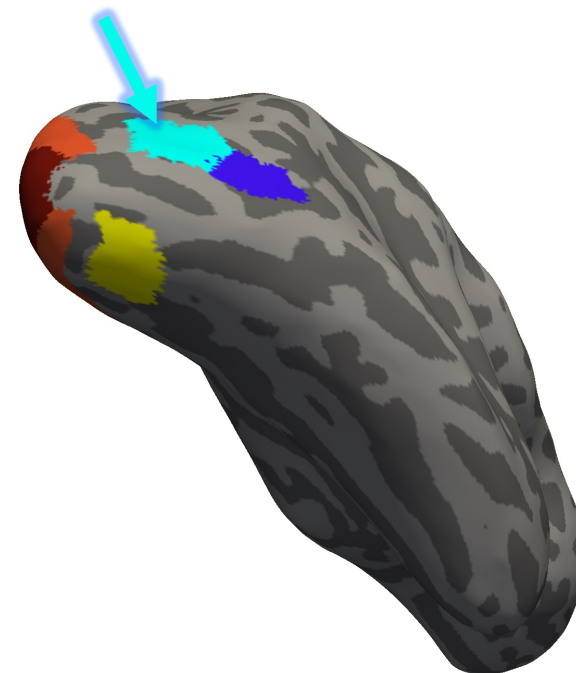Hierarchy Order

Subject 1 ROIs

# Method
## Neural data

- Each human subject viewed ~30,000 images (~9,000 unique).

- Brain activities were recorded with 7T fMRI.

- 5 bilateral ROIs were used:



TO

LO

V4

V2

V1

Ventral Visual Stream
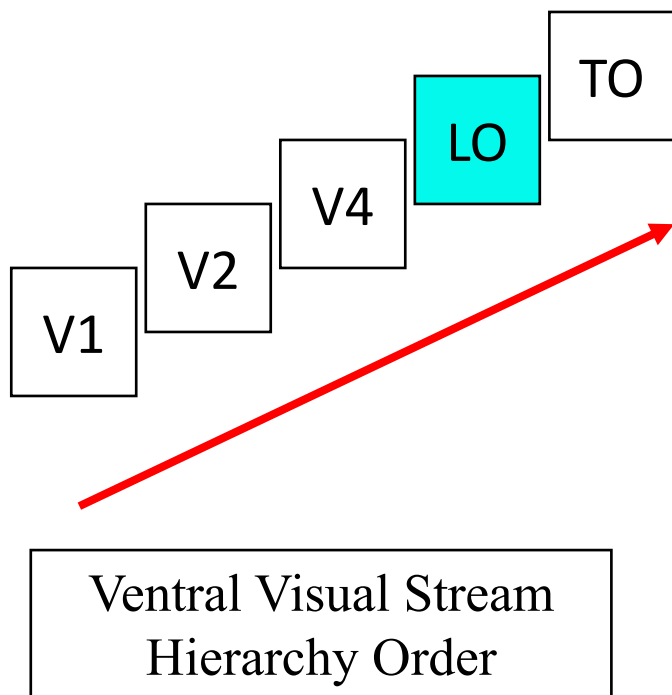Hierarchy Order

Subject 1 ROIs

# Method
## Neural data

- Each human subject viewed ~30,000 images (~9,000 unique).

- Brain activities were recorded with 7T fMRI.

- 5 bilateral ROIs were used:



