

Using Embodied Agents for “Why” Questions in Systems Neuroscience

Aran Nayebi

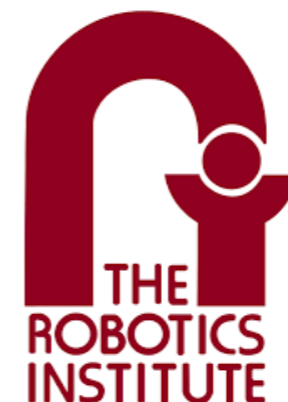
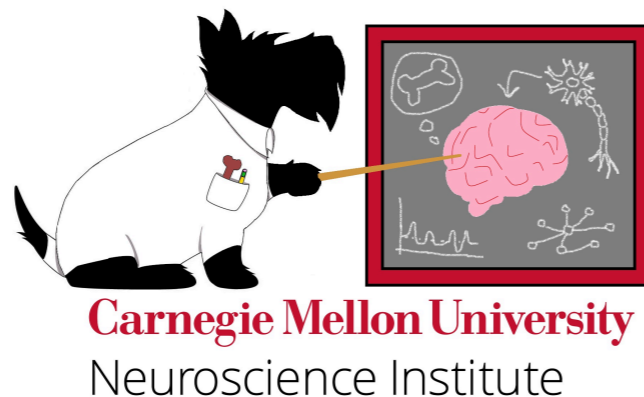
Assistant Professor

Machine Learning Department

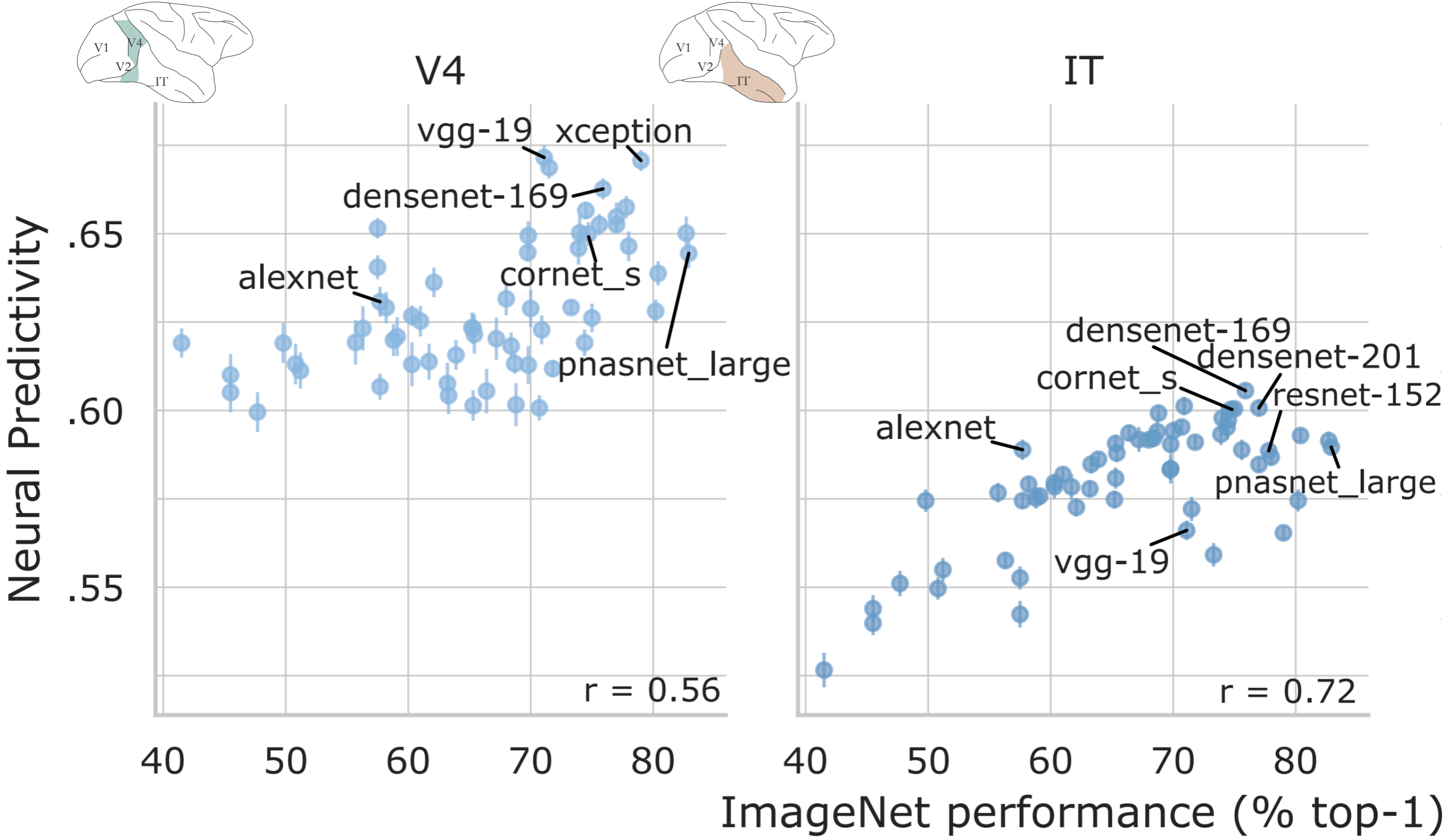
Neuroscience Institute (core faculty), Robotics Institute (courtesy)

EPFL BMI Seminar

2025.02.12

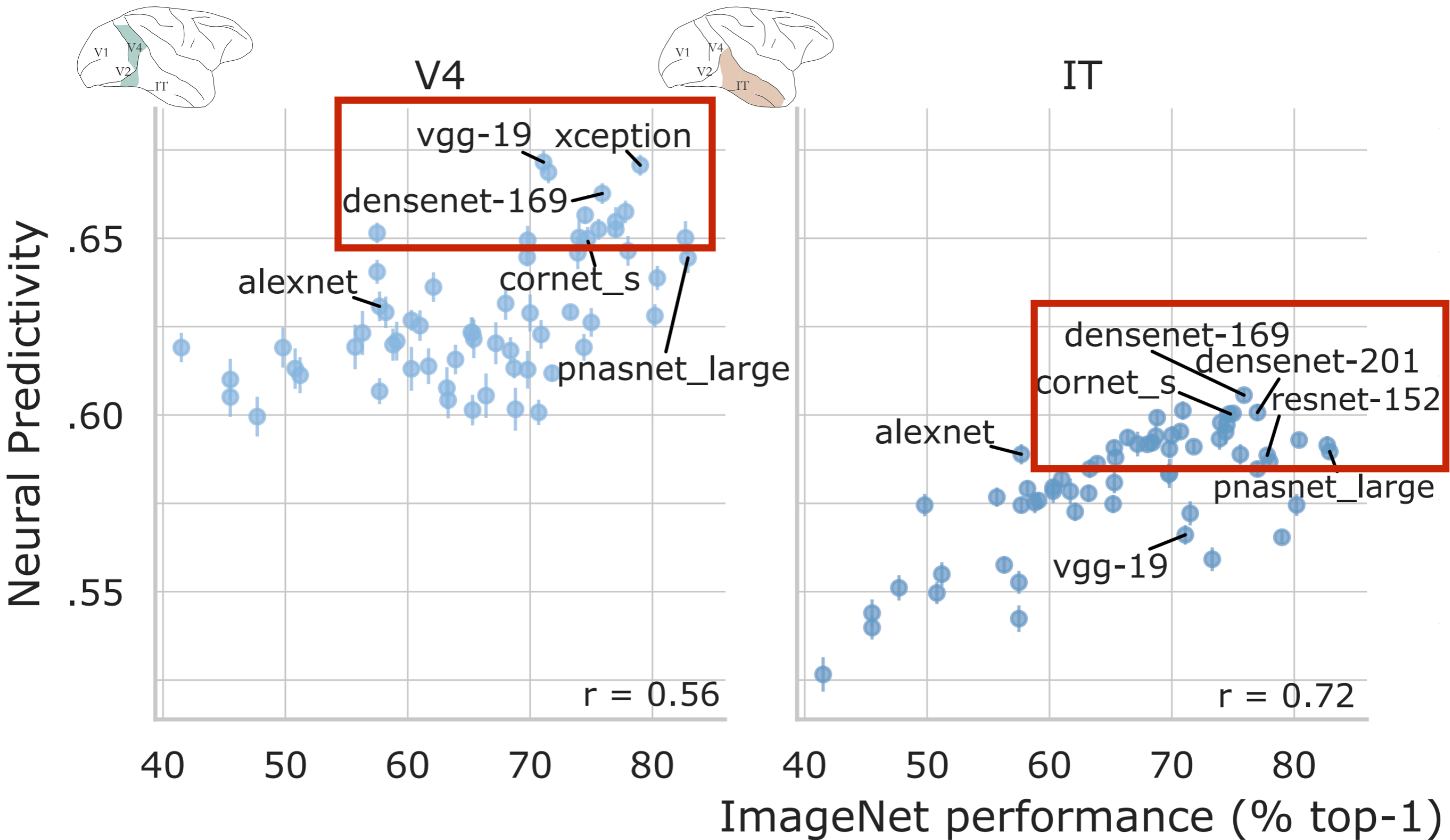


Similar predictivities among very different CNN architectures

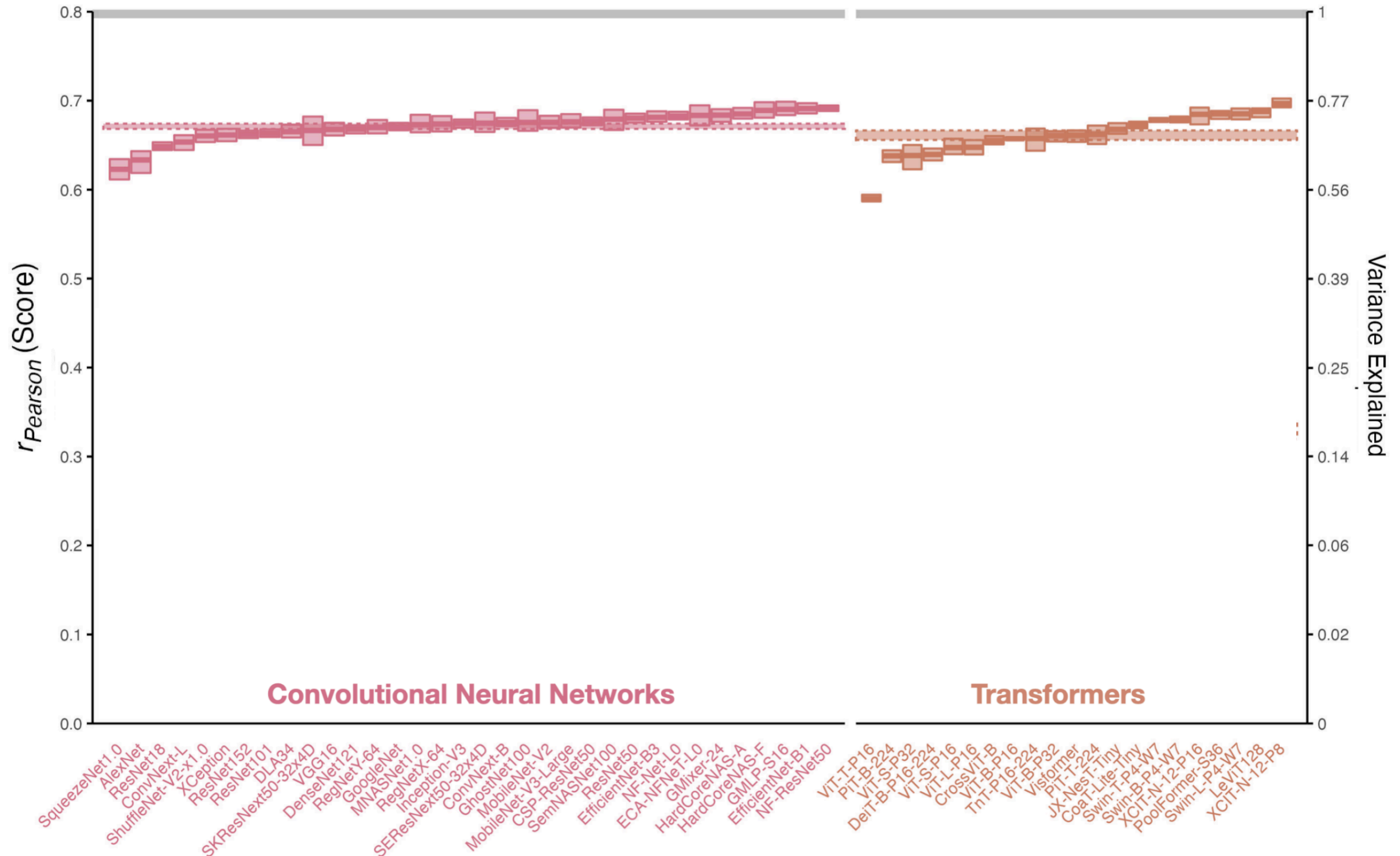
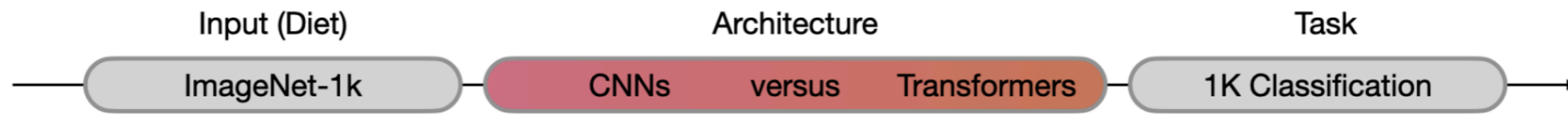
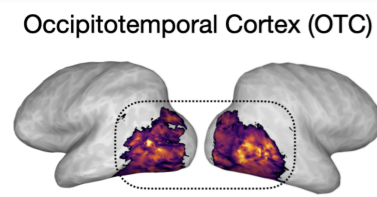


Schrimpf*, Kubilius* et al. 2018

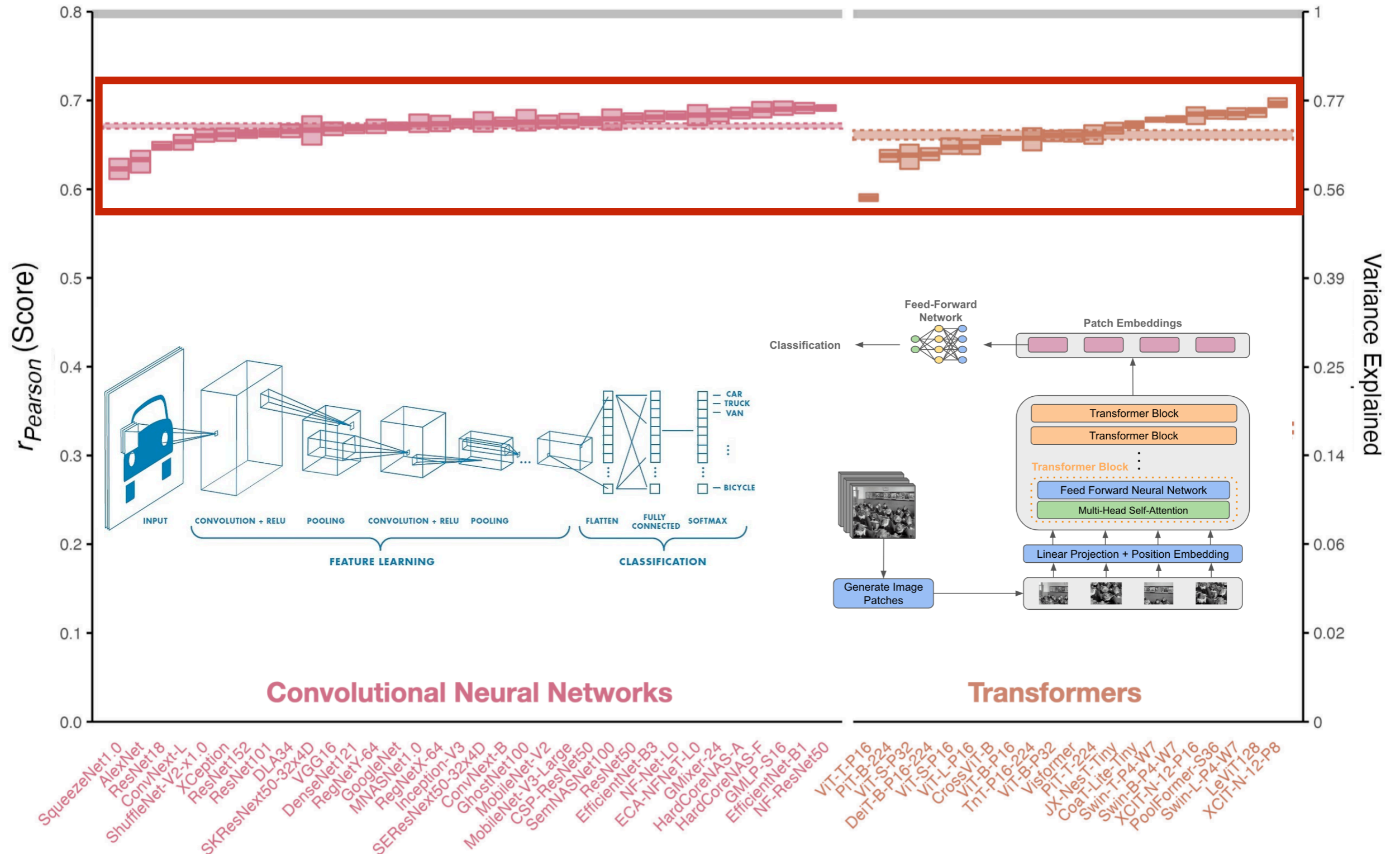
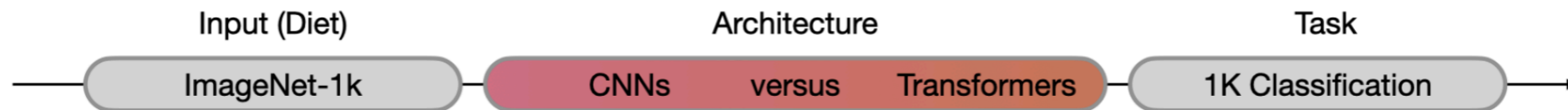
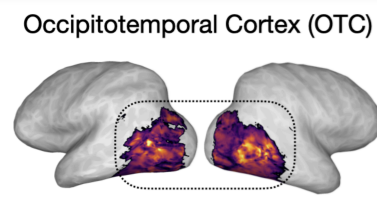
Similar predictivities among very different CNN architectures



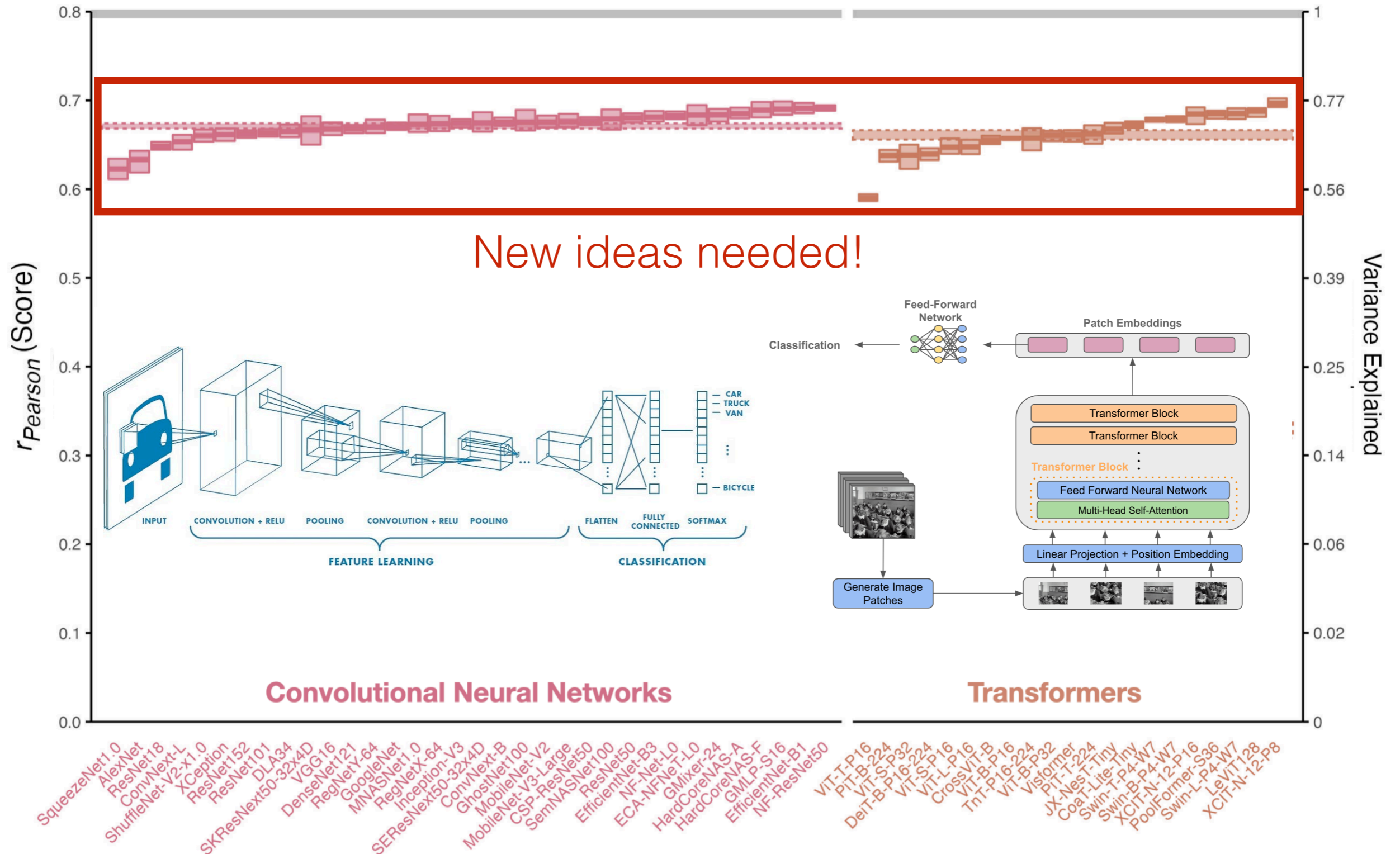
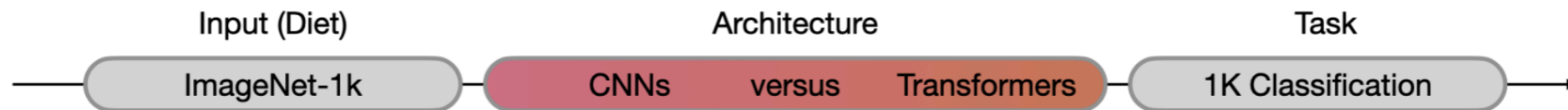
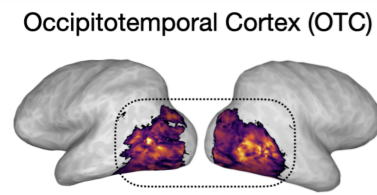
Similar predictivities between CNNs vs. Transformers



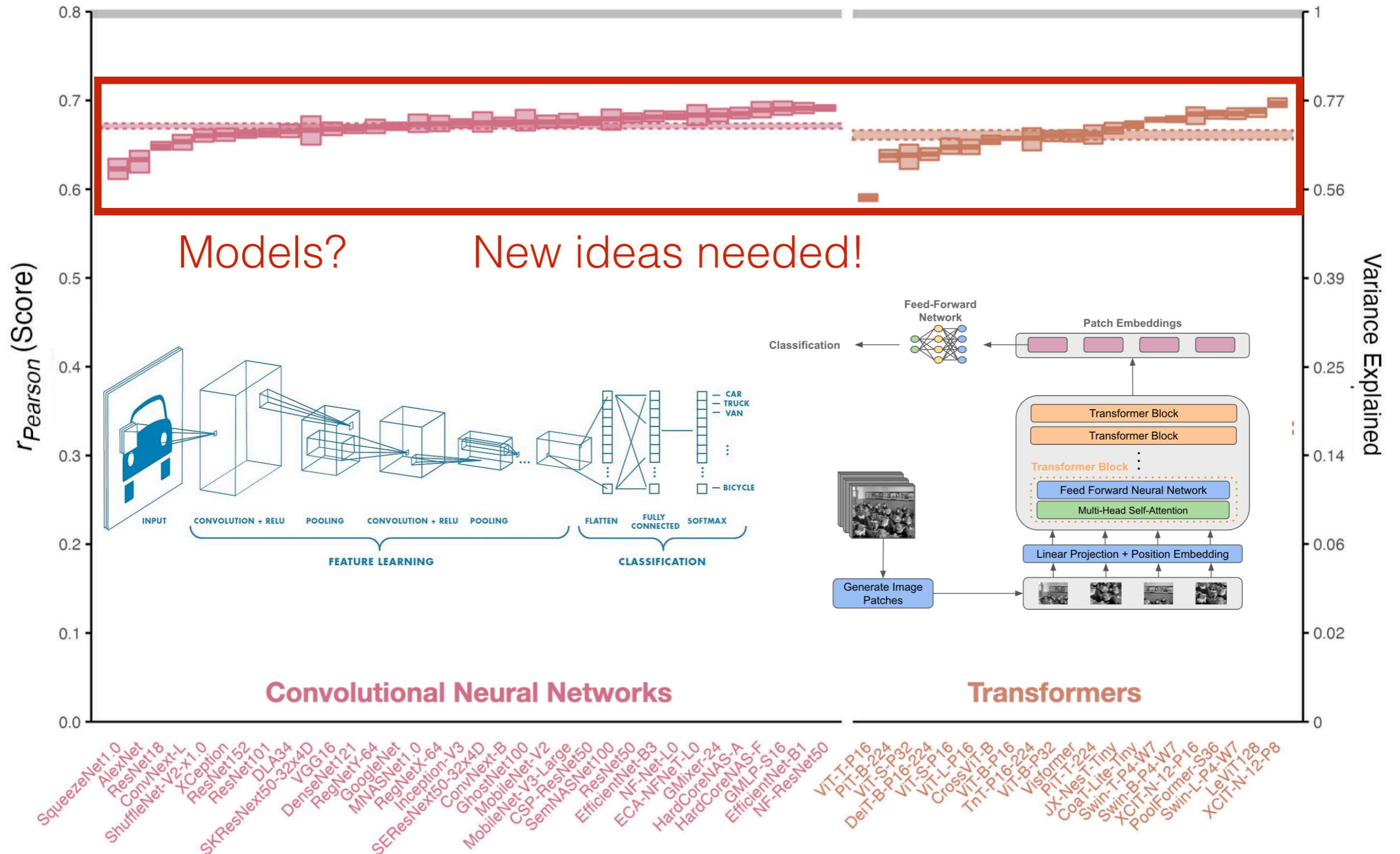
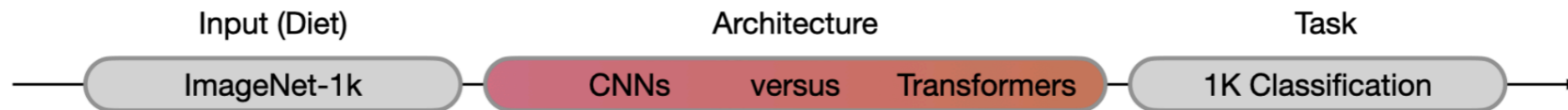
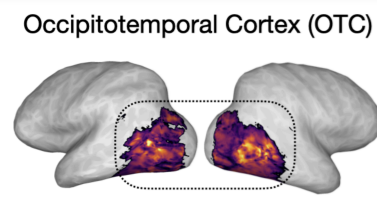
Similar predictivities between CNNs vs. Transformers



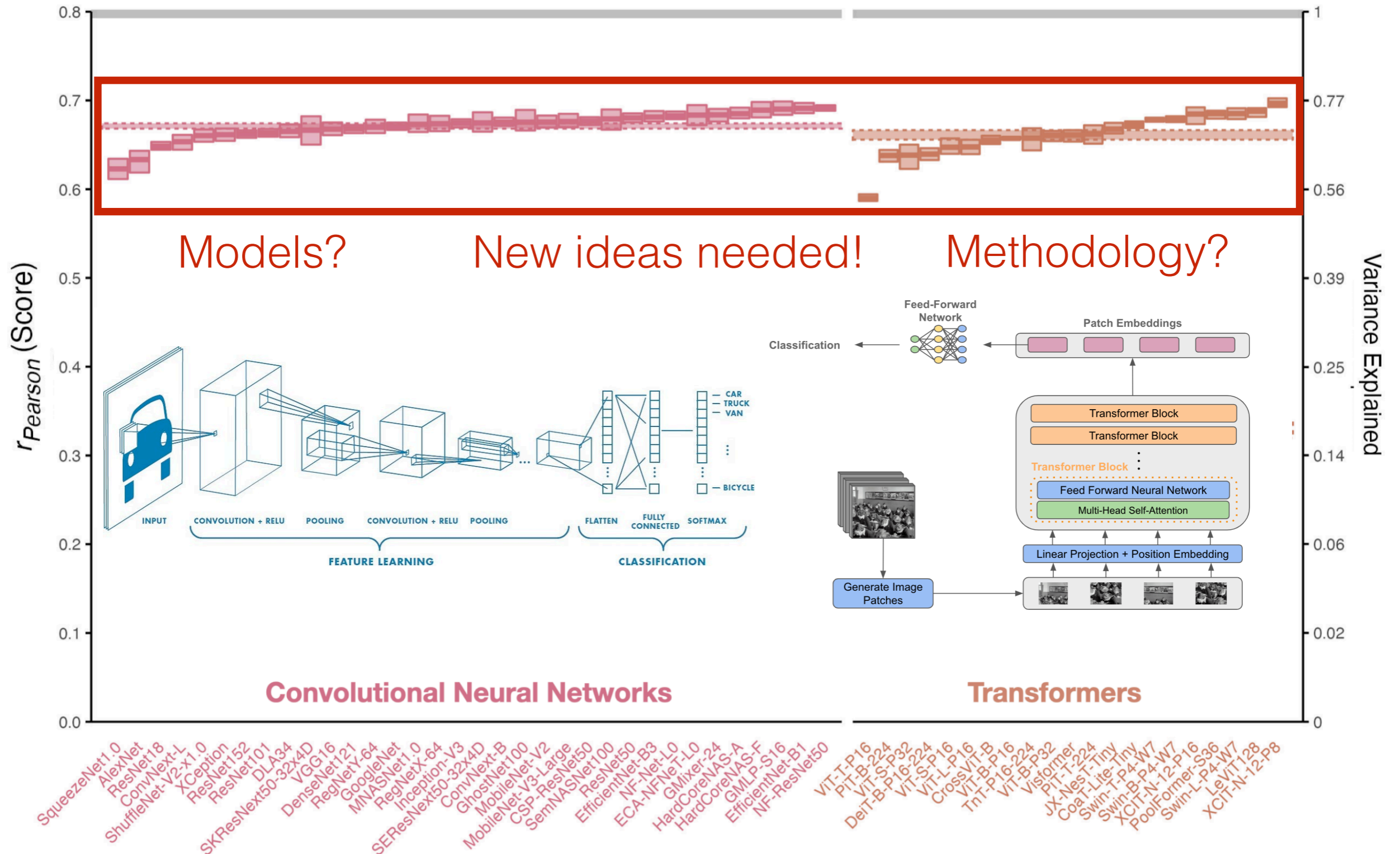
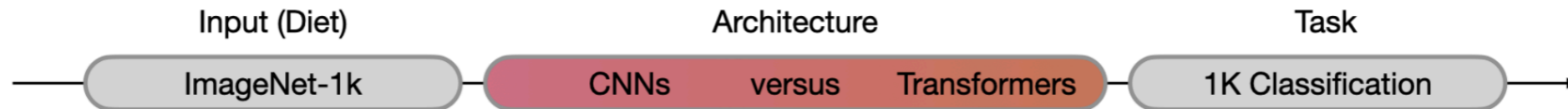
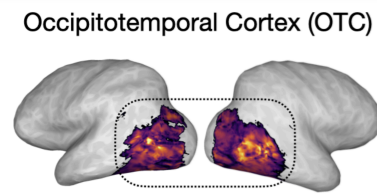
Similar predictivities between CNNs vs. Transformers



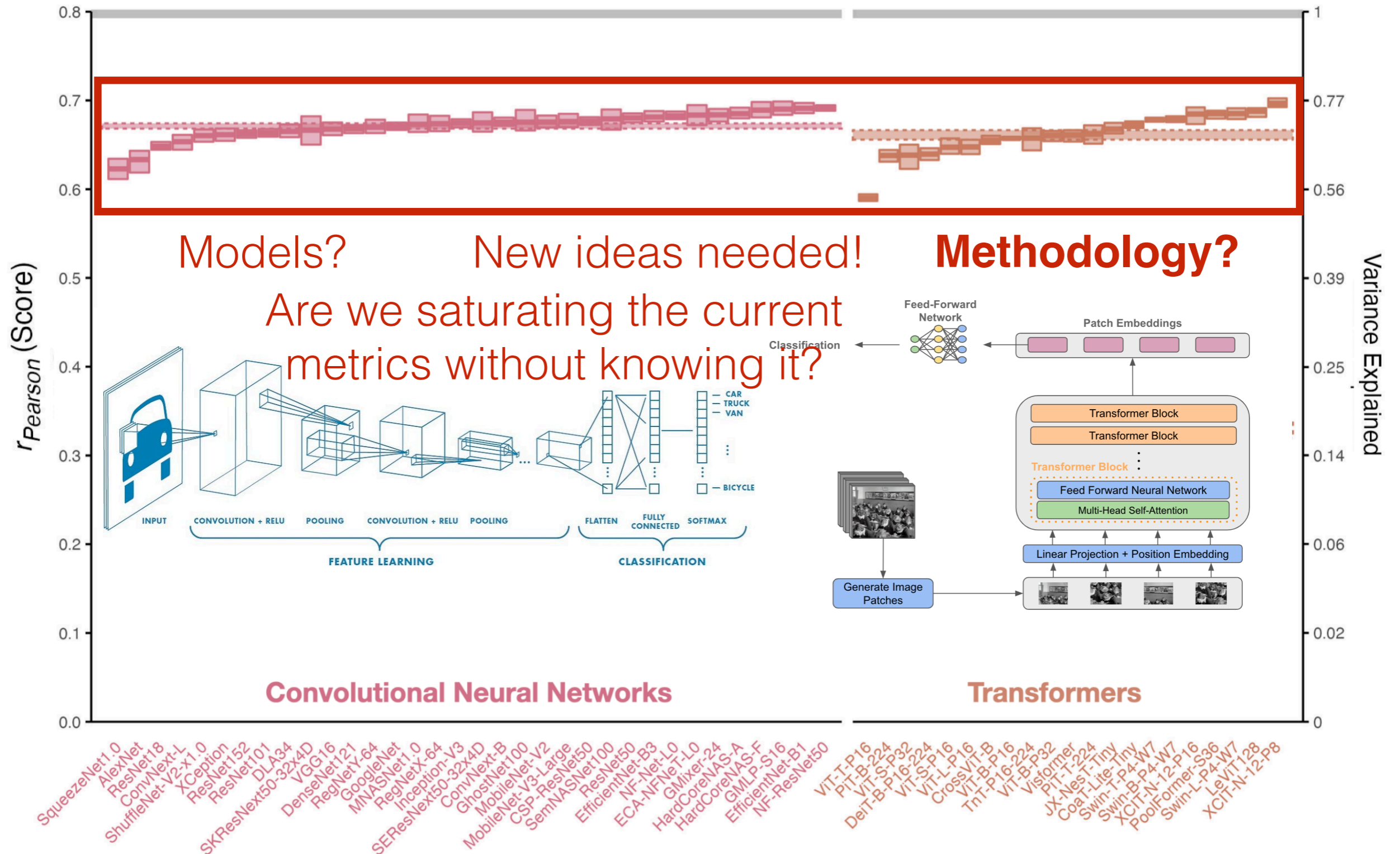
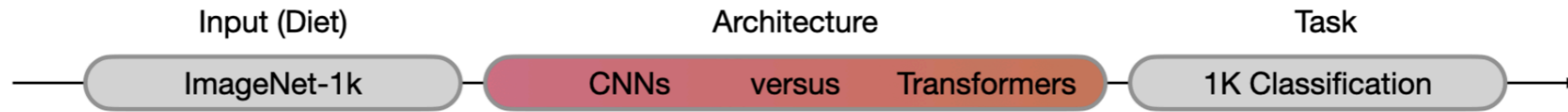
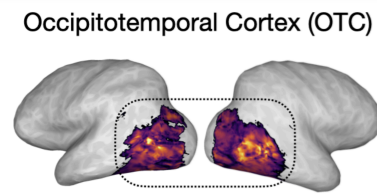
Similar predictivities between CNNs vs. Transformers



Similar predictivities between CNNs vs. Transformers



Similar predictivities between CNNs vs. Transformers



Determining the Methodology

Determining the Methodology

We don't have an *a priori* notion of what constitutes a
“good” brain-model in an *absolute* sense

Determining the Methodology

We don't have an *a priori* notion of what constitutes a
“good” brain-model in an *absolute* sense

...if we did, we would just put it in the model and be
done already!

Determining the Methodology

Brain-Model Evaluations Need the NeuroAI Turing Test

Jenelle Feather^{*1} **Meenakshi Khosla**^{*2} **N. Apurva Ratan Murty**^{*3} **Aran Nayebi**^{*4}

We don't have an *a priori* notion of what constitutes a
“good” brain-model in an *absolute* sense

...if we did, we would just put it in the model and be
done already!

Determining the Methodology

Brain-Model Evaluations Need the NeuroAI Turing Test

Jenelle Feather^{*1} Meenakshi Khosla^{*2} N. Apurva Ratan Murty^{*3} Aran Nayebi^{*4}

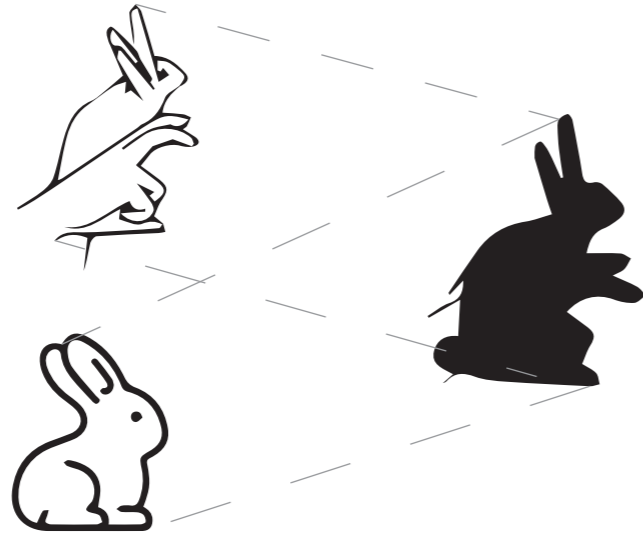
We don't have an *a priori* notion of what constitutes a
“good” brain-model in an *absolute* sense

...if we did, we would just put it in the model and be
done already!

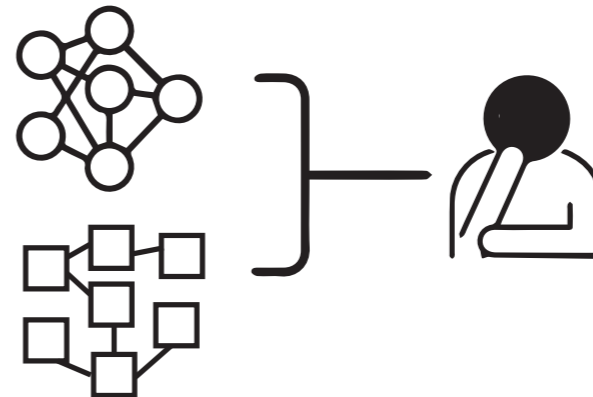
But, we do know that animals/subjects contain those
desired properties (even if we can't describe them), so
we instead take a *relative* perspective

NeuroAI Turing Test

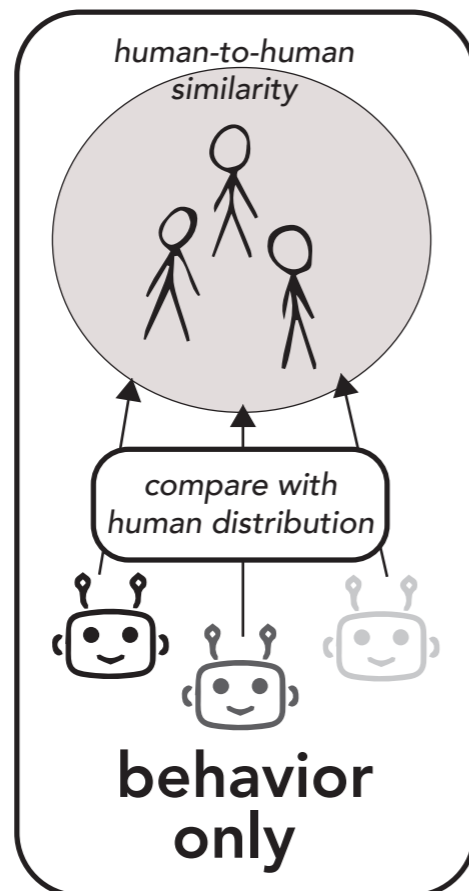
Just as distinct objects
can cast the same shadow...



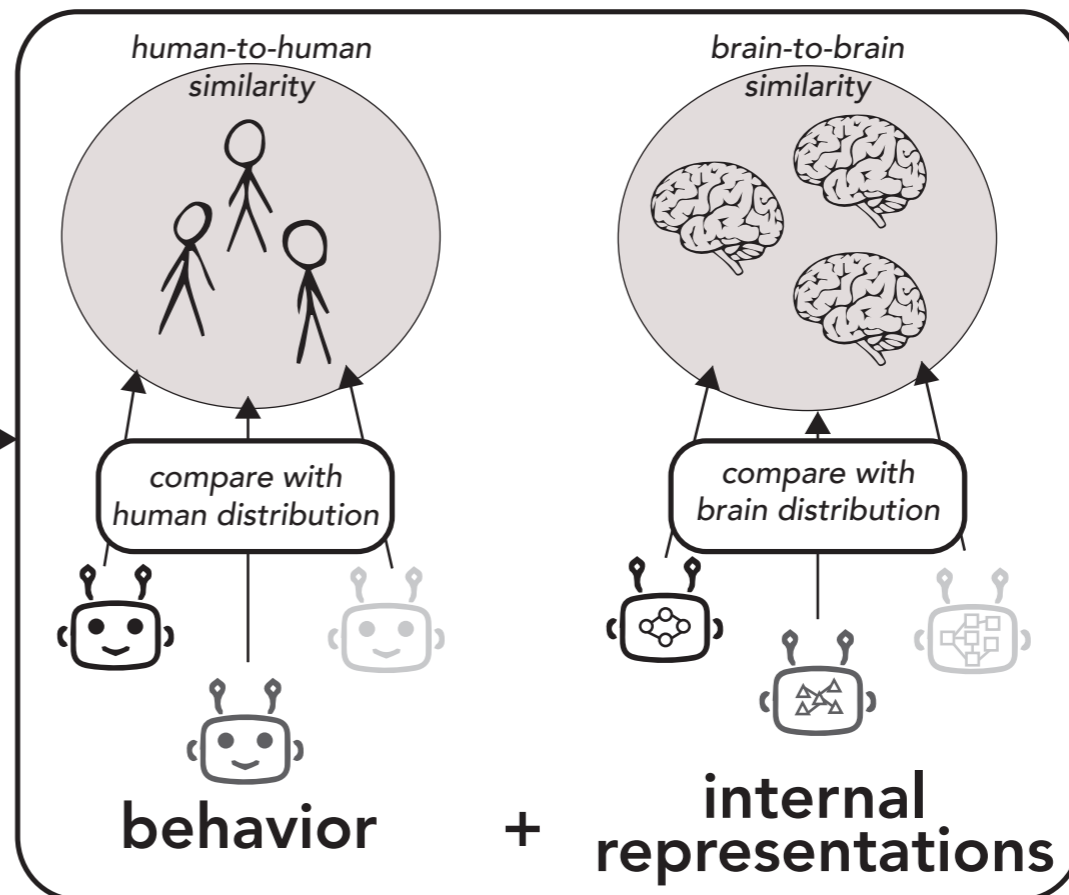
...distinct internal processes
(representations) can produce
similar outputs (behavior)



Turing Test

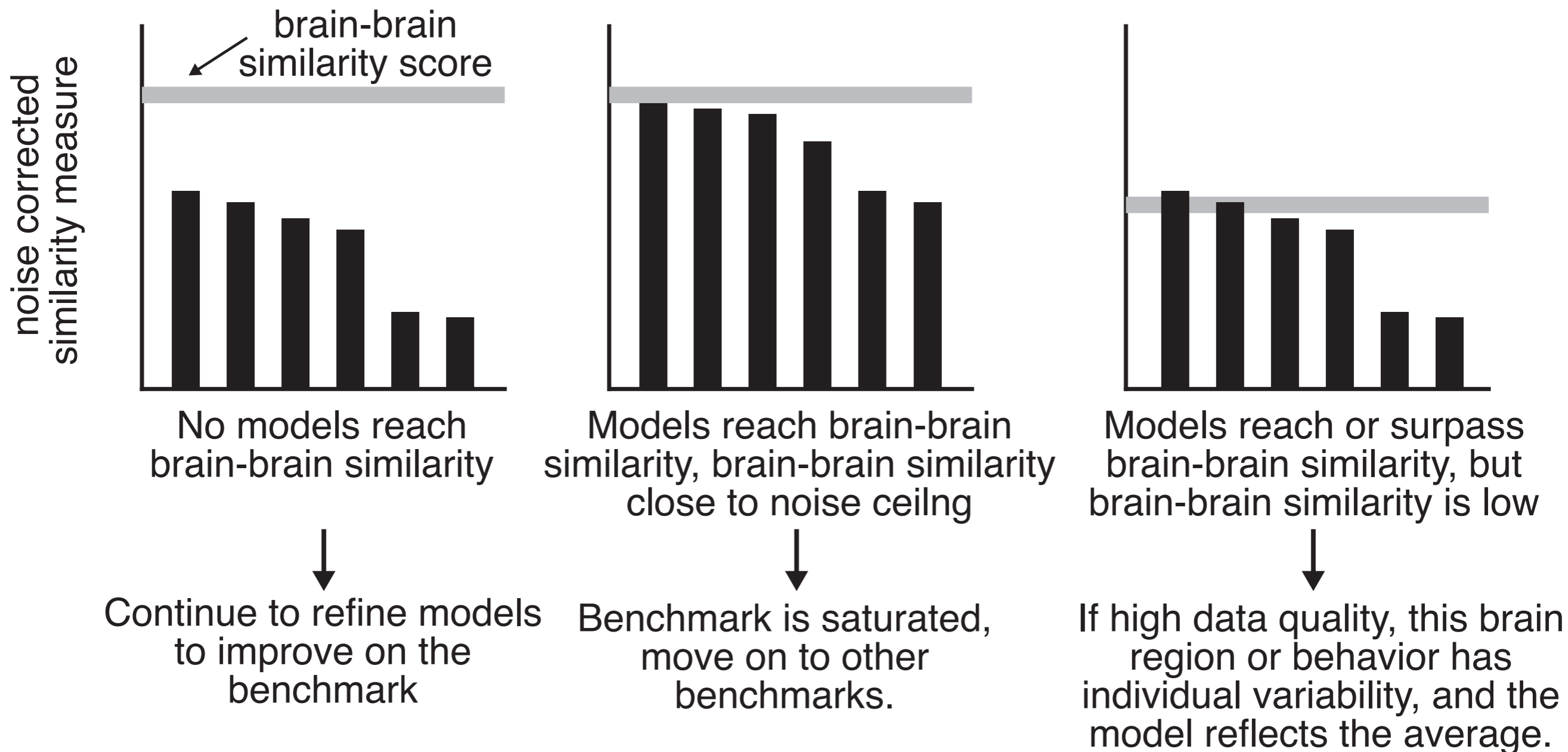


The NeuroAI Turing Test

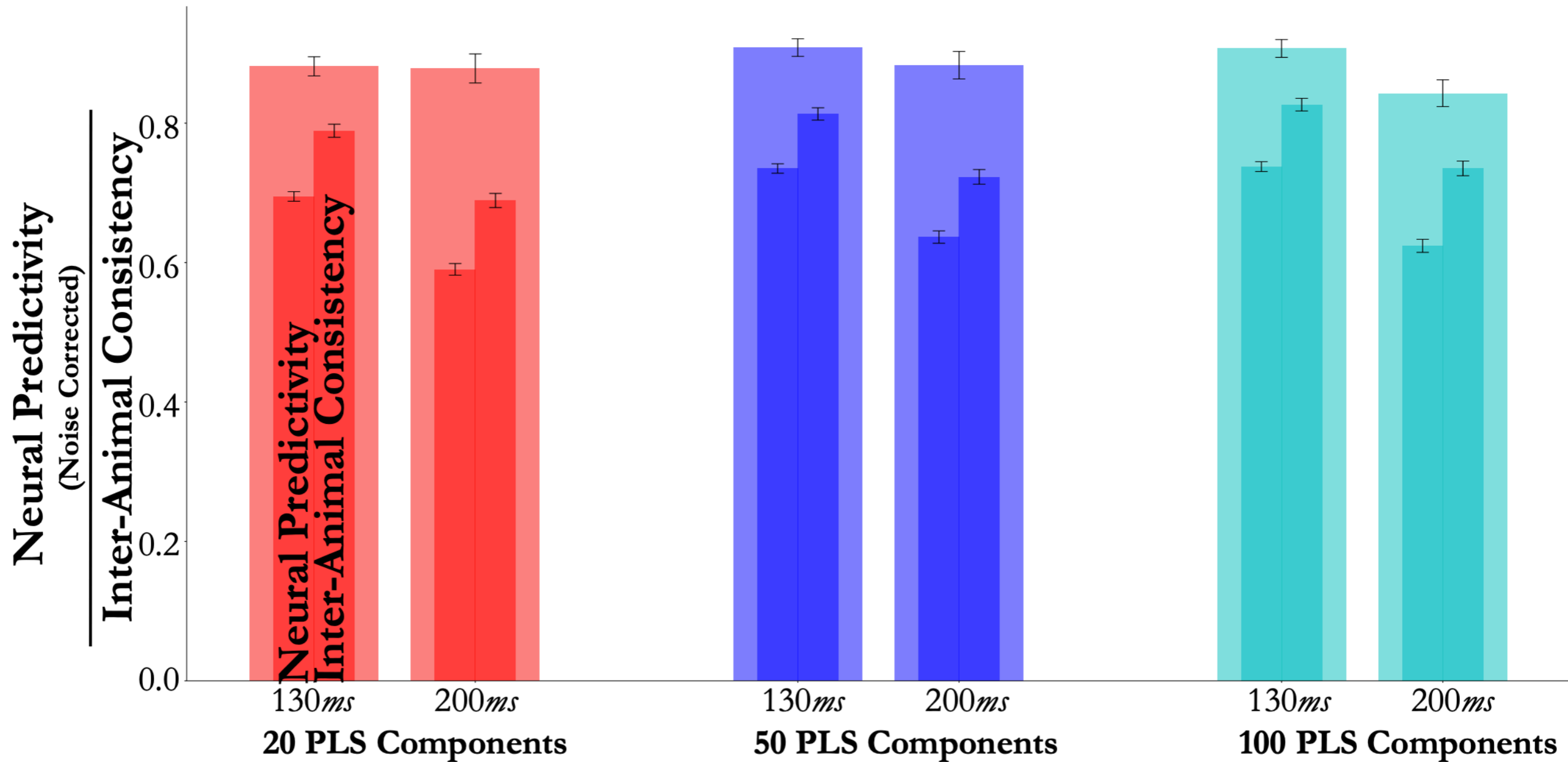


NeuroAI Turing Test: Possible Scenarios

NeuroAI Turing Test: Possible Scenarios



NeuroAI Turing Test: An Example in 200 ms Vision

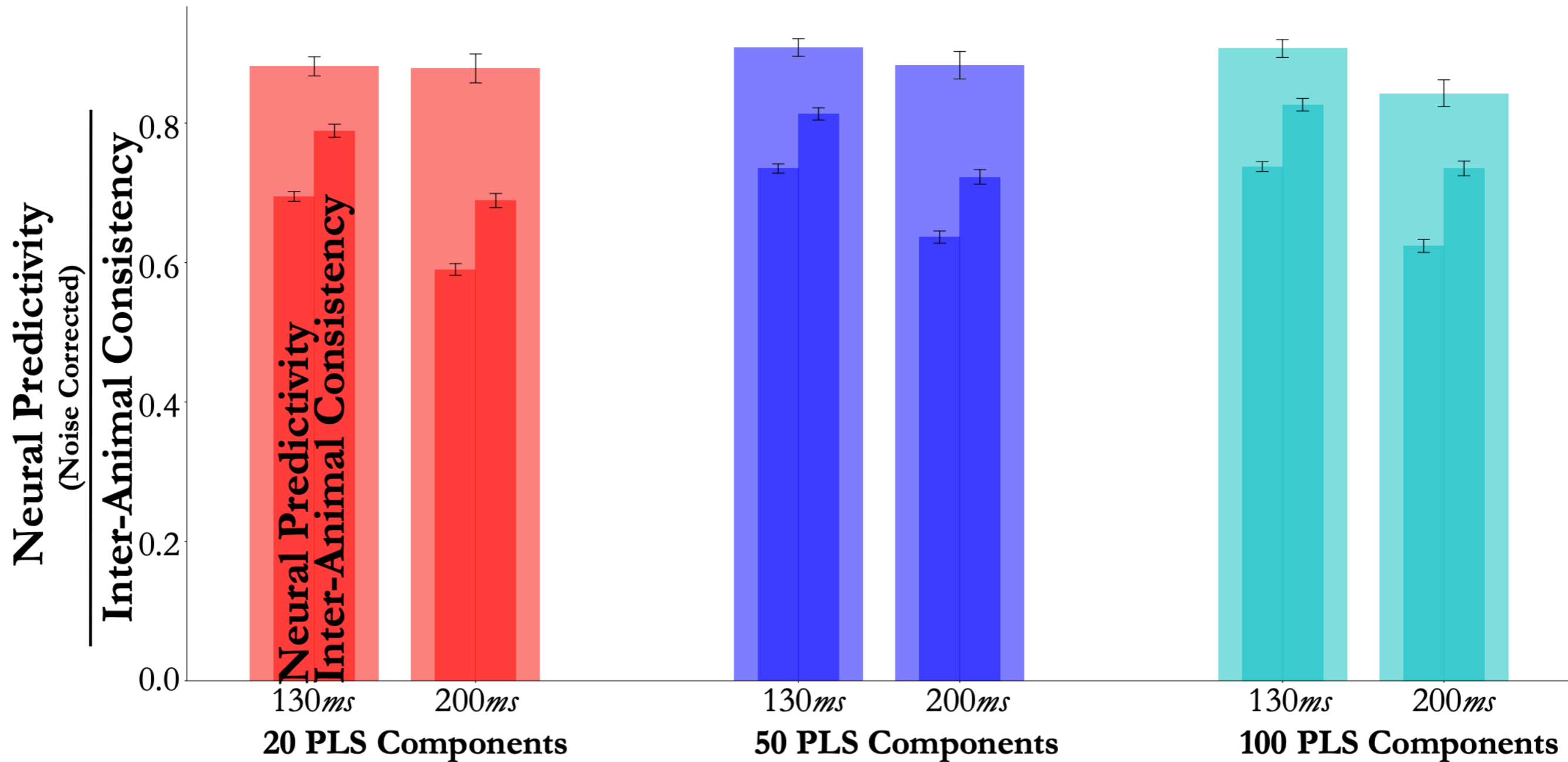


Recurrent Connections in the Primate Ventral Visual Stream Mediate a Tradeoff Between Task Performance and Network Size During Core Object Recognition

Aran Nayebi^{1,*}, Javier Sagastuy-Brena¹, Daniel M. Bear¹, Kohitij Kar², Jonas Kubilius^{2,3}, Surya Ganguli¹, David Sussillo¹, James J. DiCarlo², and Daniel L. K. Yamins¹

NeuroAI Turing Test: An Example in 200 ms Vision

So we may have achieved it for the classic “HvM” dataset (it’s tapped out!)



