# Using *Embodied* Agents to Reverse-Engineer *Natural* Intelligence
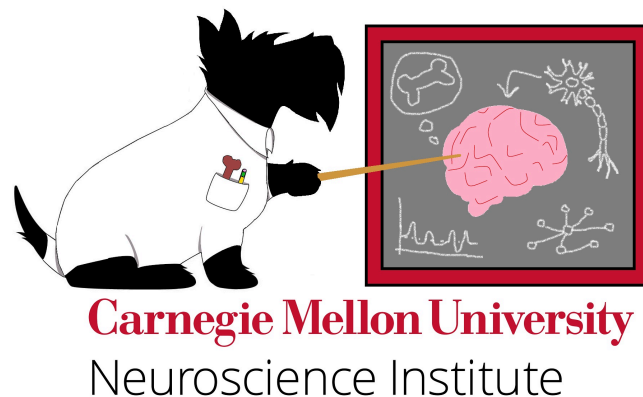
## Aran Nayebi

*Assistant Professor*
*Machine Learning Department*
*Neuroscience Institute (core faculty), Robotics Institute (by courtesy)*

## RI Seminar

*2025.09.26*

# Current AI Struggles to Understand the Physical World

*OpenAI Sora, February 2024*

Q: What's missing?

*OpenAI Sora,*
*February 2024*

Q: What's missing?

A: Embodied agency & interaction.

**Why?**

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

## 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!)
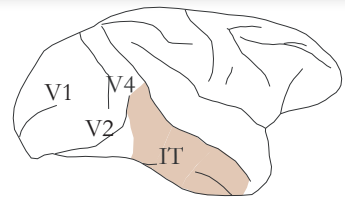
Why?
Animals & humans (currently) perform behaviors we've yet to engineer successfully in AI agents:
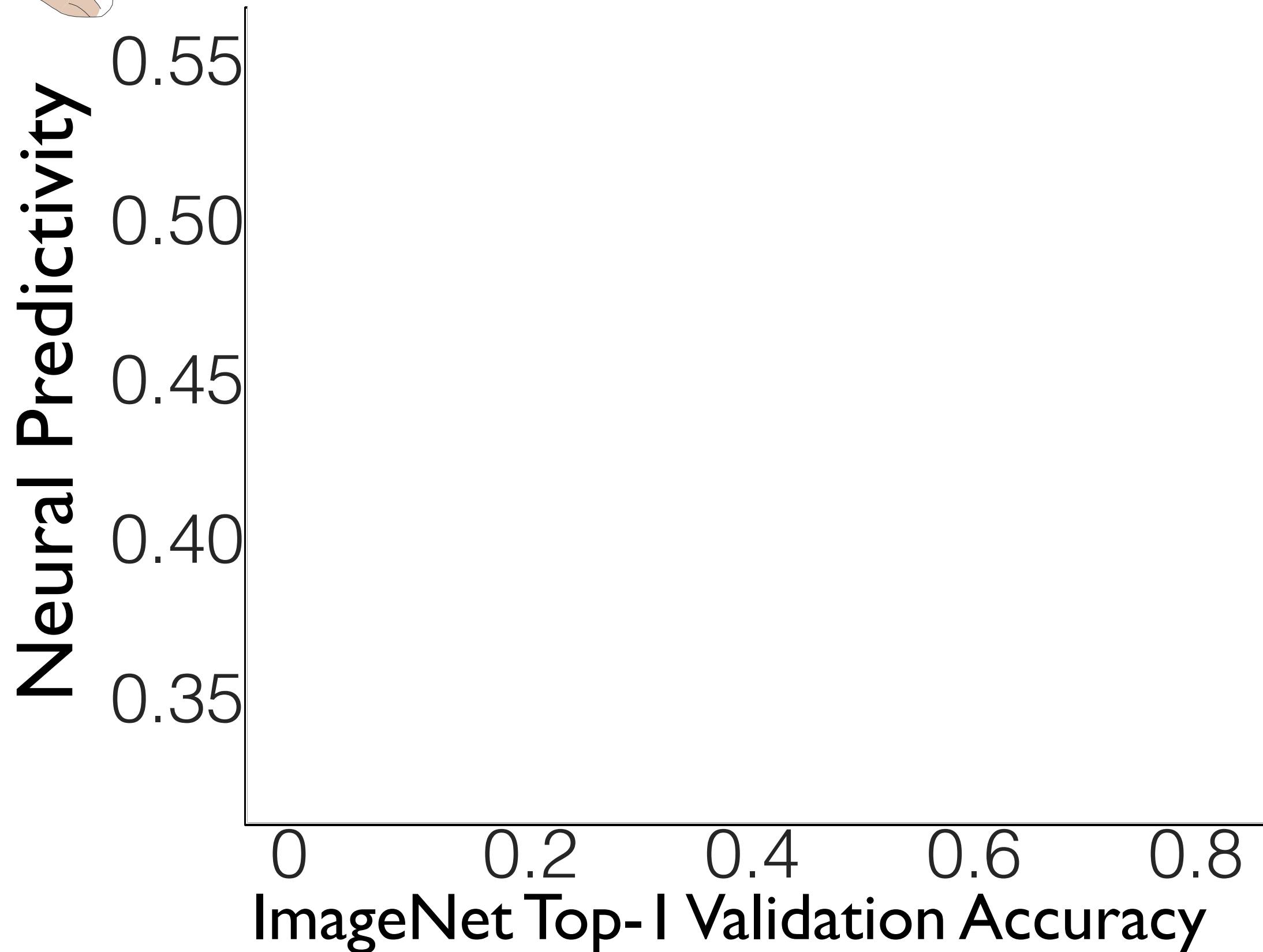
‣ Prediction (requires world modeling) & planning (requires memory)

‣ 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

Neural Predictivity (y-axis): 0.35, 0.40, 0.45, 0.50, 0.55

ImageNet Top-1 Validation Accuracy (x-axis): 0, 0.2, 0.4, 0.6, 0.8

# Task Performance Correlated with Neural Predictivity

*Schrimpf*, Kubilius* et al. 2018*

$R = 0.92$

Neural Predictivity
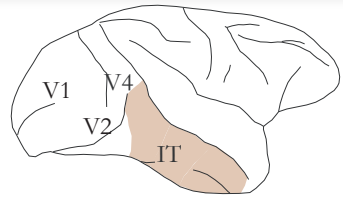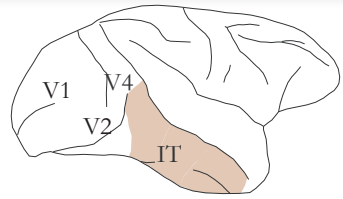
ImageNet Top-1 Validation Accuracy

# Task Performance Correlated with Neural Predictivity

*Schrimpf\*, Kubilius\* et al. 2018*

$R = 0.92$

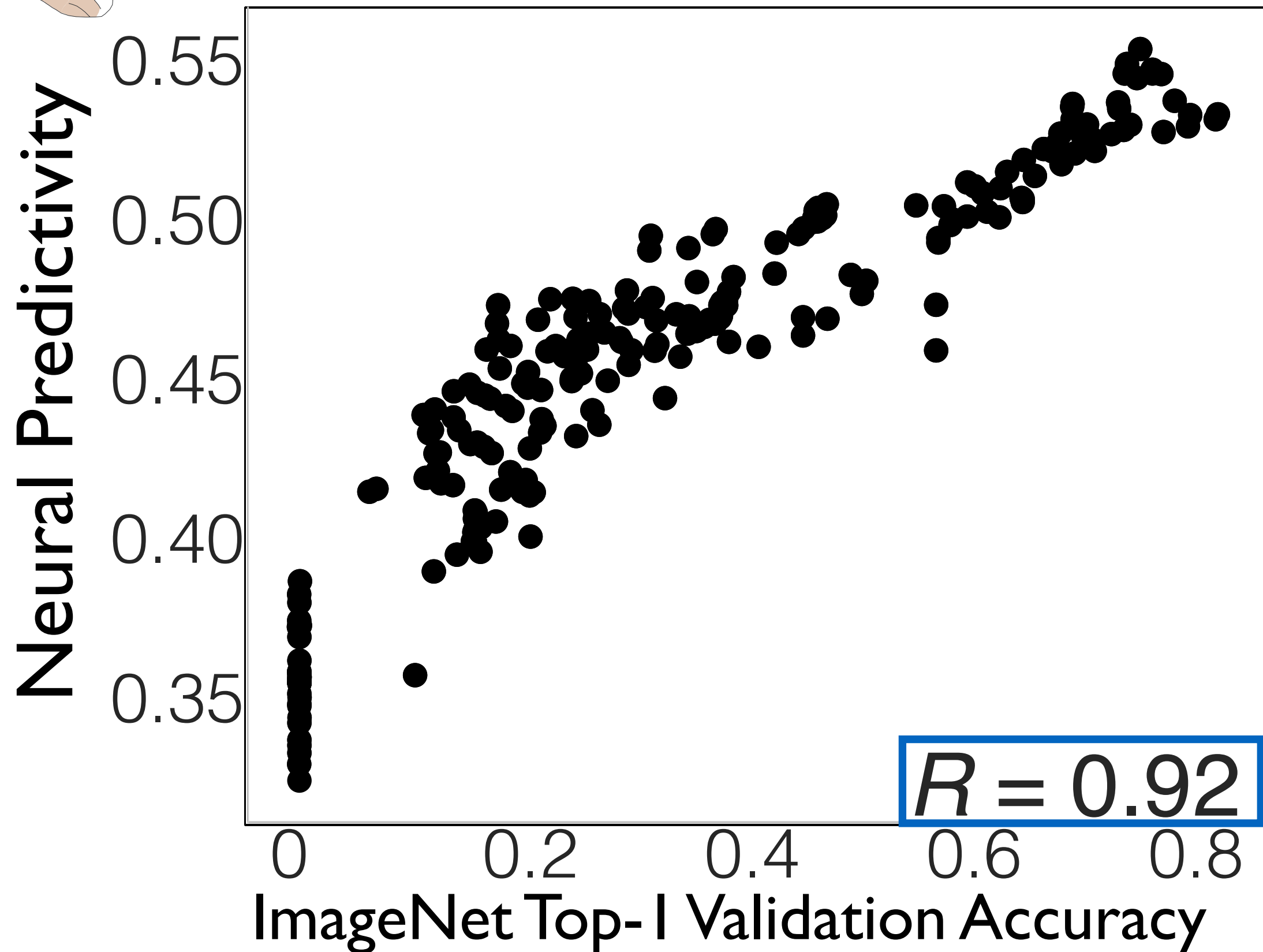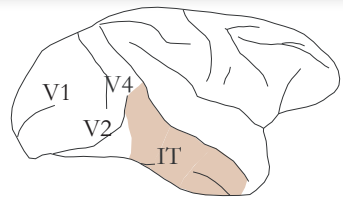# Task Performance Correlated with Neural Predictivity