| V1 | V2 | V4 | LO | TO |

Ventral Visual Stream
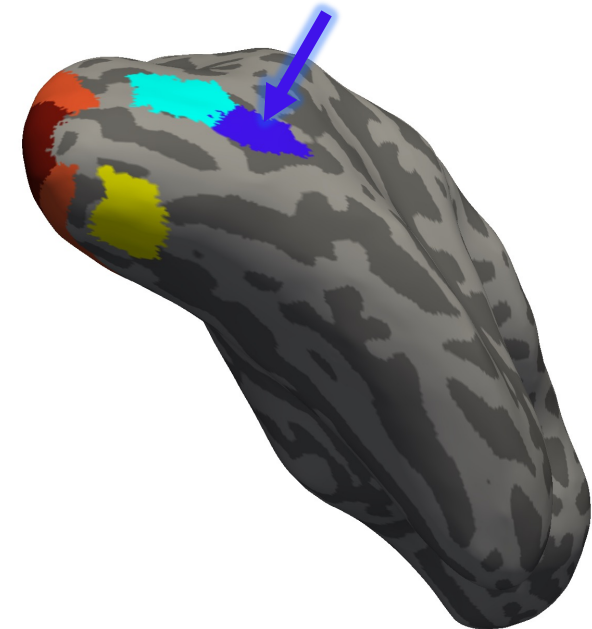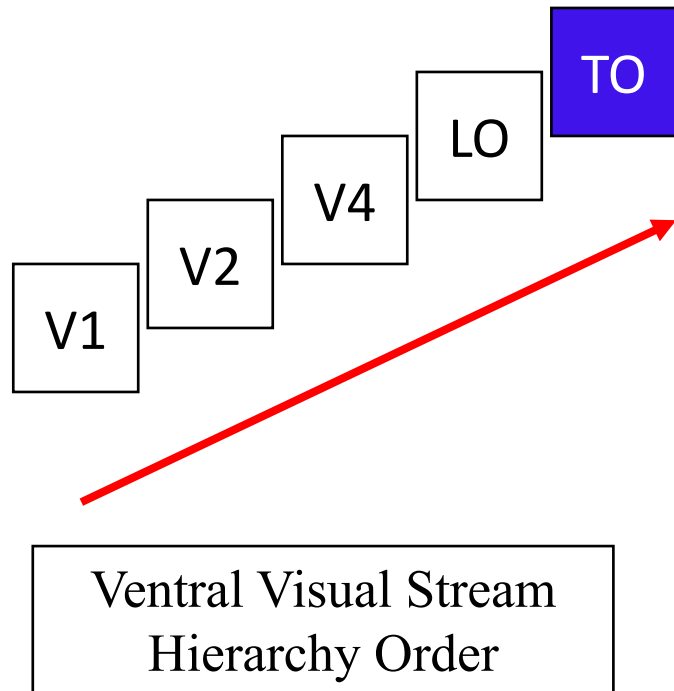Hierarchy Order

Subject 1 ROIs

# Method
## Neural data

- Each human subject viewed ~30,000 images (~9,000 unique).

- Brain activities were recorded with 7T fMRI.

- 5 bilateral ROIs were used:

TO

LO

V4

V2

V1

Ventral Visual Stream
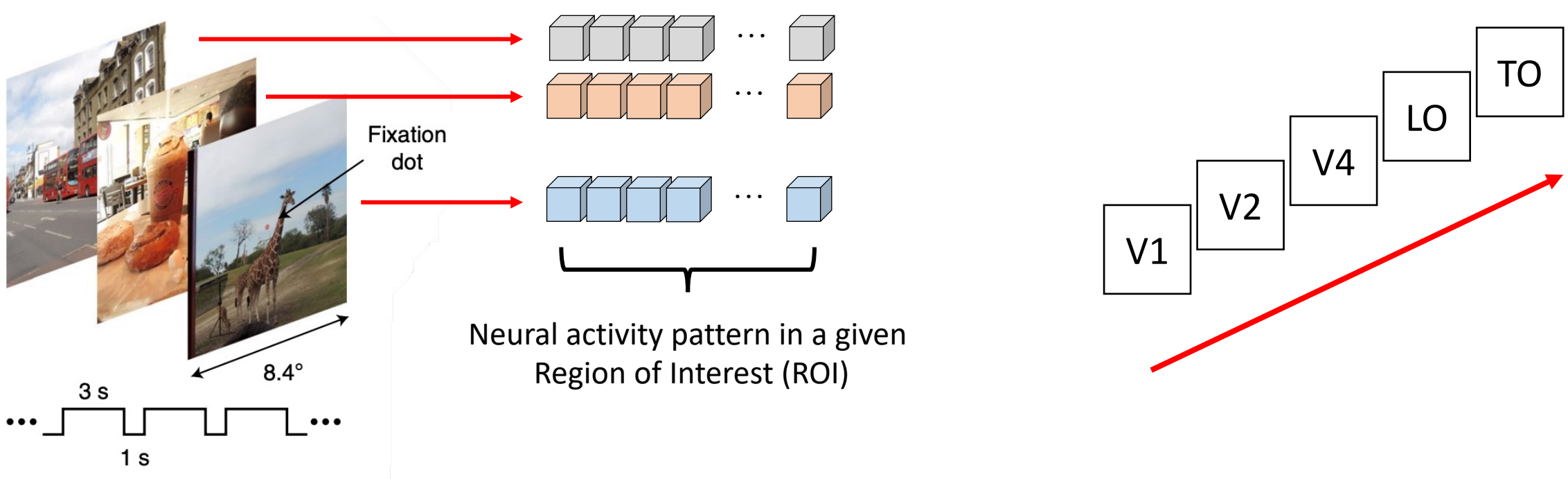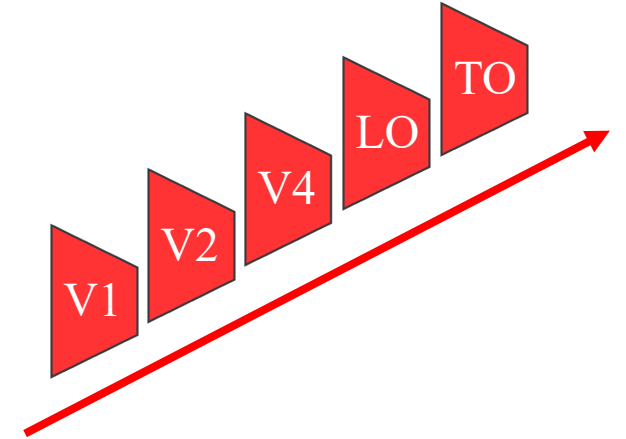Hierarchy Order

