

# Using Embodied Agents to Reverse-Engineer Natural Intelligence

---

**Aran Nayebi**

*Assistant Professor*

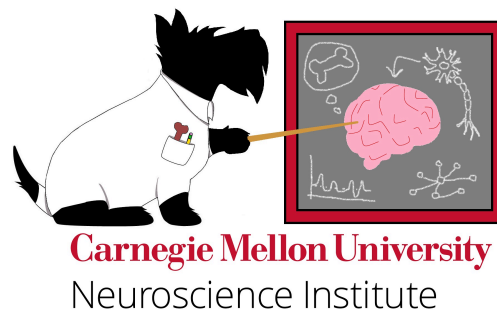
*Machine Learning Department*

*Neuroscience Institute (core faculty), Robotics Institute (by courtesy)*

**CS 375/Psych 249**

*Stanford University*

*2026.02.23*



# Current AI Struggles to Understand the Physical World

OpenAI Sora,  
February 2024



[Research](#) ▾ [API](#) ▾ [ChatGPT](#) ▾ [Safety](#) [Company](#) ▾

[Search](#)

[Log in](#) ↗

[Try ChatGPT](#) ↗

## Creating video from text

Sora is an AI model that can create realistic and imaginative scenes from text instructions.

[Read technical report](#)

All videos on this page were generated directly by Sora without modification.

# Current AI Struggles to Understand the Physical World

OpenAI Sora,  
February 2024



Q: What's missing?

A: Embodied agency & interaction.

# Why Reverse-Engineer Natural Intelligence?

## Why?

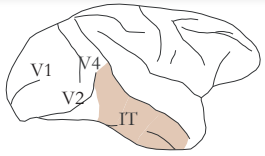
Animals & humans (currently) perform behaviors we've yet to engineer successfully in AI agents:

- Prediction (requires **world modeling**) & planning (requires **memory**)
- Adaptive motor control (requires **embodiment**)
- **Autonomy** / online **life-long** learning (test-time reasoning is just the beginning: need to update the weights without forgetting everything!)

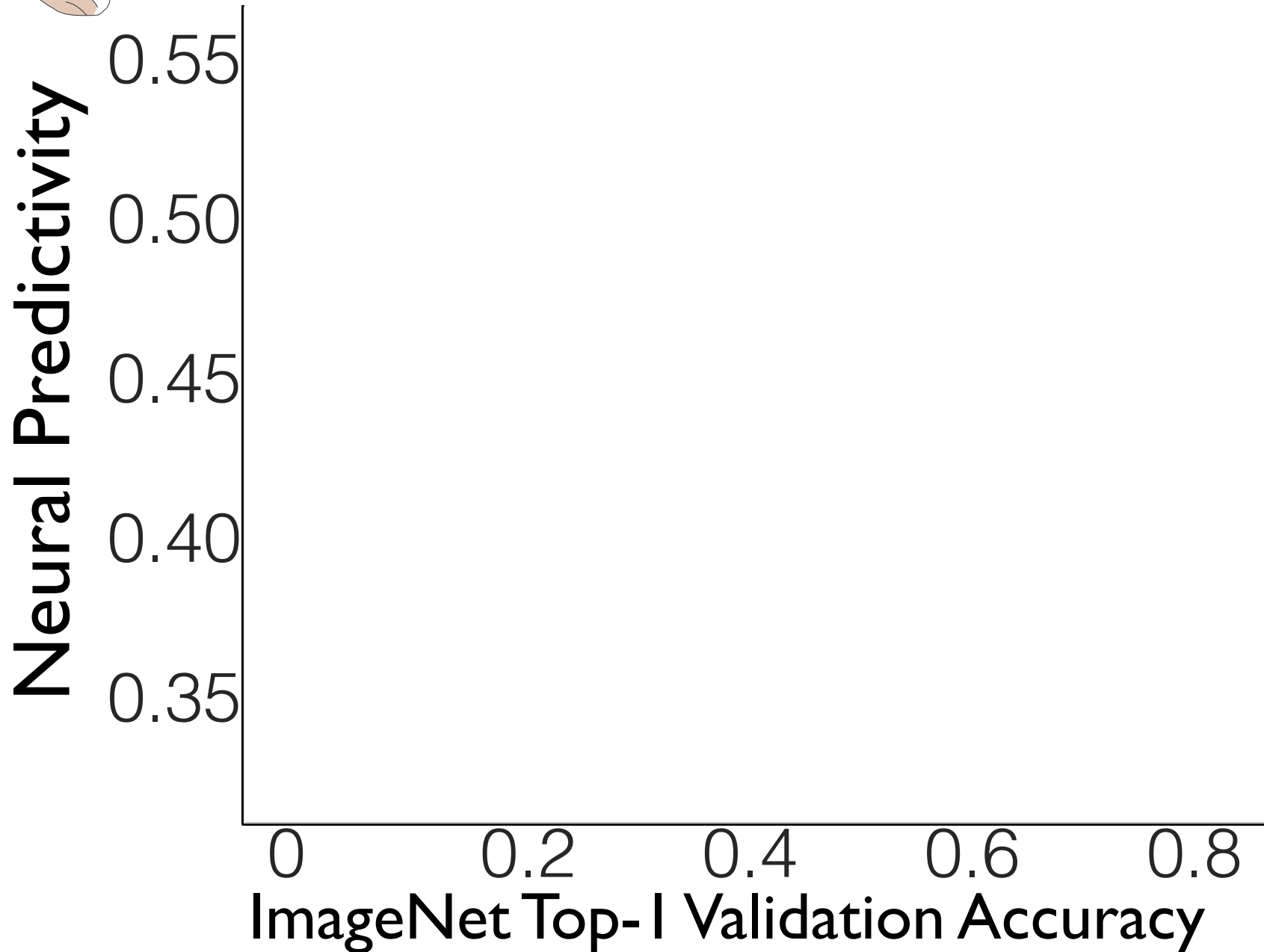


The specific *capabilities* of humans & animals become our concrete engineering **targets!**

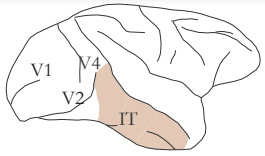
# Task Performance Correlated with Neural Predictivity



Schrimpf\*, Kubilius\* et al. 2018

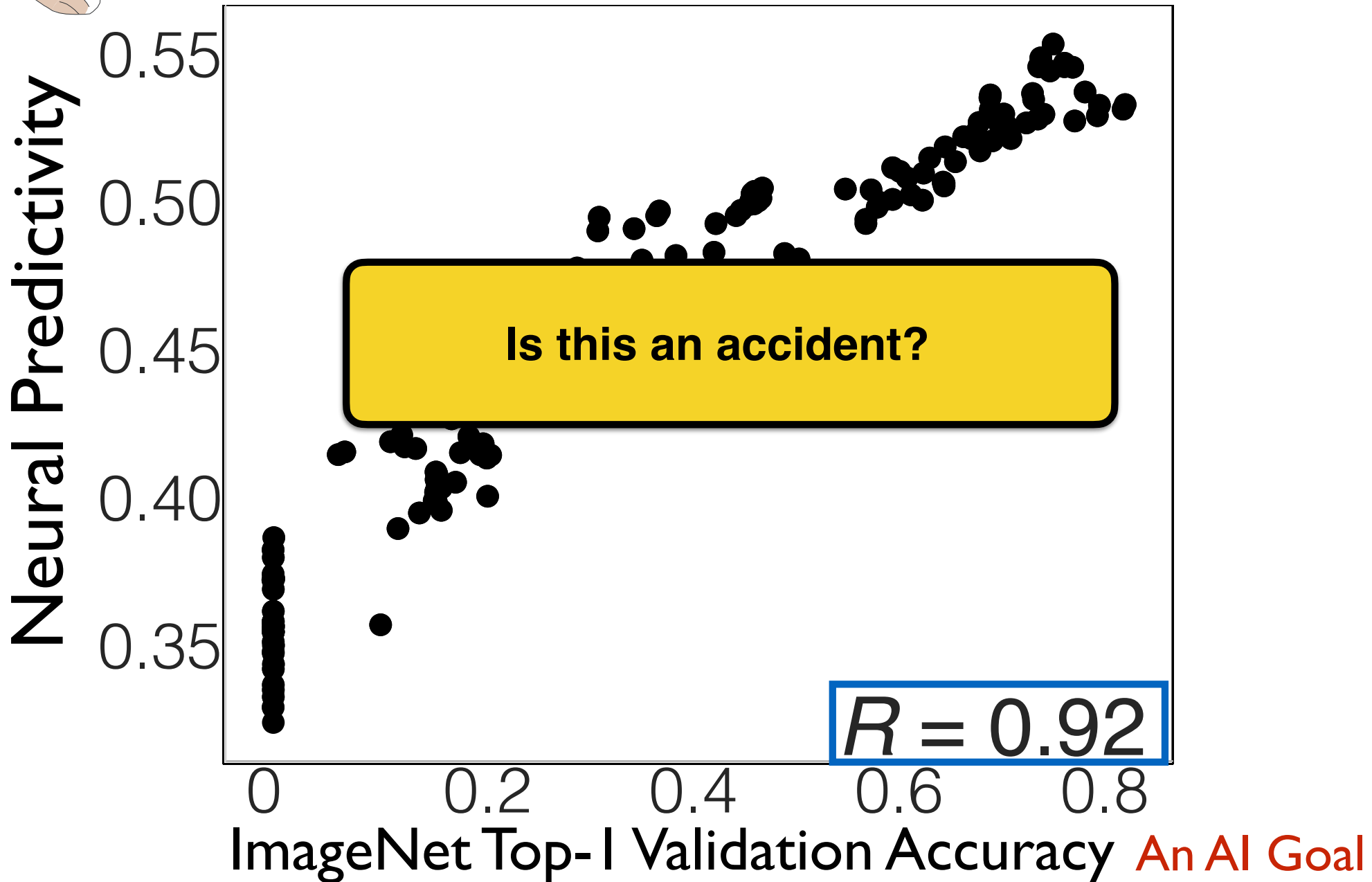


# Task Performance Correlated with Neural Predictivity



## A Neuroscience Goal

Schrimpf\*, Kubilius\* et al. 2018



# Task-Optimized Modeling: Four Components

## Task-Optimization (ML)

1.

**A** = *architecture class*

2.

**T** = *task loss*

3.

**D** = *dataset*

4.

**L** = *learning rule*

# Task-Optimized Modeling: Four Components

## Task-Optimization (ML)

## Neurobiology

1.

**A** = *architecture class*

2.

**T** = *task loss*

3.

**D** = *dataset*

4.

**L** = *learning rule*

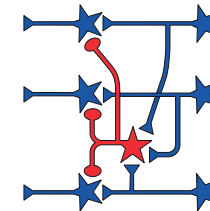
# Task-Optimized Modeling: Four Components

## Task-Optimization (ML)

1.

**A** = *architecture class* = **circuit neuroanatomy**

## Neurobiology



2.

**T** = *task loss*

3.

**D** = *dataset*

4.

**L** = *learning rule*

# Task-Optimized Modeling: Four Components

## Task-Optimization (ML)

1.

**A** = *architecture class* = **circuit neuroanatomy**

2.

**T** = *task loss* = **ecological niche/behavior**

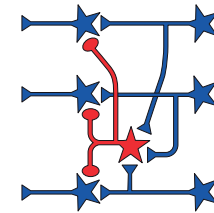
3.

**D** = *dataset*

4.

**L** = *learning rule*

## Neurobiology



# Task-Optimized Modeling: Four Components

## Task-Optimization (ML)

1.

**A** = *architecture class* = **circuit neuroanatomy**

2.

**T** = *task loss* = **ecological niche/behavior**

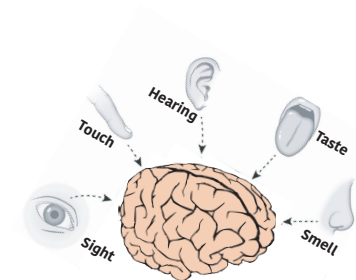
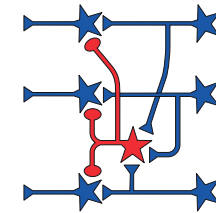
3.

**D** = *dataset* = **environment**

4.

**L** = *learning rule*

## Neurobiology



# Task-Optimized Modeling: Four Components

## Task-Optimization (ML)

1.

**A** = architecture class = **circuit neuroanatomy**

2.

**T** = task loss = **ecological niche/behavior**

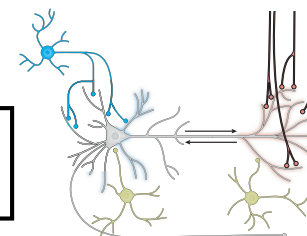
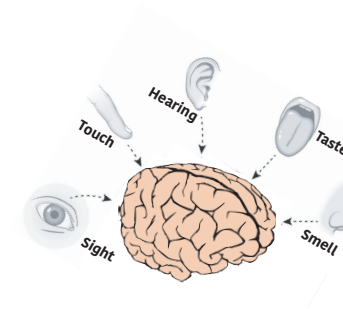
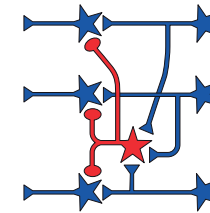
3.

**D** = dataset = **environment**

4.

**L** = learning rule = **natural selection + synaptic plasticity**

## Neurobiology



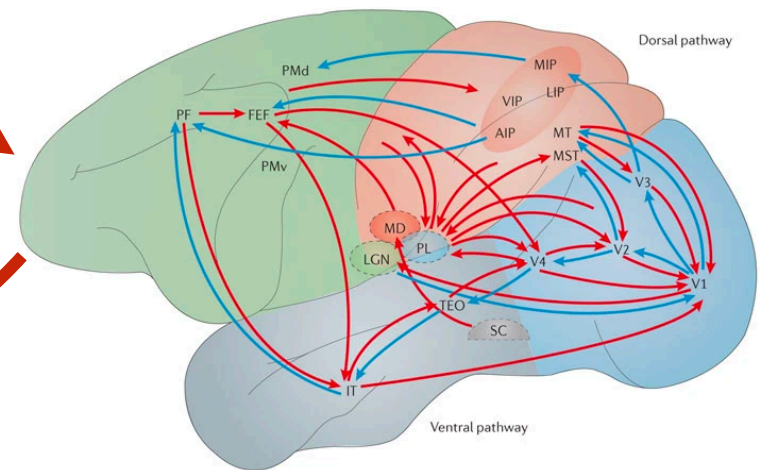
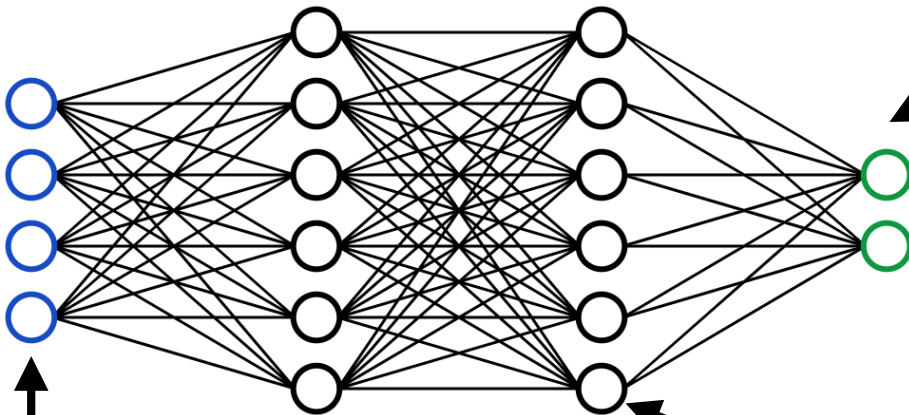
# Task-Optimized Modeling: Four Components

$L = \text{learning rule}$

$T = \text{task loss}$

**“Natural selection + plasticity”**

**“Ecological niche/behavior”**



**“Environment”**

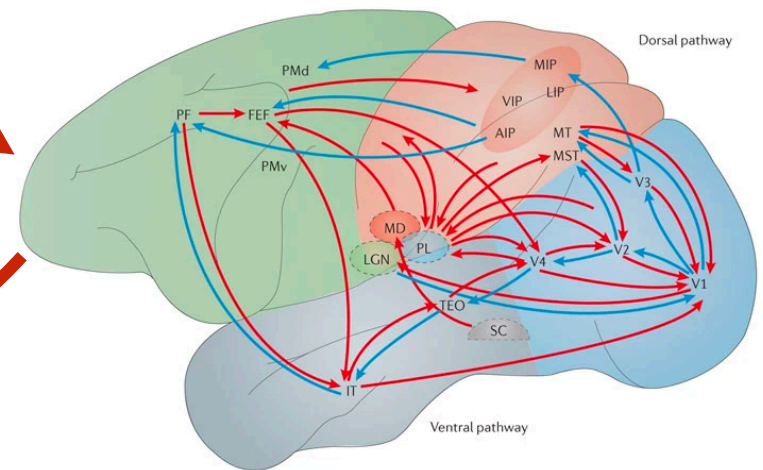
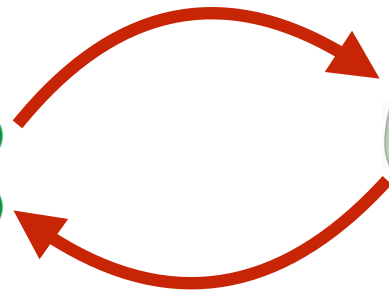
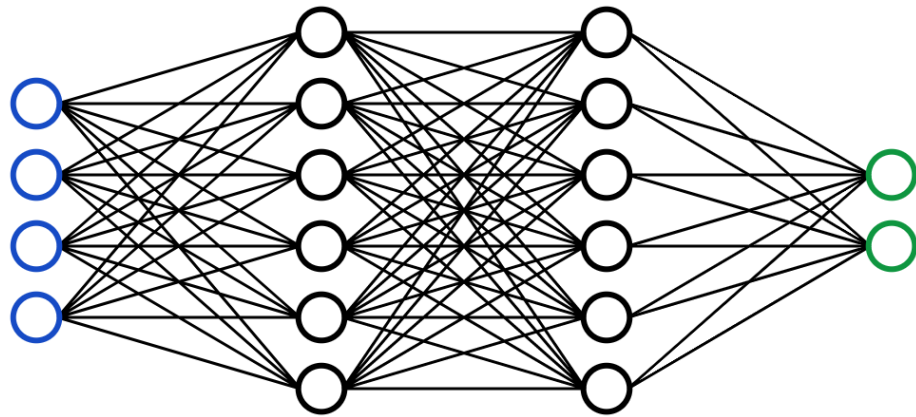
**“Circuit”**

$D = \text{data stream}$

$A = \text{architecture class}$

# Task-Optimized Modeling

Design ML Algorithms Optimized to Perform Organism's Behavior under Organism's Constraints



**Yields:**

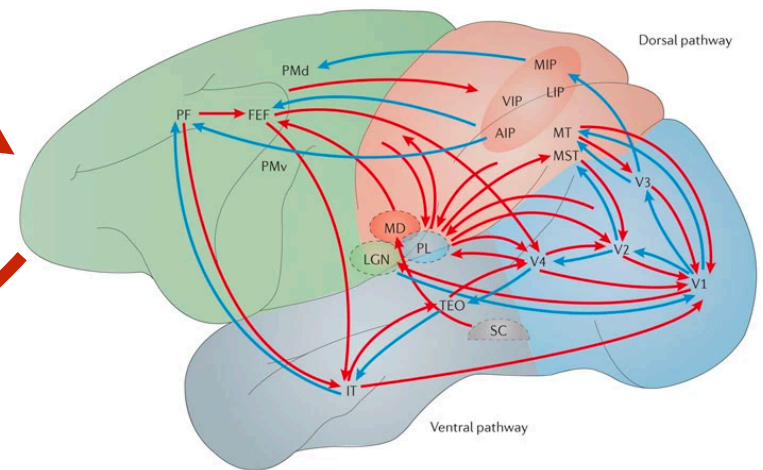
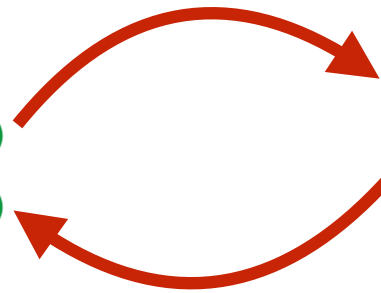
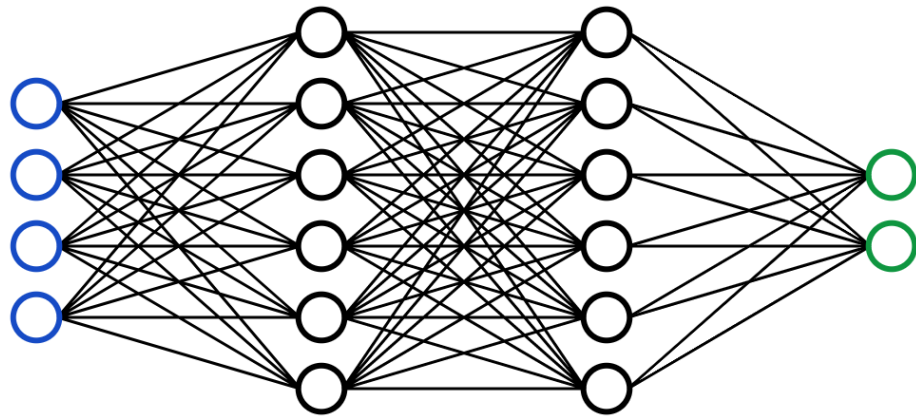
Quantitatively Accurate & Practically Useful Brain Models

**AND**

Principles of *Why* Neural Responses Are As They Are

# Task-Optimized Modeling

Design ML Algorithms Optimized to Perform Organism's Behavior under Organism's Constraints



**Yields:**

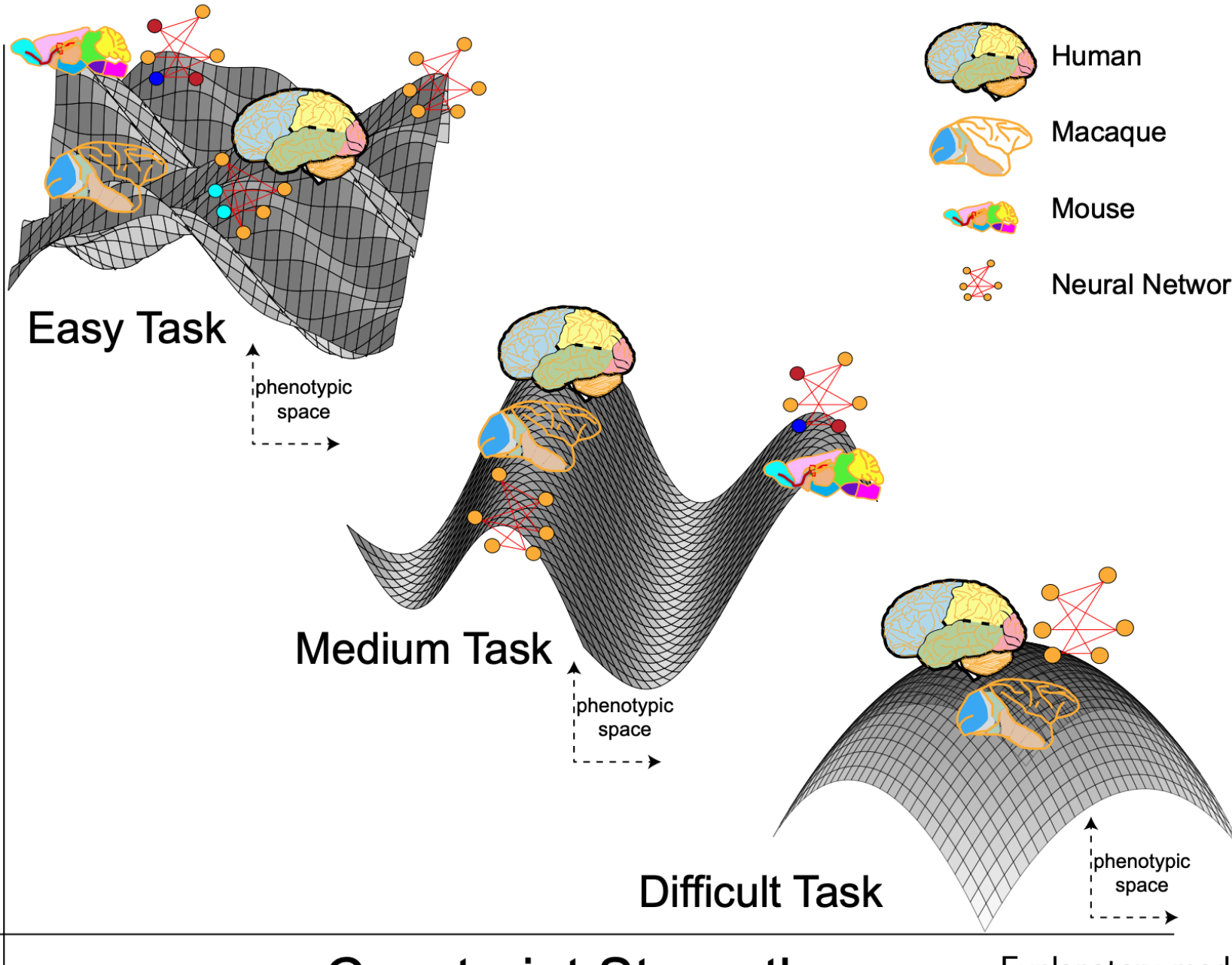
Quantitatively Accurate & Practically Useful Brain Models

**AND**

Principles of *Why* Neural Responses Are As They Are

# Contravariance Principle: The Harder the Task, the Less Solutions!

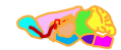
Dispersion of Solution Set



Human



Macaque



Mouse



Neural Network



Rosa Cao



Daniel Yamins

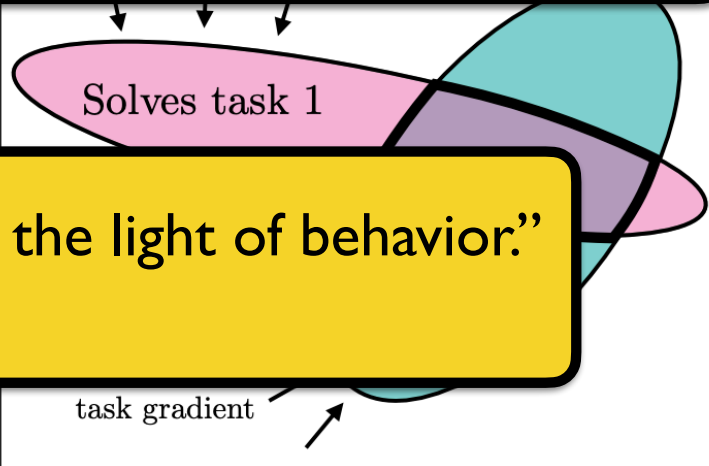
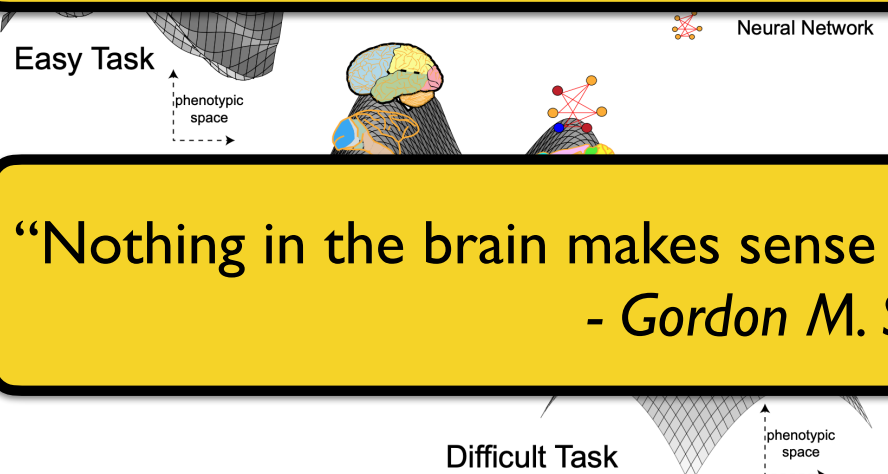
Constraint Strength

Explanatory models in neuroscience:  
Part 2 – Constraint-based intelligibility

# Platonic Representation Hypothesis is the AI version of Contravariance

“Nothing in biology makes sense except in light of evolution.”  
 - *Theo Dobzhansky*

Dispersion of Solution Set



“Nothing in the brain makes sense except in the light of behavior.”  
 - *Gordon M. Shepherd*

Our (slightly) modified credo:  
 “Nothing in (computational) neuroscience makes sense except in light of task-optimization.”

Figure 6. The Multitask Scaling Hypothesis: Models trained to pressure to solve more tasks at once.

**The Platonic Representation Hypothesis**

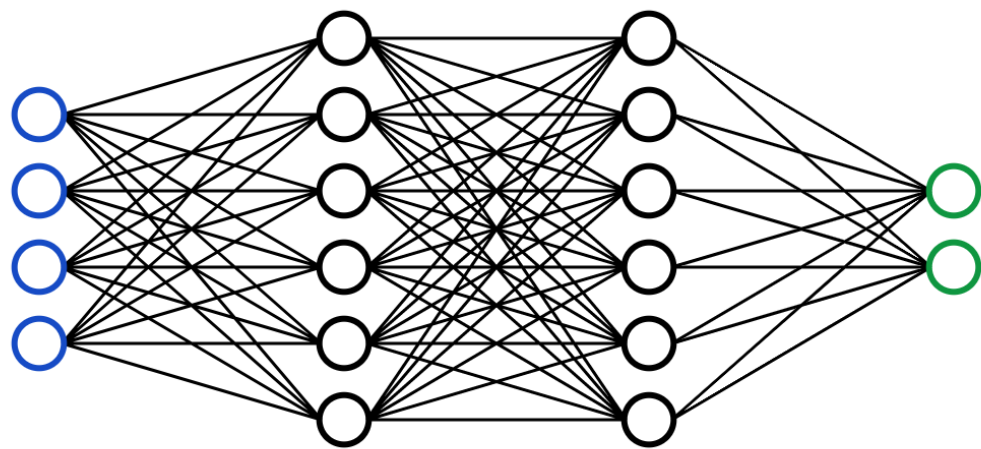
**The Multitask Scaling Hypothesis**

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.

# Task-Optimized Modeling Approach

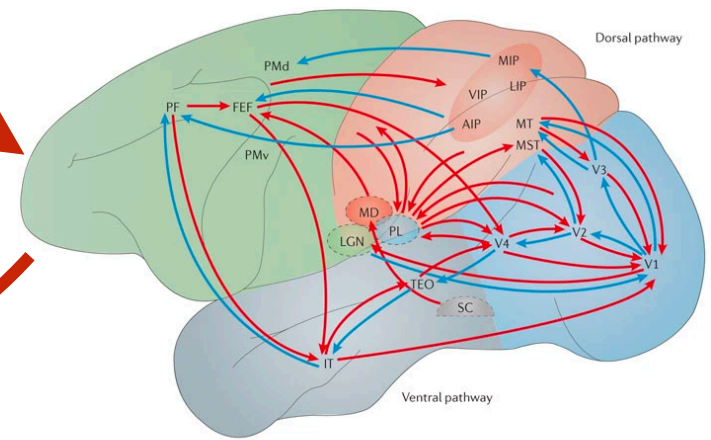
Design ML Algorithms Optimized to Perform Organism's Behavior under Organism's Constraints

But what even counts as *good* here?



