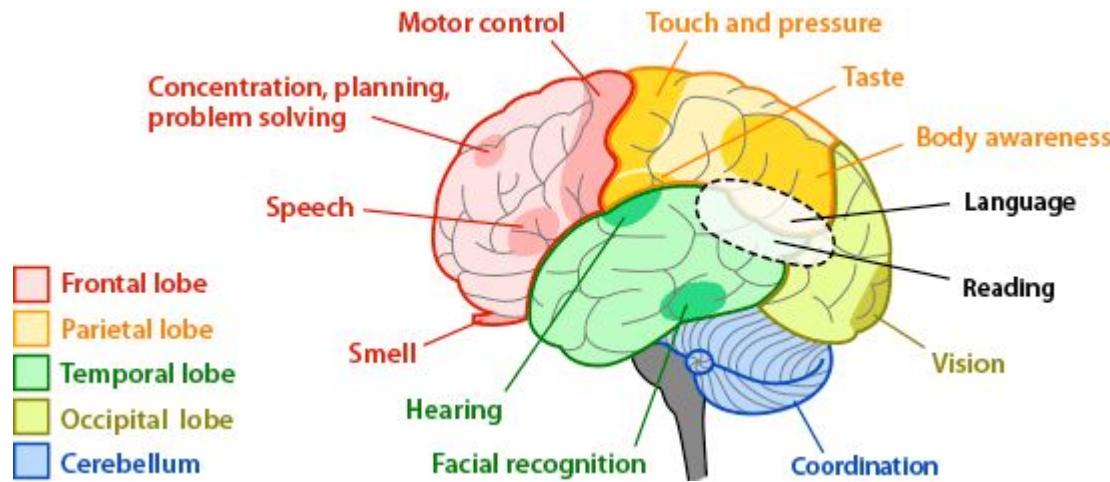
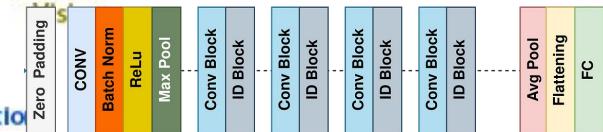
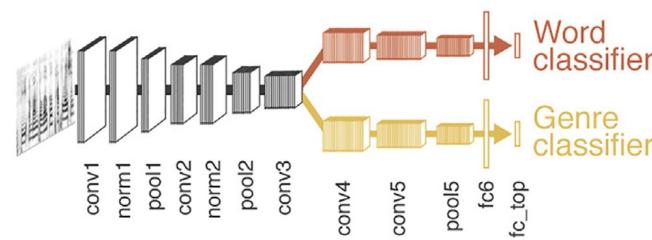
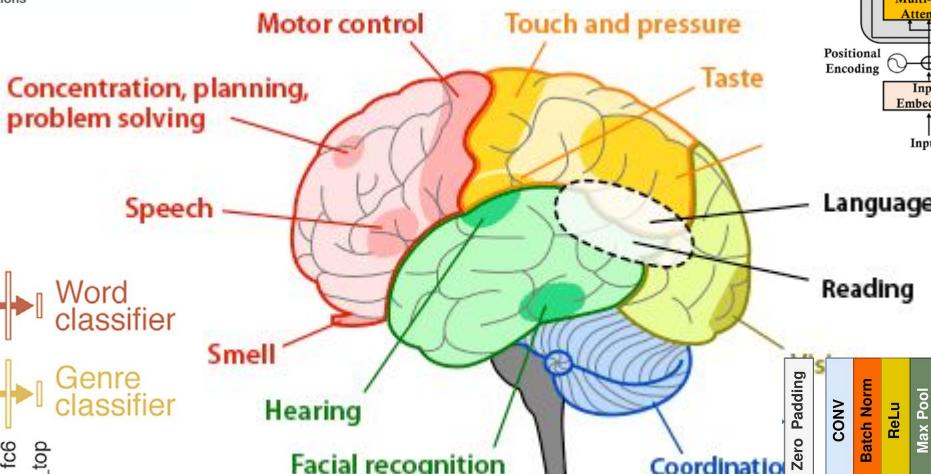
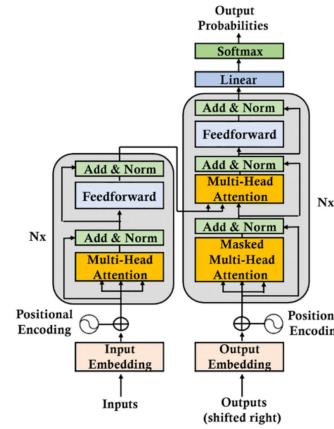
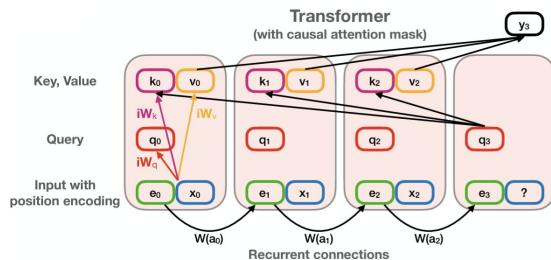


# **Bridging Theories: What Task-Optimized Models Reveal About Brain Computations**

Tahereh Toosi

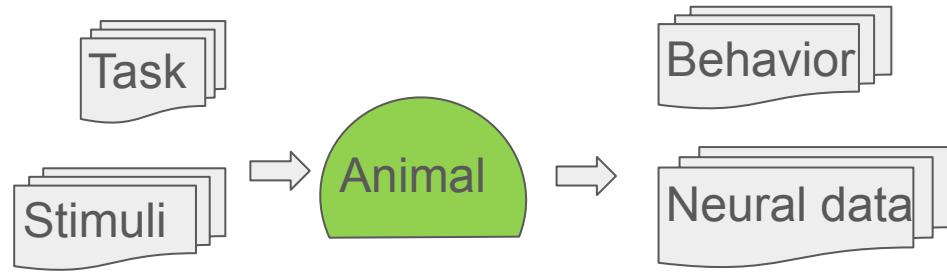
Advanced Topics in Theoretical Neuroscience  
March 2025



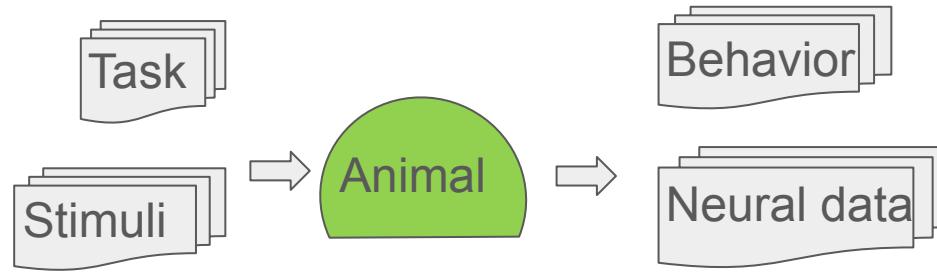


There is more to the task-optimized modeling framework than to replace every brain region with a network!

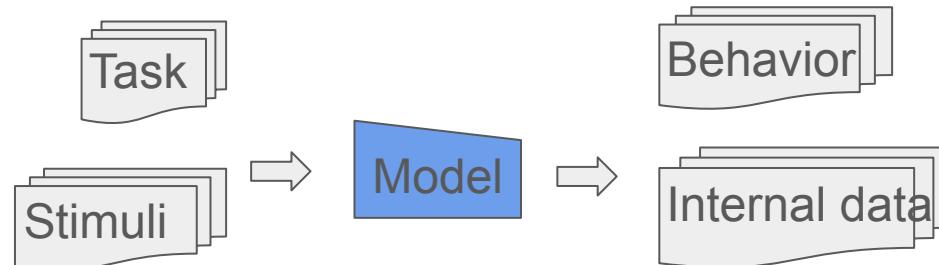
# Systems Neuroscience



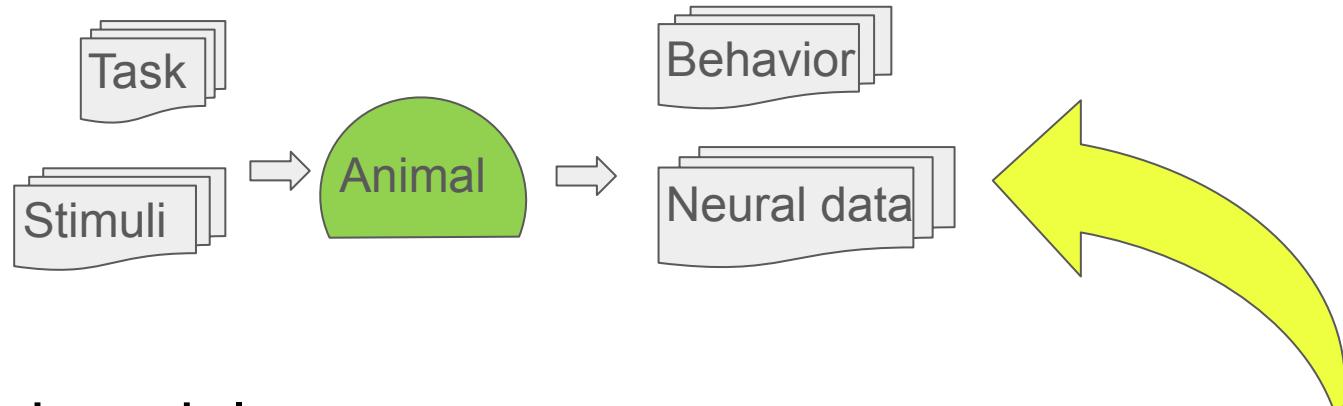
# Systems Neuroscience



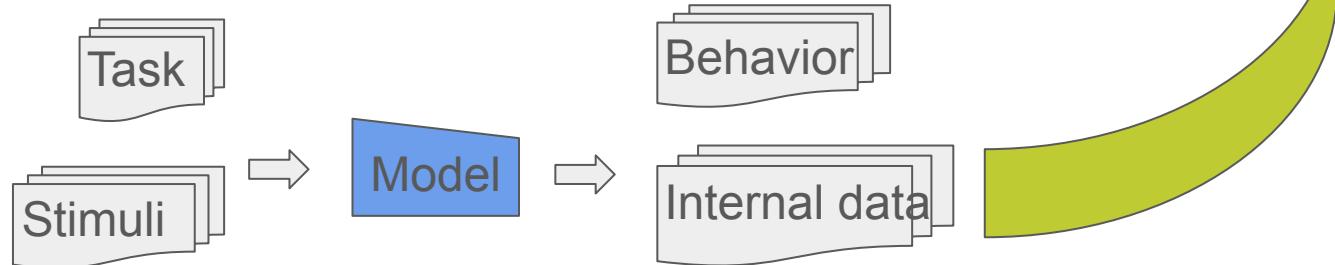
## Task-optimized models



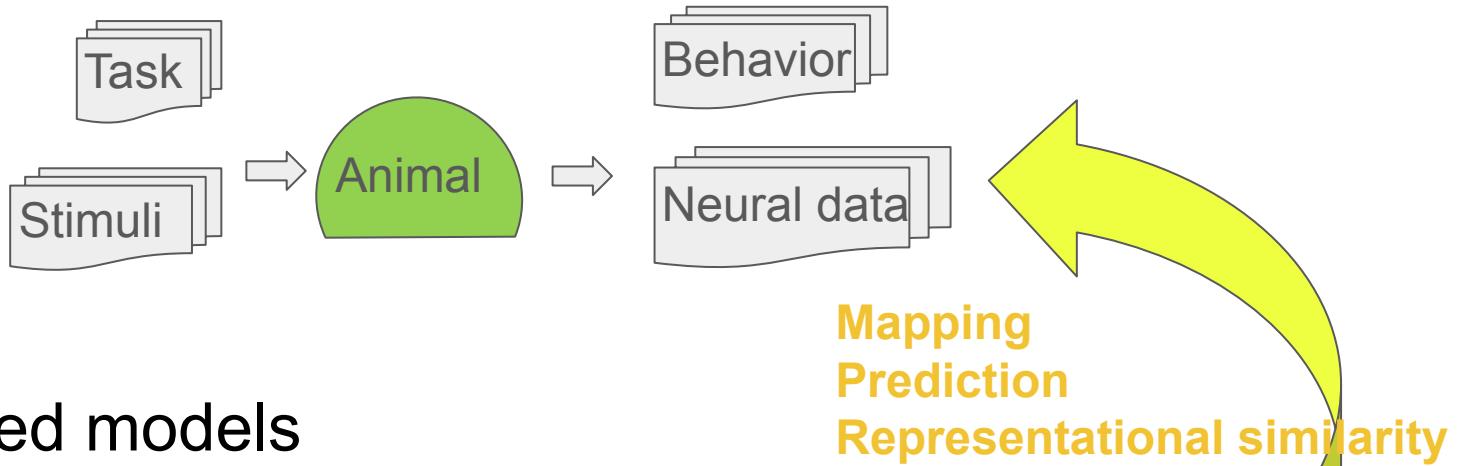
# Systems Neuroscience



## Task-optimized models



# Systems Neuroscience



## Task-optimized models



# Roadmap

- Comparing models to data
  - Compare model 1 to model 2: which one maps to data better? (**Case study 1**)
  - However, the platonic representation hypothesis!
- Closing the loop by synthesizing stimuli
  - Once a model predicts neural data well, it can be inverted, to generate new stimuli (**Case study 2**)
  - Still a very new field
- Computational tricks
  - Contrastive learning
  - Regularizations (**Case study 3**)
- Bridging theories (different perspectives, same math!)
  - Recognition | Generation
  - Learning | Attention

# Roadmap

- Comparing models to data
  - Compare model 1 to model 2: which one maps to data better? (**Case study 1**)
  - However, the platonic representation hypothesis!
- Closing the loop by synthesizing stimuli
  - Once a model predicts well, it can be inverted, to give stimuli (**Case study 2**)
  - However, Mechanistic interpretability hasn't really been useful in deep learning
- Computational tricks
  - Contrastive learning
  - Regularizations (**Case study 3**)
- Bridging theories (different perspectives, same math!)
  - Attention mechanism to Memory
  - Regularization to Generativity (**Case study 4**)

# Performance-optimized hierarchical models predict neural responses in higher visual cortex

Daniel L. K. Yamins<sup>a,1</sup>, Ha Hong<sup>a,b,1</sup>, Charles F. Cadieu<sup>a</sup>, Ethan A. Solomon<sup>a</sup>, Darren Seibert<sup>a</sup>, and James J. DiCarlo<sup>a,2</sup>

<sup>a</sup>Department of Brain and Cognitive Sciences and McGovern Institute for Brain Research, Massachusetts Institute of Technology, Cambridge, MA 02139; and <sup>b</sup>Harvard-MIT Division of Health Sciences and Technology, Institute for Medical Engineering and Science, Massachusetts Institute of Technology, Cambridge, MA 02139

Edited by Terrence J. Sejnowski, Salk Institute for Biological Studies, La Jolla, CA, and approved April 8, 2014 (received for review March 3, 2014)

The ventral visual stream underlies key human visual object recognition abilities. However, neural encoding in the higher areas of the ventral stream remains poorly understood. Here, we describe

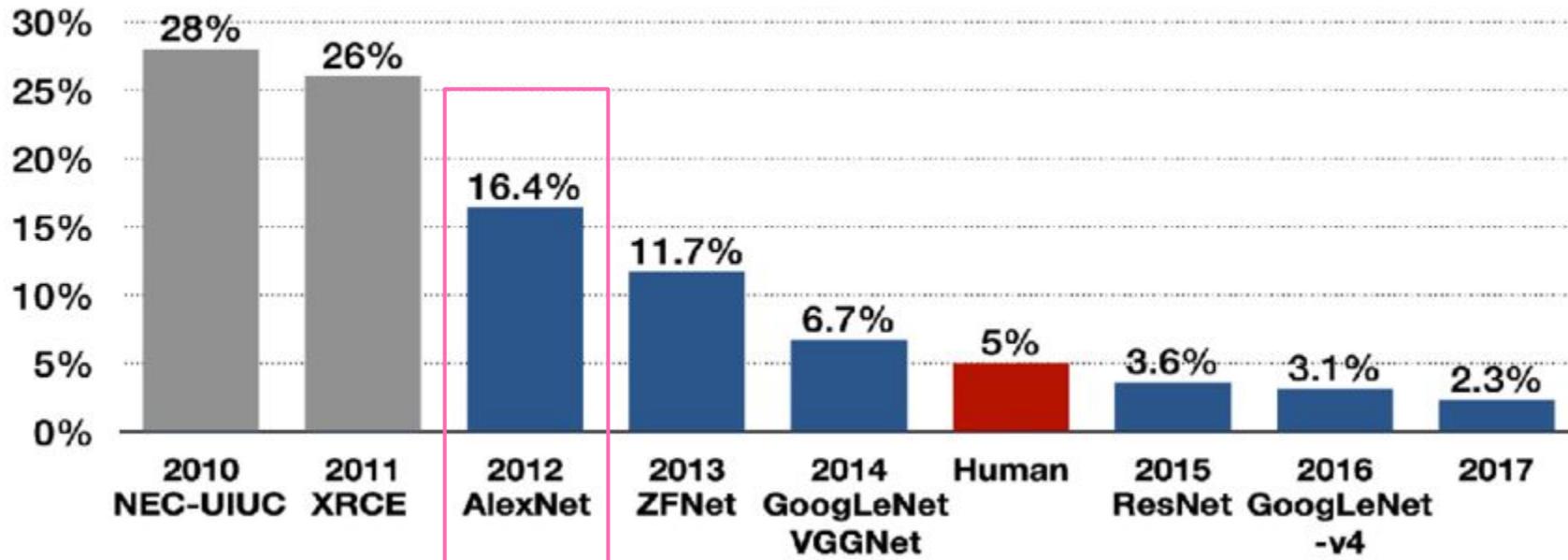
Explaining the neural encoding in these higher ventral areas thus remains a fundamental open question in systems neuroscience.