Subject 1 ROIs

# Method
## Neural data

- Each human subject viewed ~30,000 images (~9,000 unique).

- Brain activities were recorded with 7T fMRI.

- 5 bilateral ROIs were used:



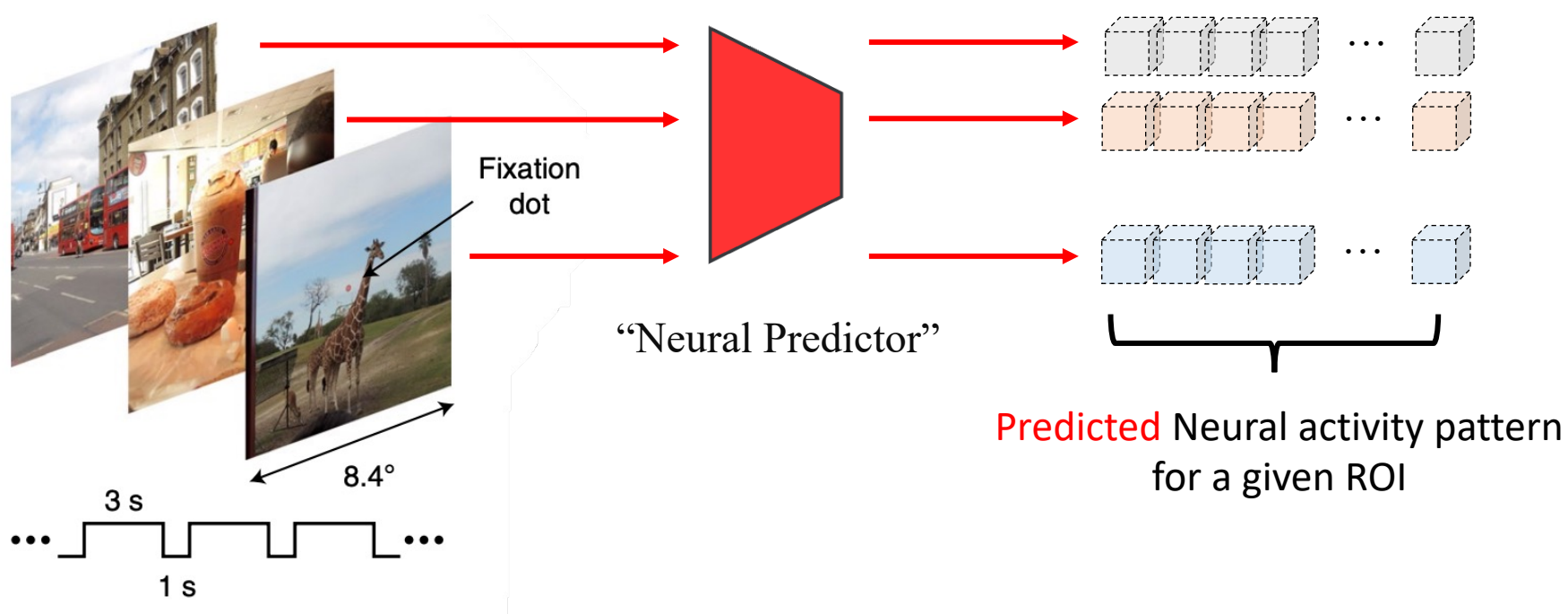Ventral Visual Stream Hierarchy Order

Subject 1 ROIs

# Method
## Neural data

- Each human subject viewed ~30,000 images (~9,000 unique).