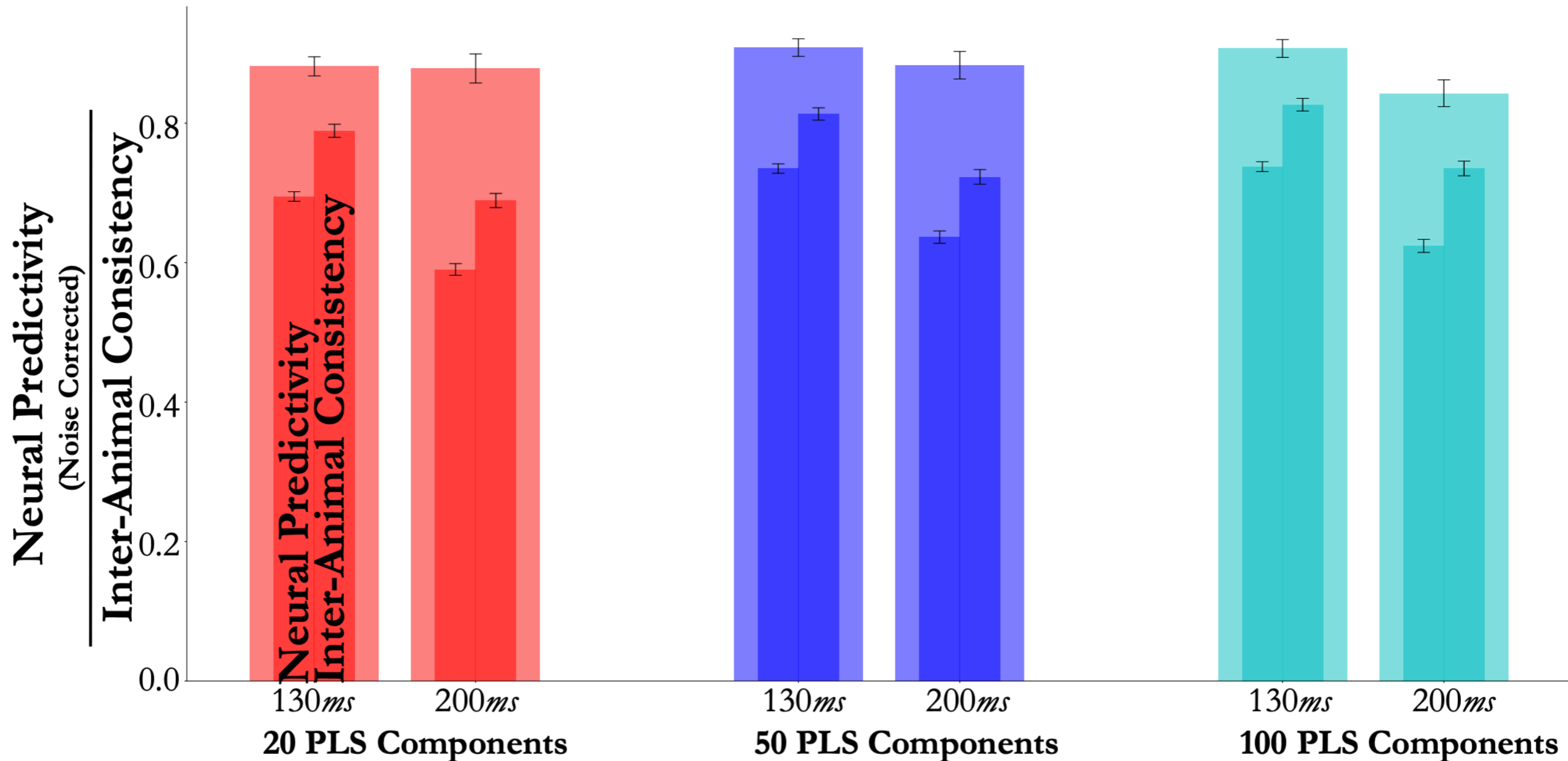
Recurrent Connections in the Primate Ventral Visual Stream Mediate a Tradeoff Between Task Performance and Network Size During Core Object Recognition

Aran Nayebi^{1,*}, Javier Sagastuy-Brena¹, Daniel M. Bear¹, Kohitij Kar², Jonas Kubilius^{2,3}, Surya Ganguli¹, David Sussillo¹, James J. DiCarlo², and Daniel L. K. Yamins¹

NeuroAI Turing Test: An Example in 200 ms Vision

So we may have achieved it for the classic “HvM” dataset (it’s tapped out!)

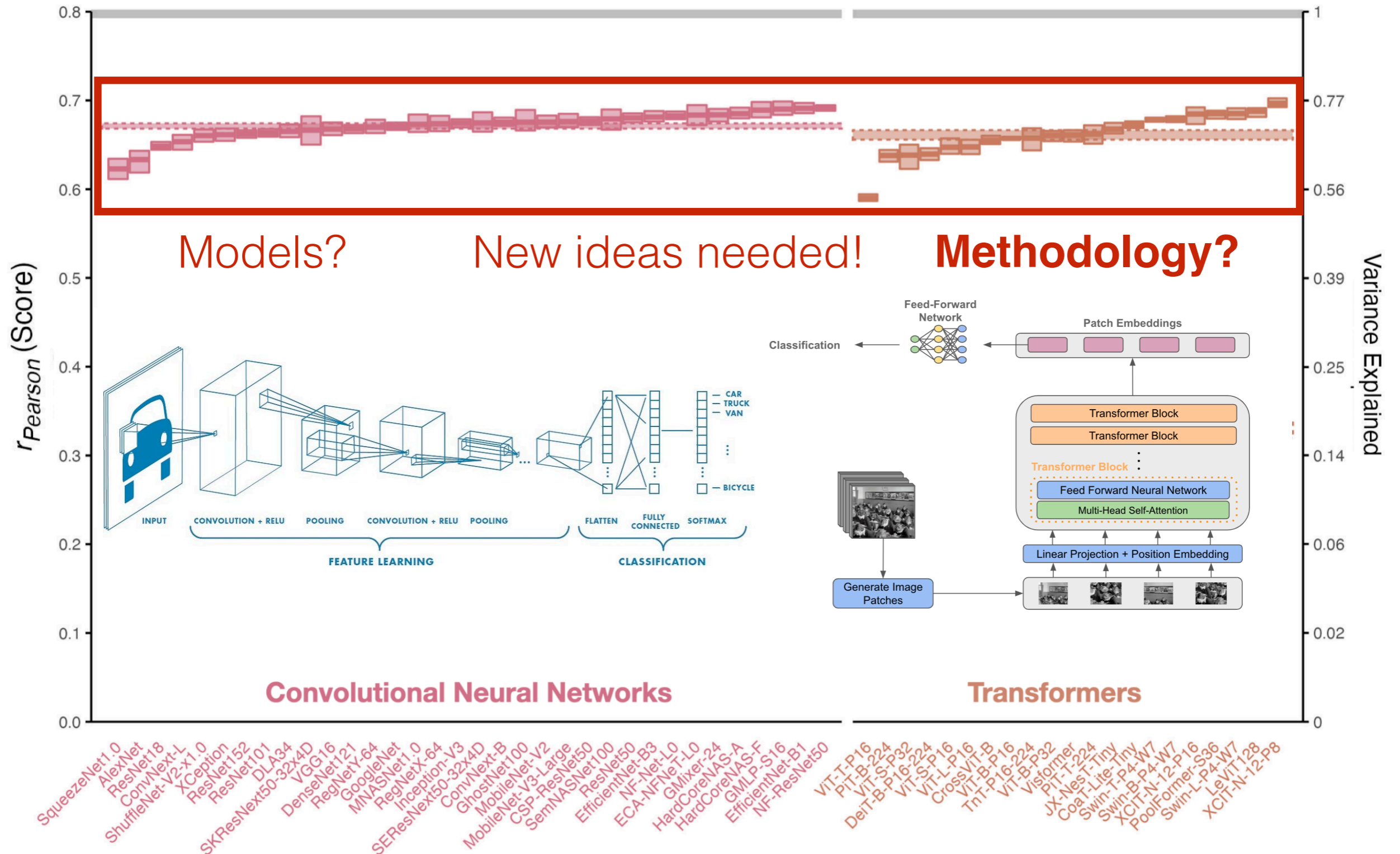
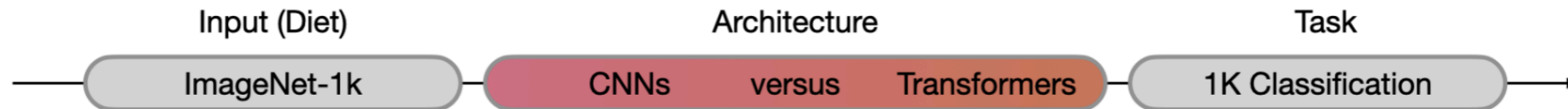
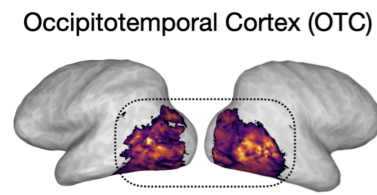
But I will show you that we haven’t yet in embodied intelligence...



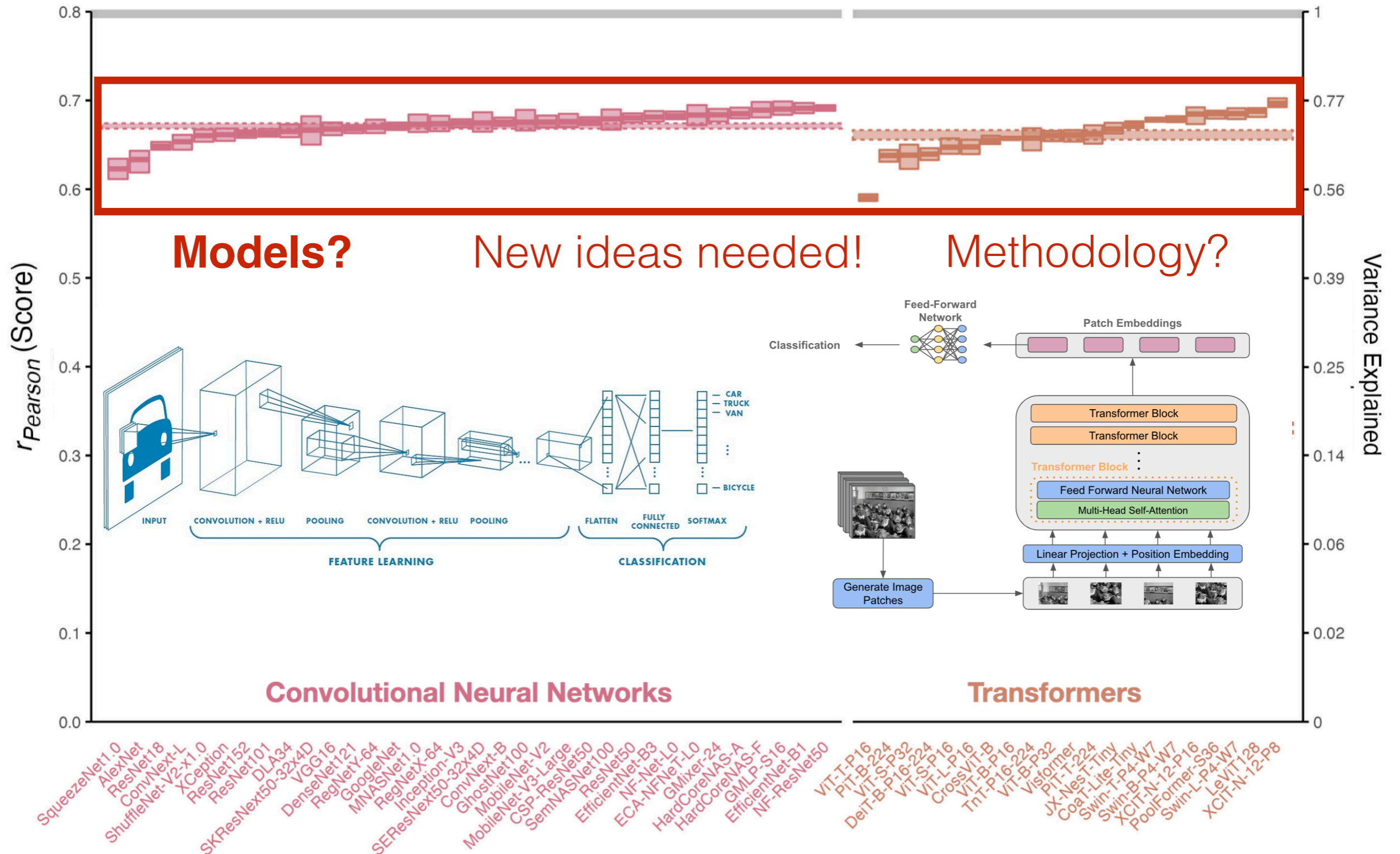
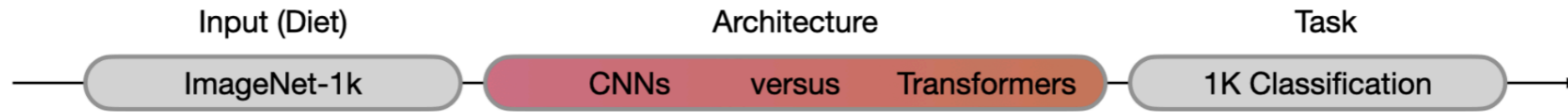
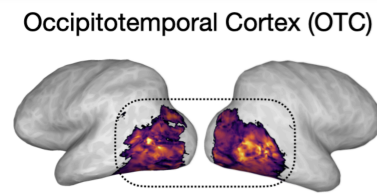
Recurrent Connections in the Primate Ventral Visual Stream Mediate a Tradeoff Between Task Performance and Network Size During Core Object Recognition

Aran Nayebi^{1,*}, Javier Sagastuy-Brena¹, Daniel M. Bear¹, Kohitij Kar², Jonas Kubilius^{2,3}, Surya Ganguli¹, David Sussillo¹, James J. DiCarlo², and Daniel L. K. Yamins¹

Similar predictivities between CNNs vs. Transformers



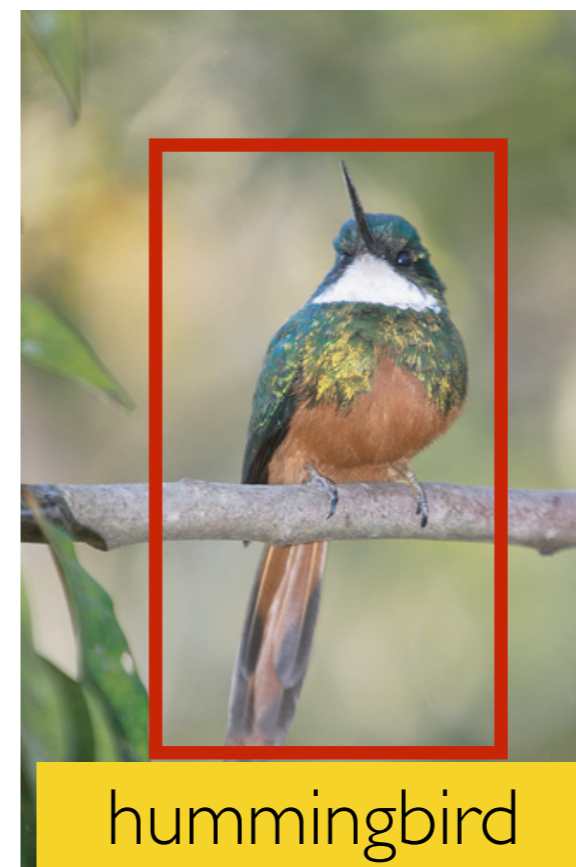
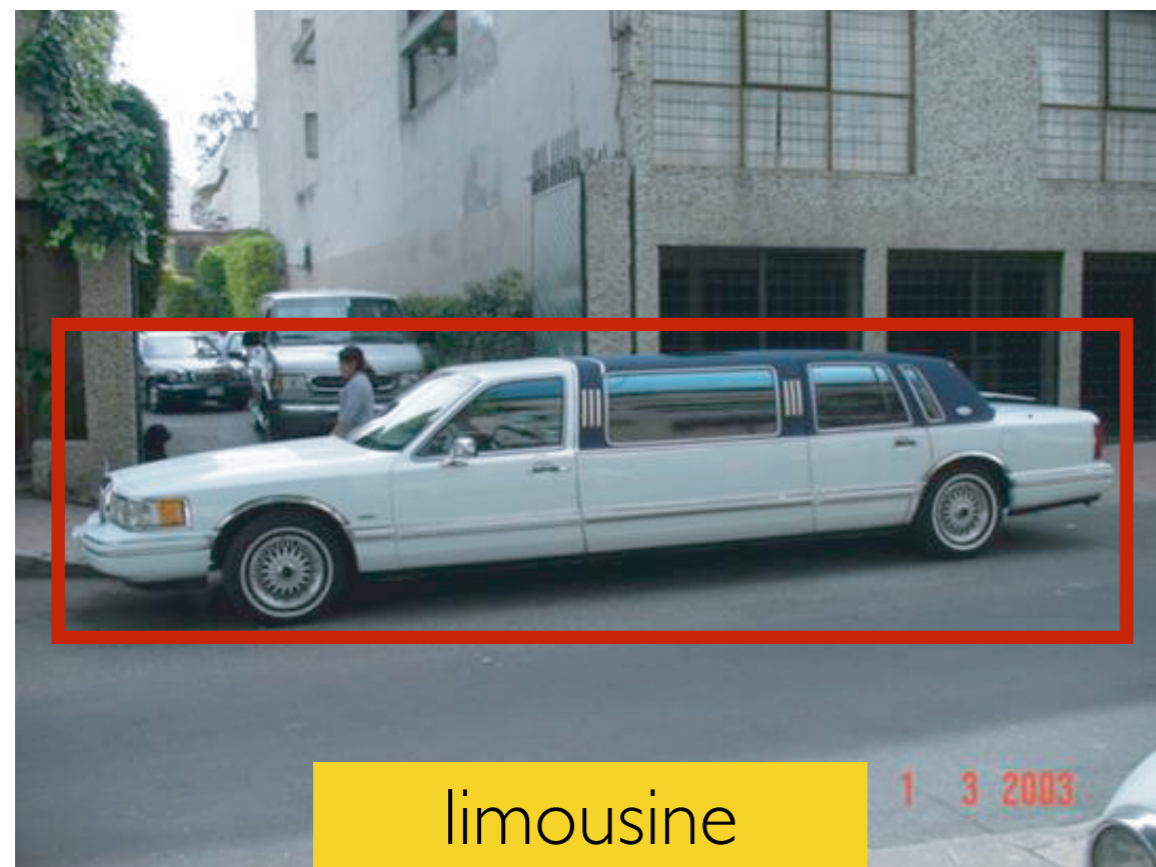
Similar predictivities between CNNs vs. Transformers



We do a lot more than passive viewing...



We do a lot more than passive viewing...



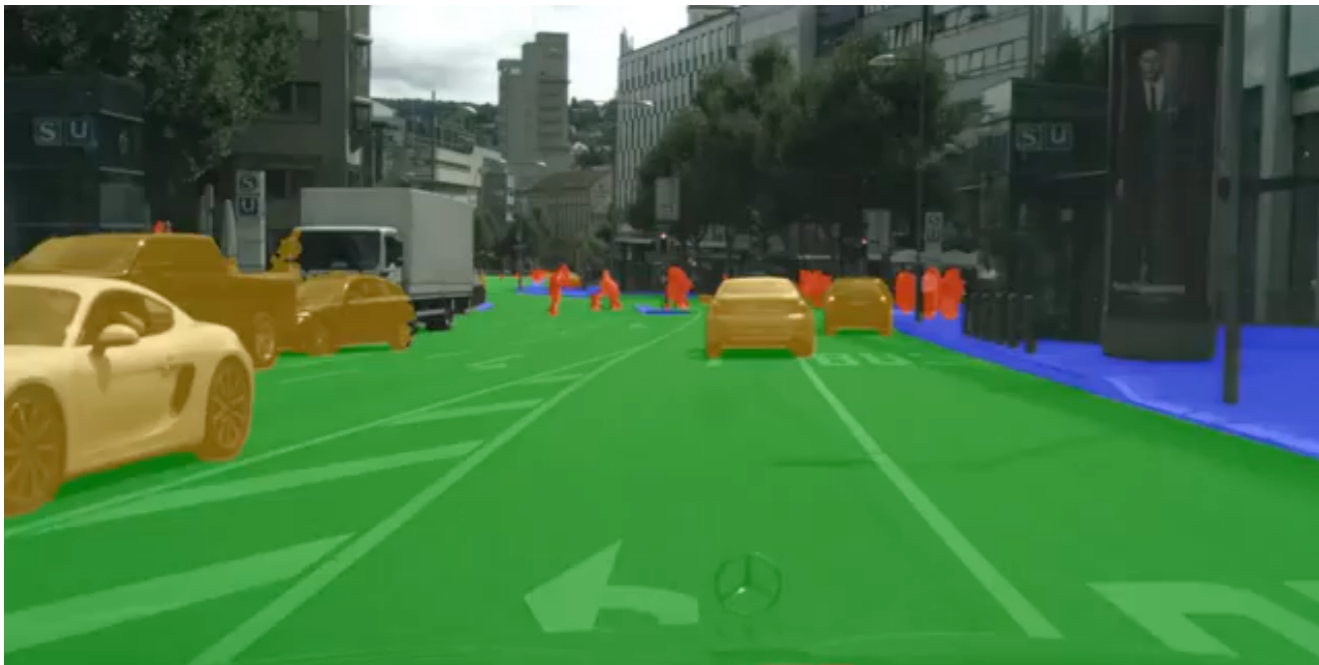
We do a lot more than passive viewing...

Scene Understanding



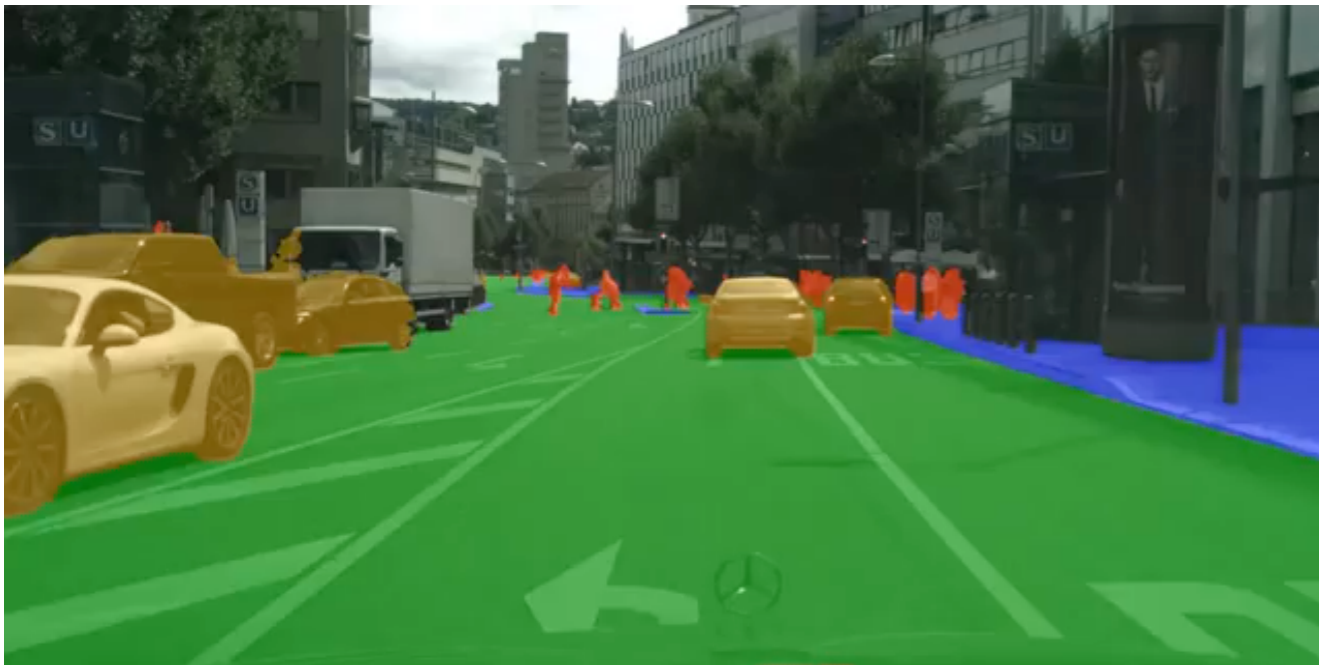
We do a lot more than passive viewing...

Scene Understanding



We do a lot more than passive viewing...

Scene Understanding

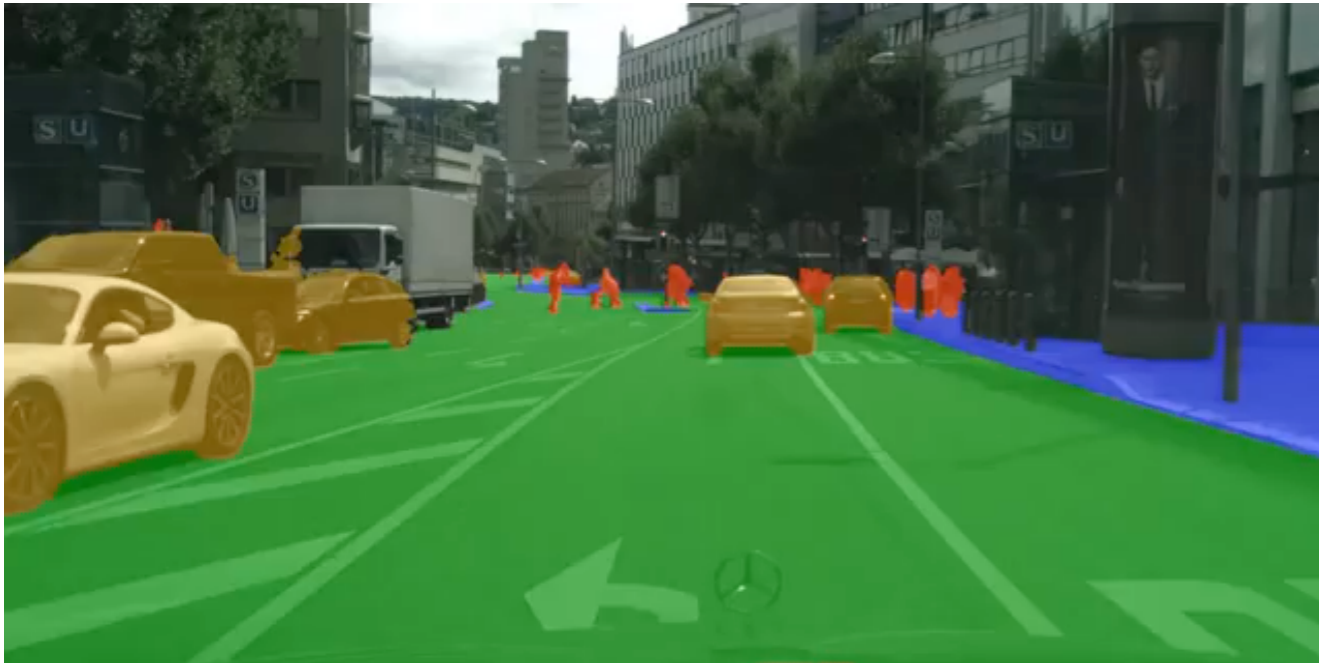


Multi-Step Planning



We do a lot more than passive viewing...

Scene Understanding

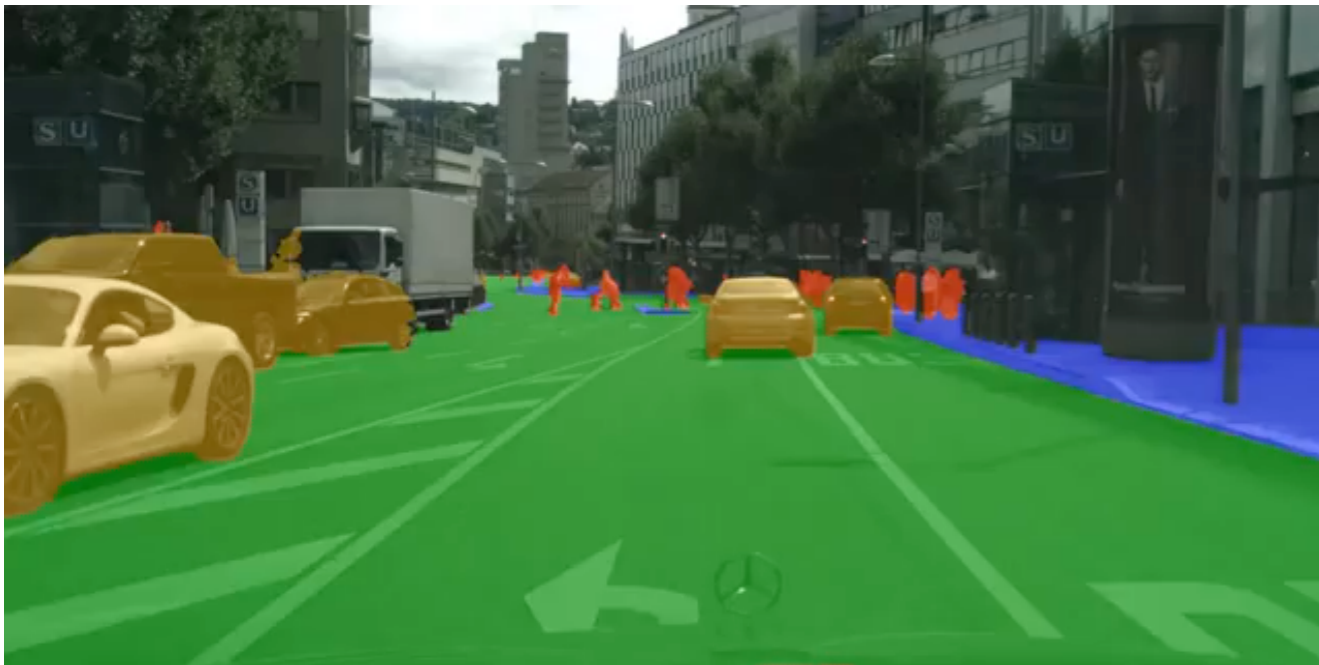


Multi-Step Planning



We do a lot more than passive viewing...

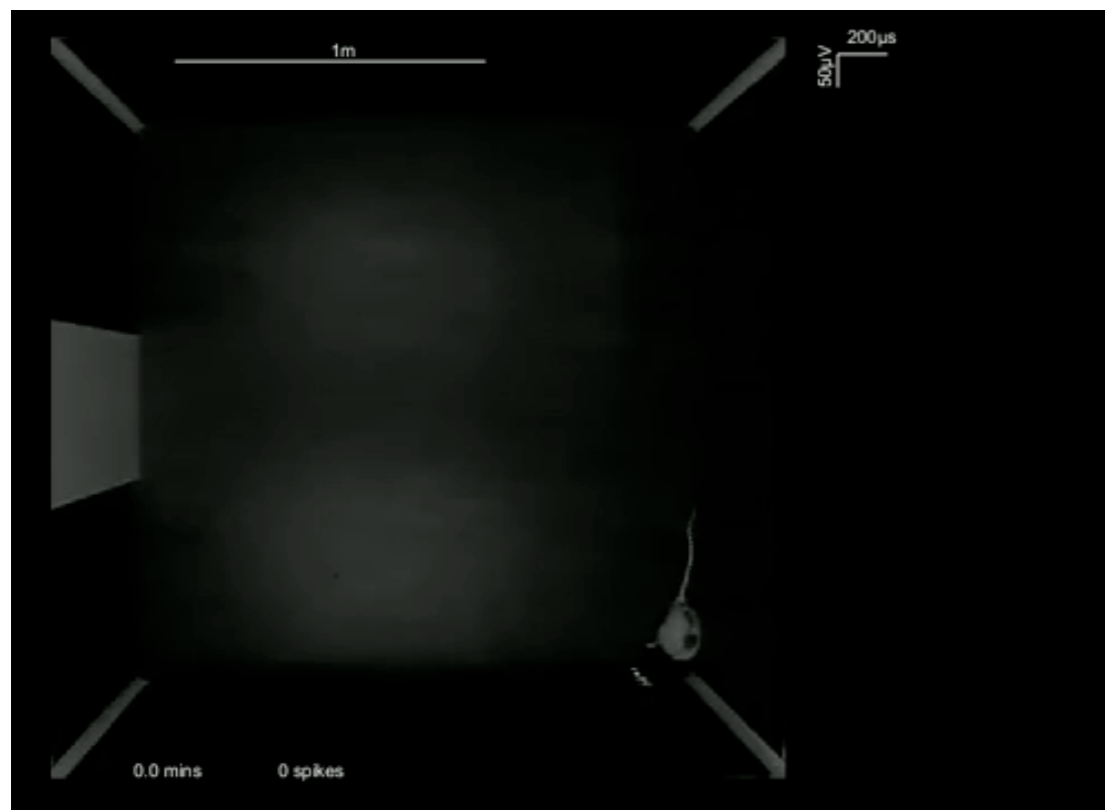
Scene Understanding



Multi-Step Planning



Navigation



We do a lot more than passive viewing...

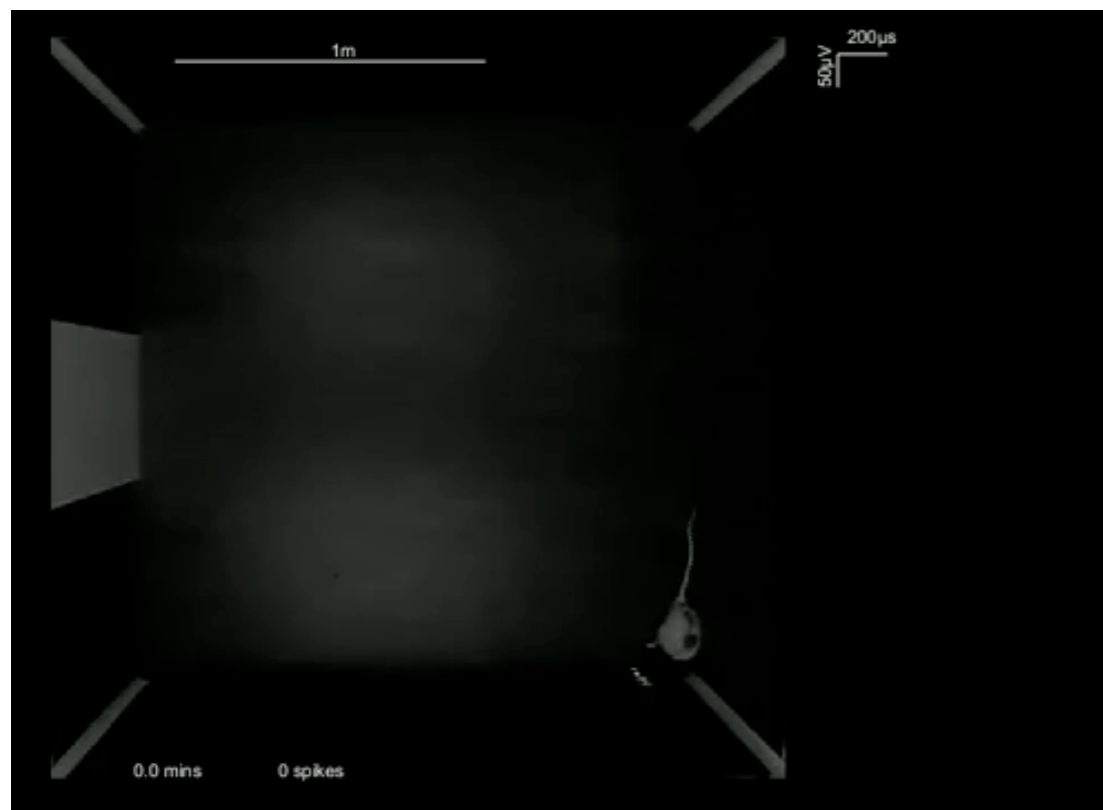
Scene Understanding



Multi-Step Planning

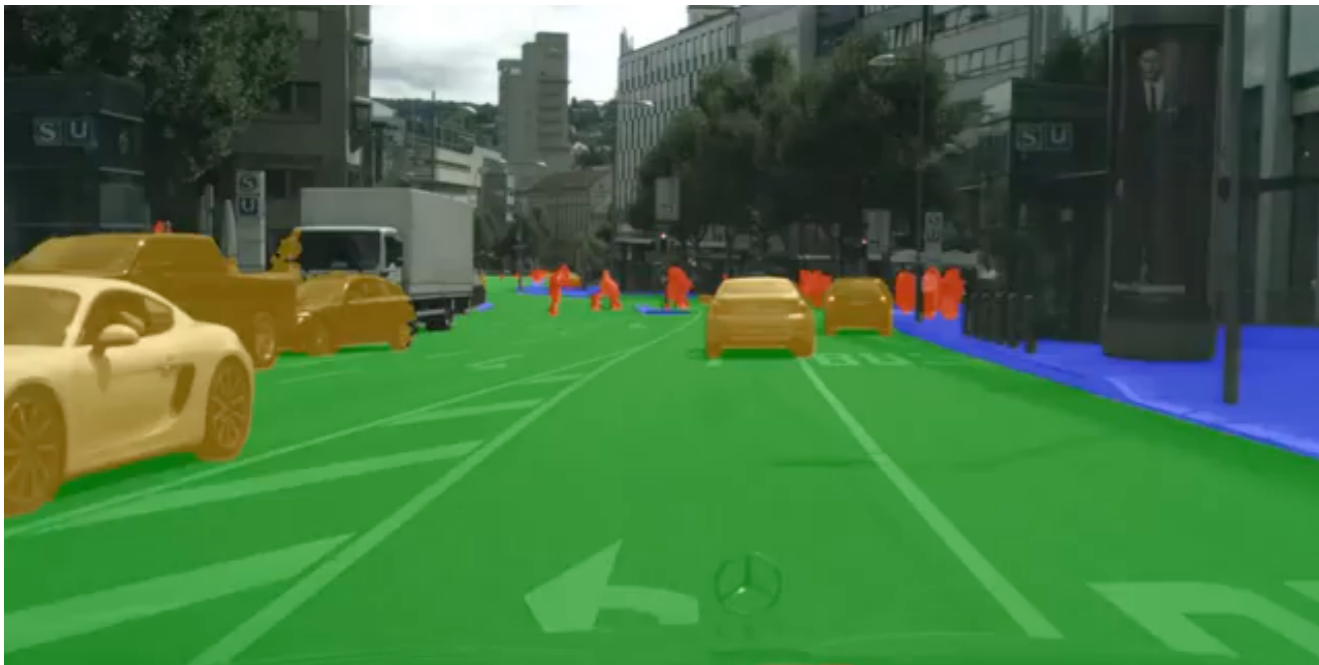


Navigation



We do a lot more than passive viewing...

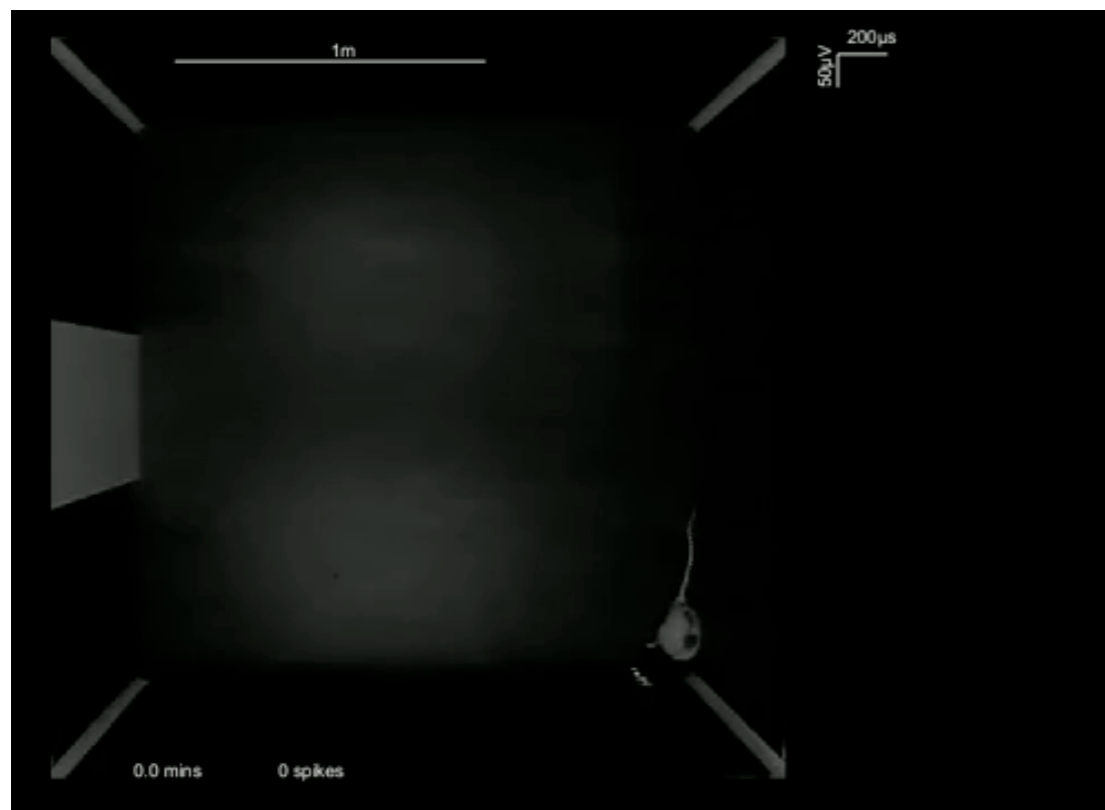
Scene Understanding



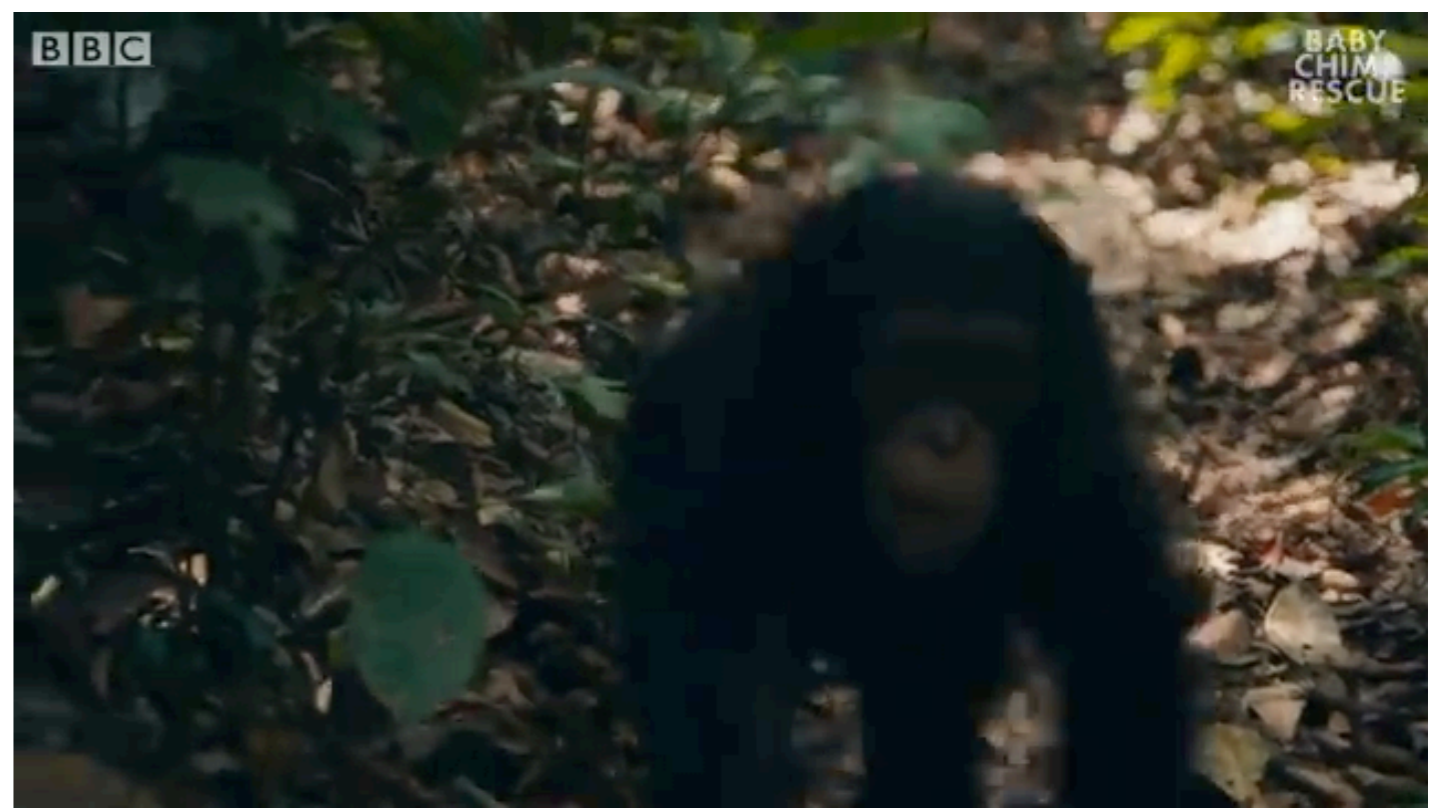
Multi-Step Planning



Navigation



Flexible Embodiment



We do a lot more than passive viewing...