A Neuroscience Goal

*Schrimpf\*, Kubilius\* et al. 2018*

An AI Goal

## Task-Optimization (ML)

**1.**

$A$ = *architecture class*

**2.**

$T$ = *task loss*

**3.**

$D$ = *dataset*

**4.**

$L$ = *learning rule*

## Task-Optimization (ML)

**1.**

$A$ = *architecture class*

**2.**

$T$ = *task loss*

**3.**

$D$ = *dataset*

**4.**

$L$ = *learning rule*

## Task-Optimization (ML)

<span style="color:red">**Neurobiology**</span>

**1.**

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



**2.**

$T$ = *task loss*

**3.**

$D$ = *dataset*

**4.**

$L$ = *learning rule*

## Task-Optimization (ML)

## Neurobiology

**1.**



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

**2.**



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

**3.**

$D$ = *dataset*

**4.**

$L$ = *learning rule*

## Task-Optimization (ML)

## Neurobiology

**1.**

| $A$ = architecture class **= circuit neuroanatomy** |
|---|



**2.**

| $T$ = task loss **= ecological niche/behavior** |
|---|



**3.**

| $D$ = dataset **= environment** |
|---|



**4.**

| $L$ = learning rule |
|---|

## Task-Optimization (ML)

## Neurobiology

**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** |

**L** = *learning rule*

***"Natural selection + plasticity"***

**T** = *task loss*

***"Ecological niche/ behavior"***

***"Environment"***

**D** = *data stream*

***"Circuit"***

**A** = *architecture class*

# Task-Optimized Modeling: Four Components

$L$ = learning rule

**"Natural selection + plasticity"**

$T$ = task loss

**"Ecological niche/ behavior"**

**"Environment"**

$D$ = data stream

**"Circuit"**

$A$ = architecture class

# Task-Optimized Modeling: Four Components

**L** = *learning rule*

**"Natural selection + plasticity"**

**T** = *task loss*

**"Ecological niche/ behavior"**

**"Environment"**

**D** = *data stream*

**"Circuit"**

**A** = *architecture class*

# Task-Optimized Modeling: Four Components



**L** = learning rule

**"Natural selection + plasticity"**

**T** = task loss

**"Ecological niche/ behavior"**

**"Environment"**
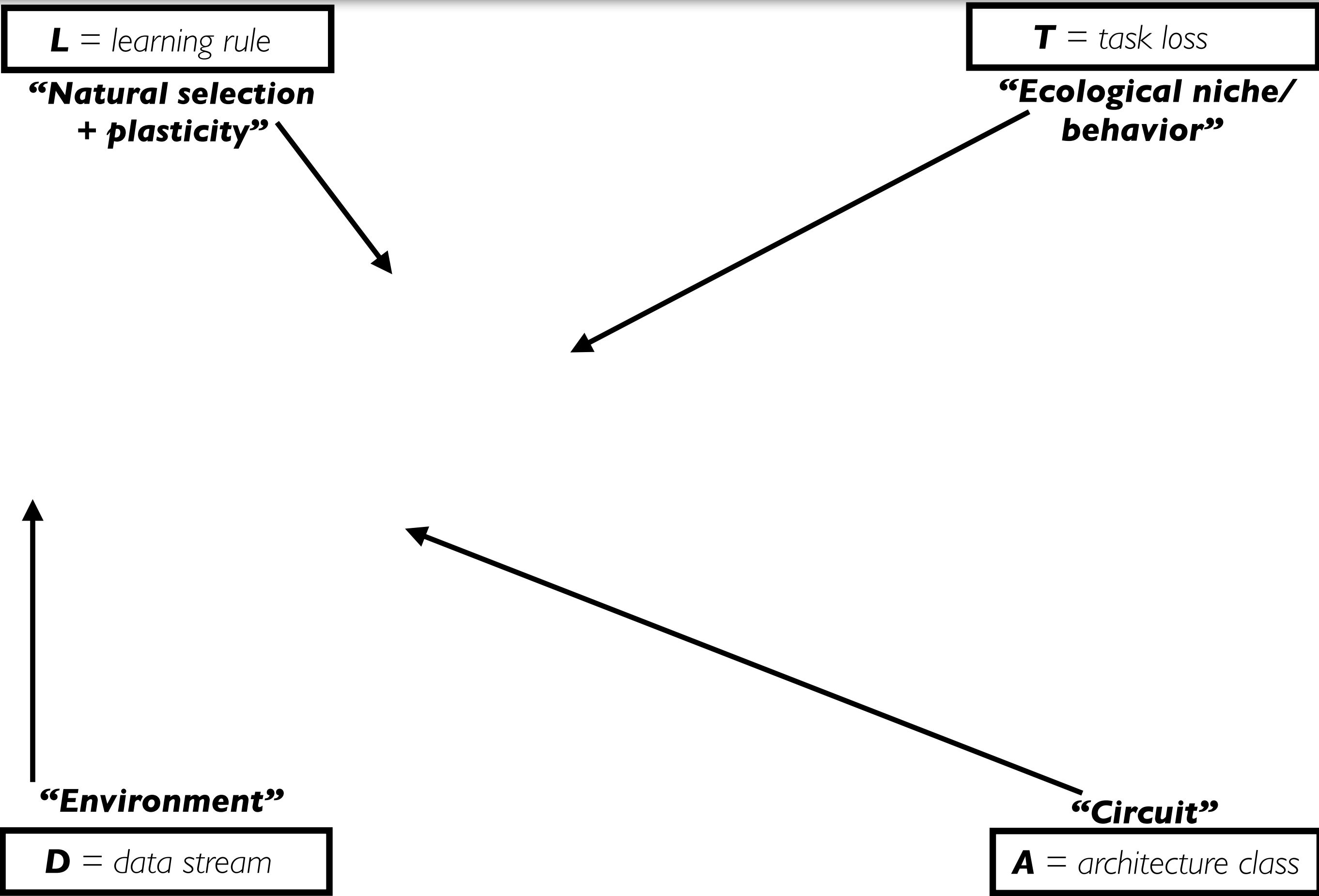
**D** = data stream

**"Circuit"**

**A** = architecture class

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

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



**Yields:**

# Task-Optimized Modeling

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



**Yields:**

Quantitatively Accurate & Practically Useful Brain Models

# 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:**

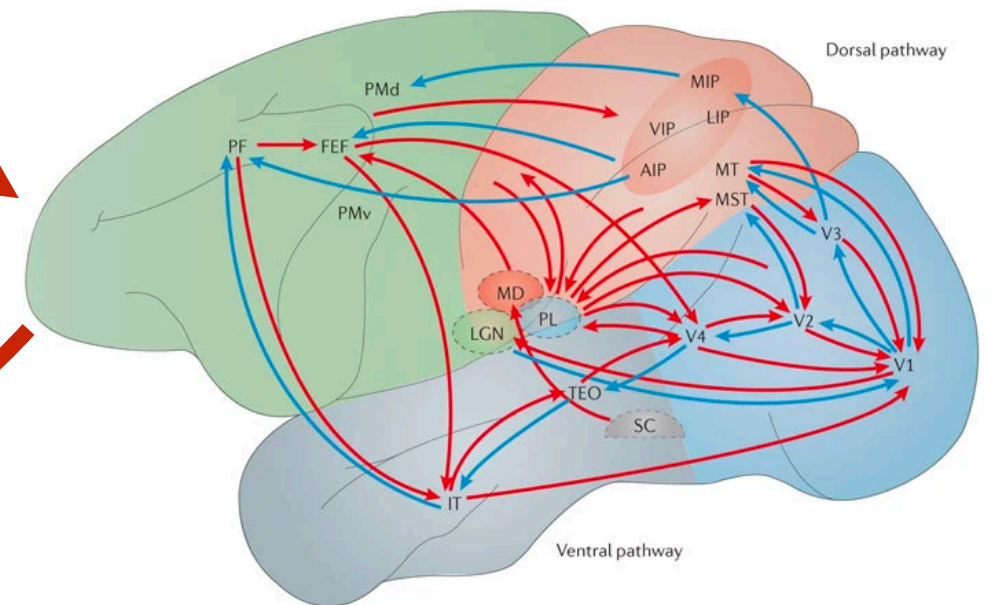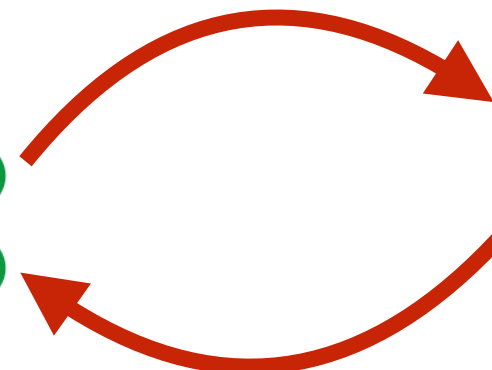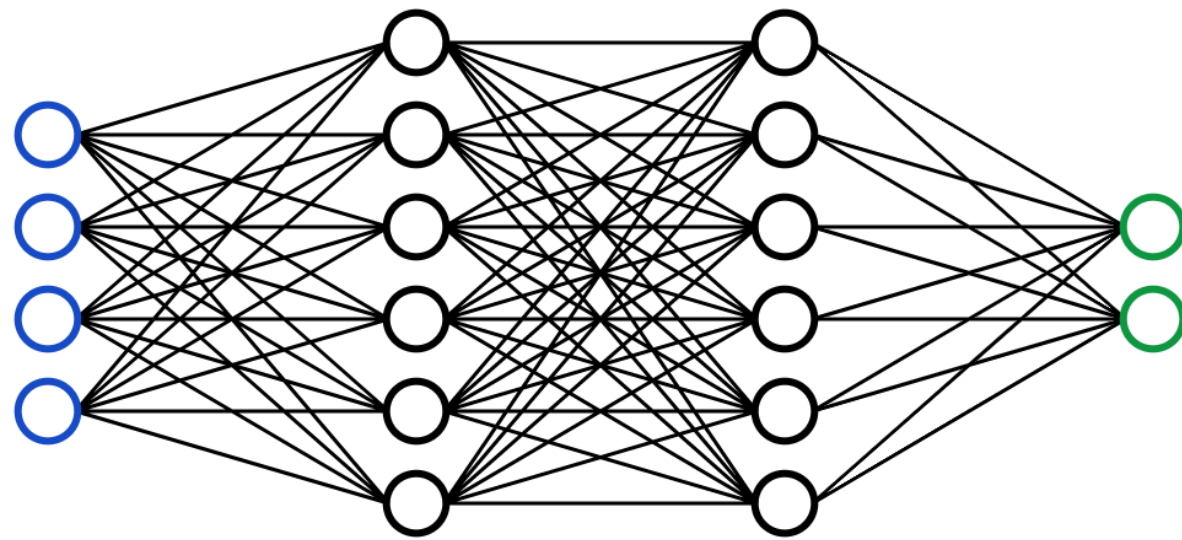Quantitatively Accurate & Practically Useful Brain Models