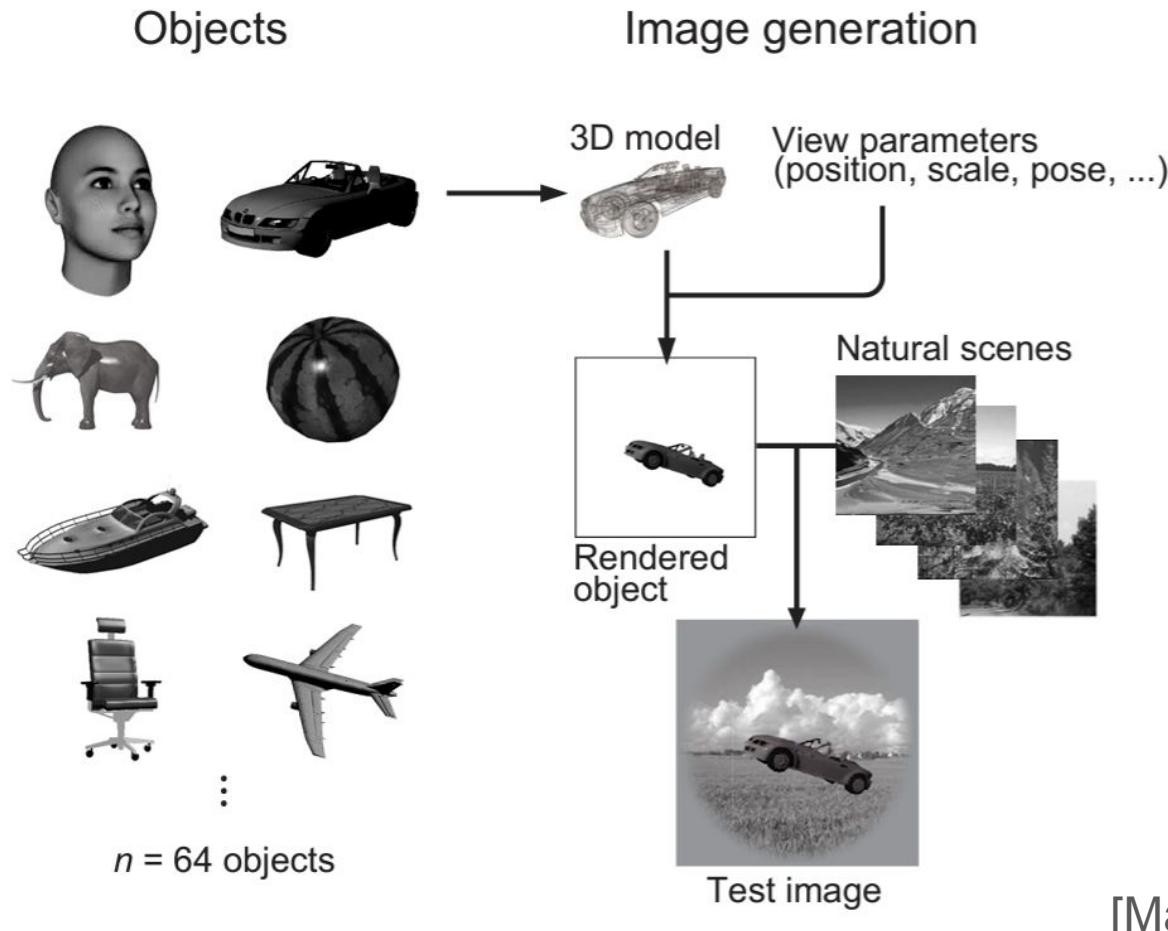
As with V1, models of higher ventral areas should be neurally plausible. However, unlike the higher ventral stream visual

## 2012 Imagenet Challenge

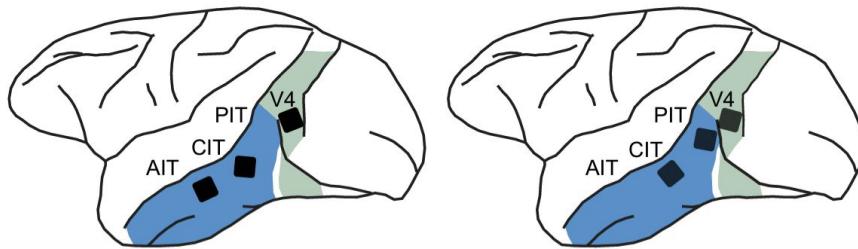
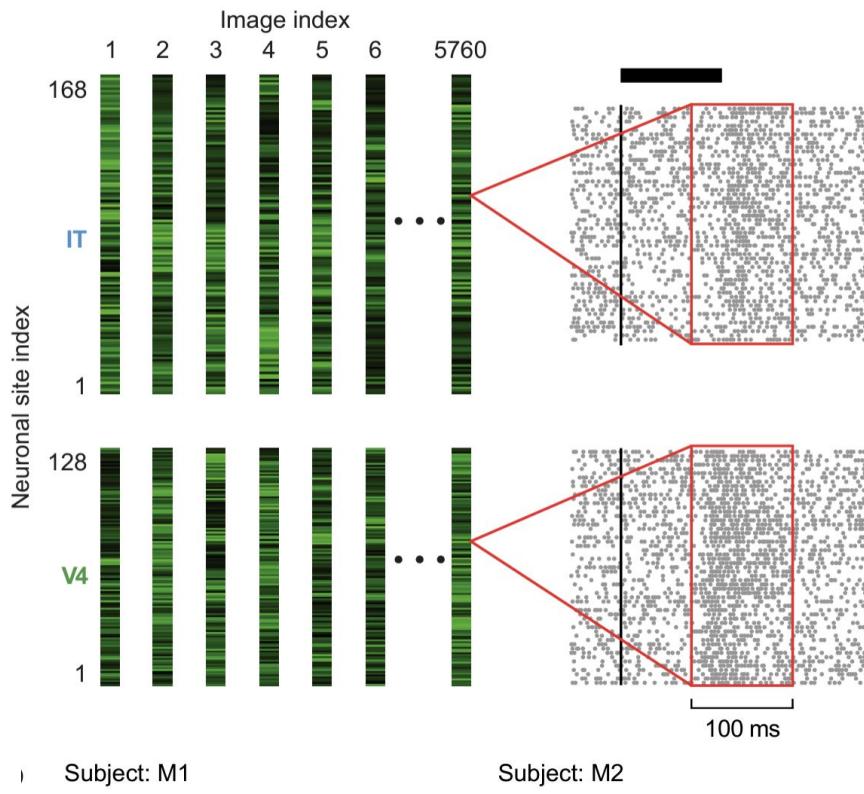
## 2014 Neuroscience papers

### Top-5 error

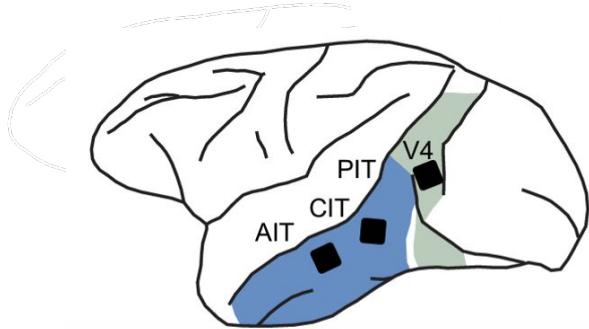
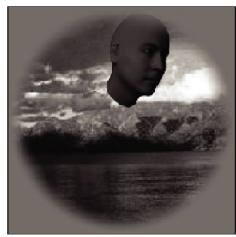
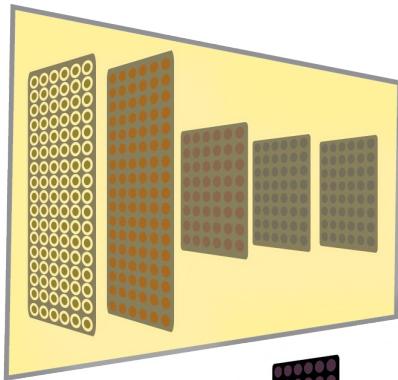
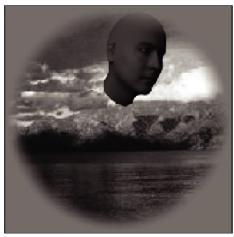


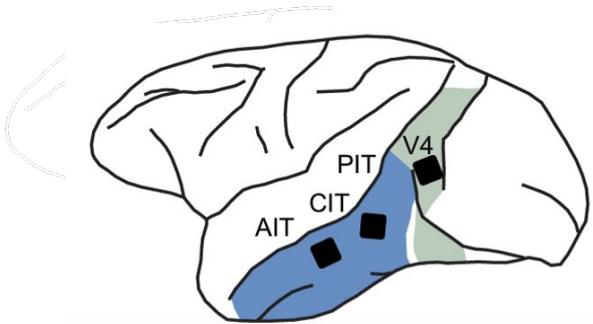
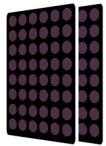
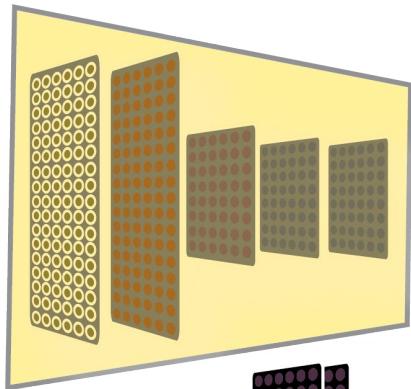


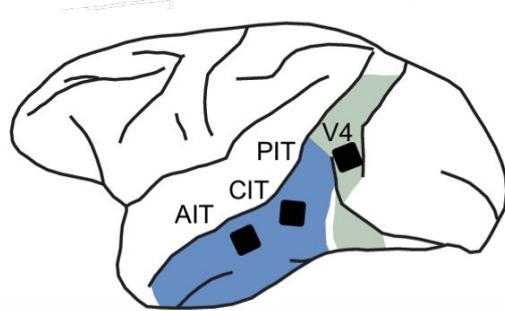
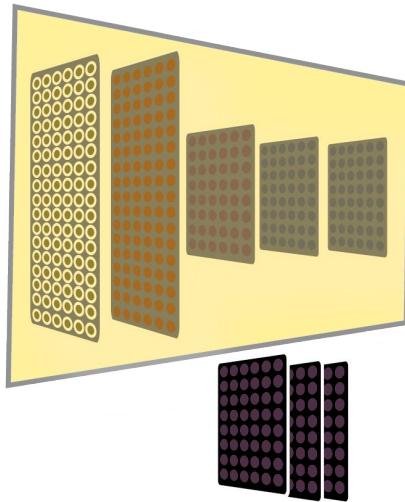
[Majaj et al 2015]

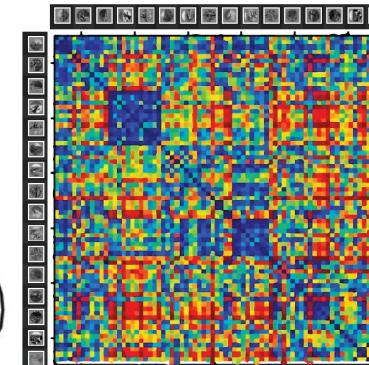
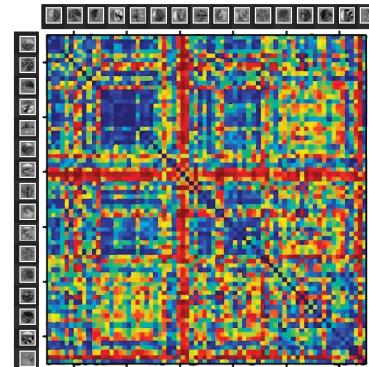
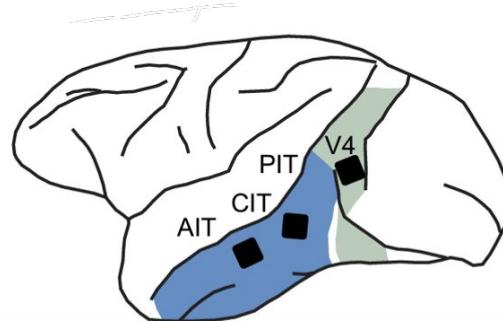
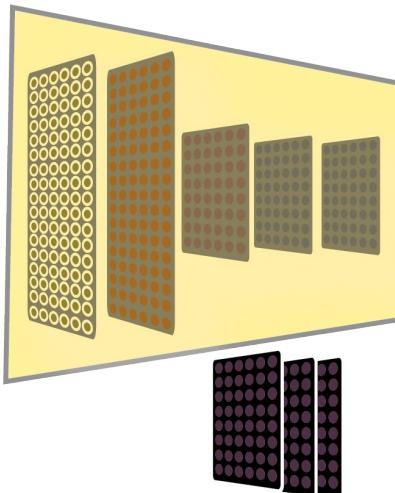


[Majaj et al 2015]





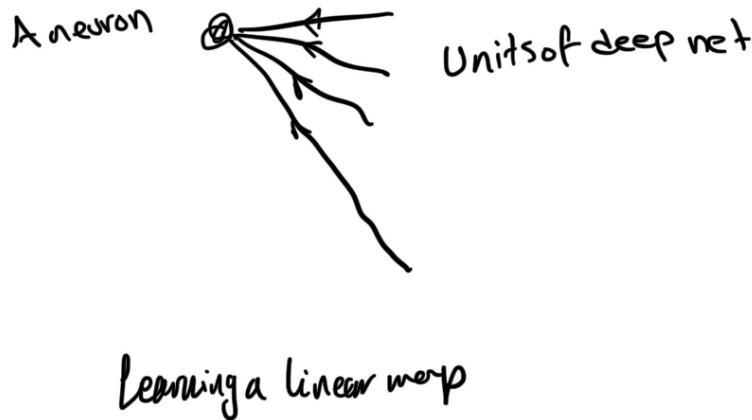




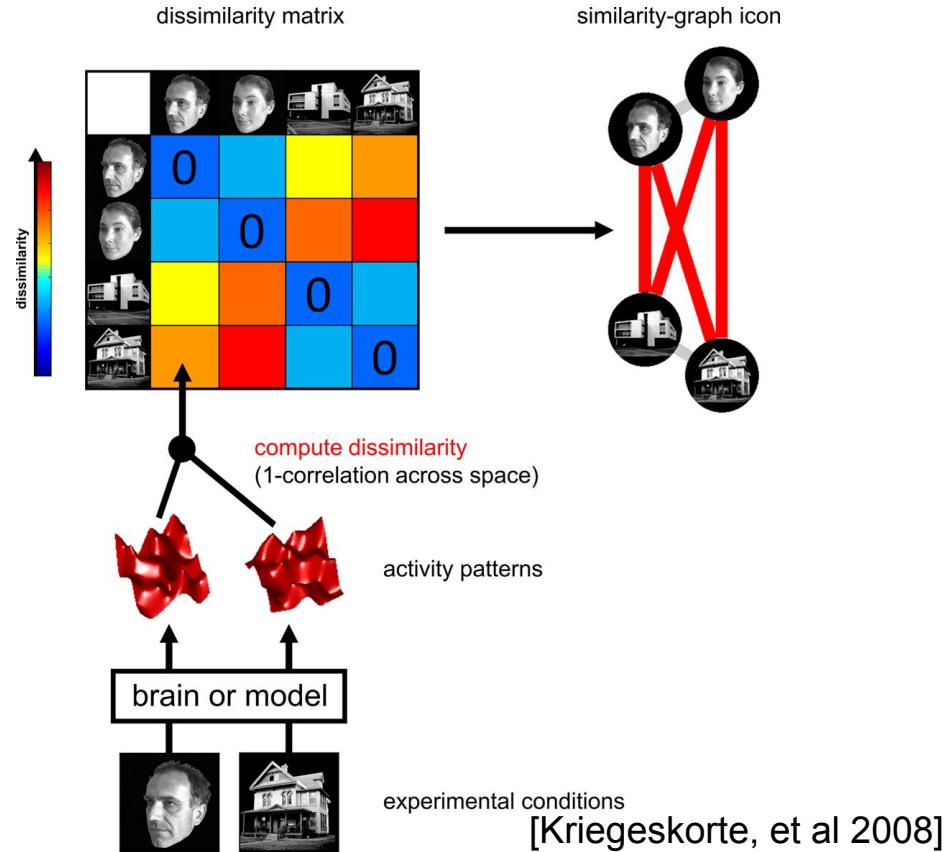
Representational  
Similarity

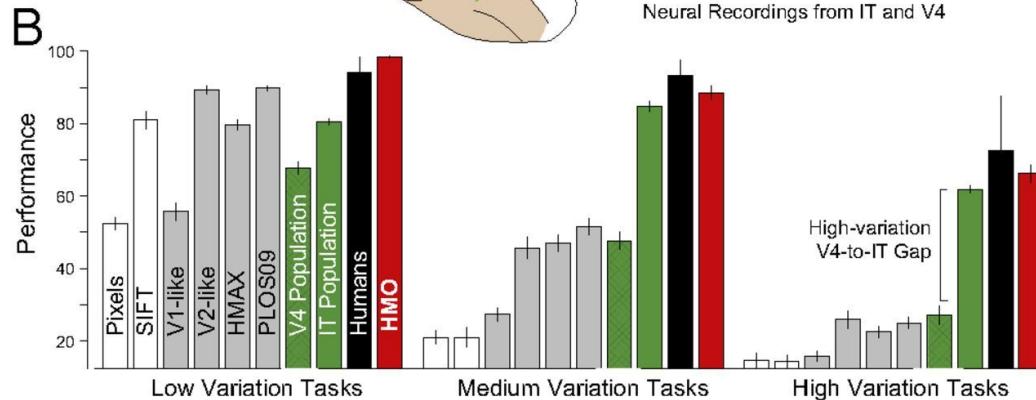
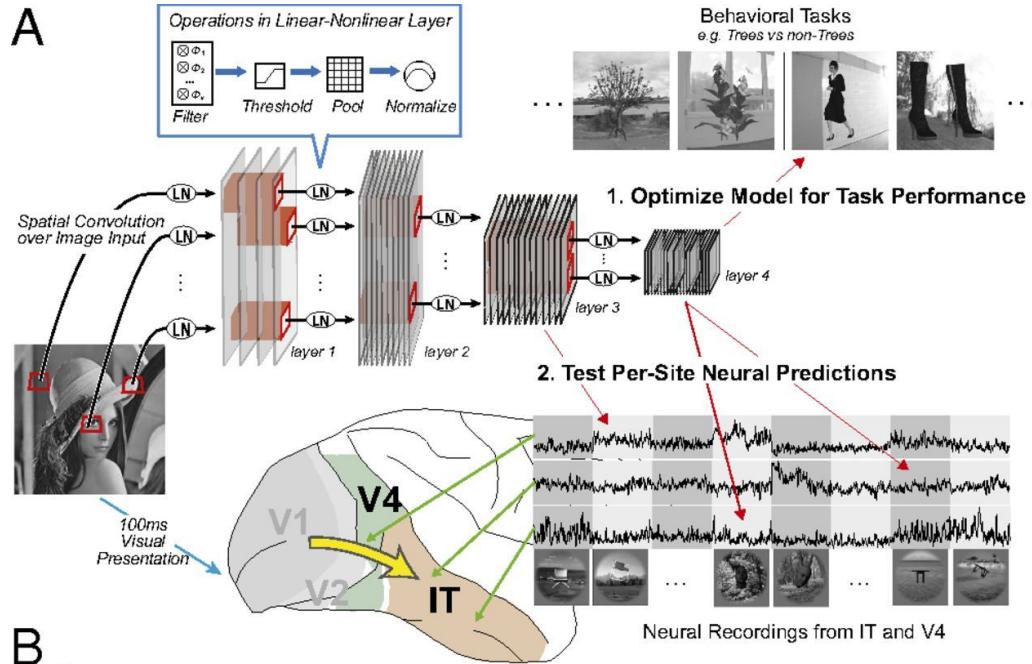
# Evaluation of similarity between DNNs and Brain