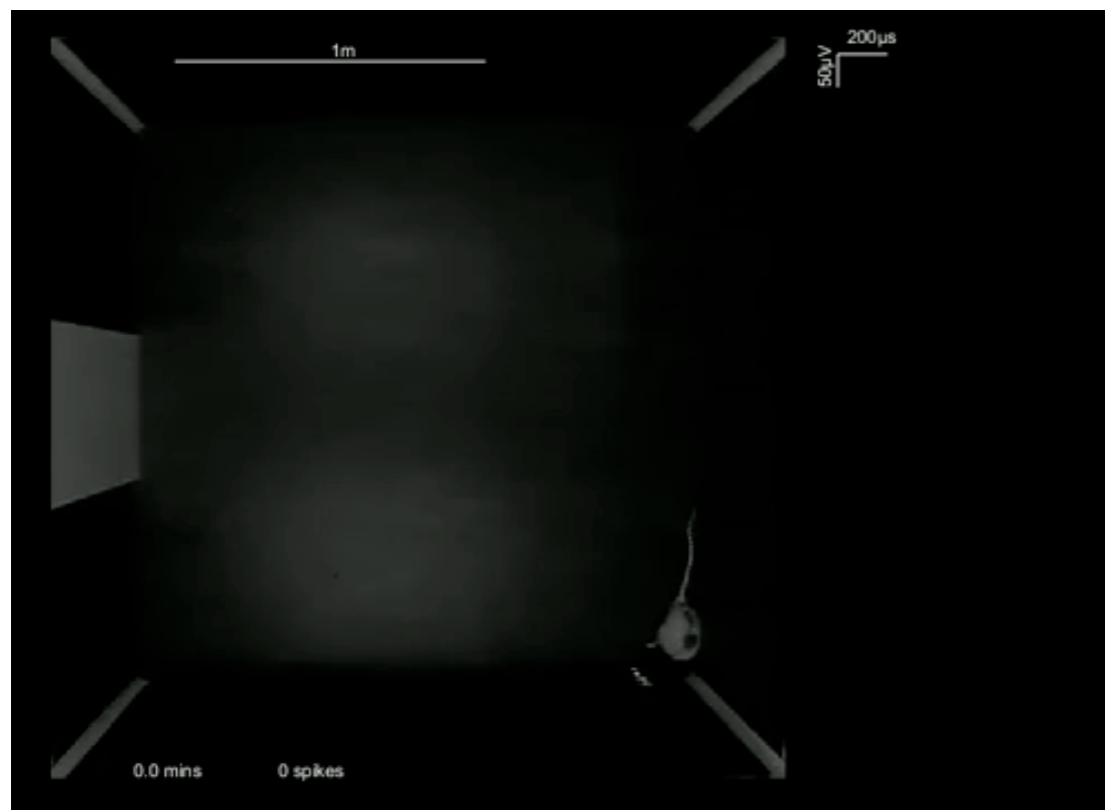
Scene Understanding



Multi-Step Planning



Navigation



Flexible Embodiment



We do a lot more than passive viewing...

Scene Understanding

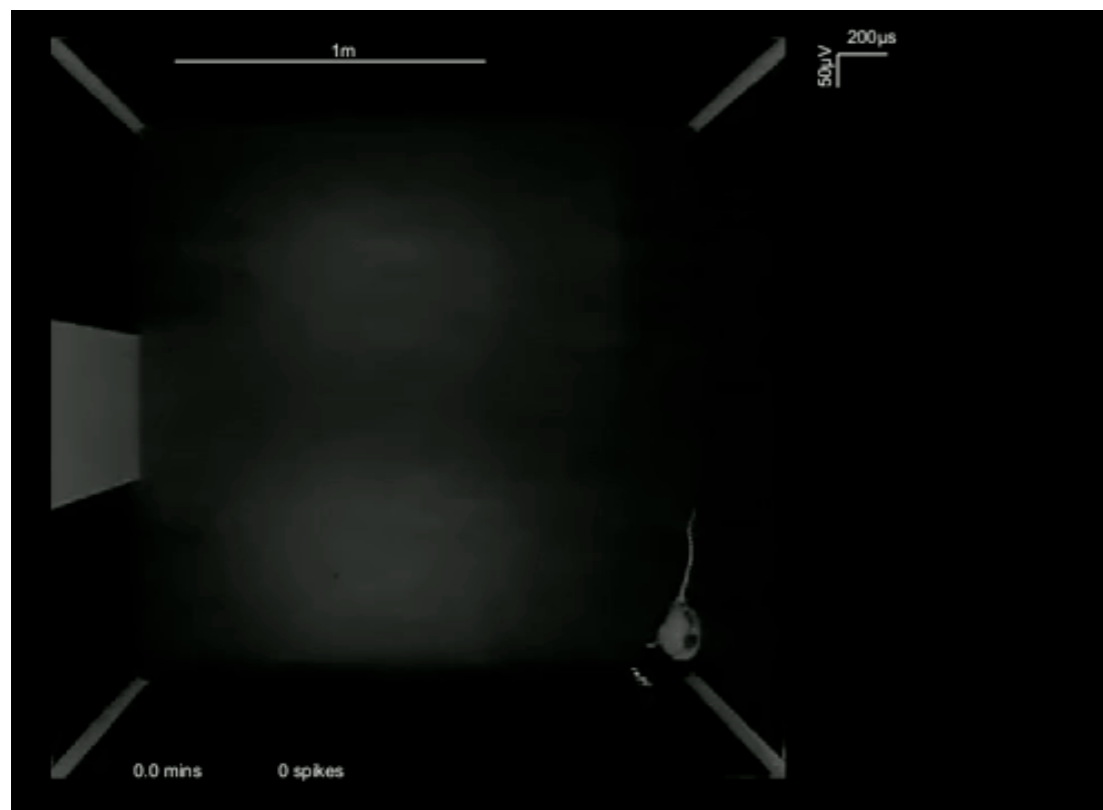


Multi-Step Planning

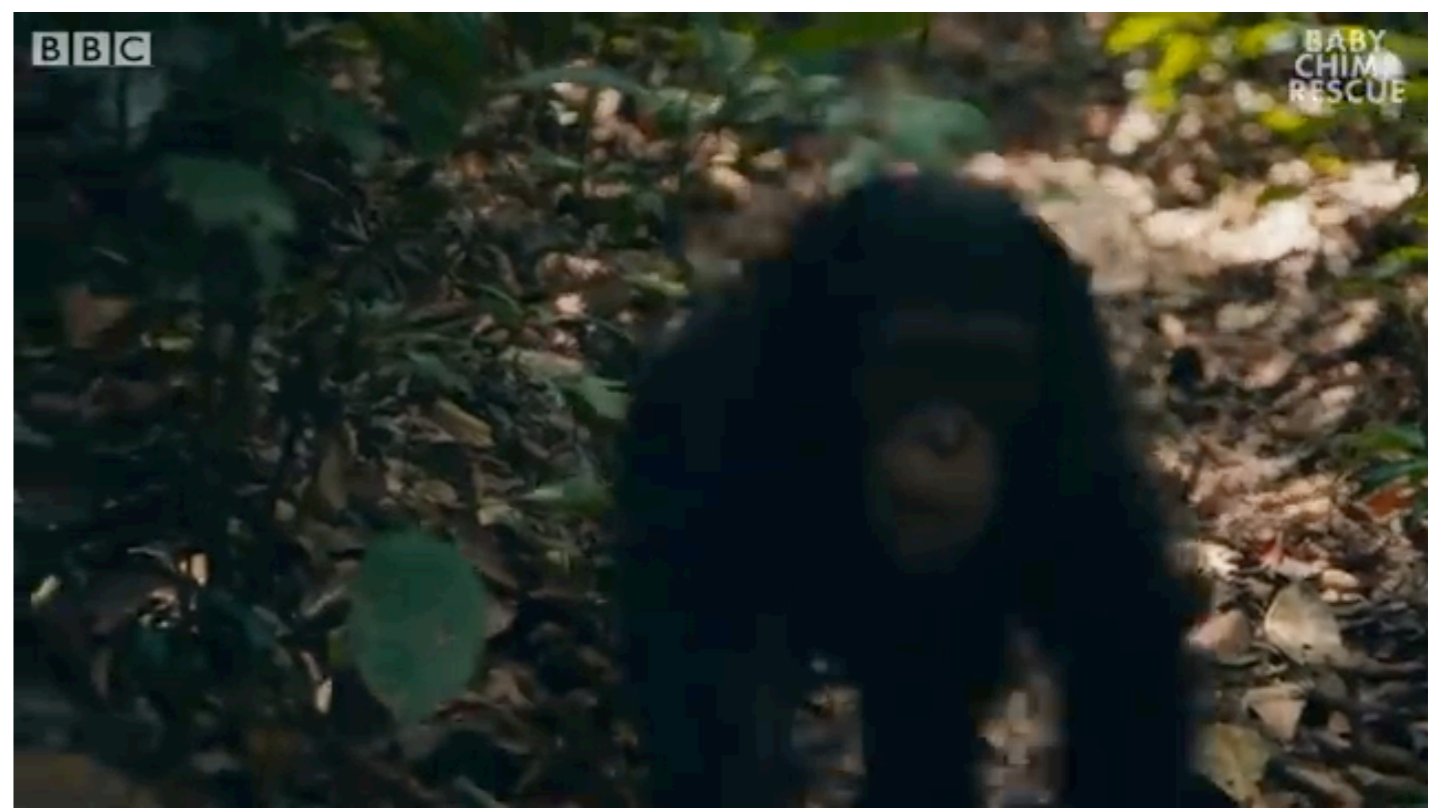


All of these behaviors are done in a body!

Navigation



Flexible Embodiment



We do a lot more than passive viewing...

Scene Understanding

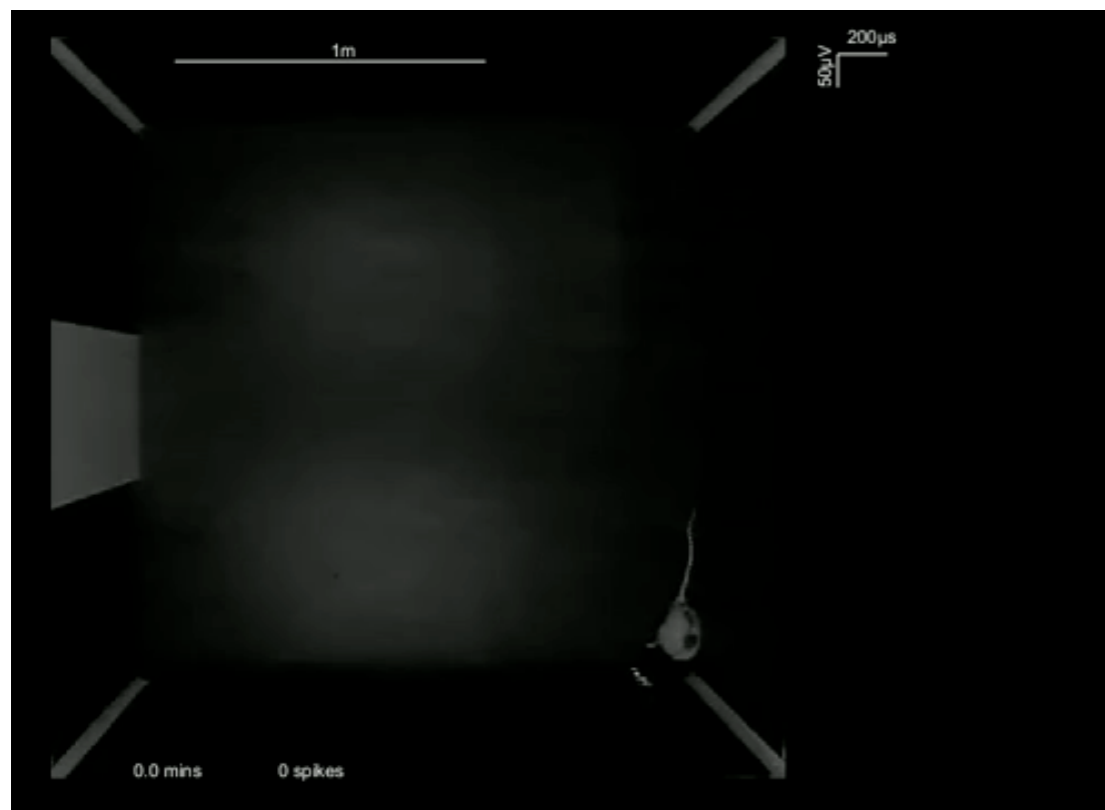


Multi-Step Planning



All of these behaviors are done in a body!

Navigation



Flexible Embodiment



We do a lot more than passive viewing...

Scene Understanding

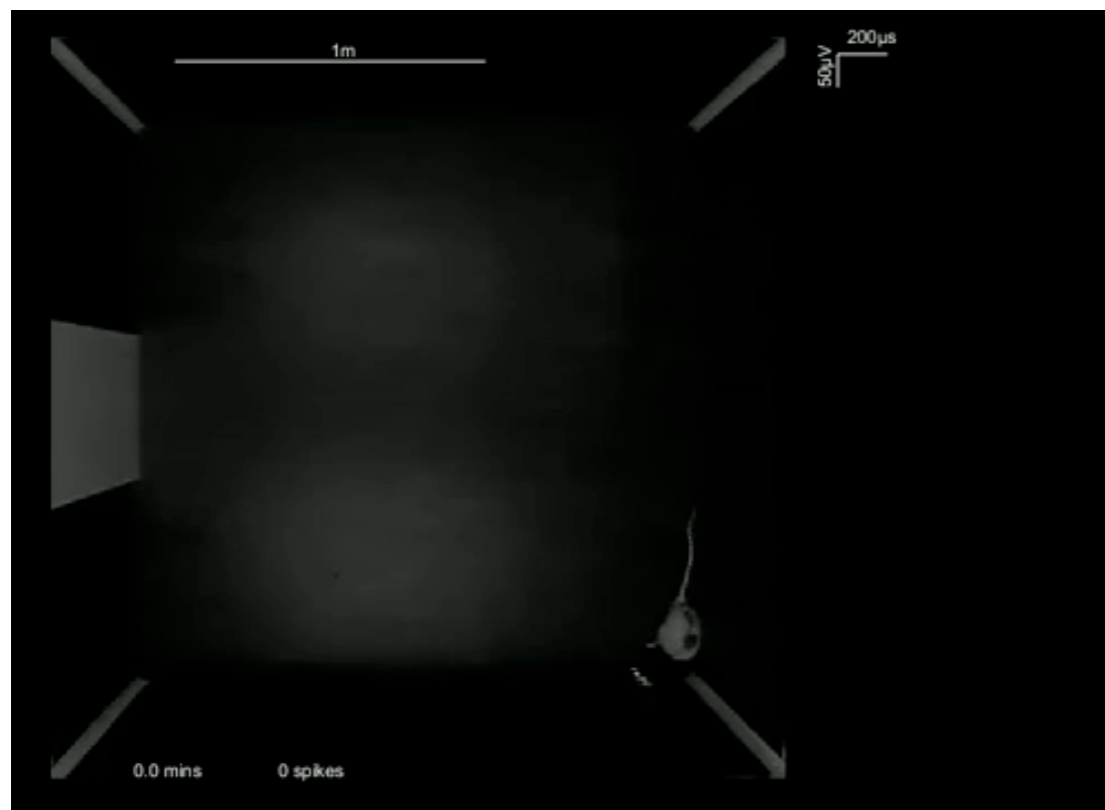


Multi-Step Planning



What are the core design principles that give rise to these abilities?

Navigation

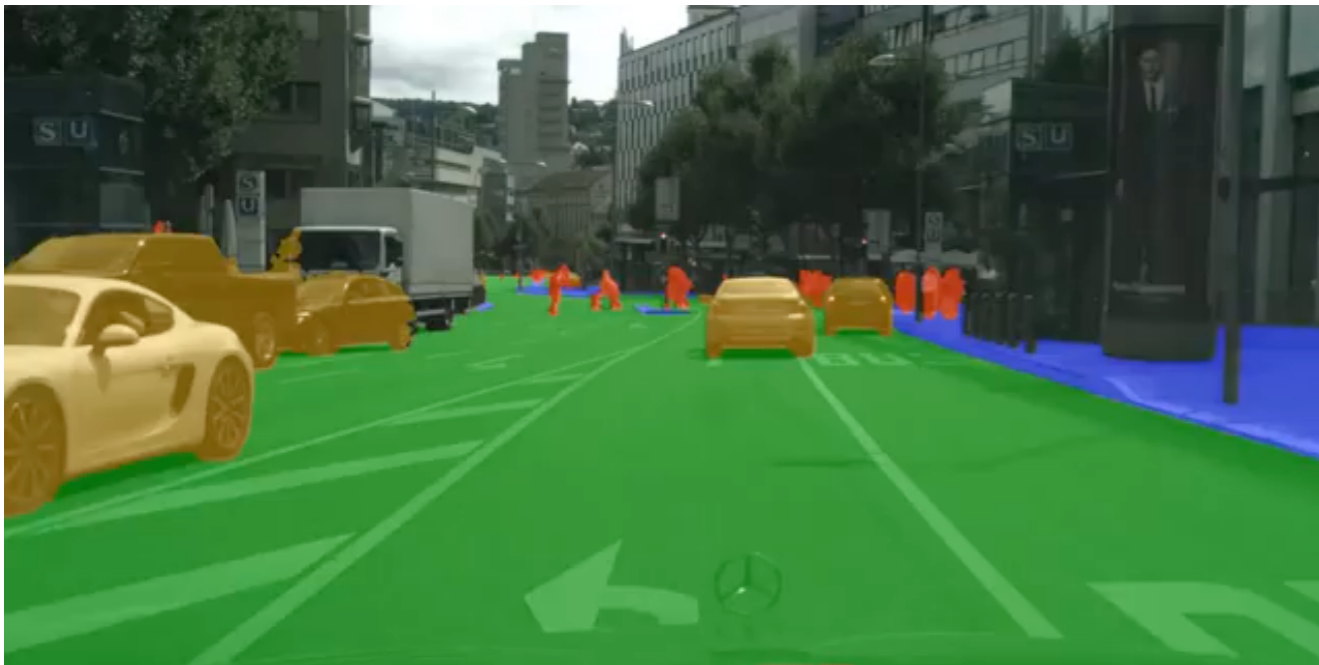


Flexible Embodiment



We do a lot more than passive viewing...

Scene Understanding

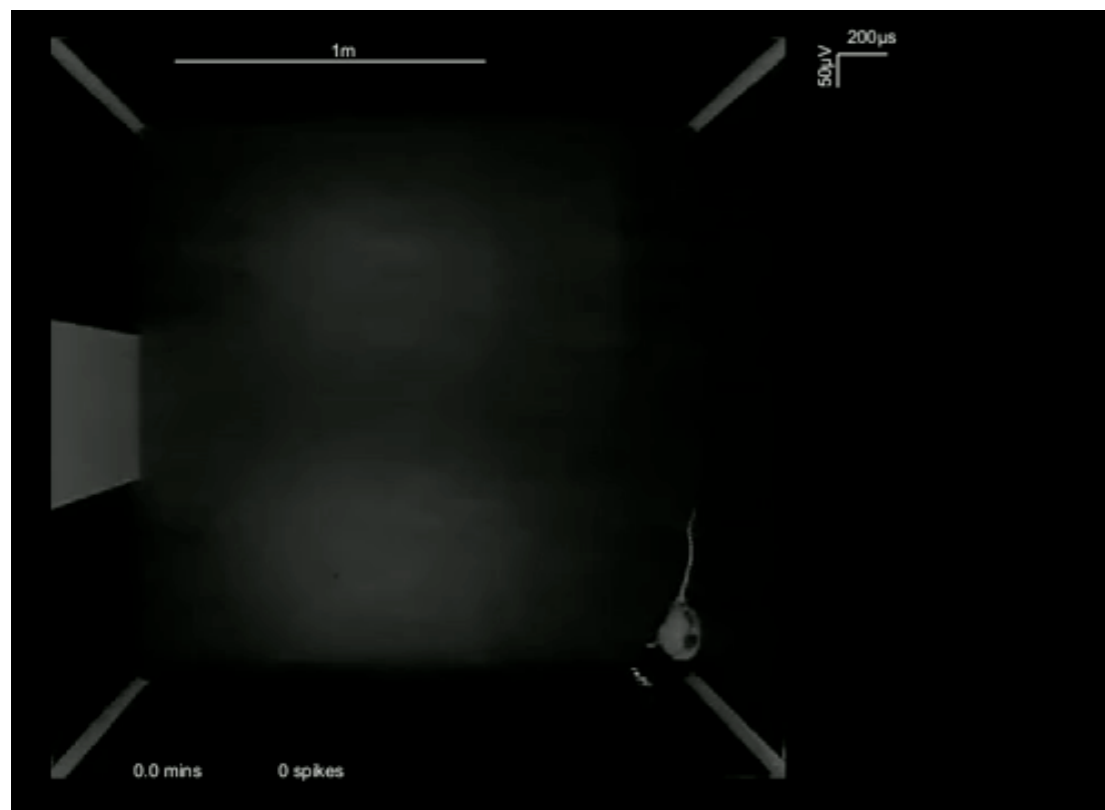


Multi-Step Planning



What are the core design principles that give rise to these abilities?

Navigation



Flexible Embodiment



We do a lot more than passive viewing...

Scene Understanding

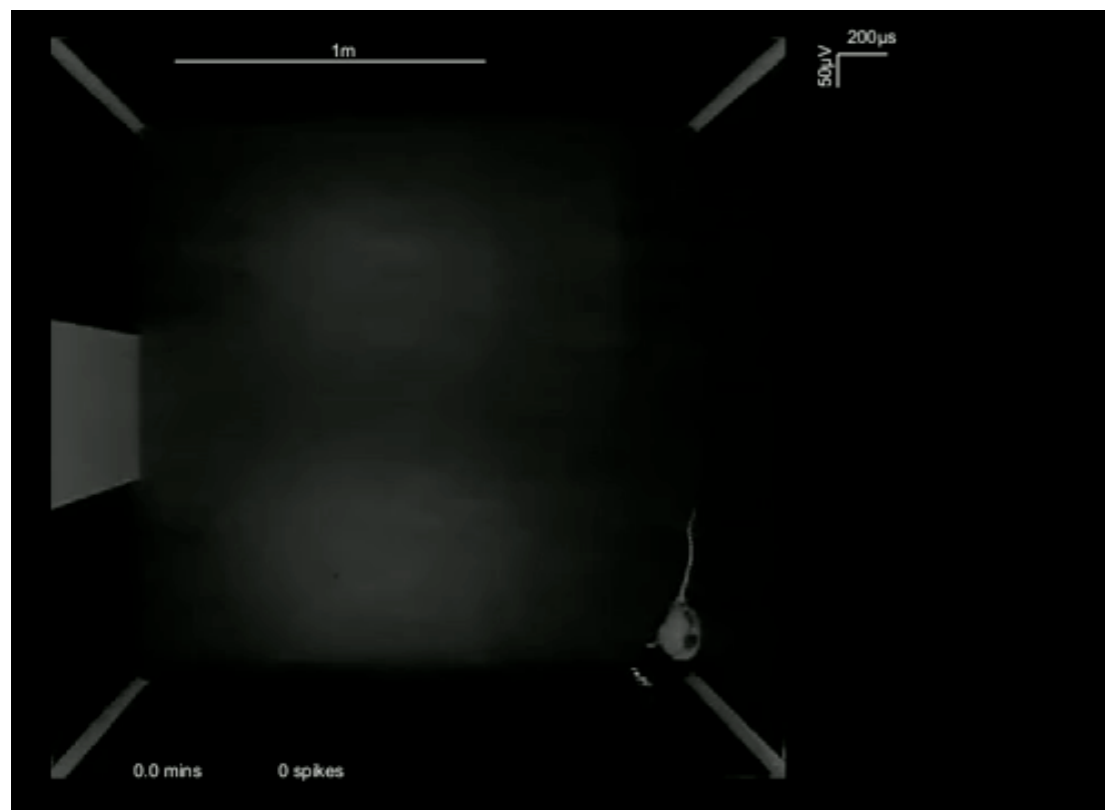


Multi-Step Planning



What are the core design principles that give rise to these abilities?

Navigation



Flexible Embodiment



We do a lot more than passive viewing...

Scene Understanding

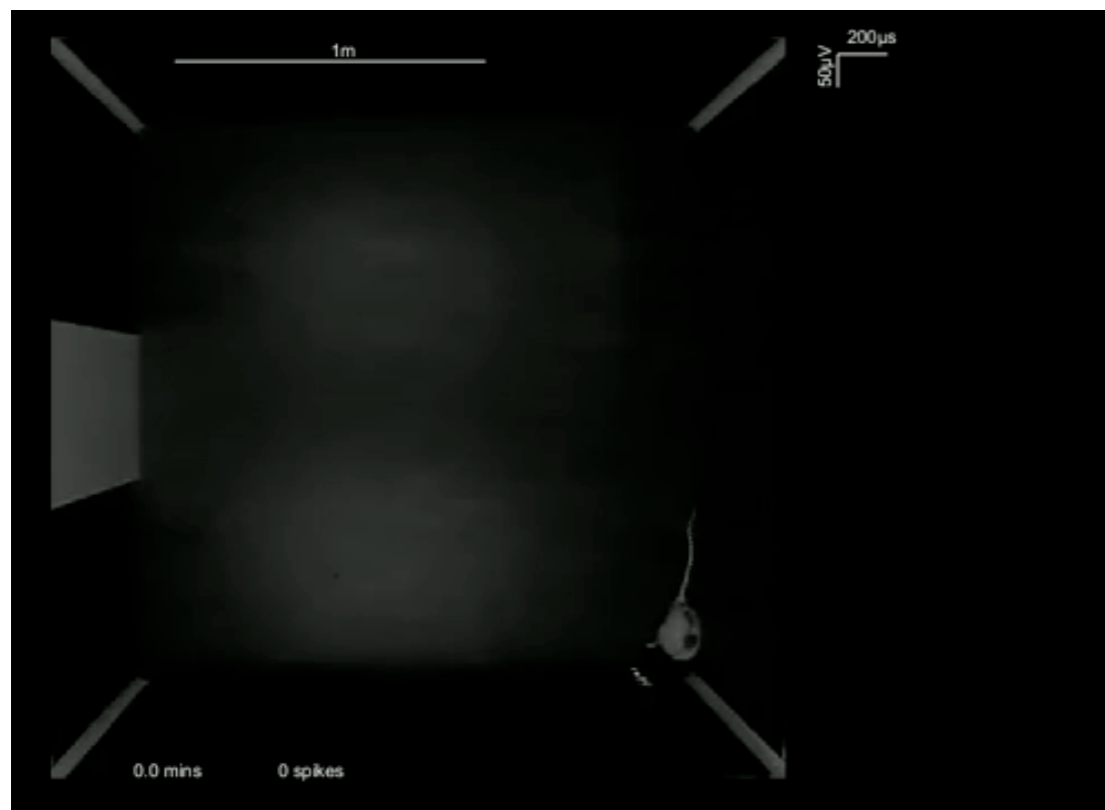


Multi-Step Planning

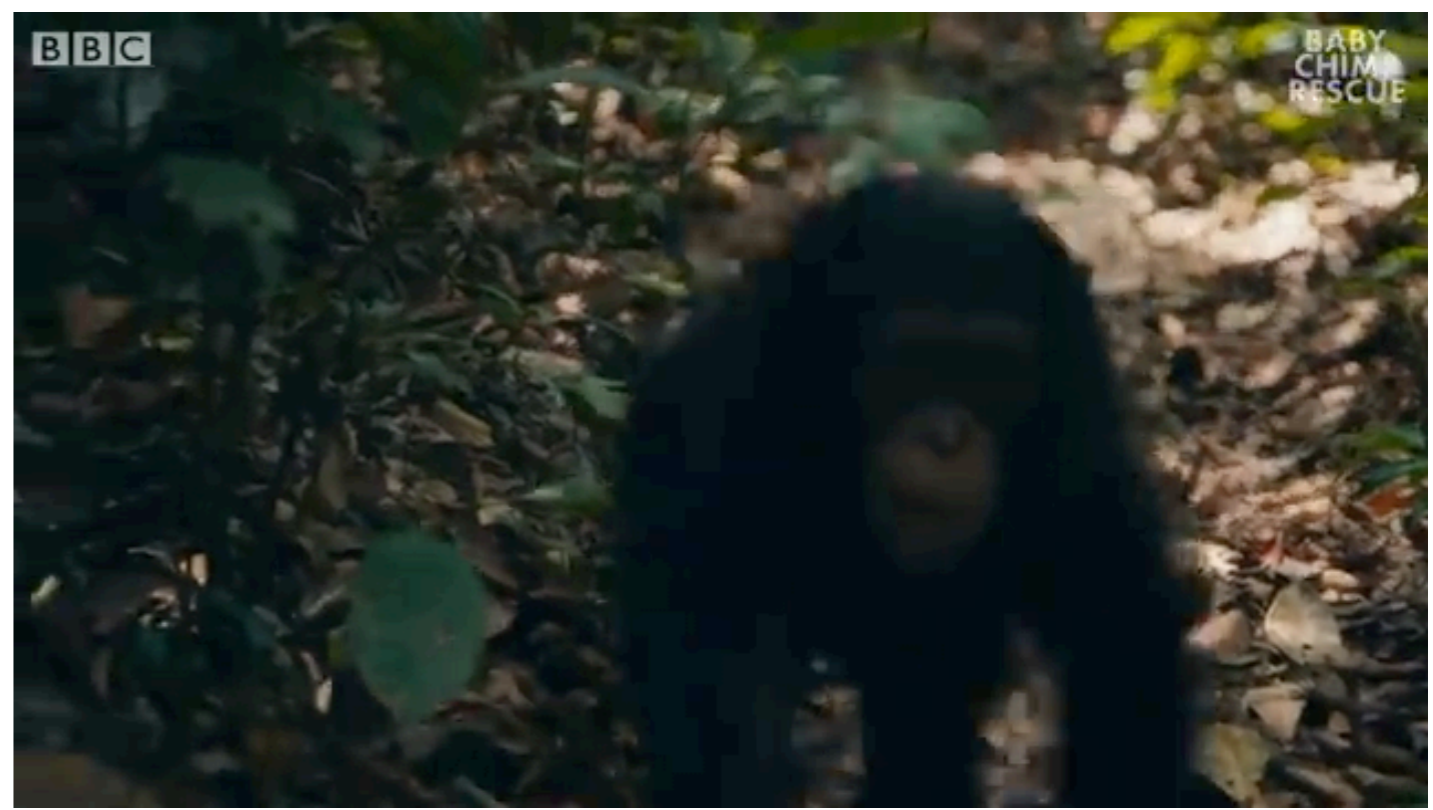


What are the core design principles that give rise to these abilities?

Navigation

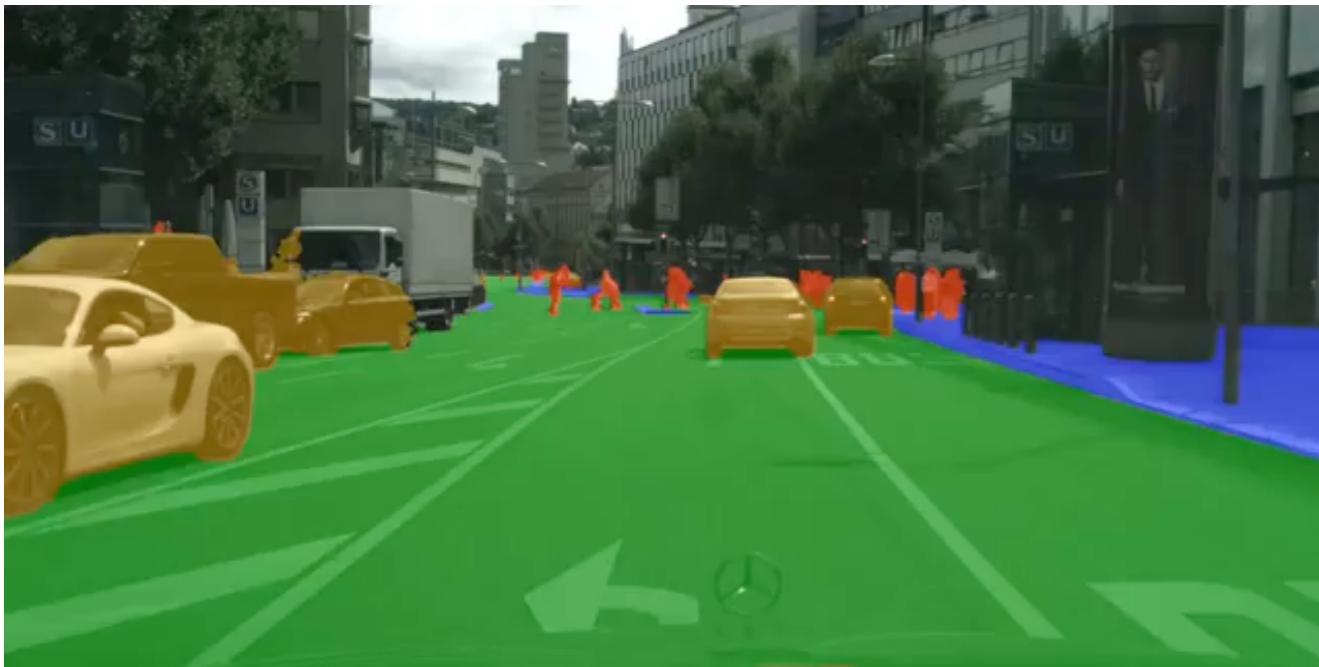


Flexible Embodiment



We do a lot more than passive viewing...

Scene Understanding

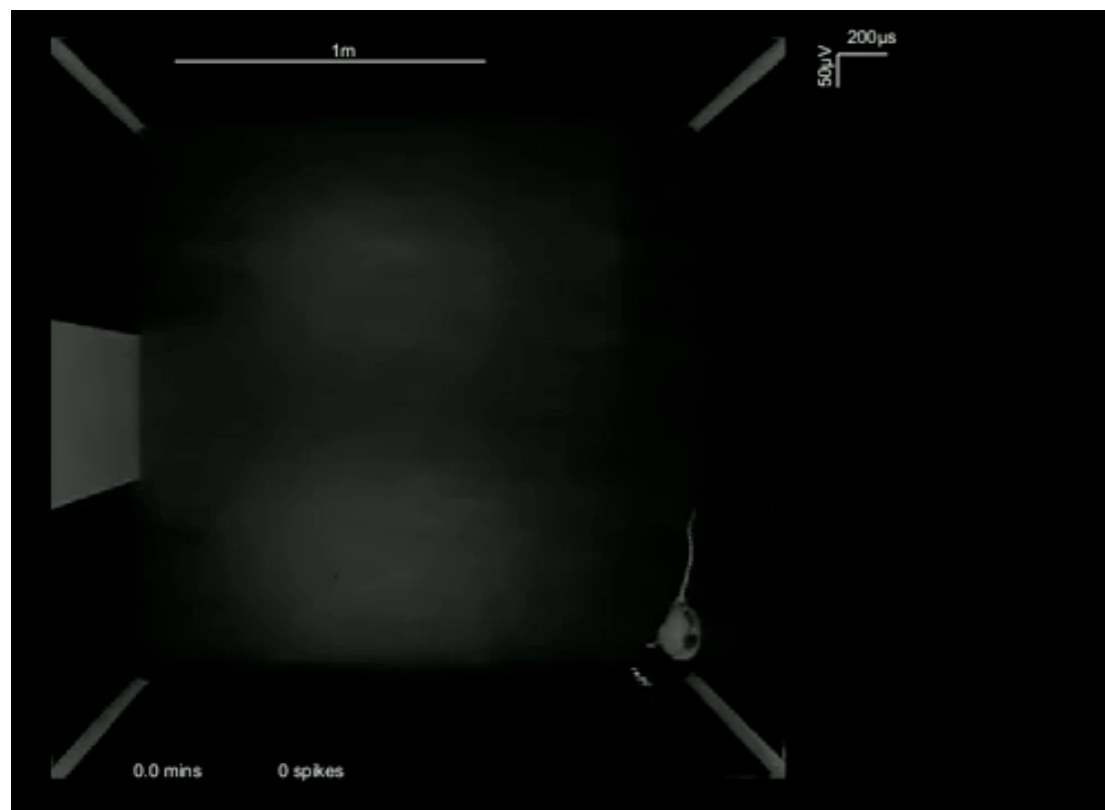


Multi-Step Planning

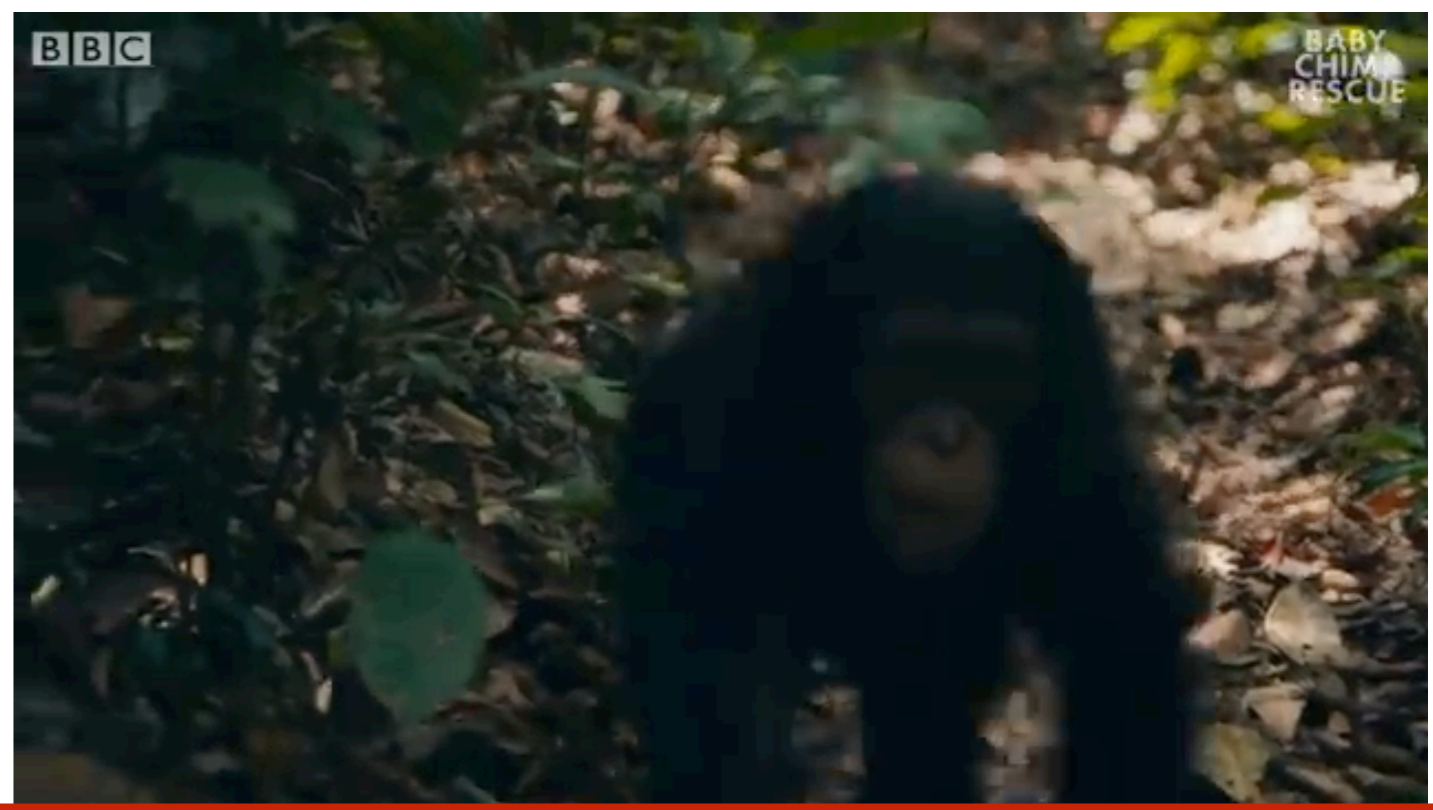


What are the core design principles that give rise to these abilities?

Navigation

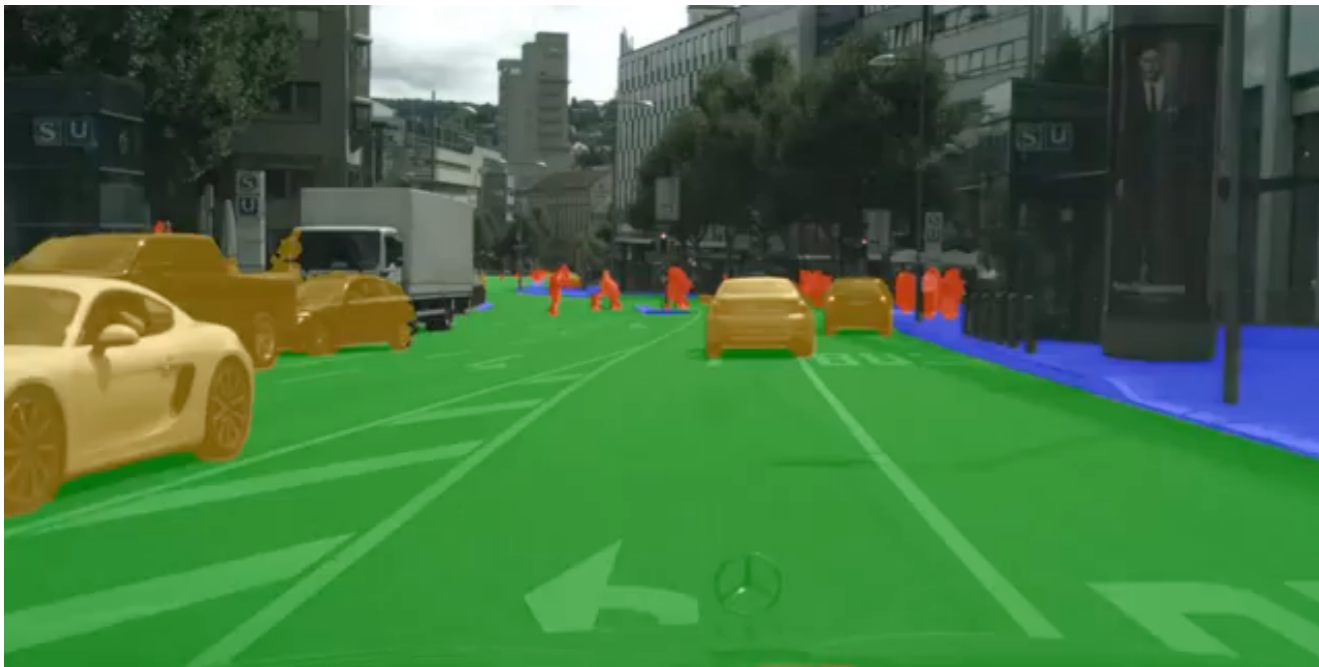


Flexible Embodiment



We do a lot more than passive viewing...

Scene Understanding

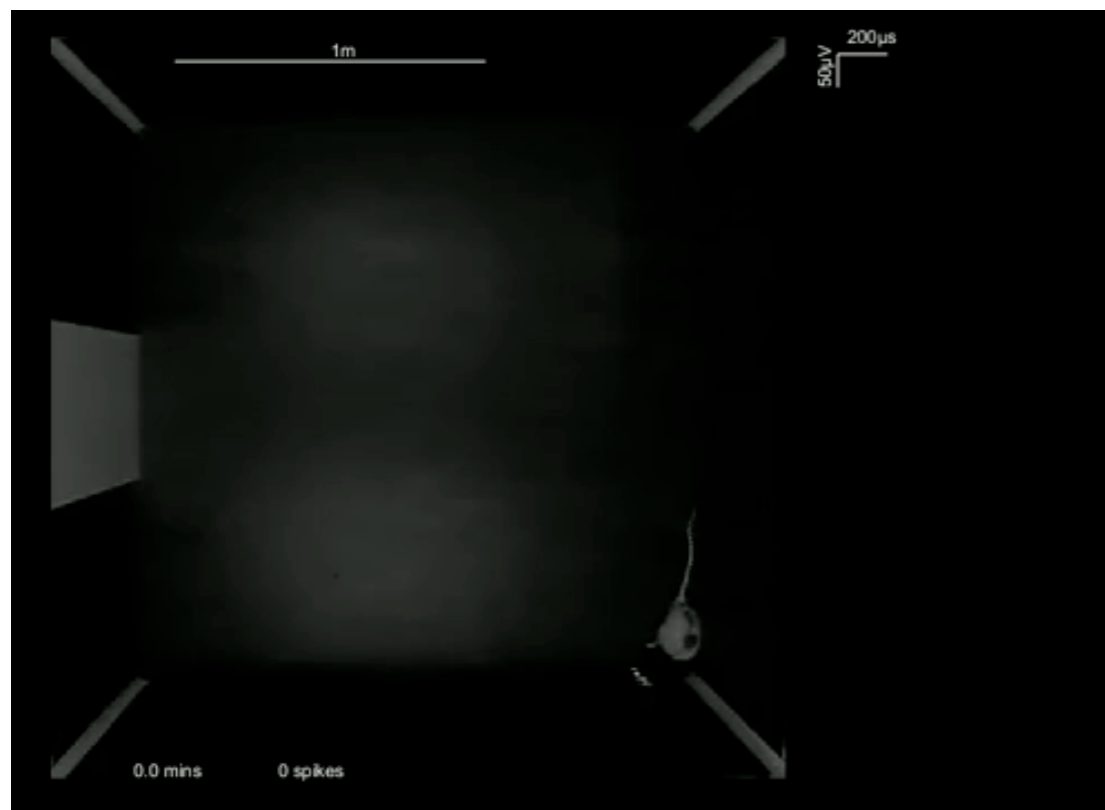


Multi-Step Planning

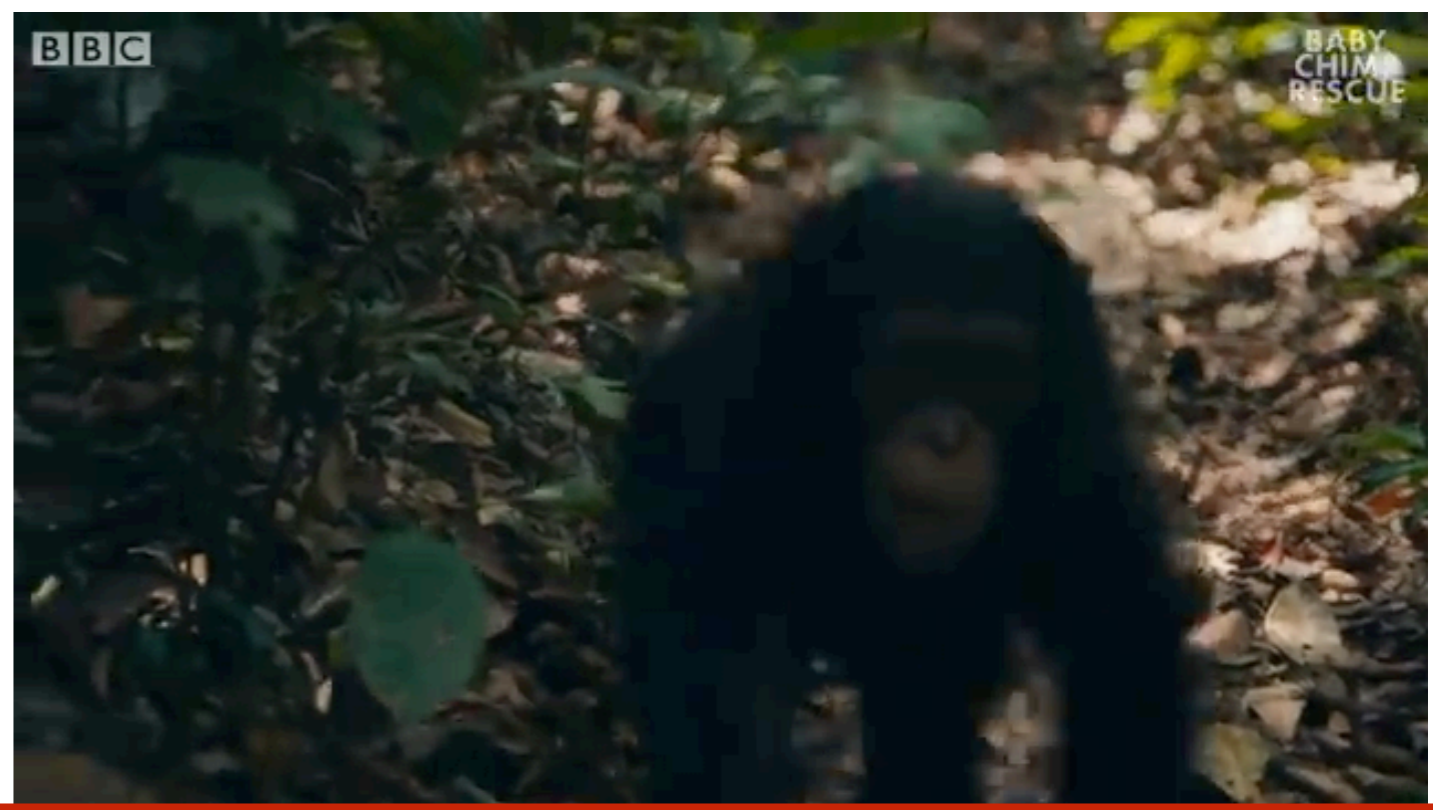


What are the core design principles that give rise to these abilities?