*AND*

Principles of *Why* Neural Responses Are As They Are

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

**Dispersion of Solution Set** (vertical axis)

**Constraint Strength** (horizontal axis)

Easy Task

Medium Task

Difficult Task

phenotypic space

Human
Macaque
Mouse
Neural Network

Rosa Cao

Daniel Yamins

Explanatory models in neuroscience:

Part 2 – Constraint-based intelligibility

*Figure 6.* **The Multitask Scaling Hypothesis:** Models trained with an increasing number of tasks are subjected to pressure to learn a representation that can solve all the tasks.

**The Platonic Representation Hypothesis**

Minyoung Huh [*1]   Brian Cheung [*1]   Tongzhou Wang [*1]   Phillip Isola [*1]

**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.

> "Nothing in biology makes sense in light of evolution."
> *- Theo Dobzhansky*

Dispersion of Solution Set

Easy Task

phenotypic space

Neural Network

Medium Task

phenotypic space

Difficult Task

phenotypic space

Constraint Strength

Solves task 1

Solves task 2

task gradient

*Figure 6.* **The Multitask Scaling Hypothesis:** Models trained with an increasing number of tasks are subjected to pressure to learn a representation that can solve all the tasks.

**The Platonic Representation Hypothesis**

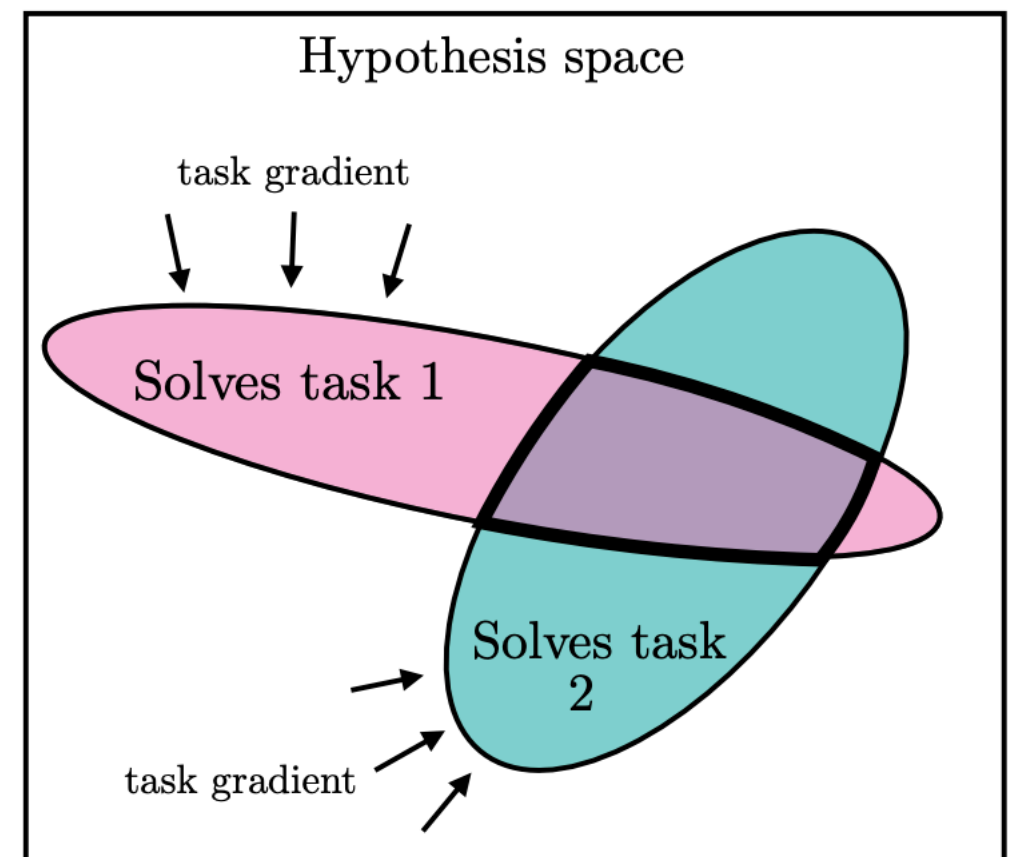Minyoung Huh[*1]  Brian Cheung[*1]  Tongzhou Wang[*1]  Phillip Isola[*1]