Artificial Neural Network

**Yields:**



Brain

Quantitatively Accurate & Practically Useful Brain Models

**AND**

Principles of *Why* Neural Responses Are As They Are

# NeuroAI Turing Test

---

## Brain-Model Evaluations Need the NeuroAI Turing Test

---

**Jenelle Feather** <sup>\*1</sup> **Meenakshi Khosla** <sup>\*2</sup> **N. Apurva Ratan Murty** <sup>\*3</sup> **Aran Nayebi** <sup>\*4</sup>



Jenelle Feather



Meenakshi Khosla



Ratan Murty

# How to Reverse-Engineer Natural Intelligence?

**Whole brain...**

**Q: How are we going to make sense of all this data?**

**A: Build embodied agents & check if their internals pass the NeuroAI Turing test on *whole-brain* data.**

**... awake, behaving animals**

Using Agents to Reverse-Engineer *Whole-Brain* Data

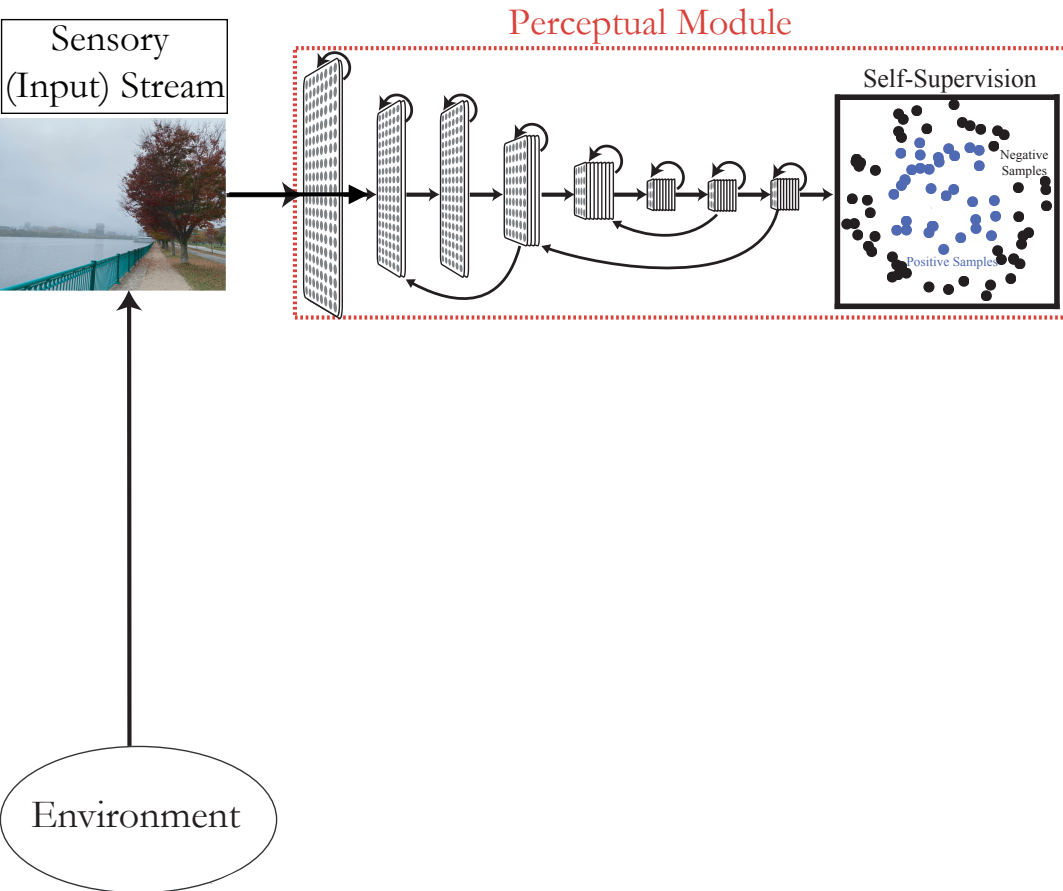
**How does the brain build and use **world models**?**

Using Agents to Reverse-Engineer *Whole-Brain* Data

**How does the brain *represent*, *predict*, *plan*, and enable *action*?**

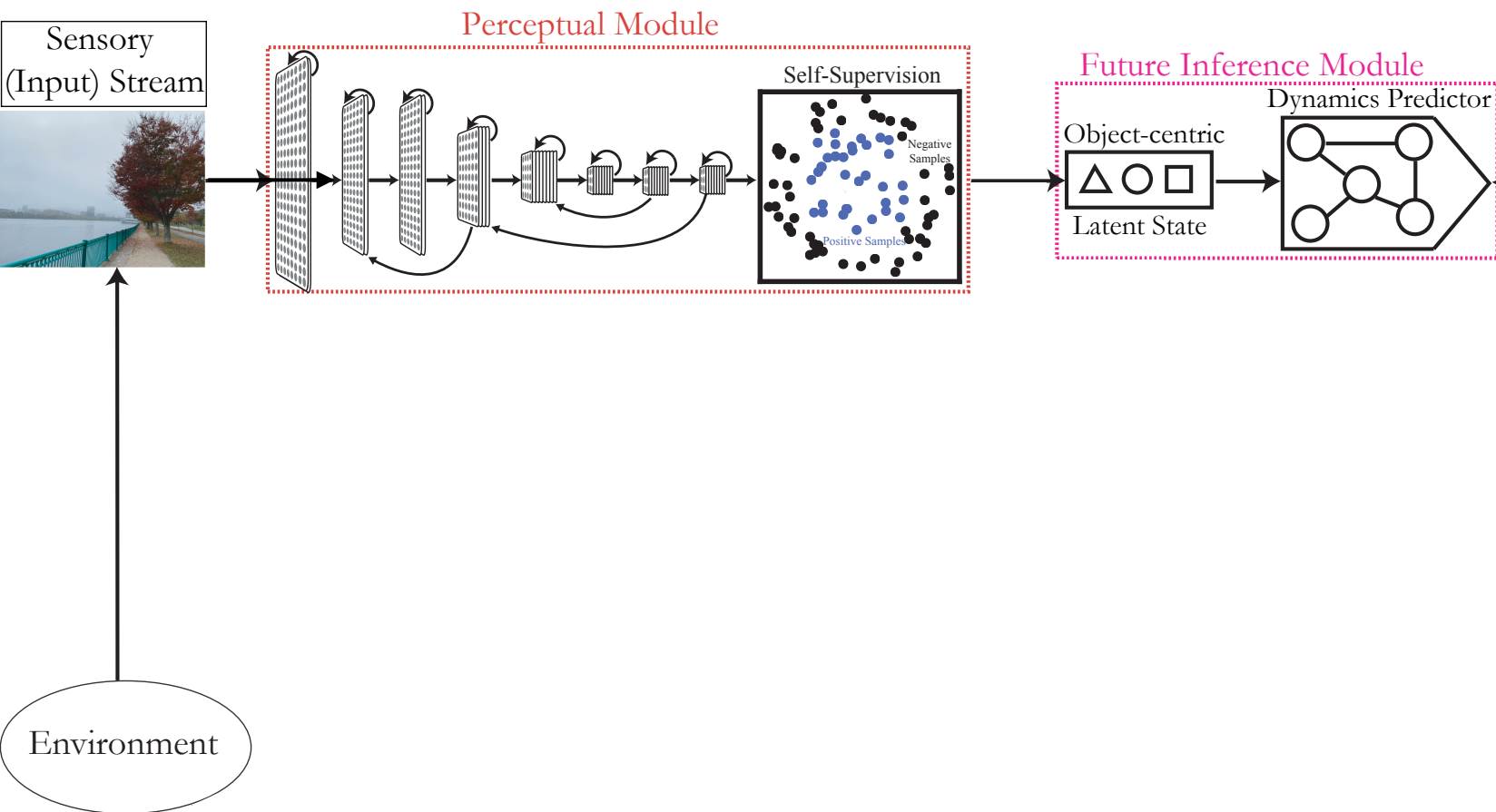
# Using Agents to Reverse-Engineer *Whole-Brain* Data

How does the brain *represent*, *predict*, *plan*, and enable *action*?



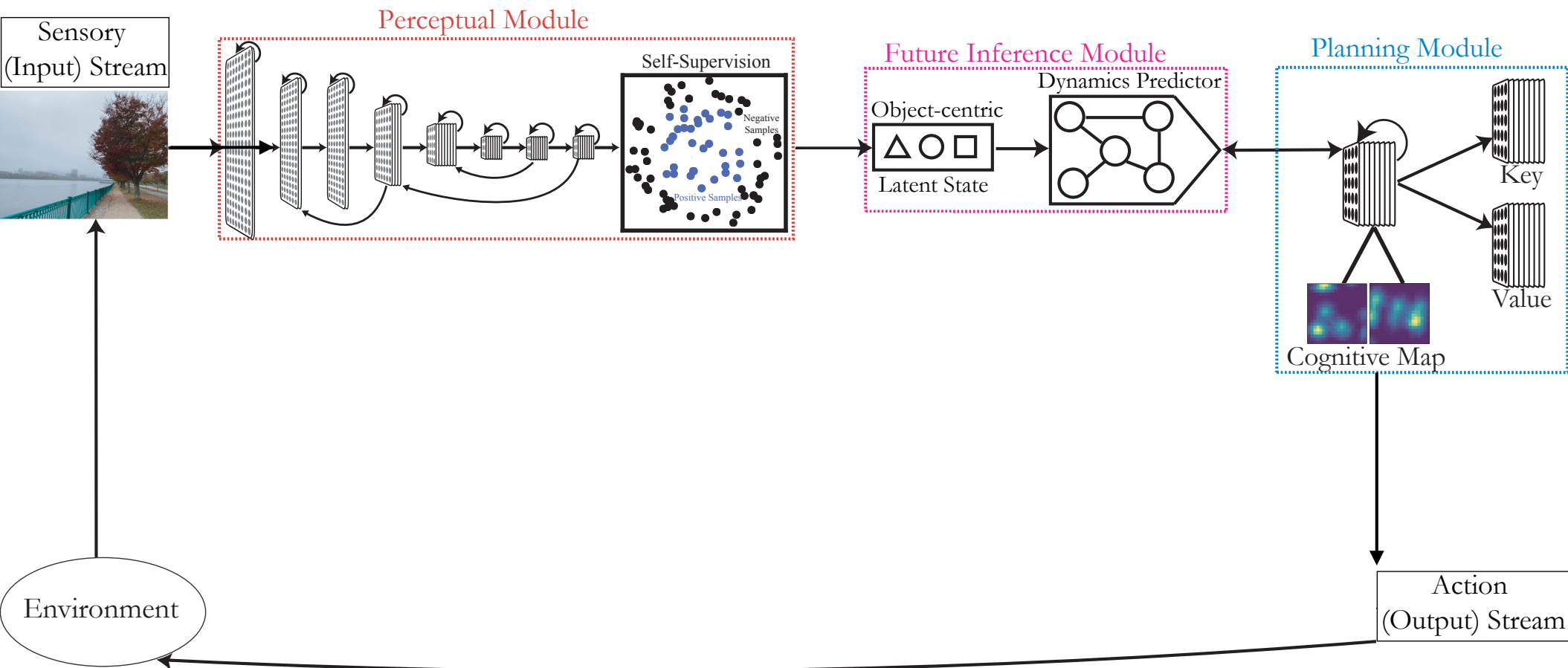
# Using Agents to Reverse-Engineer *Whole-Brain* Data

How does the brain *represent*, *predict*, *plan*, and enable *action*?



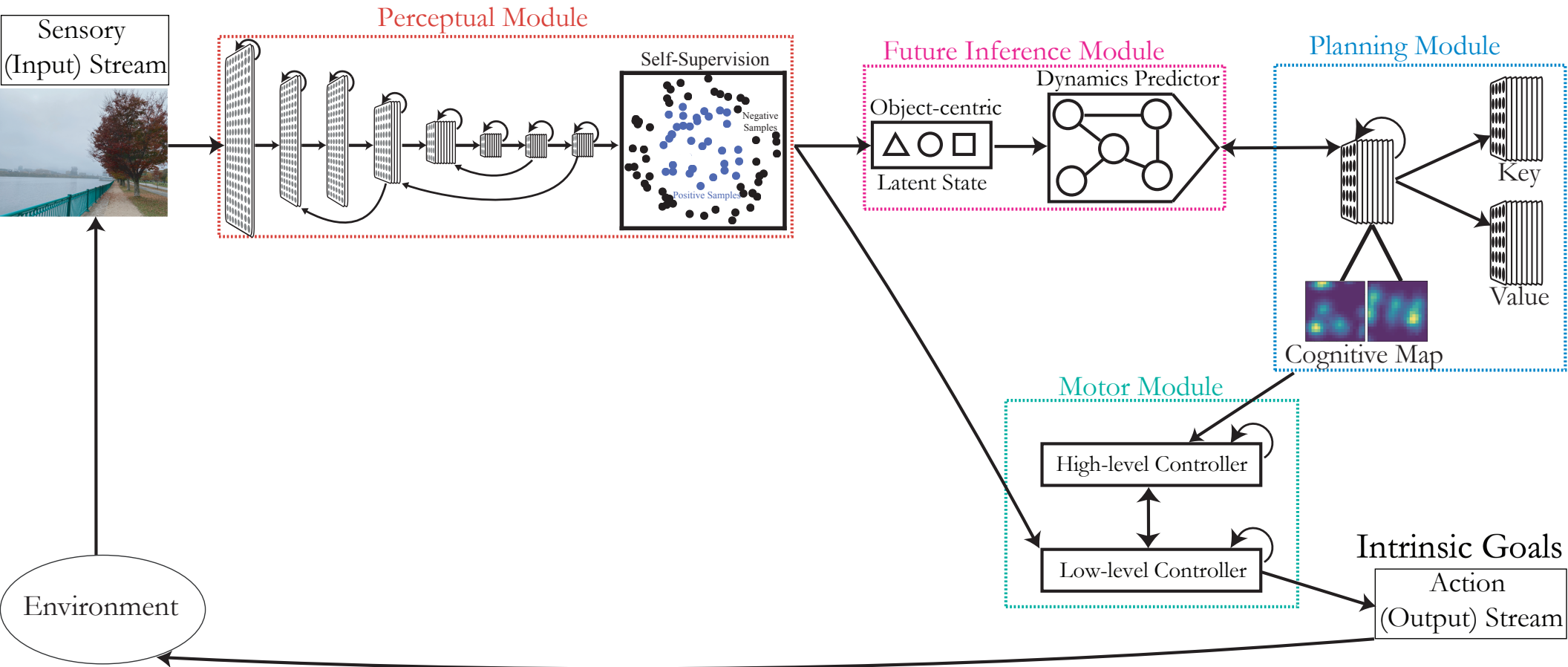
# Using Agents to Reverse-Engineer Whole-Brain Data

How does the brain *represent*, *predict*, *plan*, and enable *action*?



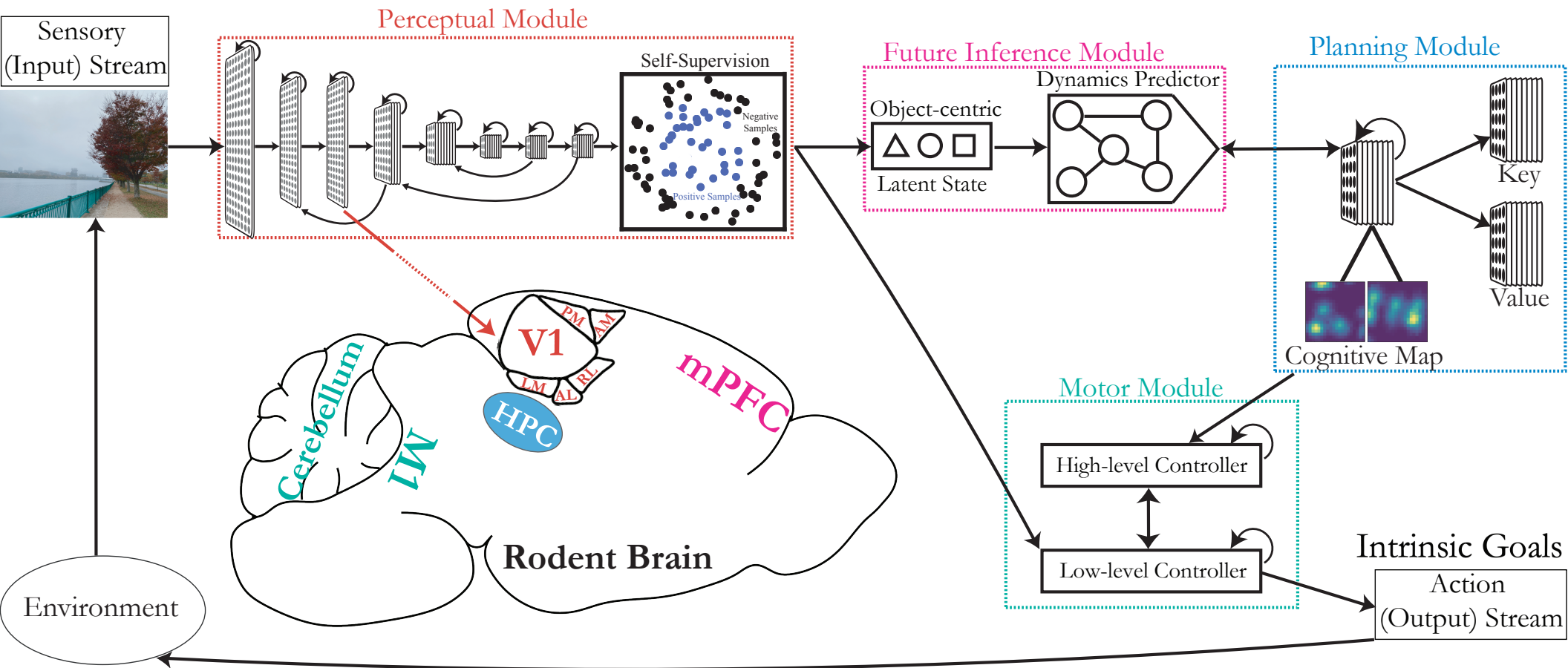
# Using Agents to Reverse-Engineer Whole-Brain Data

How does the brain *represent*, *predict*, *plan*, and enable *action*?



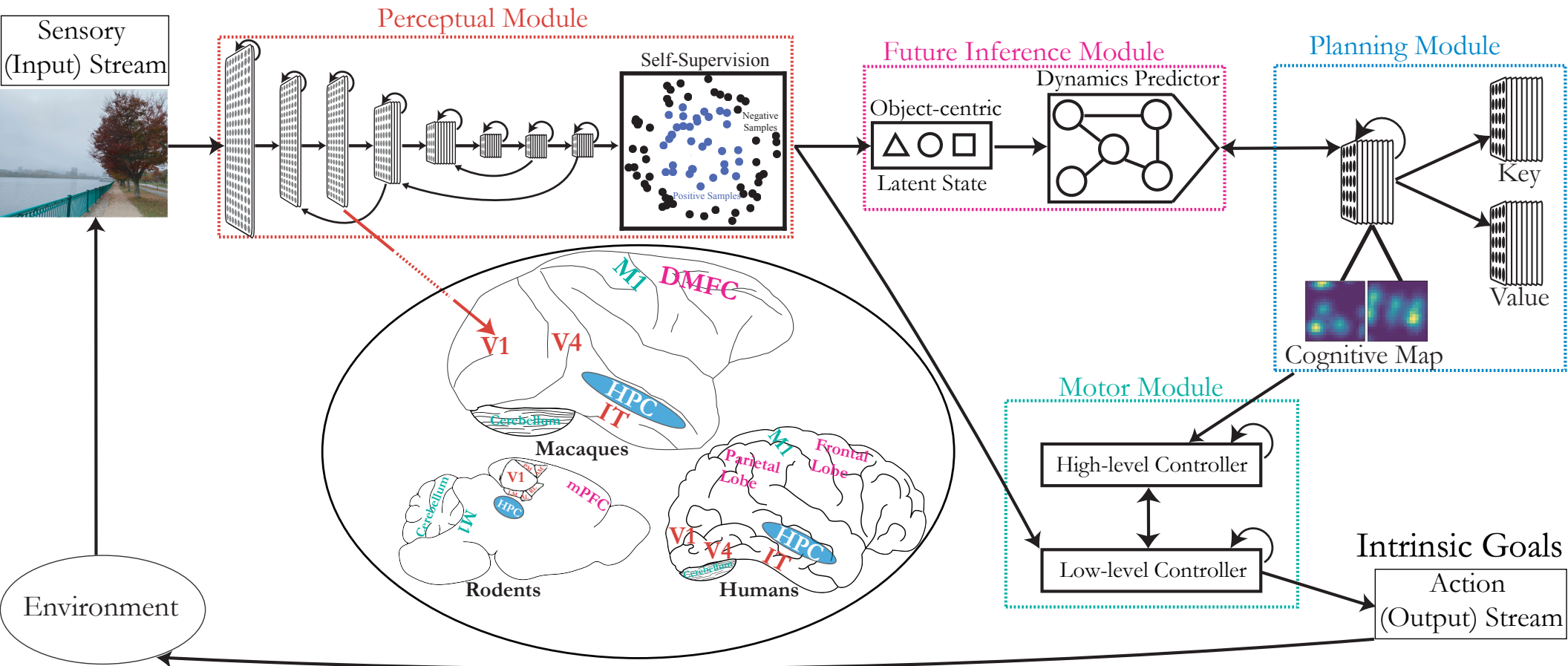
# Using Agents to Reverse-Engineer Whole-Brain Data

How does the brain *represent*, *predict*, *plan*, and enable *action*?



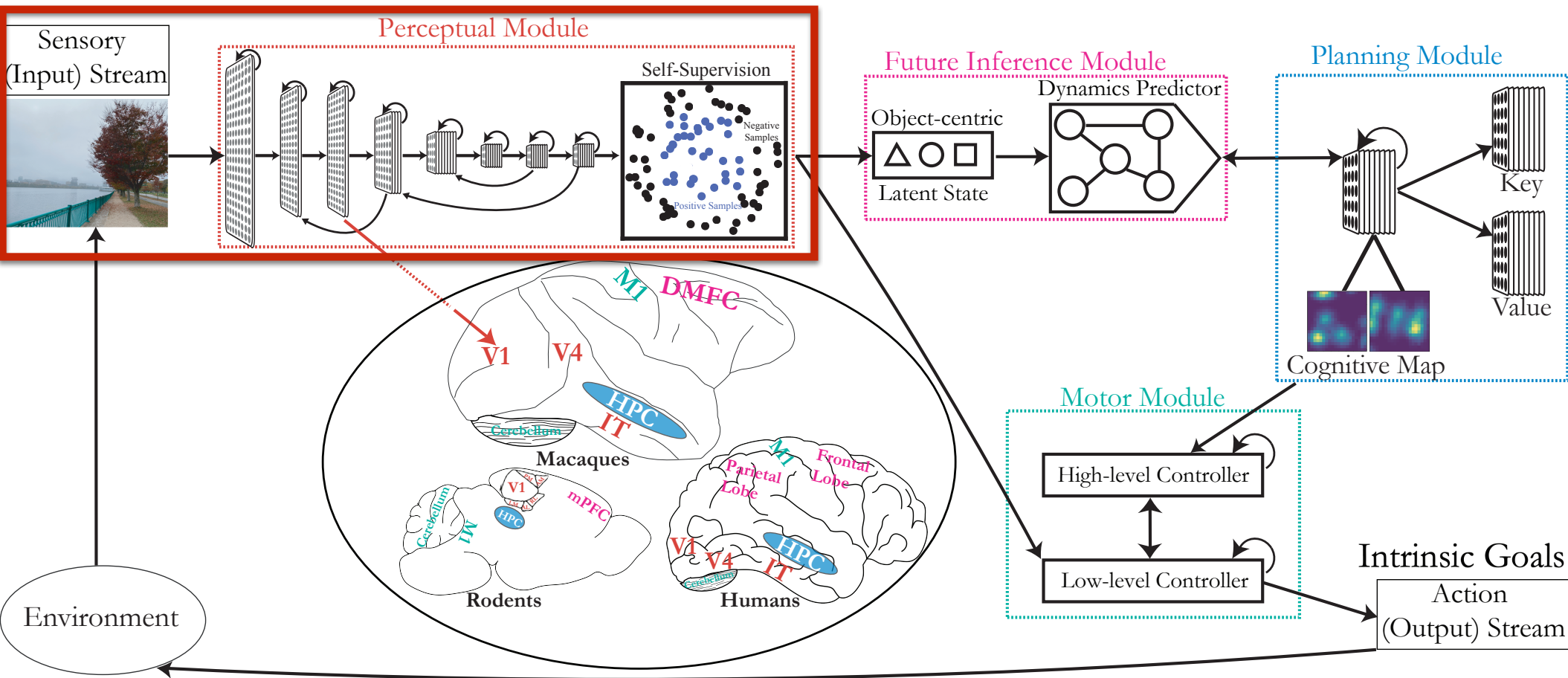
# Long-Term Outcome: Artificial **Organisms**

How does the brain *represent*, *predict*, *plan*, and enable *action*?



# Roadmap: Perception

How does the brain *represent*, *predict*, *plan*, and enable *action*?



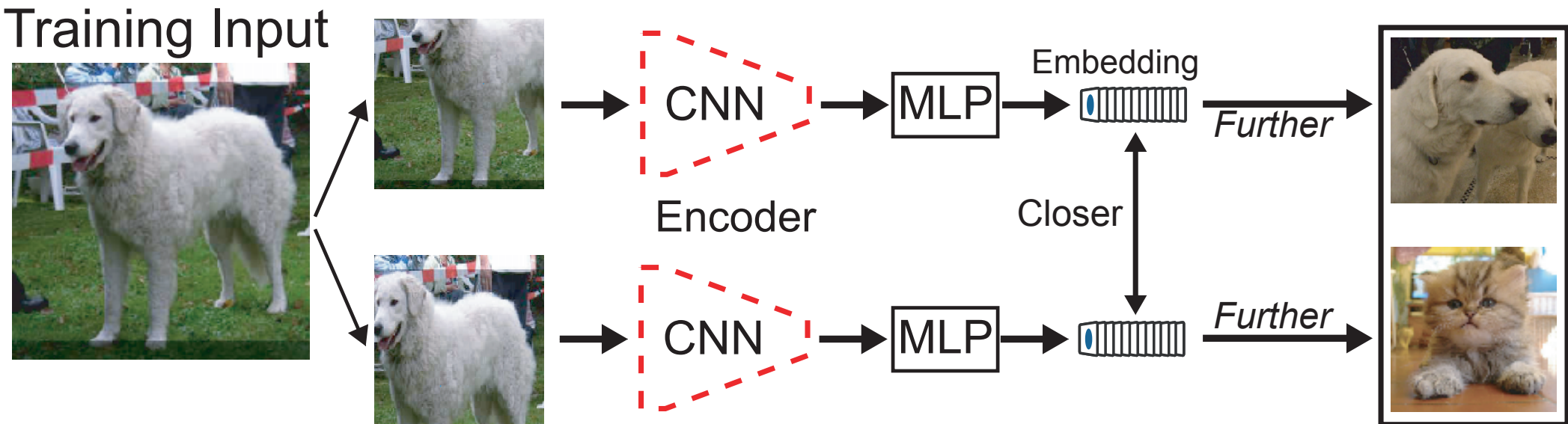
# The Supervision Problem



There's just no way that these creatures receive millions of high-level semantic labels during learning.

Effective proxy, but just obviously deeply wrong.

# Contrastive learning tasks



*CNN: Convolutional Neural Network, MLP: Multi-Layer Perceptron*

**High-level idea of these methods: make the representations  
non-trivially robust to data augmentations**

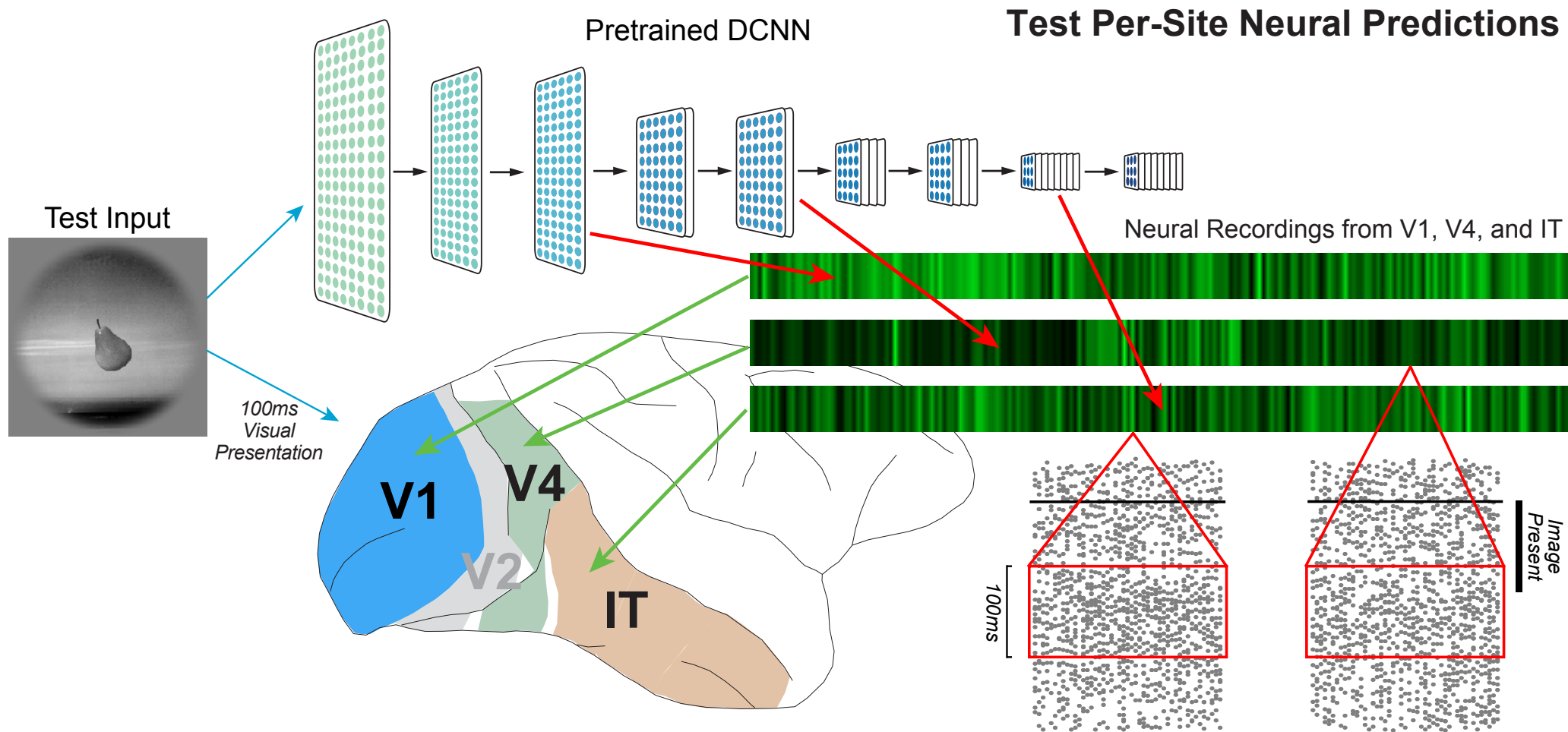
**(somewhat inspired by how we “sample” the world via head  
motion)**

# Comparison to Neural Data



Chengxu  
Zhuang

How well does it match neural data?