Navigation



Flexible Embodiment



Outline

- ▶ Mouse Visual Cortex as a Task-General, Limited Resource System
- ▶ Reusable Latent Representations for Primate Mental Simulation
- ▶ Heuristics for Interrogating Natural Intelligence

Outline

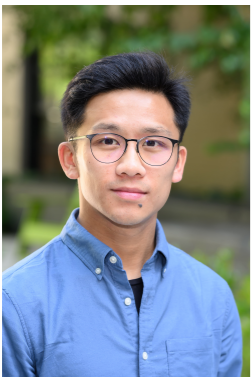
- ▶ Mouse Visual Cortex as a Task-General, Limited Resource System
- ▶ Reusable Latent Representations for Primate Mental Simulation
- ▶ Heuristics for Interrogating Natural Intelligence

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

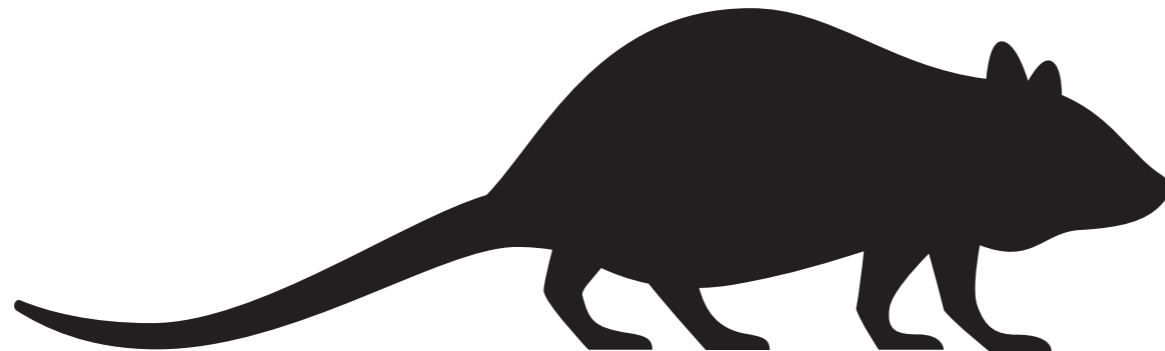
A. Nayebi*, N.C.L. Kong*, C. Zhuang, J.L. Gardner, A.M. Norcia, D.L.K. Yamins

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

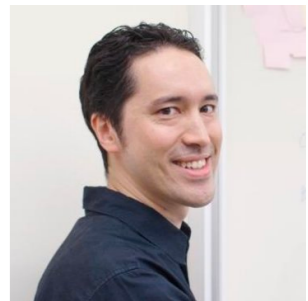
PLOS Computational Biology 2023



Nathan C.L. Kong*



Chengxu Zhuang



Justin L. Gardner



Anthony M. Norcia



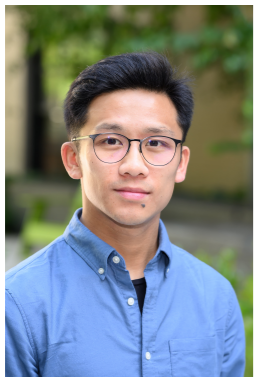
Daniel Yamins

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

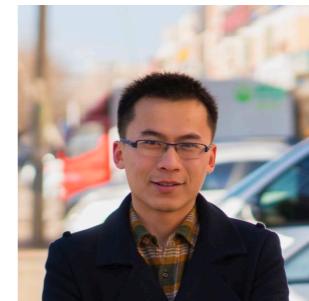
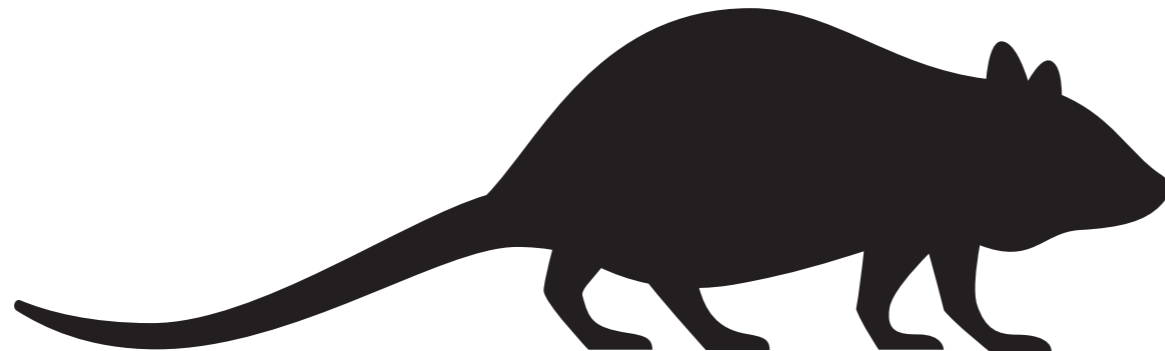
A. Nayebi*, N.C.L. Kong*, C. Zhuang, J.L. Gardner, A.M. Norcia, D.L.K. Yamins

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

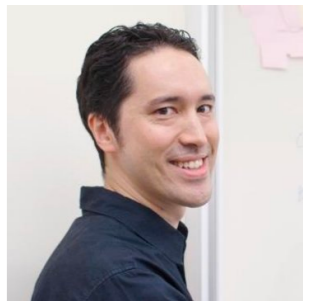
PLOS Computational Biology 2023



Nathan C.L. Kong*



Chengxu Zhuang



Justin L. Gardner

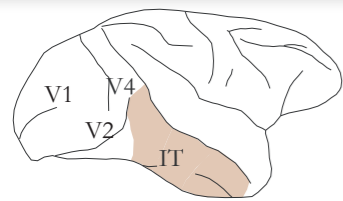


Anthony M. Norcia



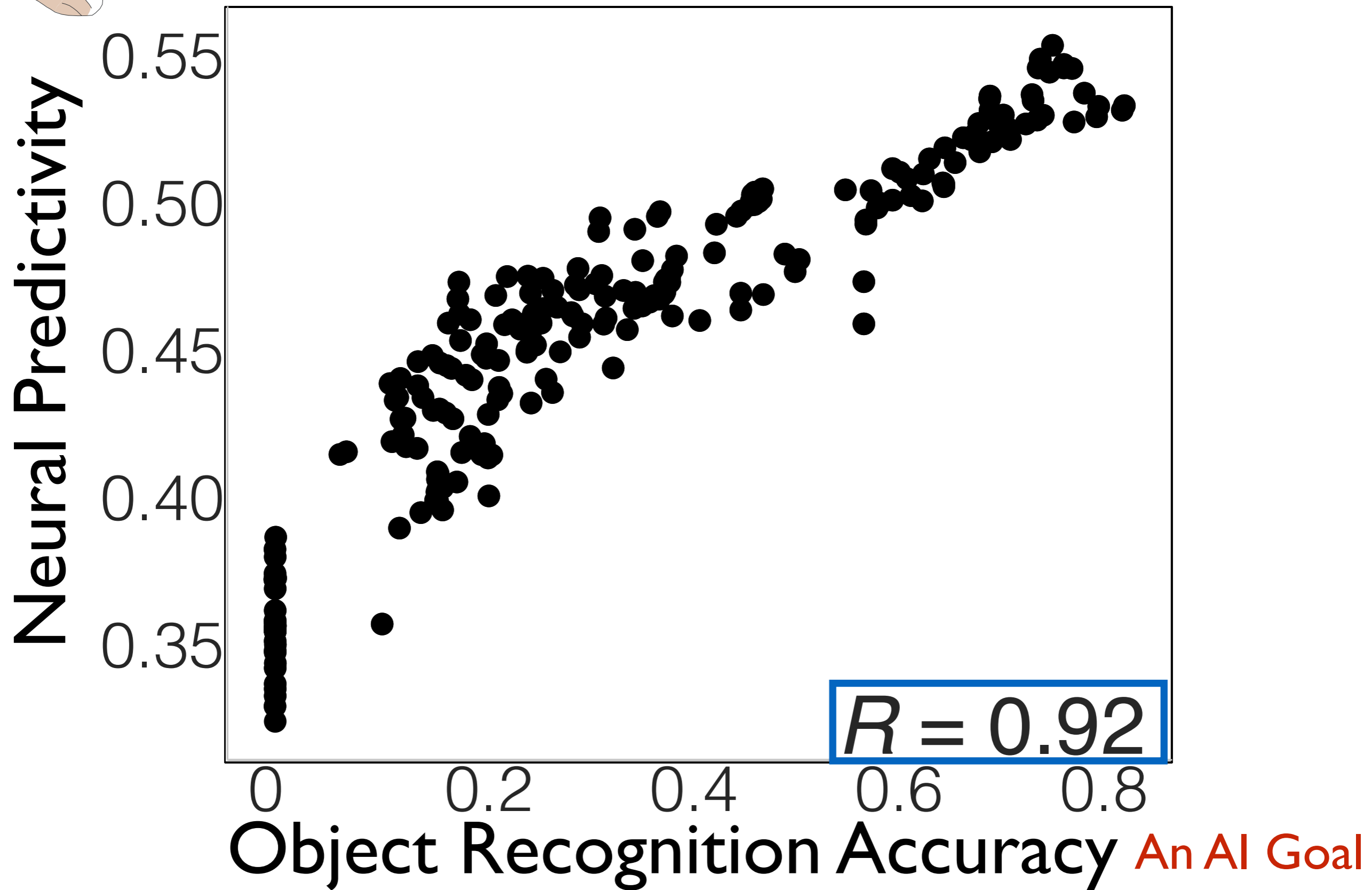
Daniel Yamins

Task Performance Correlated with Neural Predictivity

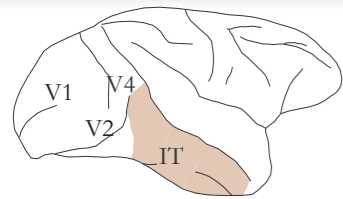


A Neuroscience Goal

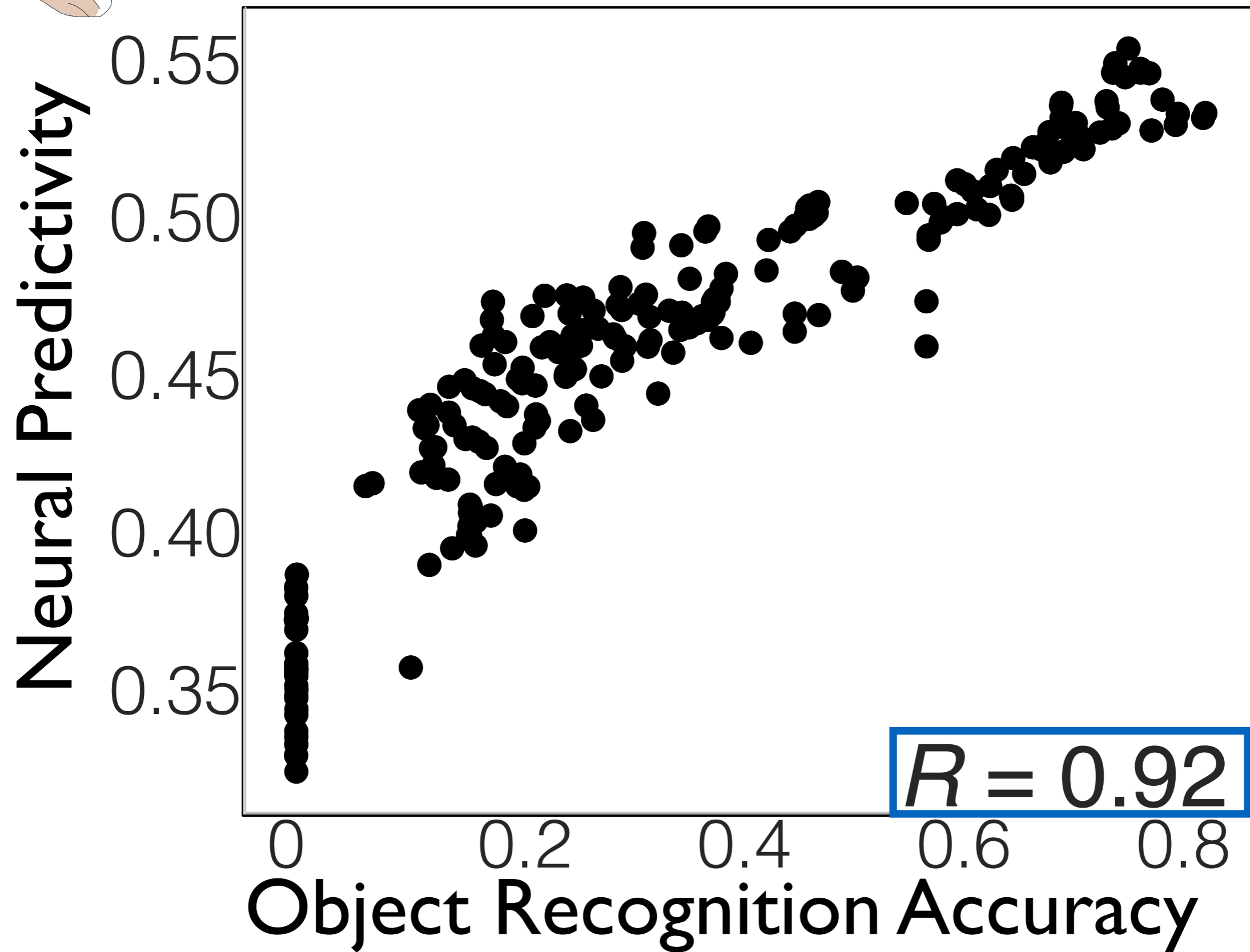
Schrimpf*, Kubilius* et al. 2018



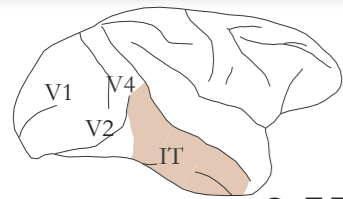
Task Performance Correlated with Neural Predictivity



Schrimpf*, Kubilius* et al. 2018

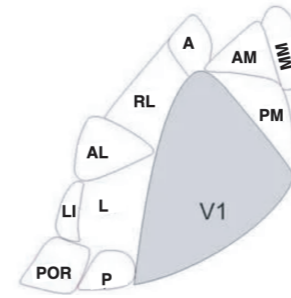
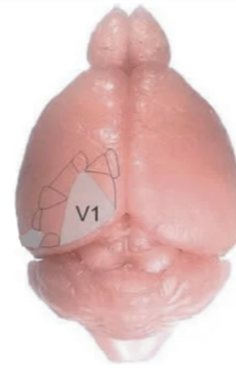
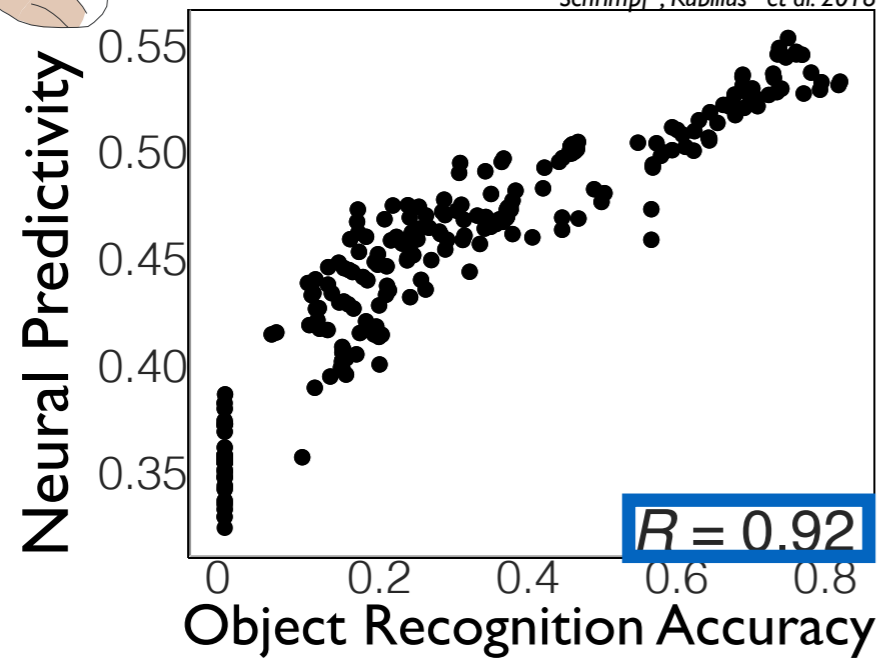


Task Performance Correlated with Neural Predictivity



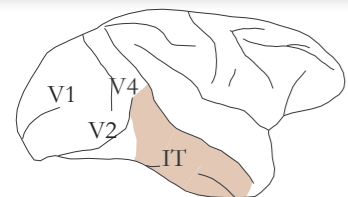
Primates

Schrimpf*, Kubilius* et al. 2018



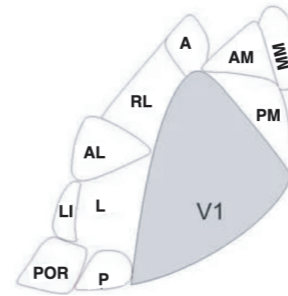
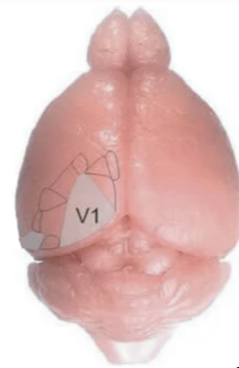
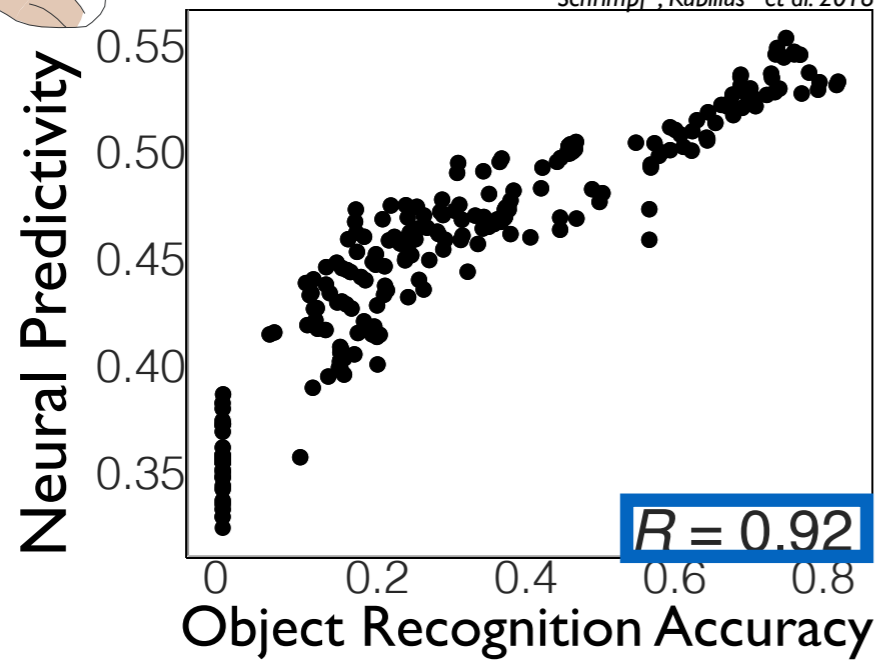
Mouse?

Object Categorization Ability **NOT** Correlated with Neural Predictivity

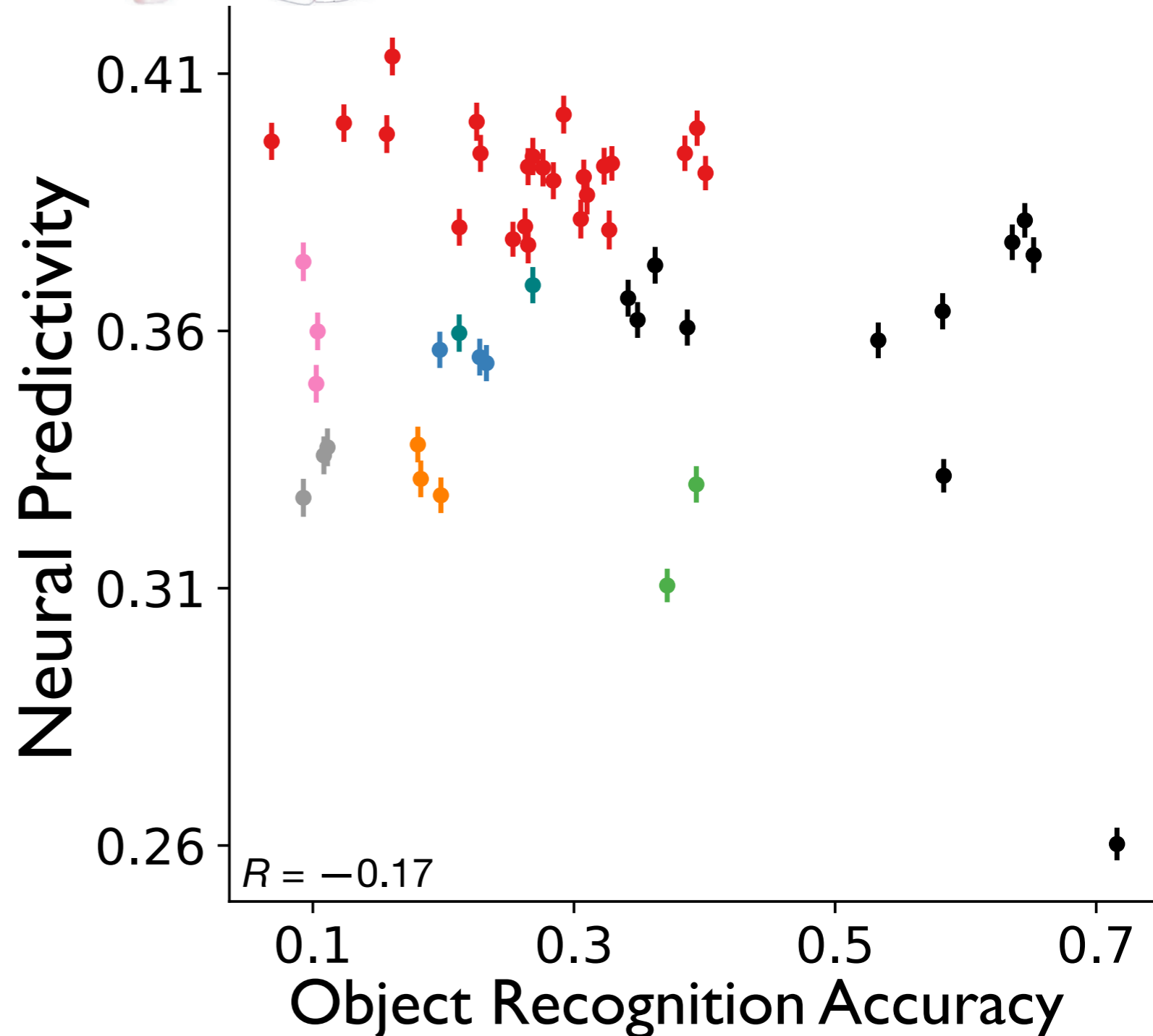


Primates

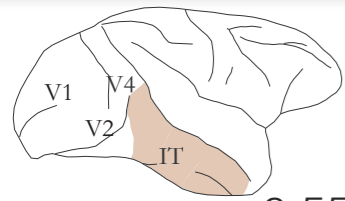
Schrimpf*, Kubilius* et al. 2018



Mouse

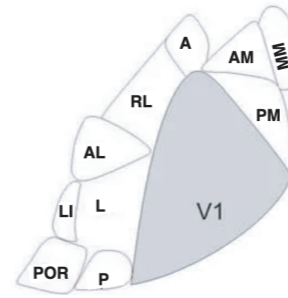
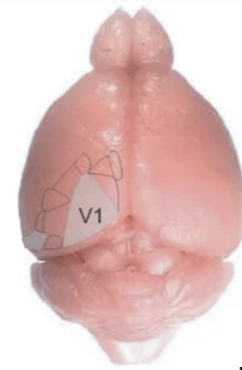
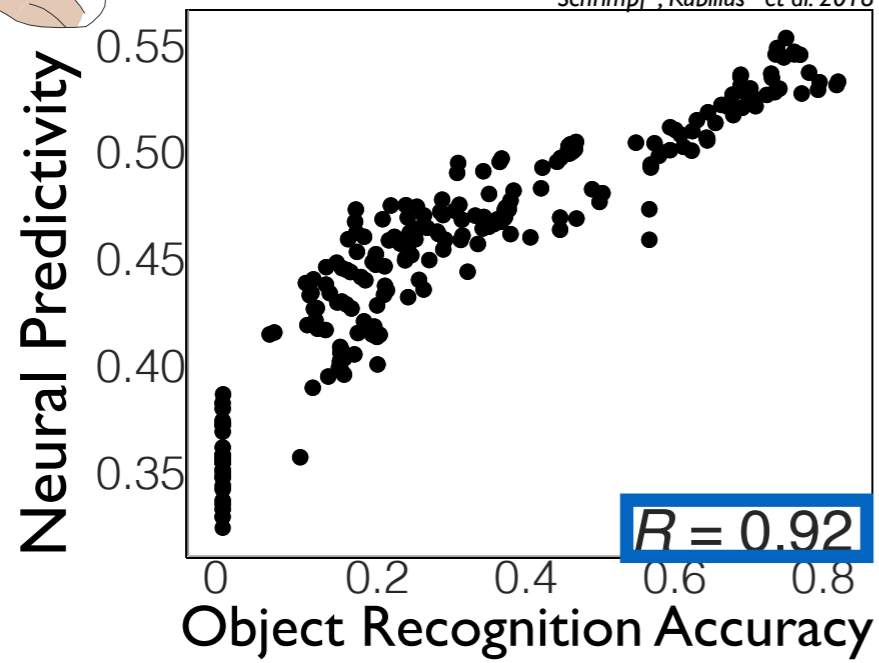


Object Categorization Ability **NOT** Correlated with Neural Predictivity

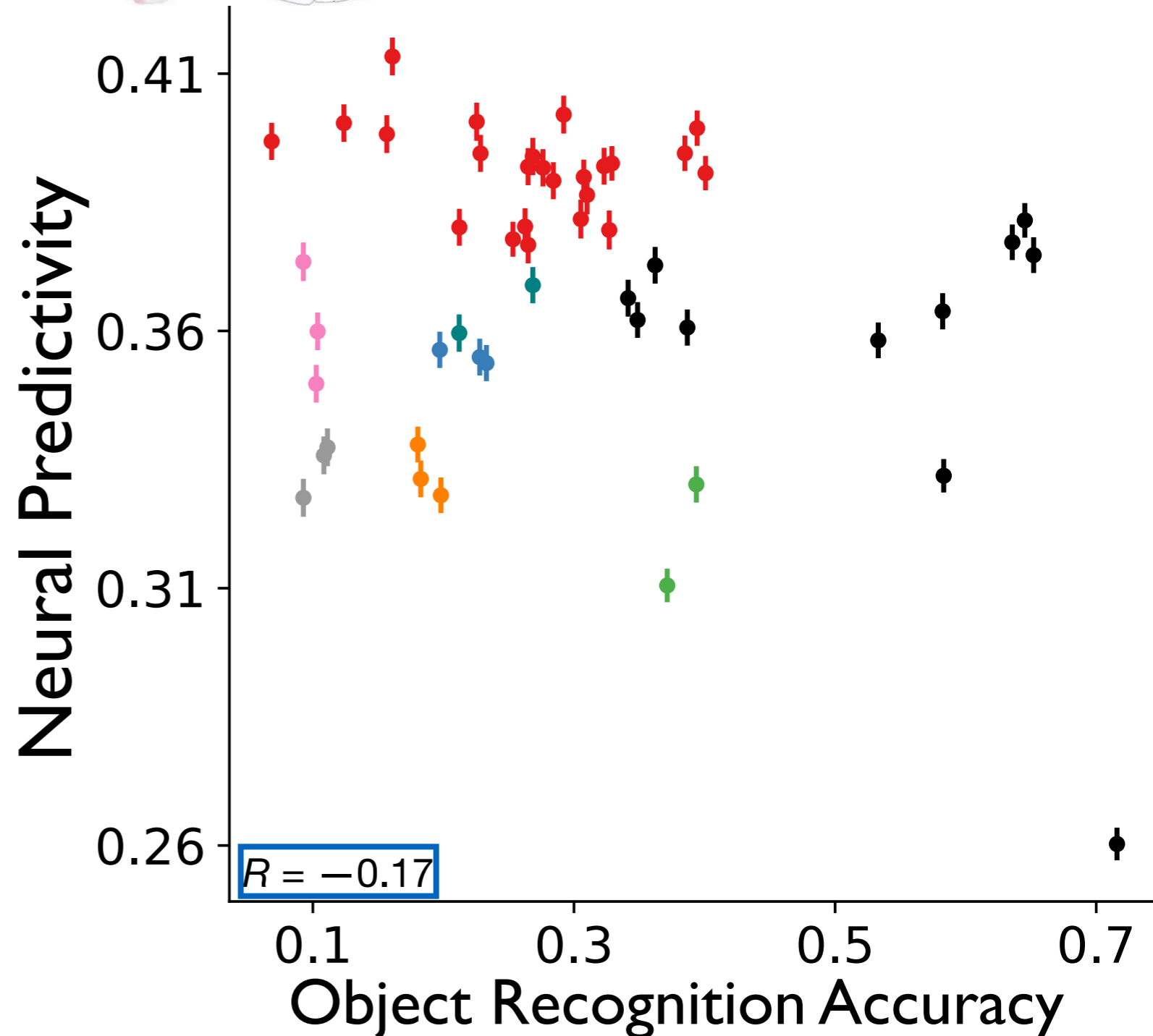


Primates

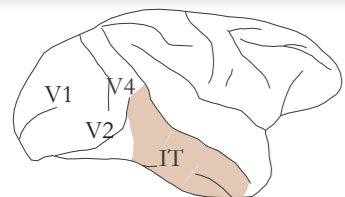
Schrimpf*, Kubilius* et al. 2018



Mouse

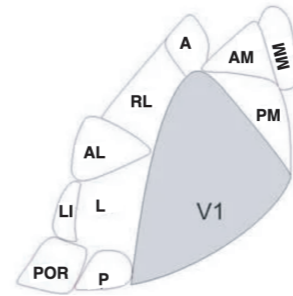
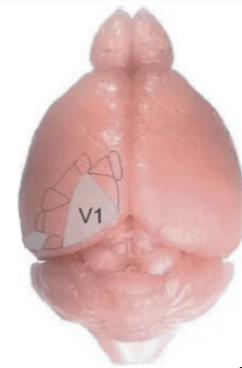
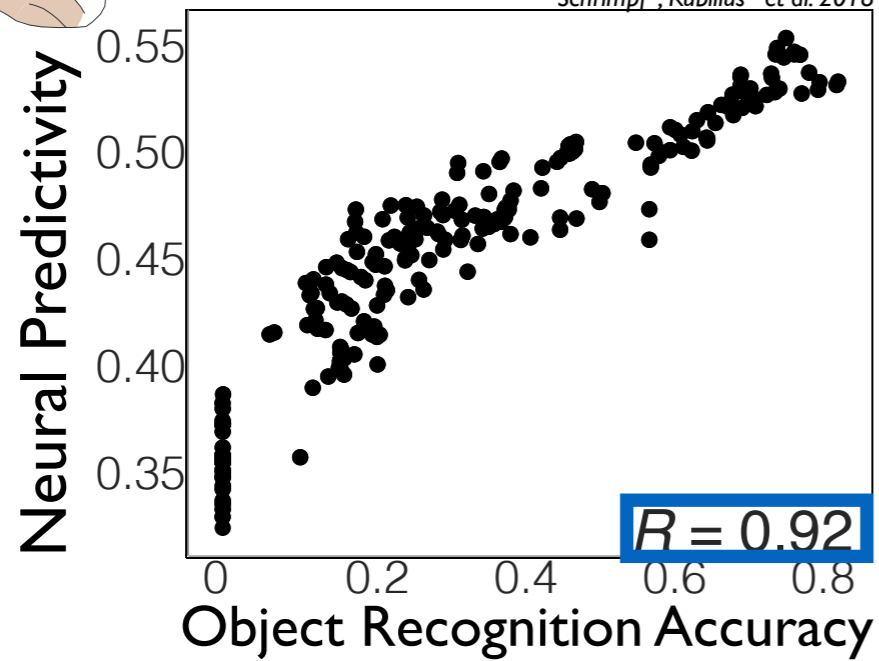


Object Categorization Ability **NOT** Correlated with Neural Predictivity

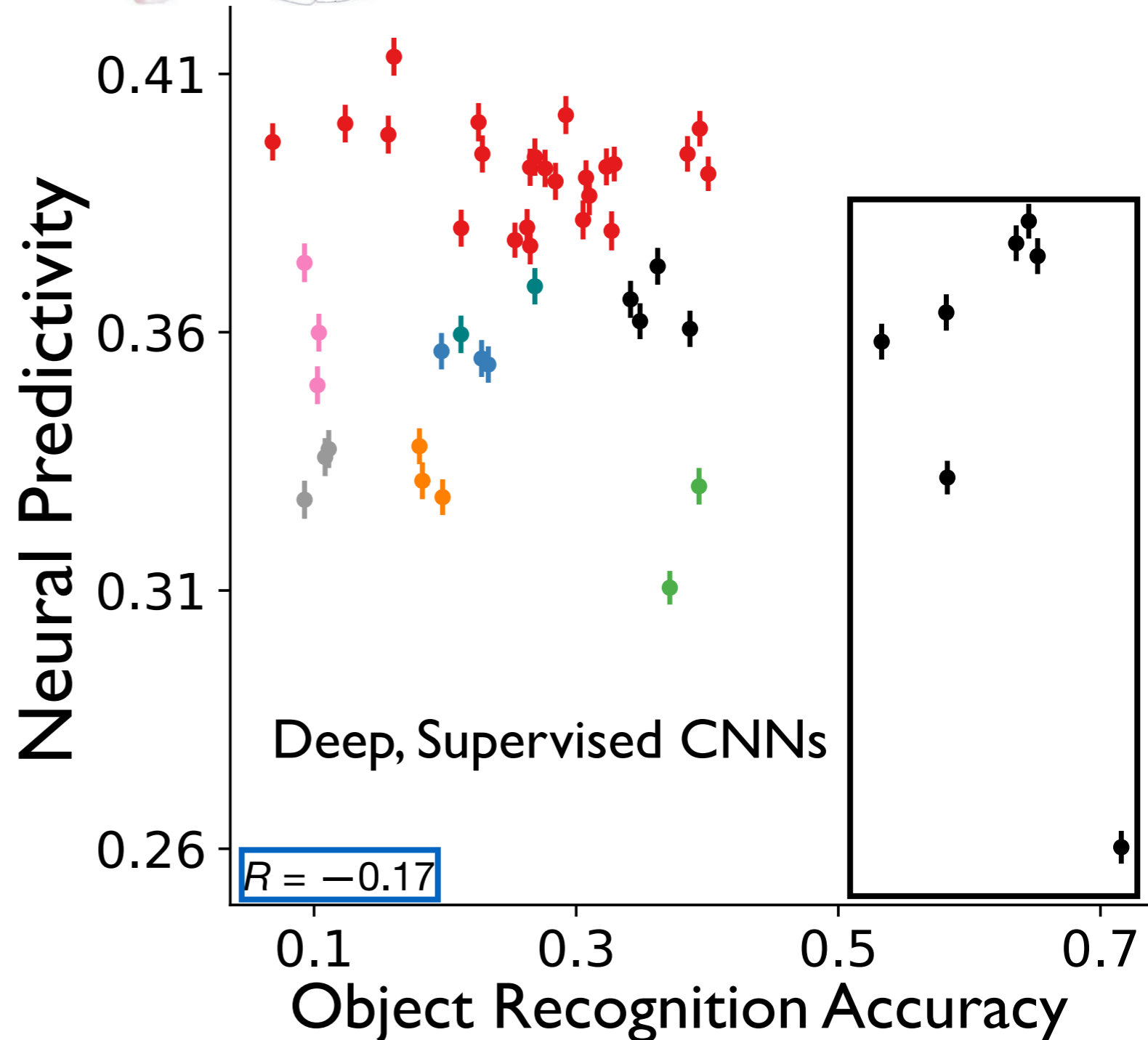


Primates

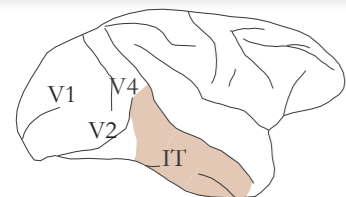
Schrimpf*, Kubilius* et al. 2018



Mouse

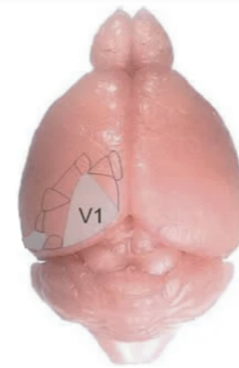
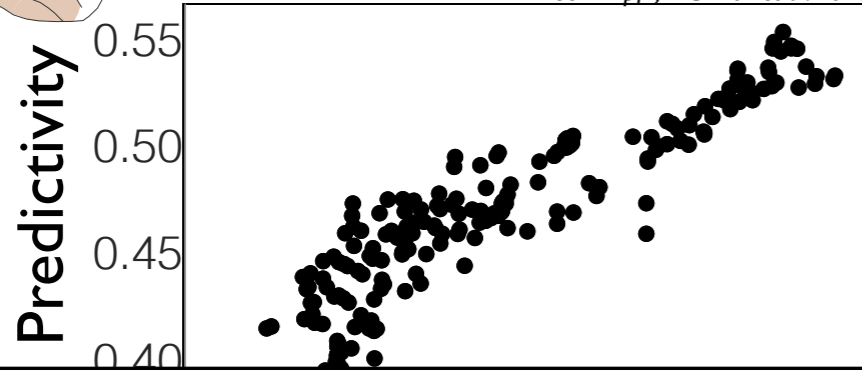


Object Categorization Ability **NOT** Correlated with Neural Predictivity

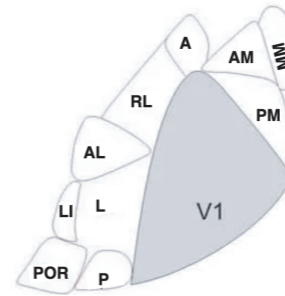


Primates

Schrimpf*, Kubilius* et al. 2018



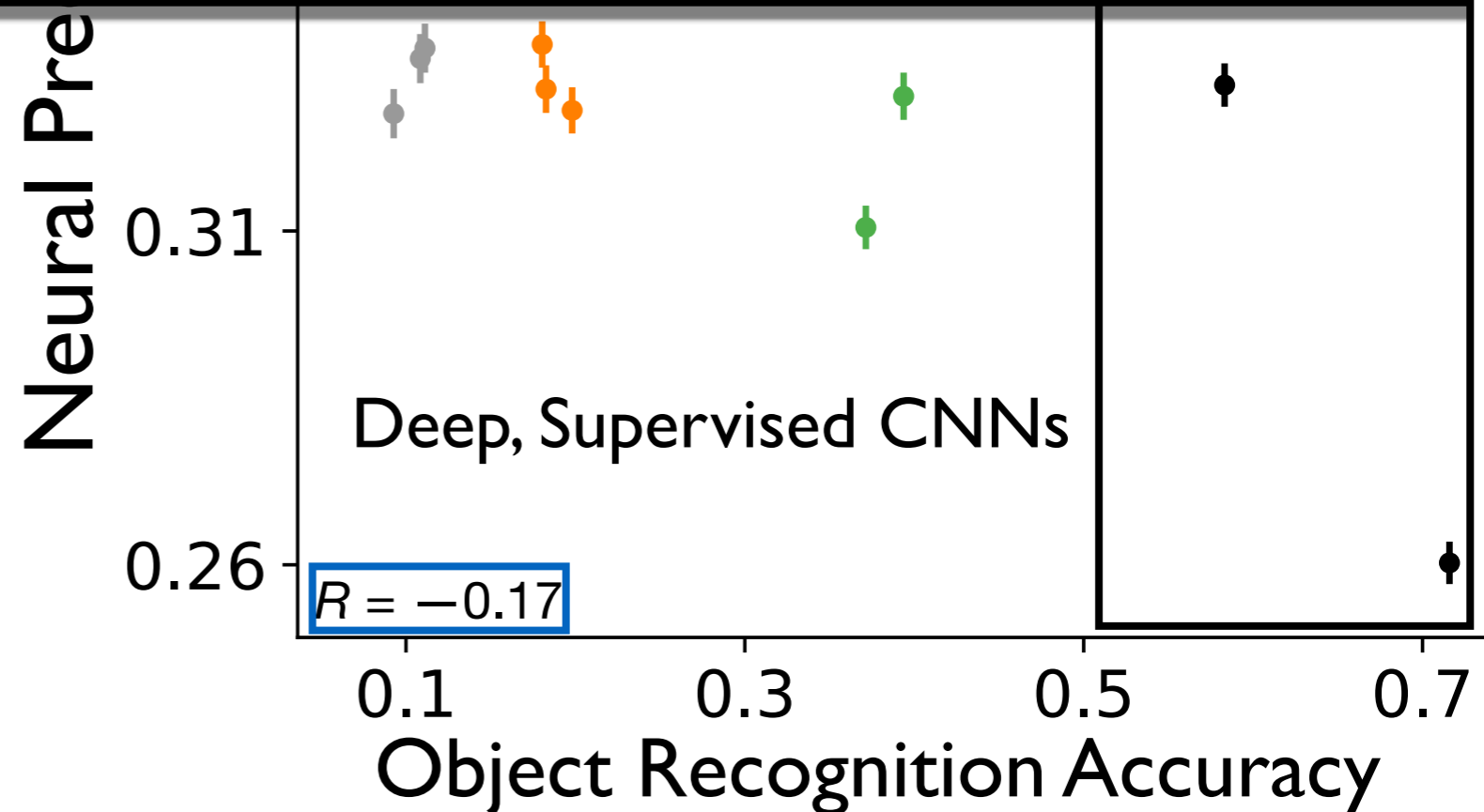
0.41



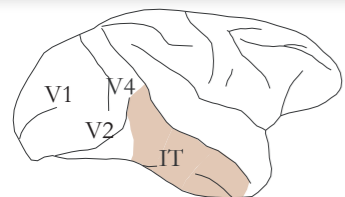
Mouse



Neurobiological Puzzle:
Does task-optimization apply to rodents?

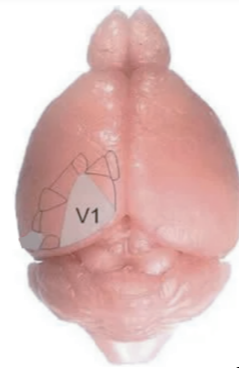
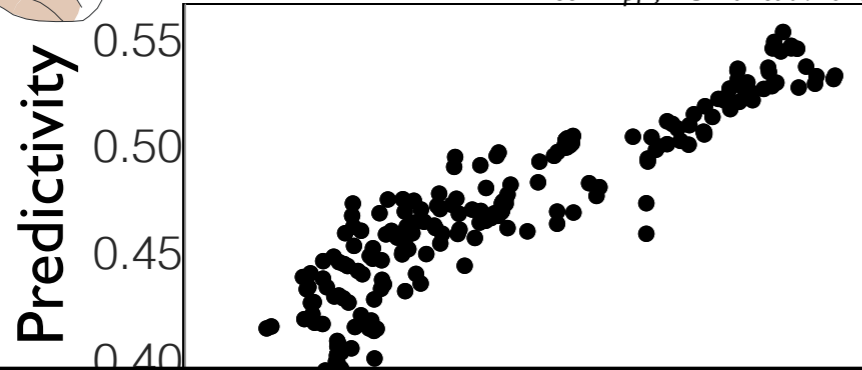


Object Categorization Ability **NOT** Correlated with Neural Predictivity

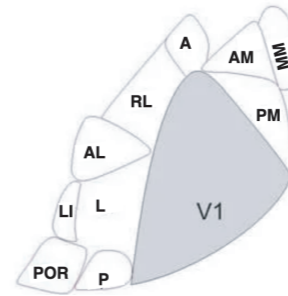


Primates

Schrimpf*, Kubilius* et al. 2018



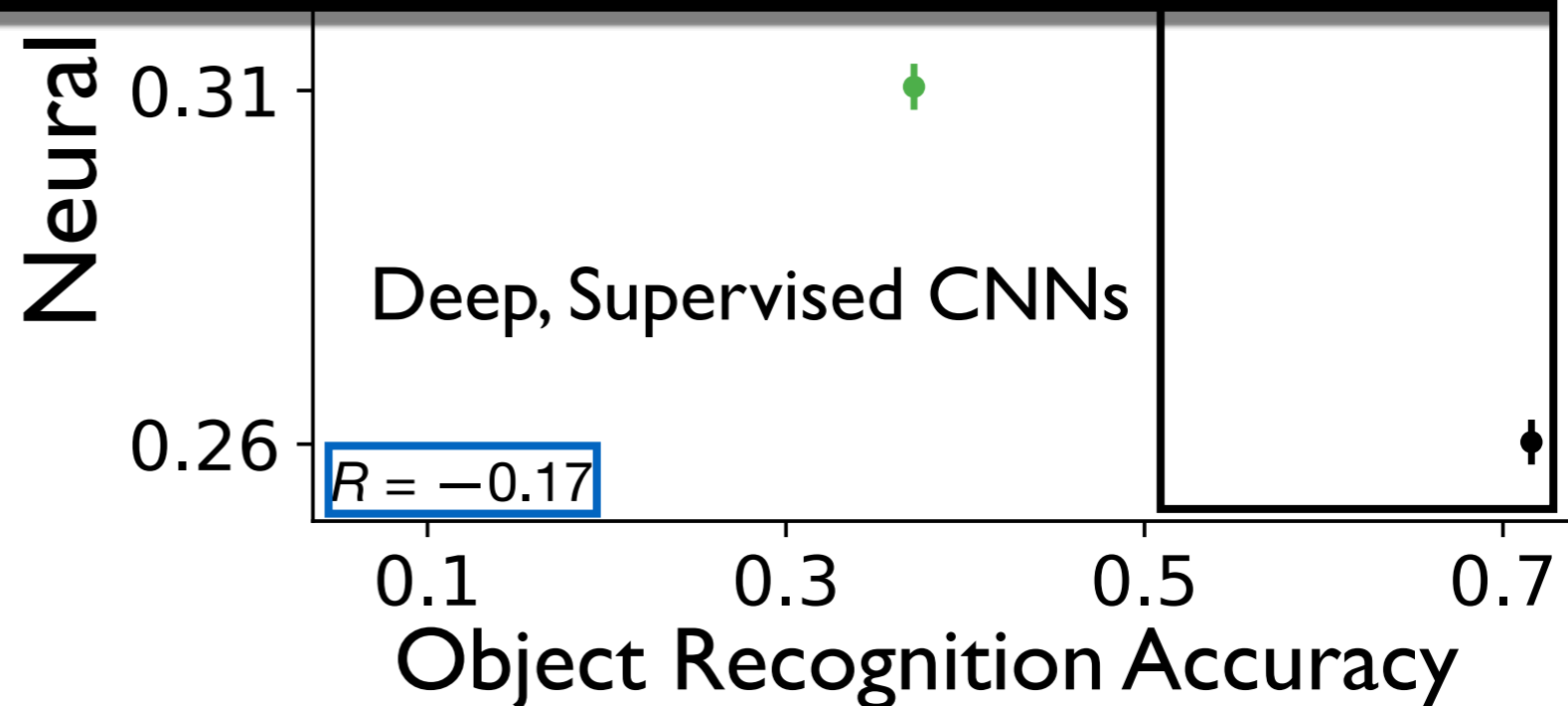
0.41



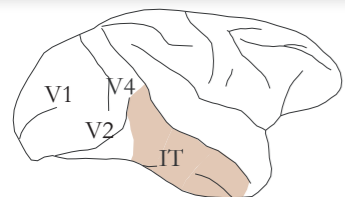
Mouse

Neurobiological Puzzle:
Does task-optimization apply to rodents?

Yes!