Linear mapping



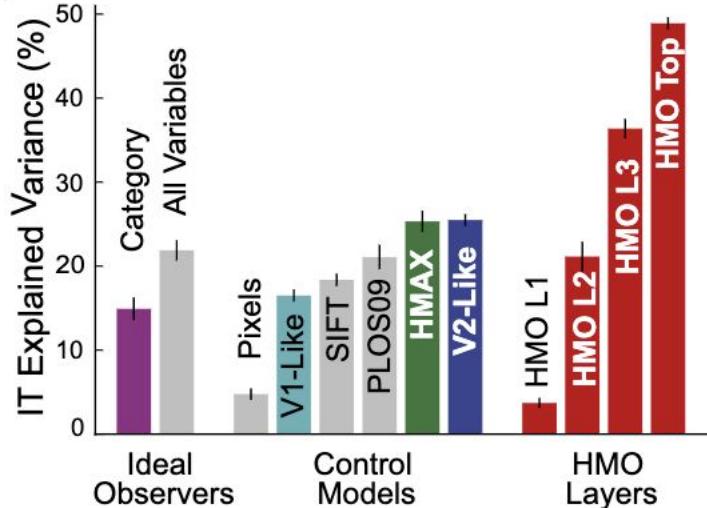
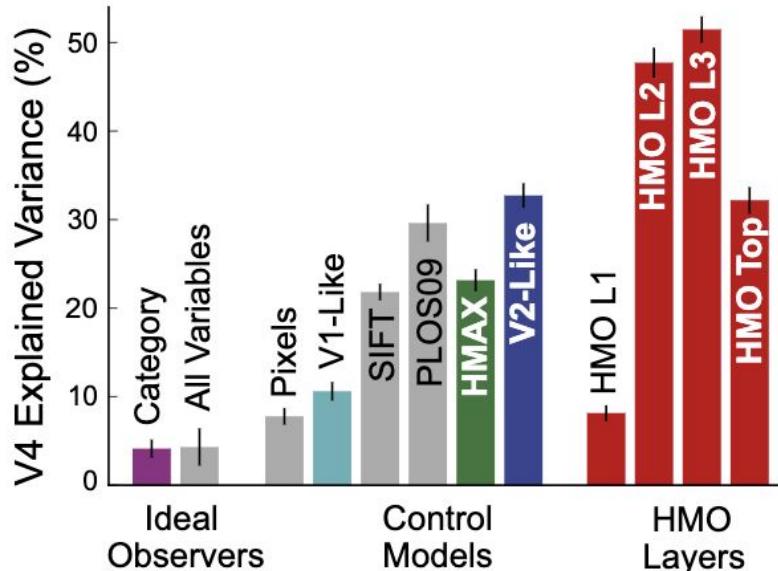
Representational (dis)similarity



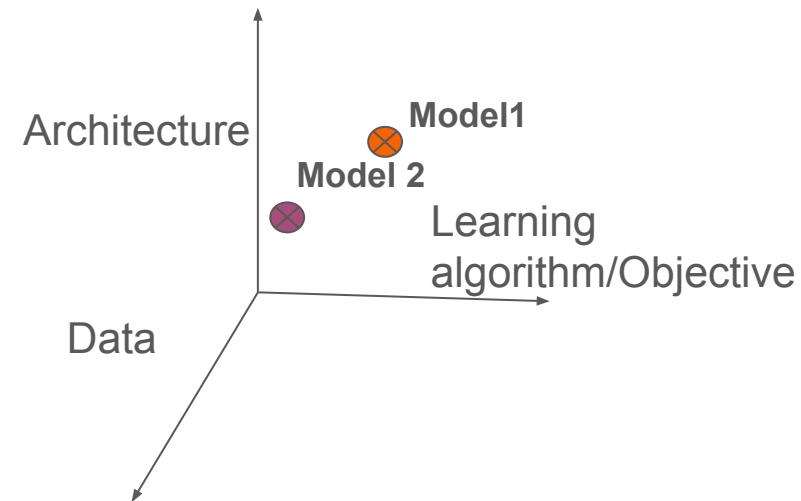
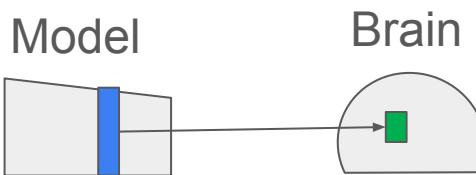


[Yamins, Hong, et al 2014]

# Mapping model units to neurons



# Model selection

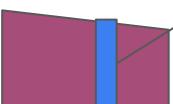


---

Model1



score 1 > score 2?



Model2



► Proc Natl Acad Sci U S A. 2021 Jan 11;118(3):e2014196118. doi: [10.1073/pnas.2014196118](https://doi.org/10.1073/pnas.2014196118)

## Unsupervised neural network models of the ventral visual stream

Chengxu Zhuang<sup>a,1</sup>, Siming Yan<sup>b</sup>, Aran Nayebi<sup>c</sup>, Martin Schrimpf<sup>d</sup>, Michael C Frank<sup>a</sup>, James J DiCarlo<sup>d</sup>, Daniel L K Yamins<sup>a,e,f</sup>



► PLoS Comput Biol. 2022 Jan 7;18(1):e1009739. doi: [10.1371/journal.pcbi.1009739](https://doi.org/10.1371/journal.pcbi.1009739)

## Increasing neural network robustness improves match to macaque V1 eigenspectrum, spatial frequency preference and predictivity

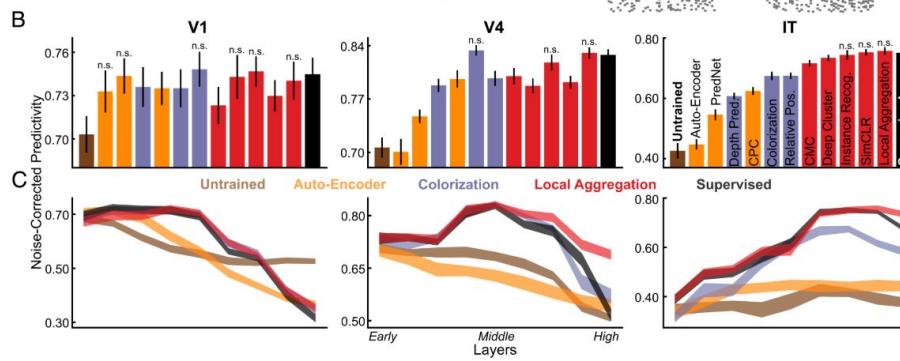
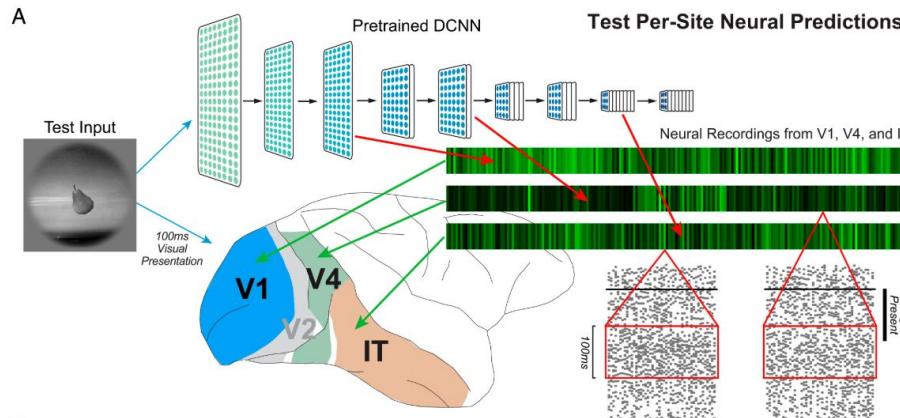
Nathan C L Kong<sup>1,3,\*</sup>, Eshed Margalit<sup>2,3</sup>, Justin L Gardner<sup>1,3</sup>, Anthony M Norcia<sup>1,3</sup>

## Neural Foundations of Mental Simulation: Future Prediction of Latent Representations on Dynamic Scenes

Aran Nayebi<sup>1,\*</sup>, Rishi Rajalingham<sup>1,4</sup>, Mehrdad Jazayeri<sup>1,2</sup>, and Guangyu Robert Yang<sup>1,2,3</sup>

being explicitly trained to do so. Finally, we find that not all foundation model latents are equal. Notably, models that future predict in the latent space of video foundation models that are optimized to support a diverse range of egocentric sensorimotor tasks, reasonably match both human behavioral error patterns and neural dynamics across all environmental scenarios that we were able to test. Overall,

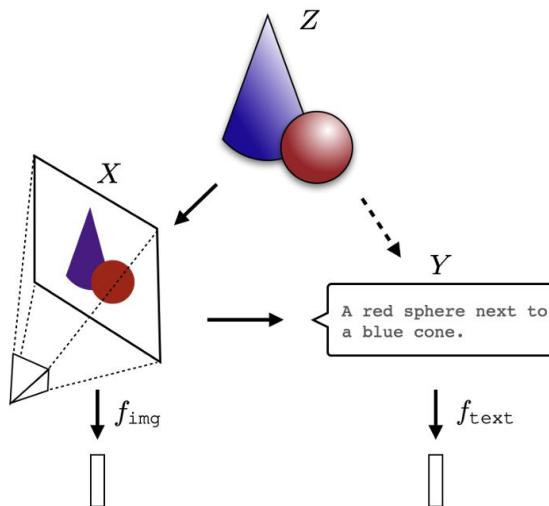
# self-supervised models compared to supervised models



# There is a catch!

## The Platonic Representation Hypothesis

Neural networks, trained with different objectives on different data and modalities, are converging to a shared statistical model of reality in their representation spaces.

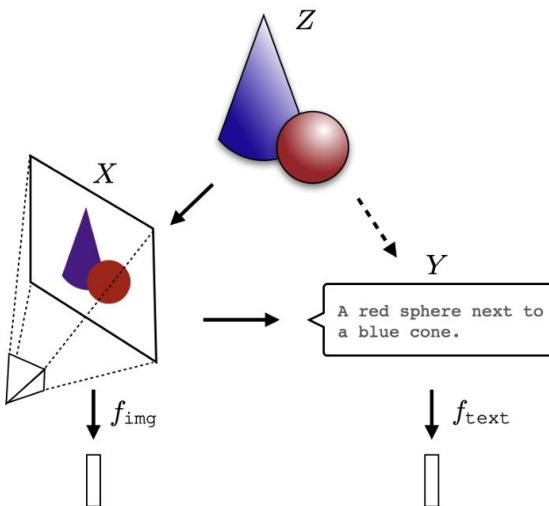


## The Platonic Representation Hypothesis

# There is a catch!

## The Platonic Representation Hypothesis

Neural networks, trained with different objectives on different data and modalities, are converging to a shared statistical model of reality in their representation spaces.



## The Multitask Scaling Hypothesis

There are fewer representations that are competent for  $N$  tasks than there are for  $M < N$  tasks. As we train more general models that solve more tasks at once, we should expect fewer possible solutions.

## The Capacity Hypothesis

Bigger models are more likely to converge to a shared representation than smaller models.

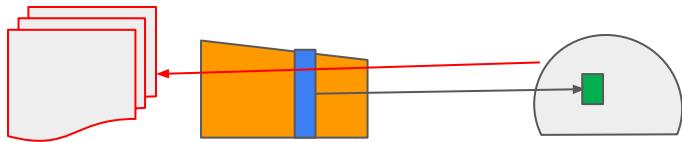
## The Simplicity Bias Hypothesis

Deep networks are biased toward finding simple fits to the data, and the bigger the model, the stronger the bias. Therefore, as models get bigger, we should expect convergence to a smaller solution space.

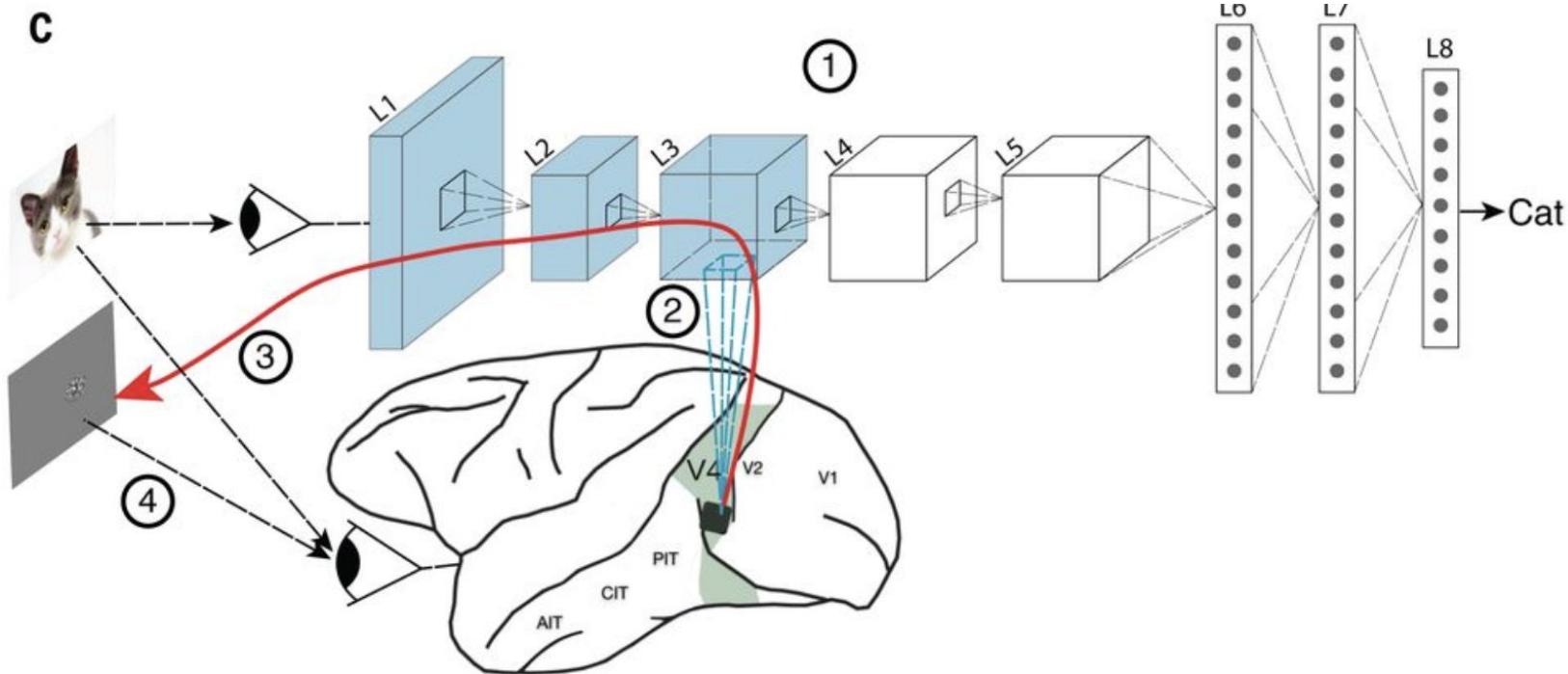
# Roadmap

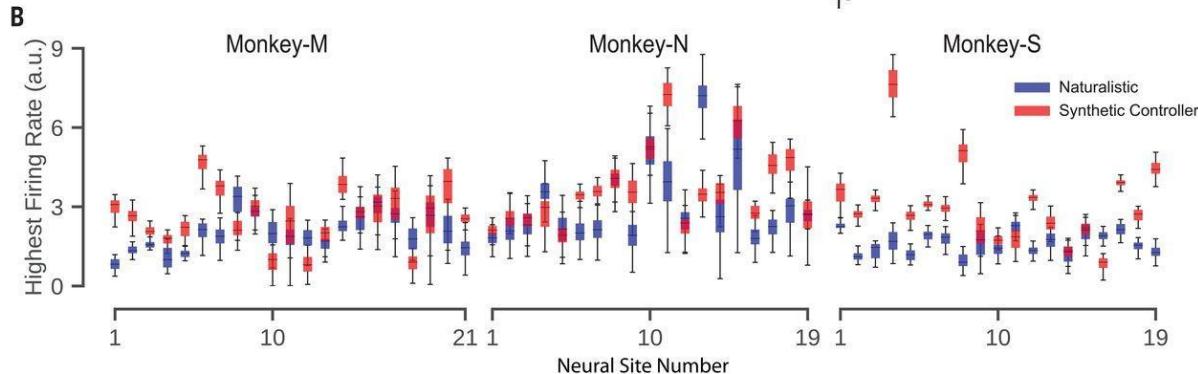
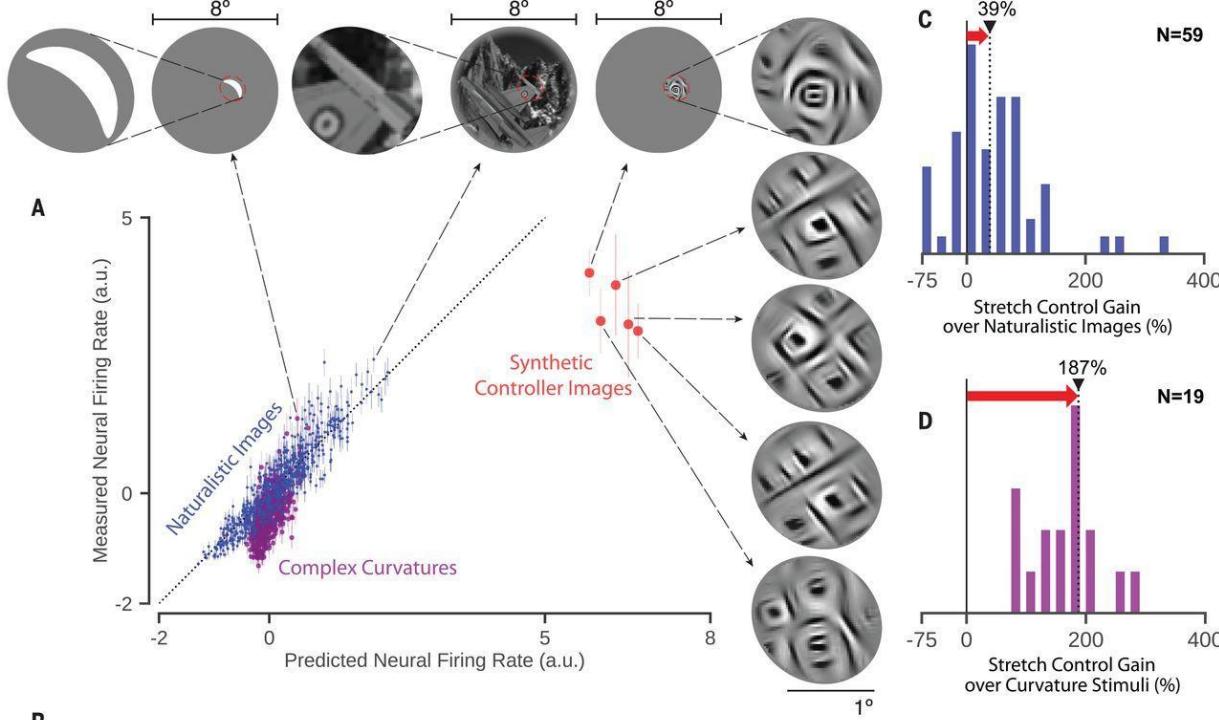
- Comparing models to data
  - Compare model 1 to model 2: which one maps to data better? (**Case study 1**)
  - However, the platonic representation hypothesis!
- Closing the loop by synthesizing stimuli
  - Once a model predicts neural data well, it can be inverted, to give new stimuli (**Case study 2**)
  - Still a very new field
- Computational tricks
  - Contrastive learning
  - Regularizations (**Case study 3**)
- Bridging theories (different perspectives, same math!)
  - Attention mechanism to Memory
  - Regularization to Generativity (**Case study 4**)