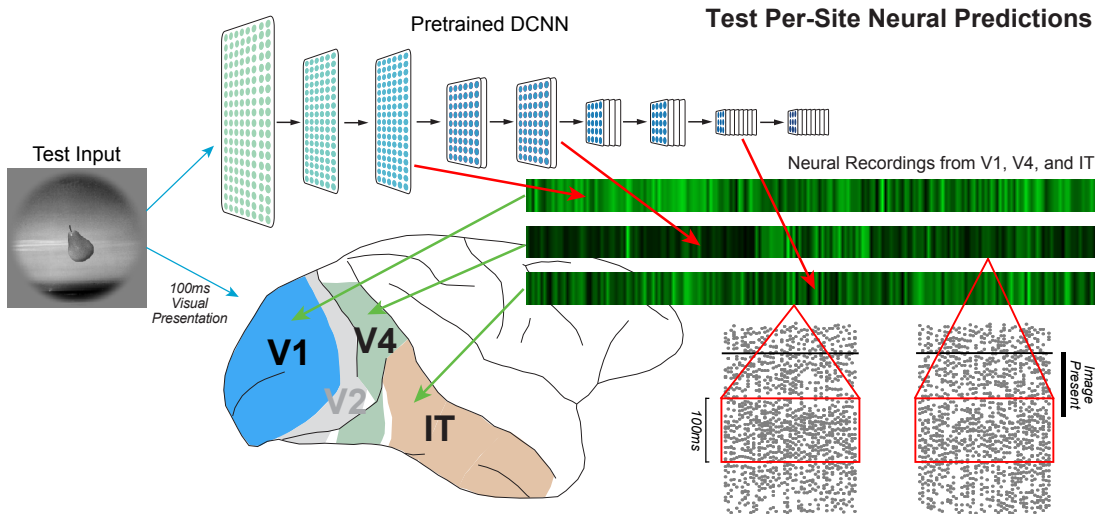


V1 data from Cadena et al. [Deep convolutional models improve predictions of macaque V1 responses to natural images](#) *PLoS Comp. Bio.*, (2019)

V4 & IT data from Majaj et al. [Simple Learned Weighted Sums of Inferior Temporal Neuronal Firing Rates Accurately Predict Human Core Object Recognition Performance](#) *J. Neurosci.* (2015)



Chengxu  
Zhuang

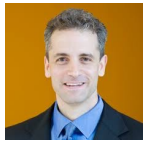
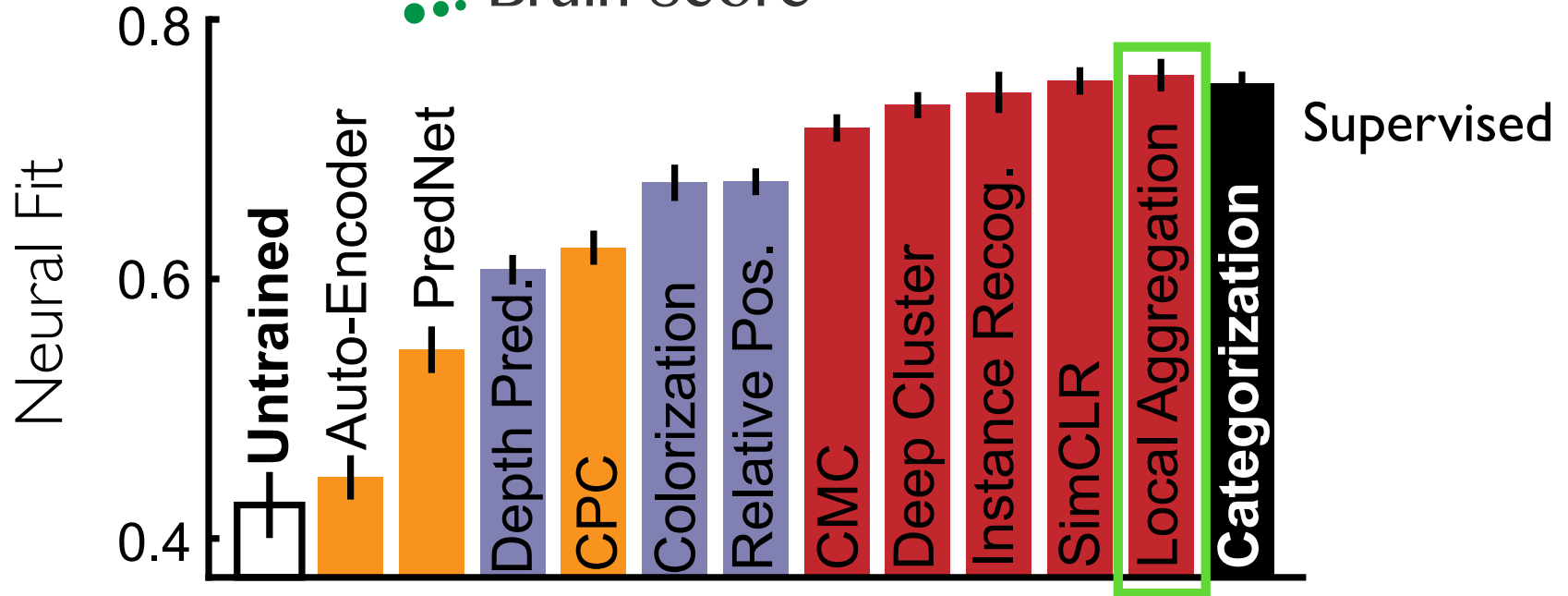


Autoencoders

Missing-Data Tasks

Deep Contrastive Embeddings

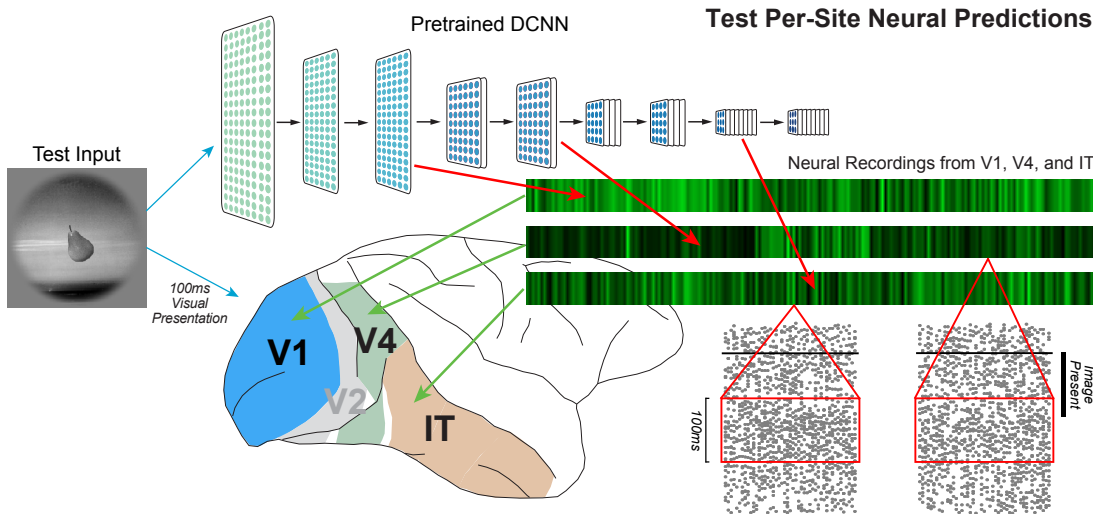
 Brain-Score



**Quantitatively accurate self-supervised model**  
**of a higher brain area.**



Chengxu Zhuang



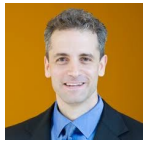
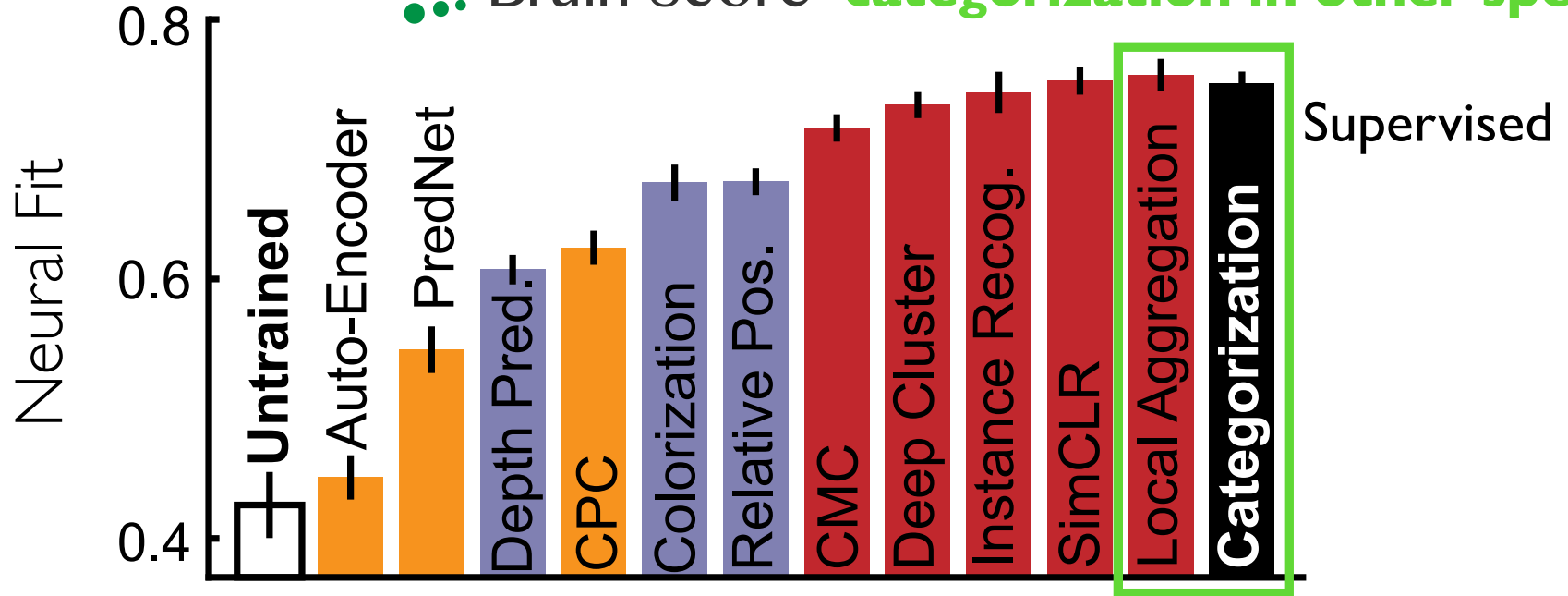
Autoencoders

Missing-Data Tasks

Deep Contrastive Embeddings

Can we do even better than categorization in other species?

Brain-Score



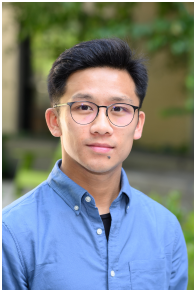
**Quantitatively accurate self-supervised model of a higher brain area.**

# Mouse Visual Cortex as a Task-General, Limited Resource System

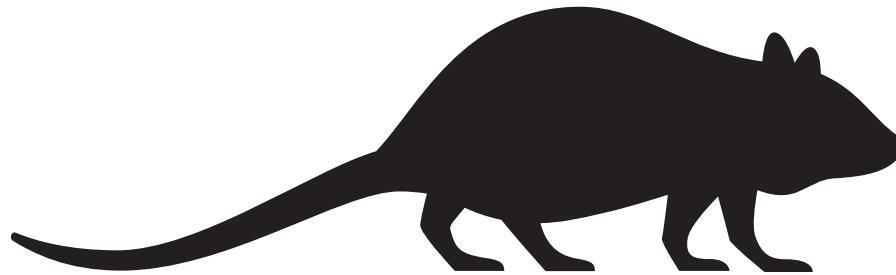
A. Nayebi\*, N.C.L. Kong\*, C. Zhuang, J.L. Gardner, A.M. Norcia, D.L.K. Yamins

Mouse visual cortex as a limited resource system that self-learns an ecologically-general representation.

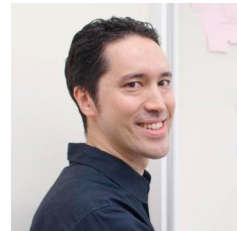
*PLOS Computational Biology* 2023



Nathan C.L. Kong\*



Chengxu Zhuang



Justin L. Gardner

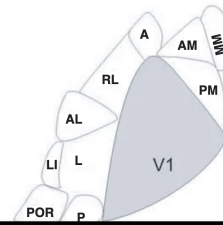
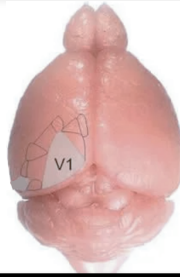
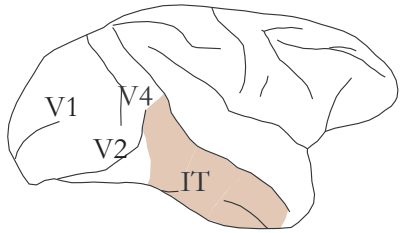


Anthony M. Norcia



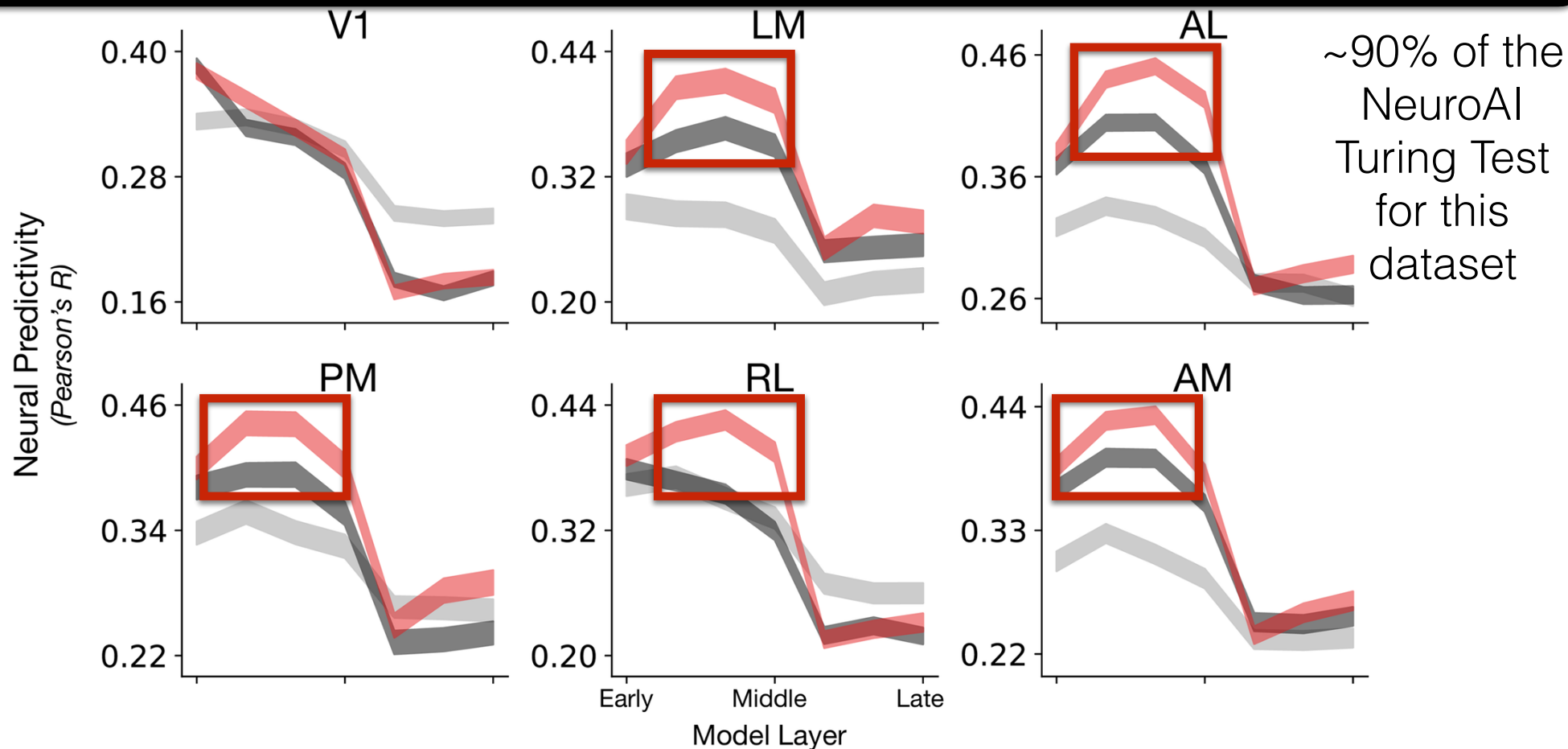
Daniel Yamins

# Contrastive Models Better Match Mouse Visual Cortex



Mouse vision is less hierarchical!

What is the ecological reason why the mouse visual system prefers *self-supervision*?  
Hypothesis: *task-generality* rather than functional specialization.



# Assessing Task-Generality

# Assessing Task-Generality

Train

*ImageNet*



# Assessing Task-Generality

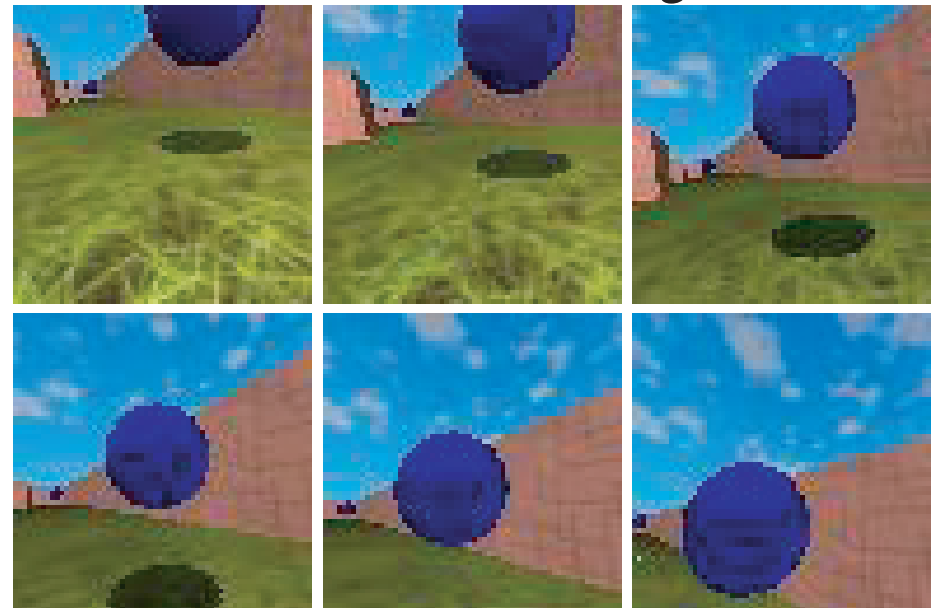
Train

*ImageNet*



Evaluate

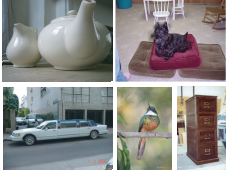
*Reward-Based Navigation*



# Assessing Task-Generality

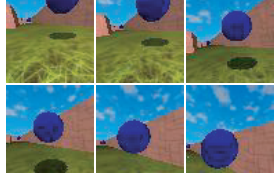
Train

*ImageNet*

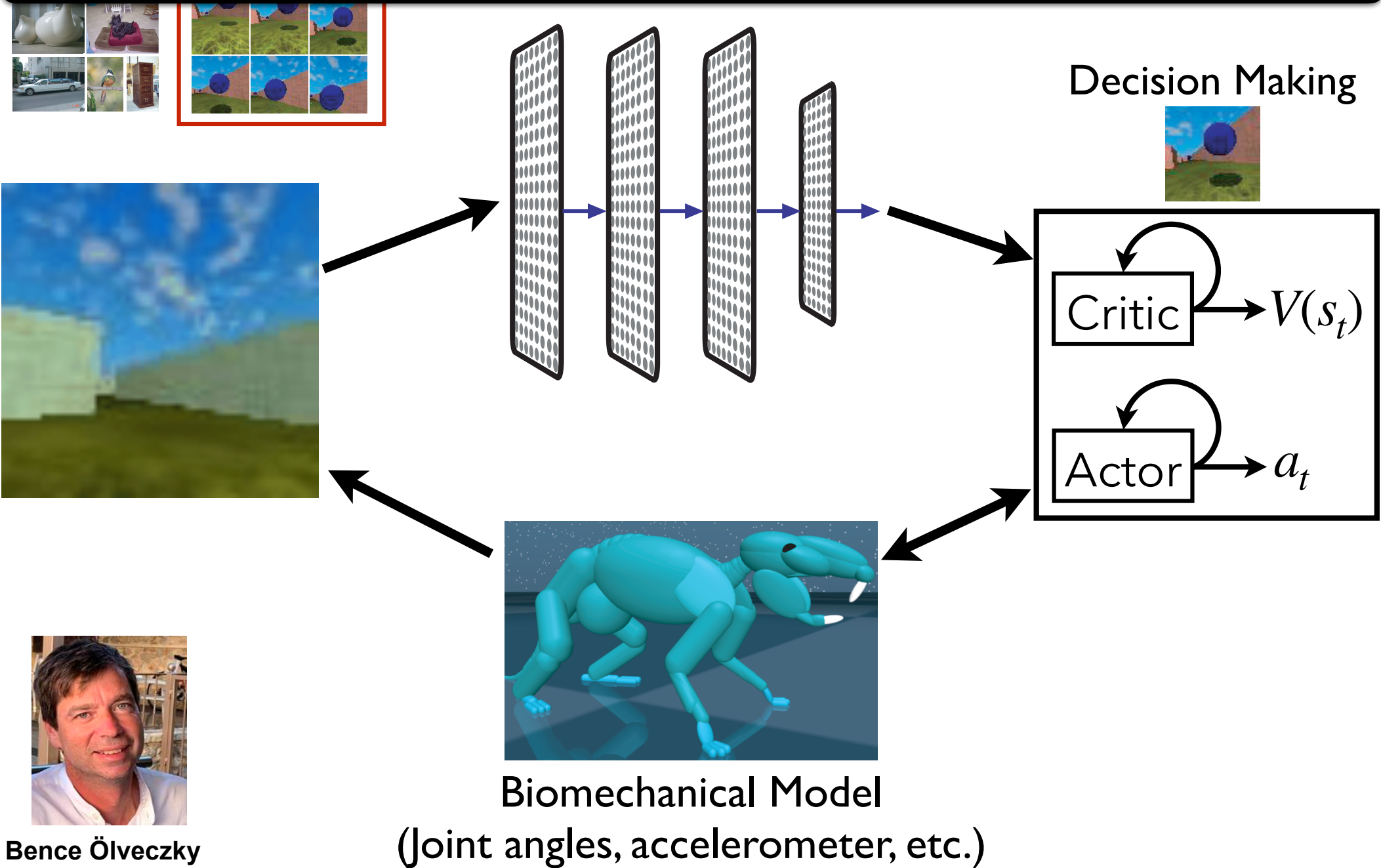


Evaluate

*Reward-Based Navigation*

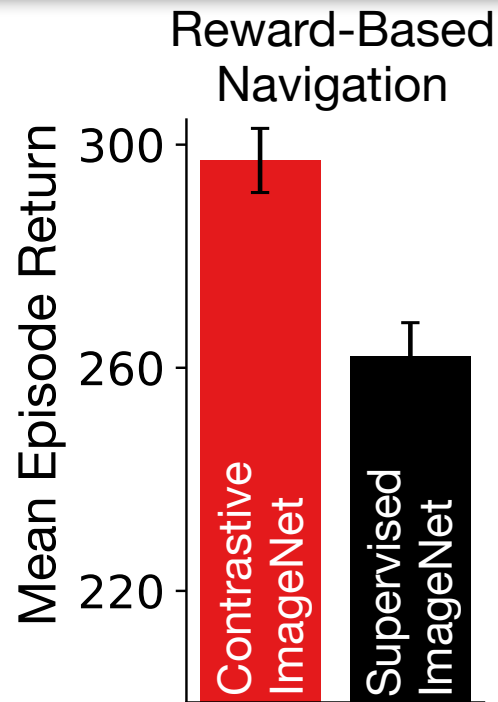


High degree-of-freedom body (38/74 controllable degrees), keeping track of history over long timescales with high-dimensional, continuous inputs

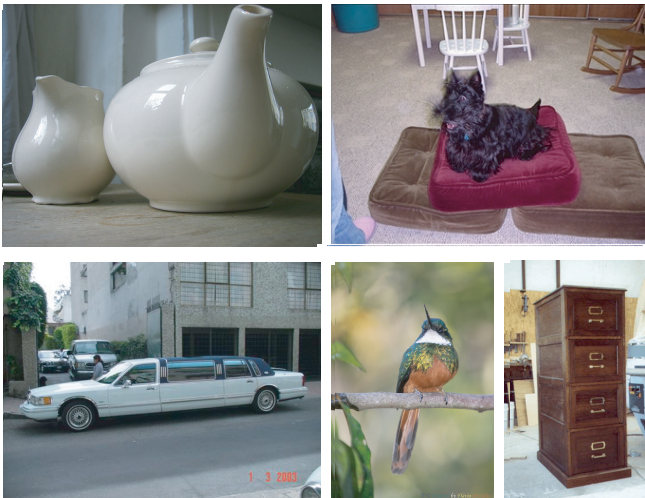


# Contrastive Models Yield Better Transfer Performance

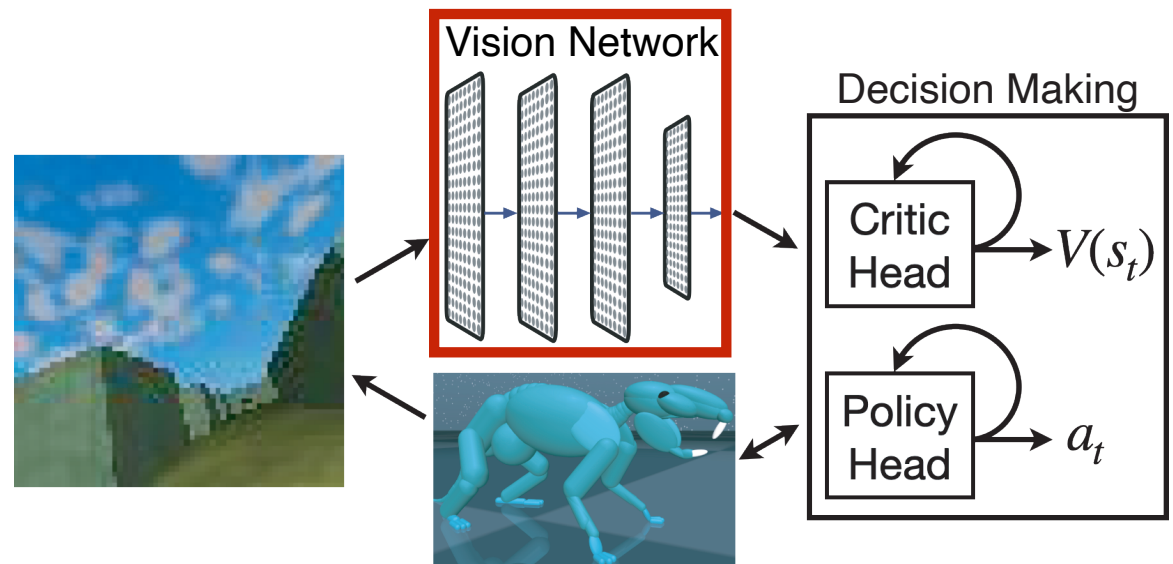
# Contrastive Models Yield Better Transfer Performance



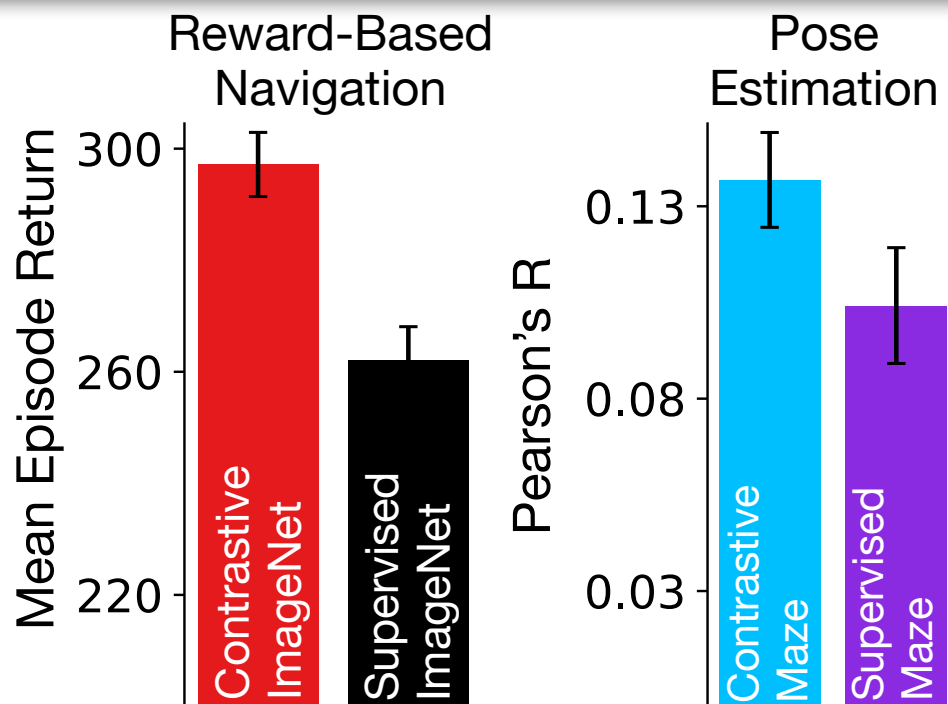
## Train ImageNet



## Evaluate

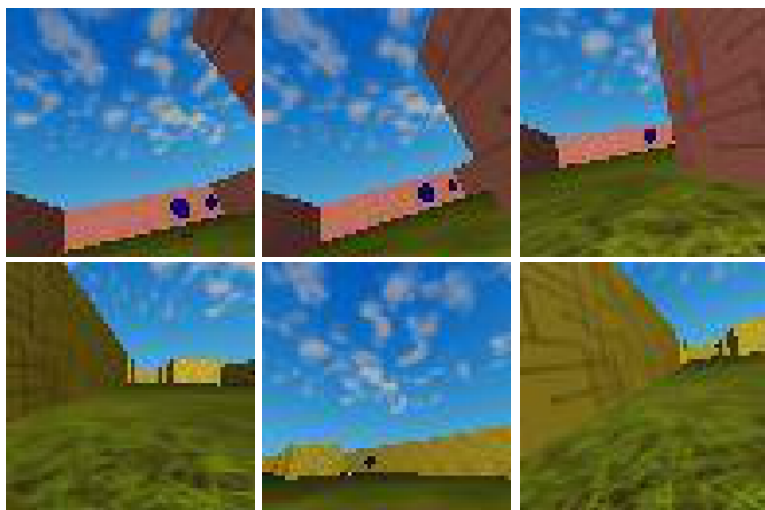


# Contrastive Models Yield Better Transfer Performance



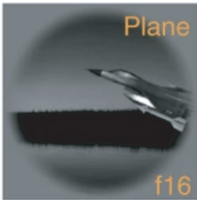
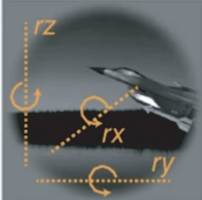

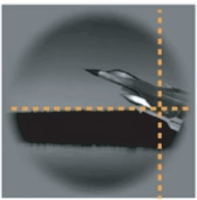
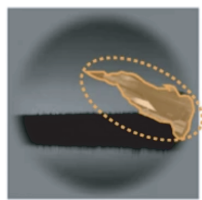
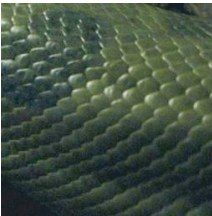
## Train