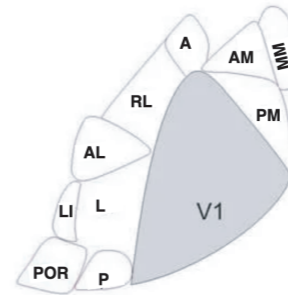
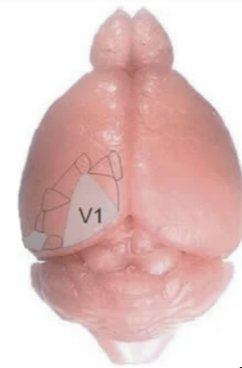
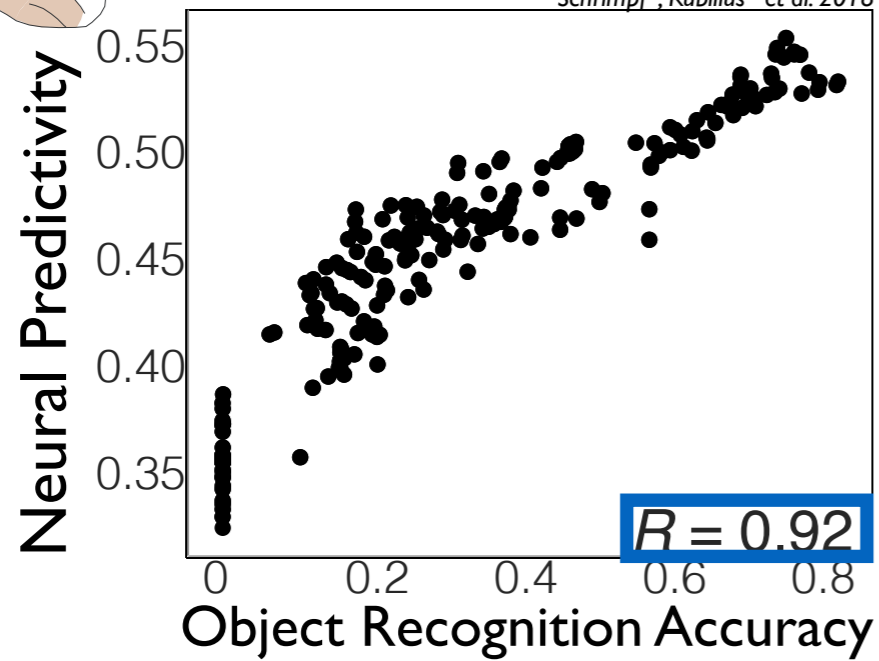


Object Categorization Ability **NOT** Correlated with Neural Predictivity

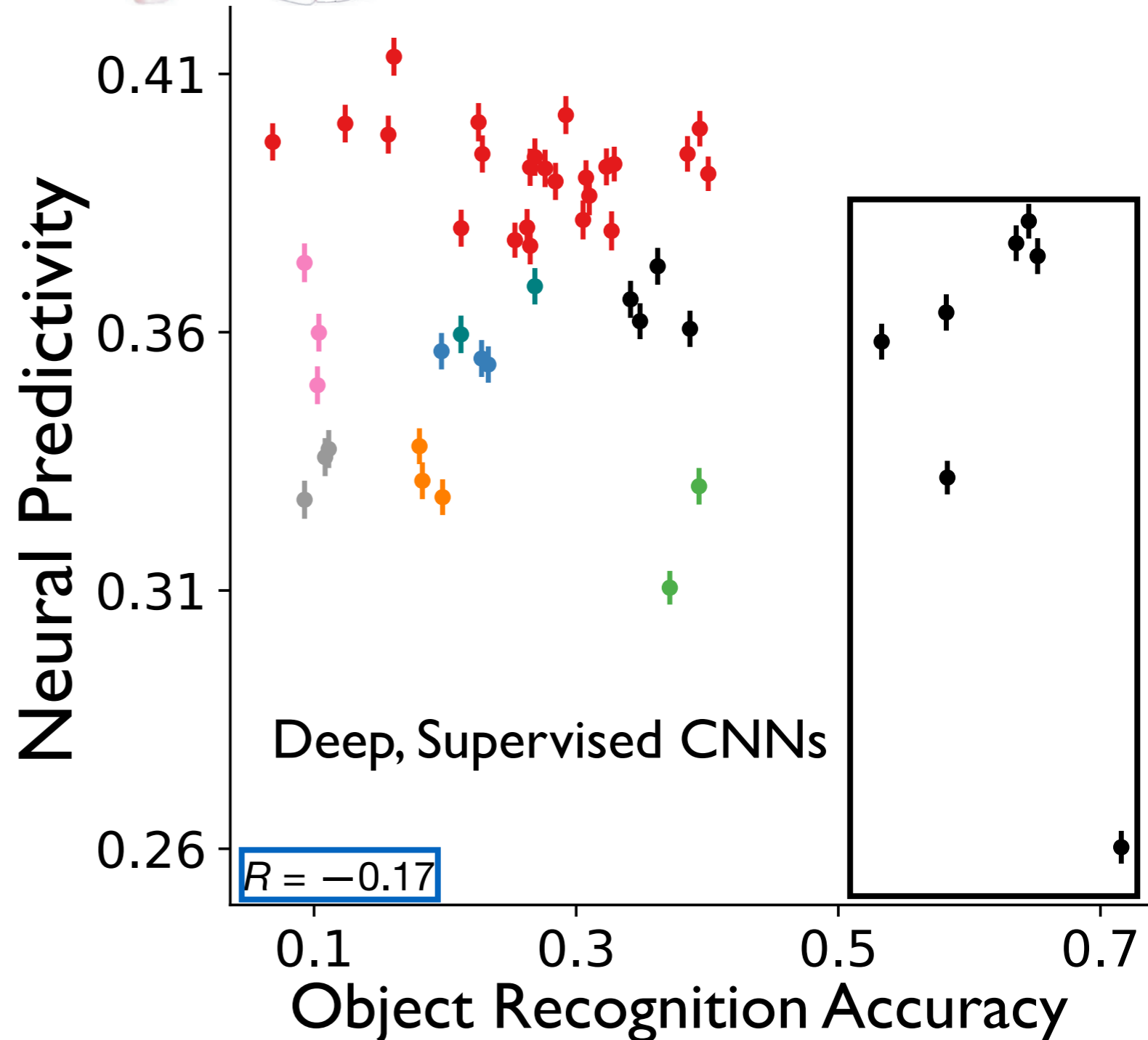


Primates

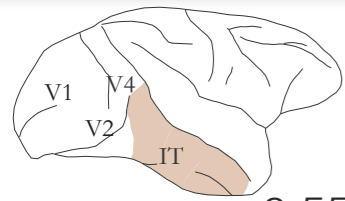
Schrimpf*, Kubilius* et al. 2018



Mouse

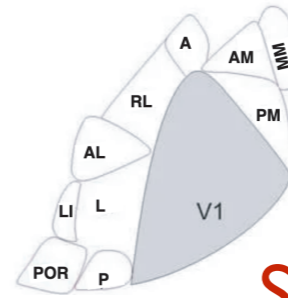
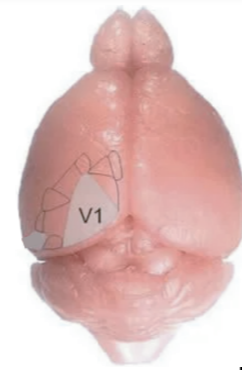
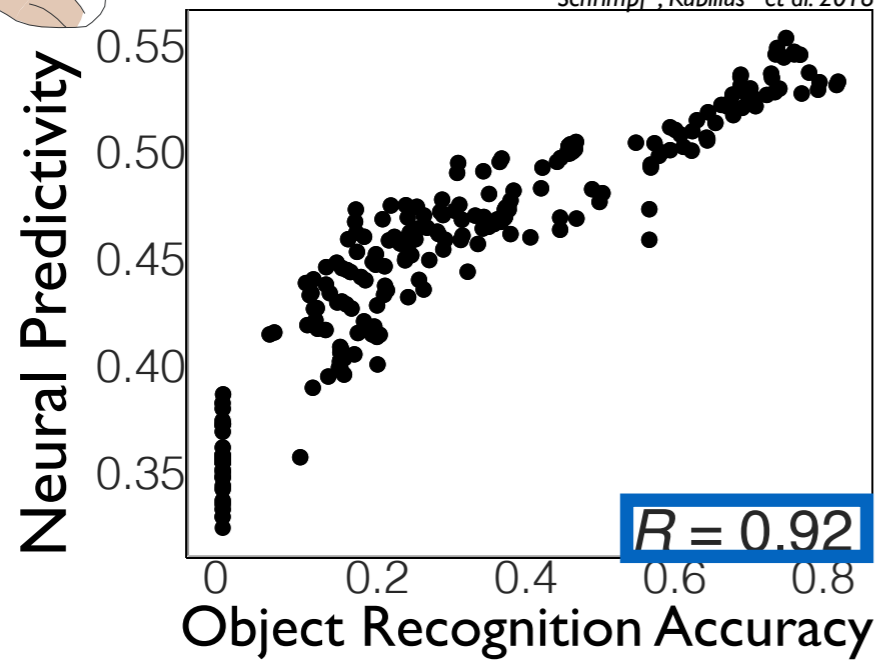


Object Categorization Ability **NOT** Correlated with Neural Predictivity



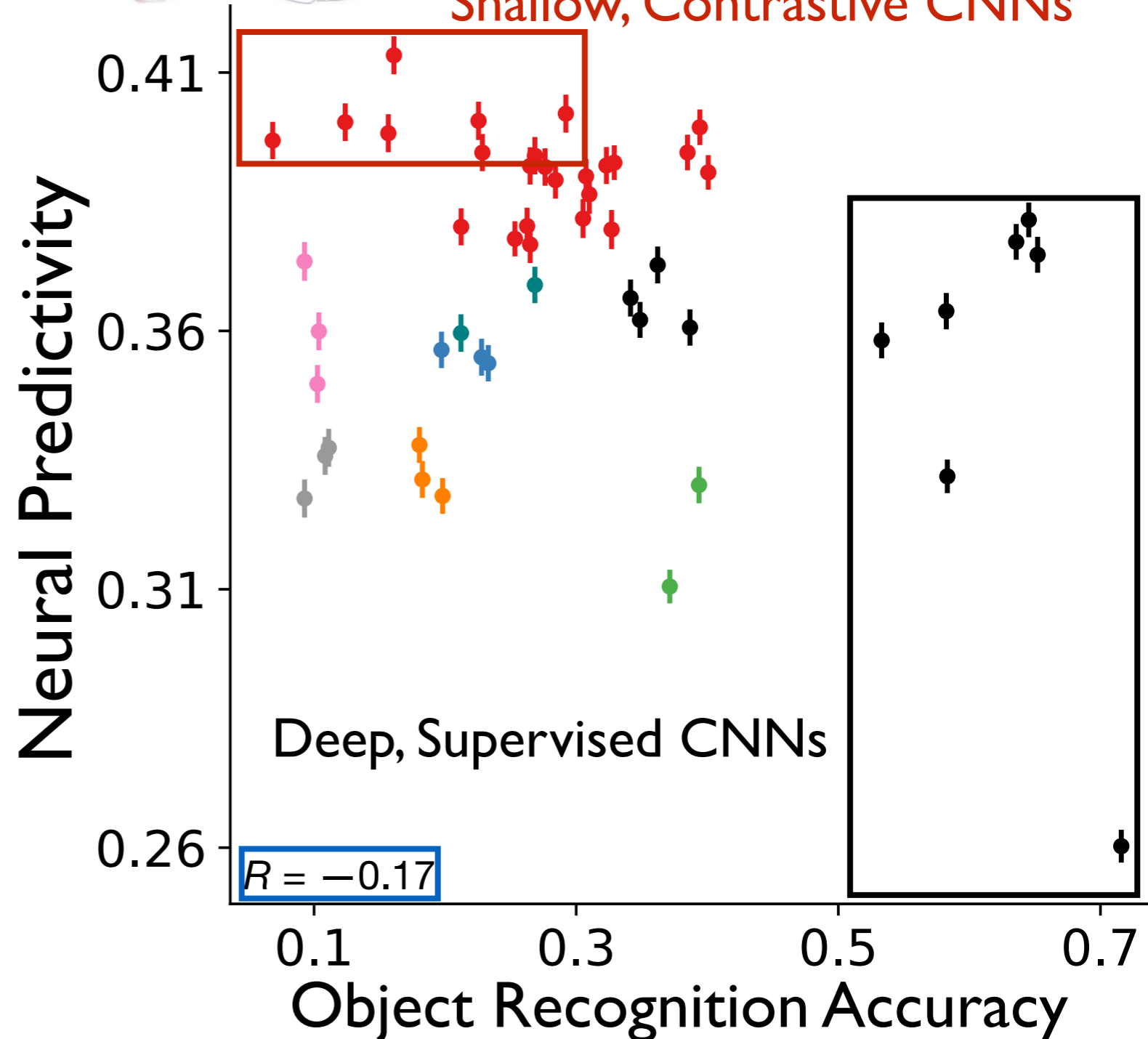
Primates

Schrimpf*, Kubilius* et al. 2018

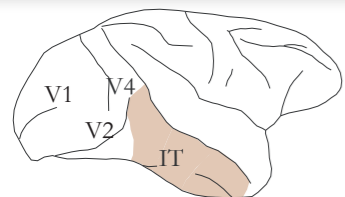


Mouse

Shallow, Contrastive CNNs

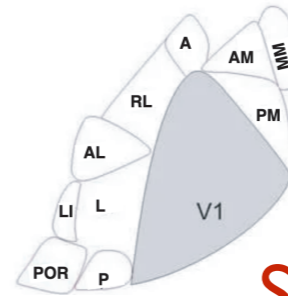
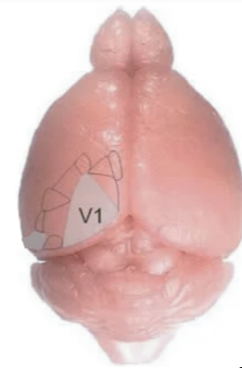
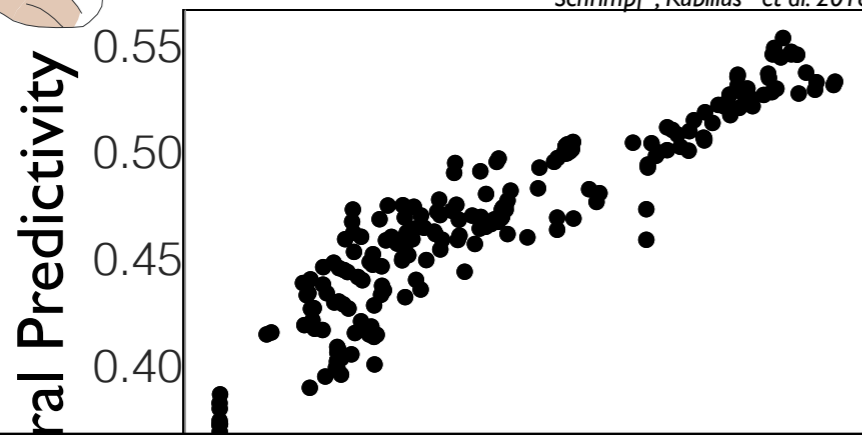


Object Categorization Ability **NOT** Correlated with Neural Predictivity



Primates

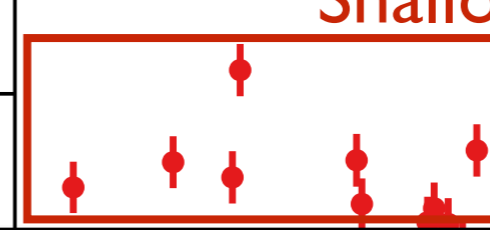
Schrimpf*, Kubilius* et al. 2018



Mouse

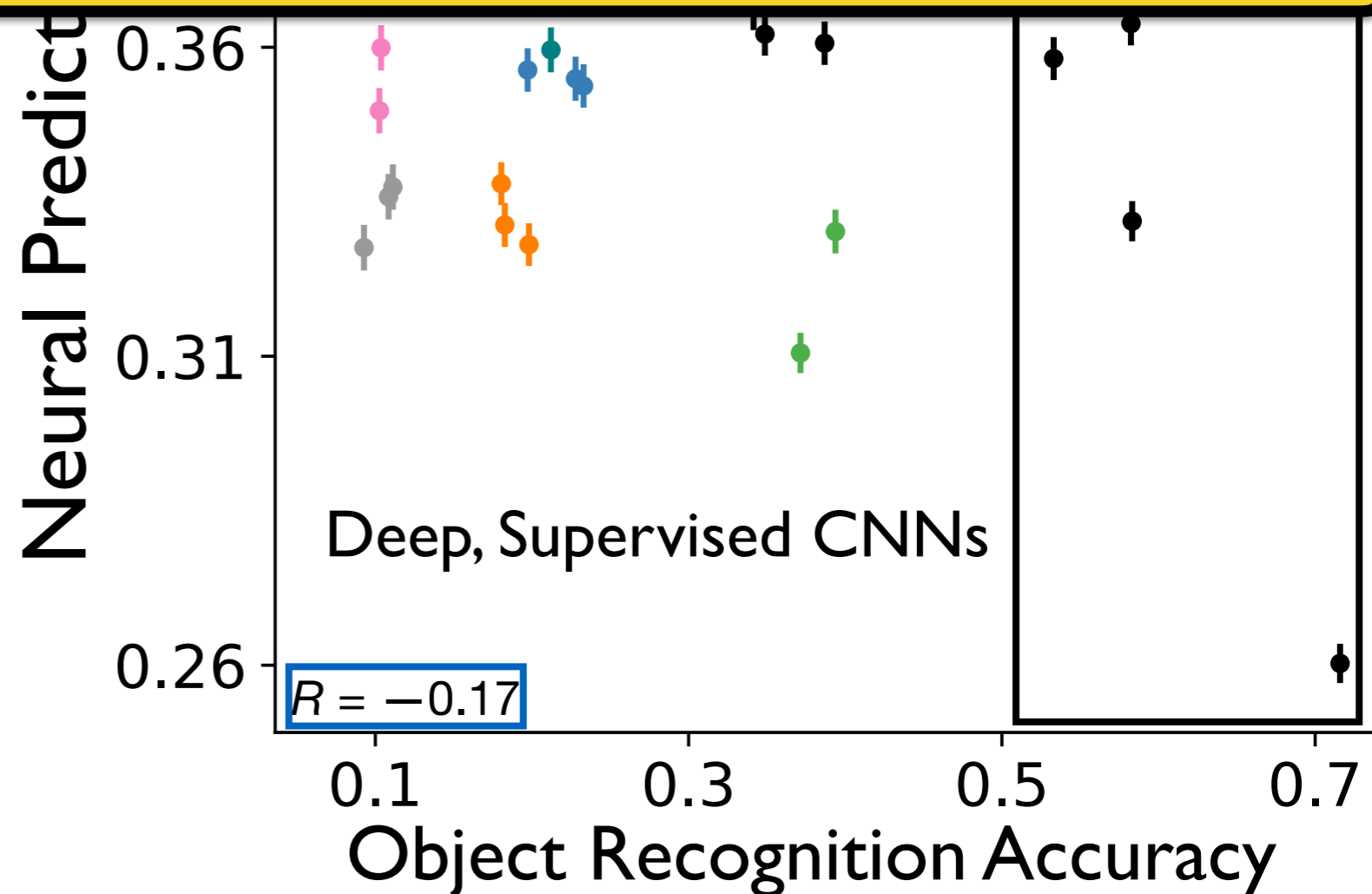
Shallow, Contrastive CNNs

0.41

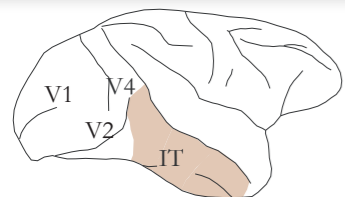


What is the ecological reason *why* the mouse visual system prefers self-supervision?

Object Recognition Accuracy

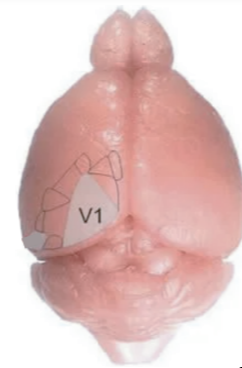
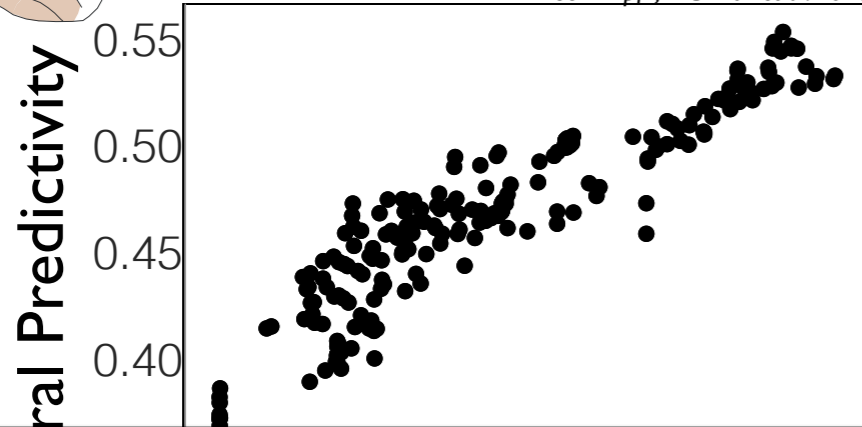


Object Categorization Ability **NOT** Correlated with Neural Predictivity

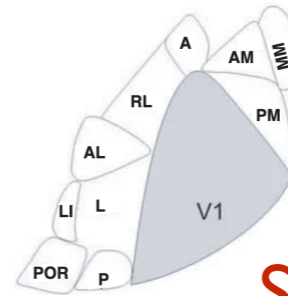


Primates

Schrimpf*, Kubilius* et al. 2018

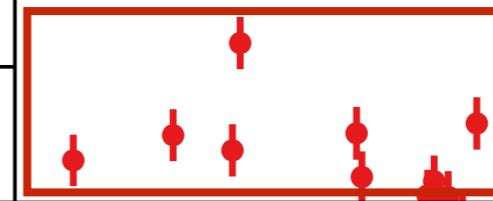


Mouse



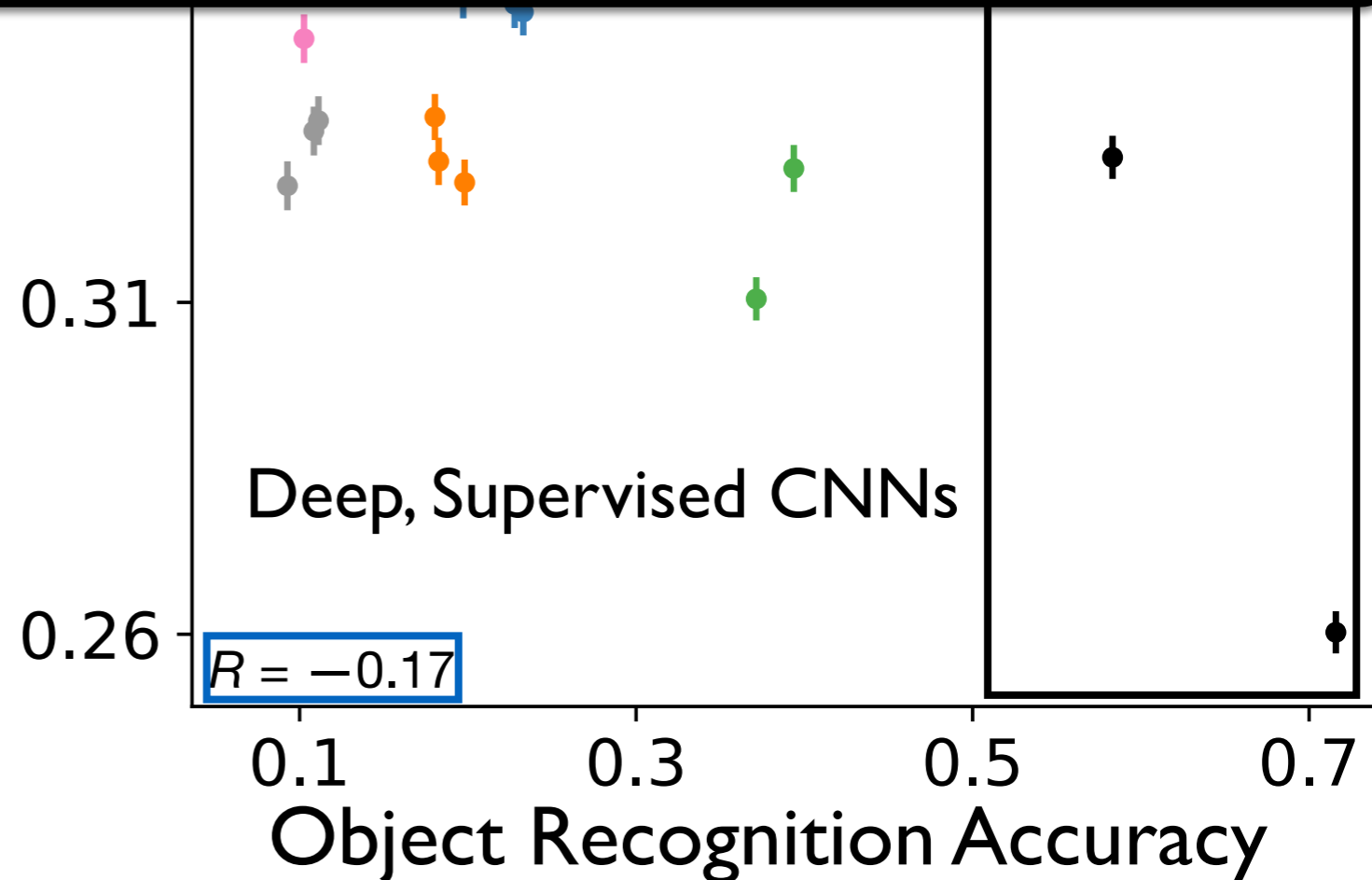
Shallow, Contrastive CNNs

0.41



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

Neural Predictivity



Deep, Supervised CNNs

$R = -0.17$

Object Recognition Accuracy

Assessing Task-Generality

Assessing Task-Generality

Train

ImageNet



Assessing Task-Generality

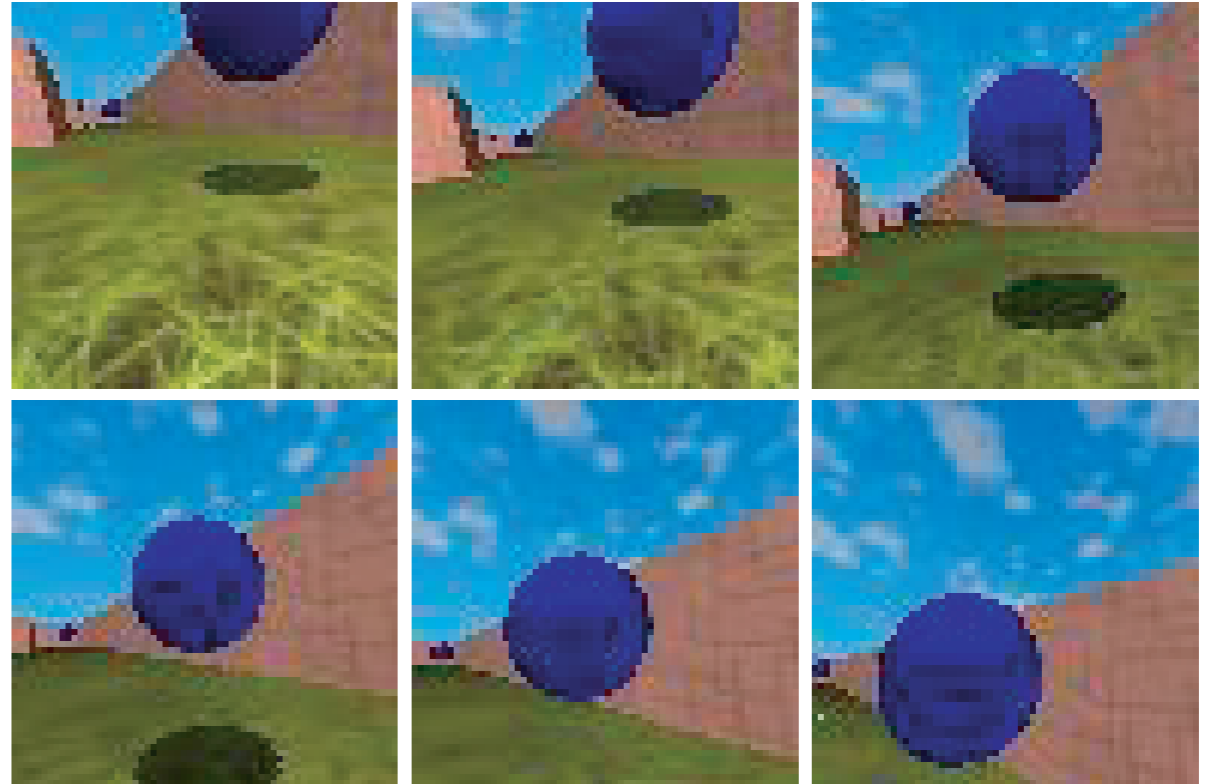
Train

ImageNet



Evaluate

Reward-Based Navigation



Assessing Task-Generality

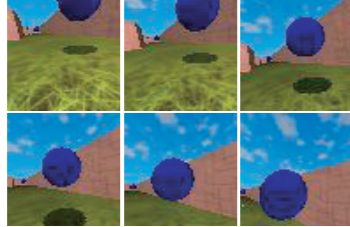
Train

ImageNet



Evaluate

Reward-Based Navigation



Assessing Task-Generality

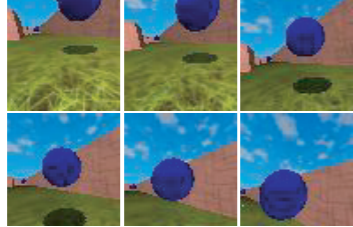
Train

ImageNet



Evaluate

Reward-Based Navigation



Embodied Virtual Rodent Navigation

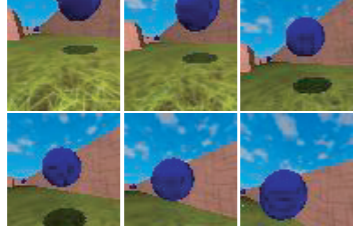
Train

ImageNet



Evaluate

Reward-Based Navigation



Embodied Virtual Rodent Navigation

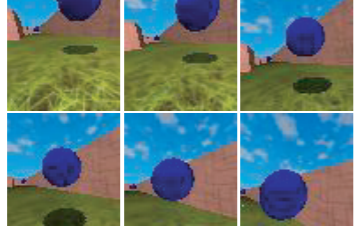
Train

ImageNet

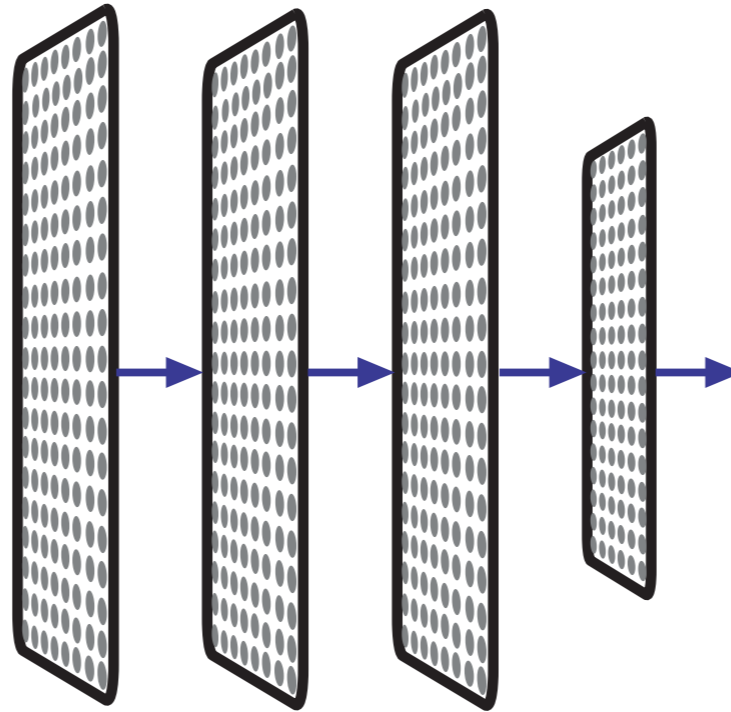
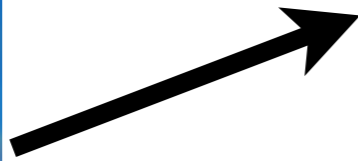


Evaluate

Reward-Based Navigation



Vision Network



Embodied Virtual Rodent Navigation

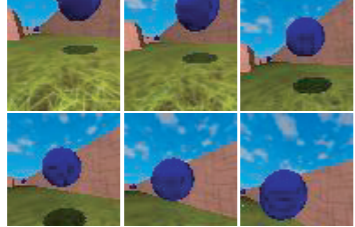
Train

ImageNet

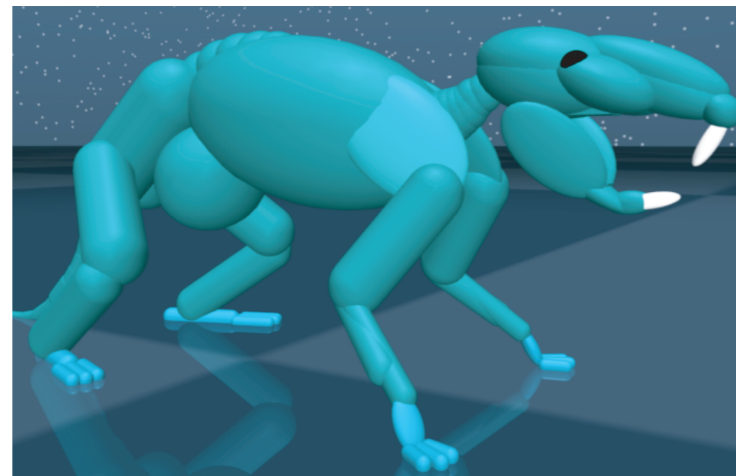
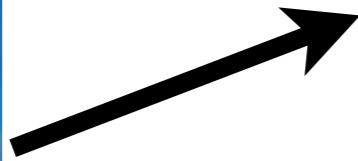
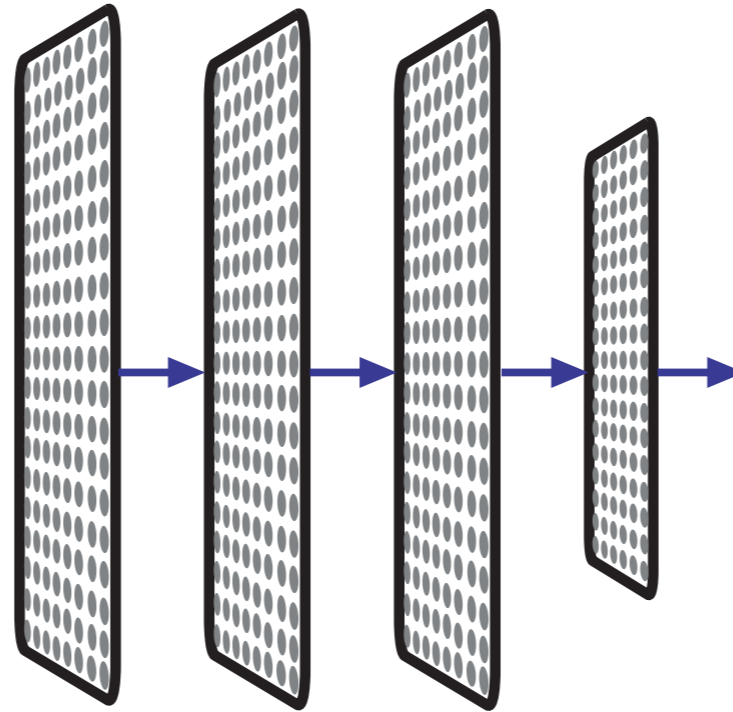


Evaluate

Reward-Based Navigation



Vision Network



Biomechanical Model

(Joint angles, accelerometer, etc.)



Bence Ölveczky

Embodied Virtual Rodent Navigation

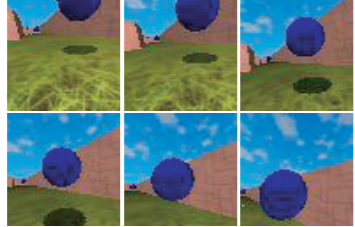
Train

ImageNet

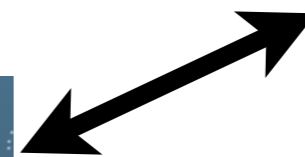
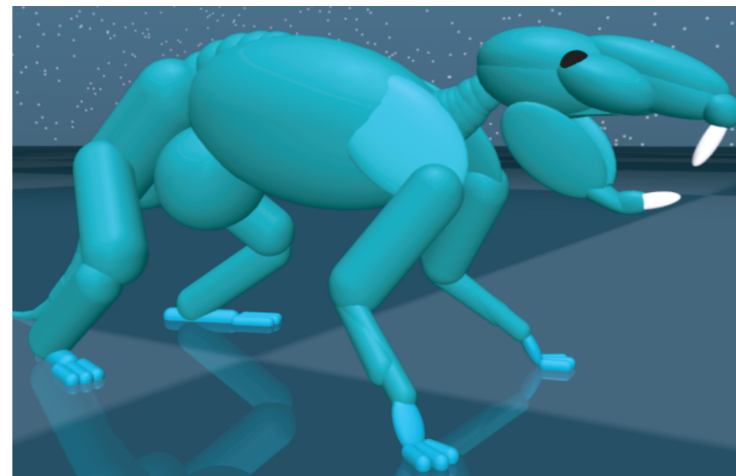
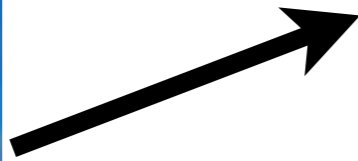
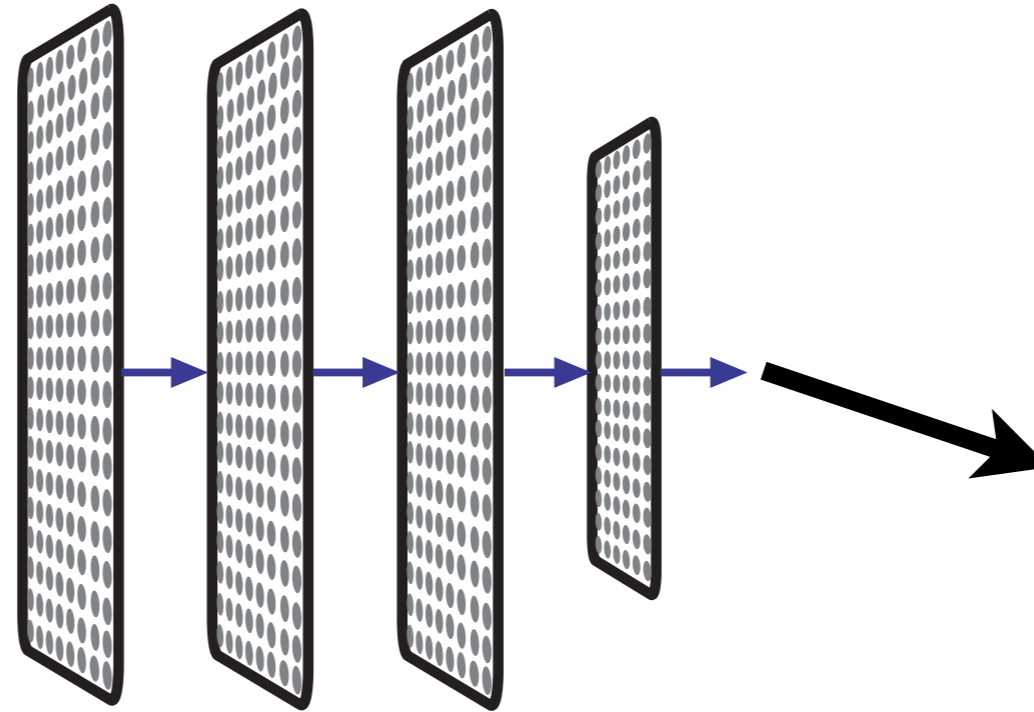


Evaluate

Reward-Based Navigation



Vision Network

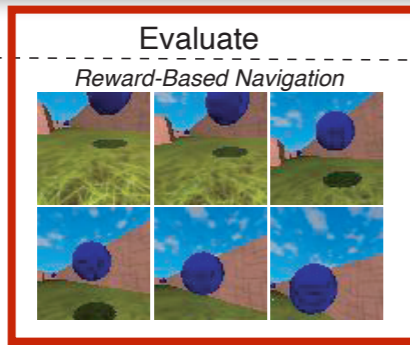
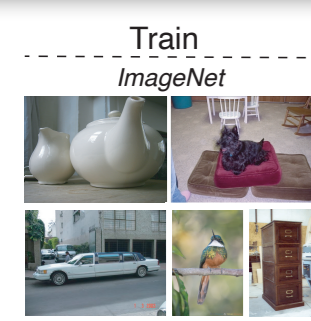


Bence Ölveczky

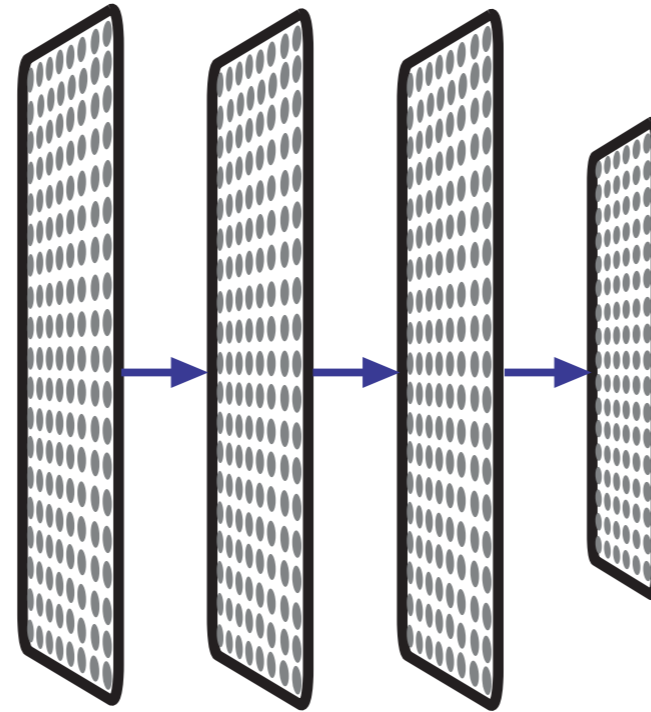
Biomechanical Model

(Joint angles, accelerometer, etc.)

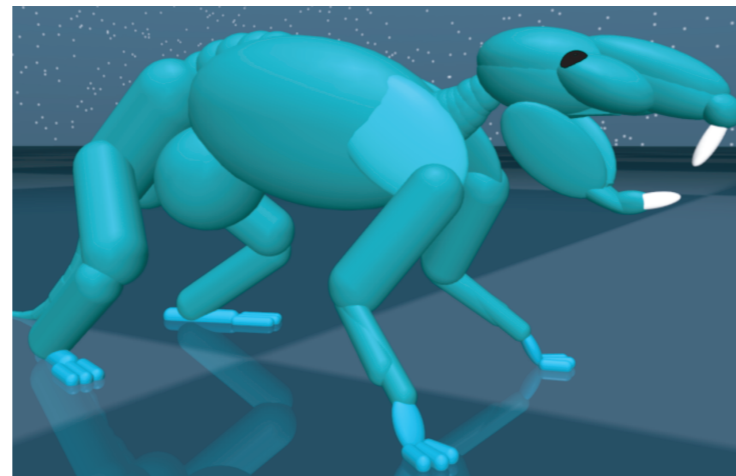
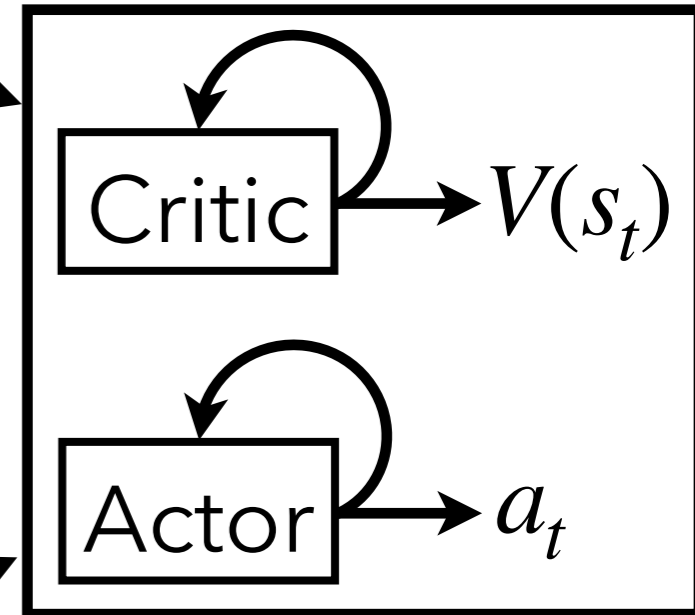
Embodied Virtual Rodent Navigation



Vision Network



Decision Making



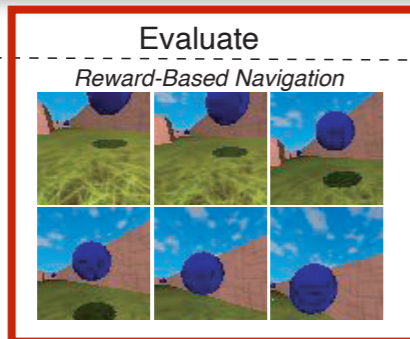
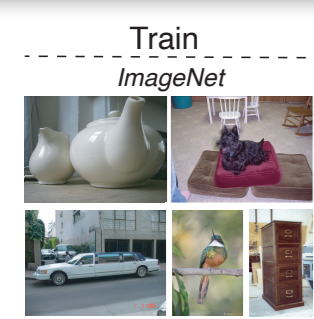
Biomechanical Model

(Joint angles, accelerometer, etc.)

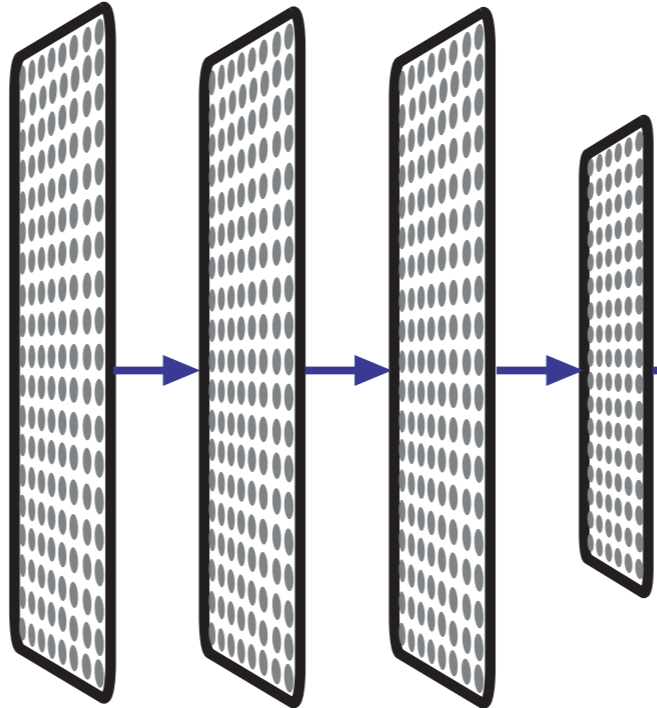


Bence Ölveczky

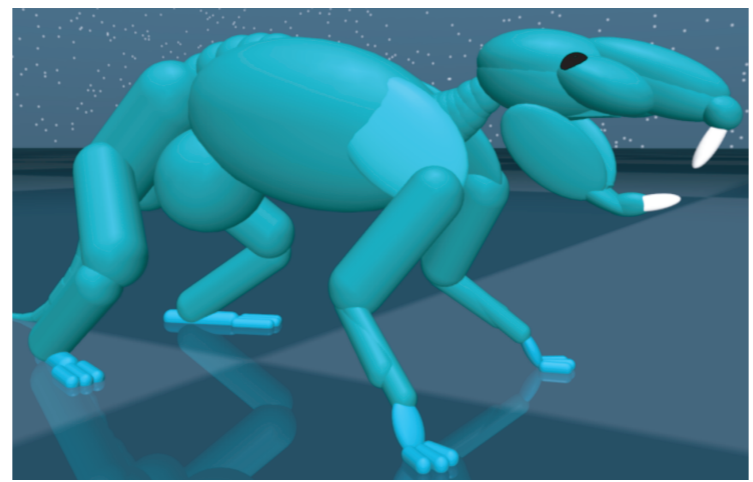
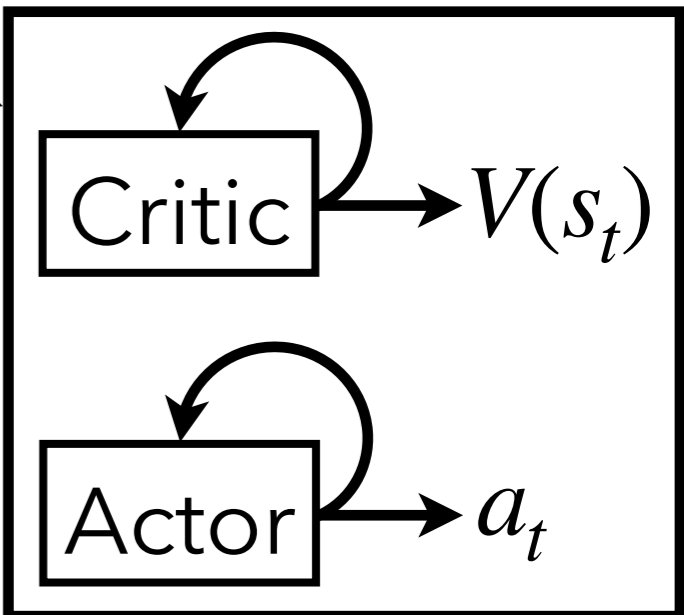
Embodied Virtual Rodent Navigation



Vision Network



Decision Making



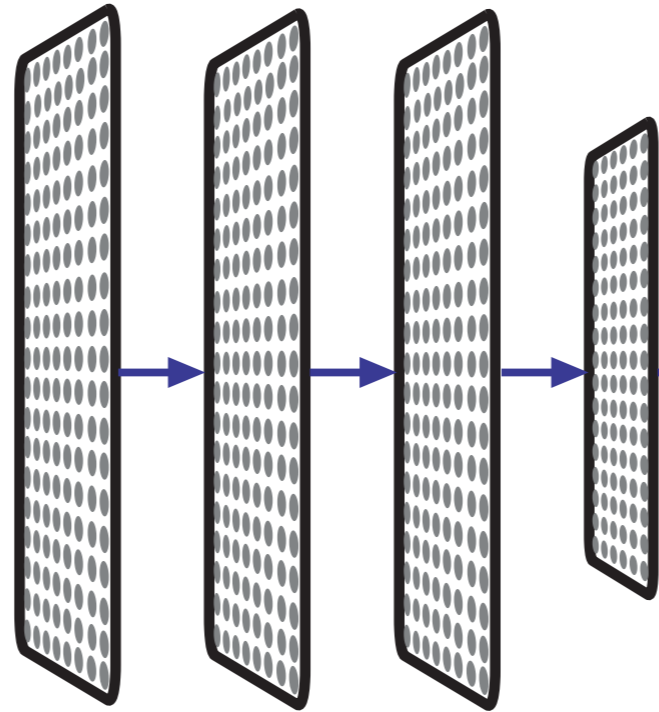
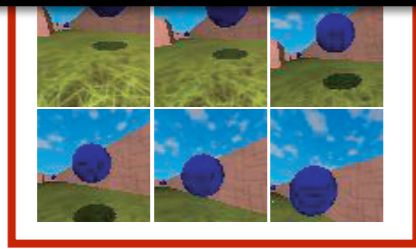
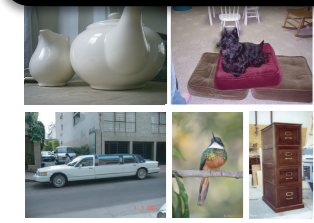
Biomechanical Model

(Joint angles, accelerometer, etc.)

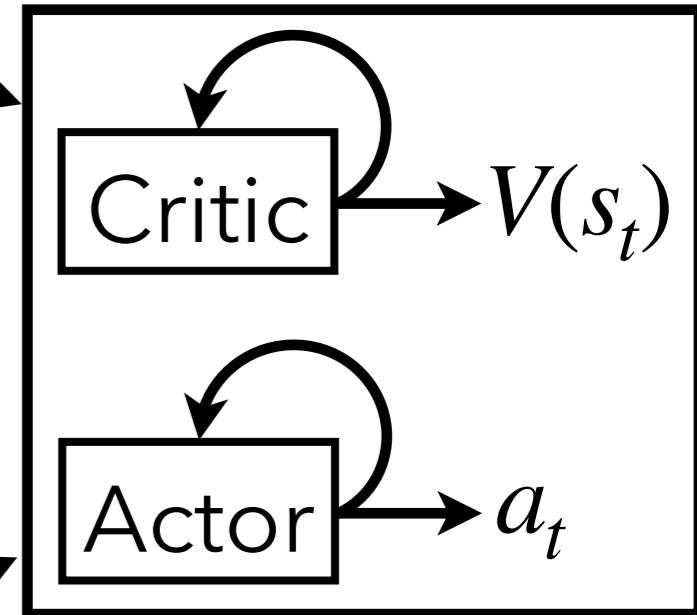


Bence Ölveczky

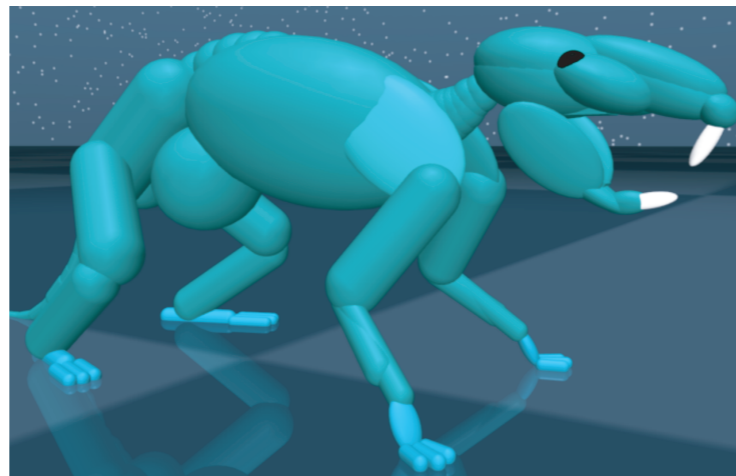
High degree-of-freedom body, keeping track of history over long timescales with high-dimensional, continuous inputs



Decision Making



Bence Ölveczky

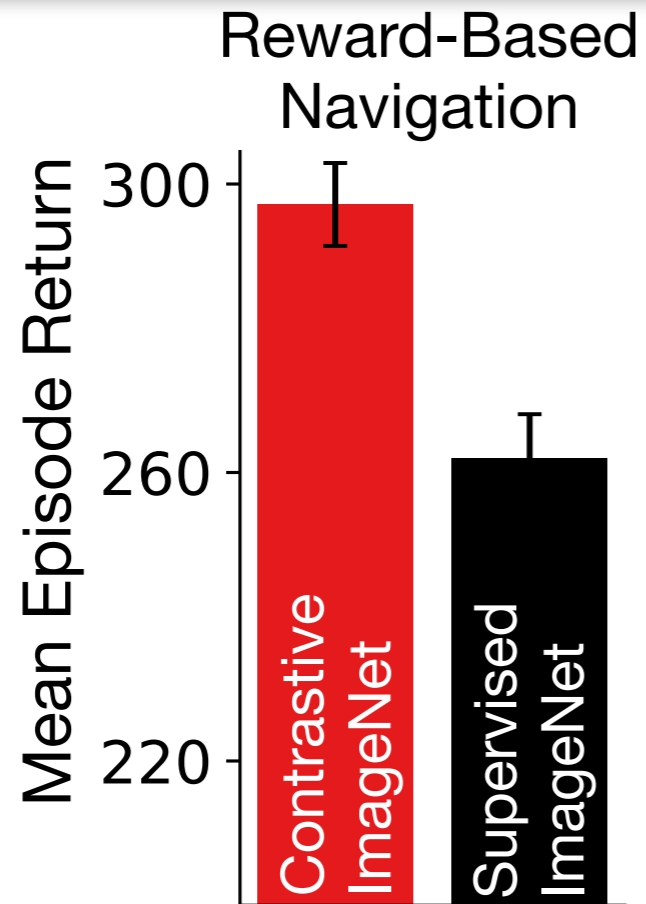


Biomechanical Model

(Joint angles, accelerometer, etc.)