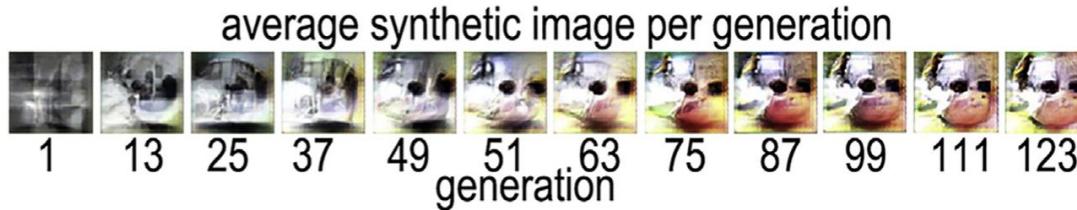
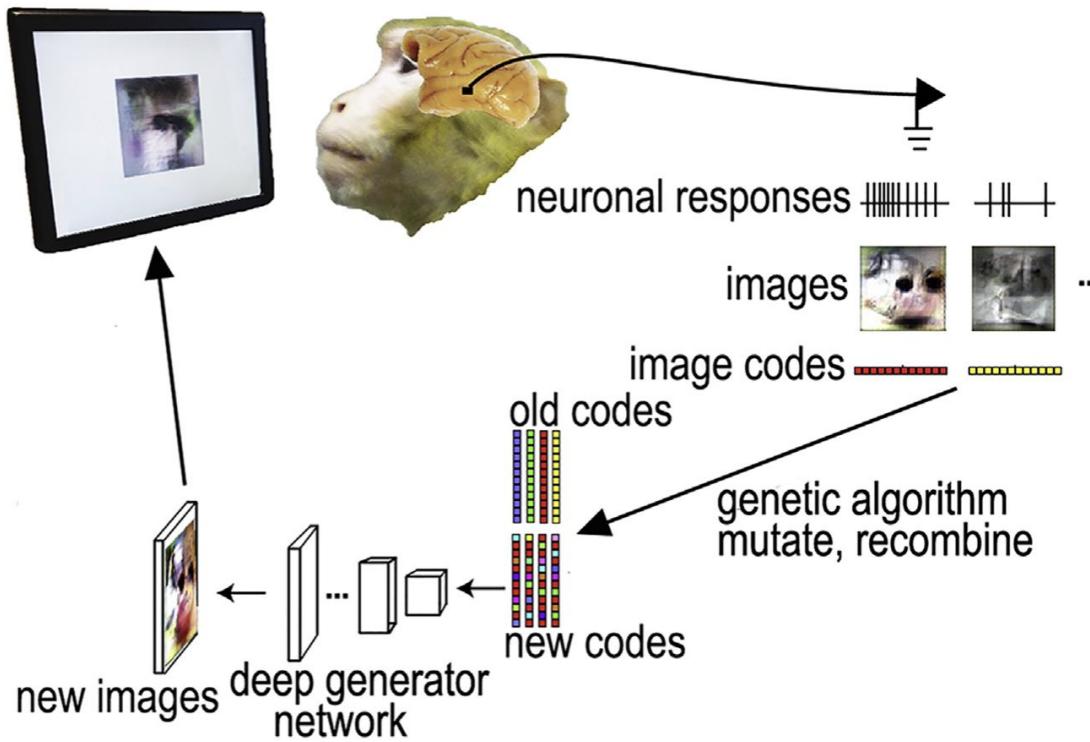
# Super stimuli or Maximally exciting inputs (MEI)



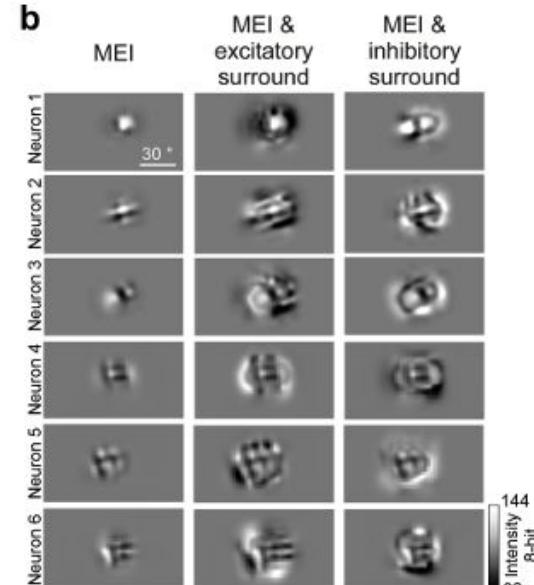
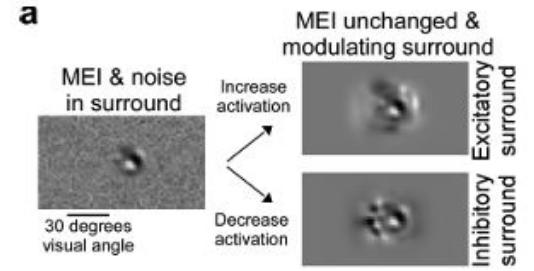
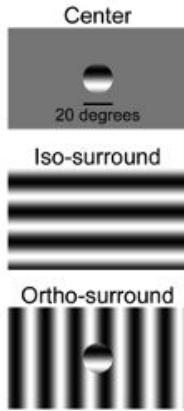
# Maximally Exciting inputs







# What can be learned from MEIs?

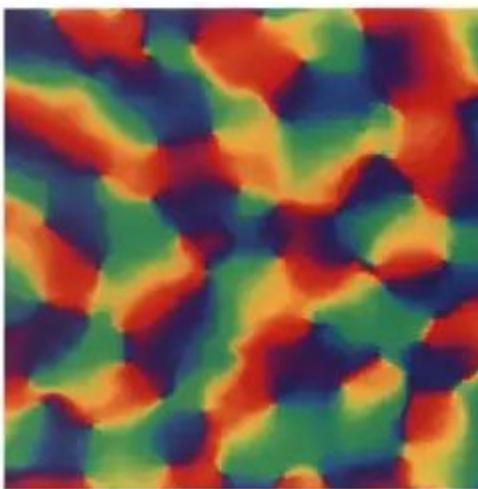


**Pattern completion and disruption characterize contextual modulation in the visual cortex**

Jiakun Fu<sup>1</sup>, Suhas Srinivasan<sup>2,\*</sup>, Luca Baroni<sup>8,\*</sup>, Zhuokun Ding<sup>1,9-11</sup>, Paul G. Fahey<sup>1,9-11</sup>, Paweł A. Pierzchlewicz<sup>2</sup>, Kayla Ponder<sup>1</sup>, Rachel Froebe<sup>1,9-11</sup>, Lydia Ntanavara<sup>1,9-11</sup>, Taliah Muhammad<sup>1</sup>, Konstantin F. Willeke<sup>2</sup>, Eric Wang<sup>1,9-11</sup>, Zhiwei Ding<sup>1,9-11</sup>, Dat T. Tran<sup>1,9-11</sup>, Stelios Papadopoulos<sup>1,9-11</sup>, Saumil Patel<sup>1,9-11</sup>, Jacob Reimer<sup>1</sup>, Alexander S. Ecker<sup>2,3</sup>, Xaq Pitkow<sup>1,6,7</sup>, Jan Antolik<sup>8</sup>, Fabian H. Sinz<sup>1,2,5</sup>, Ralf M. Häfner<sup>4</sup>, Andreas S. Tolias<sup>1,9-11,†,✉</sup>, and Katrin Franke<sup>1,9-12,†,✉</sup>

Topography matters?

# Beyond explained variance: cortical maps examples

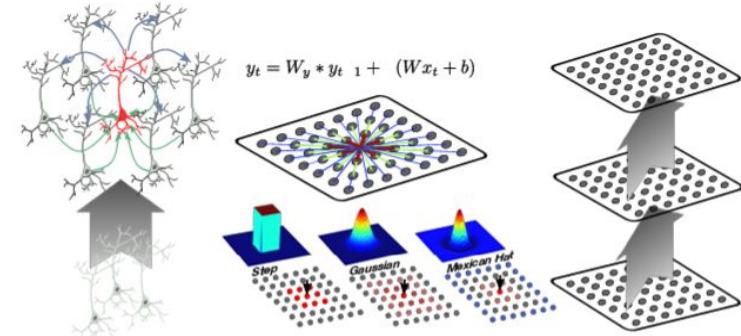
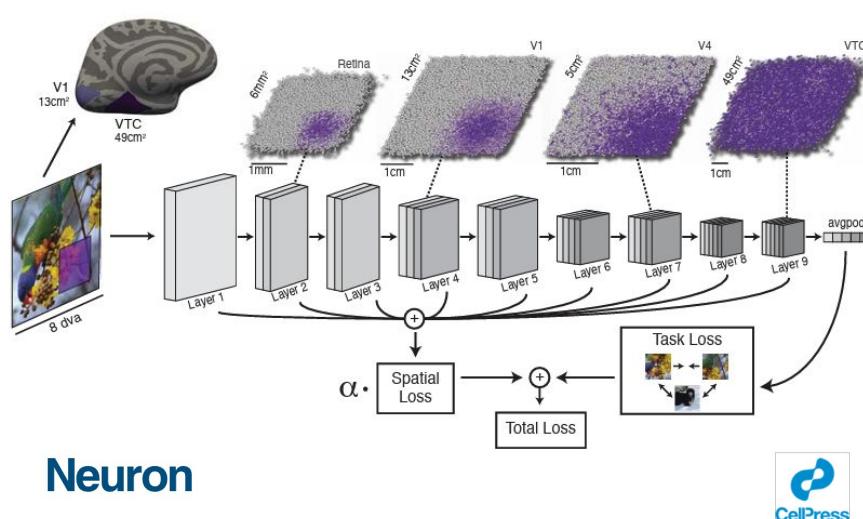


Low-level vision



High-level vision

# Two modeling approaches: direct optimization for the effect | use of bio-inspired inductive biases



Local lateral connectivity is sufficient for replicating cortex-like topographical organization in deep neural networks

Xinyu Qian<sup>1</sup>, Amir Ozhan Dehghani<sup>2</sup>, Asa Borzabadi Farahani<sup>1</sup>, Pouya Bashivan<sup>1,2,3</sup>

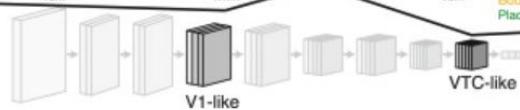
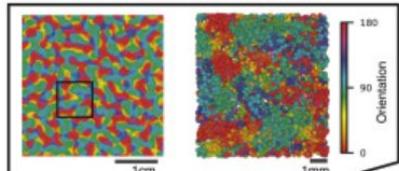
<sup>1</sup> Department of Computer Science, McGill University

<sup>2</sup> Department of Physiology, McGill University

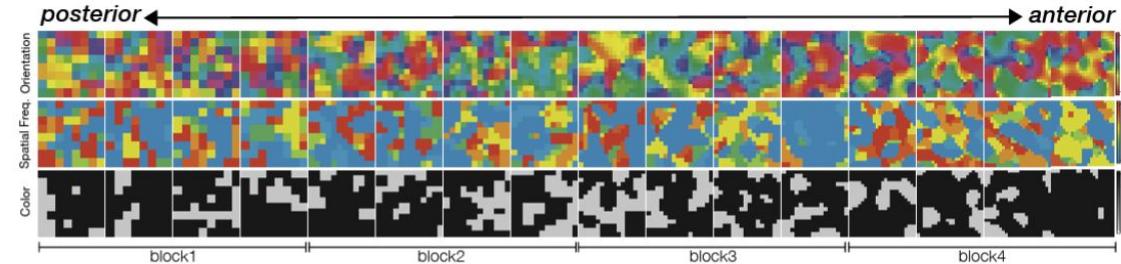
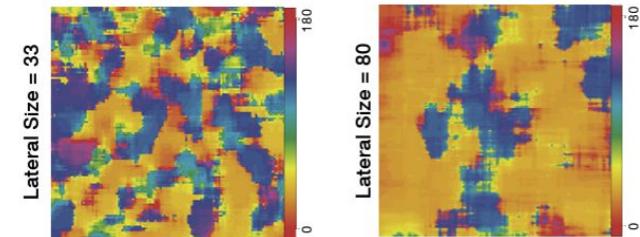
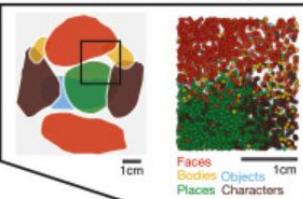
<sup>3</sup> Mila, Université de Montréal

# Both approaches work

V1: Orientation preference map smoothness, pinwheel density, smoothness of spatial frequency and chromatic maps



VTC: Category selectivity map smoothness, number of patches, patch size, colocalization of categories



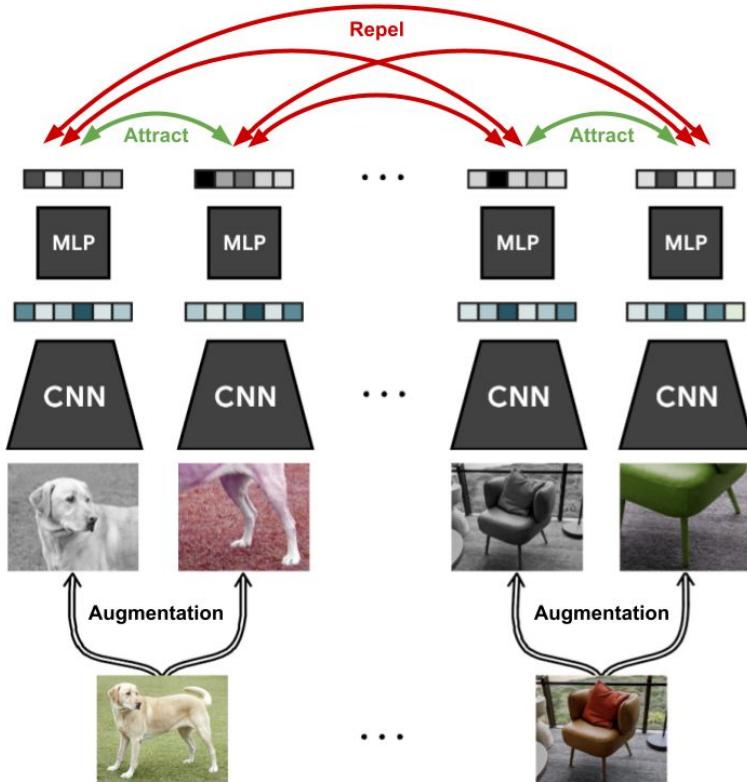
# Roadmap

- Comparing models to data
  - Compare model 1 to model 2: which one maps to data better? (**Case study 1**)
  - However, the platonic representation hypothesis!
- Closing the loop by synthesizing stimuli
  - Once a model predicts neural data well, it can be inverted, to give new stimuli (**Case study 2**)
  - Still a very new field
- Computational tricks
  - Contrastive learning
  - Regularizations (**Case study 3**)
- Bridging theories (different perspectives, same math!)
  - Attention mechanism to Memory
  - Regularization to Generativity (**Case study 4**)

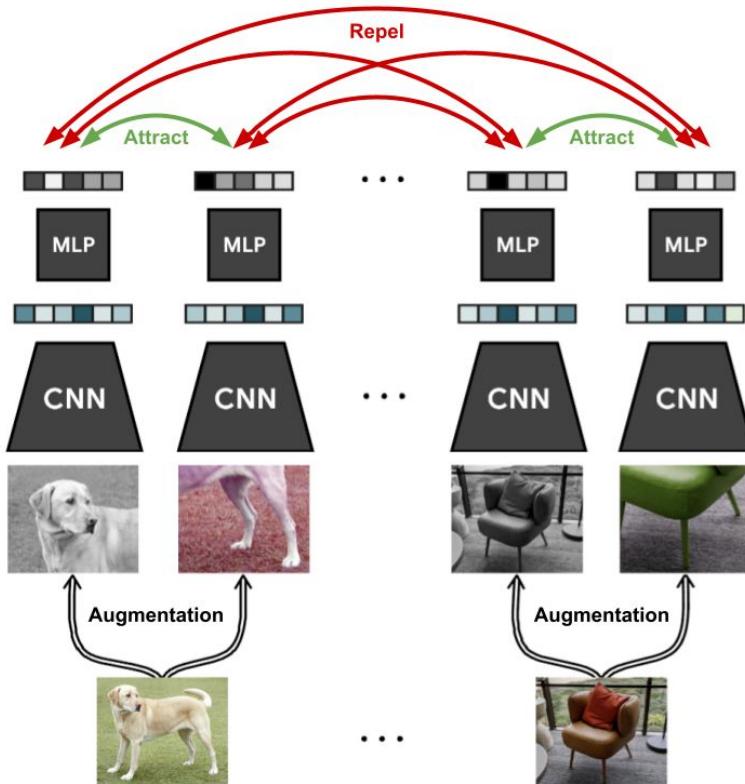
# Useful Computational hacks neuroscience learned from deep learning framework

- Contrastive learning
  - In Tools: Cebra
  - In Plasticity rules: Hebbian+predictive
- Attention
  - 3-factor plasticity rule using astrocytes
- Regularization
  - Temporal prediction can be learned from robustness to noise

# Contrastive learning



# Contrastive learning



# Regularization

## Optimization Goal

Task relevant	Task irrelevant
<b>Objective + Regularization</b>	

e.g.

Reconstruction  
Classification

e.g.