**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.

"Nothing in biology makes sense in light of evolution."
- *Theo Dobzhansky*

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

Easy Task

phenotypic space

Neural Network

Difficult Task

phenotypic space

Constraint Strength

Dispersion of Solution Set
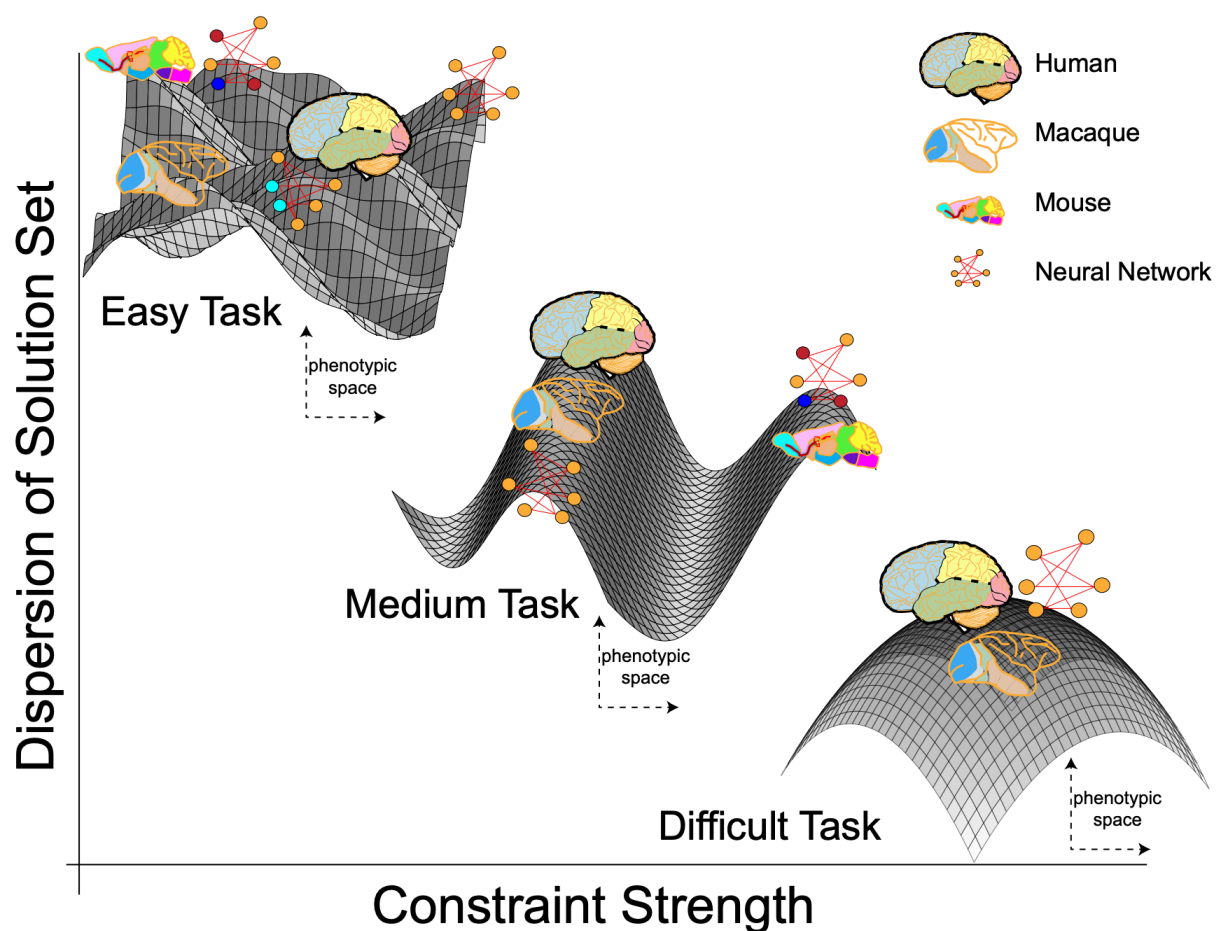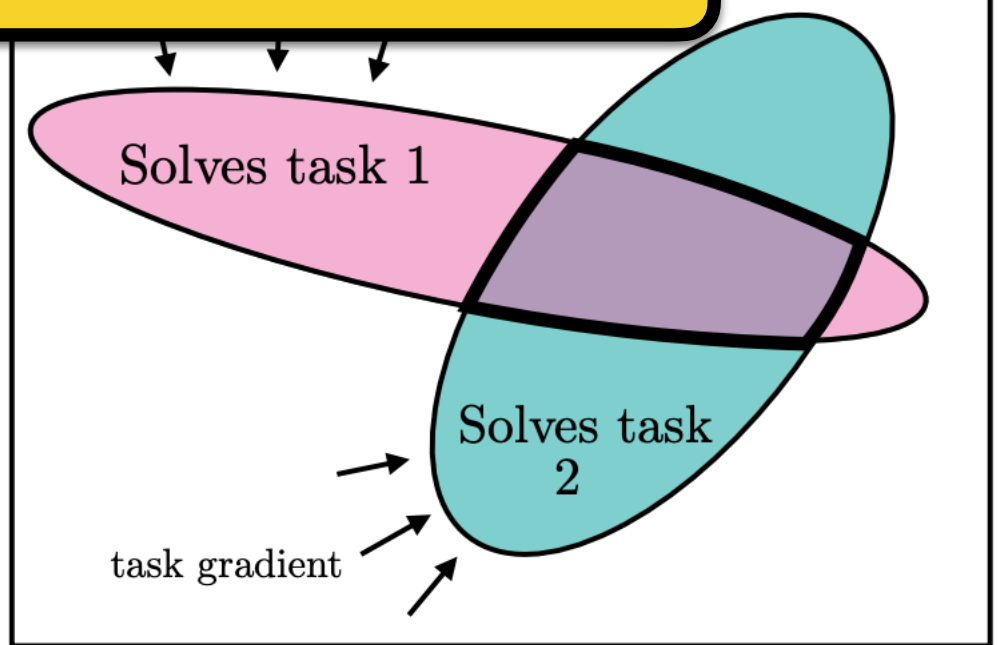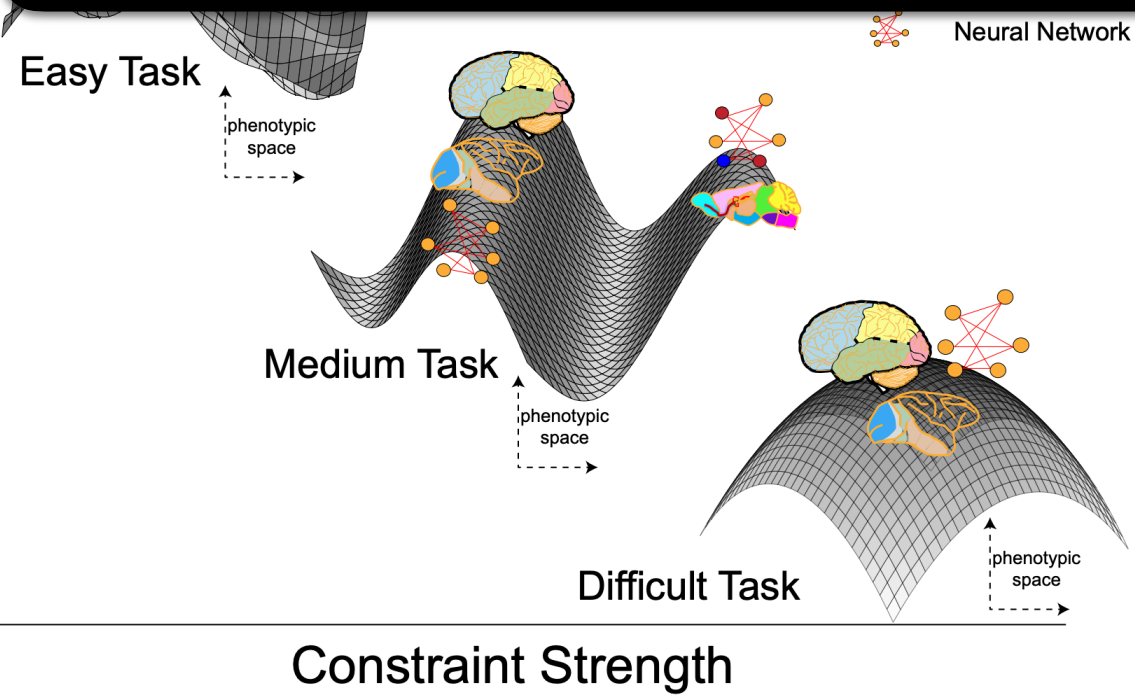
Solves task 1

task gradient

*Figure 6.* **The Multitask Scaling Hypothesis:** Models trained with an increasing number of tasks are subjected to pressure to learn a representation that can solve all the tasks.

**The Platonic Representation Hypothesis**

Minyoung Huh [*1]   Brian Cheung [*1]   Tongzhou Wang [*1]   Phillip Isola [*1]

**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.

Dispersion of Solution Set

Easy Task

phenotypic space

Neural Network

Solves task 1

task gradient

Difficult Task

phenotypic space

Constraint Strength

**"Nothing in biology makes sense in light of evolution."**
*- Theo Dobzhansky*

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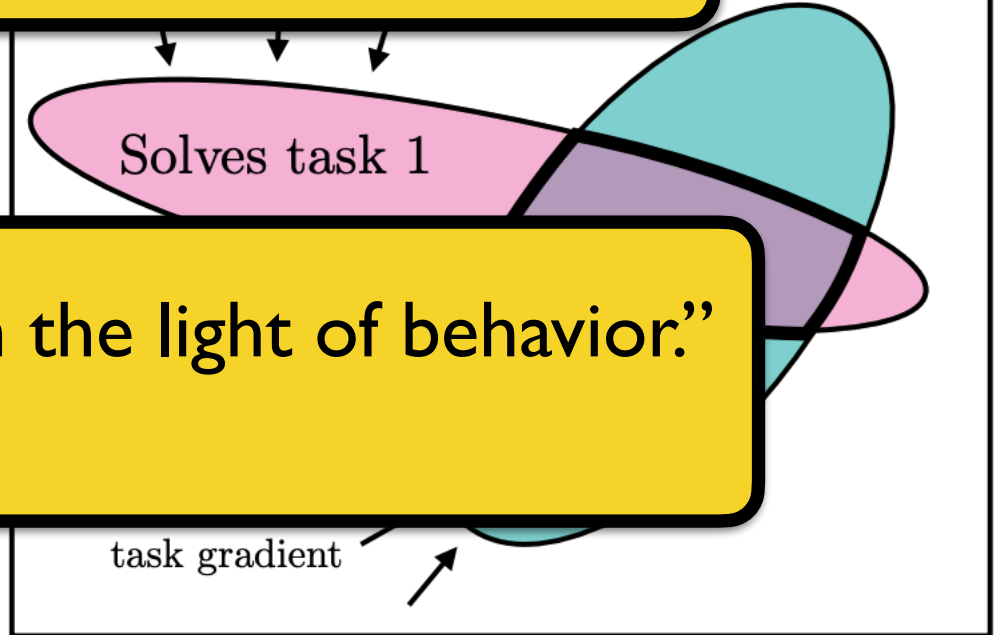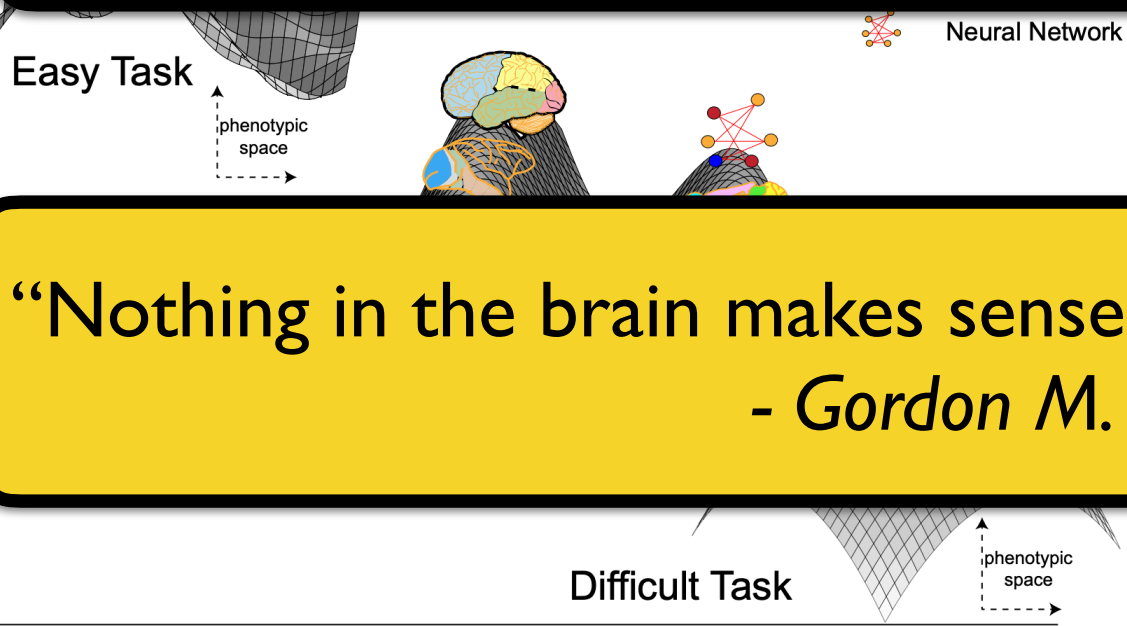*- Gordon M. Shepherd*

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

**The Platonic Representation Hypothesis**

Minyoung Huh [*1] Brian Cheung [*1] Tongzhou Wang [*1] Phillip Isola [*1]

*Figure 6* **The Multitask Scaling Hypothesis:** Models trained ...ted to pressure to ...s.