Contrastive Models Yield Better Transfer Performance

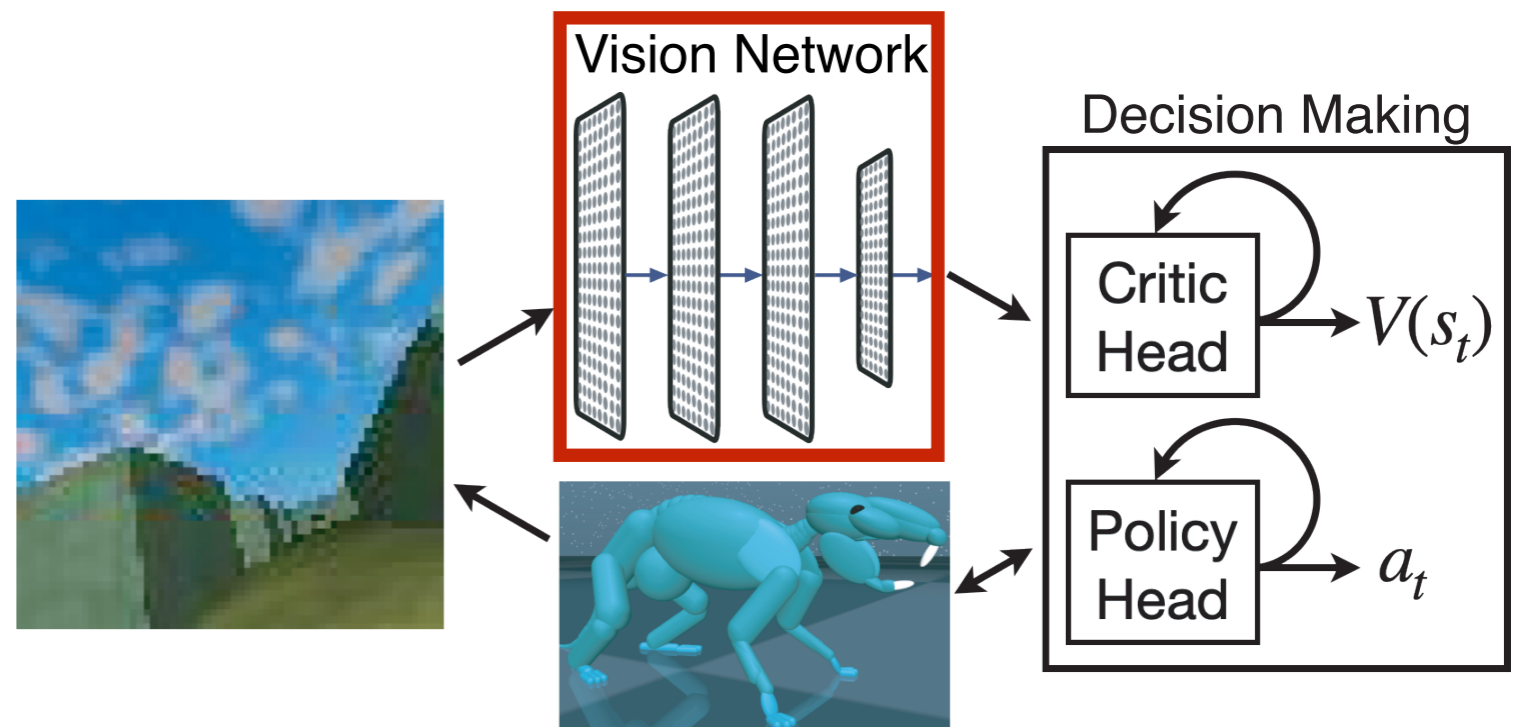
Contrastive Models Yield Better Transfer Performance



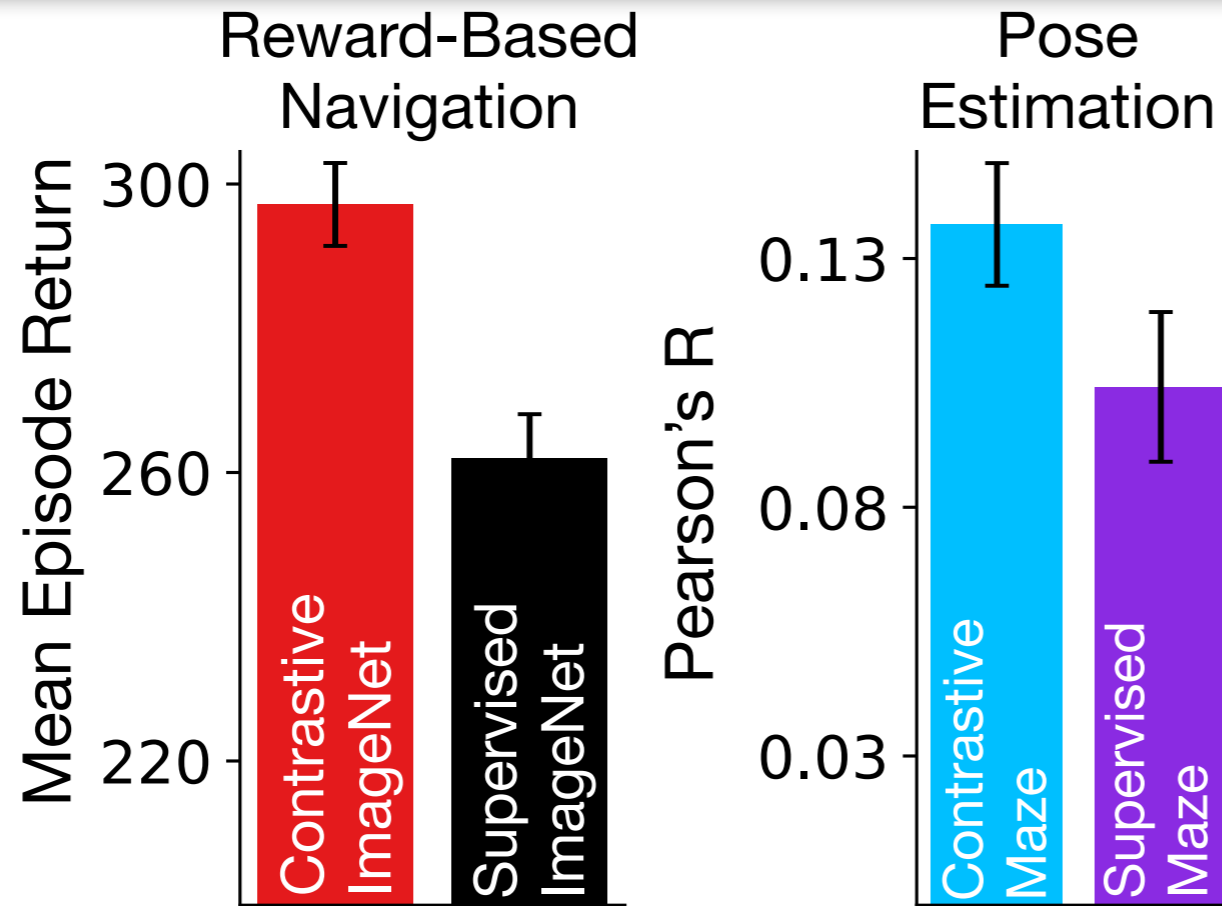
Train *ImageNet*



Evaluate

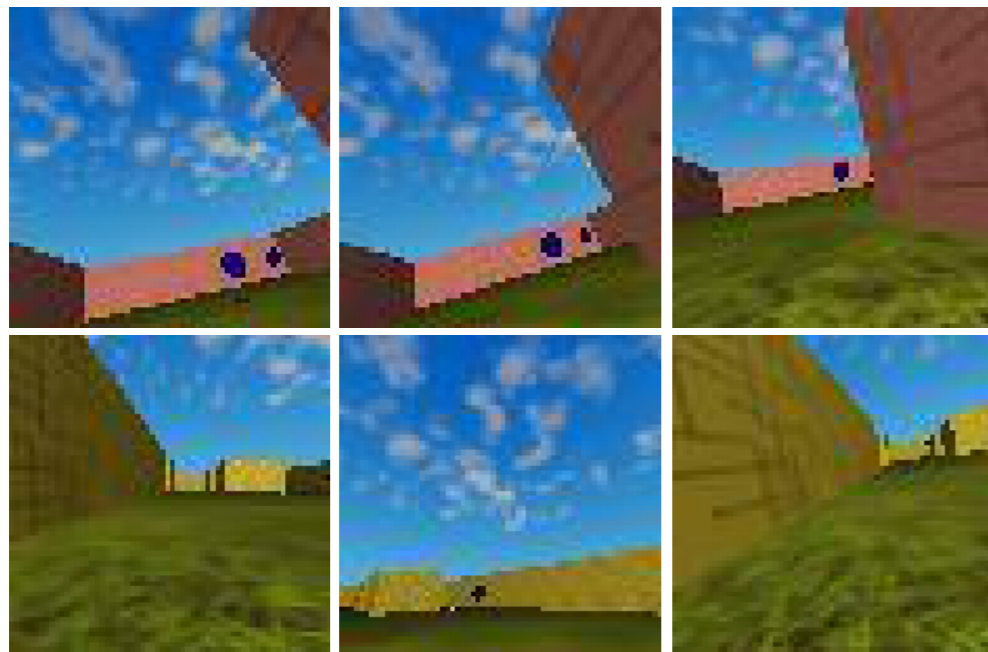


Contrastive Models Yield Better Transfer Performance



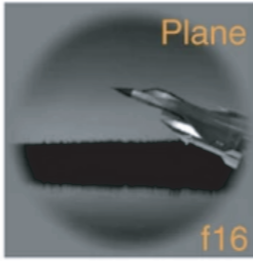
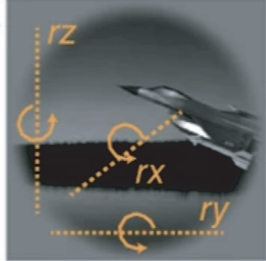

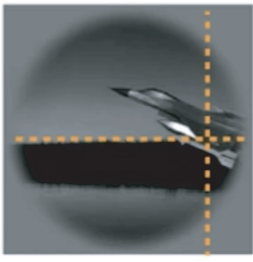
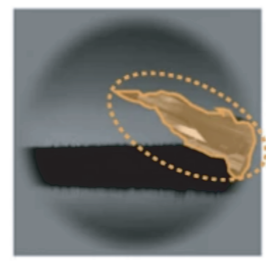
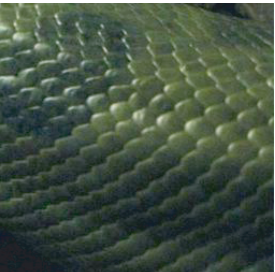
Train

Maze Environment

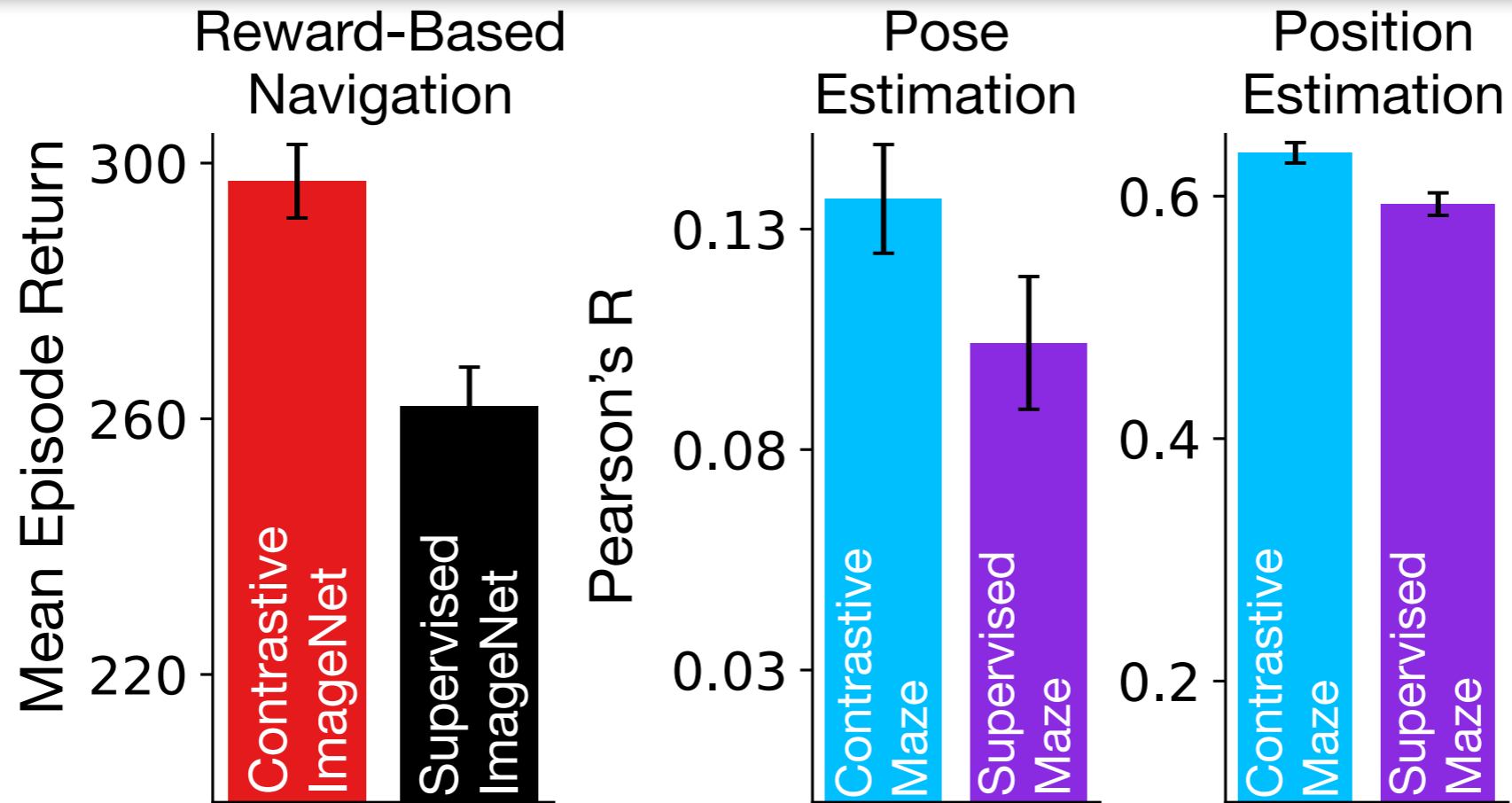


Evaluate

Visual Scene Understanding

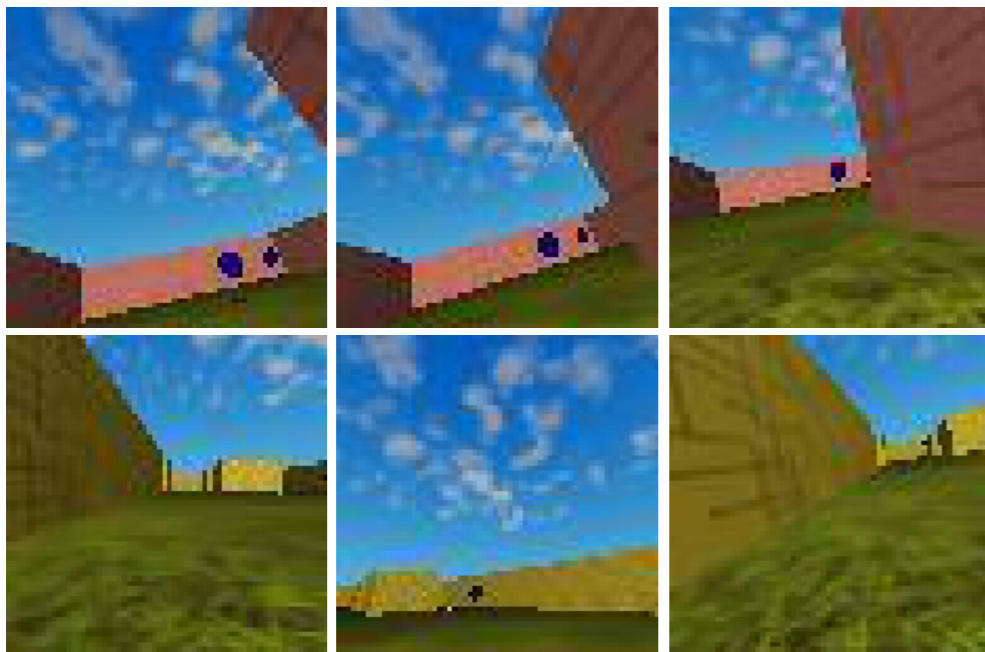
 <p>Plane</p>	Category	 <p>z axis rotation x axis rotation y axis rotation</p>	
 <p>f16</p>	Identity	 <p>Perimeter: 78 pix Two-dimensional retinal area: 146 pix Three-dimensional object scale: 1.2x</p>	
<i>Object properties</i>			<i>Texture</i>

Contrastive Models Yield Better Transfer Performance



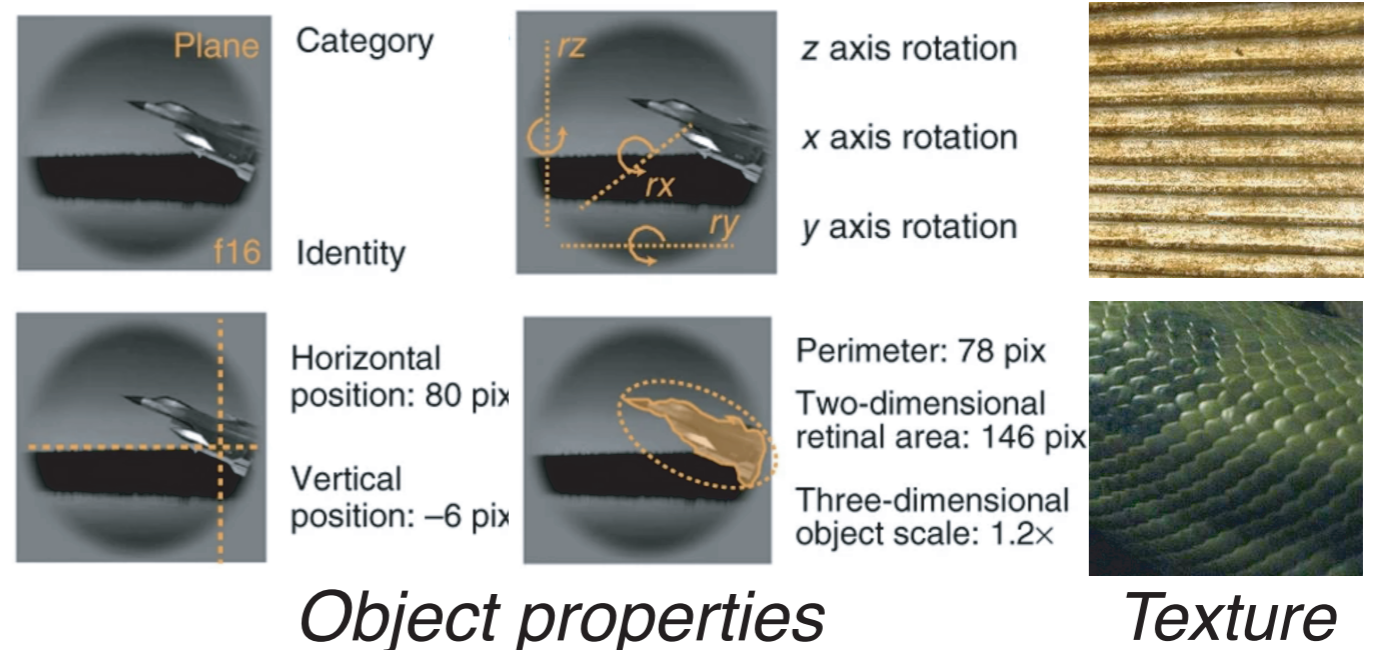
Train

Maze Environment

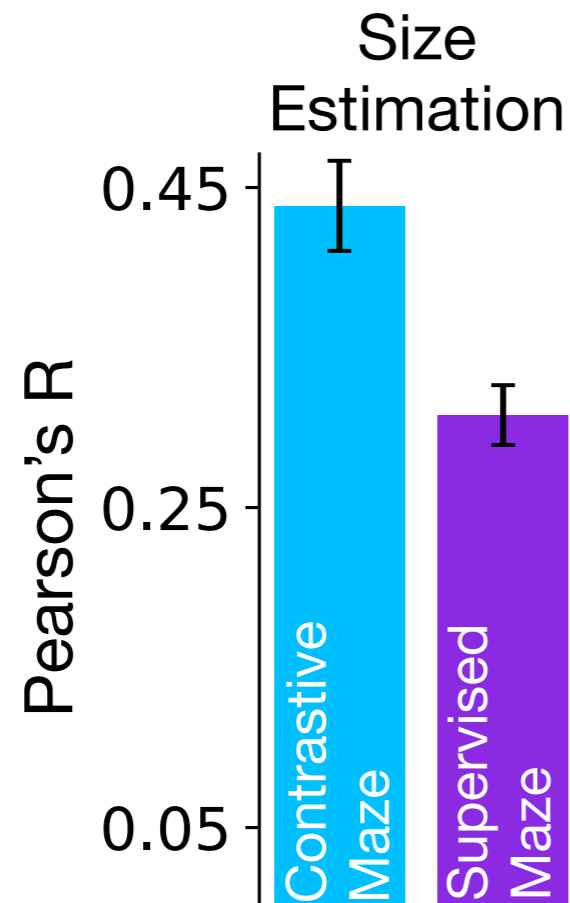
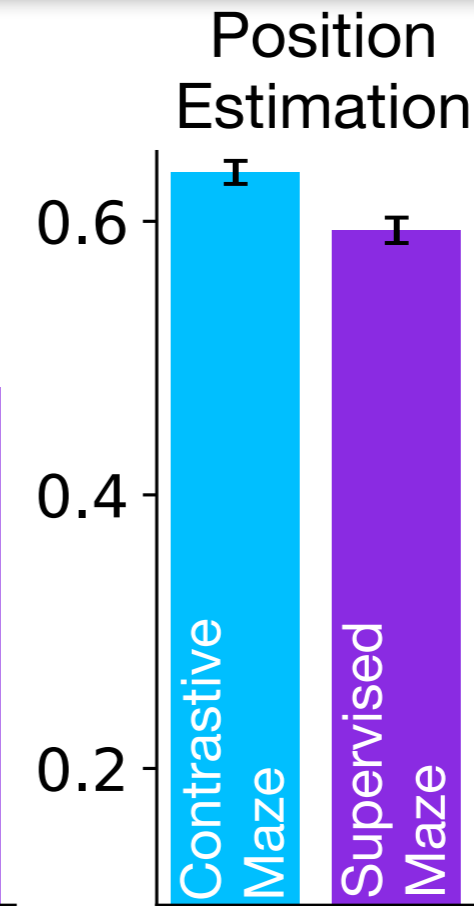
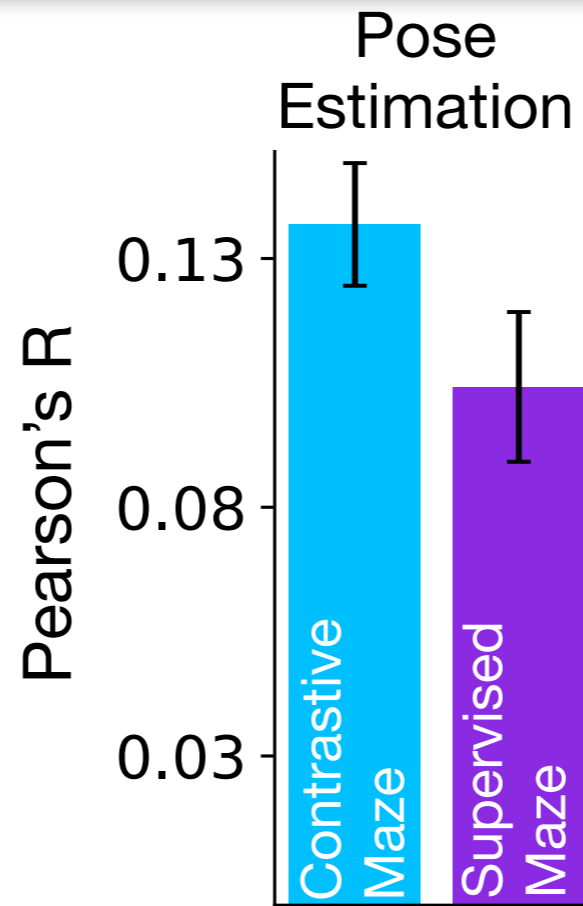
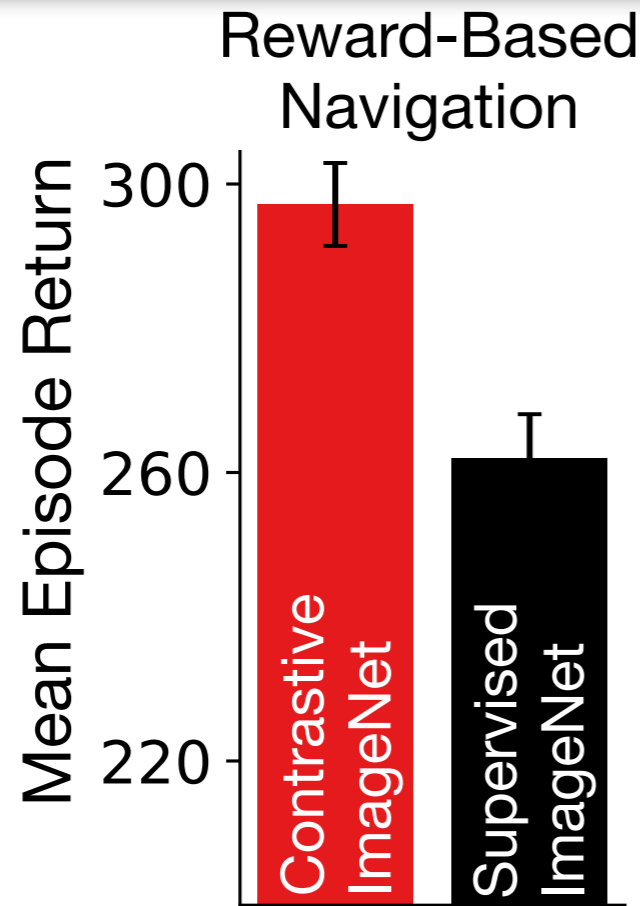


Evaluate

Visual Scene Understanding

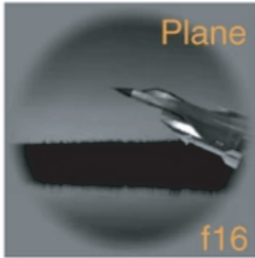
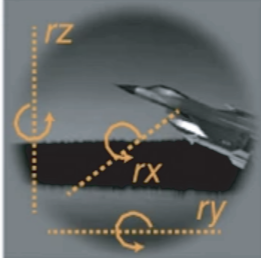

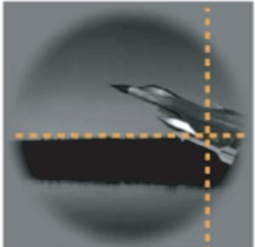
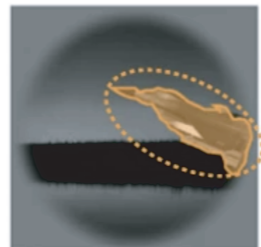
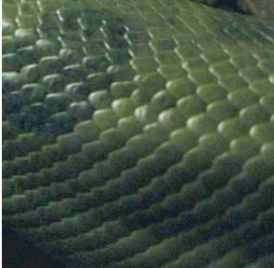


Contrastive Models Yield Better Transfer Performance

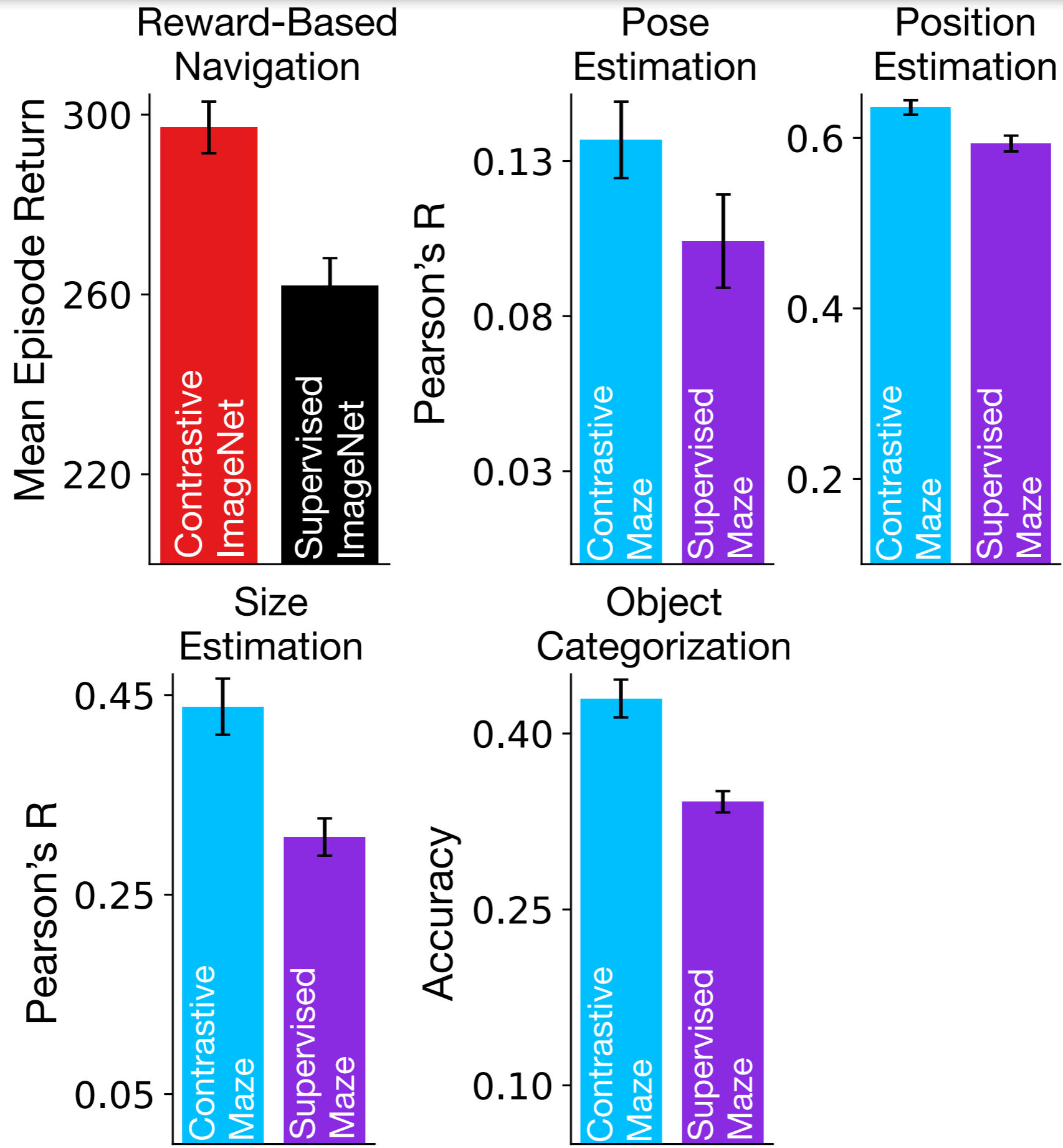


Evaluate

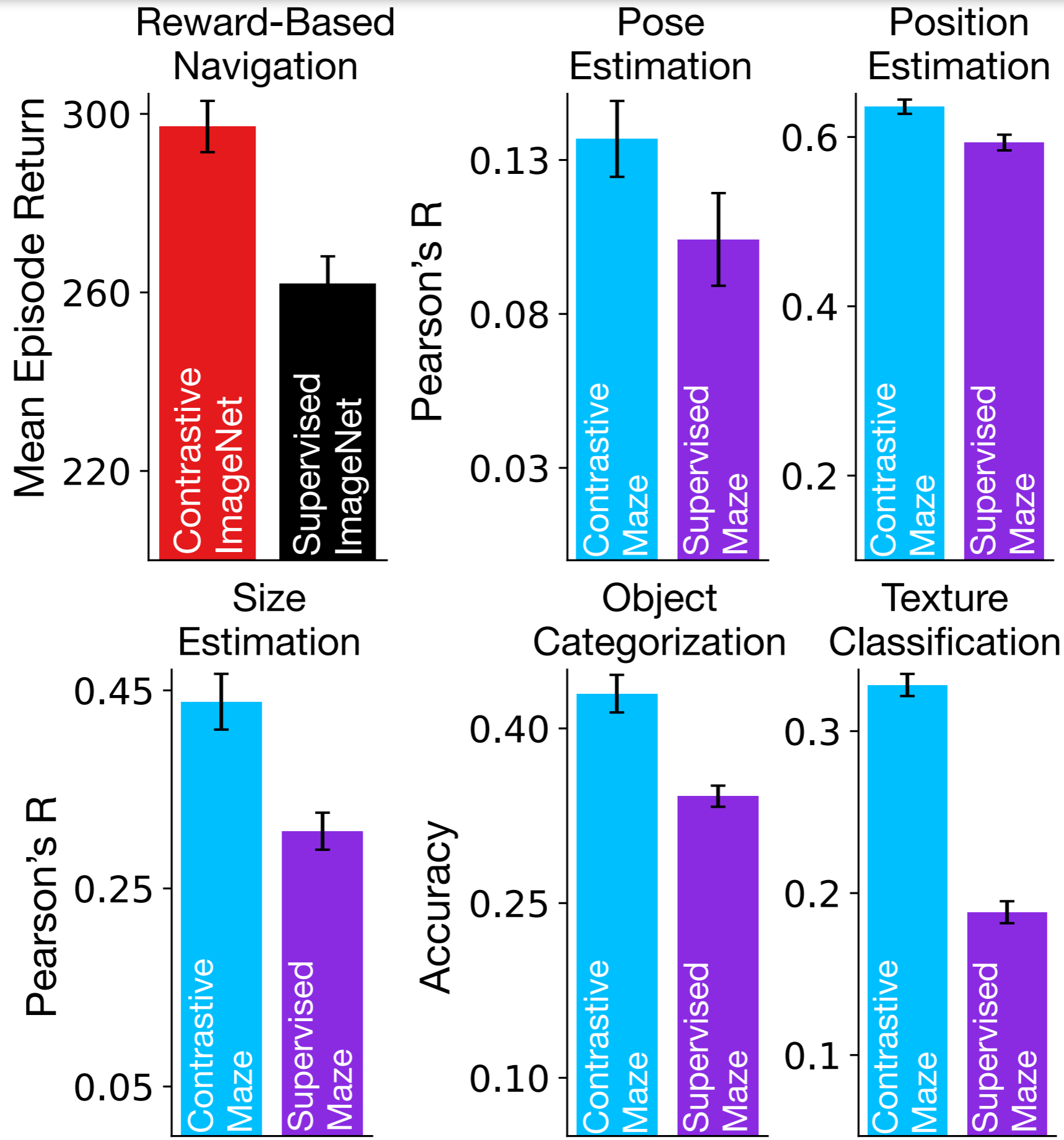
Visual Scene Understanding

 <p>Category</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>

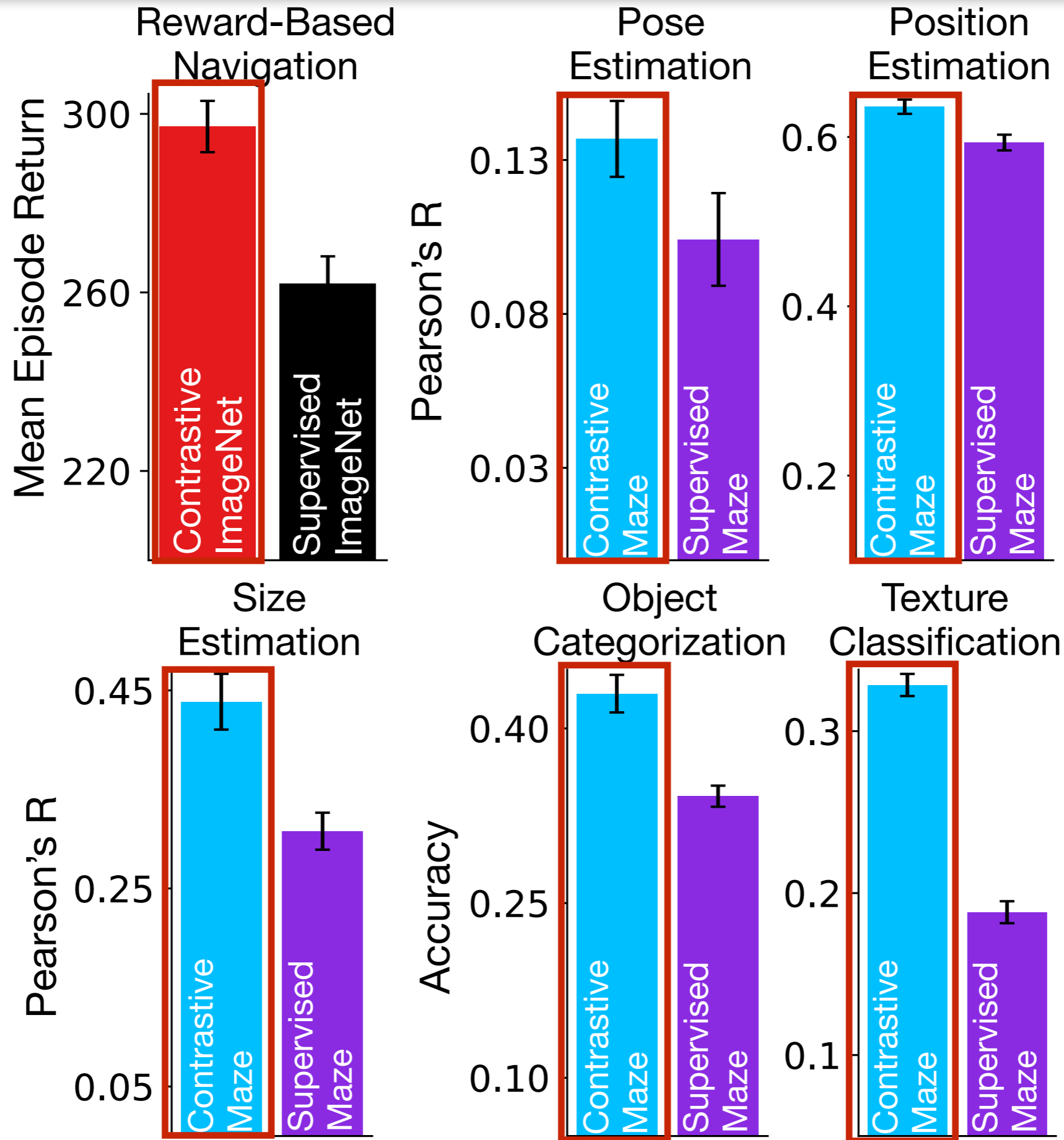
Contrastive Models Yield Better Transfer Performance



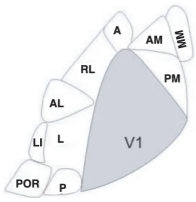
Contrastive Models Yield Better Transfer Performance



Contrastive Models Yield Better Transfer Performance



The best neural models have the best task transfer



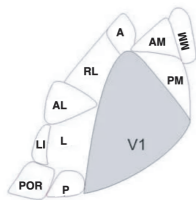
Noise-Corrected Neural
Predictivity (Pearson's R)

0.5
0.4
0.3
0.2

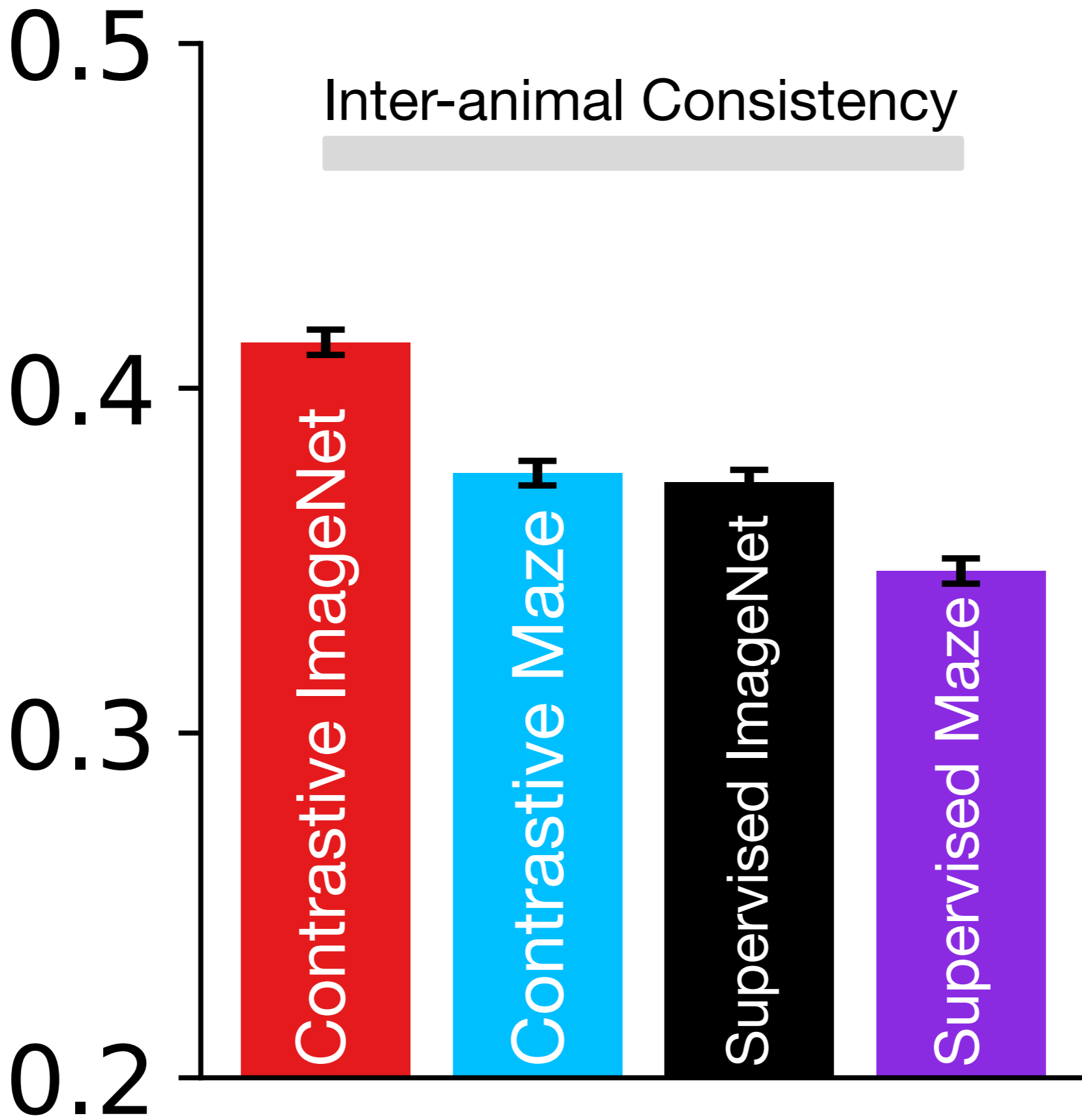
Inter-animal Consistency

Do the contrastive methods
that task generalize best,
also match the neurons better?

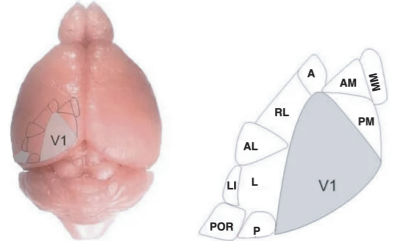
The best neural models have the best task transfer



Noise-Corrected Neural
Predictivity (Pearson's R)



The best neural models have the best task transfer



Noise-Corrected Neural
Predictivity (Pearson's R)

0.2 0.3 0.4 0.5

Contrastive ImageNet H

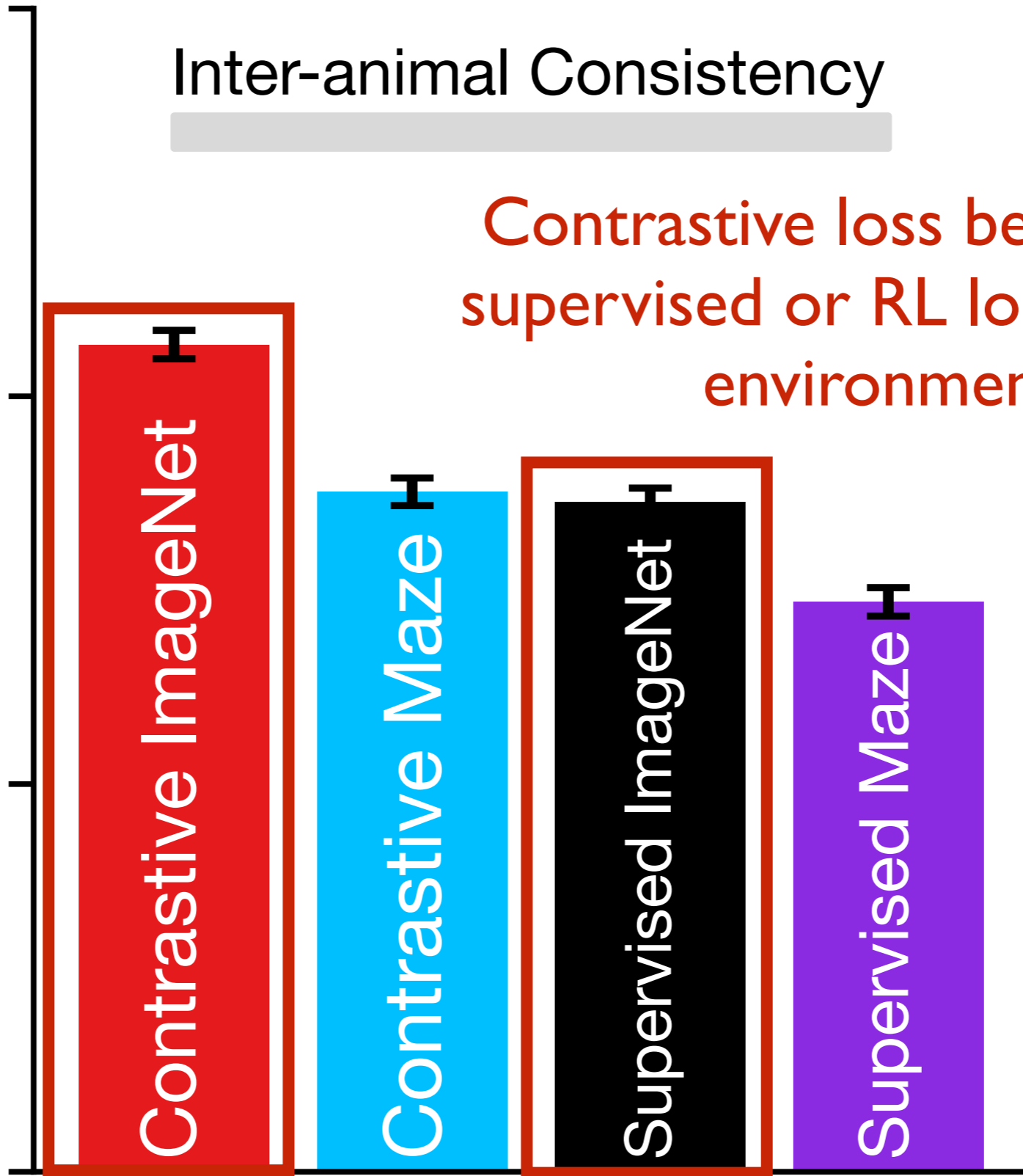
Contrastive Maze H

Supervised ImageNet H

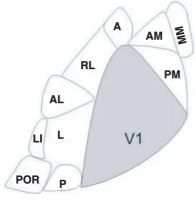
Supervised Maze H

Inter-animal Consistency

Contrastive loss better than supervised or RL loss in *same* environment



The best neural models have the best task transfer



Noise-Corrected Neural
Predictivity (Pearson's R)

0.2 0.3 0.4 0.5

Contrastive ImageNet H

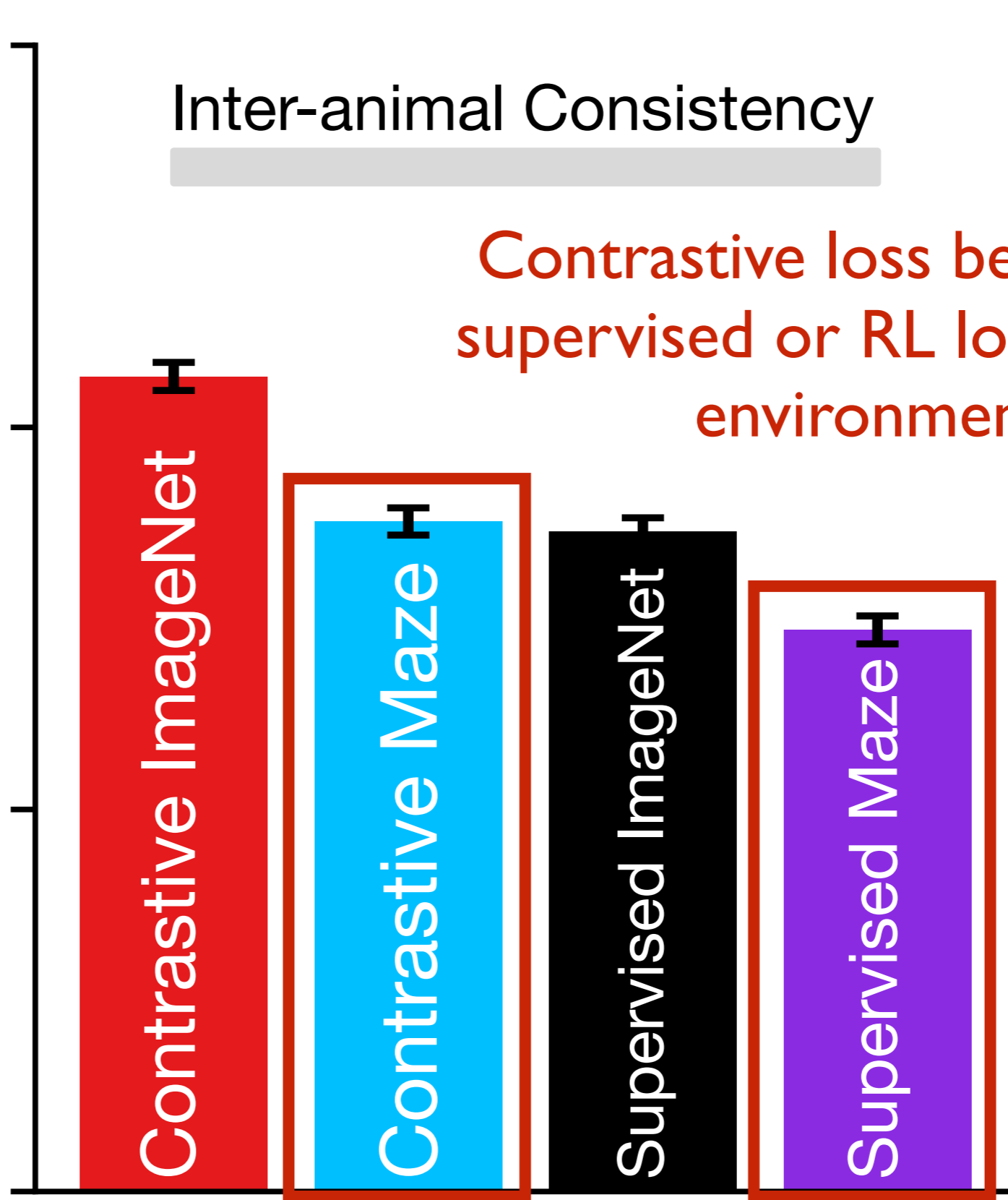
Contrastive Maze H

Supervised ImageNet H

Supervised Maze H

Inter-animal Consistency

Contrastive loss better than supervised or RL loss in *same* environment



Outline

- ▶ Mouse Visual Cortex as a Task-General, Limited Resource System
- ▶ Reusable Latent Representations for Primate Mental Simulation
- ▶ Heuristics for Interrogating Natural Intelligence

Outline

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

Mouse visual cortex (so far) is a low-acuity, shallow network that makes best use of the mouse's limited resources to create a general-purpose visual system, that can be deployed in novel environments and embodied contexts.