- Brain activities were recorded with 7T fMRI.

- 5 bilateral ROIs were used:



Fixation dot

3 s

1 s

8.4°

Neural activity pattern in a given Region of Interest (ROI)

V1  V2  V4  LO  TO

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

# Method
## Neural data

- Each human subject viewed ~30,000 images (~9,000 unique).
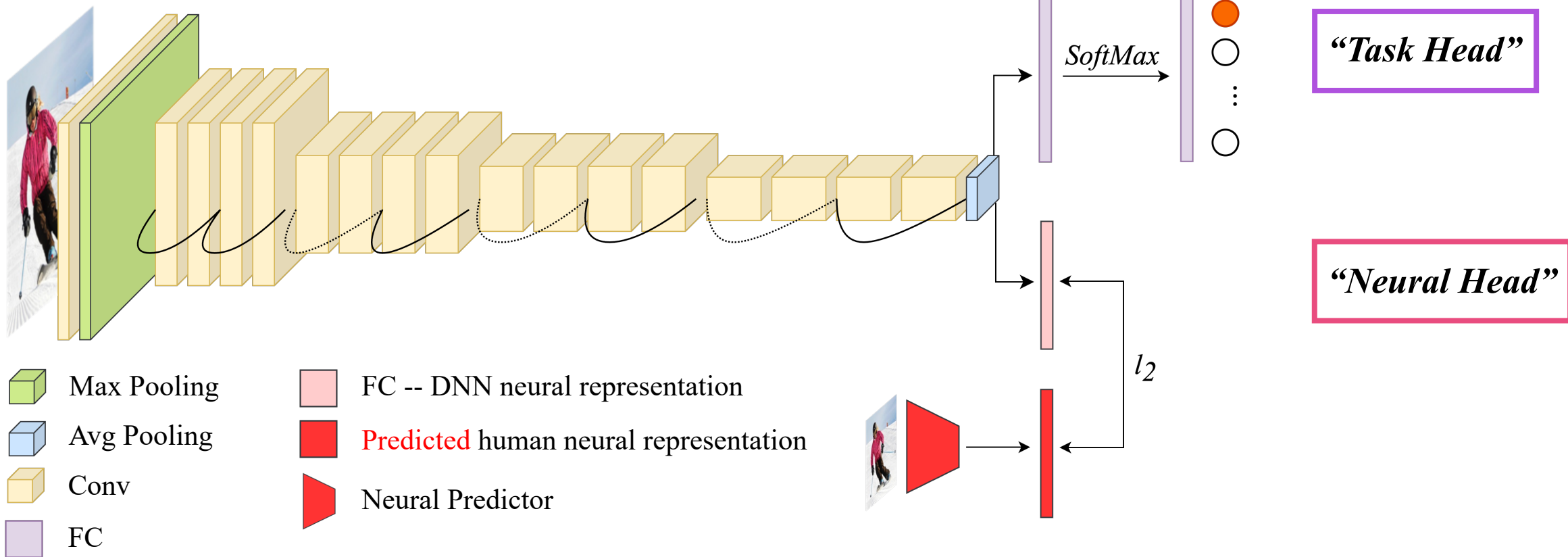- Brain activities were recorded with 7T fMRI.
- 5 bilateral ROIs were used.



"Neural Predictor"