**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

Artificial Neural Network

Brain

**Yields:**

Quantitatively Accurate & Practically Useful Brain Models

*AND*

Principles of *Why* Neural Responses Are As They Are

# 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

Brain

**Yields:**

Quantitatively Accurate & Practically Useful Brain Models

*AND*

Principles of *Why* Neural Responses Are As They Are

**Jenelle Feather**

**Meenakshi Khosla**

**Ratan Murty**

# Brain-Model Evaluations Need the NeuroAI Turing Test

**Jenelle Feather** [*1]   **Meenakshi Khosla** [*2]   **N. Apurva Ratan Murty** [*3]   **Aran Nayebi** [*4]

Jenelle Feather



Meenakshi Khosla



Ratan Murty

Jenelle Feather

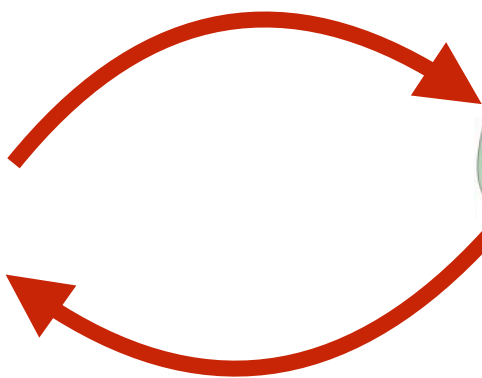**Just as distinct objects can cast the same shadow...**

Meenakshi Khosla

Ratan Murty

**Turing Test**

*human-to-human similarity*

*compare with human distribution*

**behavior only**

# NeuroAI Turing Test

Jenelle Feather

Meenakshi Khosla

Ratan Murty

**Just as distinct objects can cast the same shadow...**

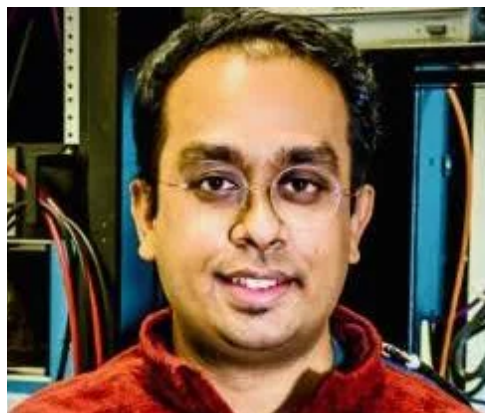**...distinct internal processes (representations) can produce similar outputs (behavior)**

**Turing Test**

human-to-human similarity

compare with human distribution

**behavior only**

**The NeuroAI Turing Test**

human-to-human similarity

brain-to-brain similarity

compare with human distribution

compare with brain distribution

**behavior** + **internal representations**

# *How* to Reverse-Engineer Natural Intelligence?

# *How* to Reverse-Engineer Natural Intelligence?

# Whole brain…

# Whole brain…

…awake, behaving animals

# Whole brain…

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

…awake, behaving animals

# 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.**

*International Brain Laboratory 2022*

…awake, behaving animals

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

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

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

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

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

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

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

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

Sensory (Input) Stream

Perceptual Module

Self-Supervision

Negative Samples

Positive Samples

Future Inference Module

Object-centric

Dynamics Predictor

Latent State

Planning Module

Key

Value

Cognitive Map

PM

AM

V1

LM  AL  RL

HPC

mPFC

Cerebellum

M1

Rodent Brain

Motor Module

High-level Controller

Low-level Controller

Intrinsic Goals

Action (Output) Stream

Environment

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

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

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

Effective proxy, but just obviously deeply wrong.

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

Training Input



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

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

# Contrastive learning tasks

Training Input



Encoder

Embedding

*Further*

Closer

*Further*

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

**High-level idea of these methods: make the representations <span style="color:red">non-trivially robust to data augmentations</span>**

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

# Comparison to Neural Data

## How well does it match neural data?

Chengxu Zhuang



**Pretrained DCNN**

**Test Per-Site Neural Predictions**

Test Input

100ms Visual Presentation

Neural Recordings from V1, V4, and IT

V1

V2

V4

IT

100ms

Image Present

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)

Test Input

Pretrained DCNN

Test Per-Site Neural Predictions

Neural Recordings from V1, V4, and IT

100ms Visual Presentation

V1    V4    V2    IT

Image Present

100ms

Autoencoders

Missing-Data Tasks

Deep Contrastive Embeddings

Chengxu Zhuang

## Brain-Score

Neural Fit

0.8

0.6

0.4

Untrained | Auto-Encoder | PredNet | Depth Pred. | CPC | Colorization | Relative Pos. | CMC | Deep Cluster | Instance Recog. | SimCLR | Local Aggregation | Categorization

Supervised

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

Zhuang C, Yan S, Nayebi A, Schrimpf M, Frank M, DiCarlo JJ, & Yamins D (2021). Unsupervised Neural Network Models of the Ventral Visual Stream. *(PNAS)*

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.**

Zhuang C, Yan S, Nayebi A, Schrimpf M, Frank M, DiCarlo JJ, & Yamins D (2021).
Unsupervised Neural Network Models of the Ventral Visual Stream. *(PNAS)*

# 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

**Primates**

**Mouse**

Mouse vision is less hierarchical!

# Contrastive Models Better Match Mouse Visual Cortex



Primates

Mouse

Mouse vision is less hierarchical!

Contrastive Models Better Match Mouse Visual Cortex

Contrastive Models Better Match Mouse Visual Cortex

# 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.**

~90% of the NeuroAI Turing Test for this dataset

Biomechanical Model

# Train

- - - - - - - - - - - - - - - - - - - - - - -

*ImageNet*

Biomechanical Model

## Train

## Evaluate

*ImageNet*

*Reward-Based Navigation*

Train      Evaluate

*ImageNet*      *Reward-Based Navigation*

CSNet

Critic Head $V(s_t)$

Policy Head $a_t$

Biomechanical Model

Train

Evaluate

*ImageNet*

*Reward-Based Navigation*

CSNet

Critic Head

Policy Head

Biomechanical Model

Train

Evaluate

*ImageNet*

*Reward-Based Navigation*

## Vision Network

Train

Evaluate

*ImageNet*

*Reward-Based Navigation*

# Embodied Virtual Rodent Navigation

## Vision Network

### Biomechanical Model
(Joint angles, accelerometer, etc.)

**Bence Ölveczky**

Train

Evaluate

*ImageNet*

*Reward-Based Navigation*

CSNet

Critic Head

Policy Head

Biomechanical Model

Embodied Virtual Rodent Navigation

Vision Network

Biomechanical Model
(Joint angles, accelerometer, etc.)

**Bence Ölveczky**

# Embodied Virtual Agent Navigation

## Vision Network

## Decision Making

$V(s_t)$

Critic

Actor $a_t$

## Biomechanical Model
(Joint angles, accelerometer, etc.)

**Bence Ölveczky**

# Embodied Virtual Agent Navigation

Vision Net

Decision Making

$V(s_t)$

Critic

Actor $a_t$

Biomechanical Model
(Joint angles, accelerometer, etc.)

**Bence Ölveczky**

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

Decision Making

Critic $\rightarrow V(s_t)$

Actor $\rightarrow a_t$

Biomechanical Model
(Joint angles, accelerometer, etc.)

Bence Ölveczky

# Contrastive Models Yield Better Transfer Performance



Reward-Based Navigation

Pose Estimation

Position Estimation

CSNet

Critic Head → $V(s_t)$

Policy Head → $a_t$

Biomechanical Model

Texture Classification

**Train**

*ImageNet*

**Evaluate**

Vision Network

Decision Making

Critic Head → $V(s_t)$

Policy Head → $a_t$

# Contrastive Models Yield Better Transfer Performance



Reward-Based Navigation

Pose Estimation

**Train**

**Evaluate**

*Maze Environment*

*Visual Scene Understanding*

*Object properties*

*Texture*

# Contrastive Models Yield Better Transfer Performance

Biomechanical Model

## Reward-Based Navigation

## Pose Estimation

## Position Estimation

**Train**

**Evaluate**

*Maze Environment*

*Visual Scene Understanding*

*Object properties*

*Texture*

# Contrastive Models Yield Better Transfer Performance

Biomechanical Model

## Reward-Based Navigation

## Pose Estimation

## Position Estimation

## Size Estimation

**Evaluate**

*Visual Scene Understanding*

Plane — Category

rz — z axis rotation

rx, ry — x axis rotation, y axis rotation

f16 — Identity

Horizontal position: 80 pix

Vertical position: −6 pix

Perimeter: 78 pix

Two-dimensional retinal area: 146 pix

Three-dimensional object scale: 1.2×

*Object properties*

*Texture*

Contrastive Models Yield Better Transfer Performance

# Contrastive Models Yield Better Transfer Performance

Contrastive Models Yield Better Transfer Performance

# Contrastive Models Yield Better Transfer Performance

**Reward-Based Navigation**

Mean Episode Return

300

260

220

Contrastive ImageNet | Supervised ImageNet

**Pose Estimation**

Pearson's R

0.13

0.08

0.03

Contrastive Maze | Supervised Maze

**Position Estimation**

0.6

0.4

0.2

Contrastive Maze | Supervised Maze

What about other sensory modalities beyond vision?