### *Maze Environment*

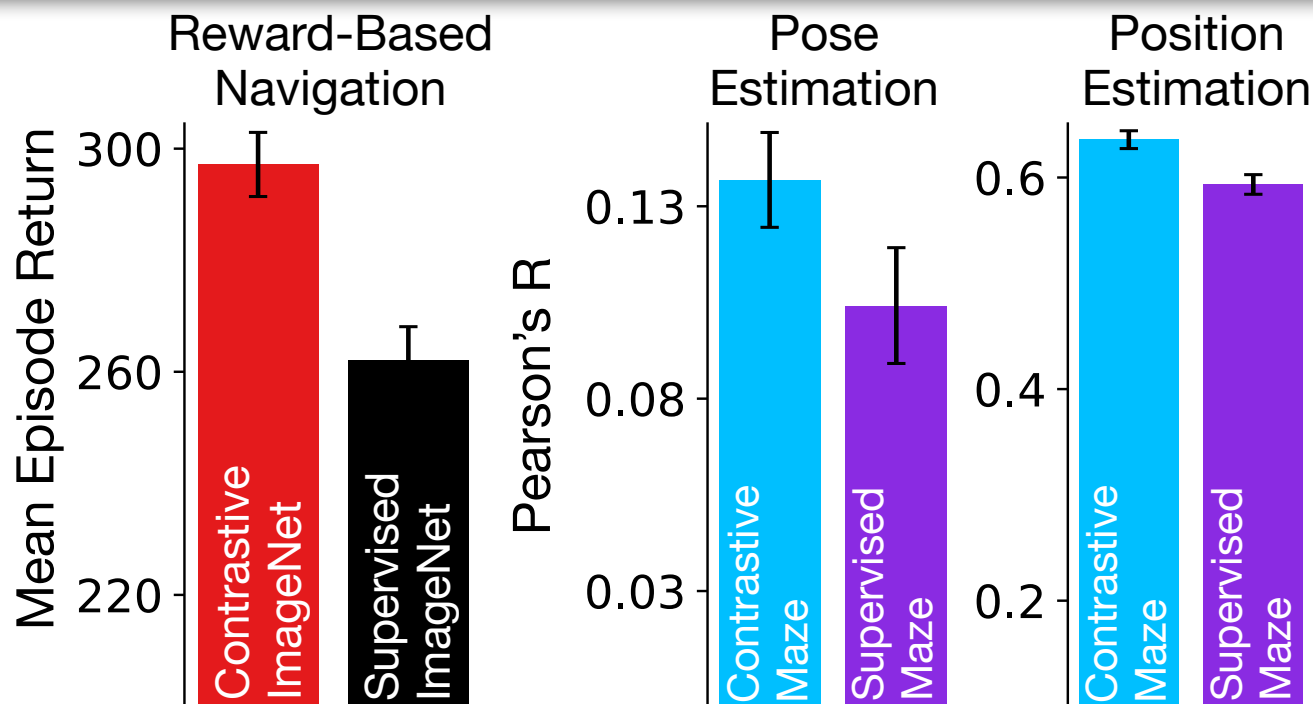


## Evaluate

### *Visual Scene Understanding*

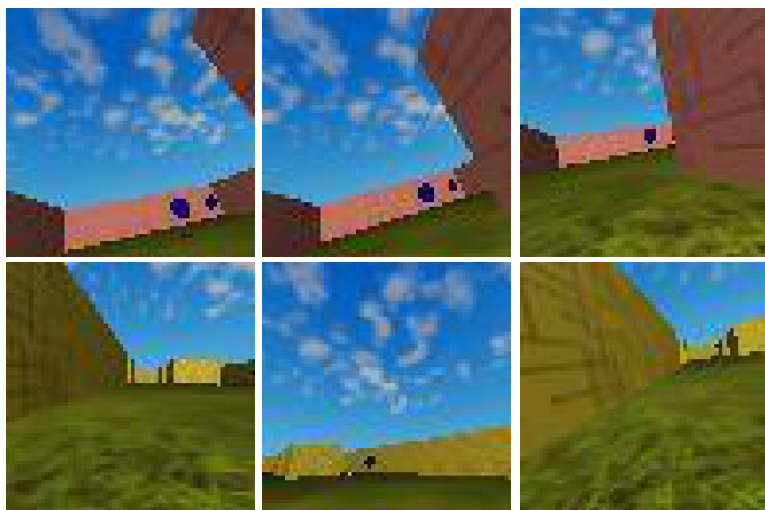
 <p>Plane</p>	Category	 <p><math>r_z</math> <math>r_x</math> <math>r_y</math></p>	z axis rotation x axis rotation y axis rotation	
 <p>f16</p>	Identity	 <p>Perimeter: 78 pix Two-dimensional retinal area: 146 pix Three-dimensional object scale: 1.2x</p>		
	Horizontal position: 80 pix Vertical position: -6 pix	<i>Object properties</i>		<i>Texture</i>

# Contrastive Models Yield Better Transfer Performance



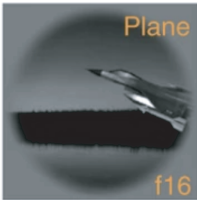
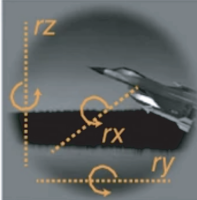

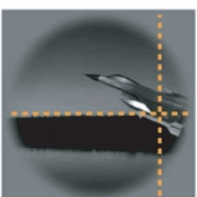
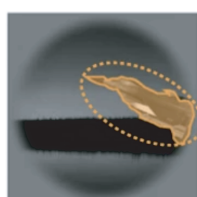
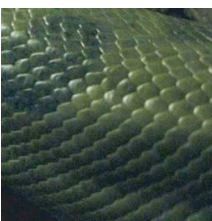
## Train

### *Maze Environment*

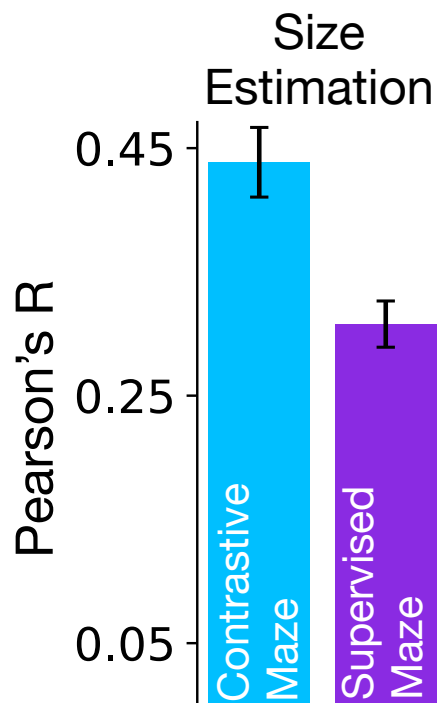
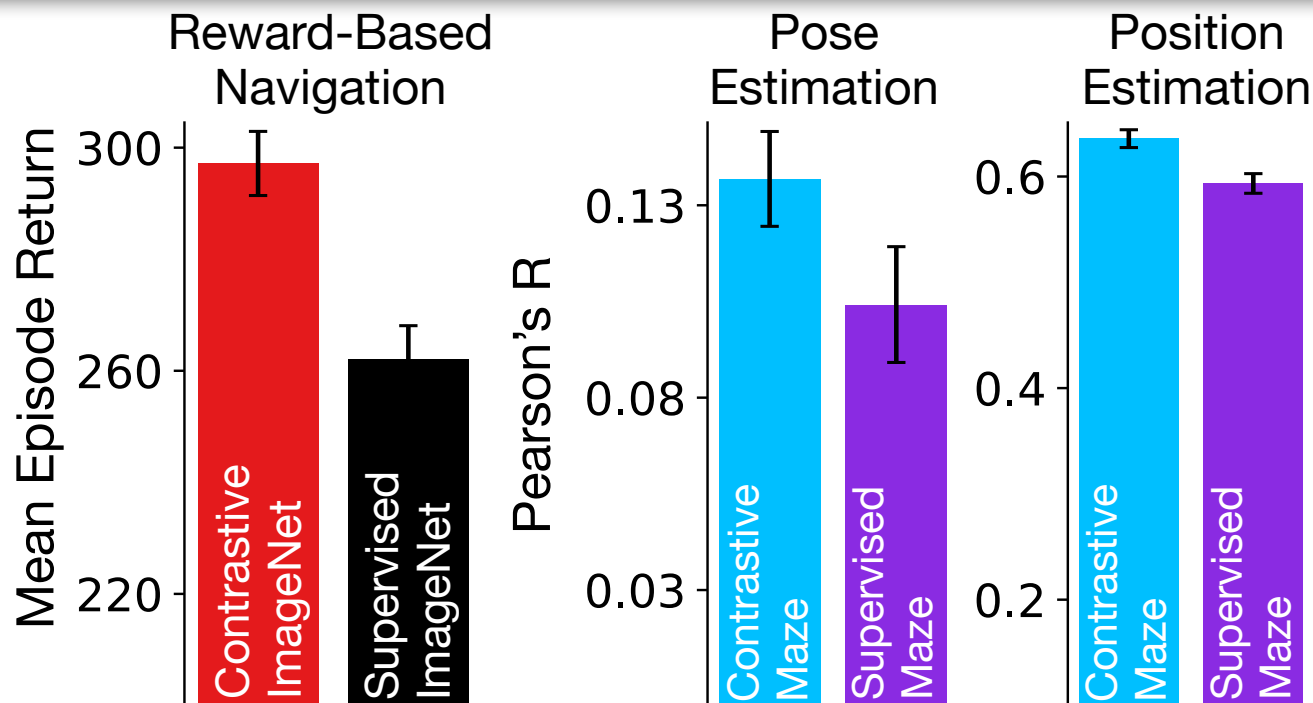


## Evaluate

### *Visual Scene Understanding*

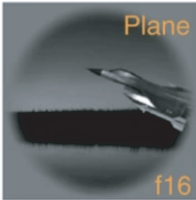
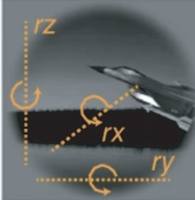

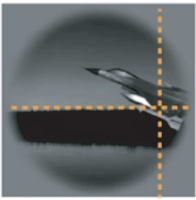
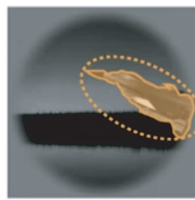
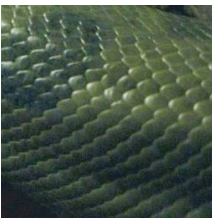
 <p>Plane</p>	Category	 <p><math>r_z</math> <math>r_x</math> <math>r_y</math></p>	z axis rotation x axis rotation y axis rotation	
 <p>f16</p>	Identity	 <p>Horizontal position: 80 pix Vertical position: -6 pix</p>	Perimeter: 78 pix Two-dimensional retinal area: 146 pix Three-dimensional object scale: 1.2x	
<i>Object properties</i>				<i>Texture</i>

# Contrastive Models Yield Better Transfer Performance

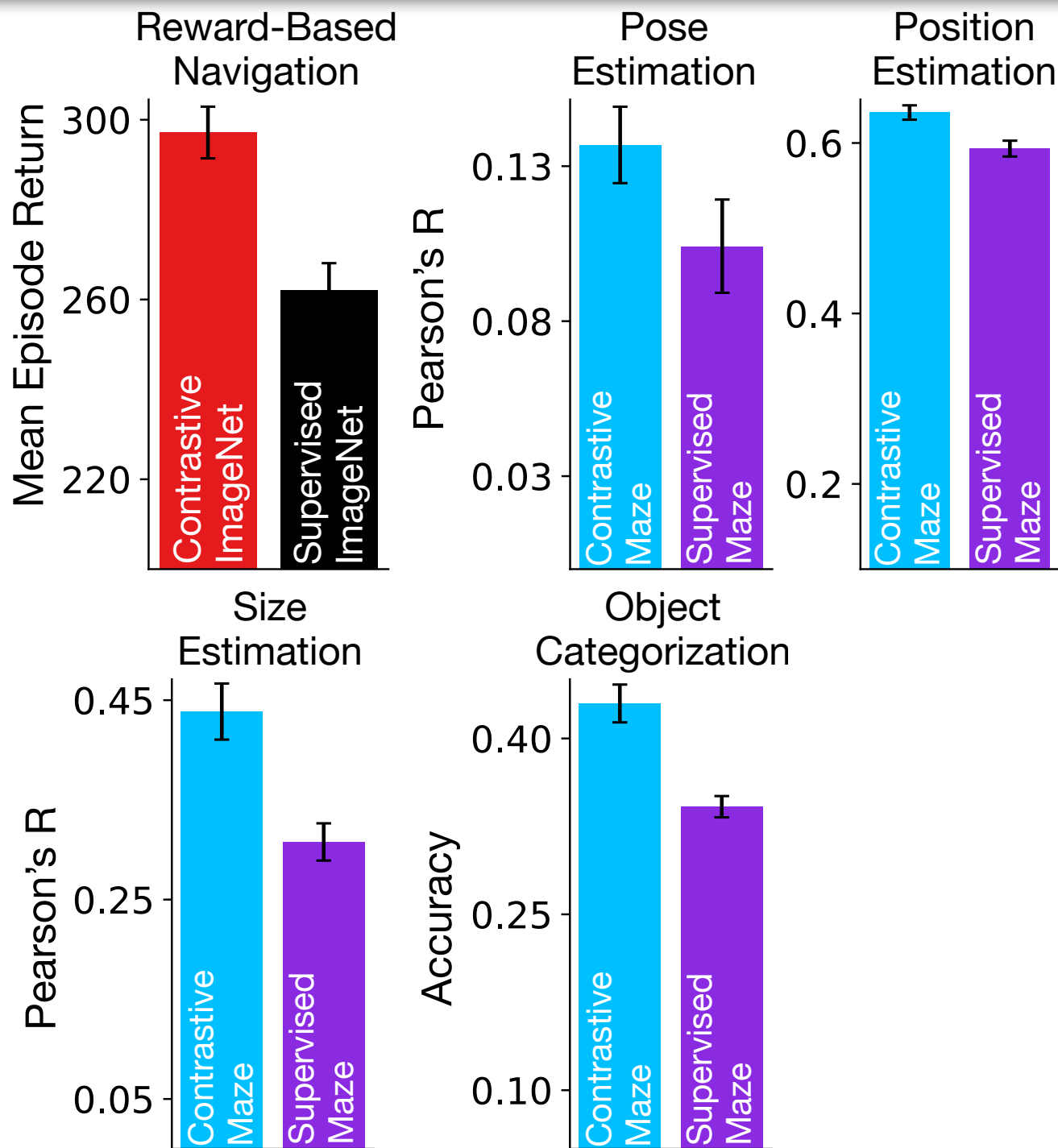


## Evaluate

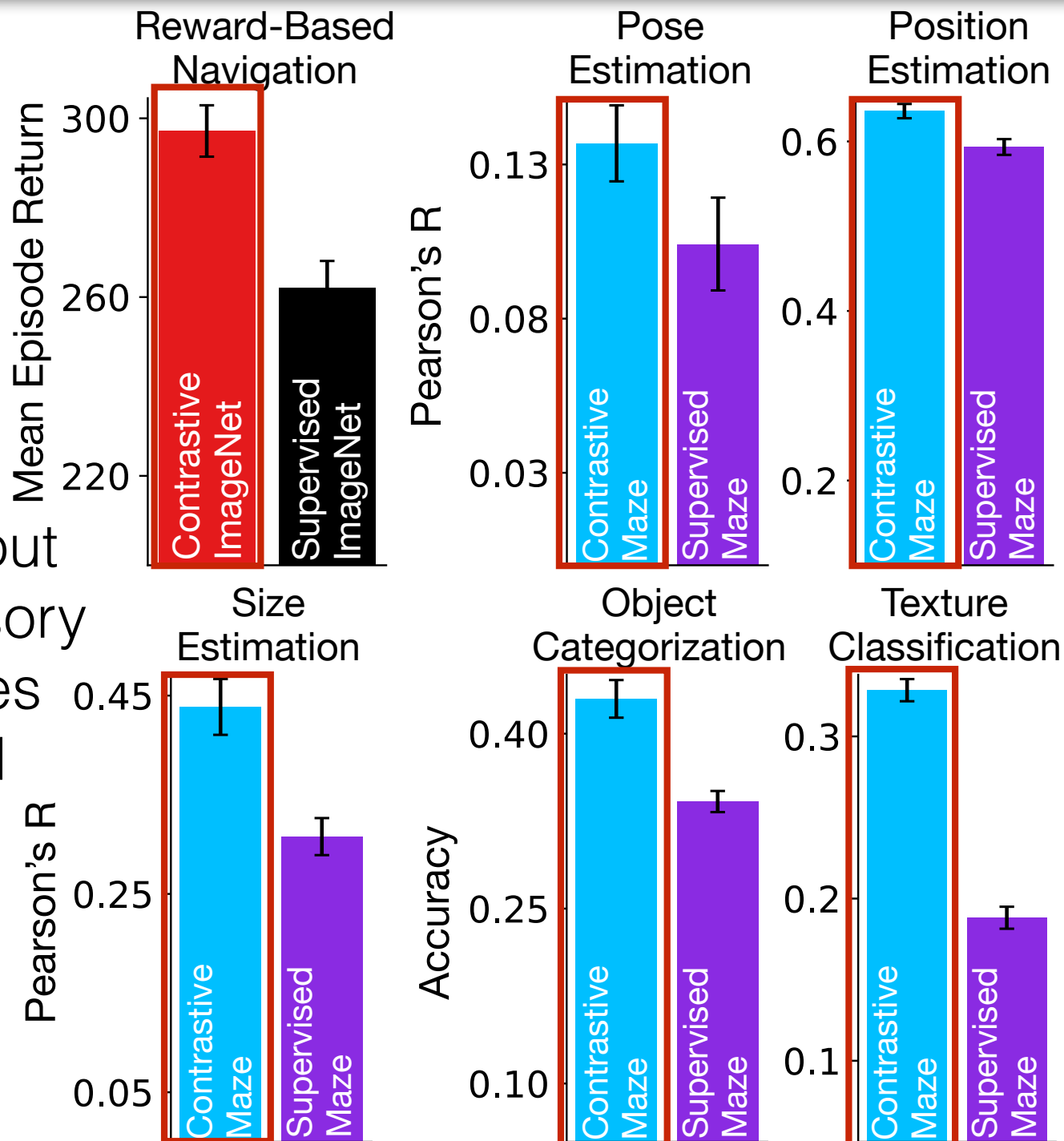
### Visual Scene Understanding

 <p>Plane</p> <p>Category</p> <p>f16</p> <p>Identity</p>	 <p>z axis rotation</p> <p>x axis rotation</p> <p>y axis rotation</p>	
 <p>Horizontal position: 80 pix</p> <p>Vertical position: -6 pix</p>	 <p>Perimeter: 78 pix</p> <p>Two-dimensional retinal area: 146 pix</p> <p>Three-dimensional object scale: 1.2x</p>	
<i>Object properties</i>		<i>Texture</i>

# Contrastive Models Yield Better Transfer Performance



# Contrastive Models Yield Better Transfer Performance



What about other sensory modalities beyond vision?

# Tactile Processing

---

## Task-Optimized Convolutional Recurrent Networks Align with Tactile Processing in the Rodent Brain

---

Trinity Chung<sup>\*,1</sup>, Yuchen Shen<sup>\*,2</sup>, Nathan C. L. Kong<sup>4</sup>, and Aran Nayebi<sup>2,3,1</sup>

<sup>1</sup>Robotics Institute, Carnegie Mellon University; Pittsburgh, PA 15213

<sup>2</sup>Machine Learning Department, Carnegie Mellon University; Pittsburgh, PA 15213

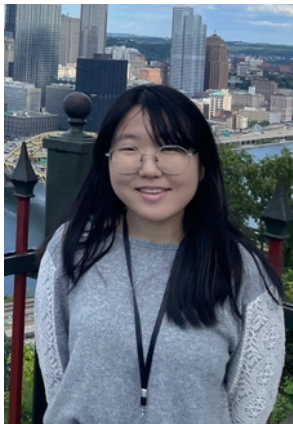
<sup>3</sup>Neuroscience Institute, Carnegie Mellon University; Pittsburgh, PA 15213

<sup>4</sup>Department of Psychology, University of Pennsylvania; Philadelphia, PA 19104

\* Equal contribution.

{trinityc, yuchens3, anayebi}@cs.cmu.edu; nclkong@sas.upenn.edu

## NeurIPS 2025 Oral



Trinity Chung\*



Yuchen Shen\*



Nathan C.L. Kong

# Why tactile?

- Tactile data is very useful in manipulating occluded and OOD objects
- Tactile hardware & sim is getting better!
- Tactile perception is still considerably under-explored in *both* neuroscience and robotics
- Many current tactile models are vision-based instead of force/torque-based

We hypothesize that model architectures that mimics brain-like processing will yield better performance for tactile data.

Trinity's search on arxiv...



# of tactile 625 results

Query: order: -announced\_date\_first; size: 50; date\_range: from 2024-06-01 to 2025-06-02; classification: Computer Science (cs), Quantitative Biology (q-bio); include\_cross\_list: True; terms: AND all=tactile; OR all=somatosensory; OR abstract=touch; NOT abstract=haptic

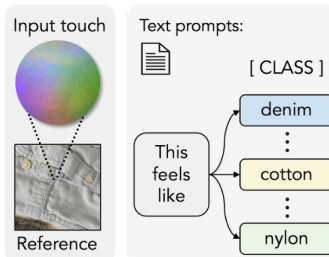
# of vision 2,577 results

Query: order: -announced\_date\_first; size: 50; date\_range: from 2024-06-01 to 2025-06-02; classification: Computer Science (cs), Quantitative Biology (q-bio); include\_cross\_list: True; terms: AND all=vision; AND title=visual

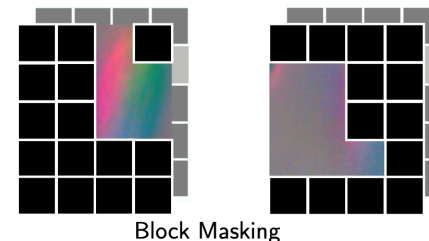
(both in the last 12 months)

e.g. UniTouch & Sparsh is trained on vision-based tactile sensors like Gelsight and DIGIT

Zero-shot Touch Understanding



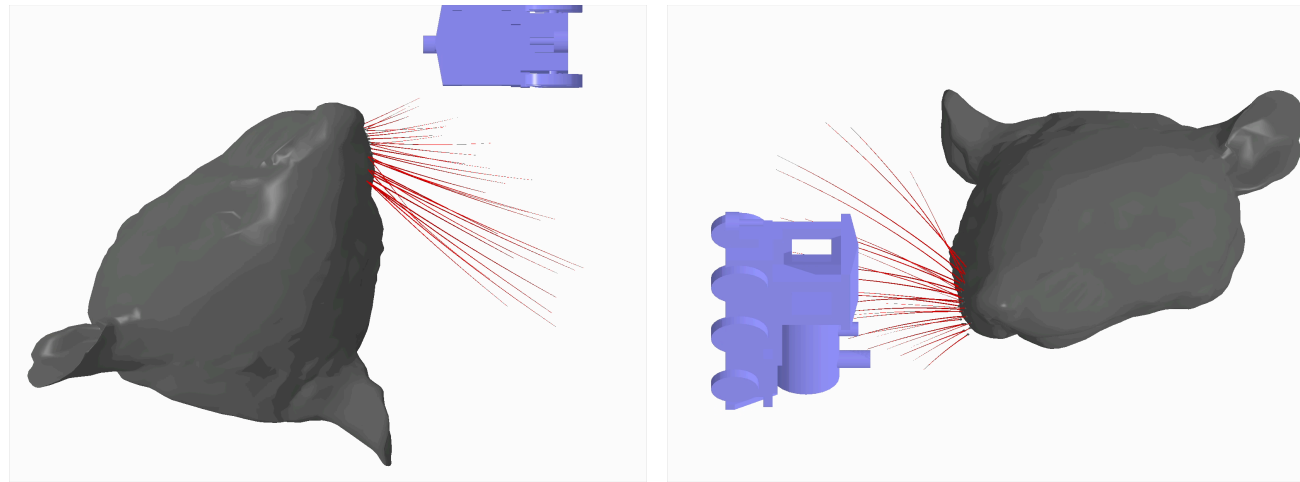
Sparsh (DINO - DINOv2)  
Self-distillation



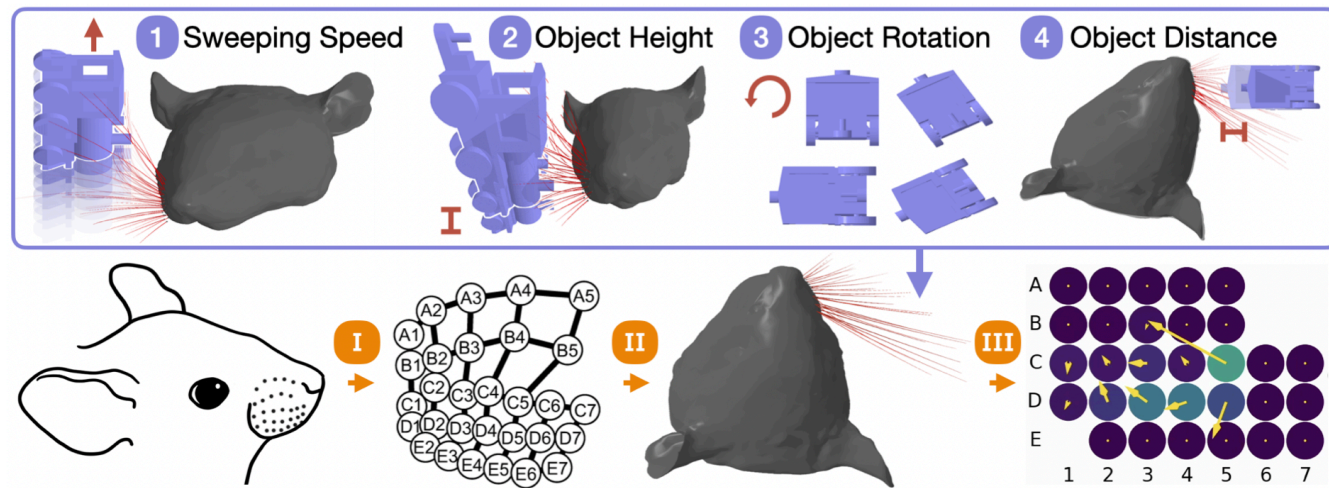
<https://arxiv.org/abs/2305.00596> <https://arxiv.org/abs/2410.24090>

# Training Data: Whisking Dataset Generation

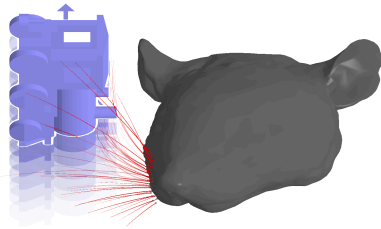
- Objects are whisked in simulation using WHISKiT [Zweifel et al., 2021], simulator based on Bullet Physics



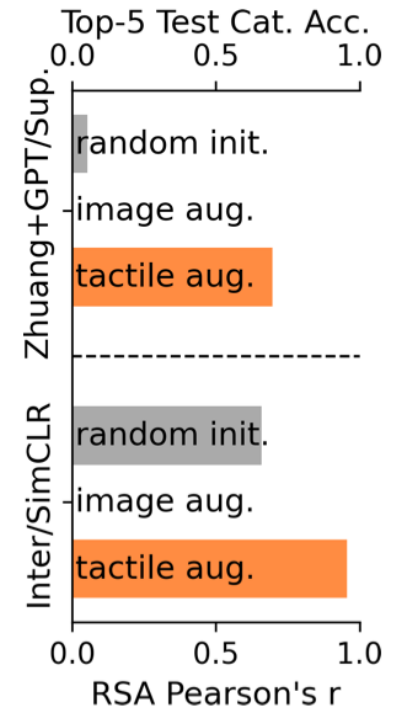
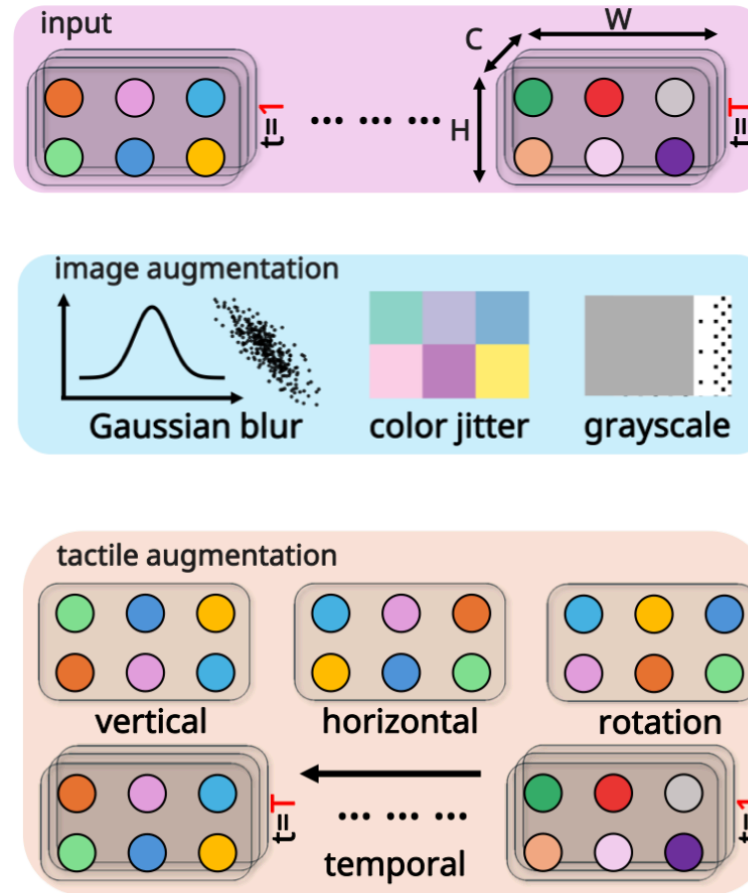
- 6-axis force/torque data for sweeping 9981 ShapeNet objects of 117 categories with various sweep augmentations



# Training Data: Tactile vs Image Augmentation

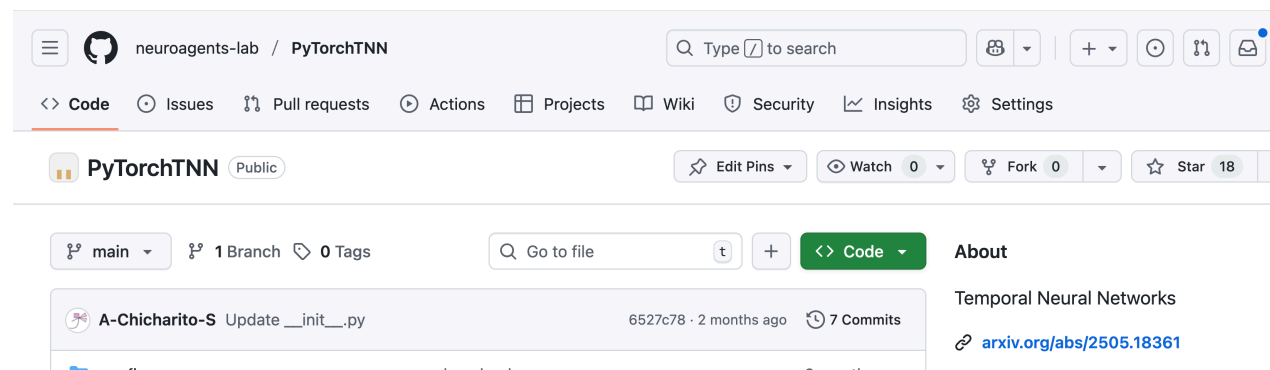
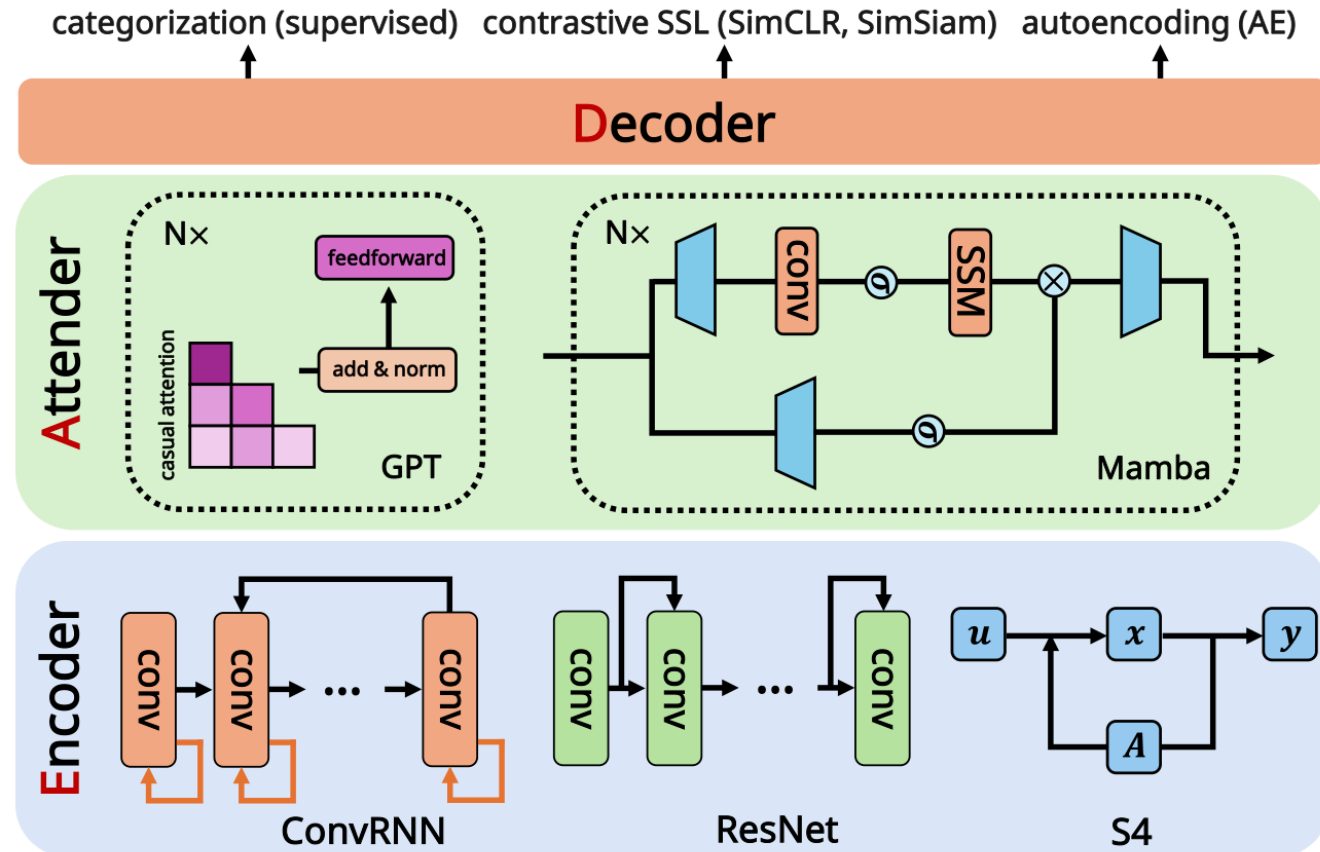