Fixation dot

8.4°

3 s

1 s

Predicted Neural activity pattern for a given ROI

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

# Method
## Neurally-guided training

- DNN visual task training with Neural Guidance:
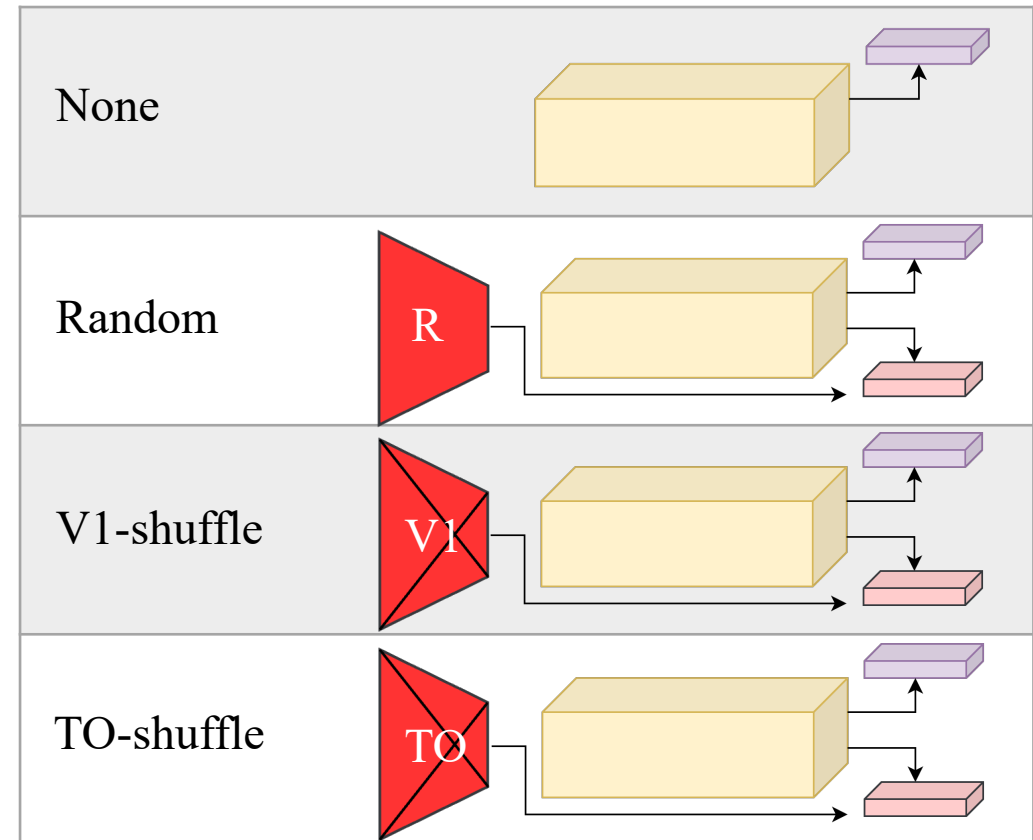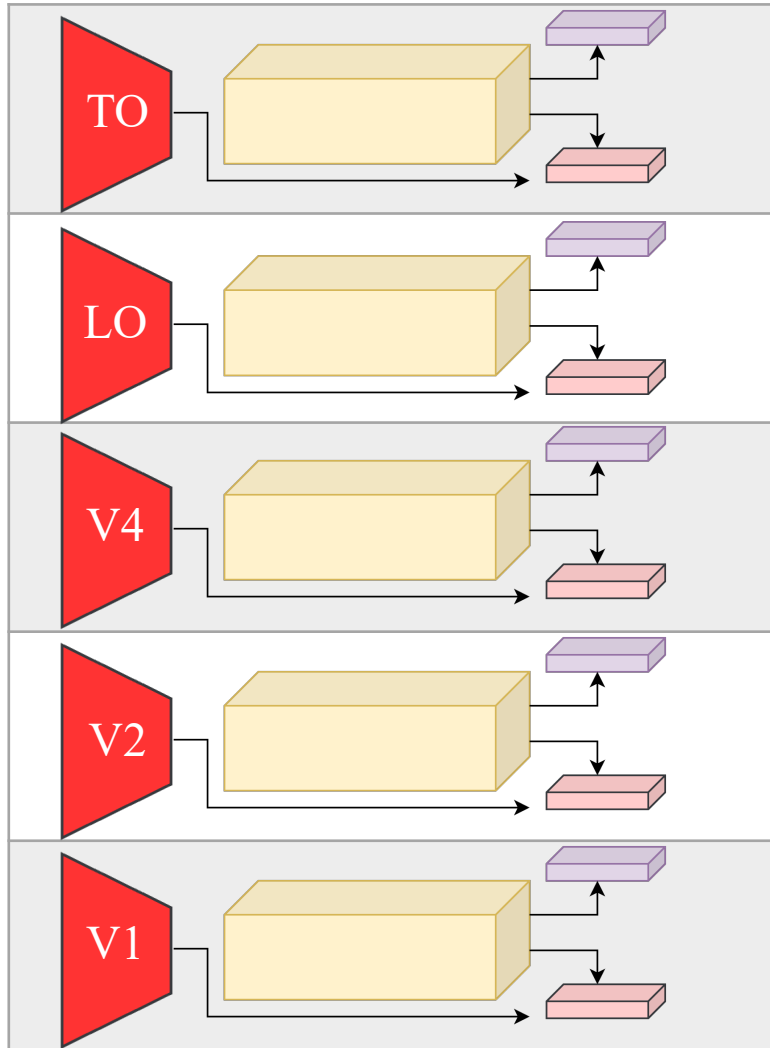
$$L = \alpha L_{task} + (1 - \alpha)||R_{DNN} - R_{human}||_2$$