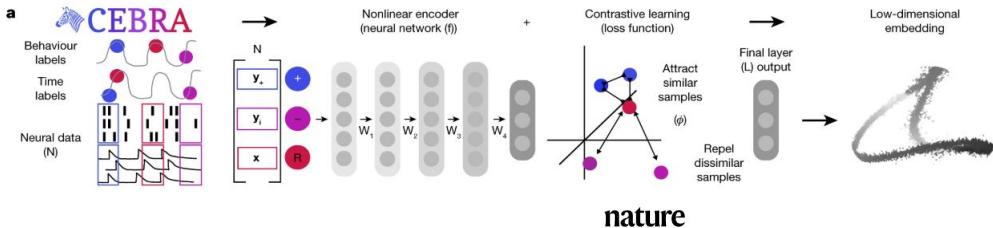
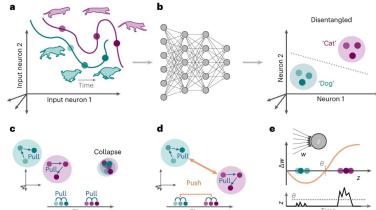
Sparse coding  
Temporal continuity  
Wiring minimization  
noise robustness

## Primary visual cortex properties

Gabor-like receptive fields  
Translation invariance  
Ocular dominance map  
Temporal prediction

# Useful Computational hacks neuroscience learned from deep learning framework

- Contrastive learning
  - In Tools: Cebra
  - In Plasticity rules: Hebbian+predictive



nature neuroscience

Explore content ▾ About the journal ▾ Publish with us ▾

nature > nature neuroscience > articles > article

Article | Open access | Published: 12 October 2023

## The combination of Hebbian and predictive plasticity learns invariant object representations in deep sensory networks

Manu Srinath Halvagal & Friedemann Zenke

nature

Explore content ▾ About the journal ▾ Publish with us ▾

nature > articles > article

Article | Open access | Published: 03 May 2023

## Learnable latent embeddings for joint behavioural and neural analysis

Steffen Schneider, Jin Hwa Lee & Mackenzie Weygandt Mathis

## Building transformers from neurons and astrocytes

Leo Kozachkov, Ksenia V. Kastanenka, and Dmitry Krotov Authors Info & Affiliations

Edited by Terrence Sejnowski, Salk Institute for Biological Studies, La Jolla, CA; received November 9, 2022; accepted June 22, 2023

- Attention
  - 3-factor plasticity rule using astrocytes
- Regularization
  - Temporal prediction can be learned from robustness to noise

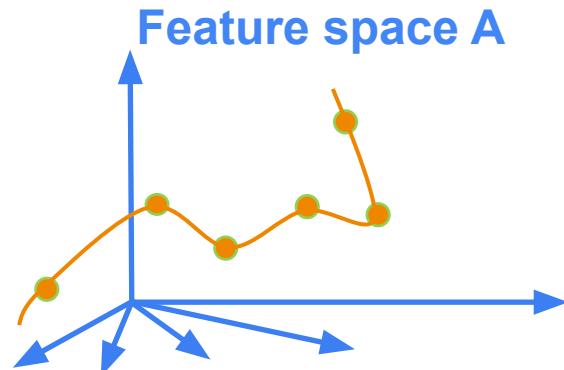
Case study 3:  
Temporal prediction can be achieved by  
regularized learning

# How to quantify the ability to predict over time?



[Henaff et al 2019]

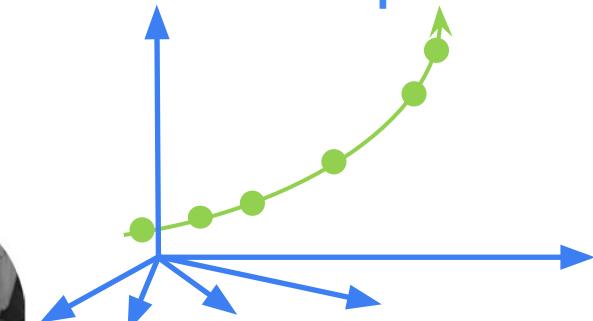
Feature space A



Pixel space

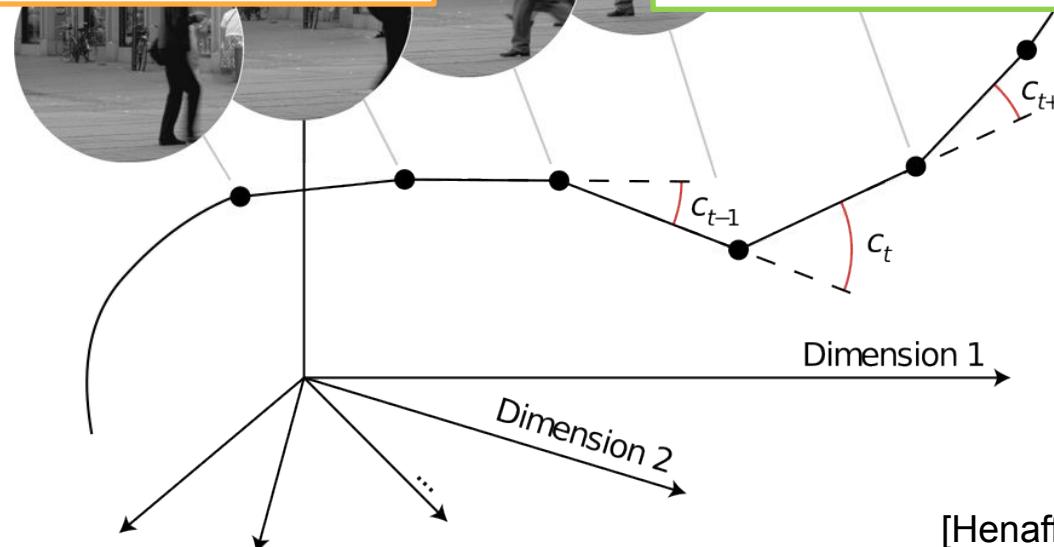


Feature space B



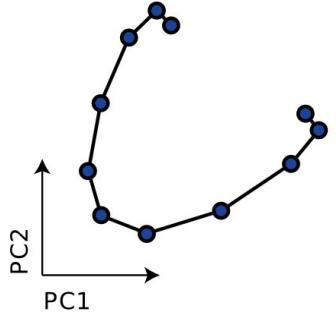
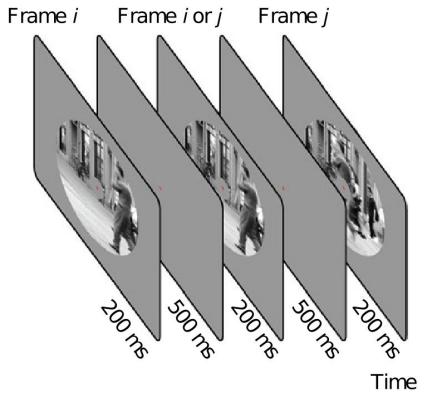
Temporal straightening X

Temporal straightening ✓



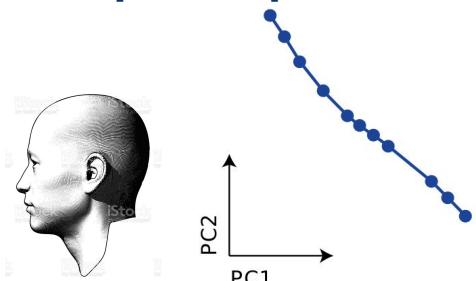
[Henaff et al 2019]

## Pixel space



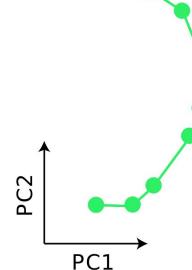
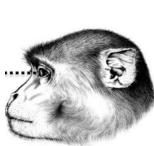
[Henaff et al 2019, 2021]

## Perceptual space of humans



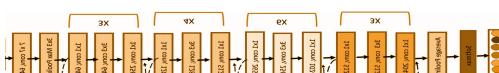
$\sim 30^\circ$

## Primary visual cortex of macaque



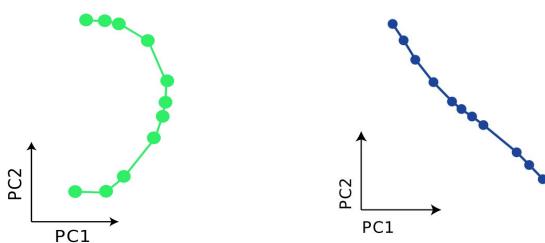
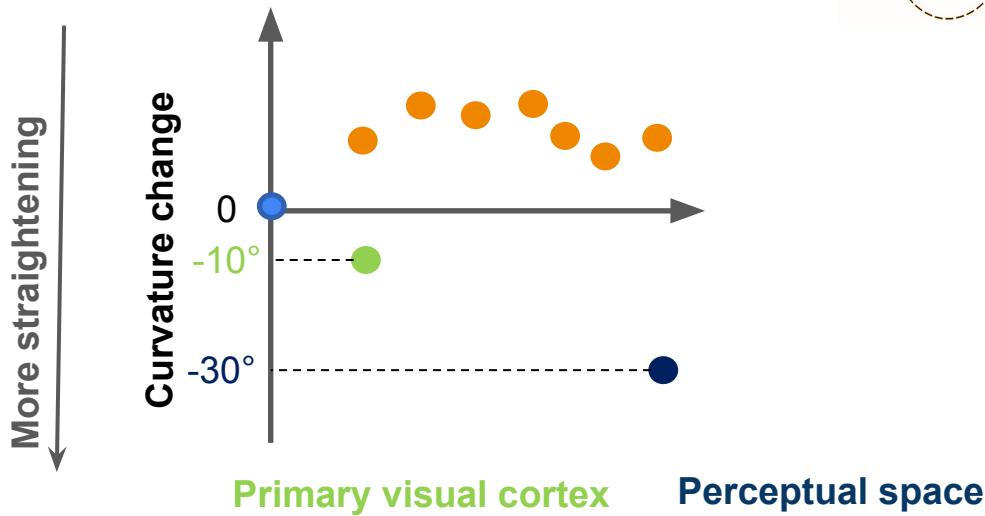
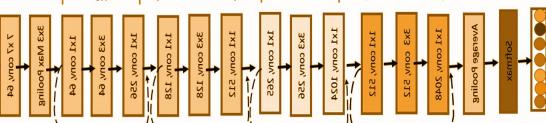
$\sim 10^\circ$

## Deep Neural Networks



?

## Deep Neural Networks

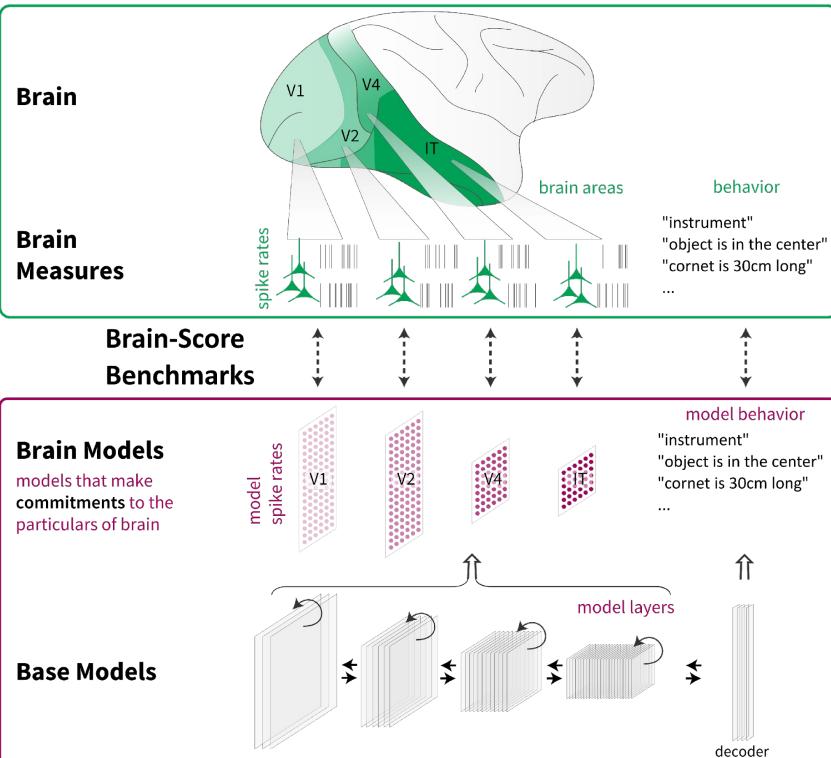
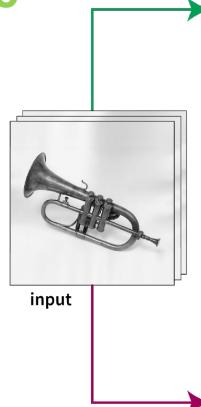


# Puzzling: Why deep neural net models of vision don't exhibit temporal straightening?

**Deep Nets explain variance in neural data**

**Deep Nets do not exhibit temporal straightening**

**Brain-like**

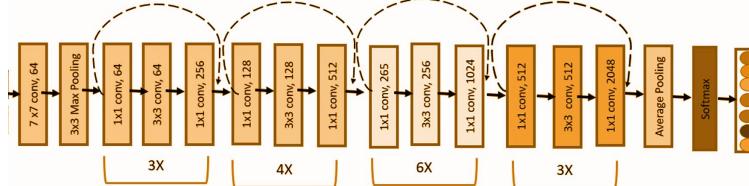


**Previous conjecture:** DNNs needs to be trained to predict the next frame to be able to straightened the temporal trajectories

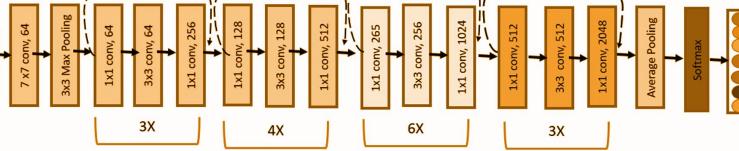
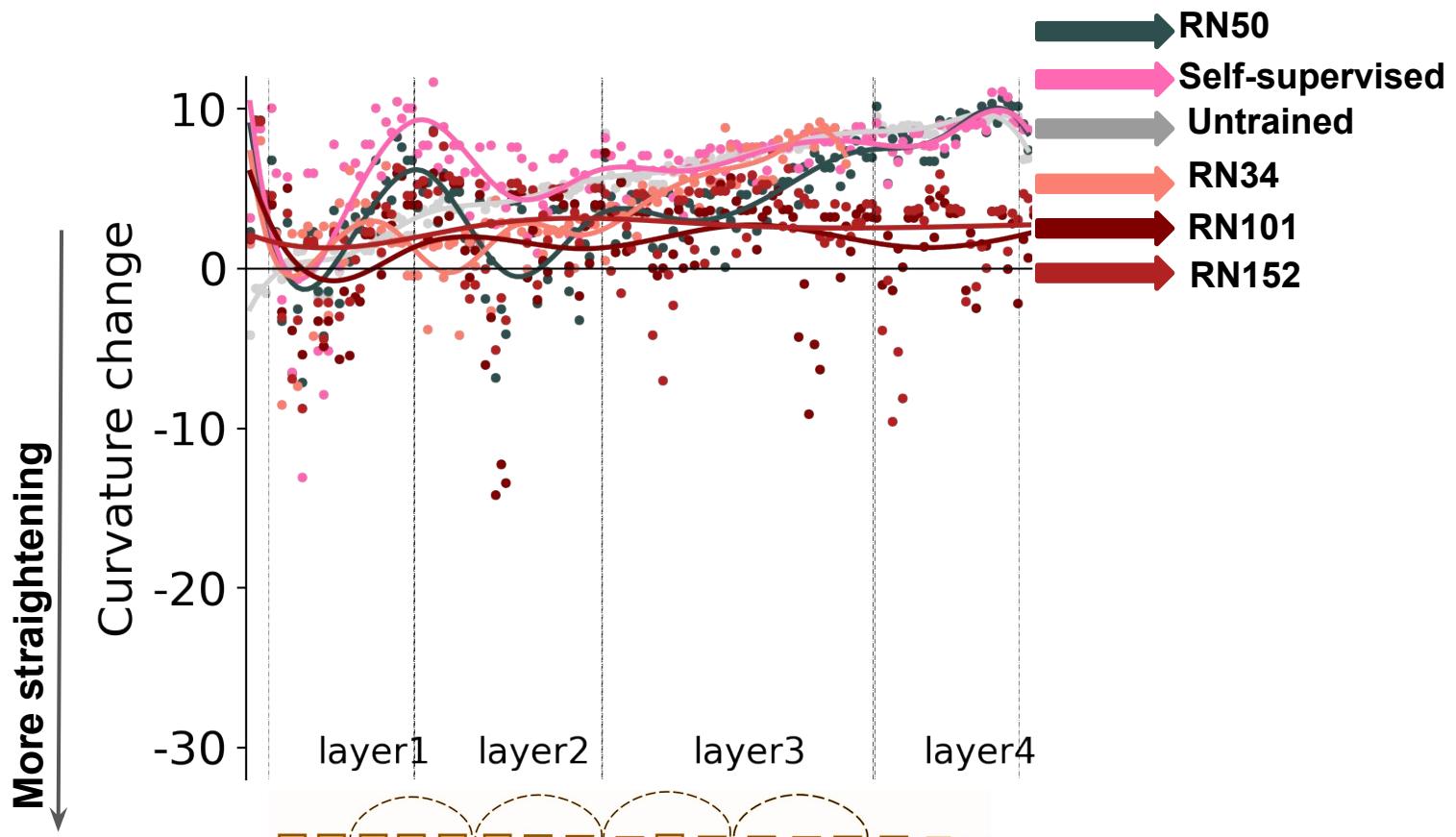
[Henaff, et al 2019, 2021]

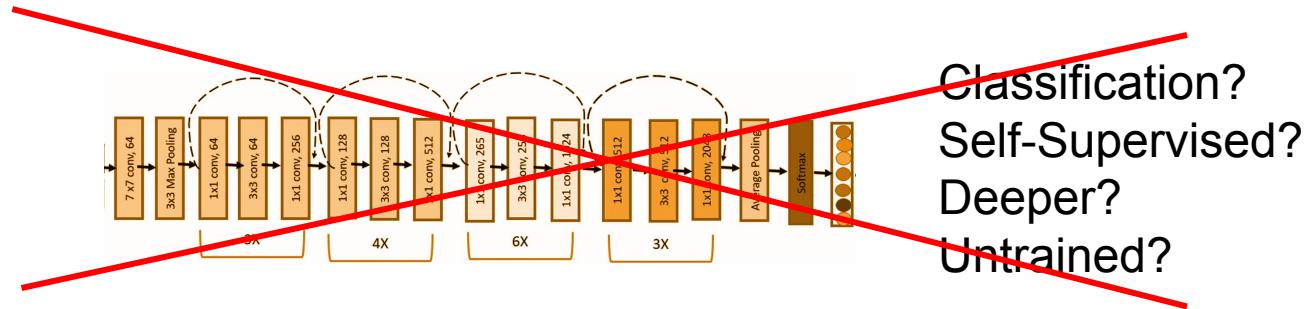
[Yamins & Hong, et al, 2014, Schrimpf et al. 2020]

Can any members of the feedforward DNN class straighten the trajectory of natural movies?



Classification?  
Self-Supervised?  
Deeper?  
Untrained?  
....