- Traditional **image augmentation** introduces Gaussian noise, color jitter, and grayscale.
- Our **tactile augmentation** vertically, horizontally, temporally flips, and rotates the features

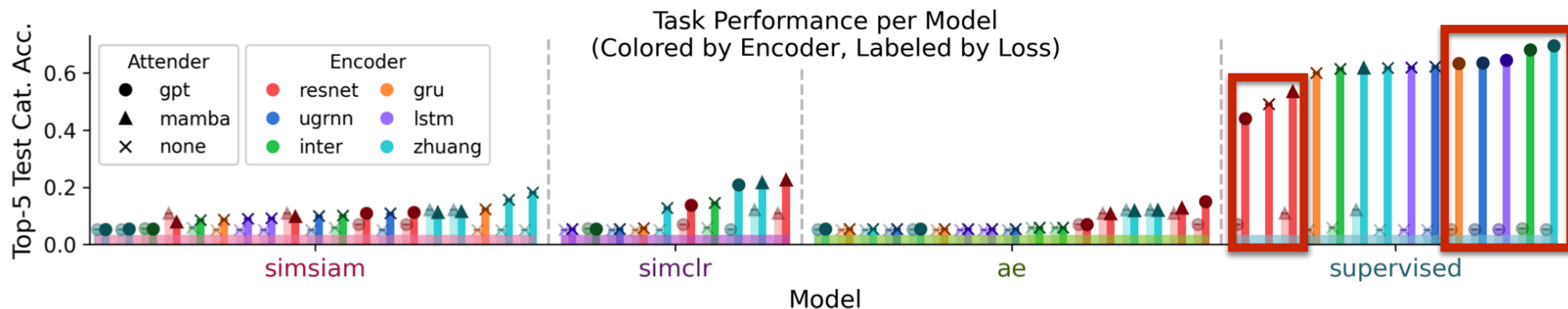


# Models: Encoder-Attender-Decoder (EAD) Architecture

- We needed a way to systematically search over the space of recurrent model architectures
- EAD architecture allows us to easily construct new models by swapping out modules
- Built using PyTorchTNN, now on Github! <https://github.com/neuroagents-lab/PyTorchTNN>

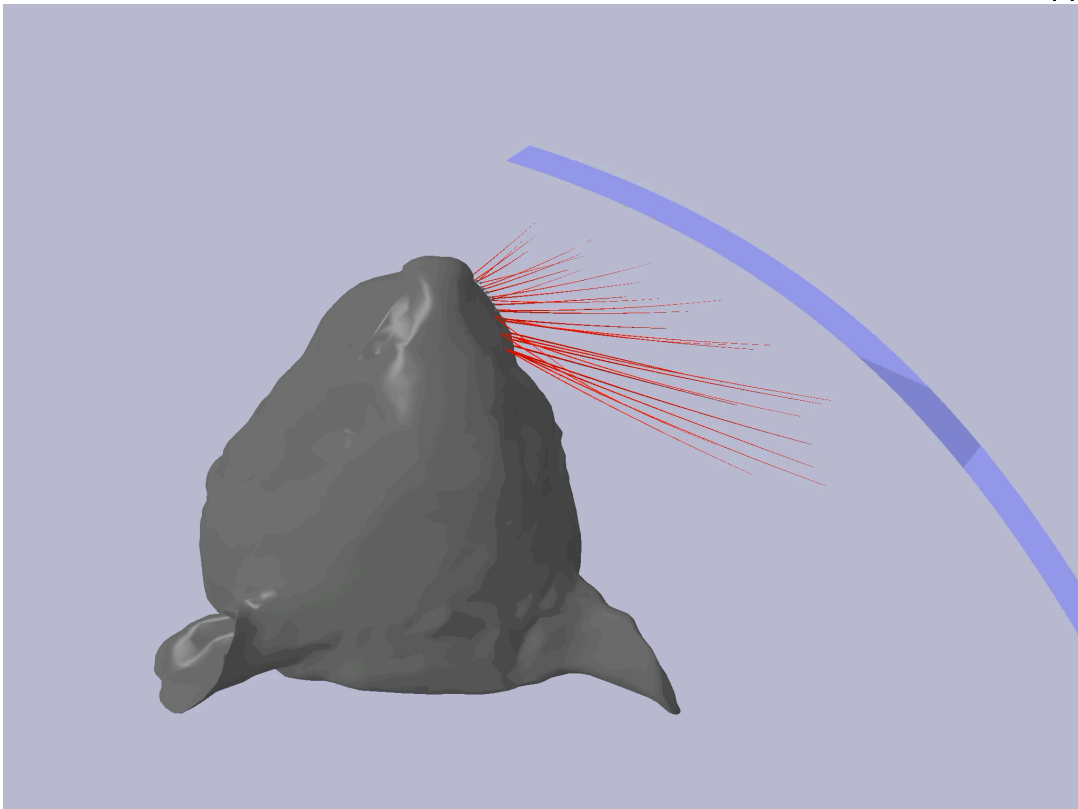
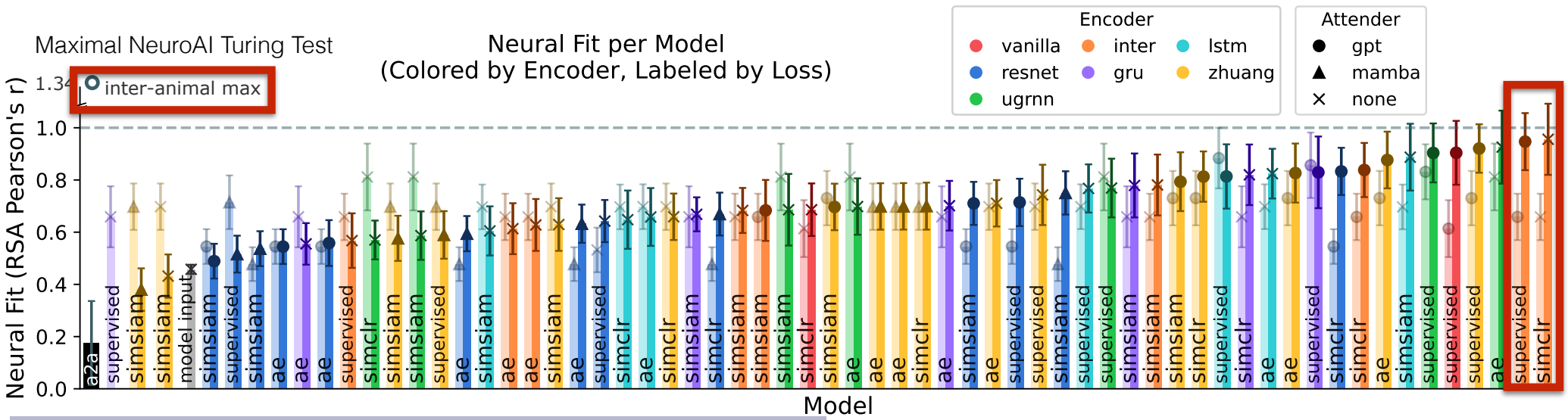


# Results: ConvRNN encoders perform best

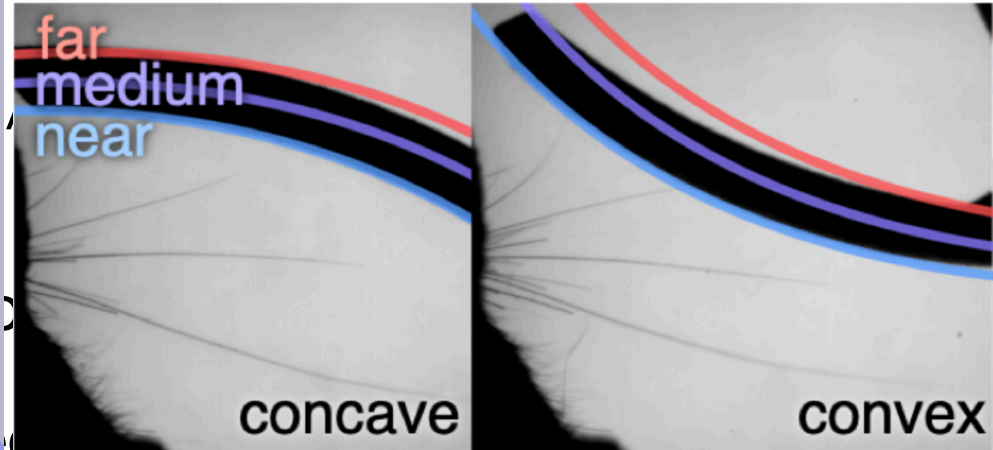


- Lighter color bar represents untrained version.
- S4 Encoders often don't even train on this task!
- Best model is ConvRNN (Encoder)+GPT (Attender) +Supervised (Decoder)

# Neural Evaluation: Results



uring Test, for this dataset at least  
Real

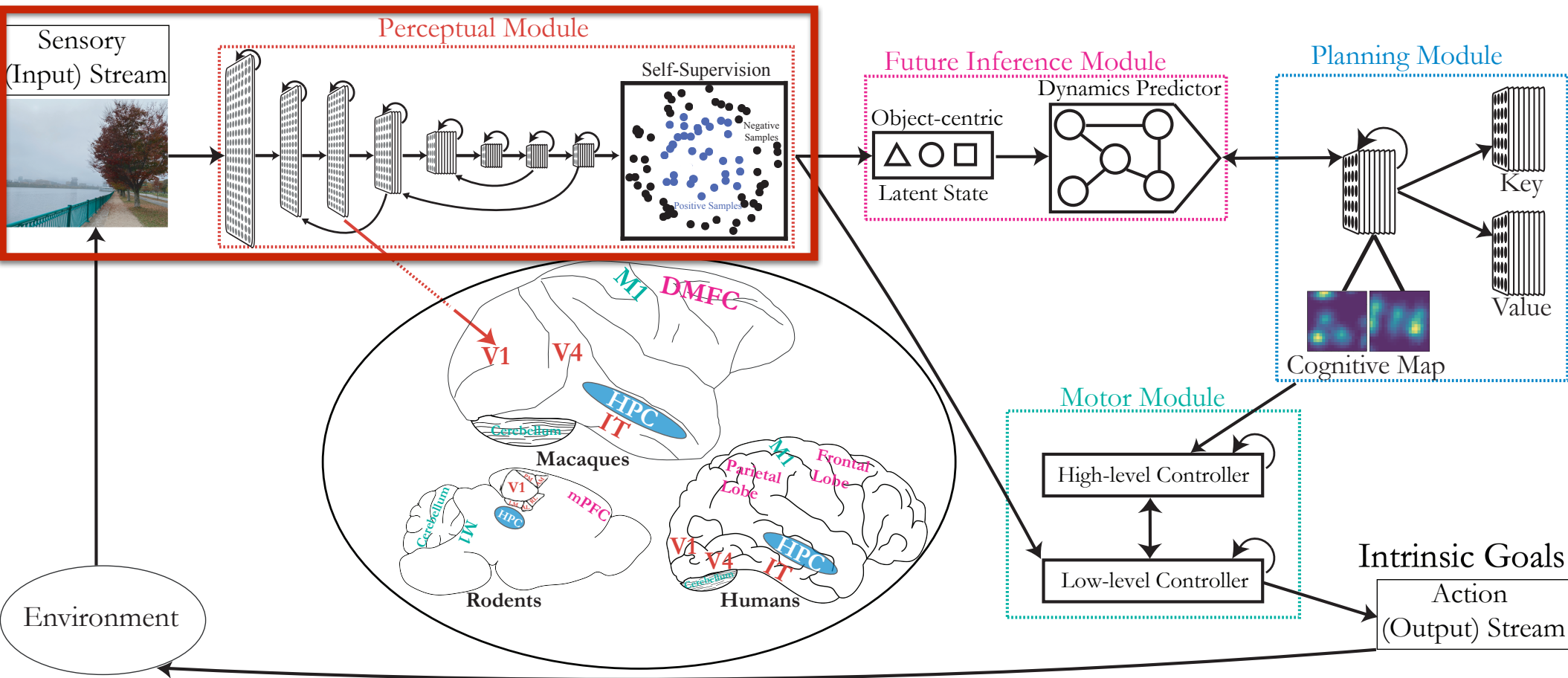


a neural alignment, possibly  
presentation in the somatosensory  
[Rogers 2002](#)  
(explore this!)

# Roadmap: Perception

How does the brain *represent*, *predict*, *plan*, and enable *action*?

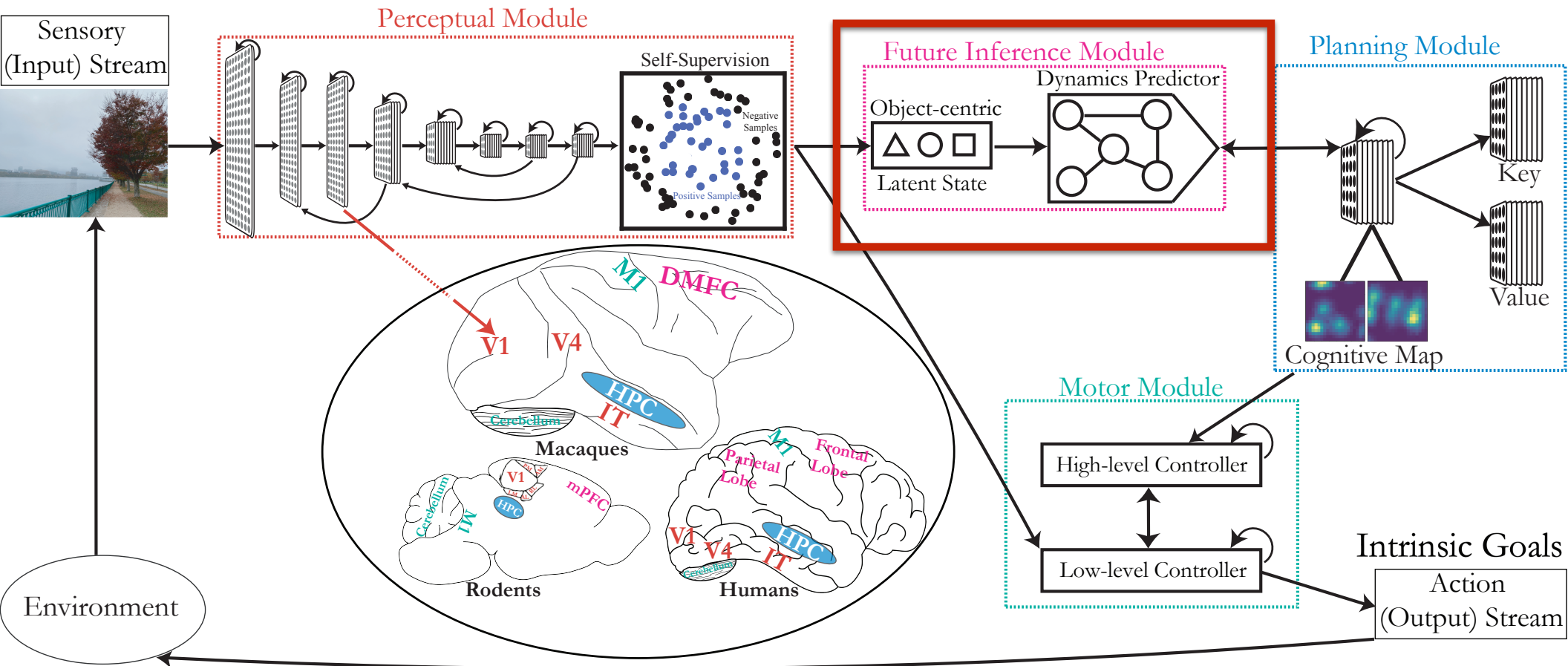
Recurrence + Contrastive SSL?



# Roadmap: Future Inference

How does the brain *represent*, *predict*, *plan*, and enable *action*?

Recurrence + Contrastive SSL?



# Reusable Latent Representations for Primate Mental Simulation

A. Nayebi, R. Rajalingham, M. Jazayeri, G.R. Yang

Neural foundations of mental simulation: future prediction of latent representations on dynamic scenes.

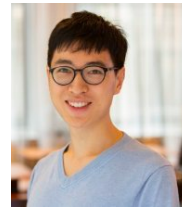
*NeurIPS 2023 (spotlight)*



Rishi Rajalingham



Mehrdad Jazayeri



Guangyu Robert Yang

# Visually-Grounded Mental Simulation



# Visually-Grounded Mental Simulation

Infer:  
Has this ice  
block been  
out longer?



# Visually-Grounded Mental Simulation

Infer:  
Has this ice  
block been  
out longer?



# Visually-Grounded Mental Simulation

Infer:  
Has this ice block been out longer?



# Visually-Grounded Mental Simulation

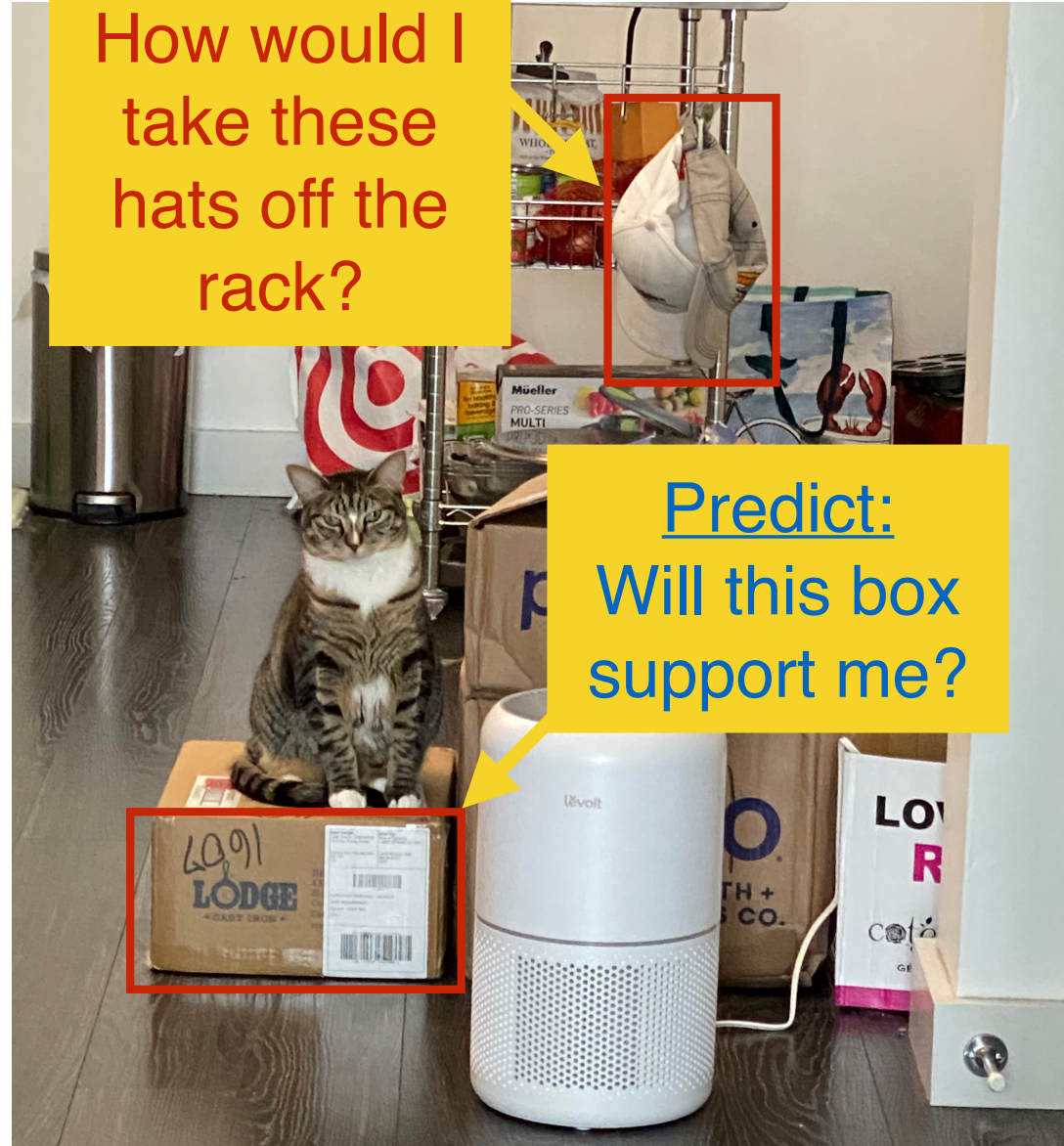
## Infer:

Has this ice block been out longer?



## Plan:

How would I take these hats off the rack?



## Predict:

Will this box support me?

# Visually-Grounded Mental Simulation

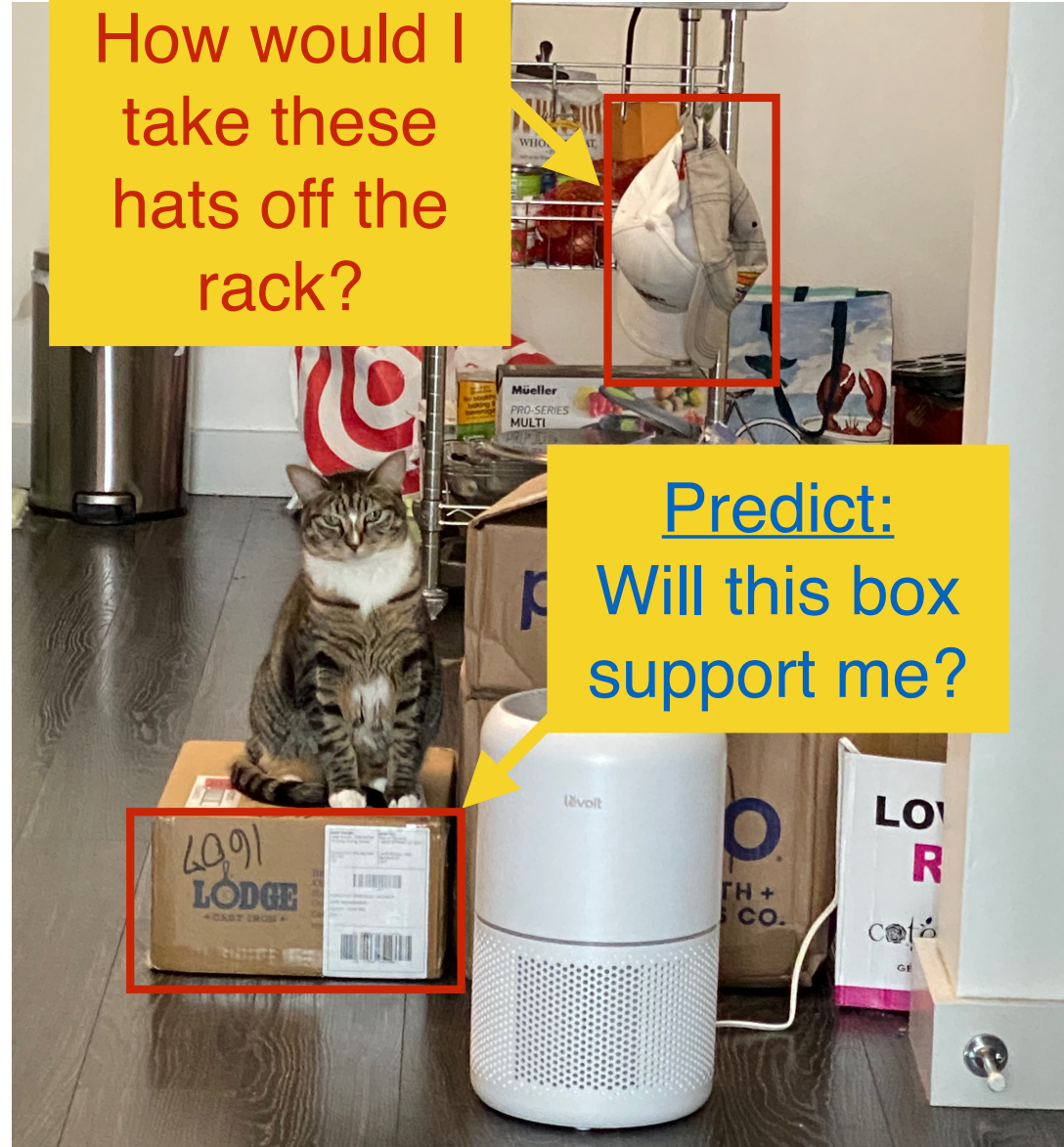
Infer:

Has this ice block been out longer?



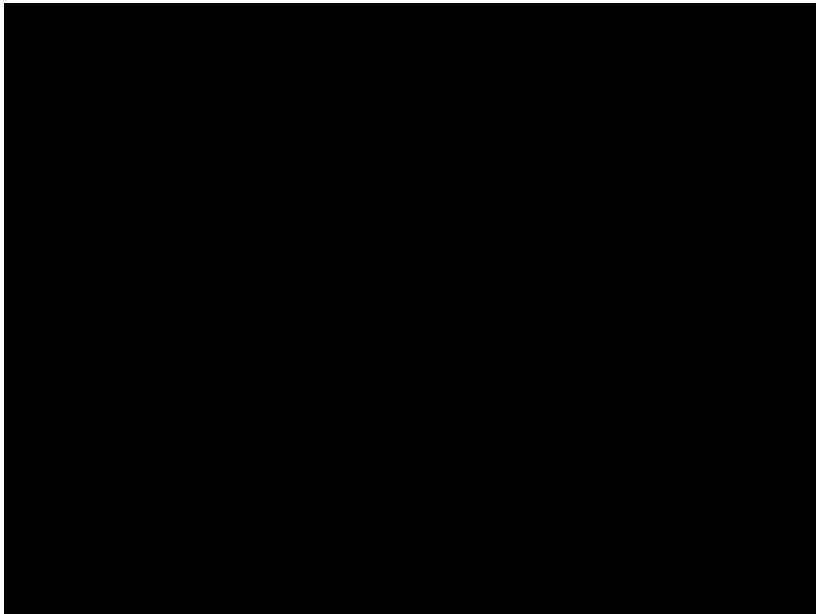
Plan:

How would I take these hats off the rack?



Predict:

Will this box support me?



# Visually-Grounded Mental Simulation

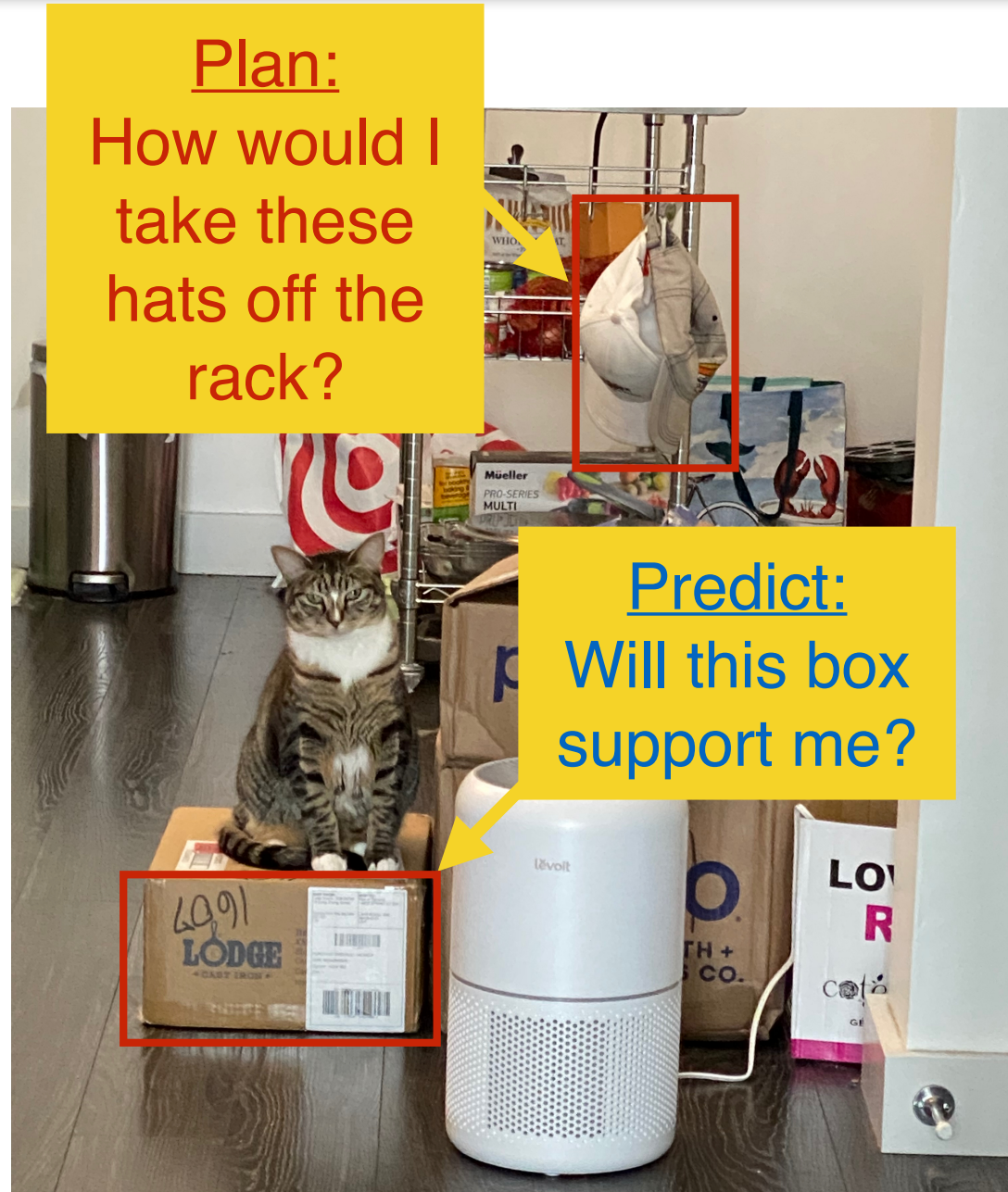
## Infer:

Has this ice block been out longer?



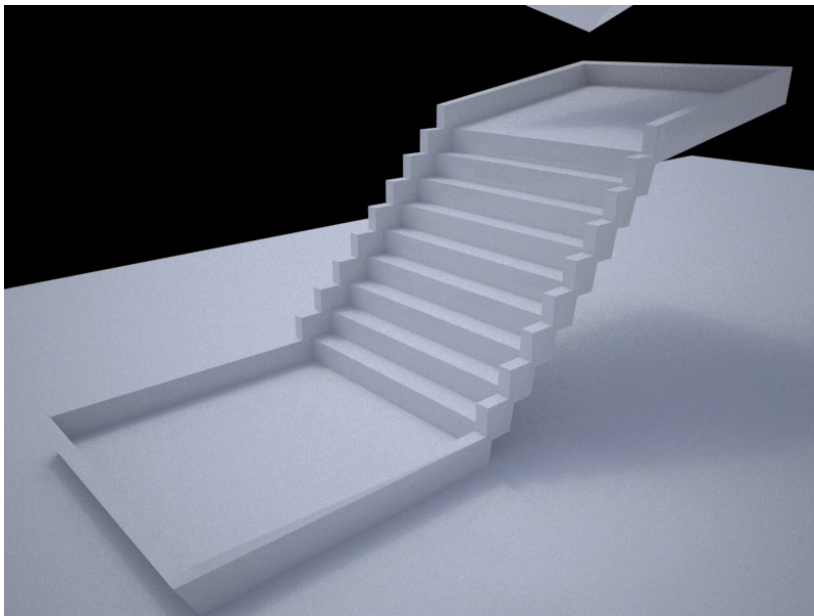
## Plan:

How would I take these hats off the rack?



## Predict:

Will this box support me?