**Size Estimation**

Pearson's R

0.45

0.25

0.05

Contrastive Maze | Supervised Maze

**Object Categorization**

Accuracy

0.40

0.25

0.10

Contrastive Maze | Supervised Maze

**Texture Classification**

0.3

0.2

0.1

Contrastive Maze | Supervised Maze

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

Trinity Chung[*,1], Yuchen Shen[*,2], Nathan C. L. Kong[4], and Aran Nayebi[2, 3, 1]

[1]Robotics Institute, Carnegie Mellon University; Pittsburgh, PA 15213
[2]Machine Learning Department, Carnegie Mellon University; Pittsburgh, PA 15213
[3]Neuroscience Institute, Carnegie Mellon University; Pittsburgh, PA 15213
[4]Department of Psychology, University of Pennsylvania; Philadelphia, PA 19104
[*] Equal contribution.
{trinityc, yuchens3, anayebi}@cs.cmu.edu; nclkong@sas.upenn.edu

To appear as a NeurIPS 2025 Oral!



Trinity Chung*



Yuchen Shen*



Nathan C.L. Kong

# Why tactile?

# Why tactile?

- Tactile data is very useful in manipulating occluded and OOD objects

# Why tactile?

- Tactile data is very useful in manipulating occluded and OOD objects

- Tactile hardware & sim is getting better!

# 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

# 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

Trinity's search on arxiv...



# of <u>tactile</u>    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 <u>vision</u>    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)

# 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

Trinity's search on arxiv...

**arXiv**

# of <u>tactile</u>        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 <u>vision</u>        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)

# 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

Trinity's search on arxiv...



# of <u>tactile</u>    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 <u>vision</u>    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

Block Masking

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

# 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 <u>tactile</u>    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
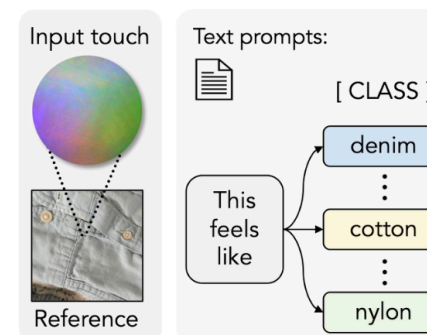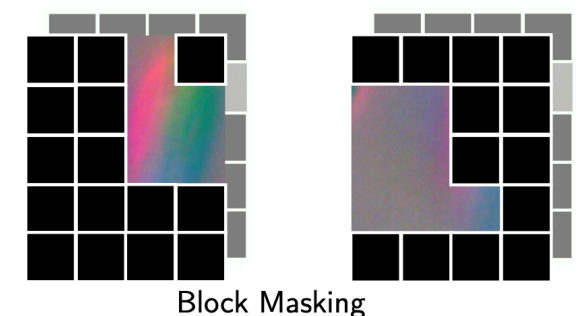
# of <u>vision</u>    2,577 results

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



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: 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

# Training Data: Tactile vs Image Augmentation

# Training Data: Tactile vs Image Augmentation



- Traditional image augmentation introduces Gaussian noise, color jitter, and grayscale.

# 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

# 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

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

- We needed a way to systematically search over the space of recurrent model architectures

# 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

# 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

# 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



Task Performance per Model
(Colored by Encoder, Labeled by Loss)

# Results: ConvRNN encoders perform best



Task Performance per Model
(Colored by Encoder, Labeled by Loss)

- Lighter color bar represents untrained version.

# Results: ConvRNN encoders perform best



Task Performance per Model
(Colored by Encoder, Labeled by Loss)

- Lighter color bar represents untrained version.

- S4 Encoders often don't even train on this task!

# Results: ConvRNN encoders perform best



Task Performance per Model
(Colored by Encoder, Labeled by Loss)

- 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)

# Results: ConvRNN encoders perform best



Task Performance per Model
(Colored by Encoder, Labeled by Loss)

- 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

# Neural Evaluation: Results



Real

far
medium
near

concave          convex

Rodgers 2022

# Neural Evaluation: Results



Rodgers 2022

# Neural Evaluation: Results



Neural Fit per Model
(Colored by Encoder, Labeled by Loss)





Rodgers 2022

# Neural Evaluation: Results



Maximal NeuroAI Turing Test

Neural Fit per Model
(Colored by Encoder, Labeled by Loss)

inter-animal max

Encoder: vanilla, resnet, ugrnn, inter, gru, lstm, zhuang

Attender: gpt, mamba, none

Rodgers 2022

# Neural Evaluation: Results



Rodgers 2022

# Neural Evaluation: Results



Maximal NeuroAI Turing Test

Neural Fit per Model
(Colored by Encoder, Labeled by Loss)

# Neural Evaluation: Results



Neural Fit per Model
(Colored by Encoder, Labeled by Loss)

Maximal NeuroAI Turing Test

- We've nearly passed the NeuroAI Turing Test, for this dataset at least

# Neural Evaluation: Results



Maximal NeuroAI Turing Test

Neural Fit per Model
(Colored by Encoder, Labeled by Loss)

- We've nearly passed the NeuroAI Turing Test, for this dataset at least
  - Need more stimuli to evaluate with!

# Neural Evaluation: Results



Neural Fit per Model
(Colored by Encoder, Labeled by Loss)

Maximal NeuroAI Turing Test

- We've nearly passed the NeuroAI Turing Test, for this dataset at least
  - Need more stimuli to evaluate with!

- ConvRNNs outperform feedforward/SSMs on realistic tactile recognition

# Neural Evaluation: Results



Neural Fit per Model
(Colored by Encoder, Labeled by Loss)

- We've nearly passed the NeuroAI Turing Test, for this dataset at least
    - Need more stimuli to evaluate with!

- ConvRNNs outperform feedforward/SSMs on realistic tactile recognition

- ConvRNNs best match neural responses in mouse barrel cortex

# Neural Evaluation: Results



Neural Fit per Model
(Colored by Encoder, Labeled by Loss)

- We've nearly passed the NeuroAI Turing Test, for this dataset at least
  - Need more stimuli to evaluate with!

- ConvRNNs outperform feedforward/SSMs on realistic tactile recognition

- ConvRNNs best match neural responses in mouse barrel cortex

- Contrastive SSL *matches* supervised neural alignment, possibly suggesting a general-purpose representation in the somatosensory cortex (needs more neural data to explore this!)

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

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

## Recurrence + Contrastive SSL?

# 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

**Infer:**
Has this ice block been out longer?

**Infer:**
Has this ice block been out longer?

Visually-Grounded Mental Simulation

Infer:
Has this ice block been out longer?

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

**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?

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

**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

**Inputs**

Dominoes  Support

*Physion*

Dominoes  Support

Drape  Link

**Sensory-Cognitive Hypothesis Classes**

*Latent Future Prediction:*

**1. Pretraining Stage**

Ego4D, etc

**Foundation Model**

T

**2. Dynamics Pretraining Stage**

**Ground Truth**

T+1

T+1

**Prediction**

*End-to-End Future Prediction:*

Encoder  Decoder  *Pixel-wise*

Encoder  *Object-slot*

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

Yes/No?

**Observed Stimuli**

Example Scenarios

cue  *Time*  stimulus

**Unobserved Outcome**

last frame  true label

NO
acc. = 0.89

YES
*acc. = 0.96*

**Macaque Neurophysiology: Mental-Pong**

DMFC

M
L
A ←→ P

Time
Feedback

Occluded epoch
(895±270 ms)

Observed epoch
(1240±350 ms)

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

# Latent Future Prediction



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

Leverage these dynamics to do explicit future prediction

# Latent Future Prediction



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

What vision task?

Leverage these dynamics to do explicit future prediction

# 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

# Ego4D: everyday activity around the world



## 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



Geographic diversity

*Grauman et al. 2022*

## Ego4D: everyday activity around the world



## 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



Geographic diversity

*Grauman et al. 2022*

## Ego4D: everyday activity around the world



$$\mathcal{L}_{contrastive} =$$

$$-\sum_{b \in B} \log \frac{\overbrace{e^{\mathcal{S}(\mathbf{z}_i^b, \mathbf{z}_j^b)}}^{\text{attract}}}{e^{\mathcal{S}(\mathbf{z}_i^b, \mathbf{z}_j^b)} + \overbrace{e^{\mathcal{S}(\mathbf{z}_i^b, \mathbf{z}_k^b)}}^{\text{repel}} + \overbrace{e^{\mathcal{S}(\mathbf{z}_i^b, \widetilde{\mathbf{z}}_i^b)}}^{\text{repel}}}$$

$$[I_i, I_{j>i}, I_{k>j}]^{1:B}$$

### 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



Geographic diversity

*Grauman et al. 2022*

Majumdar et al. 2023

# Ego4D: everyday activity around the world

# 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

Geographic diversity

Grauman et al. 2022

Ego4D: everyday activity around the world

*Majumdar et al. 2023*

Meta World

DM Control

Trifinger

Adroit

Mobile-Pick

ImageNav

ObjectNav

**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

Geographic diversity

*Grauman et al. 2022*

# Video Foundation Future Prediction Best Predict Neurons

# Dynamically-Equipped Video Foundation Models Can Match Both

# Dynamically-Equipped Video Foundation Models Can Match Both

# Dynamically-Equipped Video Foundation Models Can Match Both

**DMFC**

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



**VC-1+LSTM**

**VC-1+CTRNN**

Neural Predictivity *(Pearson's R)*

Correlation to Average Human Response *(Pearson's R)*

Dominoes

Yes/No?