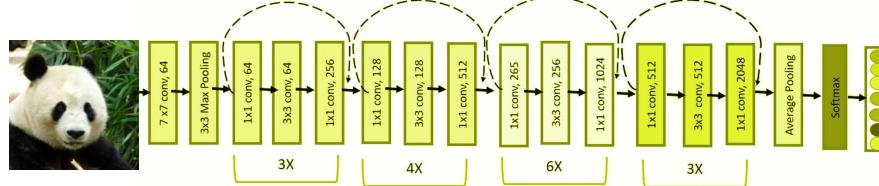




Classification?  
Self-Supervised?  
Deeper?  
Untrained?

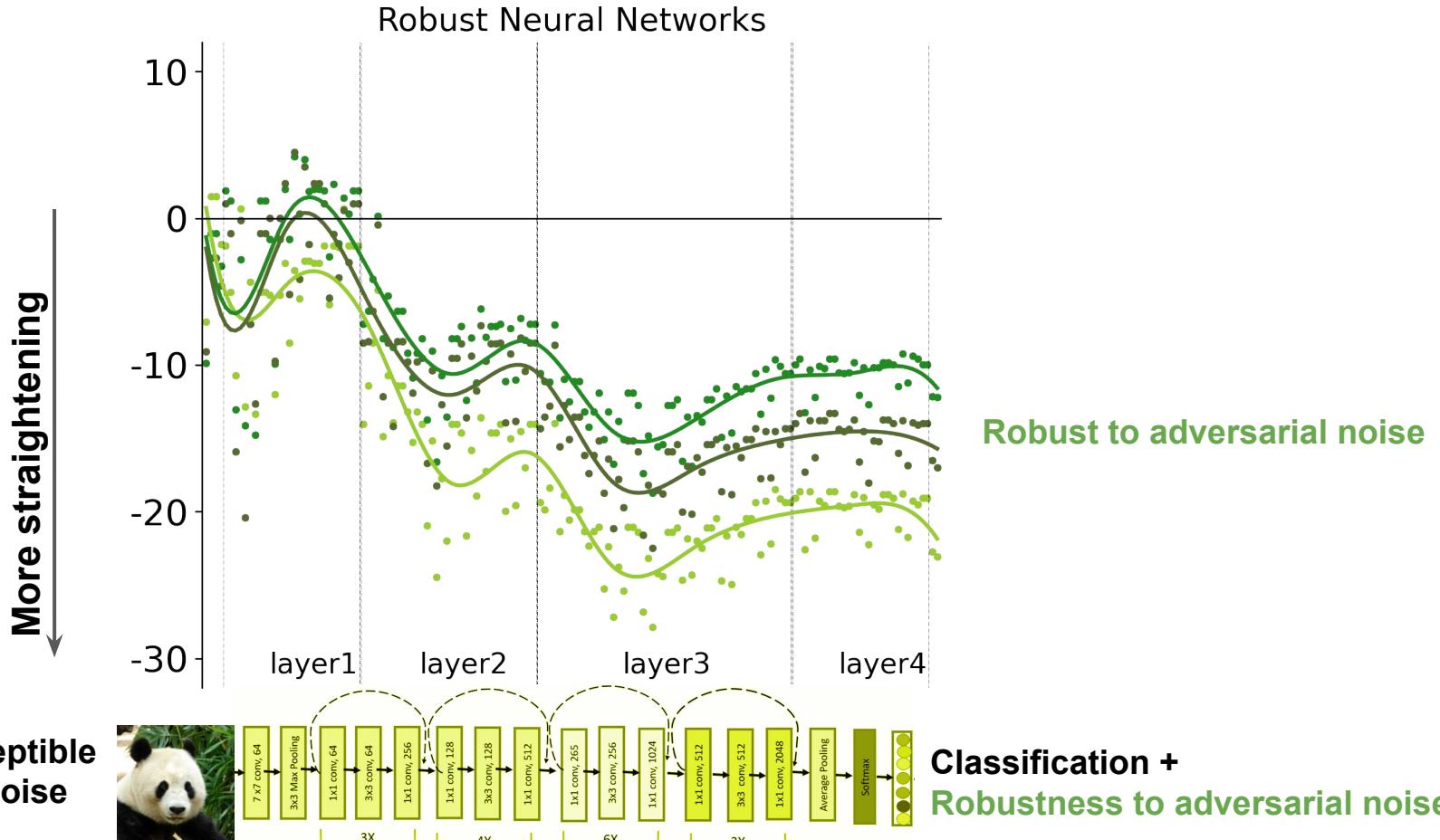
Next, we considered how ***training for noise robustness*** may affect manifold straightness in feedforward DNNs

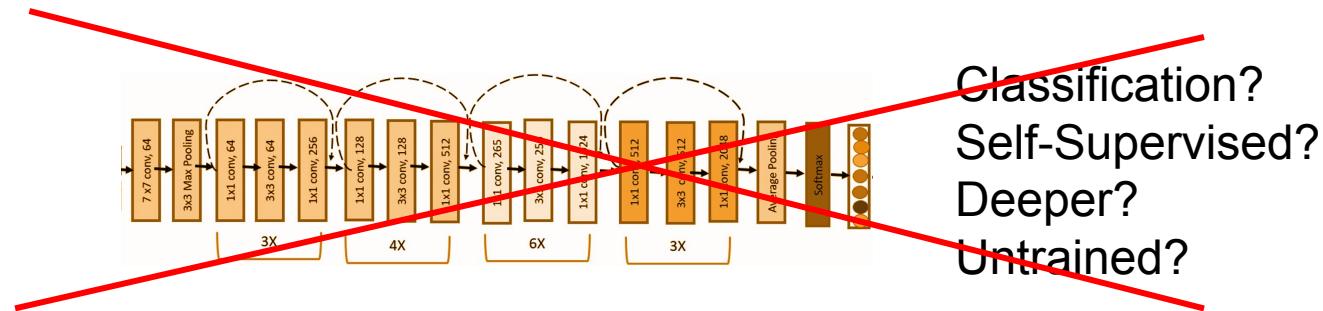
+ Imperceptible  
engineered noise



Classification +  
Robustness to adversarial noise

# Robustness to adversarial noise give rise to temporal straightening

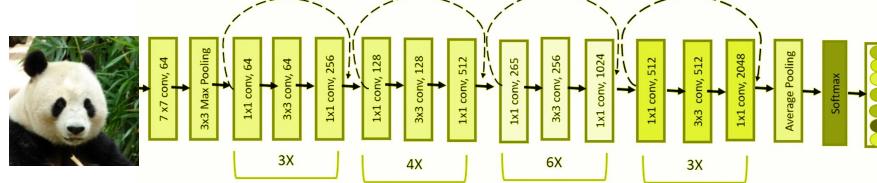




Classification?  
Self-Supervised?  
Deeper?  
Untrained?

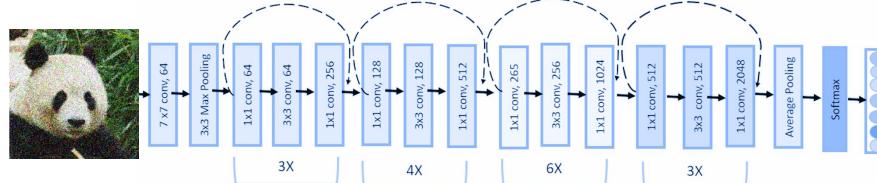
Next, we considered how ***training for noise robustness*** may affect manifold straightness in feedforward DNNs

- + Imperceptible engineered noise



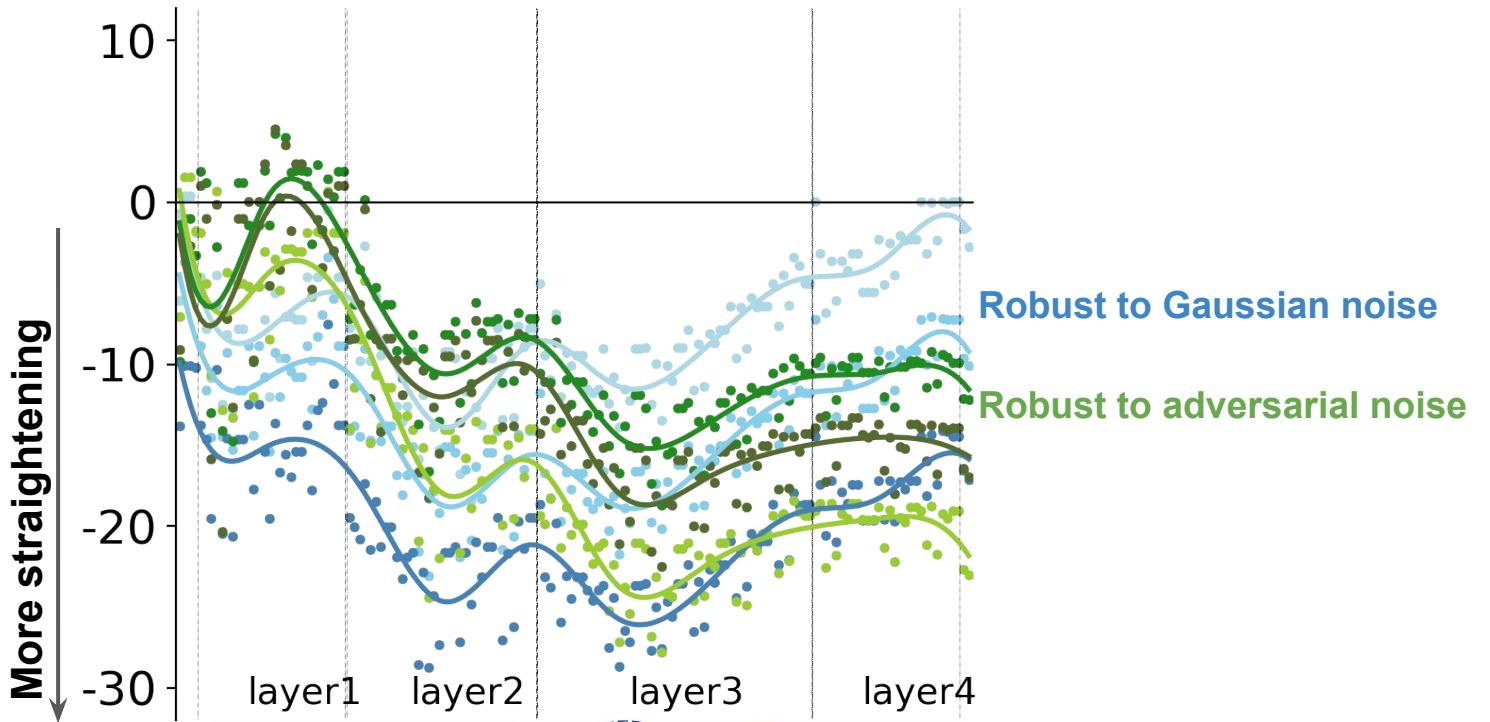
Classification +  
Robustness to adversarial noise

- + Random Gaussian noise

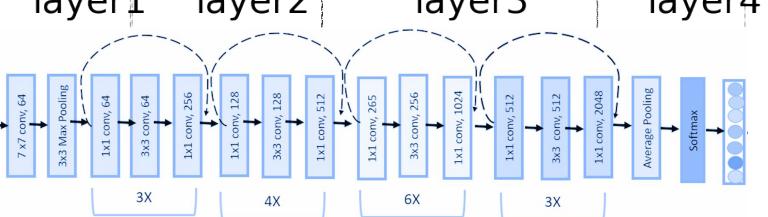


Classification +  
Robustness to Gaussian noise

# Robustness to noise give rise to temporal straightening

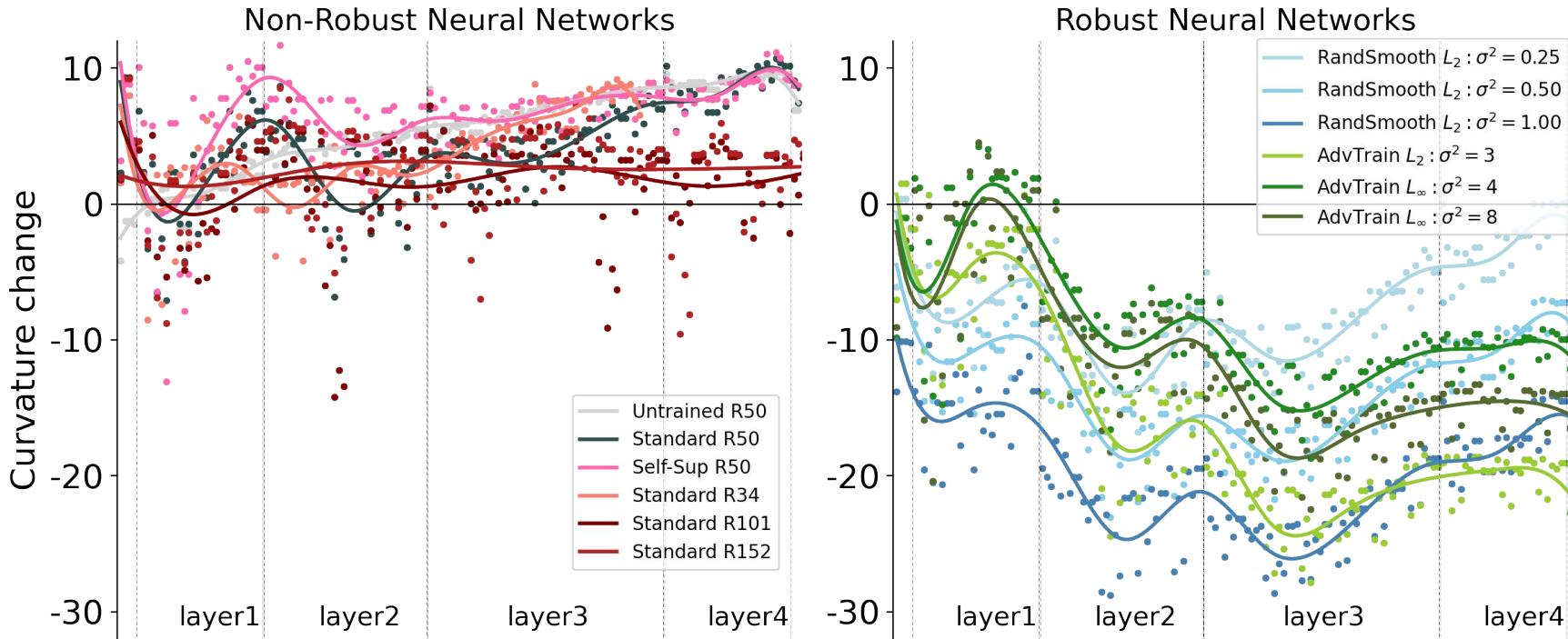


+ Random Gaussian noise



Classification +  
Robustness to Gaussian noise

# Feedforward neural networks can exhibit straightening...



If they were trained to be robust to input noise!

# How is straightening related to neural predictivity?

**Standard Deep Nets**

**Robust Deep Nets**

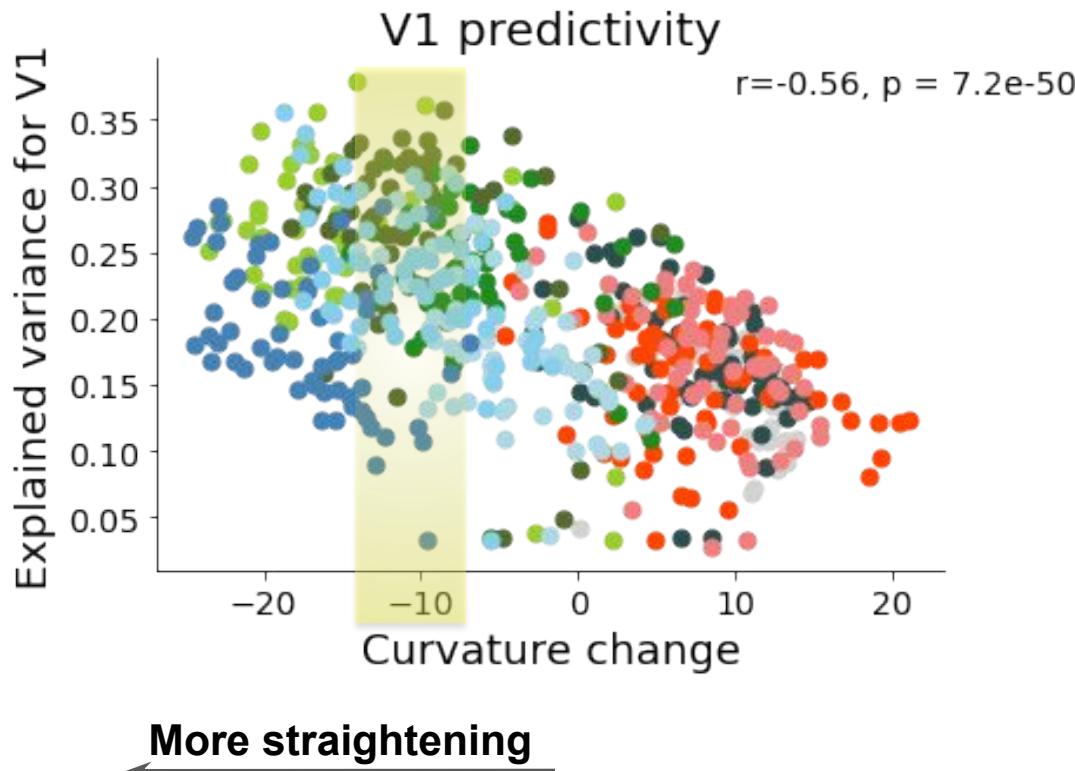
Predict neural data



temporal straightening

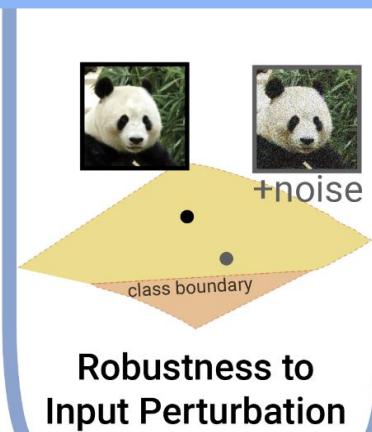
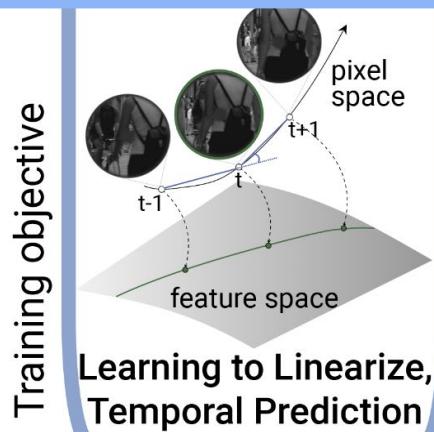


# How is straightening related to neural predictivity?



# What is the intuition (and theory) behind the link?

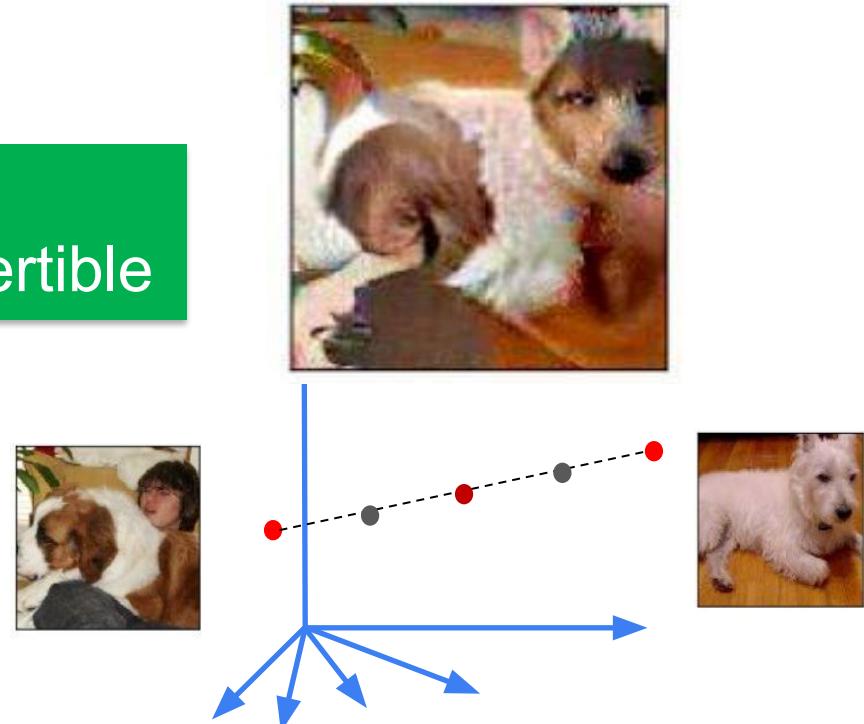
## temporal straightening & robustness to noise



**Invertibility:** Linear interpolation between the representations in the feature space corresponds to natural features

**Invertible feature space**

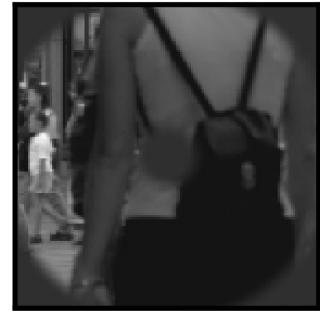
Prior finding:  
Robust neural networks are invertible



[Engstrom et al, 2018]

## Invertible feature space

Robust neural networks are invertible  
when interpolating categories but what  
about in natural movies?