# Method
## Summary of models

- 5 models with neural-guided training



- 4 baseline models for comparison

# Robustness of Neurally-guided Models
## Evaluation

- **$l_p$-based adversarial attack:**

$$\max_{||\tau||_p < \epsilon} l(f_\theta(x + \tau), y)$$
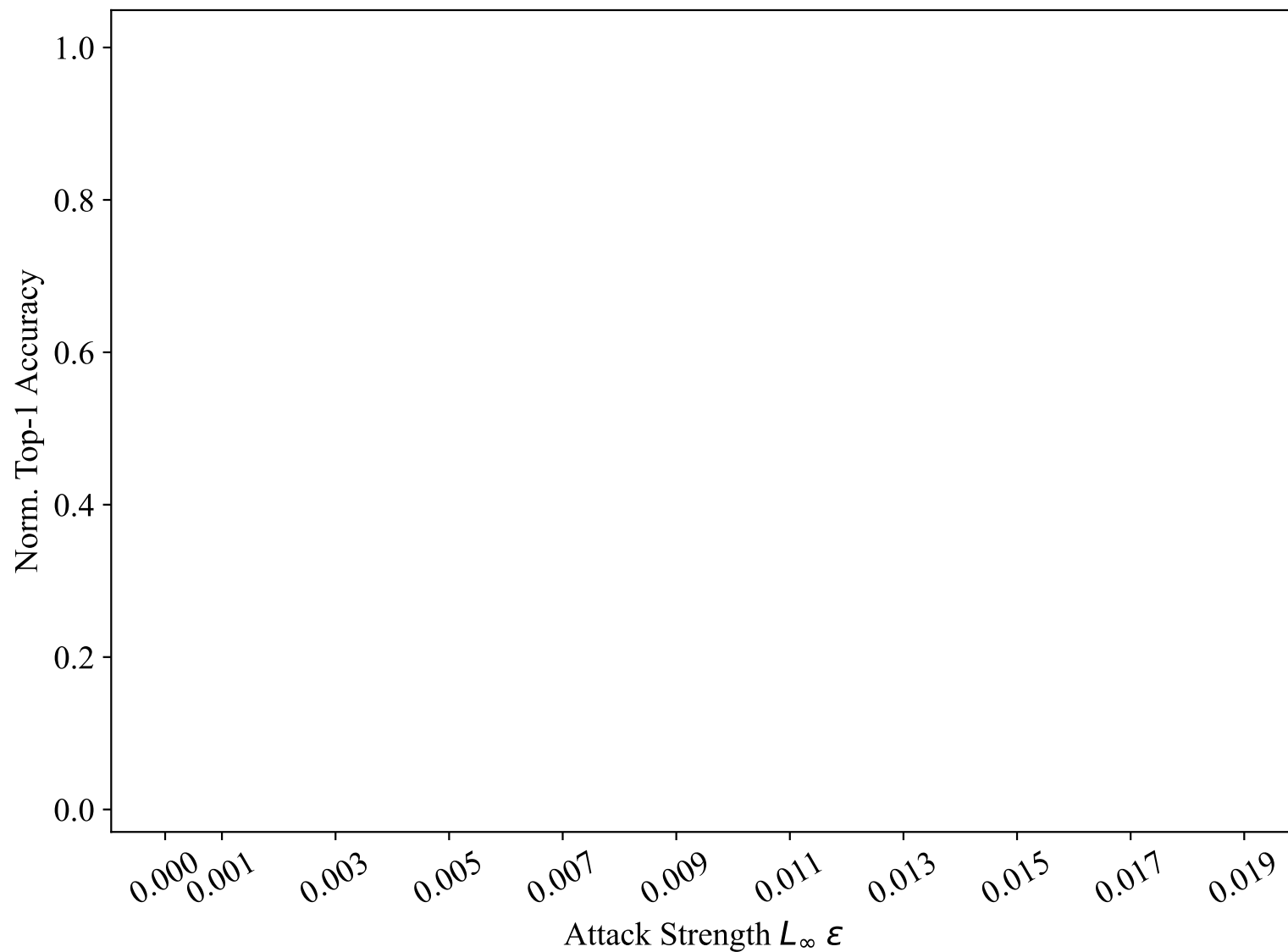
dog