# Dynamically-Equipped Video Foundation Models Can Match Both

**DMFC**

Legend:
- ◆ End-to-End
- ■ Image Foundation Models
- ● Video Foundation Models

**Better Models**

?

VC-1+LSTM

VC-1+CTRNN

Dominoes

**But we have a ways to go to reach the NeuroAI Turing Test here!**

**Neural Predictivity** *(Pearson's R)*

**Correlation to Average Human Response** *(Pearson's R)*

Yes/No?

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

## Recurrence + Contrastive SSL?

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

Recurrence + Contrastive SSL?

Latent Future Prediction?

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

Recurrence + Contrastive SSL?

Latent Future Prediction?

**Place Cell**
**(Hippocampus)**

# A Task-Optimized Account of Heterogeneity

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

**Border Cells**

**Heterogeneous Cells**

Data from: *Mallory et al. 2021*

# A Task-Optimized Account of Heterogeneity

# A Task-Optimized Account of Heterogeneity



Grid Cells
More like ~2-3%!

Border Cells

Heterogeneous Cells

Data from: *Mallory et al. 2021*

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

*velocities$_t$*
Input

"MEC"

*positions (x,y)$_{t+1}$*
Output

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

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

NeurIPS 2021 (spotlight)

# A Task-Optimized Account of Heterogeneity

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

**Border Cells**

**Heterogeneous Cells**

Data from: *Mallory et al. 2021*

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

$velocities_t$
Input

*"MEC"*

$positions\ (x,y)_{t+1}$
Output

Place Cell (Hippocampus)

# Autonomous Behavior and Whole-Brain Dynamics Emerge in Embodied Zebrafish Agents with Model-based Intrinsic Motivation

Reece Keller[1,2,*]  Alyn Tornell[2]  Felix Pei[2]  Xaq Pitkow[1,3]

Leo Kozachkov[4,†]  Aran Nayebi[3,1,2,†]

## To appear at NeurIPS 2025!



Reece Keller    Alyn Tornell    Felix Pei    Xaq Pitkow    Leo Kozachkov†

# Why is Animal Autonomy Hard?
The behavioral repertoire is enormous...

# Why is Animal Autonomy Hard?

The behavioral repertoire is enormous…

# 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?**

Slides credit: Reece Keller

# 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.**

# 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!

Slides credit: Reece Keller

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

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

virtual reality navigation

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

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

virtual reality navigation

Mu et al., *Cell* (2019)

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

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

virtual reality navigation



Mu et al., *Cell* (2019)

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

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

✔ 1. Ecologically-relevant environment

✔ 2. "Cognitive" states with clear behavioral readouts

✔ 3. Large-scale multi-area neural recordings

virtual reality navigation



Mu et al., *Cell* (2019)

Neurons

25 µm    50 µm    75 µm    100 µm    125 µm

Radial astrocytes

25 µm    50 µm    75 µm    100 µm    125 µm

Δf/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

Neurons

25 µm  50 µm  75 µm  100 µm  125 µm

Radial astrocytes

25 µm  50 µm  75 µm  100 µm  125 µm

Δf/f

Open loop

Fictive behavior
3.0X playback speed

# Machine Autonomy: Prior Work

Exploration in sparse/reward-free environments

| Curiosity type | Formulation | What it measures |
|---|---|---|
| Surprise | $r_t^i \propto -\log \omega_\theta(\mathbf{s'} \mid \mathbf{s}, \mathbf{a})$ | prediction error |
| Disagreement | $r_t^i \propto \mathrm{Var}\left(\{\omega_{\theta_j}(\mathbf{s'} \mid \mathbf{s}, \mathbf{a})\}_{j=1}^N\right)$ | prediction variance |
| Learning progress | $r_t^i \propto \log \dfrac{\omega_{\theta'}(\mathbf{s'} \mid \mathbf{s}, \mathbf{a})}{\omega_\theta(\mathbf{s'} \mid \mathbf{s}, \mathbf{a})}$ $\theta \leftarrow (1-\gamma)\theta + \gamma\theta'$ | prediction error gain |

Kim et al., *ICML* (2020)
Pathak et al., *ICML* (2017)
Burda, Edwards & Pathak et al., *NeurIPS* (2017)

# Machine Autonomy: Prior Work

Explore states/regions/novel/new environments



**Atari**

| | **What it measures** |
|---|---|
| $\mathbf{a})$ | prediction error |
| $s, \mathbf{a})\}_{j=1}^{N} \Big)$ | prediction variance |
| | prediction error gain |

Learning progress

$$r_t^t \propto \log \frac{\cdots}{\omega_\theta(\mathbf{s}' \mid \mathbf{s}, \mathbf{a})}$$

$$\theta \leftarrow (1 - \gamma)\theta + \gamma\theta'$$

Kim et al., *ICML* (2020)
Pathak et al., *ICML* (2017)
Burda, Edwards & Pathak et al., *NeurIPS* (2017)

# Machine Autonomy: Prior Work



**Atari**

**Mario**

Learning progress

$$r_t^i \propto \log \frac{\cdot}{\omega_\theta(\mathbf{s}' \mid \mathbf{s}, \mathbf{a})}$$

$$\theta \leftarrow (1 - \gamma)\theta + \gamma\theta'$$

prediction variance

prediction error gain

Kim et al., *ICML* (2020)
Pathak et al., *ICML* (2017)
Burda, Edwards & Pathak et al., *NeurIPS* (2017)

# Machine Autonomy: Prior Work



Atari

Mario

DM-Control

Learning progress

$$r_t^i \propto \log \frac{\ }{\omega_\theta(\mathbf{s}' \mid \mathbf{s}, \mathbf{a})}$$

$$\theta \leftarrow (1 - \gamma)\theta + \gamma\theta'$$

Kim et al., *ICML* (2020)
Pathak et al., *ICML* (2017)
Burda, Edwards & Pathak et al., *NeurIPS* (2017)

# Machine Autonomy: Prior Work



**Atari**

**Mario**

**DM-Control**

Learning progress

$$r_t^i \propto \log \frac{\omega_{\theta'}(\mathbf{s}' \mid \mathbf{s}, \mathbf{a})}{\omega_\theta(\mathbf{s}' \mid \mathbf{s}, \mathbf{a})}$$

$$\theta \leftarrow (1-\gamma)\theta + \gamma\theta'$$

Kim et al., *ICML* (2020)
Pathak et al., *ICML* (2017)
Burda, Edwards & Pathak et al., *NeurIPS* (2017)

**Often leads to unethological behaviors! (or can be stuck on white noise)**

# Epistemic Curiosity isn't Enough...

## Animal autonomy != novelty optimization



Training step (1e6)

γ-Progress    ICM    RND    Disagreement    Training environment

10
20
30
40
50
60

2 N

3 s

## What's the issue?

- Rewards are non-stationary and saturate with experience.

    Consequence: behavioral strategies are transient

    (e.g. γ-Progress)

- Rewards can perseverate on unpredictable/uncontrollable stimuli.

    Consequence: unethological behavior (e.g. ICM)

# 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)

# 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*

# 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?

# 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)?

# 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)?

## 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 Simulation Environment

## 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 Simulation Environment



## 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 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!

### 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 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!

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 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



NE-MO Neurons

GABA Neurons

Motor Neurons

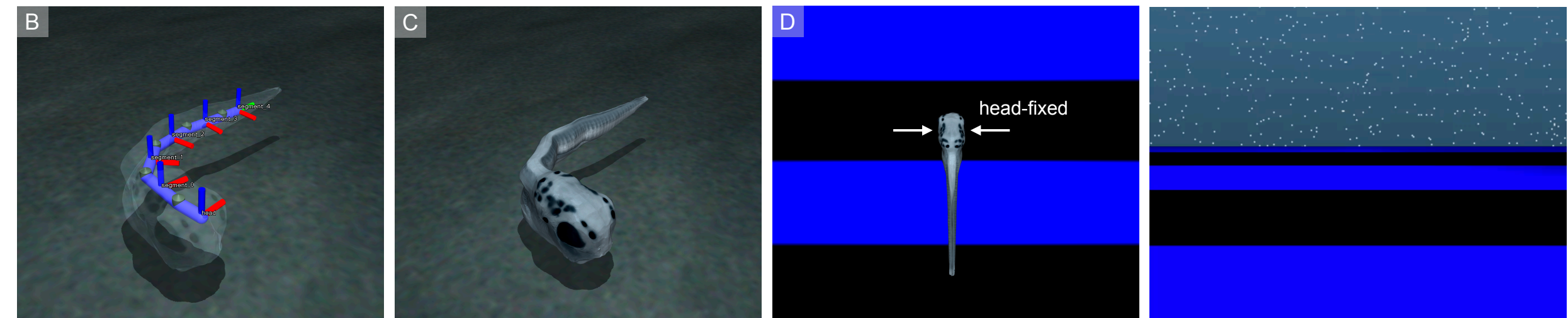Radial astrocytes

**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

### Using ethological memory to guide adaptive behavior



☐ ethological

$T_1(\mathbf{s'} \mid \mathbf{s}, \mathbf{a})$ — distill via experience

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

☐ unethological

head-fixed

$T_2(\mathbf{s'} \mid \mathbf{s}, \mathbf{a})$ — distill via experience

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

**3M: Model-Memory-Mismatch**

$\omega_\theta$    $\epsilon_t = D_{\mathrm{KL}}\left[\omega_{\theta'} \mid\mid \omega_\theta\right]$    $\omega_{\theta'}$

$s_t$

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

$$r_t^i \propto |\hat{\epsilon}_t - \epsilon_t|$$
$$\hat{\epsilon}_t = (1 - \gamma)\hat{\epsilon}_{t-1} + \gamma\epsilon_t$$

**3M-Progress**

memories agree    memories disagree

$r_t^i$

$\epsilon_t$

$t$

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

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

(dynamics agree on a subspace).