# Inverting robust vs. non-robust features of natural movies

Interpolation in  
**Robust** feature space

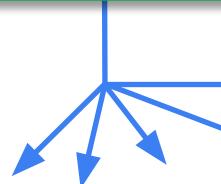
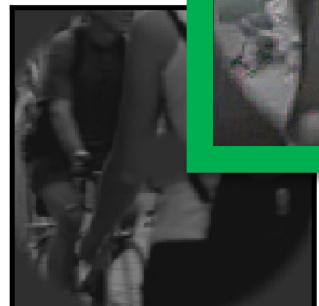
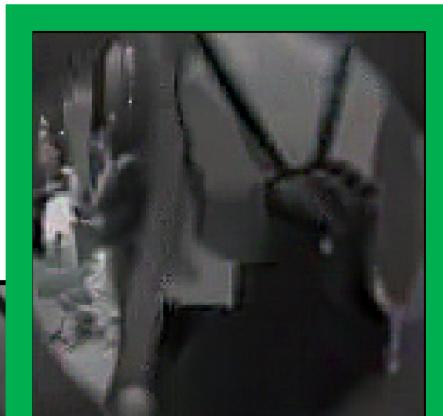


Interpolation in  
**Non-robust** feature space

# Quantitative comparison of interpolated vs. actual frame

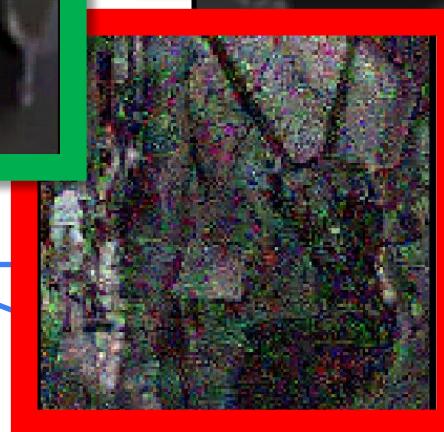
Interpolation in  
**Robust** feature space

SSIM = 0.52

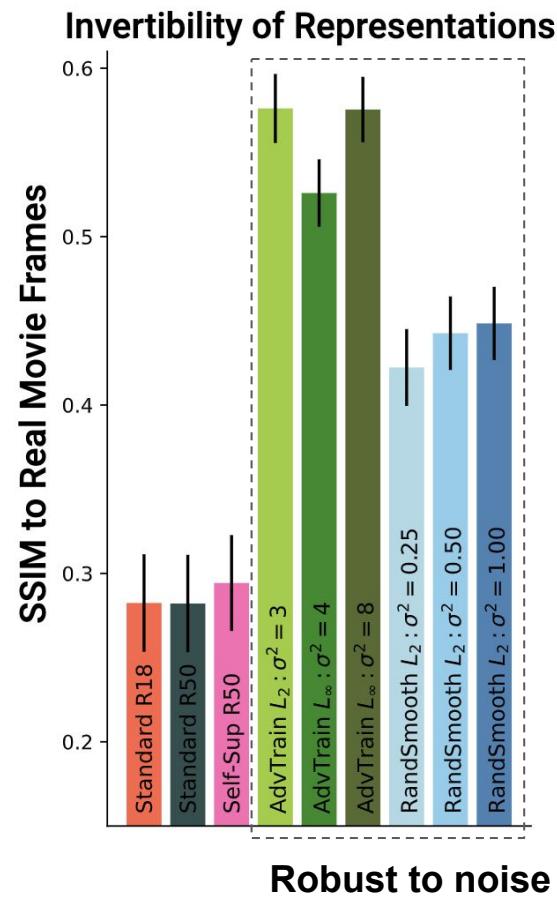
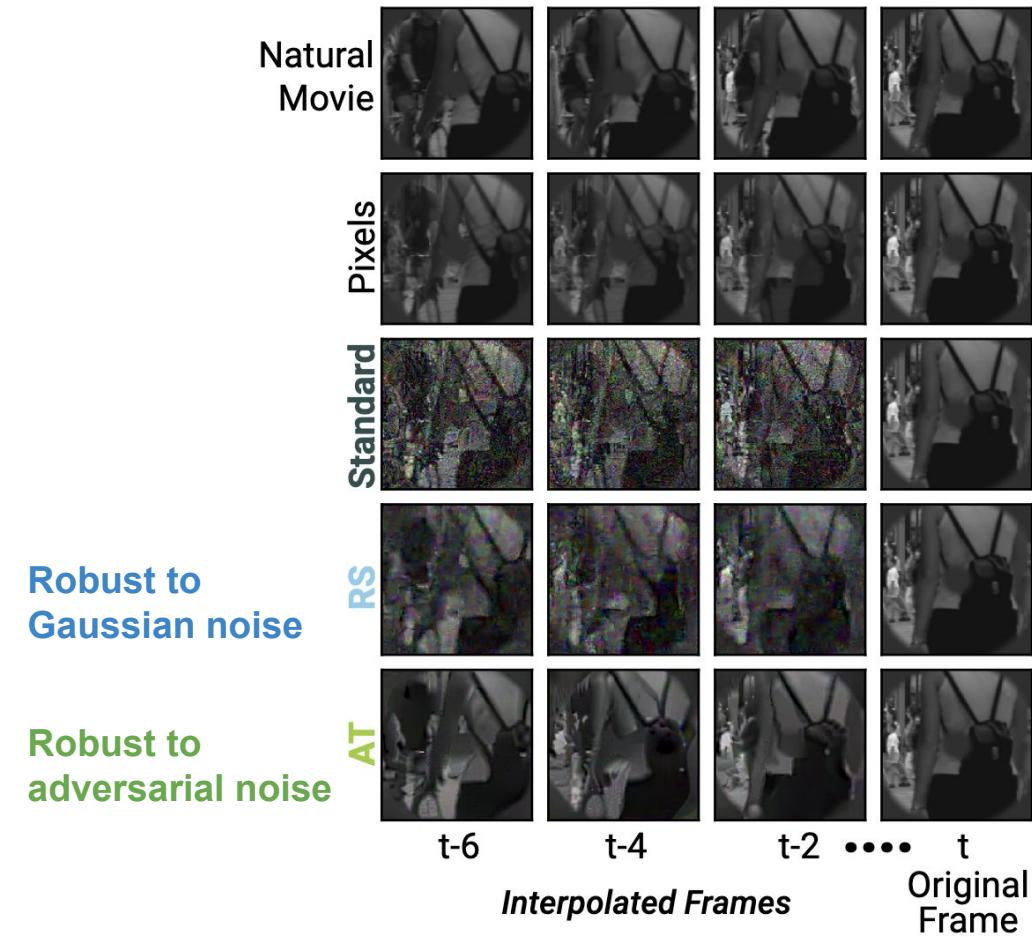


Real in-between  
frame

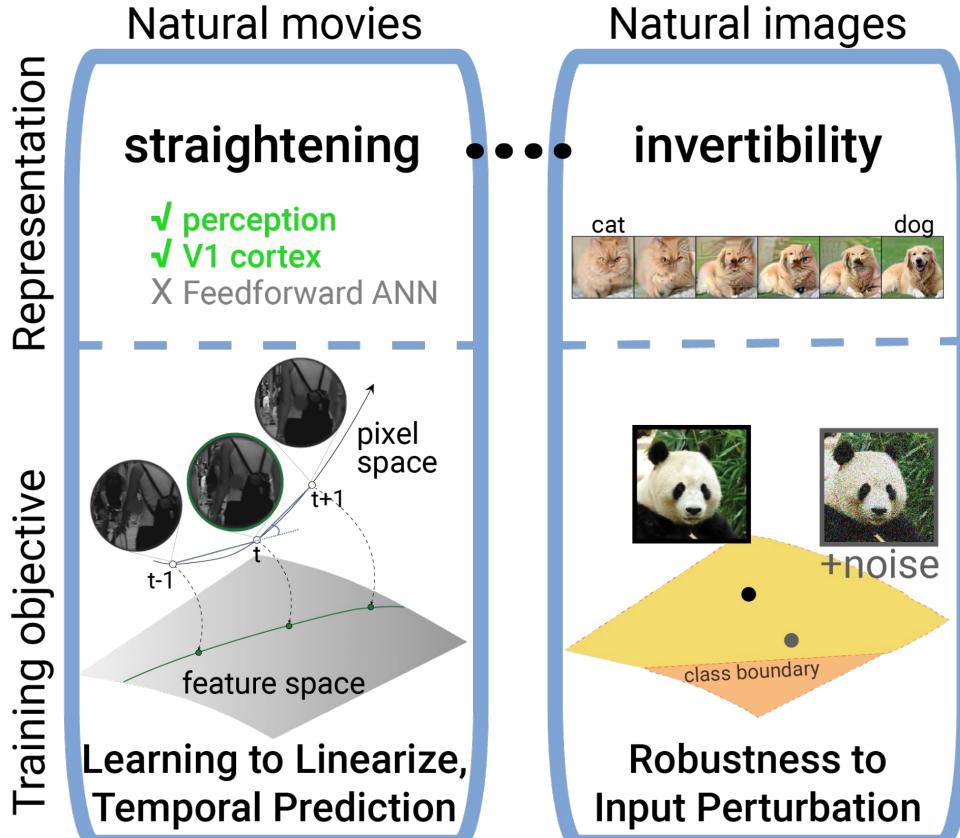
SSIM = 0.21



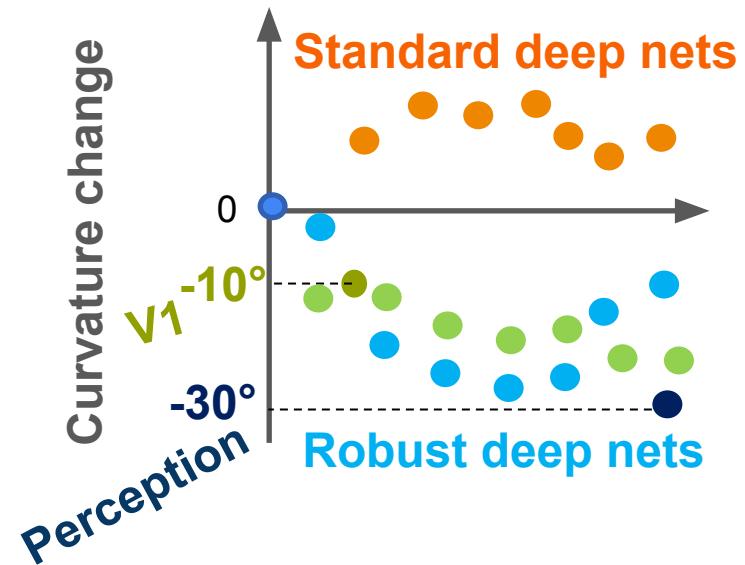
Interpolation in  
**Non-robust** feature  
space



# Summary



Robustly trained deep nets demonstrate natural movie straightening in their representations, without directly training on movies, or for prediction



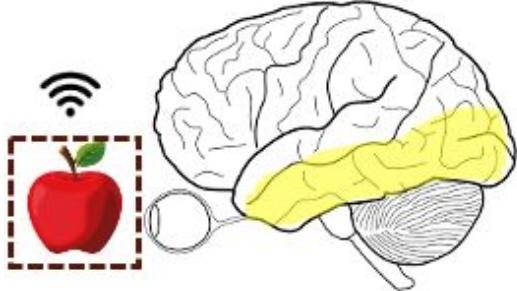
# Roadmap

- Comparing models to data
  - Compare model 1 to model 2: which one maps to data better? (**Case study 1**)
  - However, the platonic representation hypothesis!
- Closing the loop by synthesizing stimuli
  - Once a model predicts neural data well, it can be inverted, to give new stimuli (**Case study 2**)
  - Still a very new field
- Computational tricks
  - Contrastive learning
  - Regularizations (**Case study 3**)
- Bridging theories (different perspectives, same math!)
  - Recognition | Generation
  - Learning | Attention

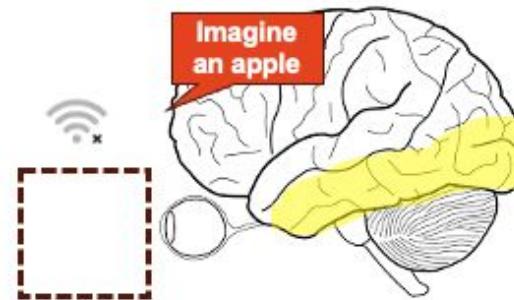
Theory bridges seemingly different computations

- How does the same neural circuitry support classification and generation?
  - (Regularized) discriminative models (i.e. classifiers) → Generative models
- Plasticity  $\longleftrightarrow$  Memory
  - Key-value attention in transformers → Memory storage and retrieval