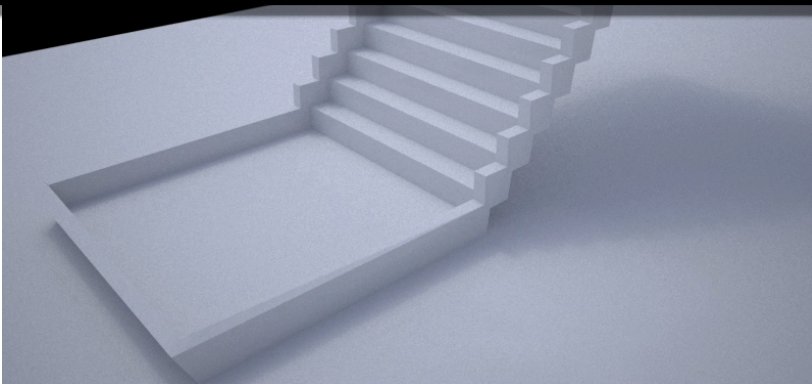


# Visually-Grounded Mental Simulation

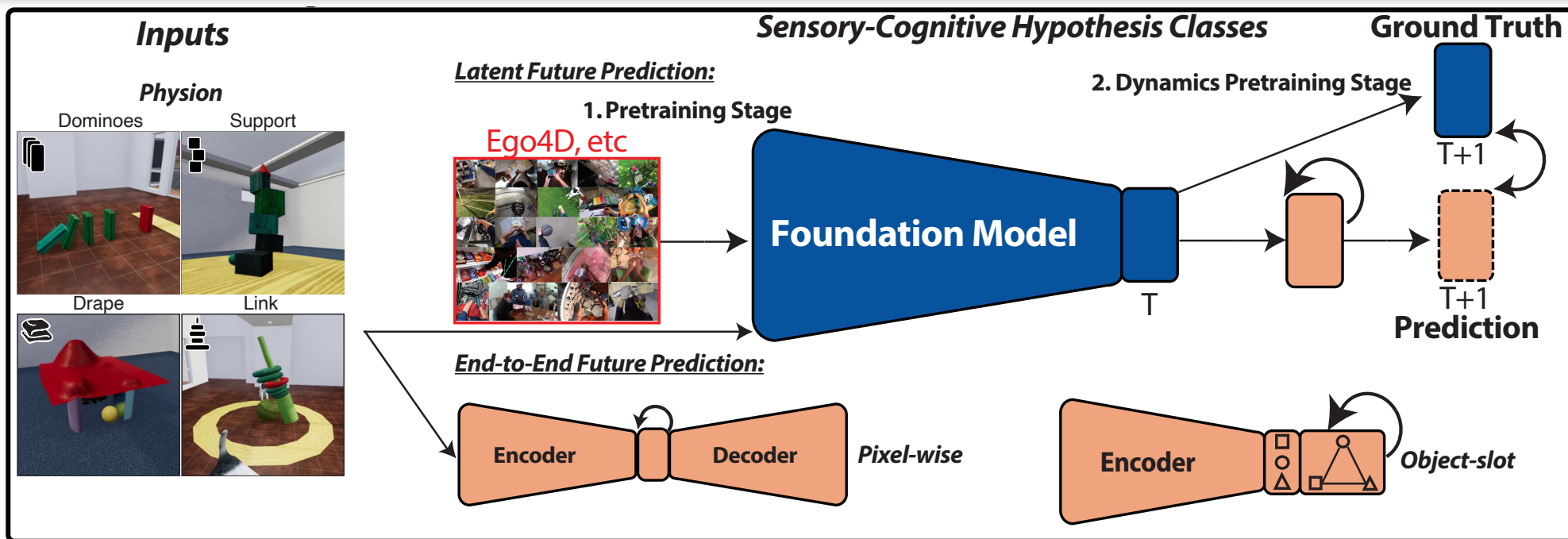


## Neurobiological Puzzle:

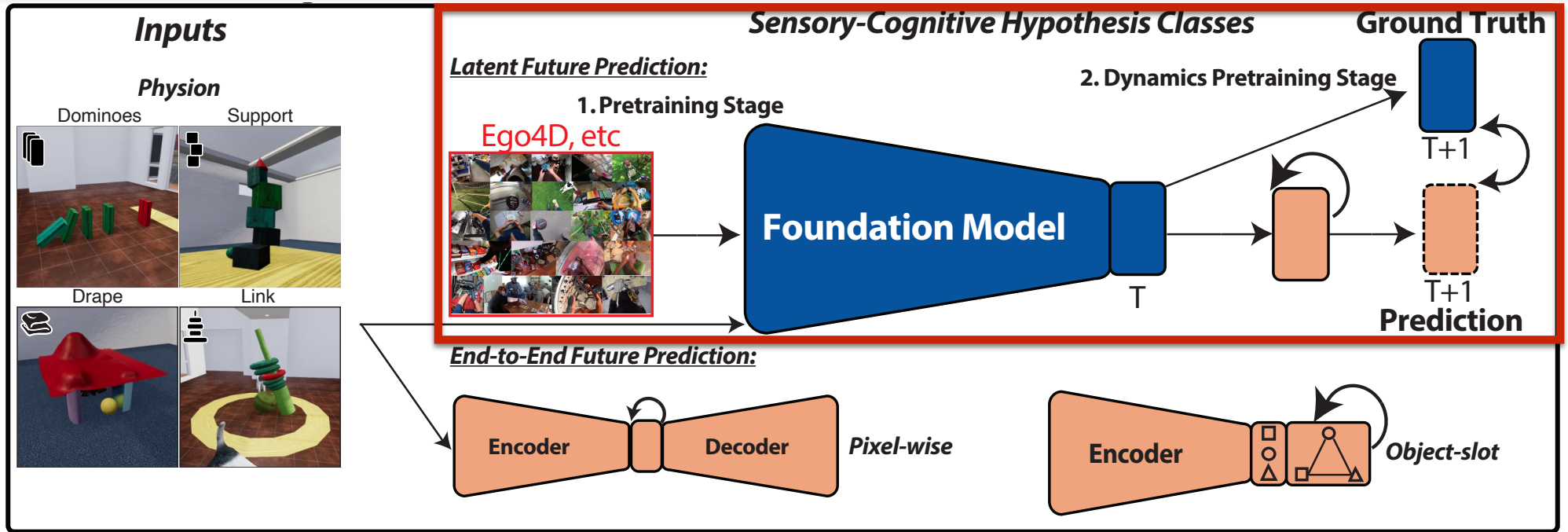
What are the functional constraints that enable us to predict the future state of our environment *across* diverse settings?



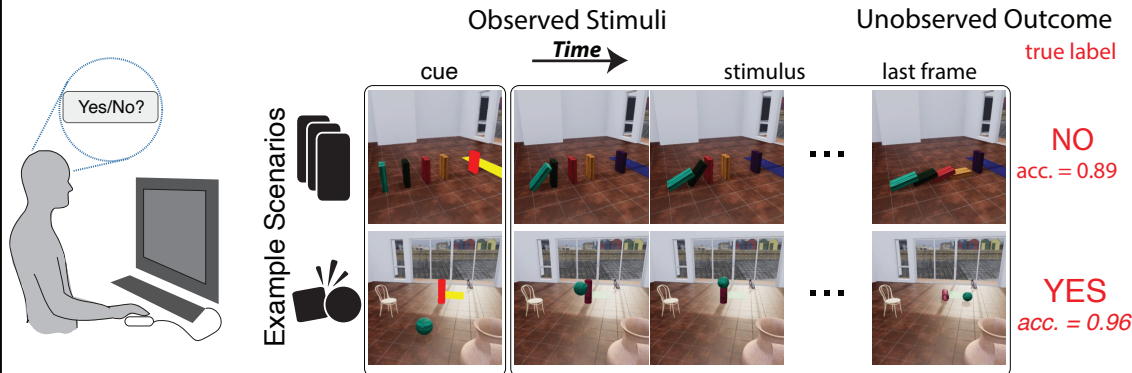
# Overall Approach: Sensory-Cognitive Hypotheses



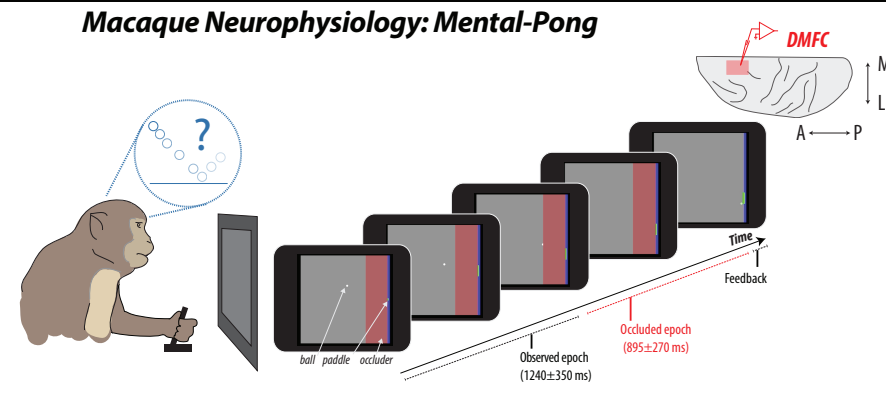
# Overall Approach: Sensory-Cognitive Hypotheses



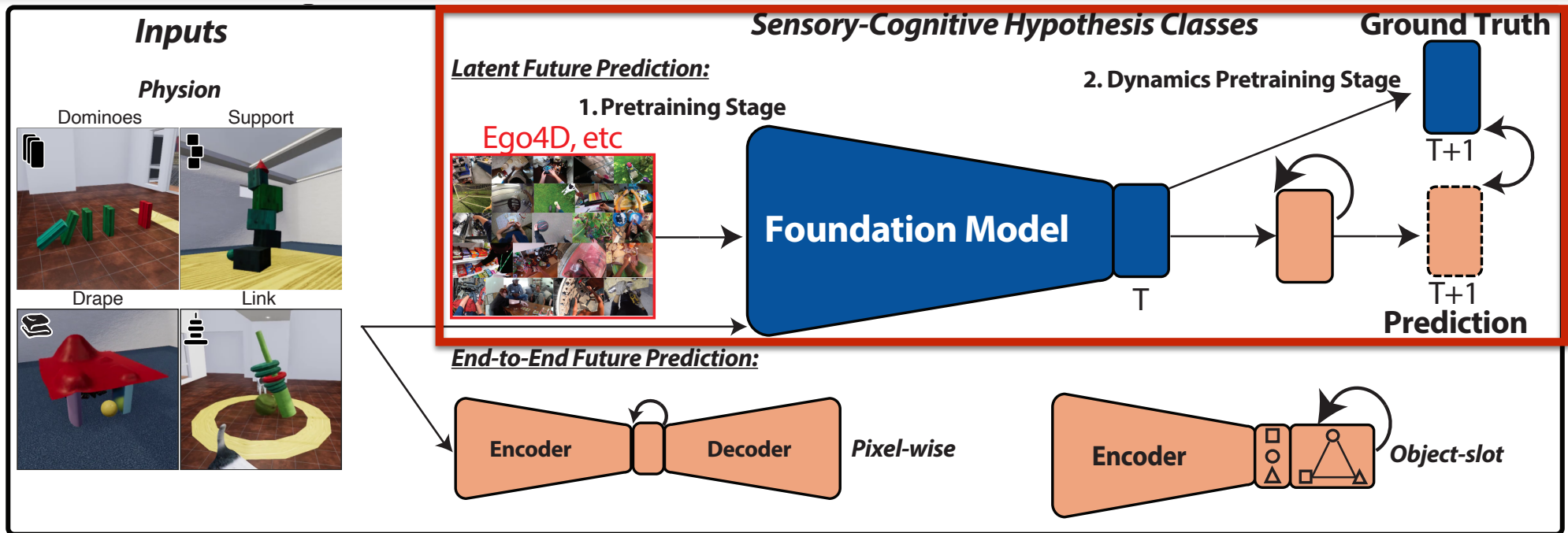
## Human Behavior: Physion Object Contact Prediction (OCP)



## Macaque Neurophysiology: Mental-Pong



# Latent Future Prediction



Learn a partial, *implicit* representation of the physical world by performing a challenging vision task (“foundation model”)

What vision task?

We do far more than engage with static images!

Leverage these dynamics to do explicit future prediction

# Video Foundation Models

## Ego4D: everyday activity around the world



$$\mathcal{L}_{contrastive} = \sum_{b \in B} \log \frac{e^{\overbrace{\mathcal{S}(\mathbf{z}_i^b, \mathbf{z}_j^b)}^{\text{attract}}}}{e^{\mathcal{S}(\mathbf{z}_i^b, \mathbf{z}_j^b)} + e^{\overbrace{\mathcal{S}(\mathbf{z}_i^b, \mathbf{z}_k^b)}^{\text{repel}}} + e^{\overbrace{\mathcal{S}(\mathbf{z}_i^b, \tilde{\mathbf{z}}_i^b)}^{\text{repel}}}}$$
$$[I_i, I_{j>i}, I_{k>j}]^{1:B}$$

## Ego4D: A massive-scale egocentric dataset

3,670 hours of in-the-wild daily life activity

931 participants from 74 worldwide locations

Multimodal: audio, 3D scans, IMU, stereo, multi-camera

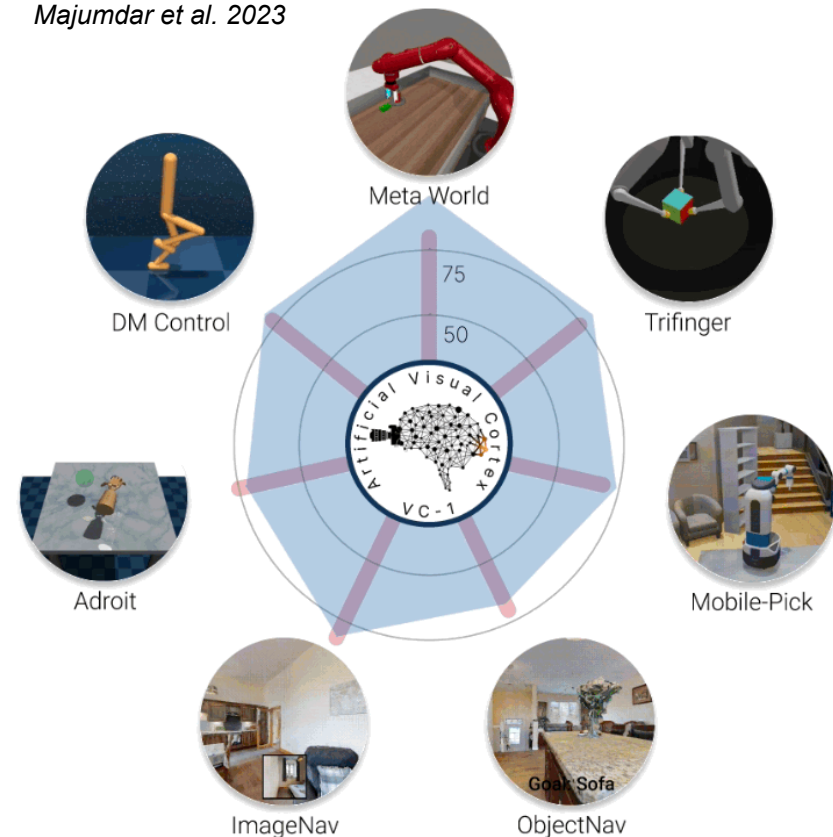


# Video Foundation Models

## Ego4D: everyday activity around the world



Majumdar et al. 2023



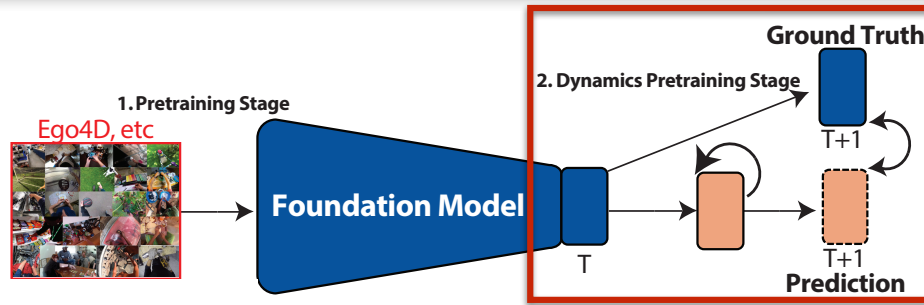
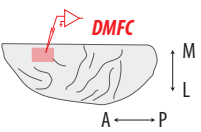
## Ego4D: A massive-scale egocentric dataset

- 3,670 hours of in-the-wild daily life activity
- 931 participants from 74 worldwide locations
- Multimodal: audio, 3D scans, IMU, stereo, multi-camera



Grauman et al. 2022

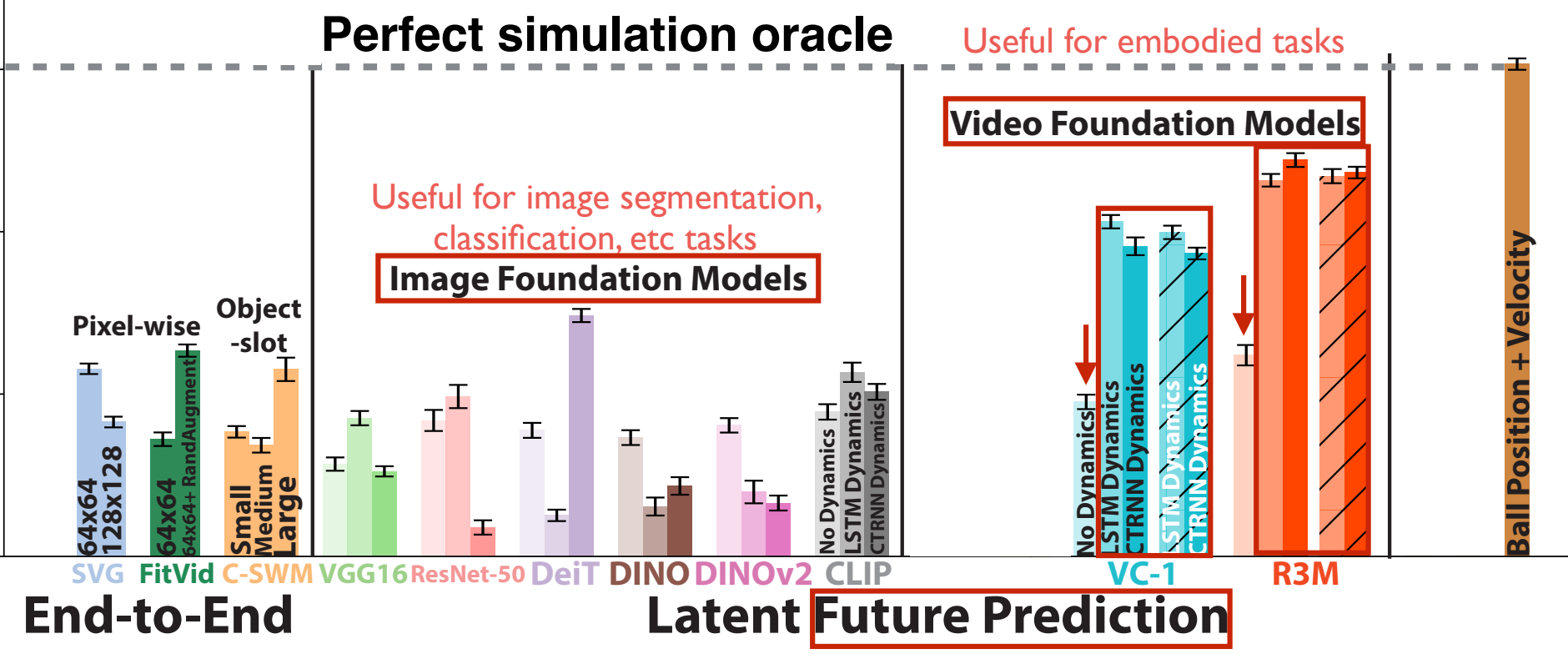
# Video Foundation Future Prediction Best Predict Neurons



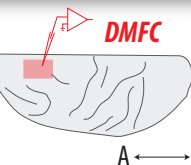
**Pretraining on Ego4D is not enough on its own:  
Need explicit future prediction!**

Neural Predictivity  
(Pearson's R)

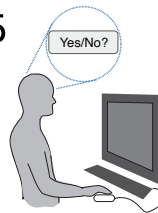
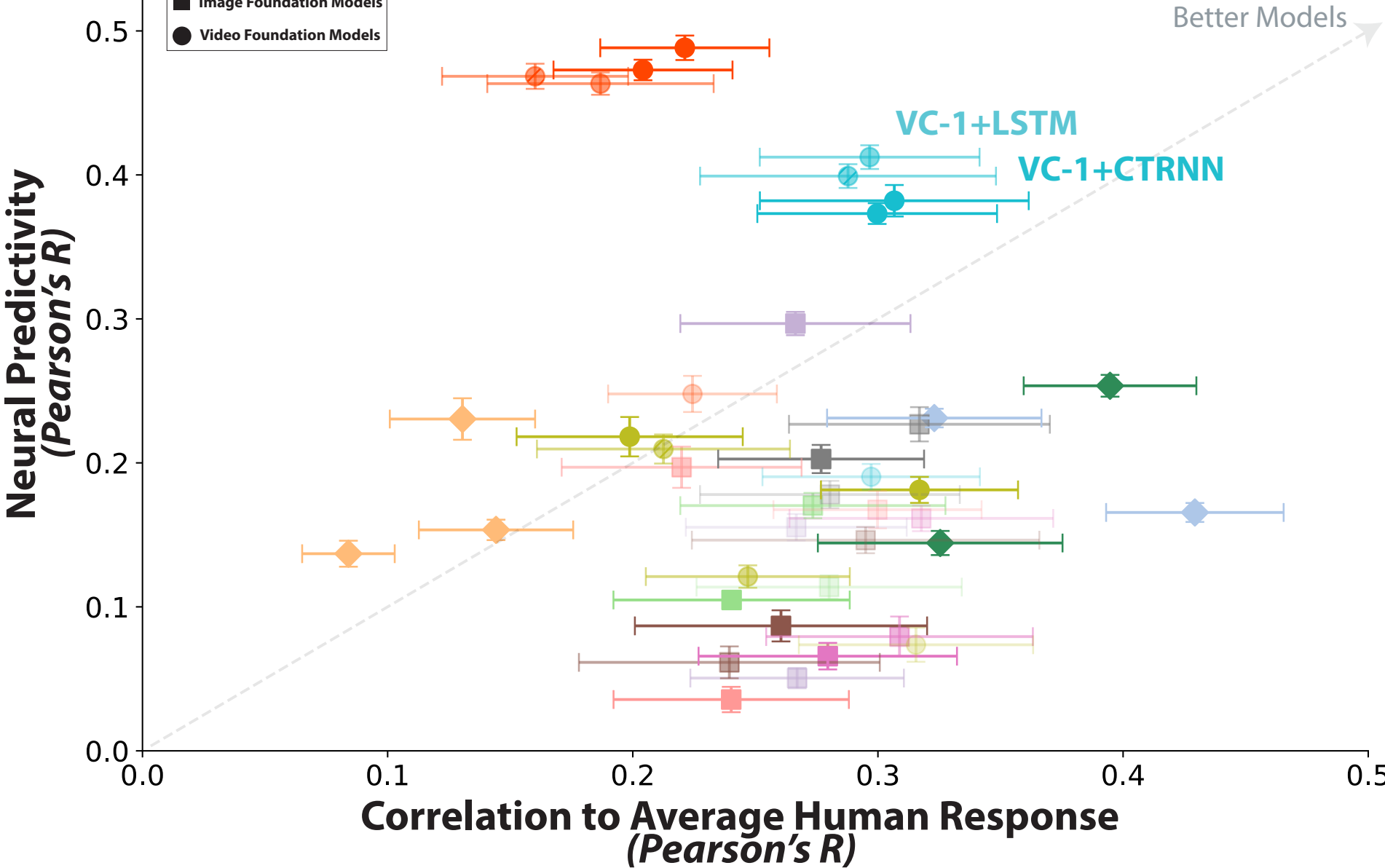
1.0  
0.8  
0.6  
0.4  
0.2  
0.0



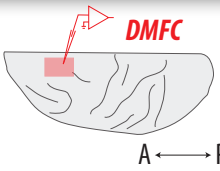
# Dynamically-Equipped Video Foundation Models Can Match Both



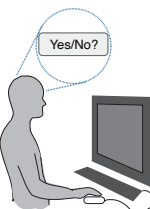
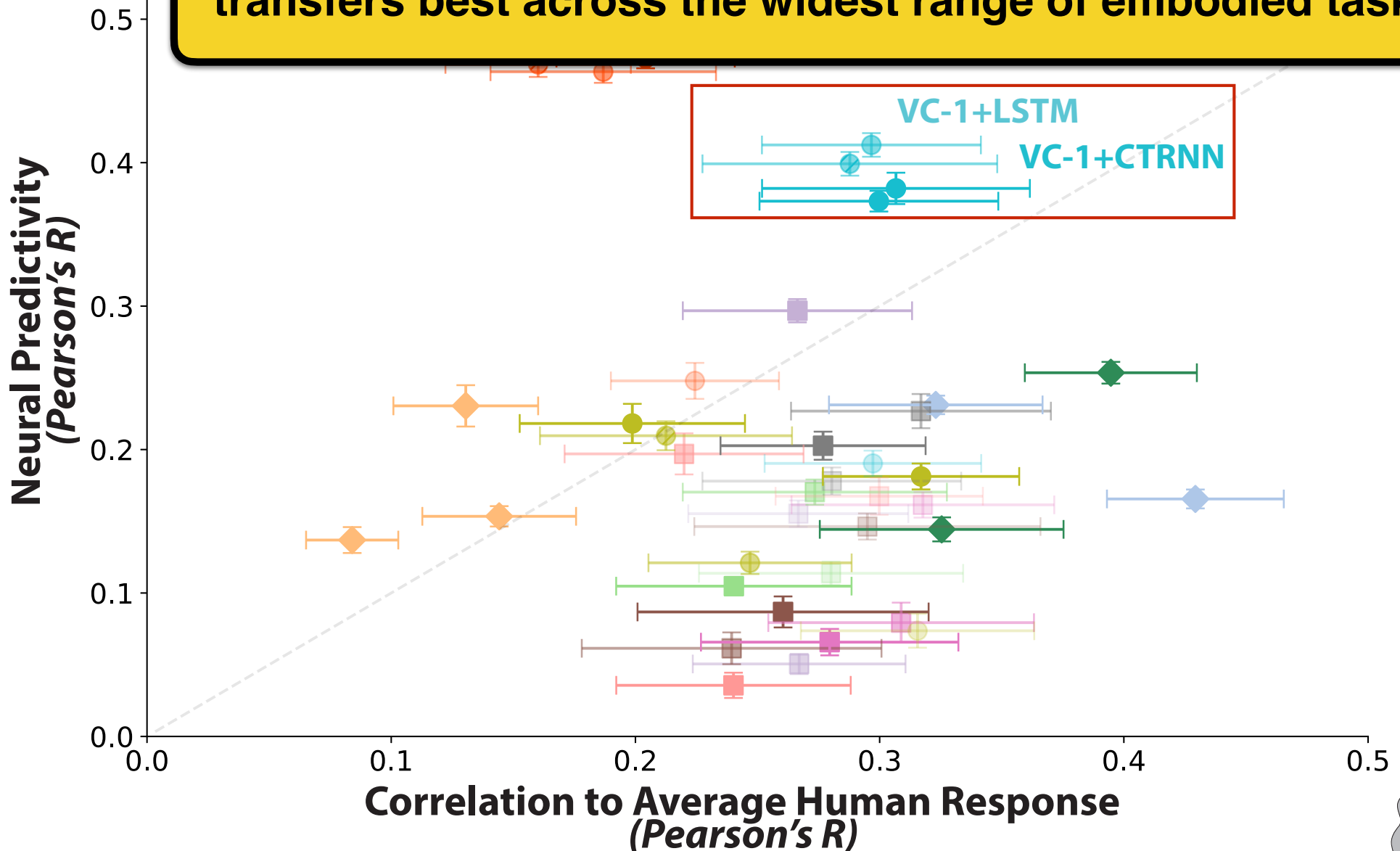
- ◆ End-to-End
- Image Foundation Models
- Video Foundation Models



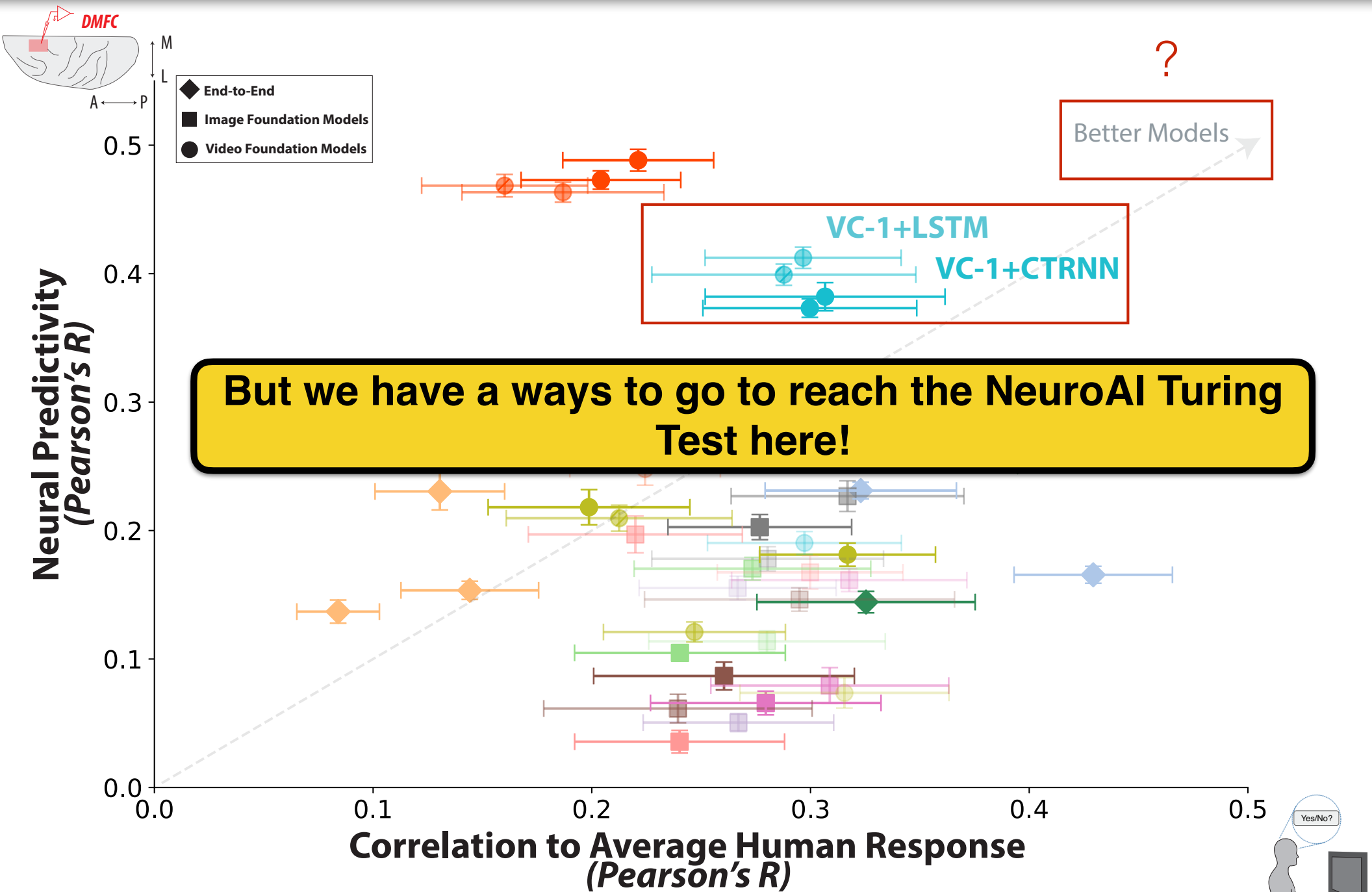
# Dynamically-Equipped Video Foundation Models Can Match Both



**Exposed to the largest variety of egocentric video sources & transfers best across the widest range of embodied tasks.**



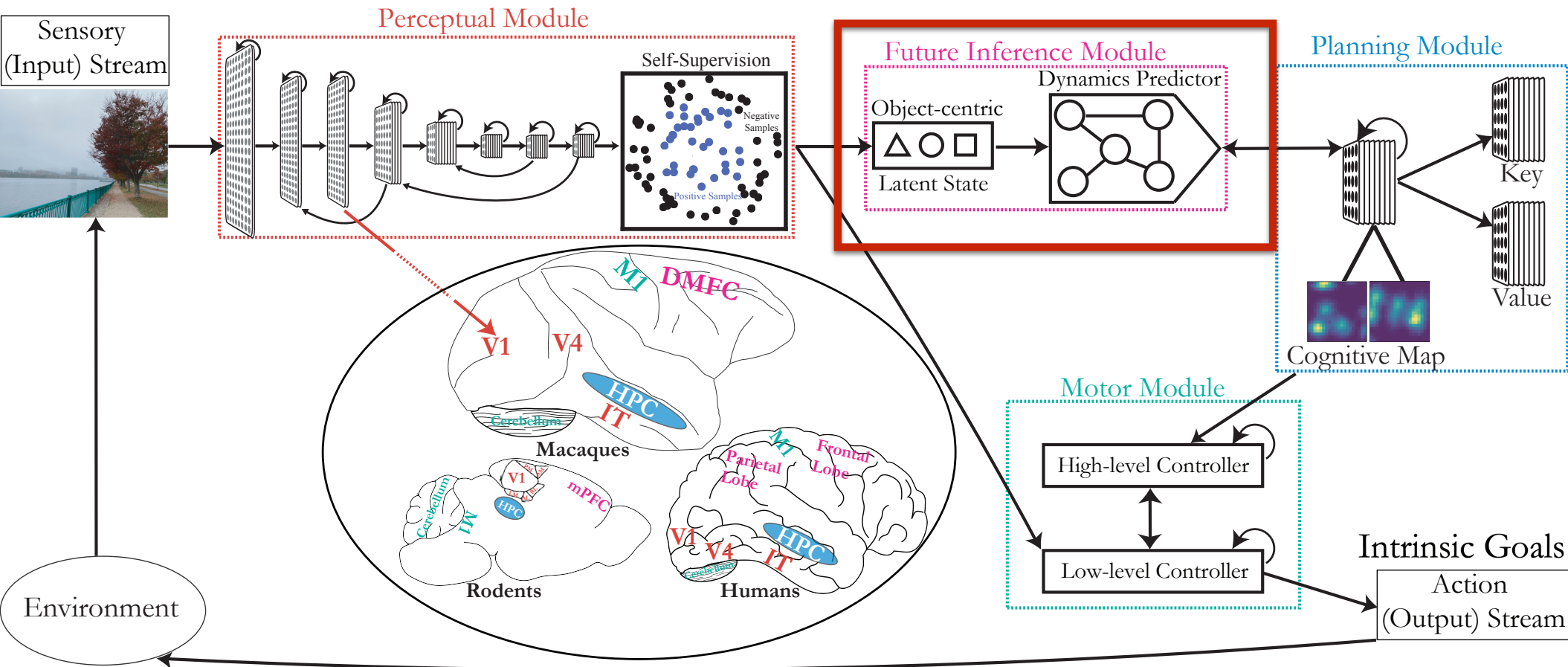
# Dynamically-Equipped Video Foundation Models Can Match Both



# Roadmap: Future Inference

How does the brain *represent*, *predict*, *plan*, and enable *action*?