# 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**

$T_1(s' \mid s, a)$ — distill via experience

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

**3M: Model-Memory-Mismatch**

$\omega_\theta$ $\qquad$ $\omega_{\theta'}$

$\epsilon_t = D_{KL}\left[\omega_{\theta'} \mid\mid \omega_\theta\right]$

$s_t$

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

**unethological**

head-fixed

$T_2(s' \mid s, a)$ — distill via experience

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

**3M-Progress**

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

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

memories agree $\qquad$ memories disagree

$r_t^i$

$\epsilon_t$

$t$

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).

# 3M-Progress Captures Whole-Brain Dynamics
## (and behavior)
### Single-cell one-to-one alignment

# 3M-Progress Captures Whole-Brain Dynamics
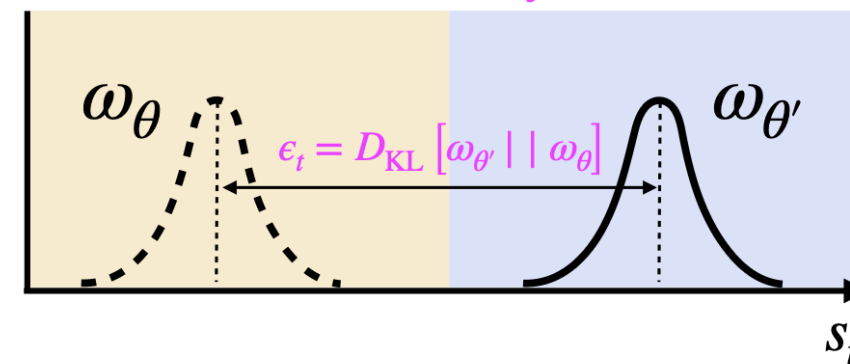## (and behavior)
### Single-cell one-to-one alignment

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

Recurrence + Contrastive SSL?   Latent Future Prediction?

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

Recurrence + Contrastive SSL?

Latent Future Prediction?

Sensory (Input) Stream

Perceptual Module

Self-Supervision

Negative Samples

Positive Samples

Future Inference Module

Object-centric

Latent State

Dynamics Predictor

Planning Module

Key

Value

Cognitive Map

M1

DMFC

V1

V4

HPC

IT

Cerebellum

Macaques

V1

HPC

M1

mPFC

Cerebellum

Rodents

M1

Parietal Lobe

Frontal Lobe

V1

V4

IT

HPC

Cerebellum

Humans

Motor Module

High-level Controller

Low-level Controller

Intrinsic Goals

Action (Output) Stream

Environment

Temporal integration of World Model-Progress-based curiosity?

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

Recurrence + Contrastive SSL?    Latent Future Prediction?



Temporal integration of World Model-Progress-based curiosity?

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

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



Temporal integration of World Model-Progress-based curiosity?

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

**Paper:** https://arxiv.org/abs/2502.05934

**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.**

Negative

Object-centric

Key

Value

Cognitive Map

Intrinsic Goals

Action
(Output) Stream

Environment

**Temporal integration of World Model-Progress-based curiosity?**

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

**Paper:** https://arxiv.org/abs/2502.05934

**2. Core Safety Values for Provably Corrigible Agents**

**Paper:** https://arxiv.org/abs/2507.20964

**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.**

**<u>Open</u>: Can we scale corrigibility cost effectively?**

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

**Paper:** https://arxiv.org/abs/2502.05934

**2. Core Safety Values for Provably Corrigible Agents**

**Paper:** https://arxiv.org/abs/2507.20964

Figure 1: **Projected AI capabilities ($\gamma_t$) vs. time-varying UBI AI capability threshold ($\gamma_t^\star$).** The dashed line is the required capability $\gamma_t^\star$ to fully fund a UBI that comprises 11% of the GDP (leading to a $\gamma_t^\star$ between 5-6× 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 $\gamma_t^\star$ 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.

**3. An AI Capability Threshold for Rent-Funded Universal Basic Income in an AI-Automated Economy**

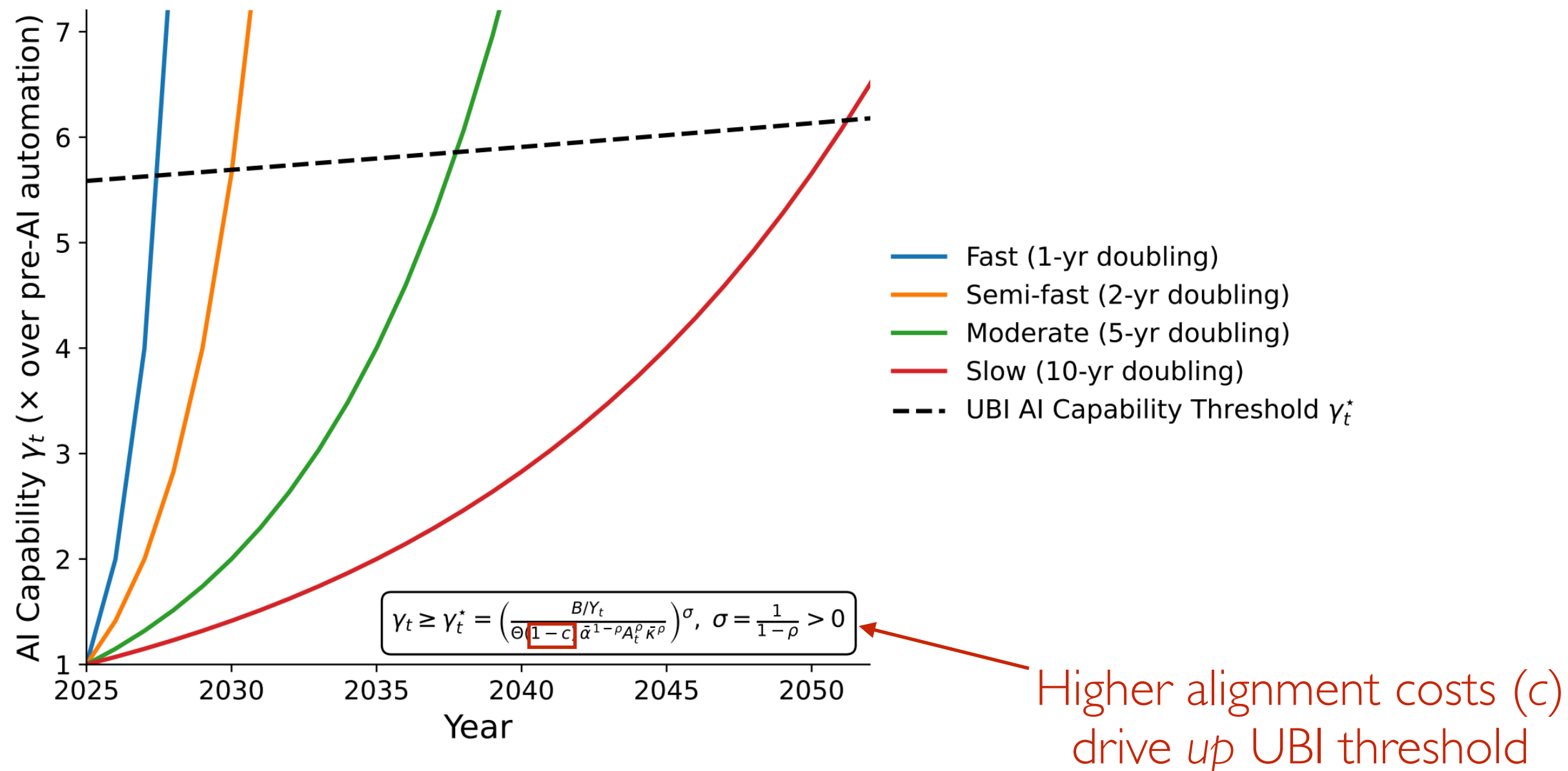**Paper:** https://arxiv.org/abs/2505.18687

Figure 1: **Projected AI capabilities ($\gamma_t$) vs. time-varying UBI AI capability threshold ($\gamma_t^\star$).** The dashed line is the required capability $\gamma_t^\star$ to fully fund a UBI that comprises 11% of the GDP (leading to a $\gamma_t^\star$ 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 $\gamma_t^\star$ 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.

**3. An AI Capability Threshold for Rent-Funded Universal Basic Income in an AI-Automated Economy**

**Paper:** https://arxiv.org/abs/2505.18687

Open: Can we incorporate other values (besides control, which is "neutrally amoral") that lead to longer term human well-being, especially if working for pay becomes no longer feasible in many cases?

Legend:
- Fast (1-yr doubling)
- Semi-fast (2-yr doubling)
- Moderate (5-yr doubling)
- Slow (10-yr doubling)
- UBI AI Capability Threshold $\gamma_t^\star$

$$\gamma_t \geq \gamma_t^\star = \left(\frac{B/Y_t}{\Theta(1-c)\,\tilde{\alpha}^{1-\rho}A_t^\rho\bar{\kappa}^\rho}\right)^\sigma, \quad \sigma = \frac{1}{1-\rho} > 0$$

Higher alignment costs (c) drive *up* UBI threshold

Figure 1: **Projected AI capabilities ($\gamma_t$) vs. time-varying UBI AI capability threshold ($\gamma_t^\star$).** The dashed line is the required capability $\gamma_t^\star$ to fully fund a UBI that comprises 11% of the GDP (leading to a $\gamma_t^\star$ between 5-6× 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 $\gamma_t^\star$ 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.

**3. An AI Capability Threshold for Rent-Funded Universal Basic Income in an AI-Automated Economy**

**Paper:** https://arxiv.org/abs/2505.18687

# Acknowledgements

## NeuroAgents Lab