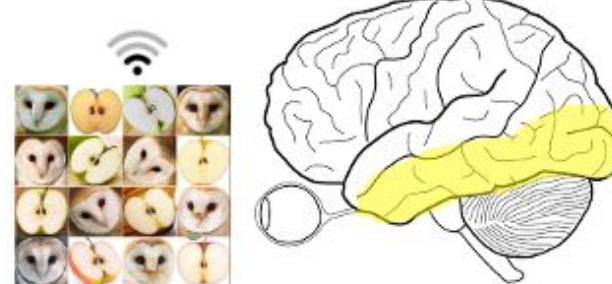
Sensory driven



Prior driven

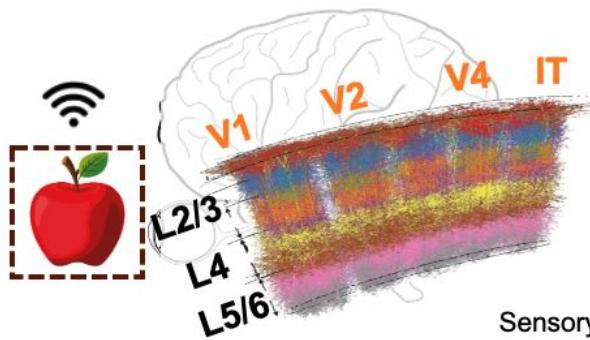


Sensory and prior integration

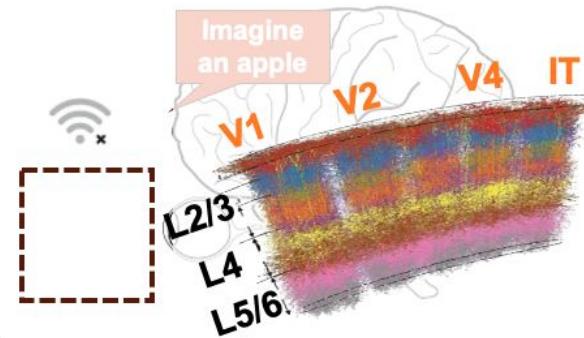


# Same visual system supports both classification and generation

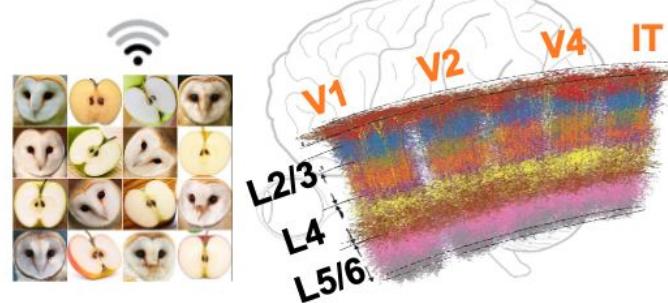
Sensory driven



Prior driven

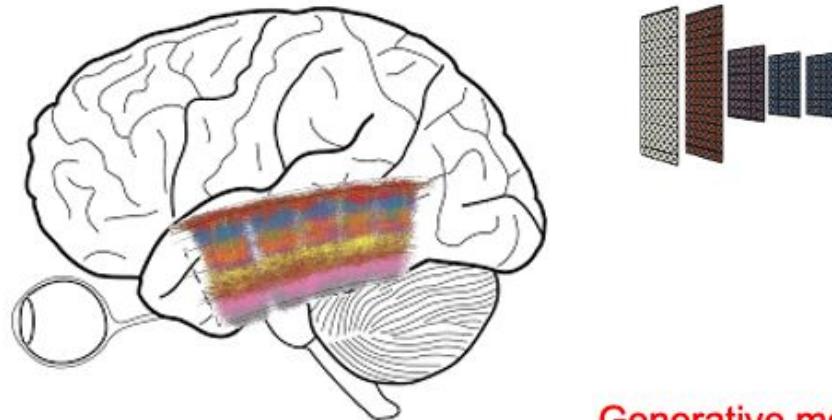


Sensory and prior integration

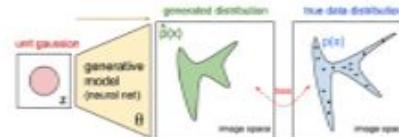


# But different classes of models for each capability...

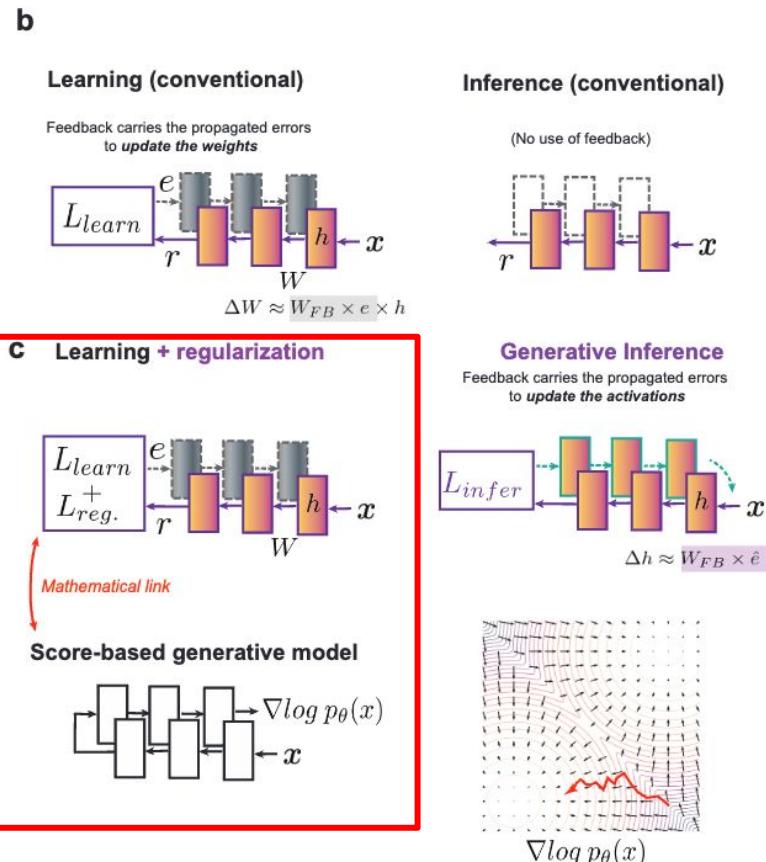
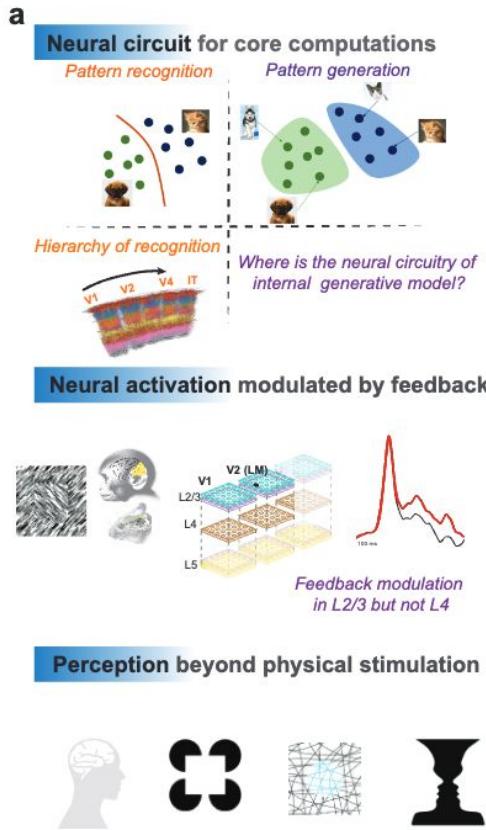
Discriminative models

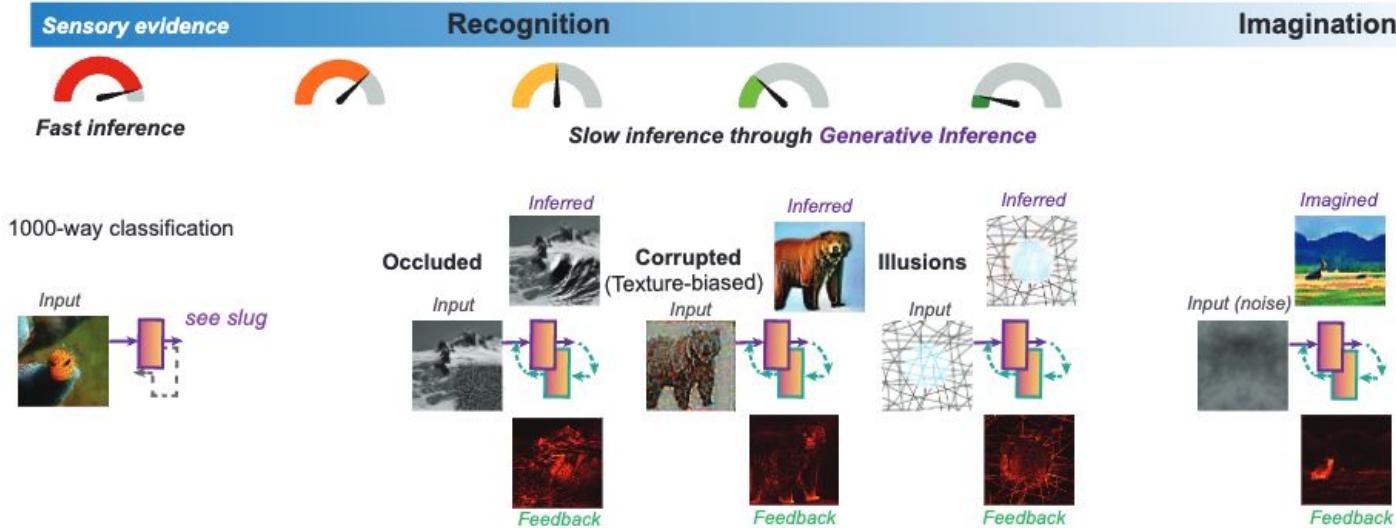


Generative models



# Recognition $\longleftrightarrow$ Generation





Duality of discrimination and generation we established using deep learning framework, helped to understand the puzzling role of feedback in vision.

Generative perceptual inference in deep neural network models of object recognition induces illusory contours and shapes

Kazuki Irie\*, Róbert Csordás\*, Jürgen Schmidhuber  
**The Dual Form of Neural Networks Revisited**  
**ICML 2022, <https://arxiv.org/abs/2202.05798>**

These systems are “equivalent”

*Arbitrary  $v$  and  $k$*

Linear layer

$$\mathbf{W} = \sum_{t=1}^T \mathbf{v}_t \otimes \mathbf{k}_t$$

$$S_1(\mathbf{x}) = \mathbf{W}\mathbf{x}$$

Key-value memory attention layer

$$\mathbf{K} = (\mathbf{k}_1, \dots, \mathbf{k}_T) \in \mathbb{R}^{d_{\text{in}} \times T}$$

$$\mathbf{V} = (\mathbf{v}_1, \dots, \mathbf{v}_T) \in \mathbb{R}^{d_{\text{out}} \times T}$$

$$S_2(\mathbf{x}) = \text{Attention}(\mathbf{K}, \mathbf{V}, \mathbf{x})$$

NB: Unnormalised dot attention  $\text{Attention}(\mathbf{K}, \mathbf{V}, \mathbf{q}) = \mathbf{V}\mathbf{K}^\top \mathbf{q}$   
 Standard attention:  $\mathbf{V}\text{softmax}(\mathbf{K}^\top \mathbf{q})$

[Slide from Katzuki Irie]

## Application to a Linear Layer Trained by Gradient Descent:

Forward computation:

$$\mathbf{y} = \mathbf{W}\mathbf{x}$$

Backward computation (gradient descent) to update  $\mathbf{W}$ :

$$\mathbf{W}_{t+1} = \mathbf{W}_t - \eta_t (\nabla_{\mathbf{y}} \mathcal{L})_t \otimes \mathbf{x}_t$$

$e_t$

*t now denotes training iteration!*

for some error function  $\mathcal{L}$  learning rate  $\eta_t$  at step  $t \in \mathbb{N}$

We can directly apply the duality from the previous slide to:

$$\mathbf{W} = \mathbf{W}_0 + \sum_{t=1}^T e_t \otimes \mathbf{x}_t$$

Linear layer trained by  
gradient descent

**Store:**  $\mathbf{W} = \mathbf{W}_0 + \sum_{t=1}^T \mathbf{e}_t \otimes \mathbf{x}_t$

**Compute:**  $S_1(\mathbf{x}) = \mathbf{W}\mathbf{x}$

Key/value-attention memory  
storing entire training experience

$$\mathbf{E} = (\mathbf{e}_1, \dots, \mathbf{e}_T) \in \mathbb{R}^{d_{\text{out}} \times T}$$

$$\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_T) \in \mathbb{R}^{d_{\text{in}} \times T}$$

$$S_2(\mathbf{x}) = \mathbf{W}_0\mathbf{x} + \text{Attention}(\mathbf{X}, \mathbf{E}, \mathbf{x})$$

## Linear layer trained by gradient descent

**Store:**  $\mathbf{W} = \mathbf{W}_0 + \sum_{t=1}^T \mathbf{e}_t \otimes \mathbf{x}_t$

## Key-value memory in the brain

Samuel J. Gershman<sup>1,2,3,\*</sup>, Ilia Fiete<sup>4</sup>, and Kazuki Irie<sup>2</sup>

<sup>1</sup>Department of Psychology

<sup>2</sup>Center for Brain Science

<sup>3</sup>Kempner Institute for the Study of Natural and Artificial Intelligence,  
Harvard University, Cambridge, MA, USA

<sup>4</sup>McGovern Institute for Brain Research and Department of Brain and Cognitive Sciences,  
Massachusetts Institute of Technology

\*Corresponding author: gershman@fas.harvard.edu

## Key/value-attention memory storing entire training experience

$$\mathbf{E} = (\mathbf{e}_1, \dots, \mathbf{e}_T) \in \mathbb{R}^{d_{\text{out}} \times T}$$

$$\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_T) \in \mathbb{R}^{d_{\text{in}} \times T}$$

$$S_2(\mathbf{x}) = \mathbf{W}_0 \mathbf{x} + \text{Attention}(\mathbf{X}, \mathbf{E}, \mathbf{x})$$

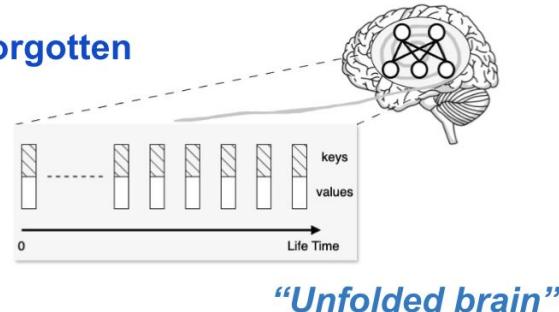
[Slide from Katzuki Irie]

# Reminiscent of some biological findings?

In the model above, key-value pairs are **never forgotten**

**BUT** some of them **become inaccessible** as  
the “query” signals to activate them are lost.

**Forgetting = “Failure to remember”**



Key-value pairs ~=  
“Memory engrams”

→ “Silent engram”?



## Silent memory engrams as the basis for retrograde amnesia

Dheeraj S. Roy<sup>a,†</sup>, Shruti Muralidhar<sup>a</sup>, Lillian M. Smith<sup>a</sup>, and Susumu Tonegawa<sup>a,b,c,2</sup>

<sup>a</sup>RIKEN-Massachusetts Institute of Technology, Department of Biology, and <sup>b</sup>Massachusetts Institute of Technology, Cambridge, MA 02139, USA; and <sup>c</sup>RIKEN Brain Science Institute, Wako, Saitama 351-0198, Japan

Contributed by Susumu Tonegawa

## Engram cells retain memory under retrograde amnesia

Tomás J. Ryan,<sup>1,2,\*</sup> Dheeraj S. Roy,<sup>1,\*</sup> Michele Pignatelli,<sup>1,\*</sup> Autumn Arons,<sup>1,2</sup> Susumu Tonegawa<sup>1,2,†</sup>

86

# Pointers to other insightful dualities

Generative diffusion models ~ Modern hopfield networks

---

**In search of dispersed memories: Generative diffusion models are associative memory networks**

---

**Luca Ambrogioni**

Radboud University, Donders Institute for Brain, Cognition and Behaviour  
[luca.ambrogioni@donders.ru.nl](mailto:luca.ambrogioni@donders.ru.nl)

Variational autoencoders ~ Generative diffusion models

---

**A Variational Perspective on Diffusion-Based Generative Models and Score Matching**

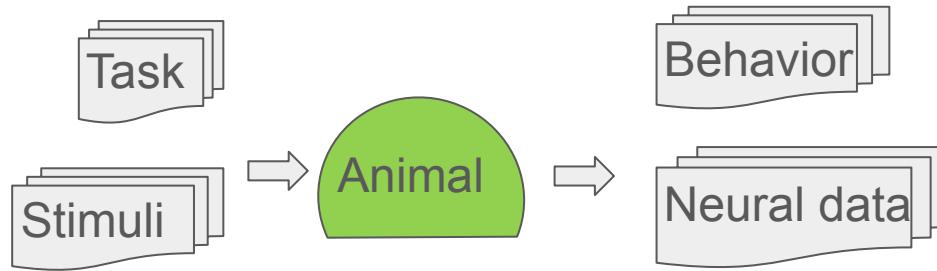
---

**Chin-Wei Huang, Jae Hyun Lim, Aaron Courville**

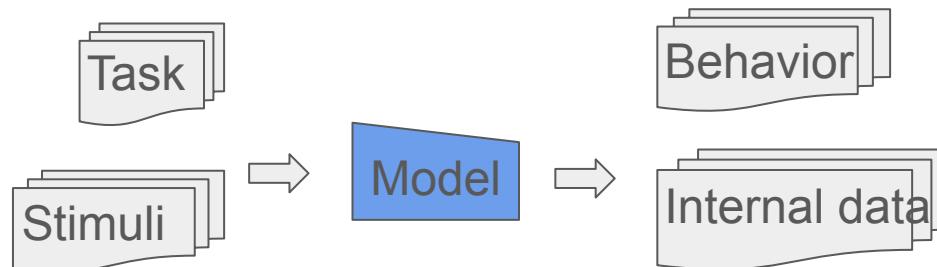
University of Montreal & Mila

{chin-wei.huang, jae.hyun.lim, aaron.courville}@umontreal.ca

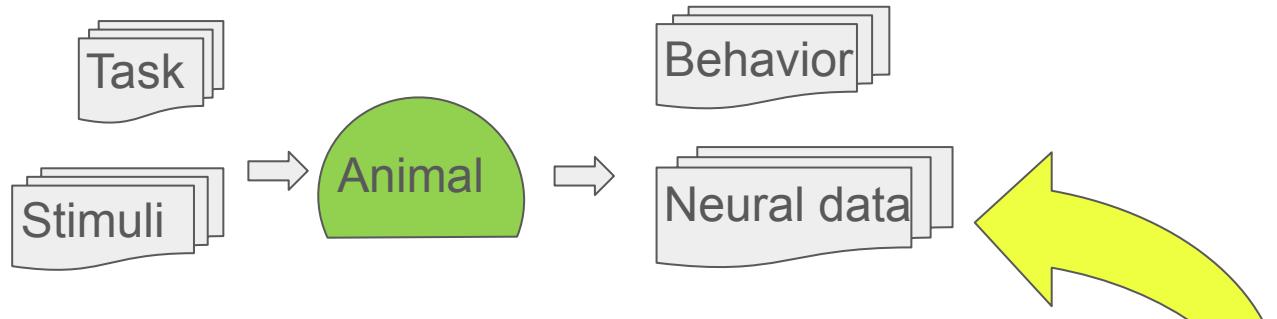
# Systems Neuroscience



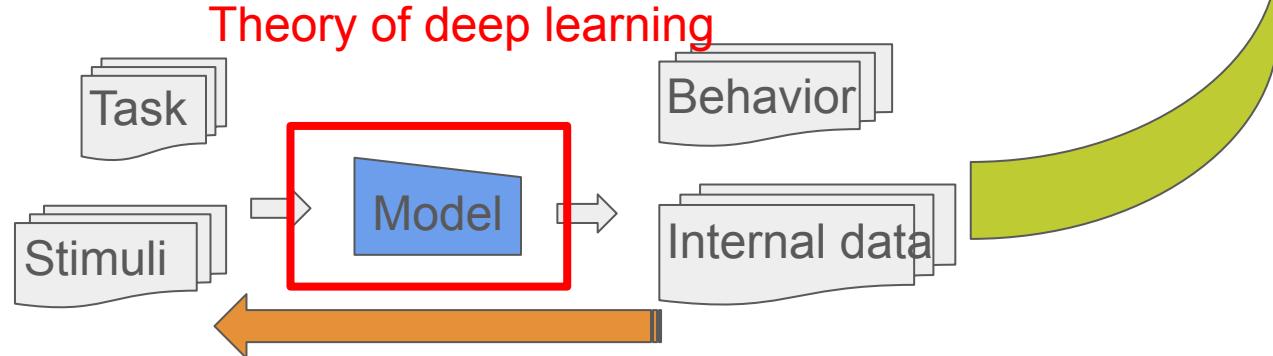
## Task-optimized models



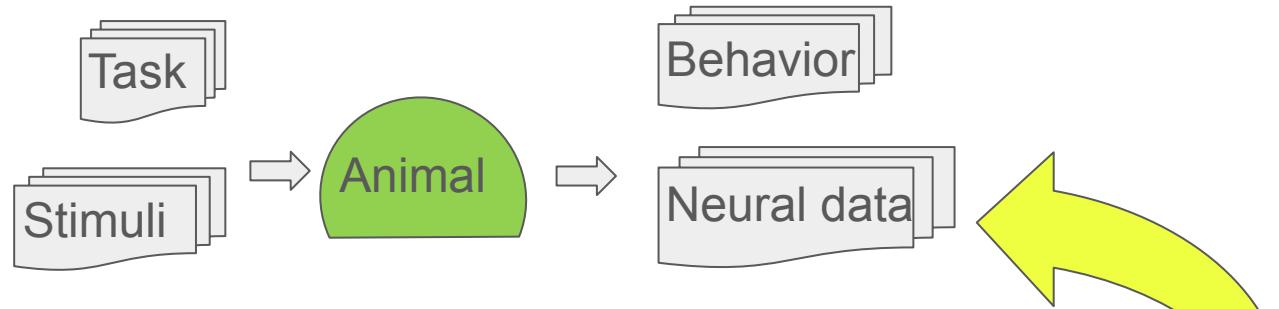
# Systems Neuroscience



## Task-optimized models



# Systems Neuroscience



3. Learning core computations from models

## Task-optimized models

Theory of deep learning