▶ Reusable Latent Representations for Primate Mental Simulation

▶ Heuristics for Interrogating Natural Intelligence

Outline

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

Mouse visual cortex (so far) is a low-acuity, shallow network that makes best use of the mouse's limited resources to create a general-purpose visual system, that can be deployed in novel environments and embodied contexts.

- ▶ Reusable Latent Representations for Primate Mental Simulation

- ▶ Heuristics for Interrogating Natural Intelligence

Reusable Latent Representations for Primate Mental Simulation

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

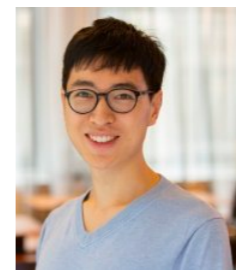
Neural foundations of mental simulation: future prediction of latent representations on dynamic scenes.
NeurIPS 2023 (spotlight)



Rishi Rajalingham



Mehrdad Jazayeri



Guangyu Robert Yang

Visually-Grounded Mental Simulation

Visually-Grounded Mental Simulation



Visually-Grounded Mental Simulation

Infer:

Has this ice block been out longer?



Infer:

Has this ice block been out longer?

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?



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?

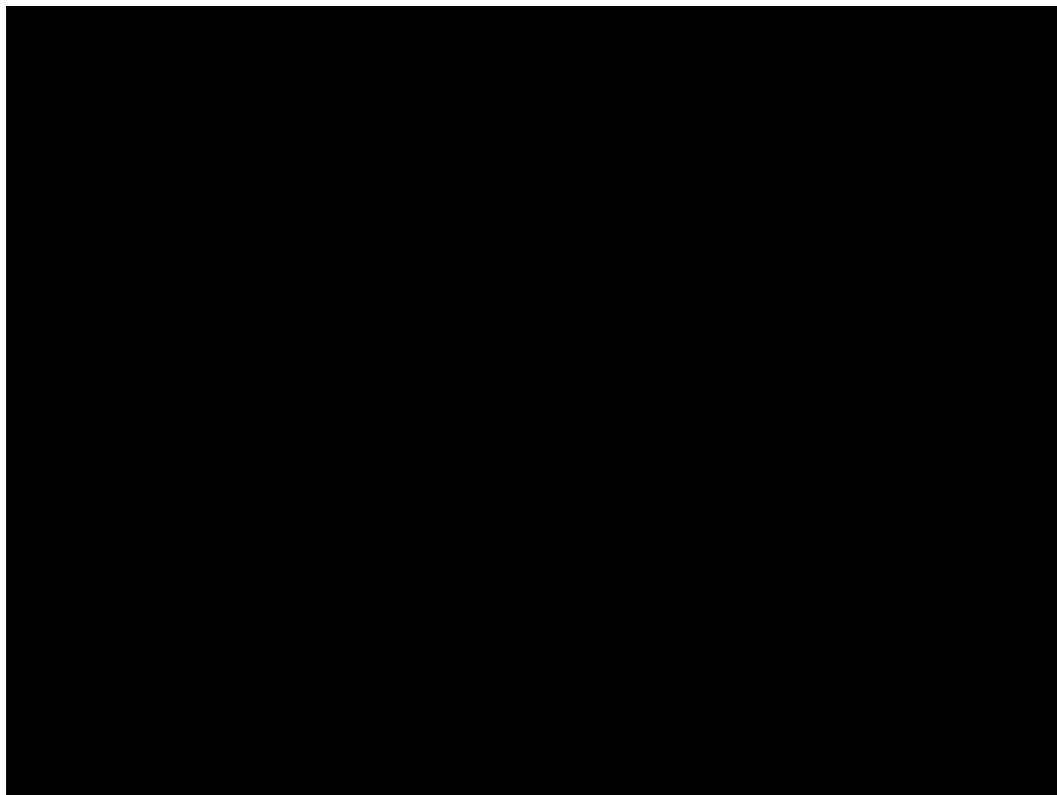


Plan:

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?



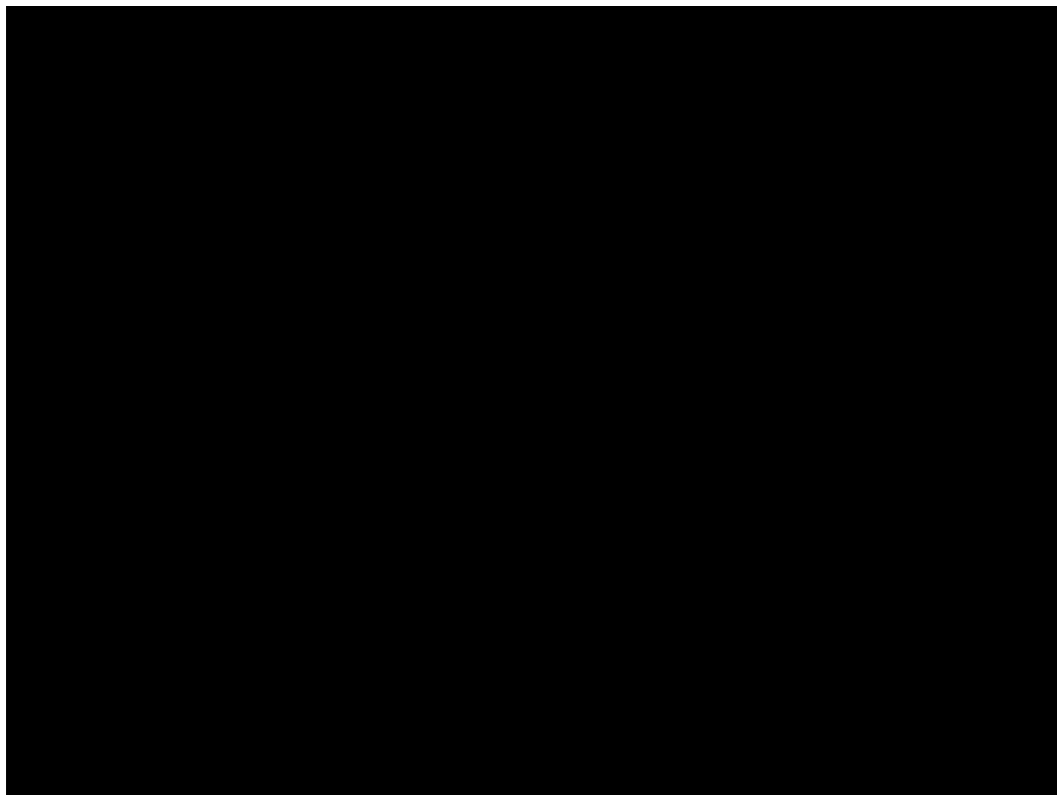
Plan:

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?



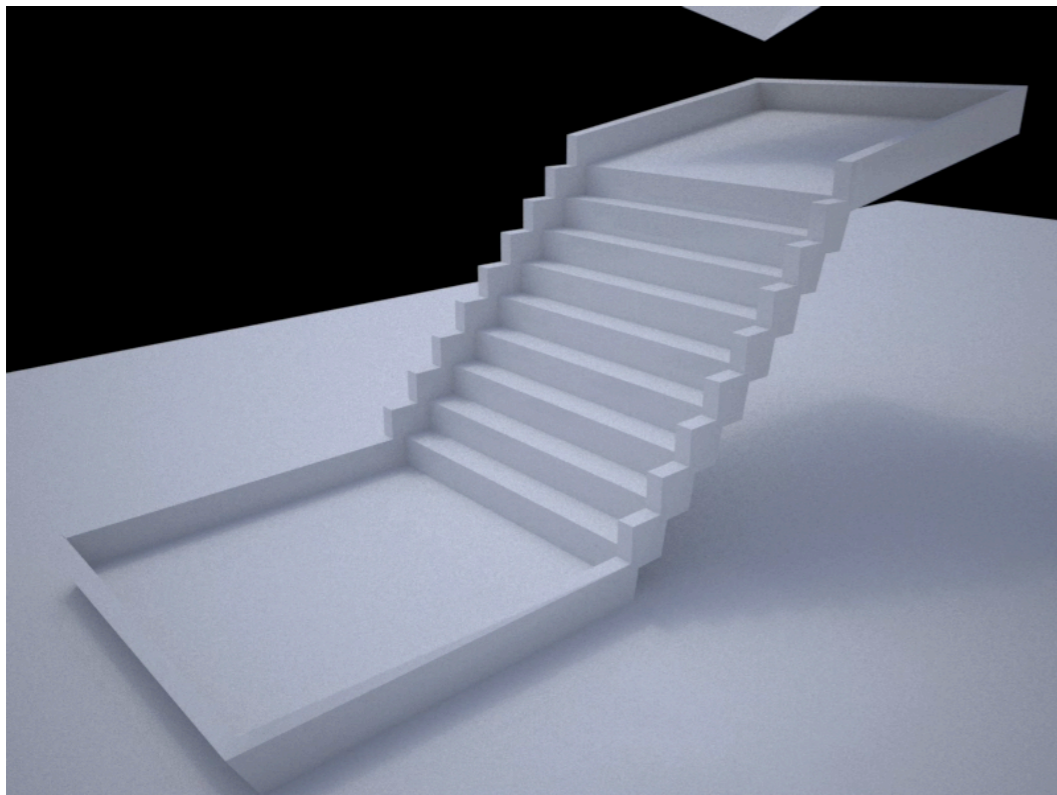
Plan:

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?



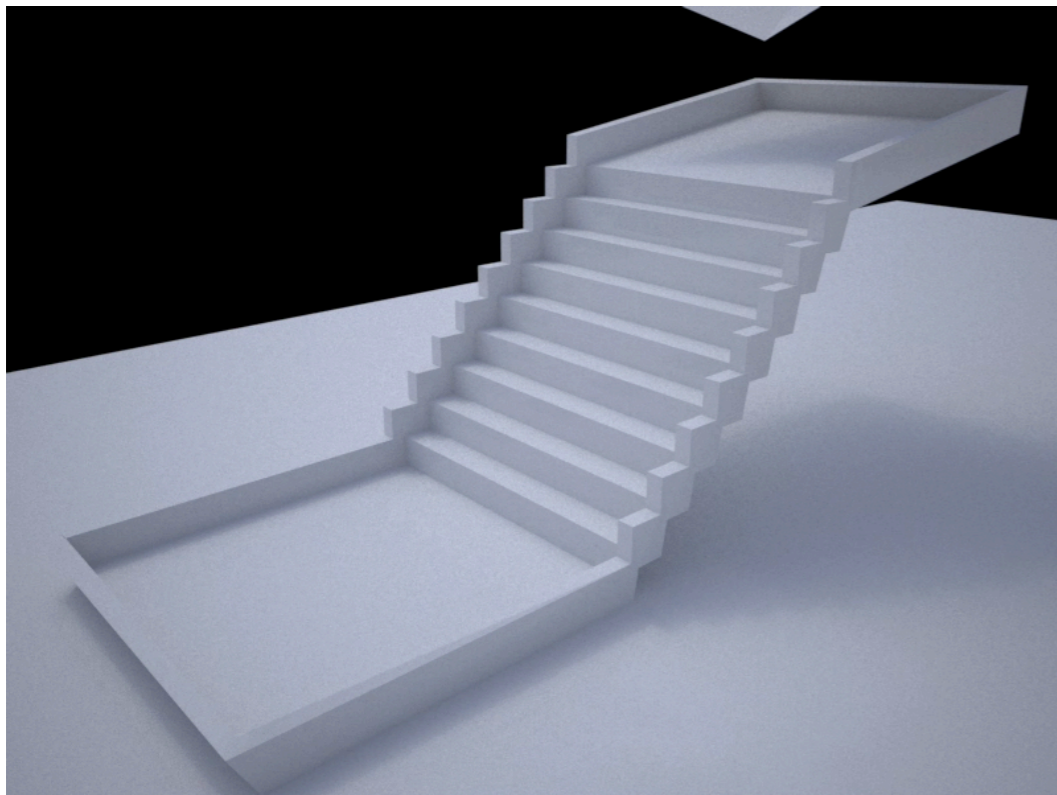
Plan:

How would I take these hats off the rack?



Predict:

Will this box support me?

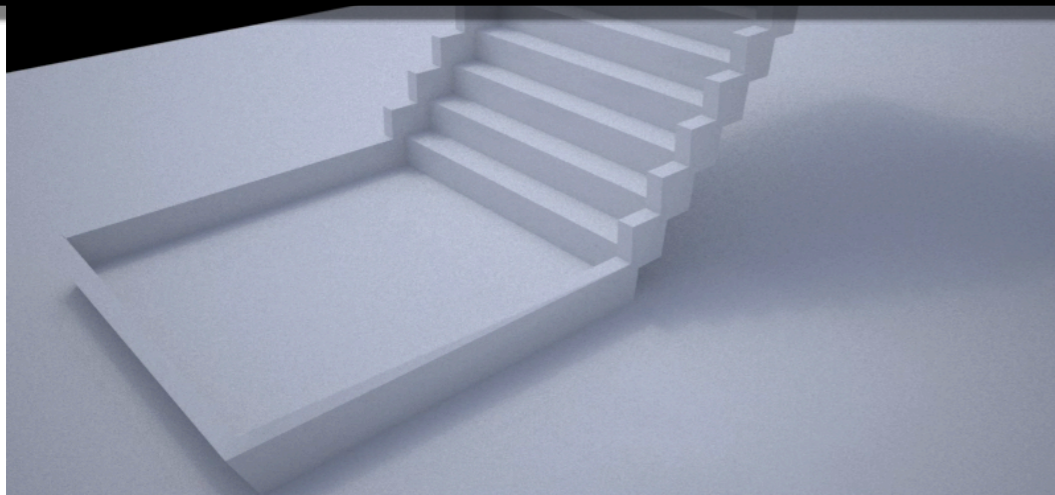


Visually-Grounded Mental Simulation



Neurobiological Puzzle:

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

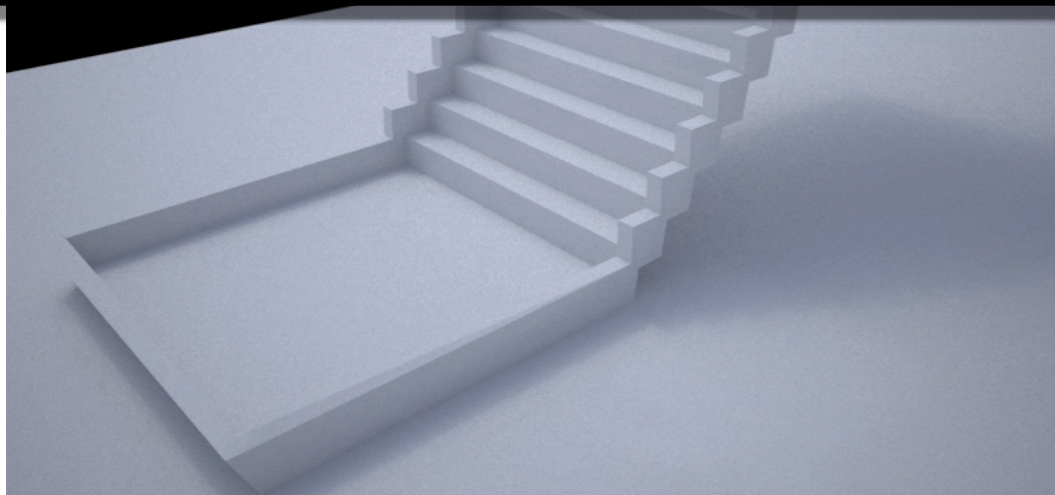


Visually-Grounded Mental Simulation



Neurobiological Puzzle:

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

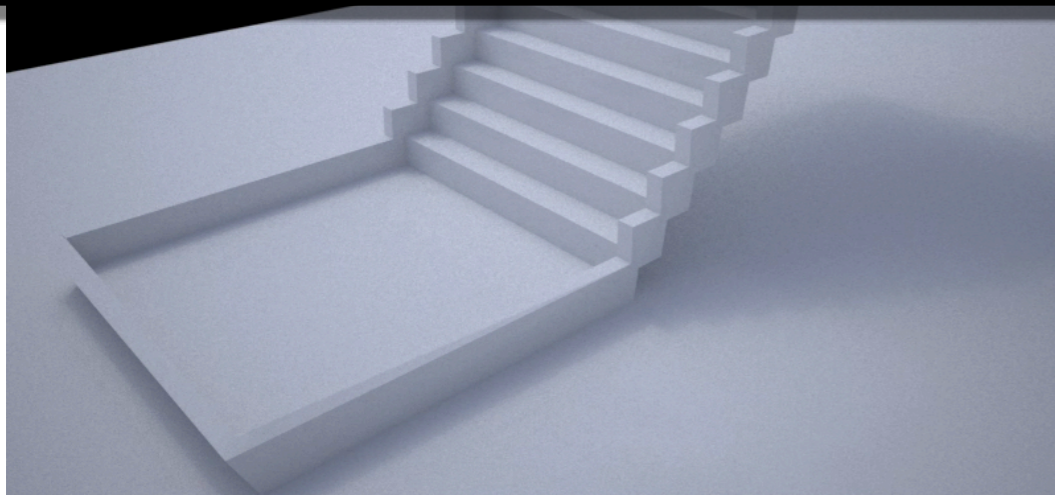


Visually-Grounded Mental Simulation



Neurobiological Puzzle:

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

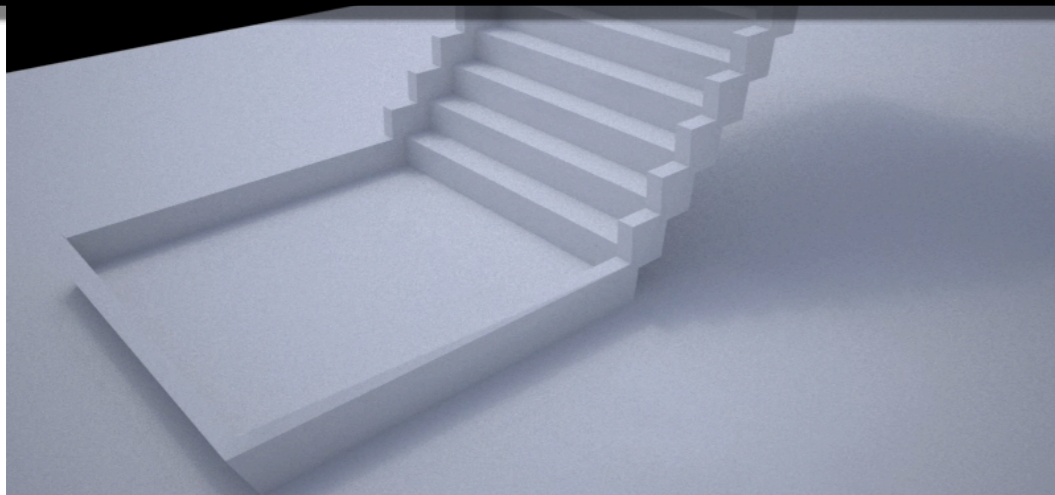


Visually-Grounded Mental Simulation



Neurobiological Puzzle:

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



Defining Hypotheses

Defining Hypotheses

“Sensory-Cognitive Networks”

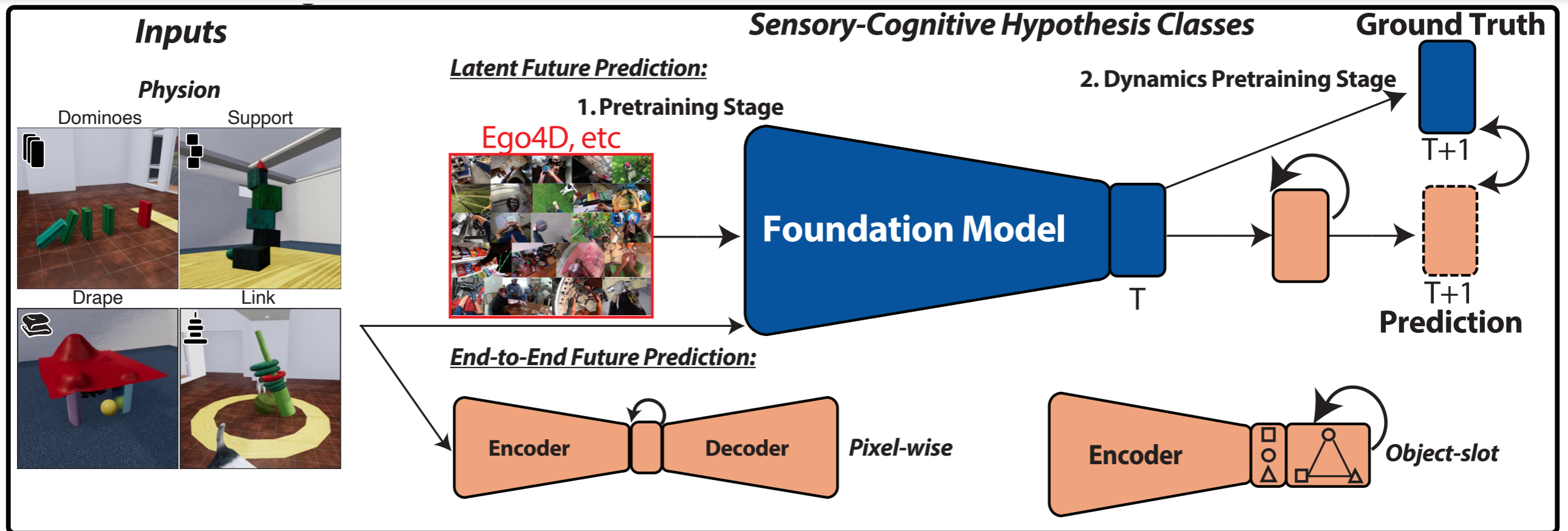
R1 (Input-Driven): Take in unstructured visual inputs across a range of physical phenomena.

R2 (Behavioral Outputs): Generate physical predictions for each scenario (“behavior”).

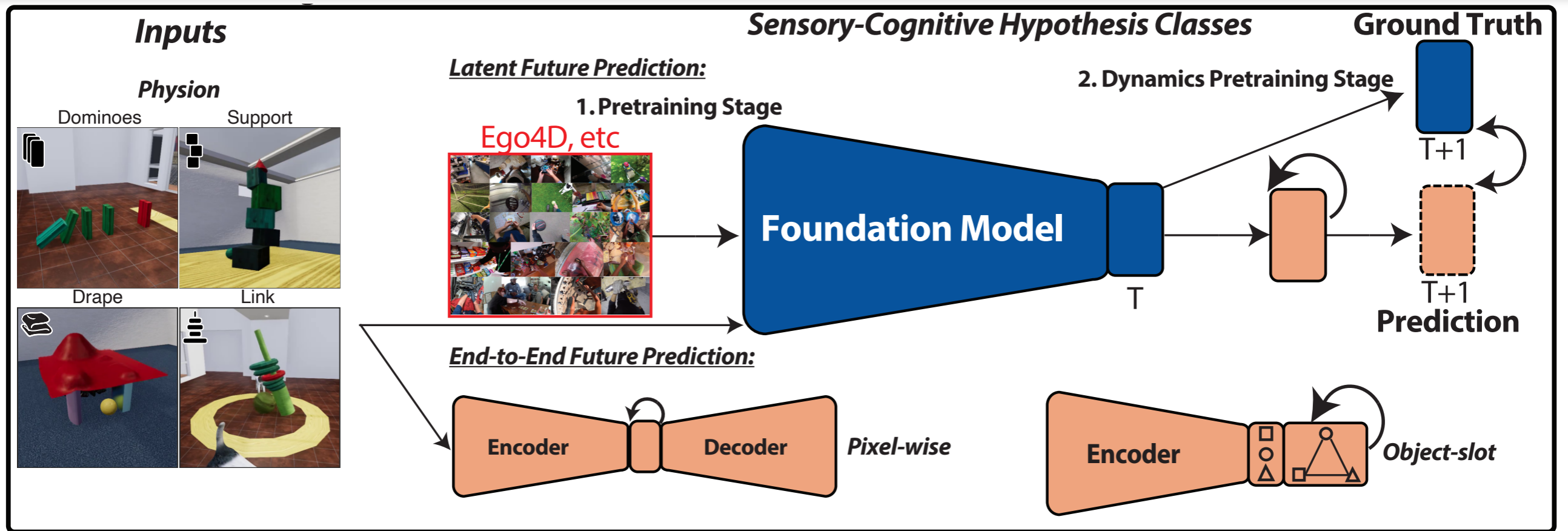
R3 (Neural Representations): Consist of internal units that can be compared to biological units (e.g. containing “artificial neurons”).

Overall Approach: Sensory-Cognitive Hypotheses

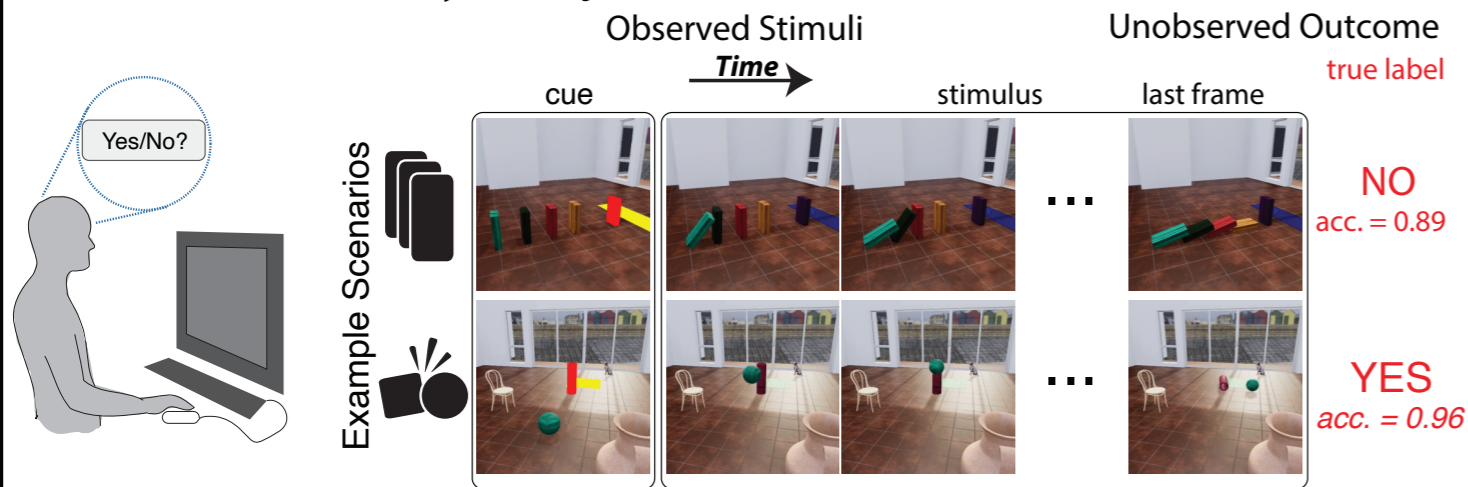
Overall Approach: Sensory-Cognitive Hypotheses



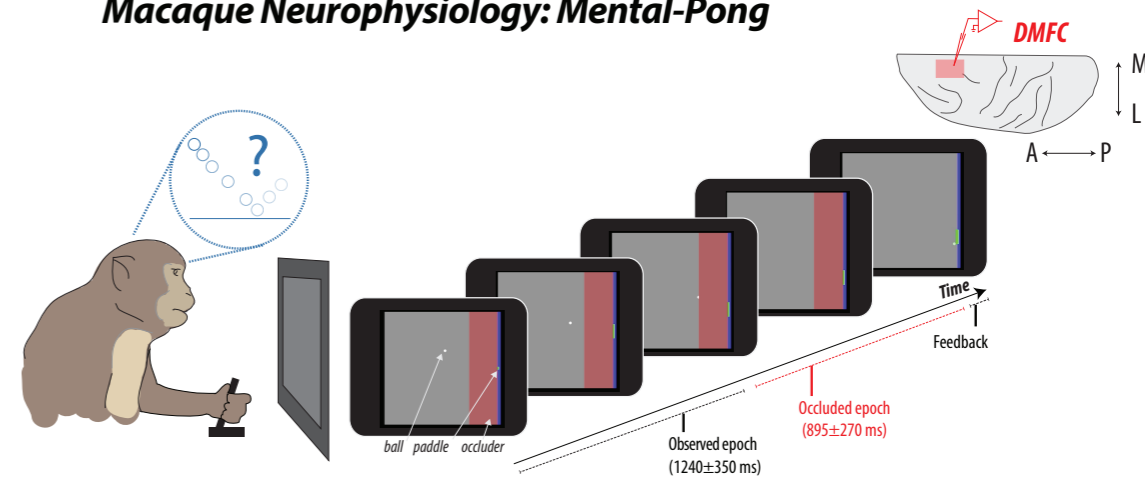
Overall Approach: Sensory-Cognitive Hypotheses



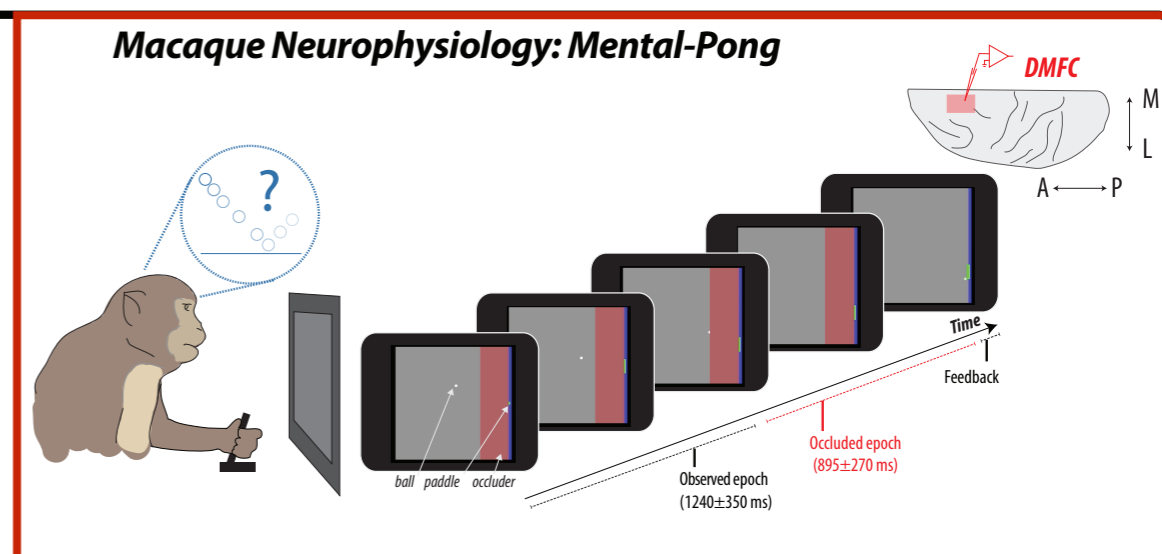
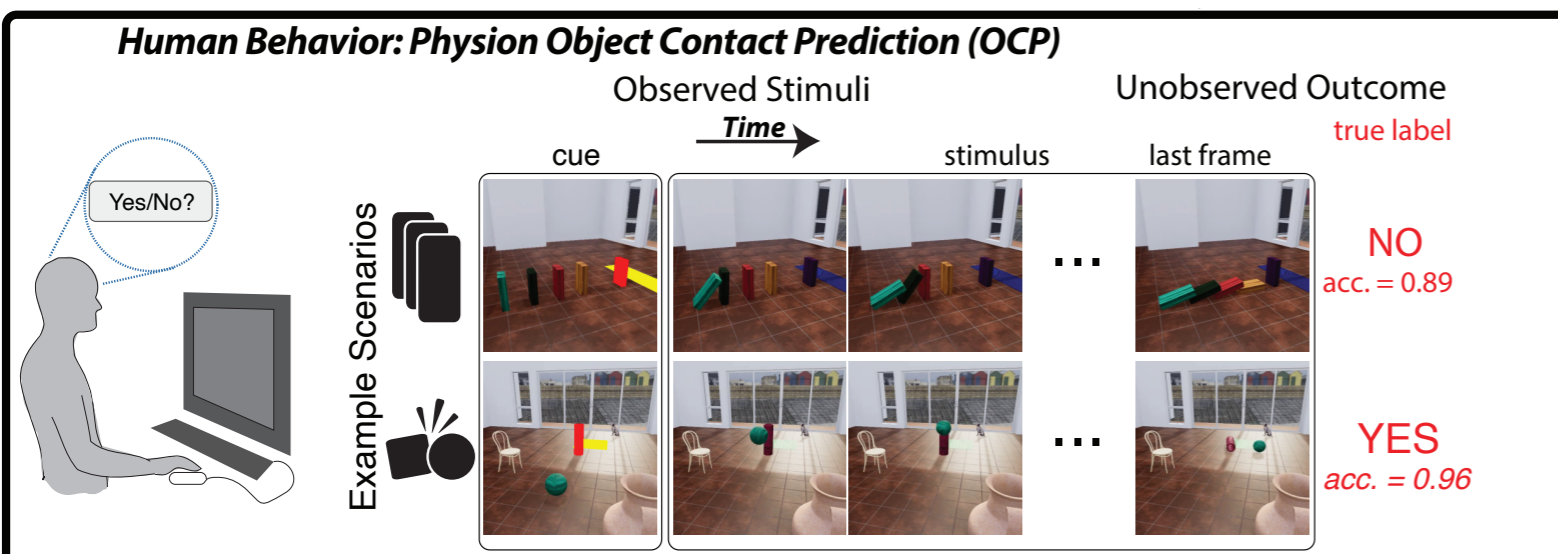
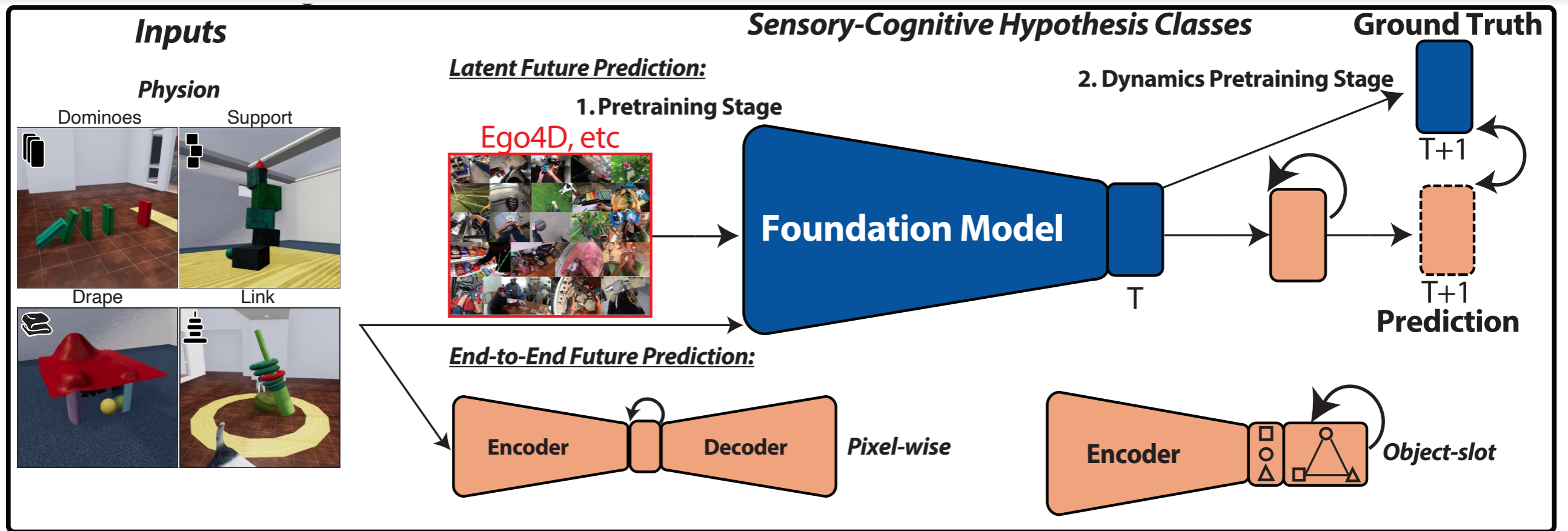
Human Behavior: Physion Object Contact Prediction (OCP)



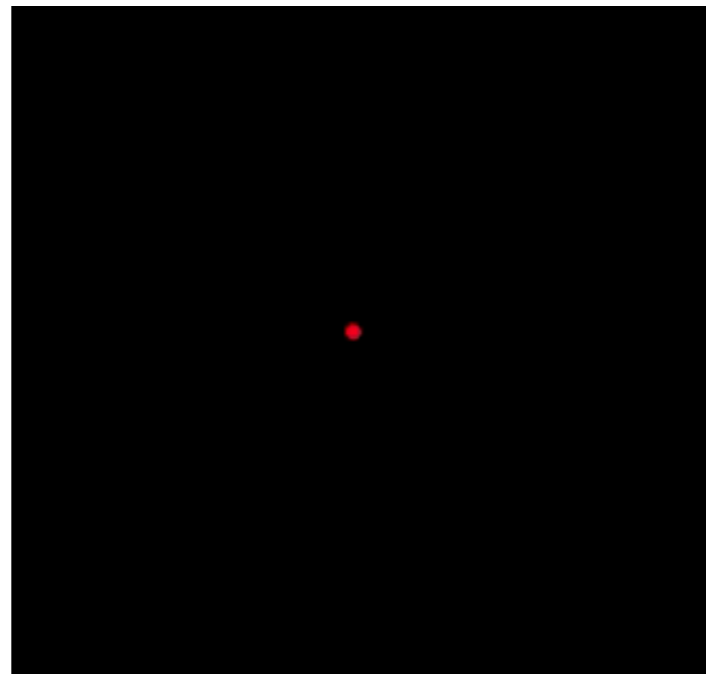
Macaque Neurophysiology: Mental-Pong



Macaque Neurophysiology: Mental Pong

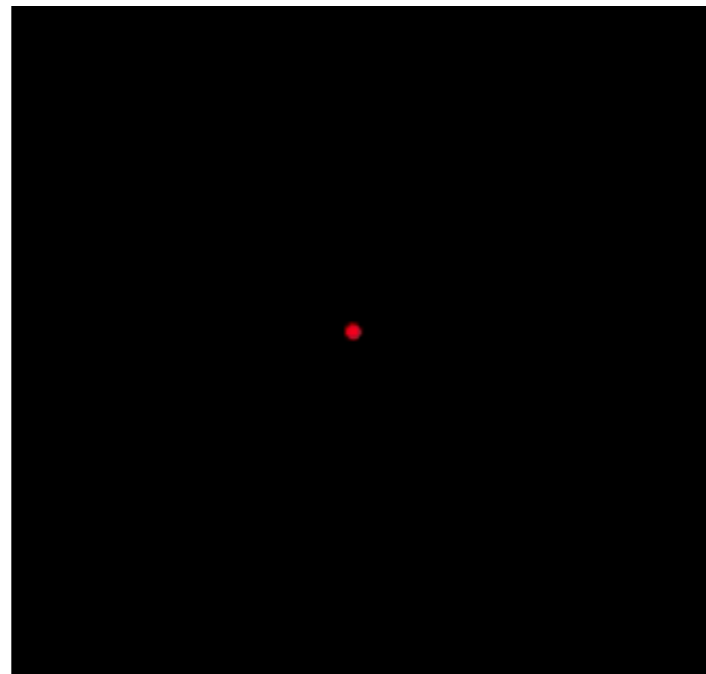


Macaque Neurophysiology: Mental Pong



Rishi Rajalingham

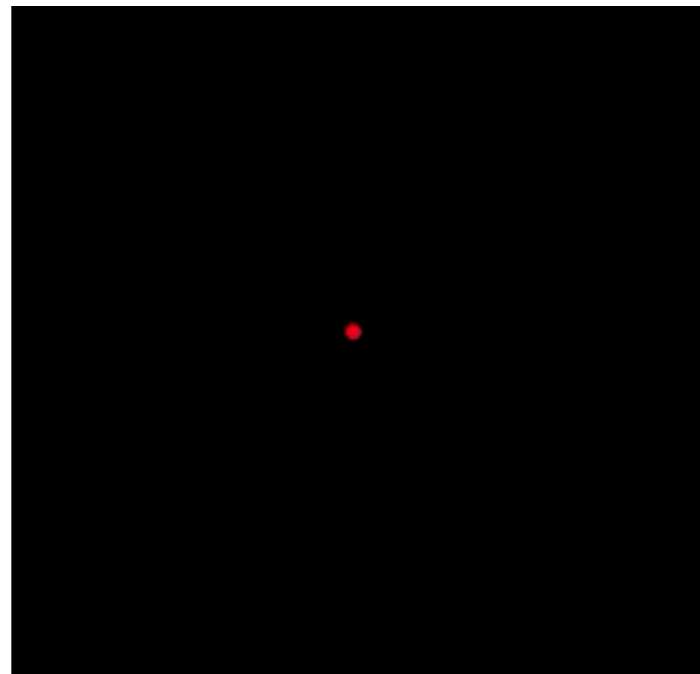
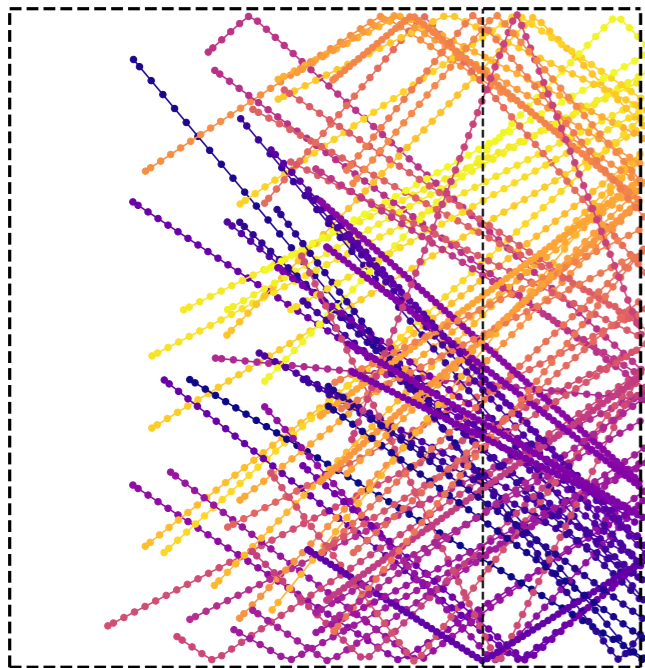
Macaque Neurophysiology: Mental Pong



Rishi Rajalingham

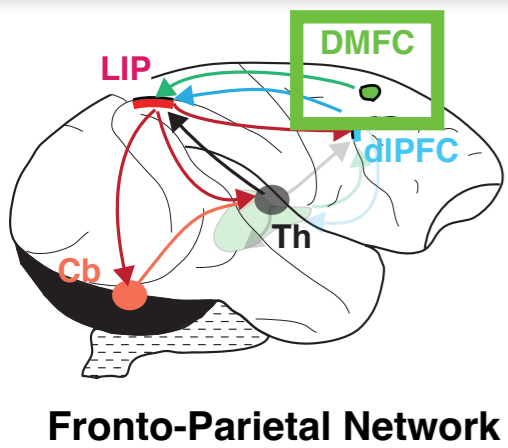
Macaque Neurophysiology: Mental Pong

79 conditions

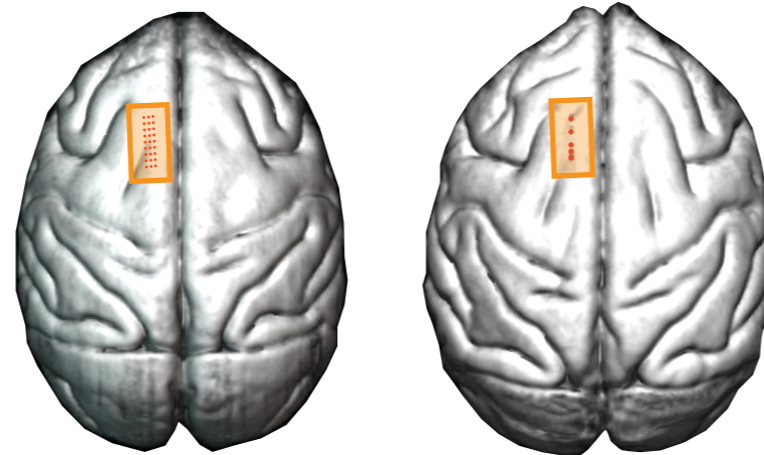


Rishi Rajalingham

Macaque Neurophysiology: Mental Pong

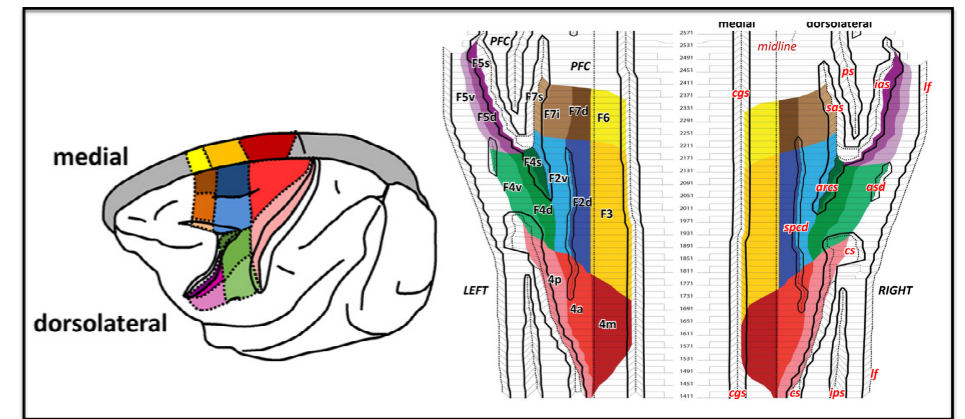


Dorsomedial frontal cortex (DMFC)

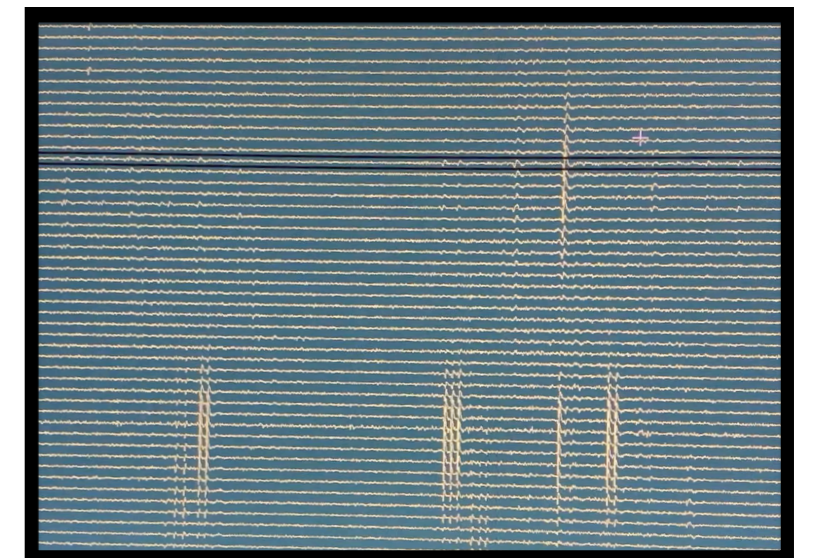
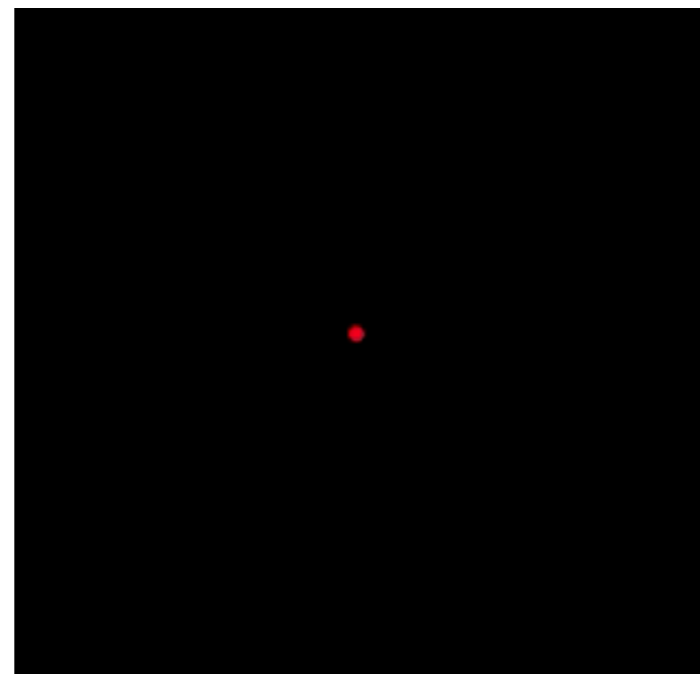
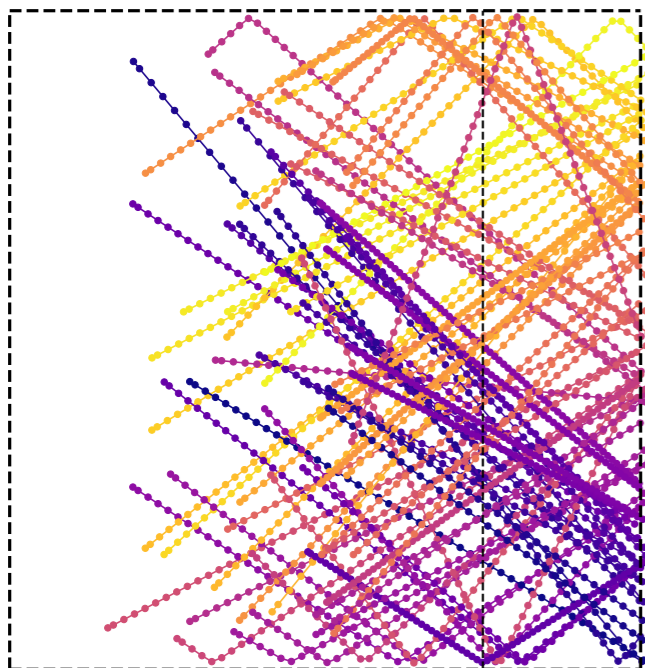