+



=

Ostrich



(Szegedy et al., 2014)

We check:
1. Does neural guidance improve neural network robustness?  ✅/❌
2. Is such improvement hierarchical?  ✅/❌

# Results
## $l_\infty$-based PGD adversarial attack

"None"

# Results

## $l_\infty$-based PGD adversarial attack

"None"

# Results
## $l_\infty$-based PGD adversarial attack

"None"  "Random"

"V1-shuffle"  "TO-shuffle"

# Results
## $l_\infty$-based PGD adversarial attack

- **Neural guidance improves robustness** (max: 22% accuracy increase) ✅

- **There exists a hierarchy of improvement's magnitude** ✅

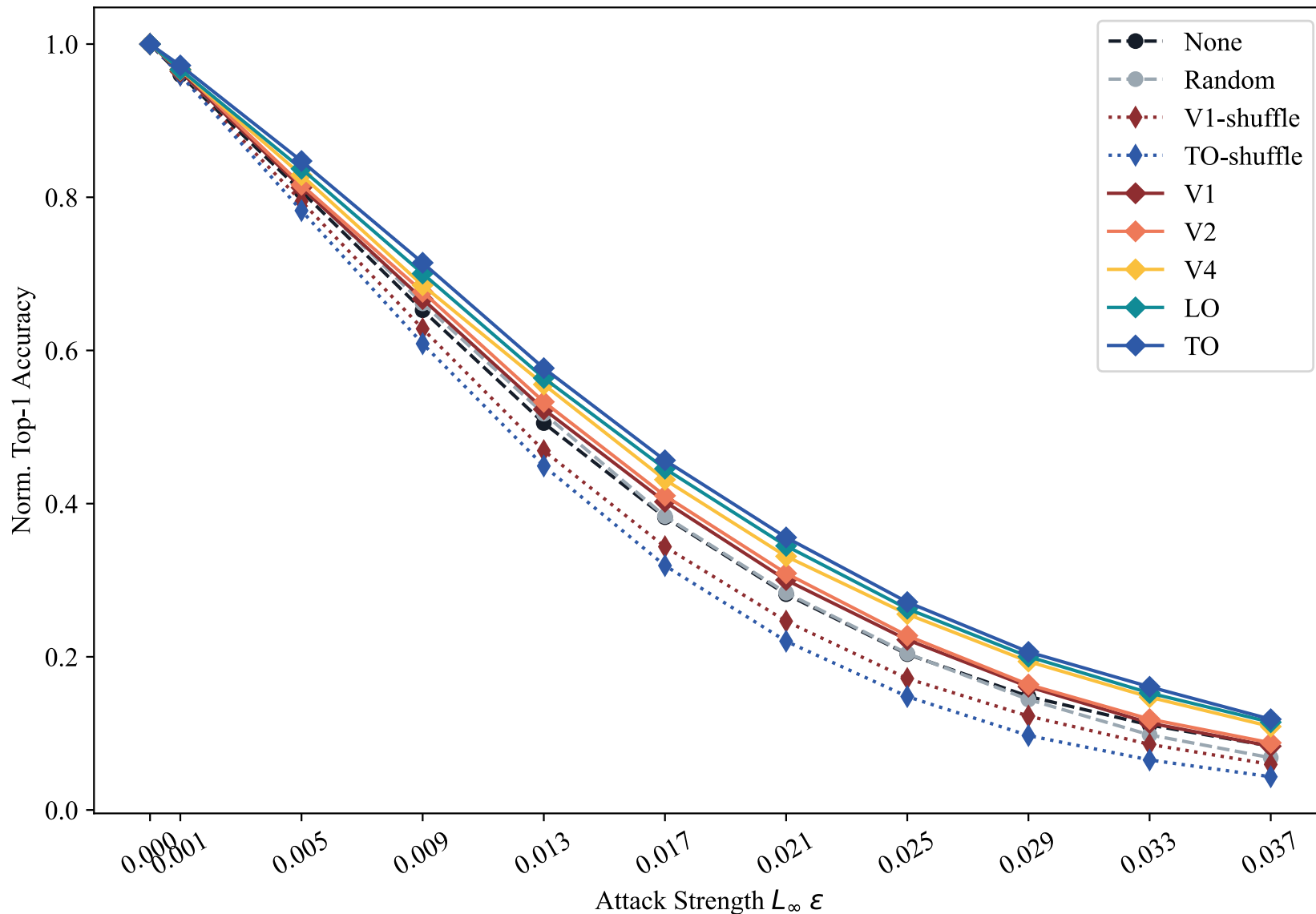- **Similar results have been observed with:**
  - $L_\infty$ FGSM
  - Auto-Attack (APGD-CE, APGD-T, FAB square)
  - $L_2$ FGM
  - $L_2$ Deepfool