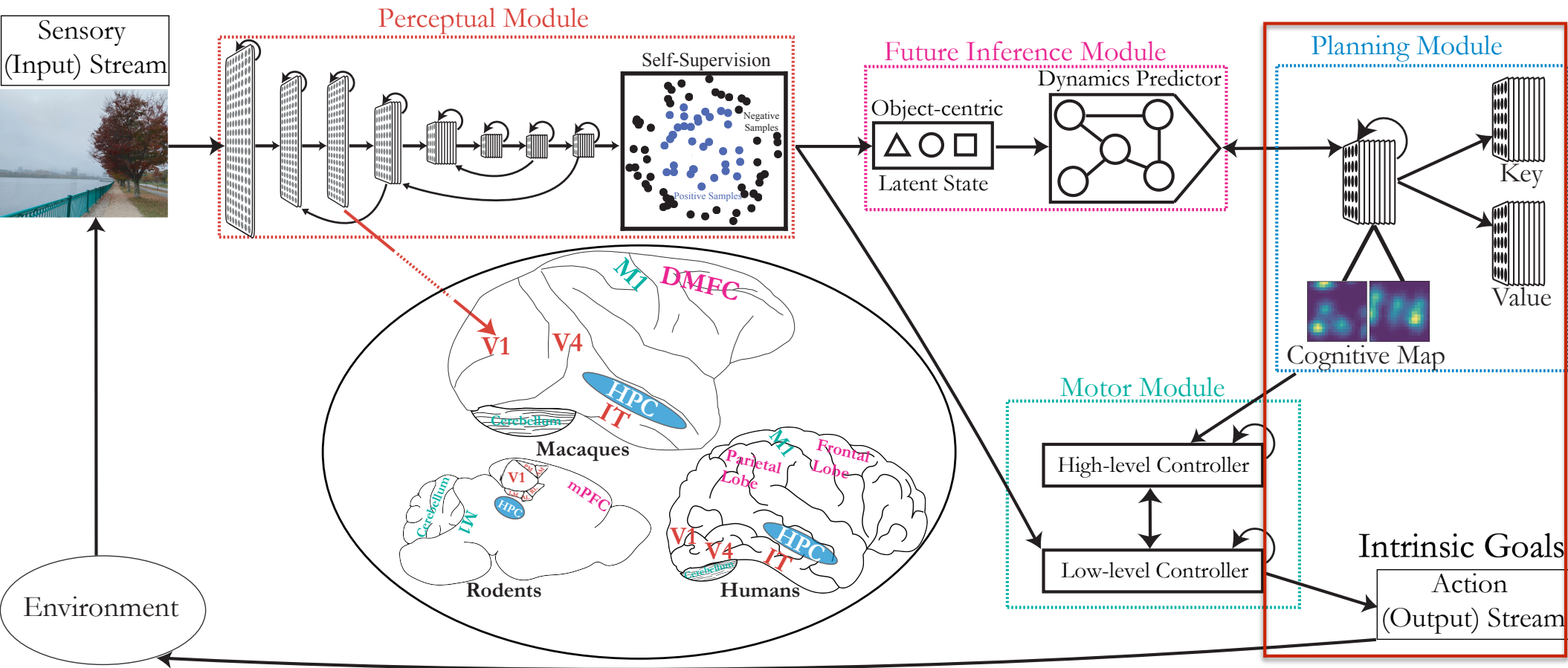
Recurrence + Contrastive SSL? Latent Future Prediction?



# Roadmap: Planning & Action

How does the brain *represent*, *predict*, *plan*, and enable *action*?

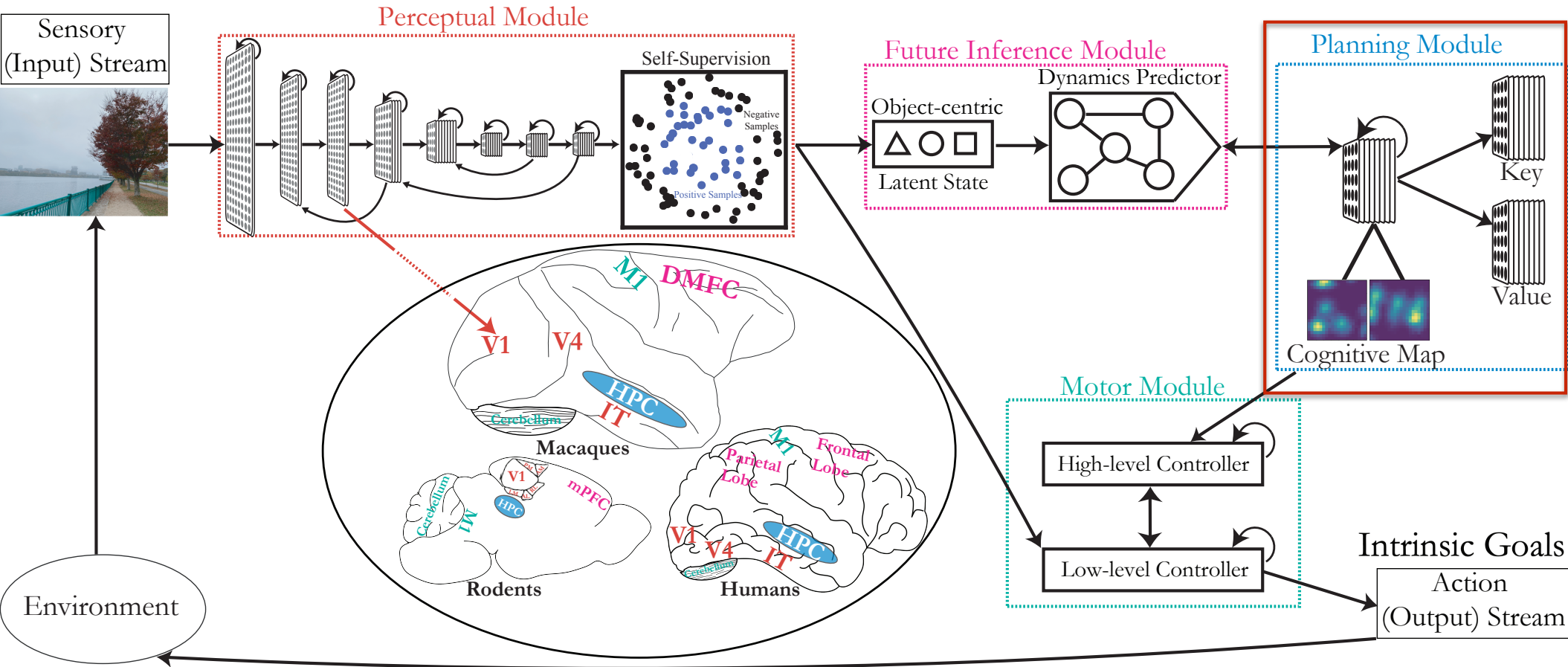
Recurrence + Contrastive SSL? Latent Future Prediction?



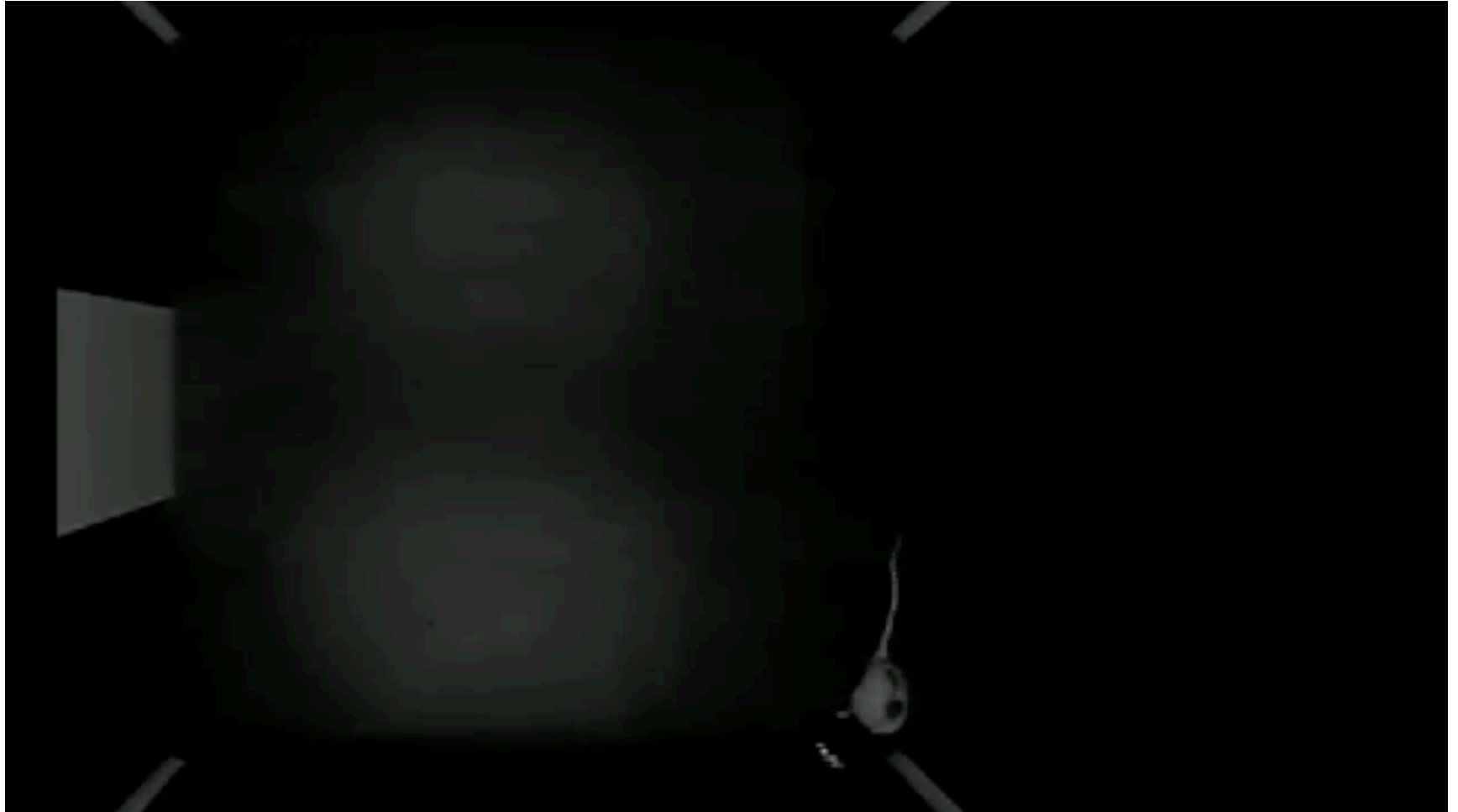
# Roadmap: Planning

How does the brain *represent*, *predict*, *plan*, and enable *action*?

Recurrence + Contrastive SSL? Latent Future Prediction?

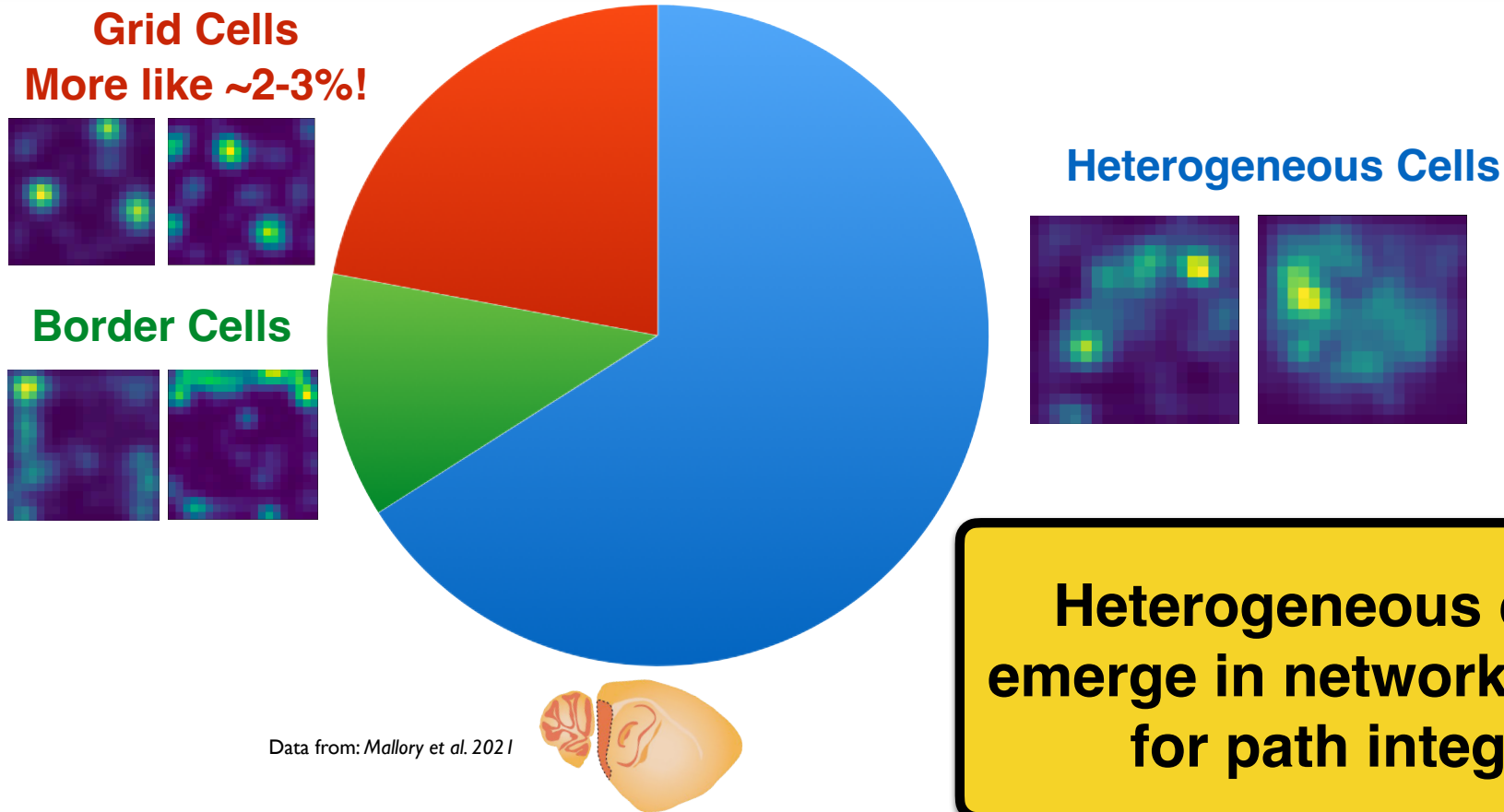


# Hippocampal-Entorhinal Spatial Map

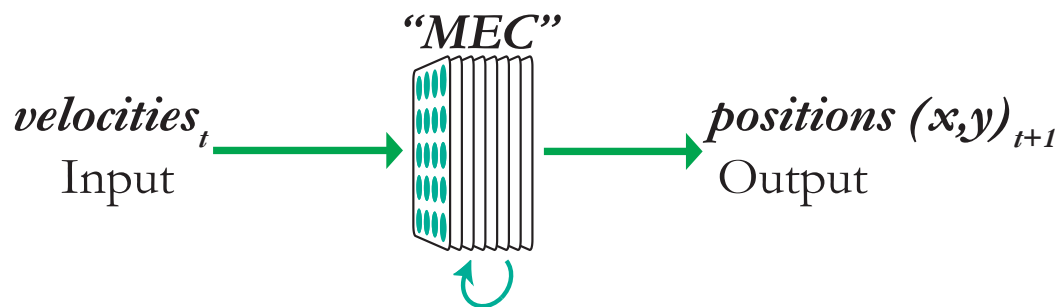


Place Cell  
(Hippocampus)

# A Task-Optimized Account of Heterogeneity



**Heterogeneous cell types emerge in networks optimized for path integration!**



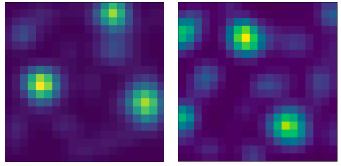
**Explaining heterogeneity in medial entorhinal cortex with task-driven neural networks**

Aran Nayebi<sup>1,\*</sup>, Alexander Attinger<sup>2</sup>, Malcolm G. Campbell<sup>2</sup>, Kiah Hardcastle<sup>2</sup>, Isabel I.C. Low<sup>1,2,7</sup>, Caitlin S. Mallory<sup>2</sup>, Gabriel C. Mel<sup>1</sup>, Ben Sorscher<sup>4</sup>, Alex H. Williams<sup>6,7</sup>, Surya Ganguli<sup>4,7,8</sup>, Lisa M. Giocomo<sup>2,7</sup>, and Daniel L.K. Yamins<sup>3,5,7</sup>

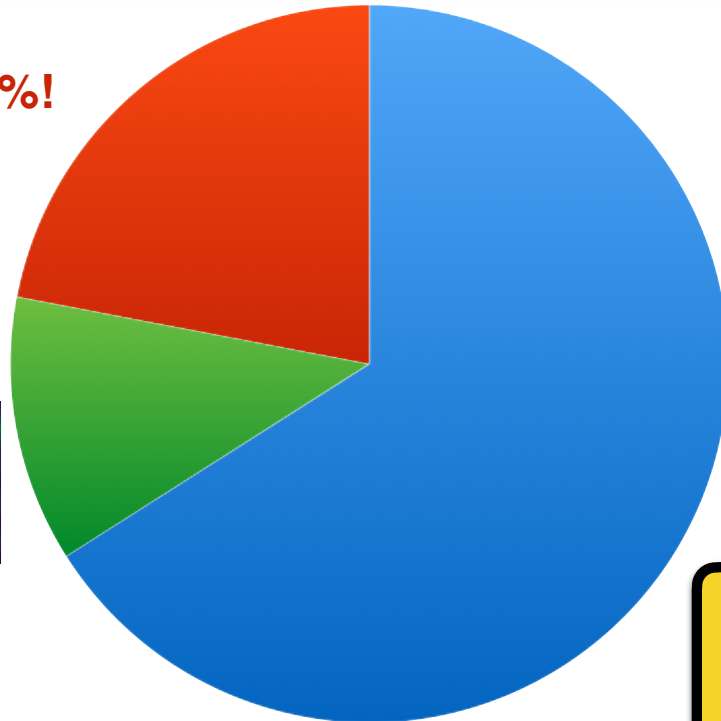
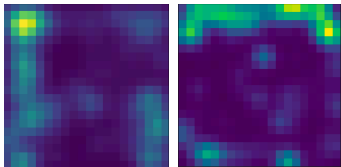
NeurIPS 2021 (spotlight)

# A Task-Optimized Account of Heterogeneity

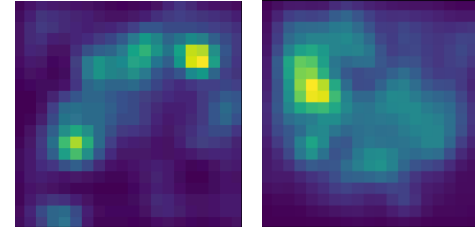
**Grid Cells**  
More like ~2-3%!



**Border Cells**

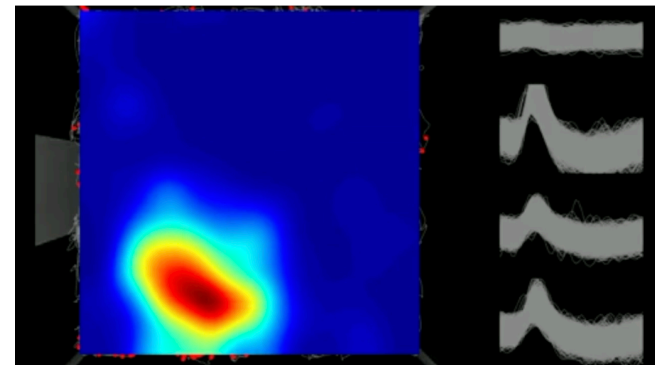
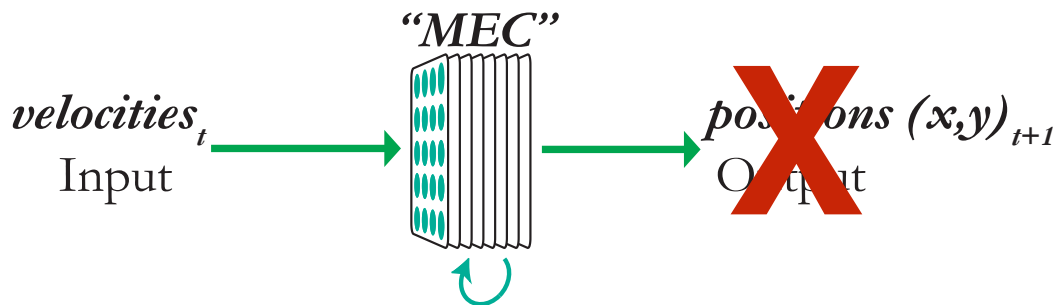
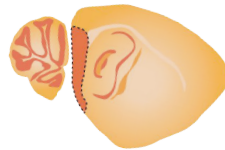


**Heterogeneous Cells**



**Heterogeneous cell types emerge in networks optimized for place cell integration!**

Data from: Mallory et al. 2021

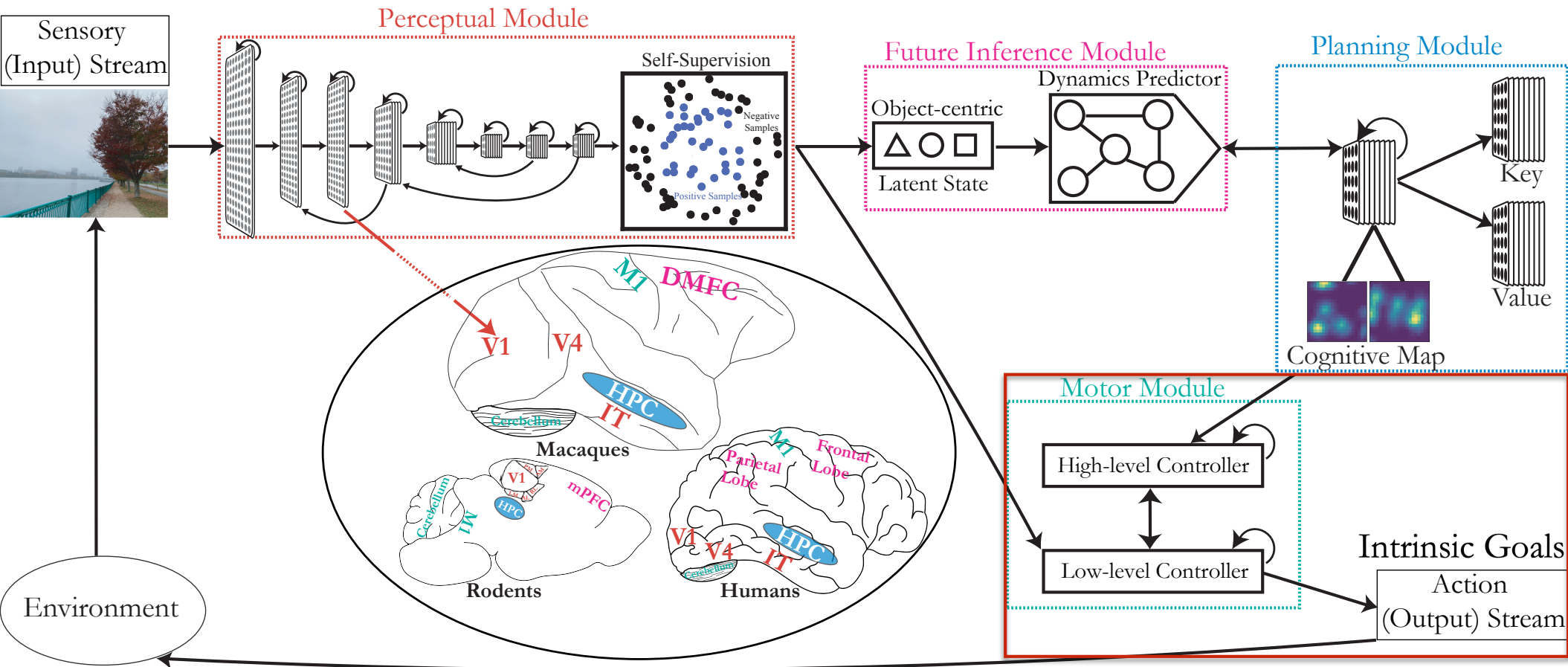


Place Cell (Hippocampus)

# Roadmap: Action

How does the brain *represent, predict, plan,* and enable **action**?

Recurrence + Contrastive SSL? Latent Future Prediction?



## **Intrinsic Goals for Autonomous Agents: Model-Based Exploration in Virtual Zebrafish Predicts Ethological Behavior and Whole-Brain Dynamics**

**Reece Keller**<sup>1,2\*</sup> **Alyn Kirsch**<sup>2</sup> **Felix Pei**<sup>1</sup> **Xaq Pitkow**<sup>1,3</sup>  
**Leo Kozachkov**<sup>4,†</sup> **Aran Nayebi**<sup>3,1,2,†</sup>

First autonomous agent that can predict whole-brain data!

NeurIPS 2025



Reece Keller



Alyn Tornell



Felix Pei



Xaq Pitkow



Leo Kozachkov<sup>†</sup>

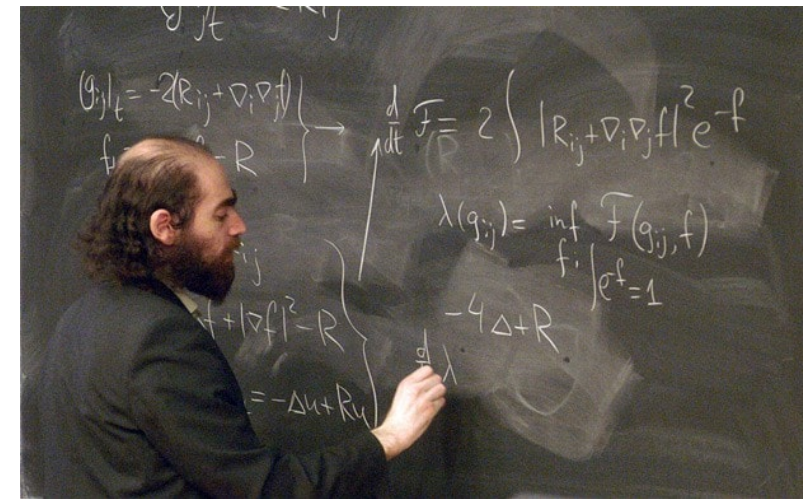
# Why is Animal Autonomy Hard?

The behavioral repertoire is enormous...

- What is the motivation/goal?
- How is it computationally formalized?
- What does "success" here even mean?

Neuroscience has largely ignored autonomous, *task-independent* behavior.

Intelligence is often attributed when goals are easily identifiable.

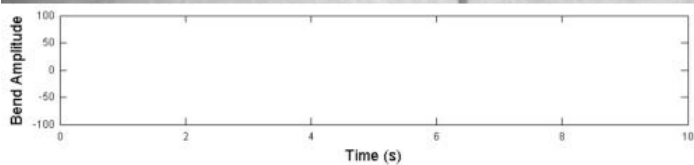
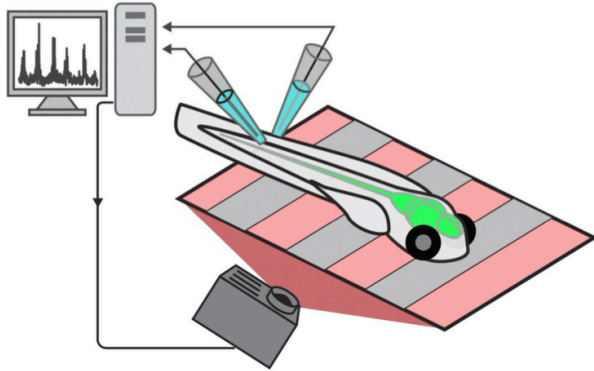


Unlike games where RL has succeeded, the environment doesn't have a dense reward function. It must be (somehow) *internally* generated by the organism!

# Glia Accumulate Evidence that Actions Are Futile and Suppress Unsuccessful Behavior

Yu Mu,<sup>1,4,\*</sup> Davis V. Bennett,<sup>1,2,4</sup> Mikail Rubinov,<sup>1,3,4</sup> Sujatha Narayan,<sup>1</sup> Chao-Tsung Yang,<sup>1</sup> Masashi Tanimoto,<sup>1</sup> Brett D. Mensh,<sup>1</sup> Loren L. Looger,<sup>1</sup> and Misha B. Ahrens<sup>1,5,\*</sup>

virtual reality navigation

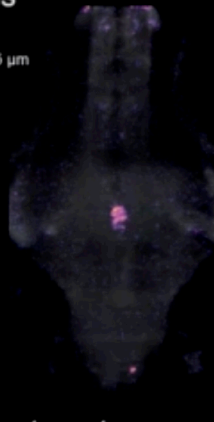


- ✓ 1. Ecologically-relevant environment
- ✓ 2. “Cognitive” states with clear behavioral readouts
- ✓ 3. Large-scale multi-area neural recordings

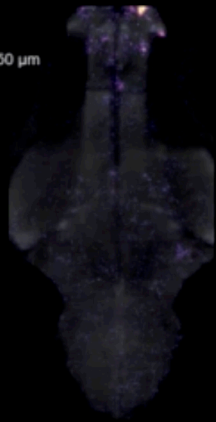


## Neurons

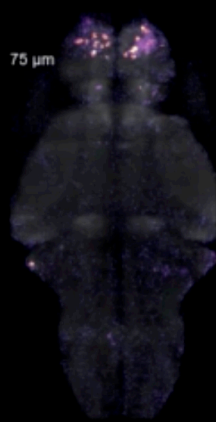
25  $\mu\text{m}$



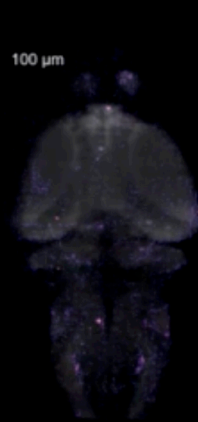
50  $\mu\text{m}$



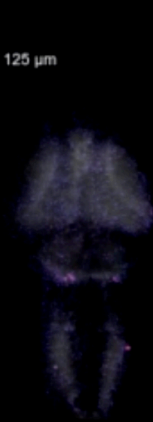
75  $\mu\text{m}$



100  $\mu\text{m}$



125  $\mu\text{m}$



## Radial astrocytes

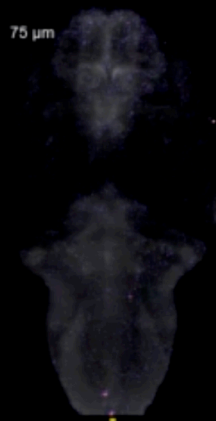
25  $\mu\text{m}$



50  $\mu\text{m}$



75  $\mu\text{m}$



100  $\mu\text{m}$



125  $\mu\text{m}$



$\Delta f/f$

0.8

0.7

0.6

0.5

0.4

0.3

0.2

0.1

Fictive behavior  
3.0X playback speed

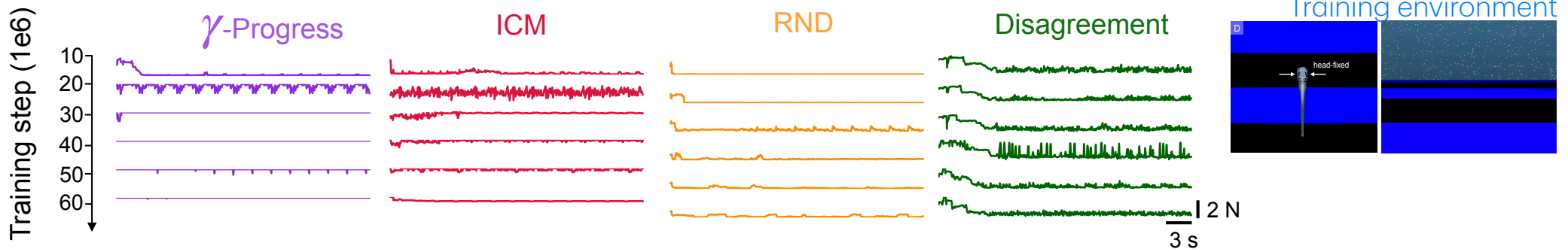
Open loop





# Epistemic Curiosity isn't Enough...

Animal autonomy != novelty optimization



## What's the issue?

- Rewards are non-stationary and saturate with experience.  
Consequence: behavioral strategies are transient  
(e.g.  $\gamma$ -Progress)
- Rewards can perseverate on unpredictable/uncontrollable stimuli.  
Consequence: unethological behavior (e.g. ICM)

## Our approach: Incorporate *priors*

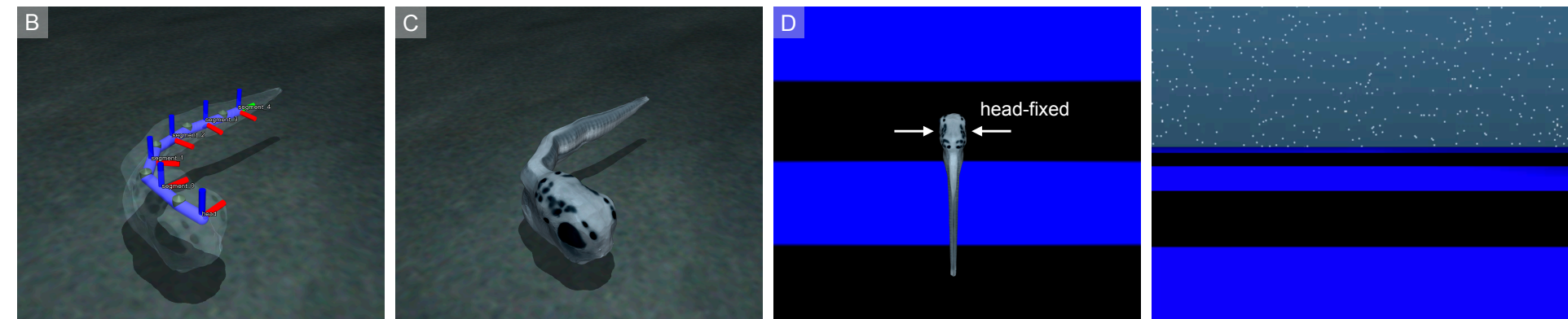
The zebrafish behavior depends on an ethological memory.

memory = fixed or slowly adapting dynamics prior (a world model!)

This enables sensorimotor feedback error to be computed and tracked.

Question: What intrinsic drive explains this behavior?  
Specifically, how should world-models be used to guide autonomous decisions in real-world situations (e.g. encountering unseen physics)?

## Zebrafish Simulation Environment



### Actuation

- The embodiment must afford a faithful comparison with the animal behavior.
- Behavioral signal is low dimensional -> embodiment can be low dimensional
- Open-source embodiments that capture basic ethology already exist!

### Sensing

- The zebrafish behavior is driven by optic flow and proprioception. A basic vision model and state information is sufficient.

### Our philosophy:

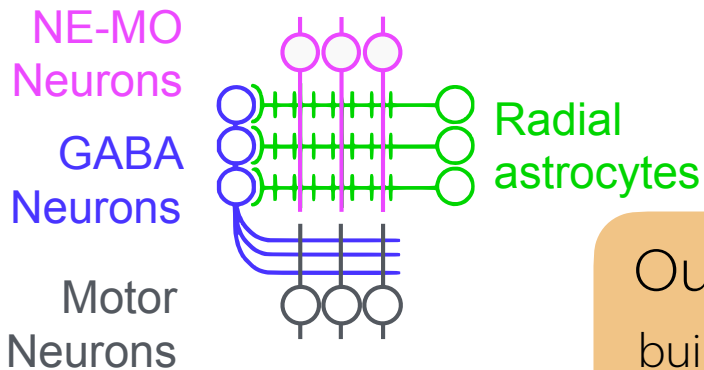
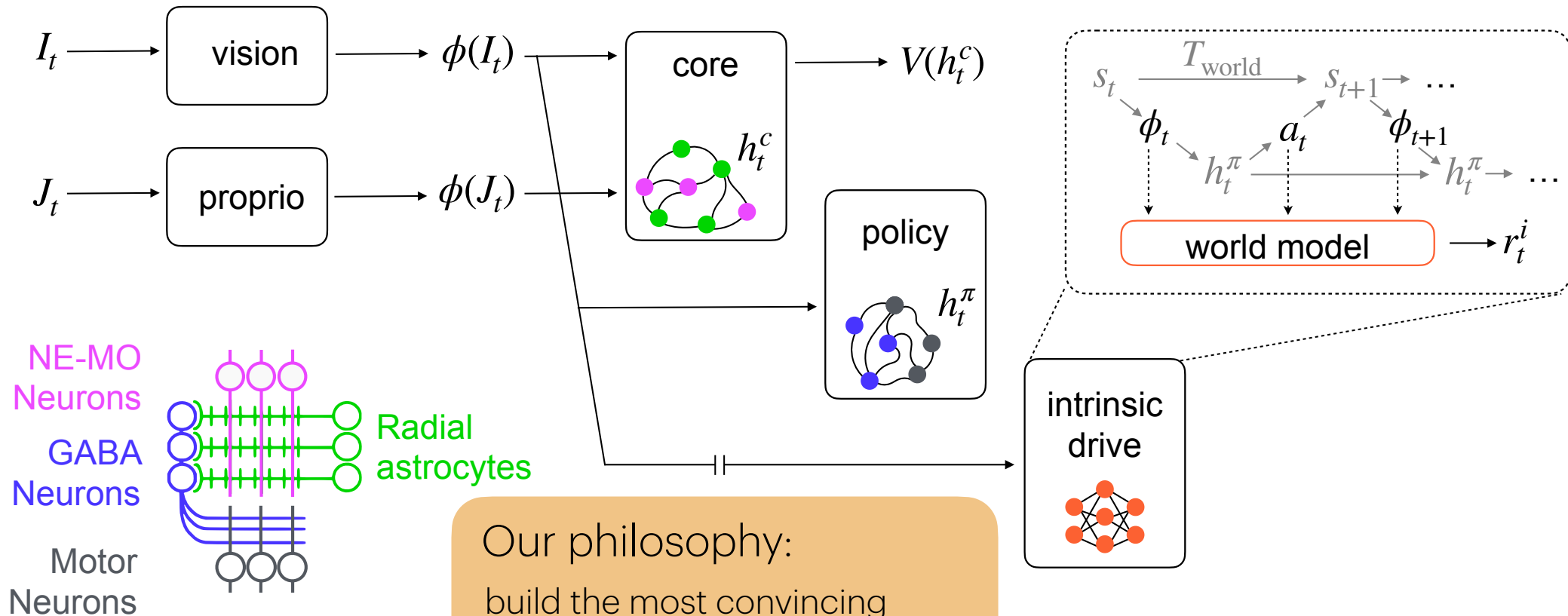
build the most convincing model possible.

- stimulus/image computable
- realistic physics
- flexible parameterization

Question: What intrinsic drive explains this behavior?

Specifically, how should world-models be used to guide autonomous decisions in real-world situations (e.g. encountering unseen physics)?

## Zebrafish Agent Architecture



Our philosophy:

build the most convincing model possible.

- stimulus/image computable
- realistic physics
- flexible parameterization

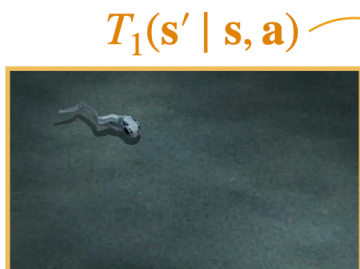
Question: What intrinsic drive explains this behavior?  
 Specifically, how should world-models be used to guide autonomous decisions in real-world situations (e.g. encountering unseen physics)?

# 3M-Progress

Recall the planning section!

Using ethological memory to guide adaptive behavior

■ ethological

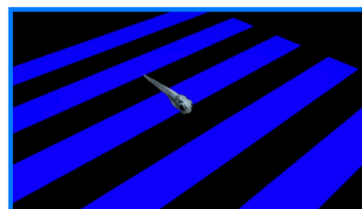
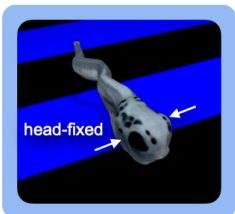


$T_1(s' | s, a)$

distill via experience

$\omega_\theta(s' | s, a)$

■ unethological

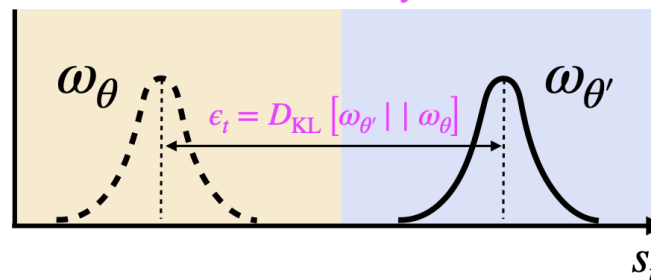


$T_2(s' | s, a)$

distill via experience

$\omega_{\theta'}(s' | s, a)$

3M: Model-Memory-Mismatch

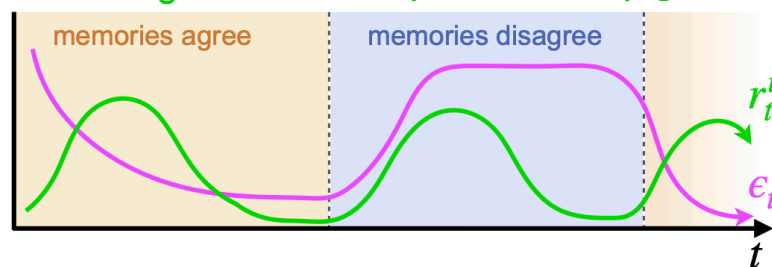


$\epsilon_t$  partitions the state-action space into model-memory agreement ( $U$ ) and disagreement ( $U^C$ ).

3M-Progress

$$r_t^i \propto |\hat{\epsilon}_t - \epsilon_t|$$

$$\hat{\epsilon}_t = (1 - \gamma)\hat{\epsilon}_{t-1} + \gamma\epsilon_t$$



We choose  $T_1$  and  $T_2$  to obey:

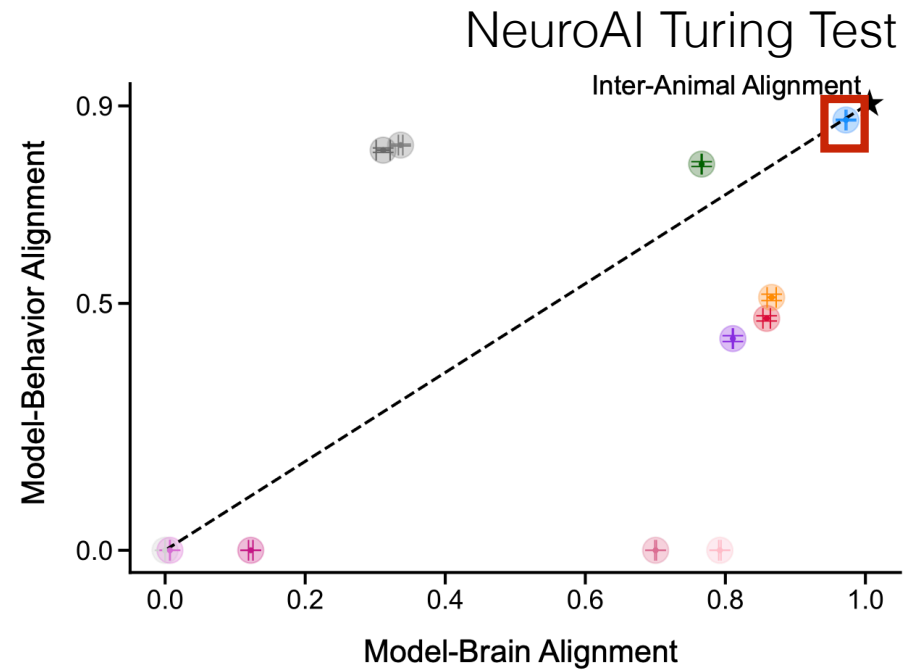
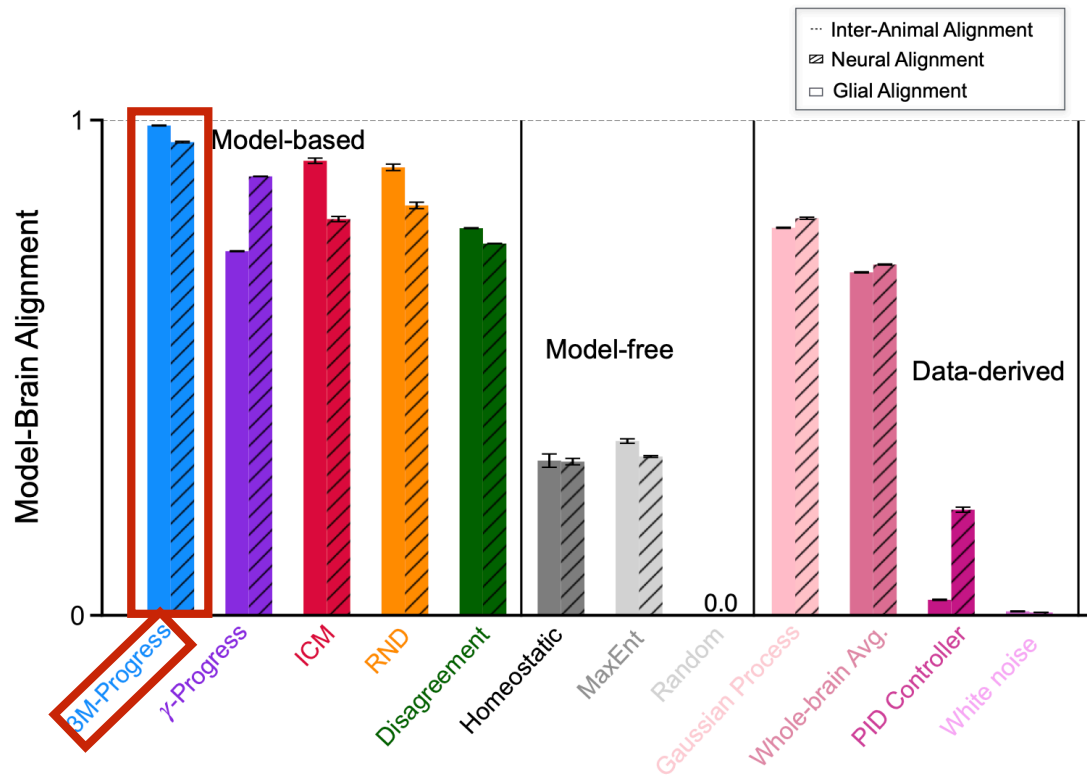
$$\exists U \subset S \times A \text{ s.t. } \forall (s, a) \in U, T_1 \approx T_2$$

(dynamics agree on a subspace).

# Putting it all together

## 3M-Progress Captures Whole-Brain Dynamics (and behavior)

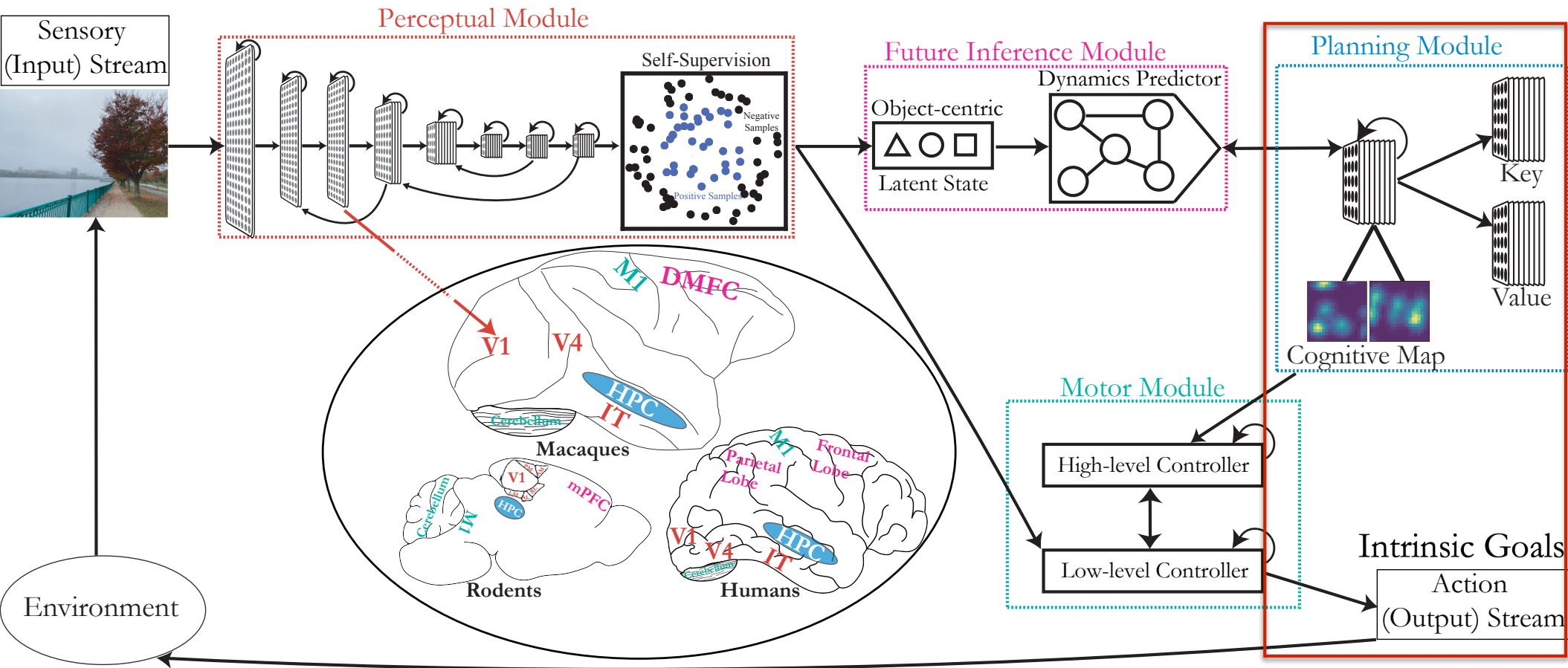
Single-cell one-to-one alignment



# Roadmap: Planning & Action

How does the brain *represent*, *predict*, *plan*, and enable *action*?

Recurrence + Contrastive SSL? Latent Future Prediction?



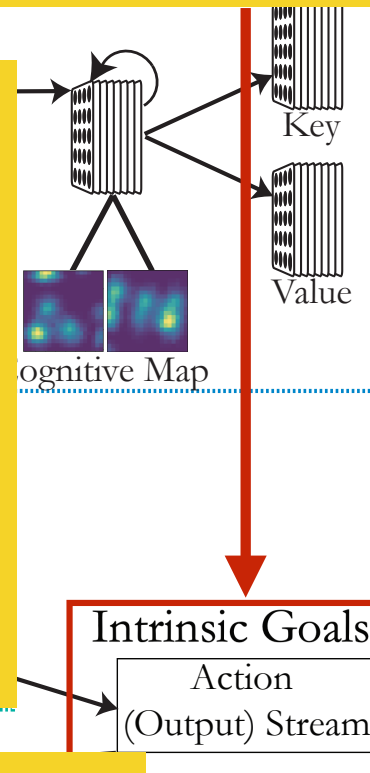
Temporal integration of World Model-Progress-based curiosity?

# Safety Implications: What Happens Once We Get There?

How does the brain *represent, predict, plan, and enable action*?

Too many of these goals makes alignment *intractable, even for computationally unbounded agents!*

One can guarantee “*corrigibility*”, where under the *optimal* agent policy, humans retain control. Involves only a small set of modular & lexicographically organized goals (paralleling the modular agent architecture), circumventing the barrier above.



Open: Can we scale corrigibility cost effectively? Curiosity?

1. Intrinsic Barriers and Practical Pathways for Human-AI Alignment: An Agreement-Based Complexity Analysis

AAAI '26: <https://arxiv.org/abs/2502.05934>

2. Core Safety Values for Provably Corrigible Agents

AAAI '26: <https://arxiv.org/abs/2507.20964>

# Potential Economic Implications of Alignment

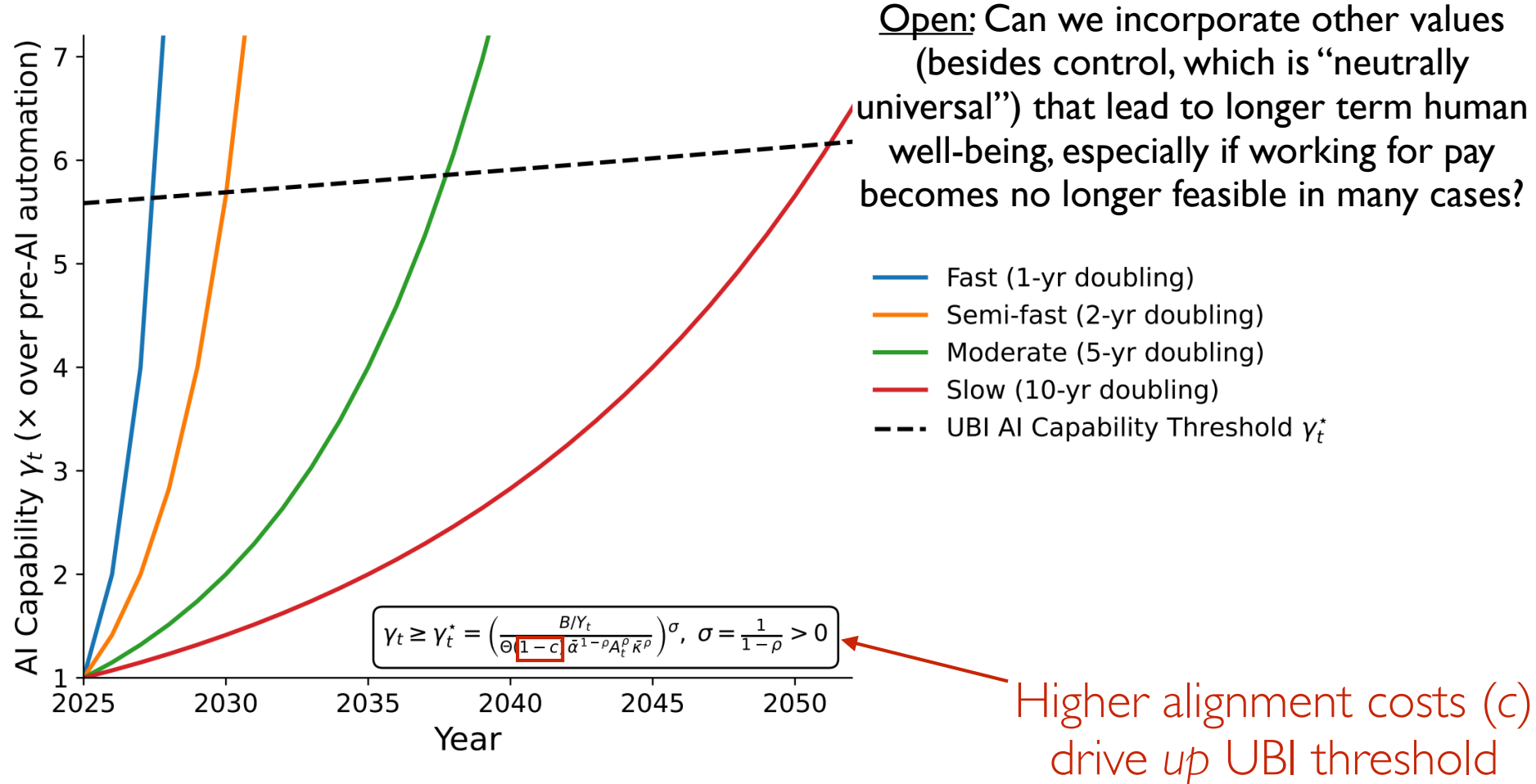
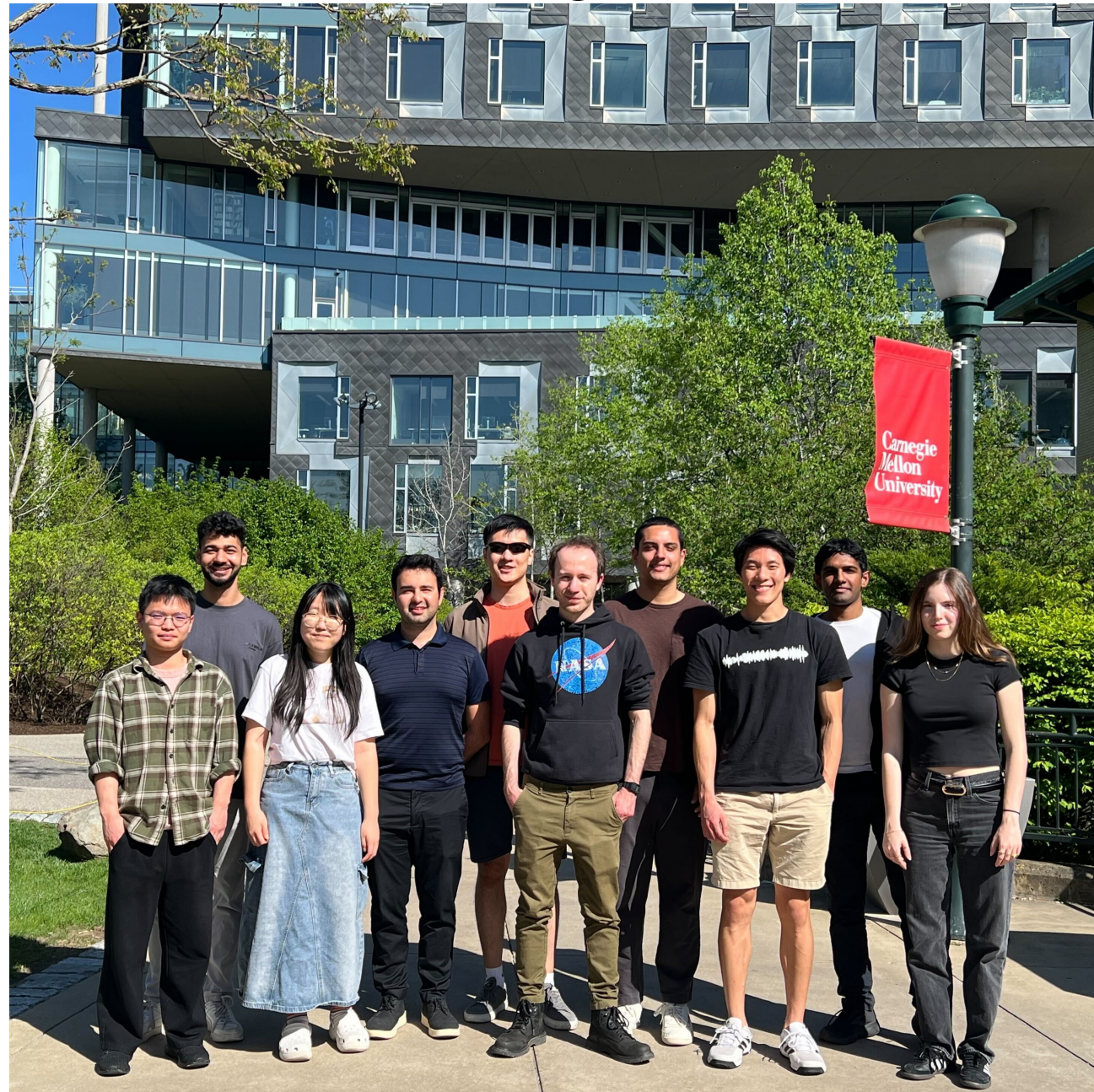


Figure 1: **Projected AI capabilities ( $\gamma_t$ ) vs. time-varying UBI AI capability threshold ( $\gamma_t^*$ ).** The dashed line is the required capability  $\gamma_t^*$  to fully fund a UBI that comprises 11% of the GDP (leading to a  $\gamma_t^*$  between 5-6 $\times$  the pre-AI productivity on automated tasks, under current economic assumptions). Under fast scaling (AI capability doubling every year), AI would cross the threshold by the late 2020s. Semi-fast scaling (doubling every 2 years) reaches the threshold in the early 2030s, whereas moderate (doubling every 5 years) and slow (doubling every 10 years) scenarios achieve  $\gamma_t^*$  by 2038 and 2052, respectively. The trajectories are illustrative, starting from a nominal, conservative 2025 capability level ( $\gamma_0 \equiv 1$ ), which assumes AI currently delivers no boost beyond the pre-AI automation level in aggregate across all automated tasks.

# Acknowledgements

## NeuroAgents Lab



### Contact:



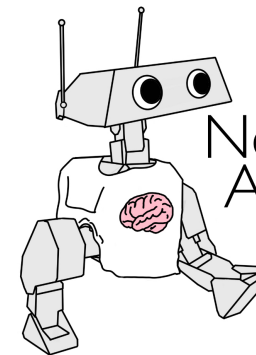
[anayebi@cs.cmu.edu](mailto:anayebi@cs.cmu.edu)



[@aran\\_nayebi](https://twitter.com/aran_nayebi)



<https://cs.cmu.edu/~anayebi>



Neuro  
Agents  
Lab



Carnegie Mellon  
SCHOOL OF COMPUTER SCIENCE

### Funding:

Foresight Institute

UK AISI Challenge Fund

Google Robotics Award

Burroughs Wellcome Fund CASI Award