# Results

## $l_\infty$-based PGD adversarial attack

- **Neural guidance improves robustness** (max: 12% accuracy increase) ✓

- **There exists a hierarchy of improvement's magnitude** ✓

# Results
## $l_\infty$-based PGD adversarial attack

# Results
## $l_\infty$-based PGD adversarial attack

- Do representations from neurally-guided DNNs benefit other visual tasks beyond basic classification?



Encoder  Decoder

*"Happy dog sitting in the bed of a pickup truck."*

*("Show, Attend, &Tell", Xu et al., 2015)*

# Results
## $l_\infty$-based PGD adversarial attack

- Do representations from neurally-guided DNNs benefit other visual tasks beyond basic classification?



Fully-trained
neurally-guided DNN backbone

Attention    LSTM

*Happy*
*Dog*
*Sitting*
....

BLEU score

*"Happy dog sitting in the bed of a pickup truck."*

*("Show, Attend, &Tell", Xu et al., 2015)*

# Results
## $l_\infty$-based PGD adversarial attack

- **Neural regularization improves robustness** (max: 0.03 BLEU-1 increase) ✅

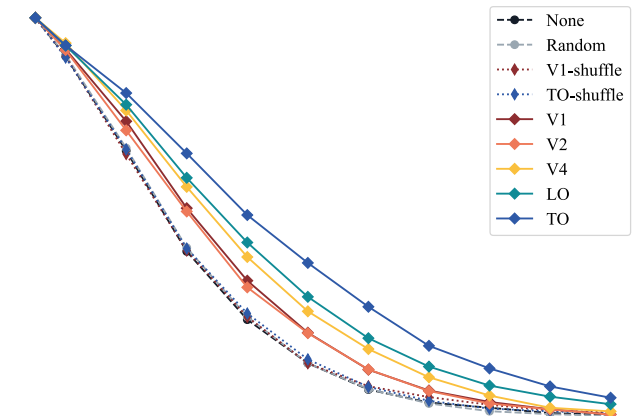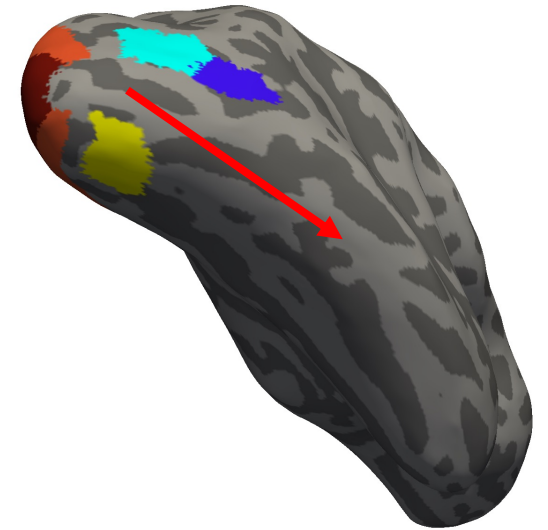- **There exists a hierarchy of improvement's magnitude** ✅

Neural-guidance
→ Robust feature extractor

# Conclusion & Discussion



- We found hierarchical improvements in DNN robustness across:
  - Datasets (ImageNet, CIFAR-100, MSCOCO)
  - Tasks (Classification, Captioning)
  - Attacks ($L_\infty$ PGD, $L_\infty$ FGSM, Autoattack, $L_2$ FGM, $L_2$ Deepfool)
- Implications:
  - Evolving representation space along ventral visual stream
  - Learnable and improvable with generic DNN structures
  - Potential for uncovering principles of building human-like representation space and advancing DNN architectural development
- Further analysis
  - Neurally-guided models are more shape-biased
  - Smoother output surface achieved in a different way from conventional solutions.
  - Neurally-guided models experience profound changes in their representation space

# Acknowledgement



*Linjian Ma*          *Bo Li*          *Diane M. Beck*