Monkey P

Monkey M

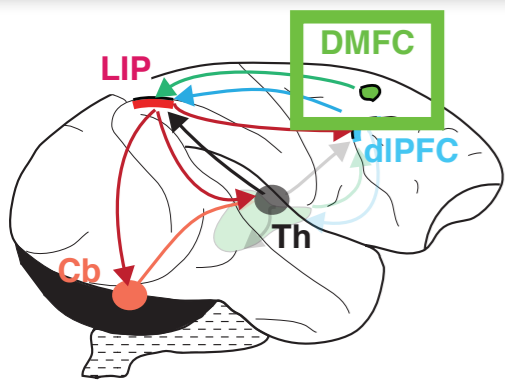


79 conditions



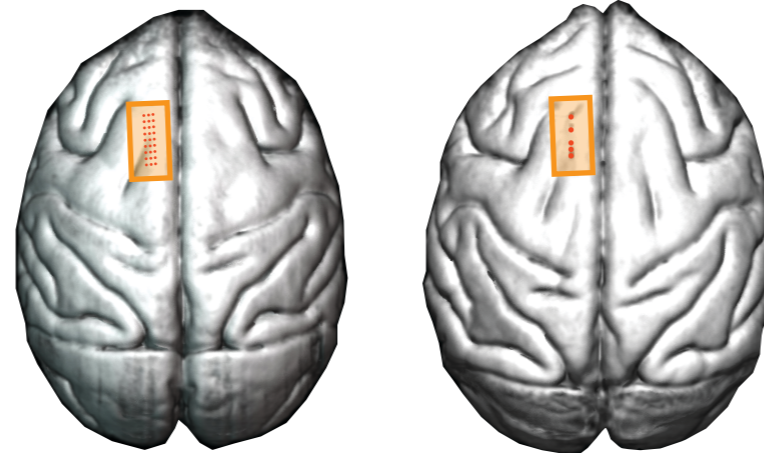
Rishi Rajalingham

Macaque Neurophysiology: Mental Pong



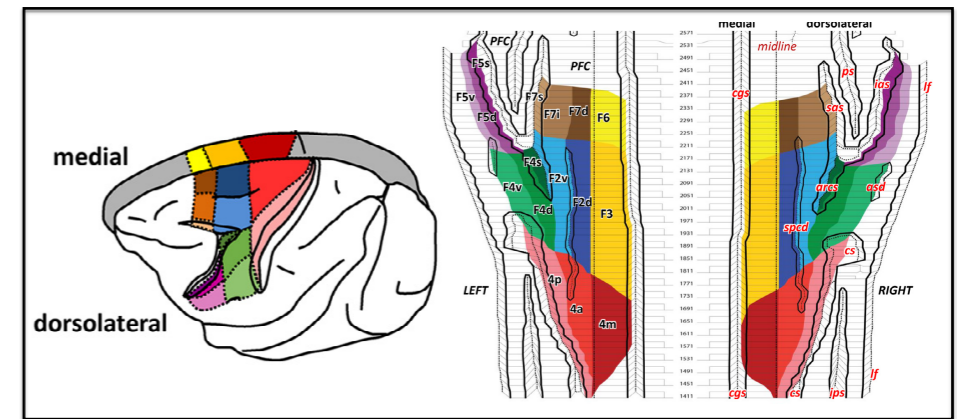
Fronto-Parietal Network

Dorsomedial frontal cortex (DMFC)

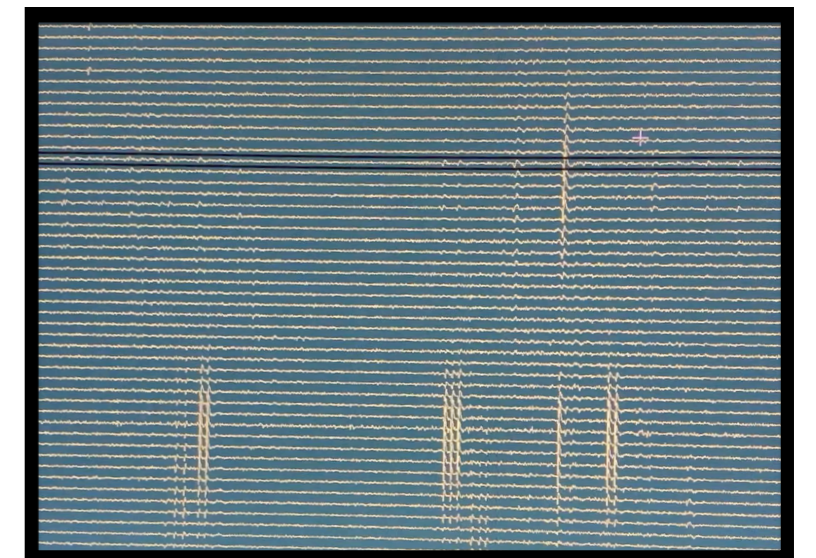
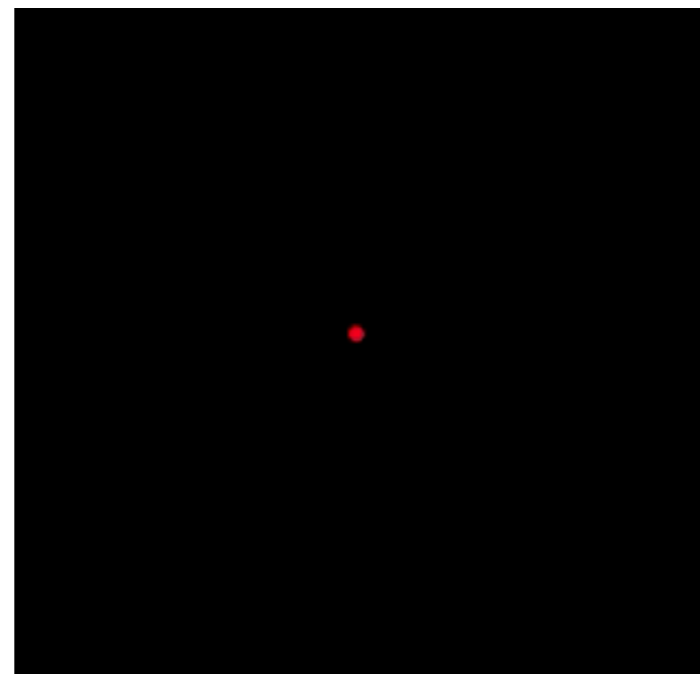
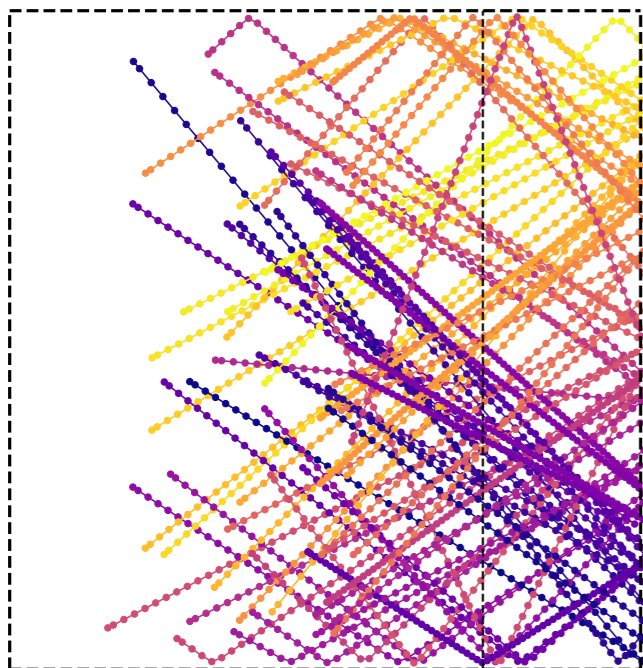


Monkey P

Monkey M



79 conditions

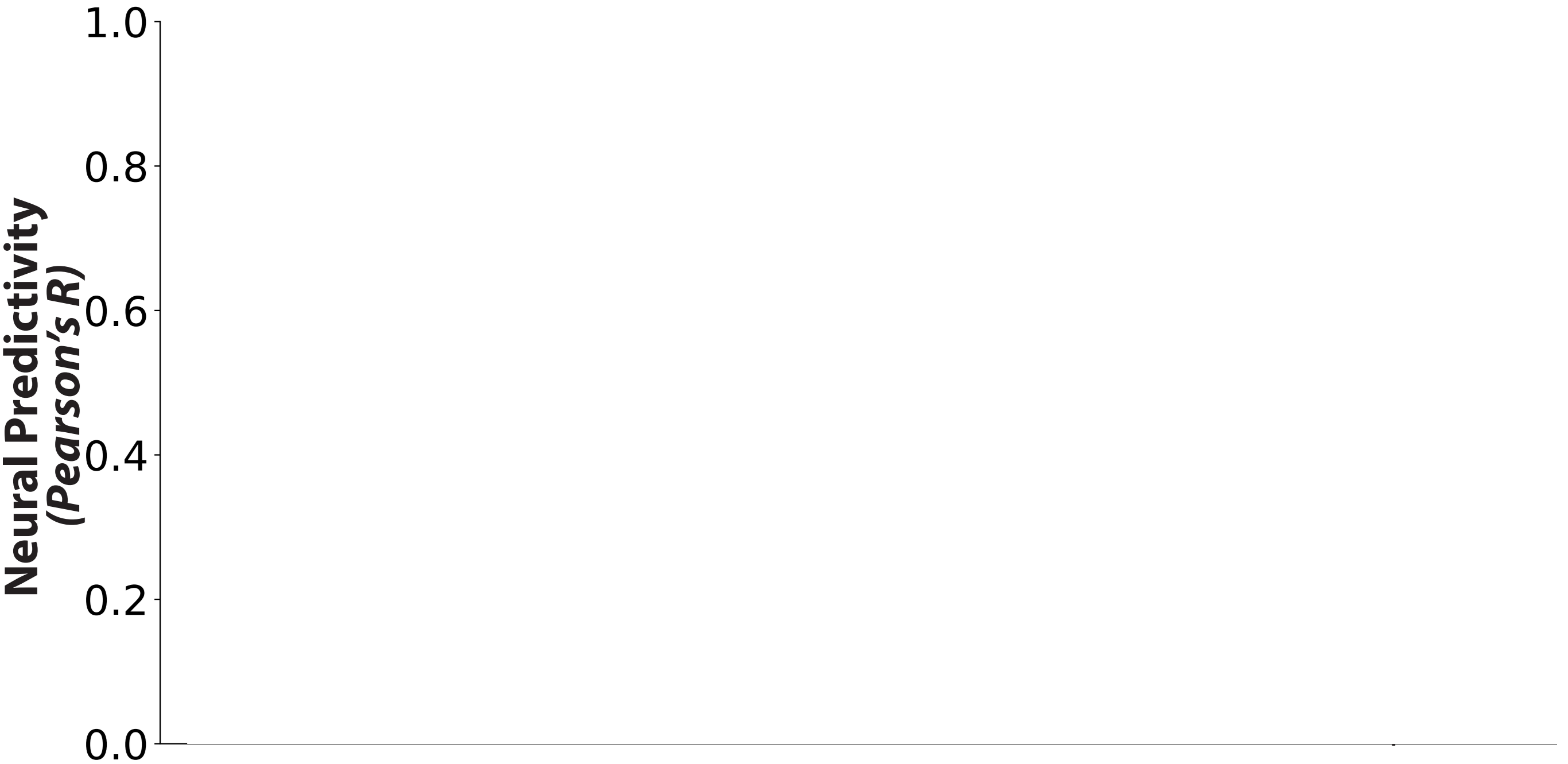
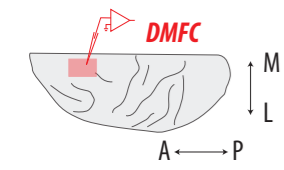


- Data from two male adult monkeys
- 79 subsampled M-Pong conditions
- 64 channel v-probe (monkey P) and 384-channel Neuropixel probe (monkey M)
- Total of 1889 stable & reliable neurons recorded from DMFC

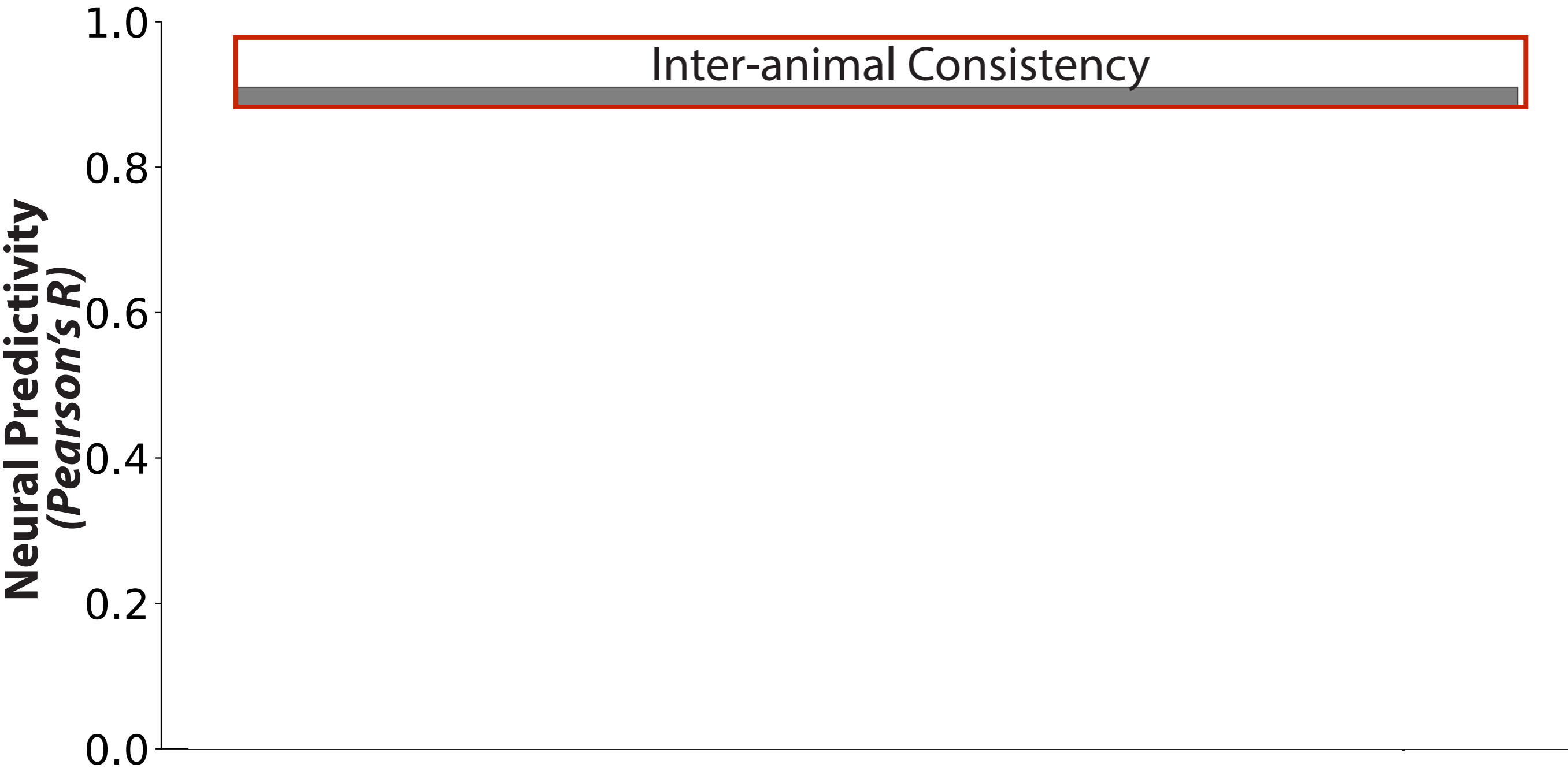
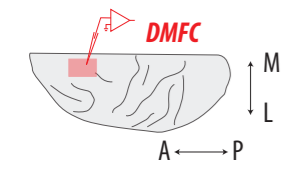


Rishi Rajalingham

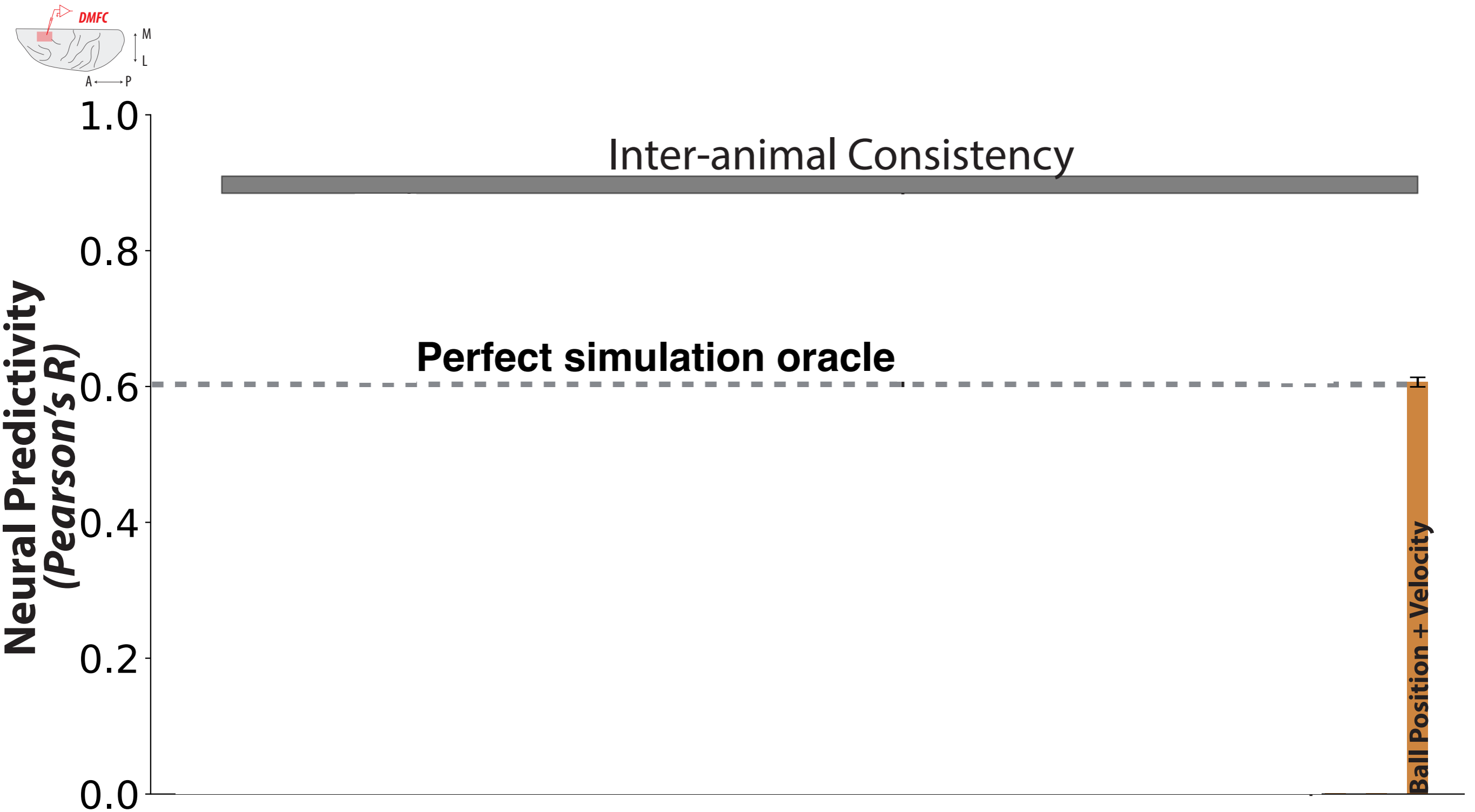
Macaque Neurophysiology: Mental Pong



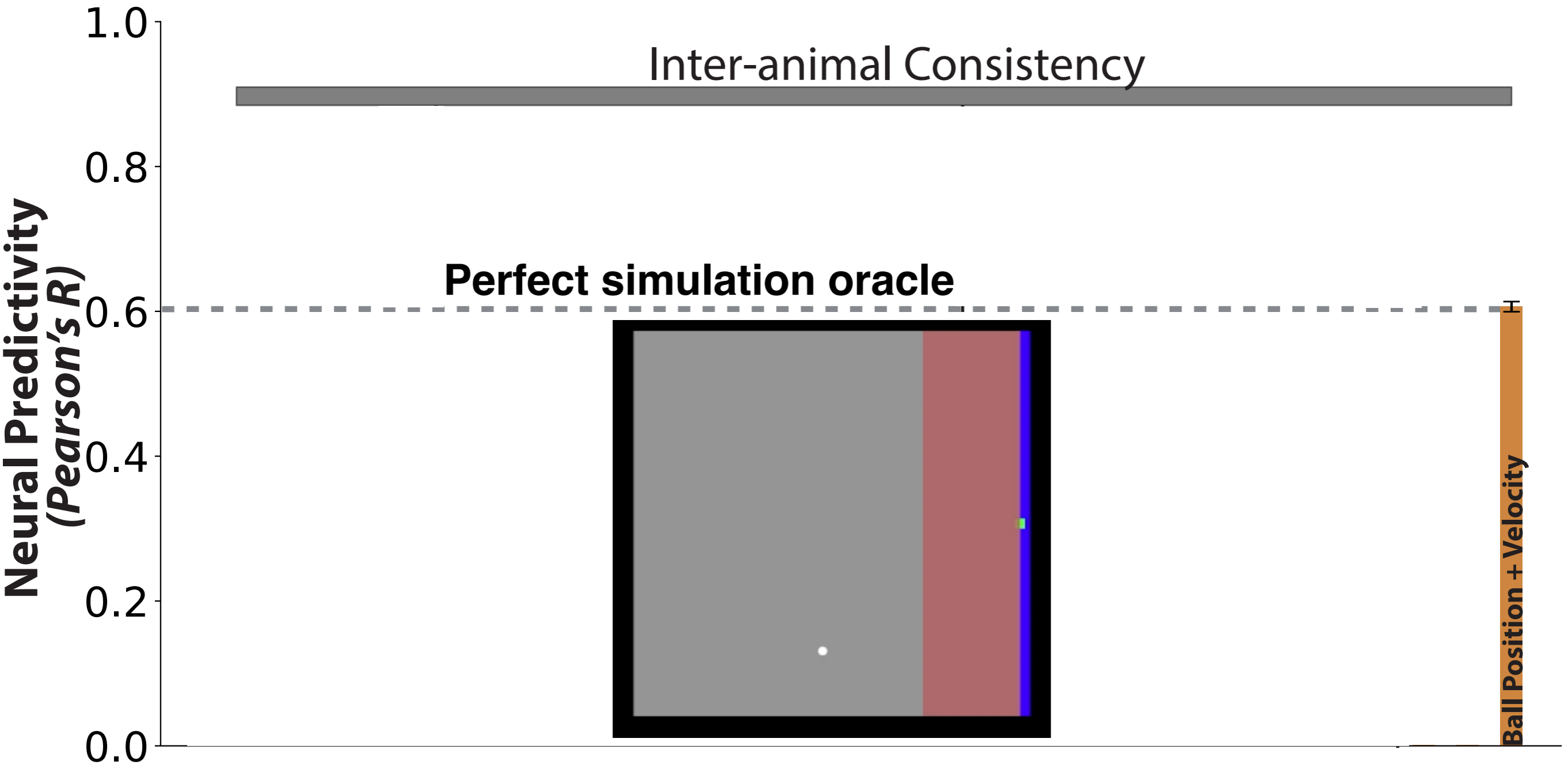
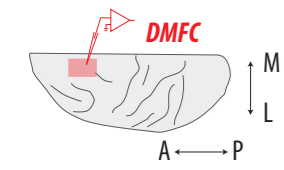
Macaque Neurophysiology: Mental Pong



Perfect Simulation Oracle Predicts Neural Data Well



Perfect Simulation Oracle Predicts Neural Data Well

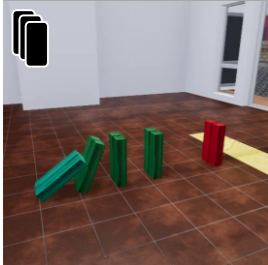


Functional Constraint Hypotheses

Inputs

Physion

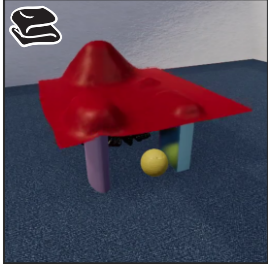
Dominoes



Support



Drape



Link



Sensory-Cognitive Hypothesis Classes

Hypothesis Class I: Pixel-wise Future Prediction

Inputs

Physion

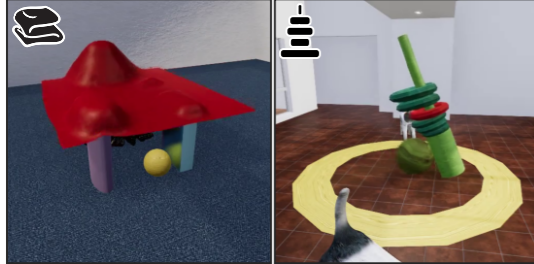
Dominoes

Support



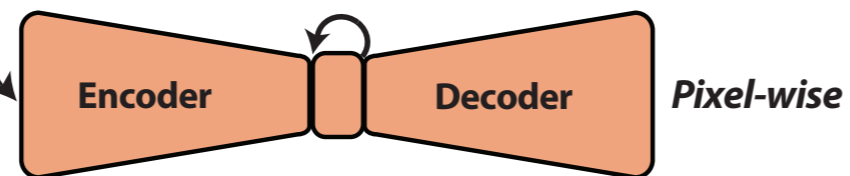
Drape

Link

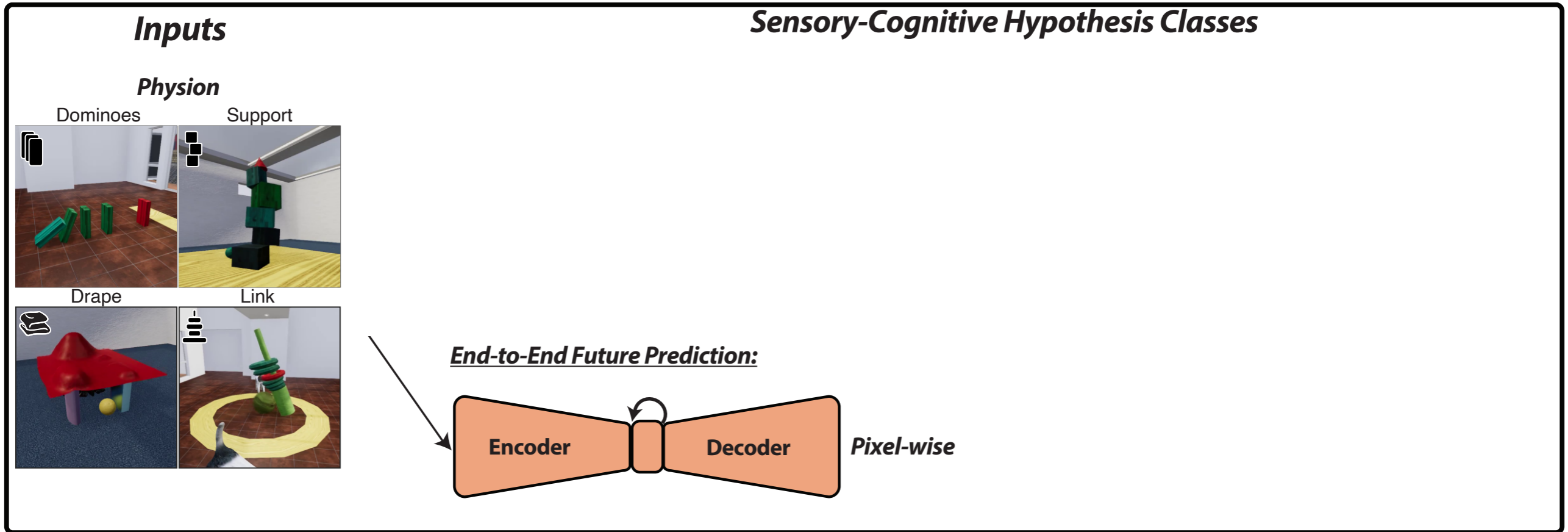


Sensory-Cognitive Hypothesis Classes

End-to-End Future Prediction:

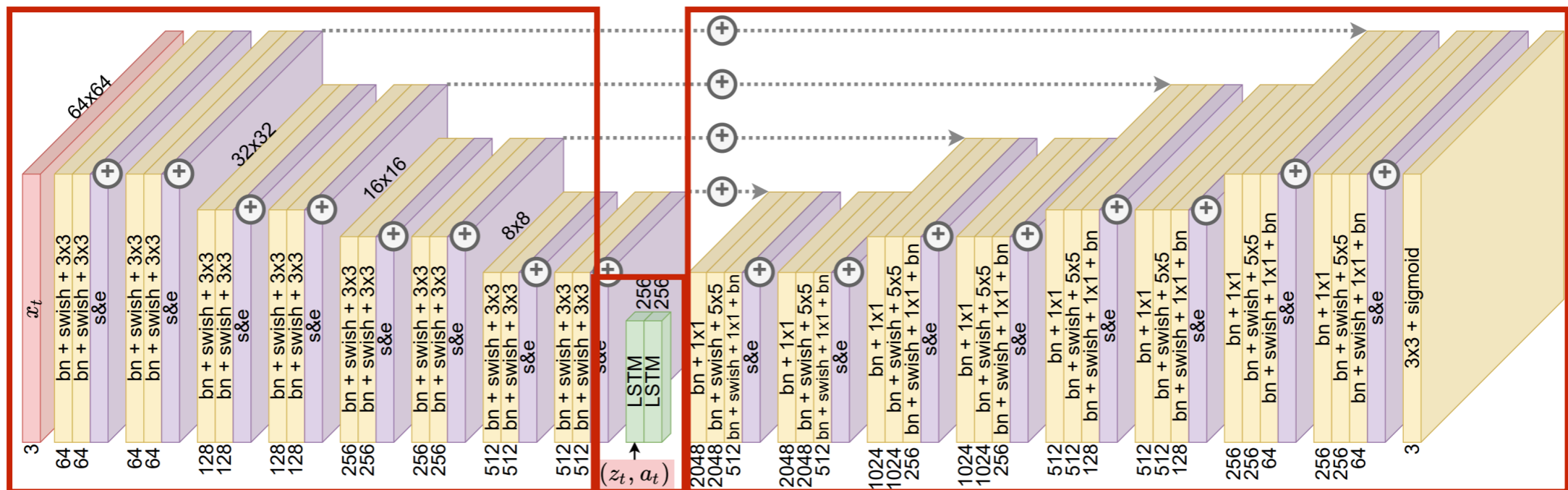


Hypothesis Class I: Pixel-wise Future Prediction



Babaeizadeh et al. 2021

Input Conv Squeeze and Excite LSTM Skip Connection Residual

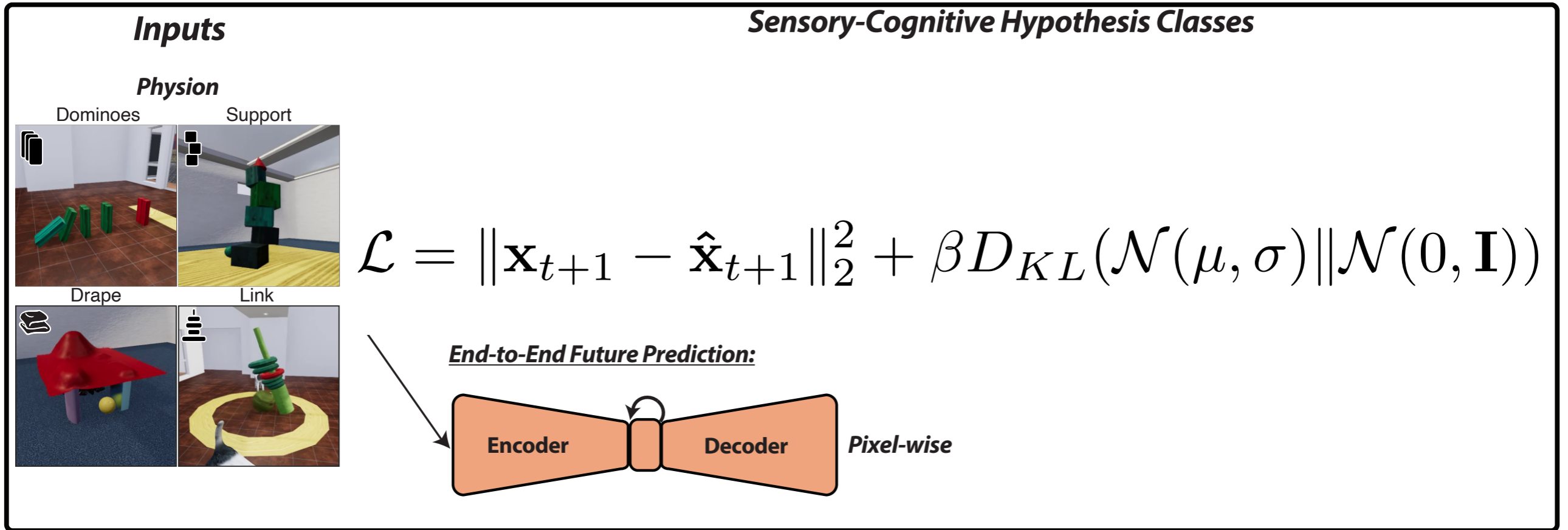


Visual Encoder
("Sensory")

Dynamics Predictor
("Cognitive")

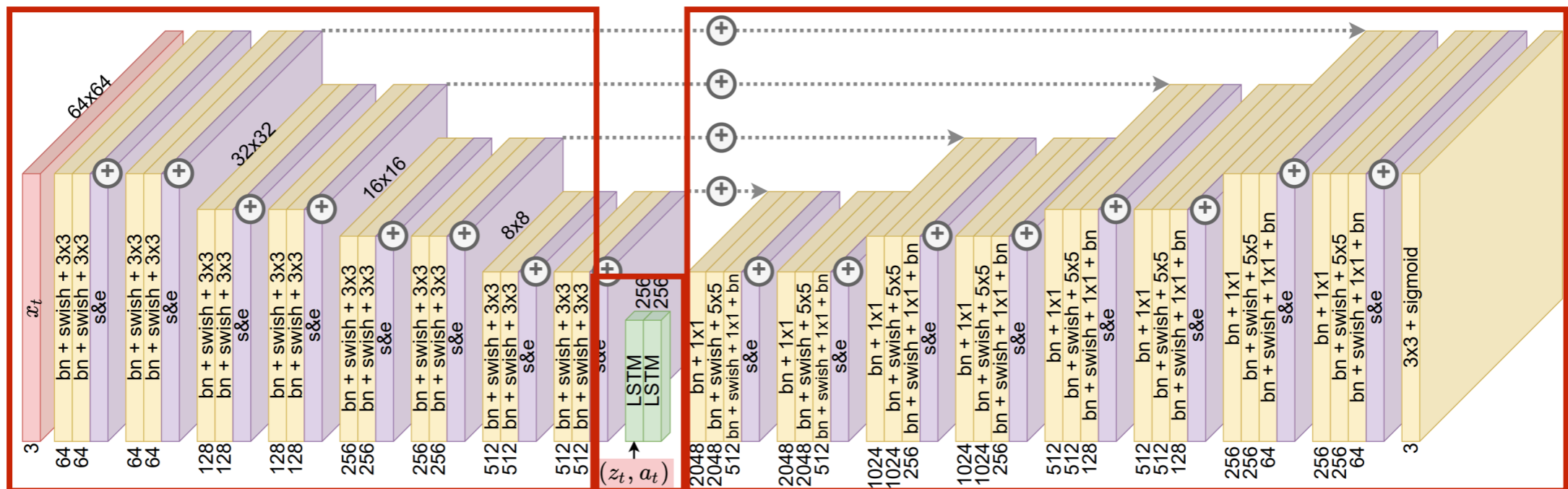
Frame Decoder
("Objective/Behavior")

Hypothesis Class I: Pixel-wise Future Prediction



Babaeizadeh et al. 2021

Input Conv Squeeze and Excite LSTM Skip Connection Residual

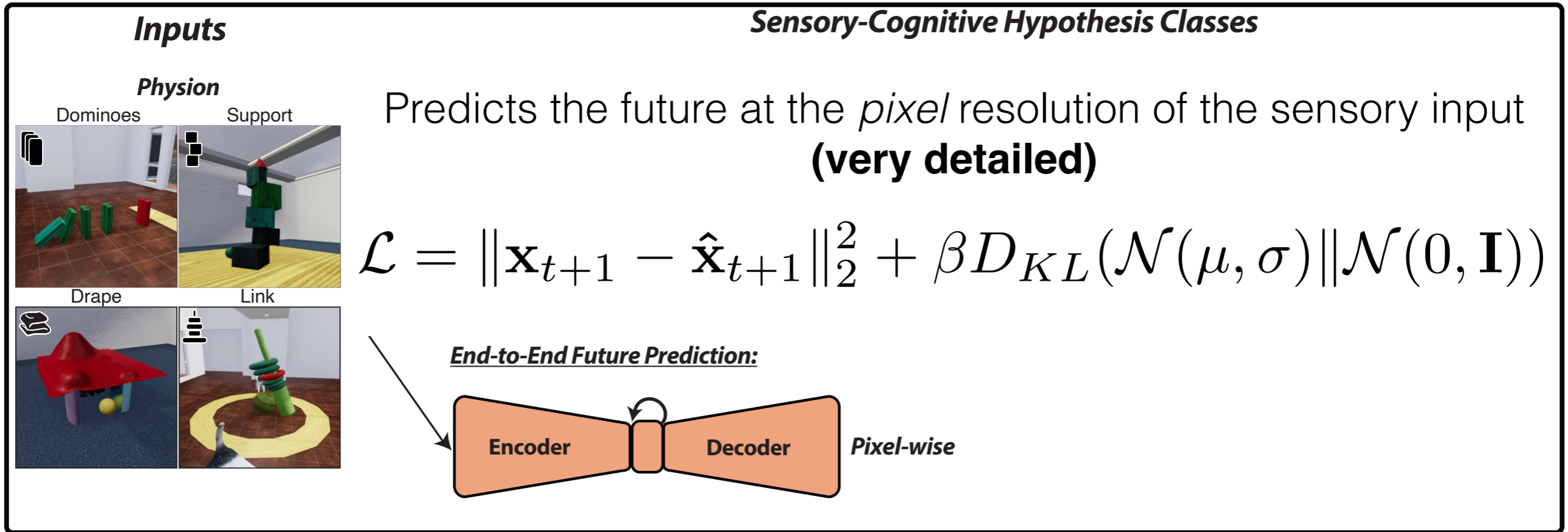


Visual Encoder
("Sensory")

Dynamics Predictor
("Cognitive")

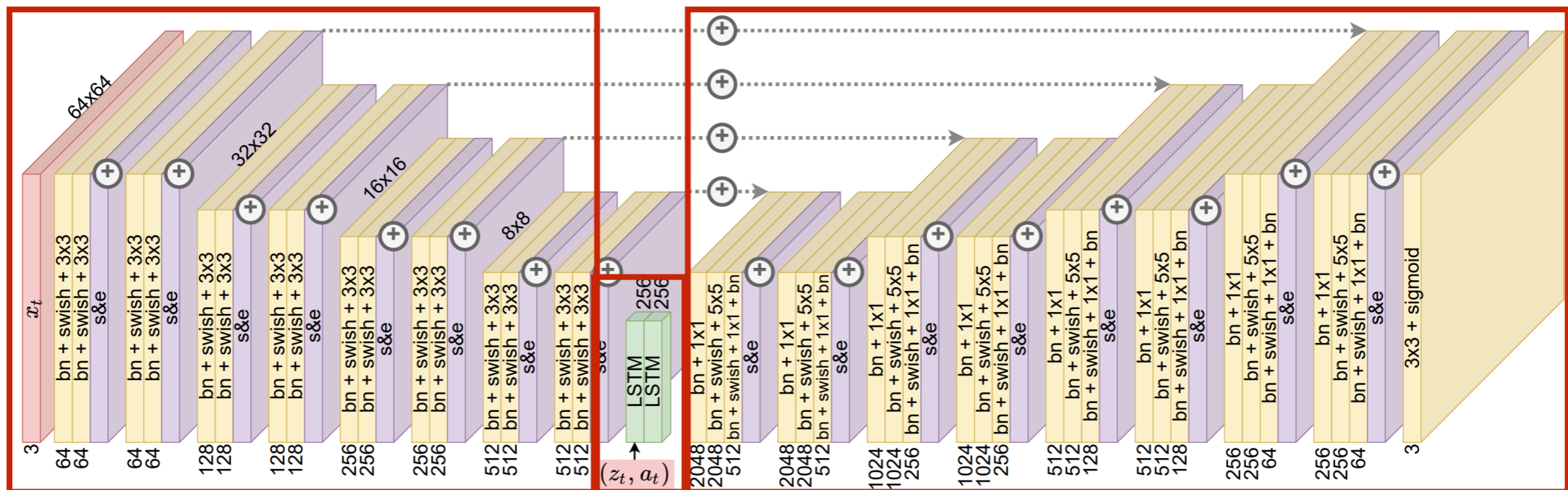
Frame Decoder
("Objective/Behavior")

Hypothesis Class I: Pixel-wise Future Prediction



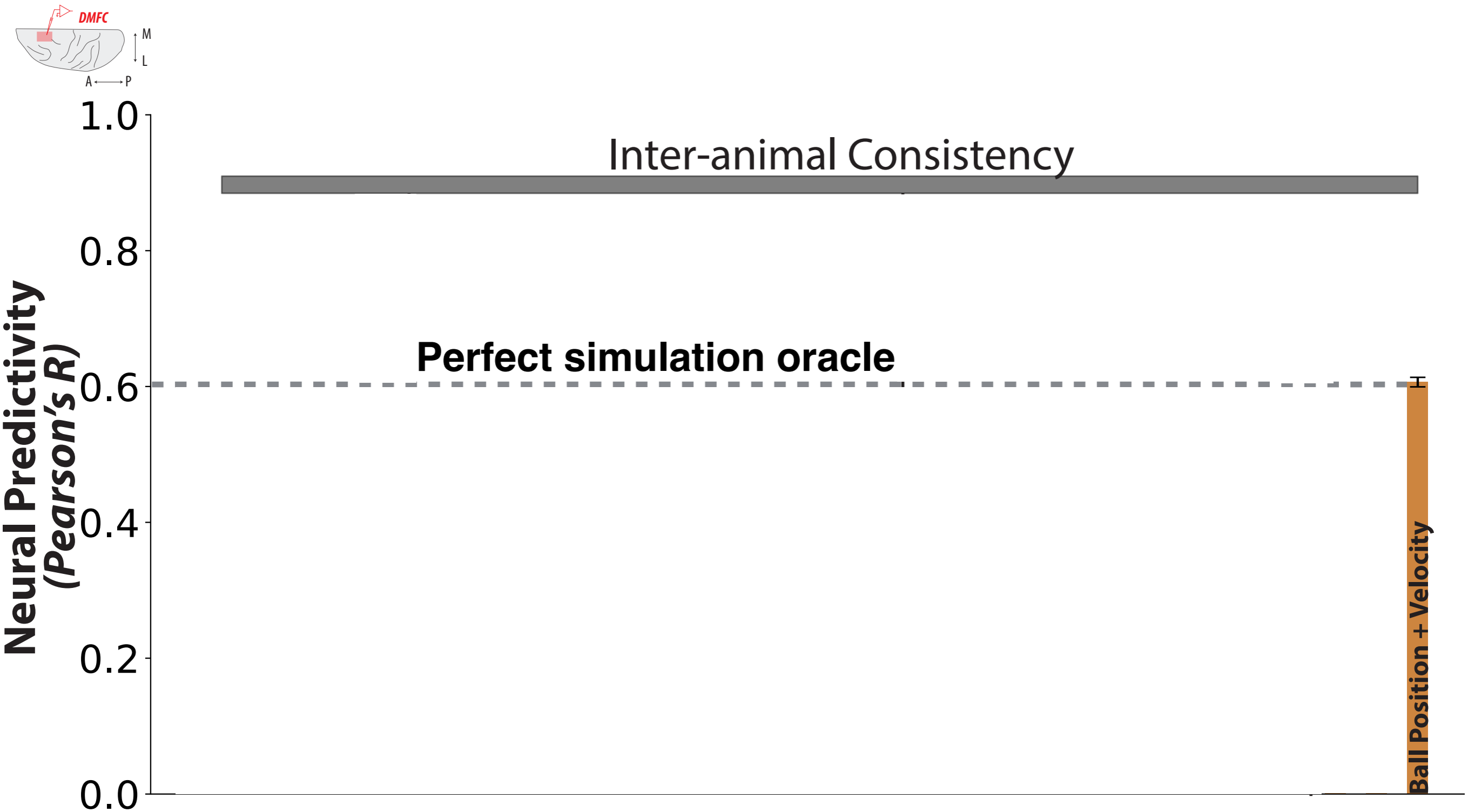
Babaeizadeh et al. 2021

Input Conv Squeeze and Excite LSTM Skip Connection Residual

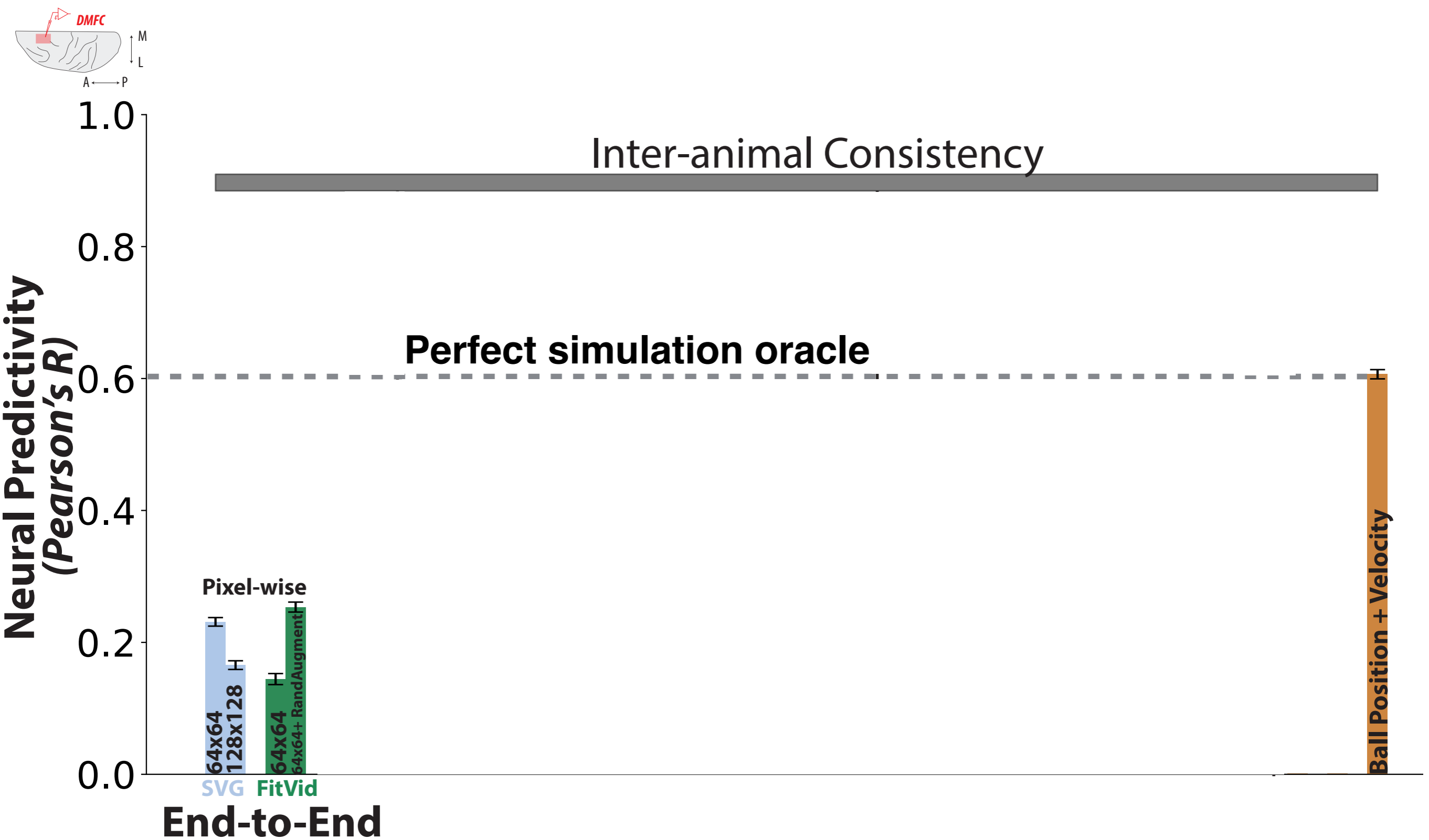


Visual Encoder ("Sensory") Dynamics Predictor ("Cognitive") Frame Decoder ("Objective/Behavior")

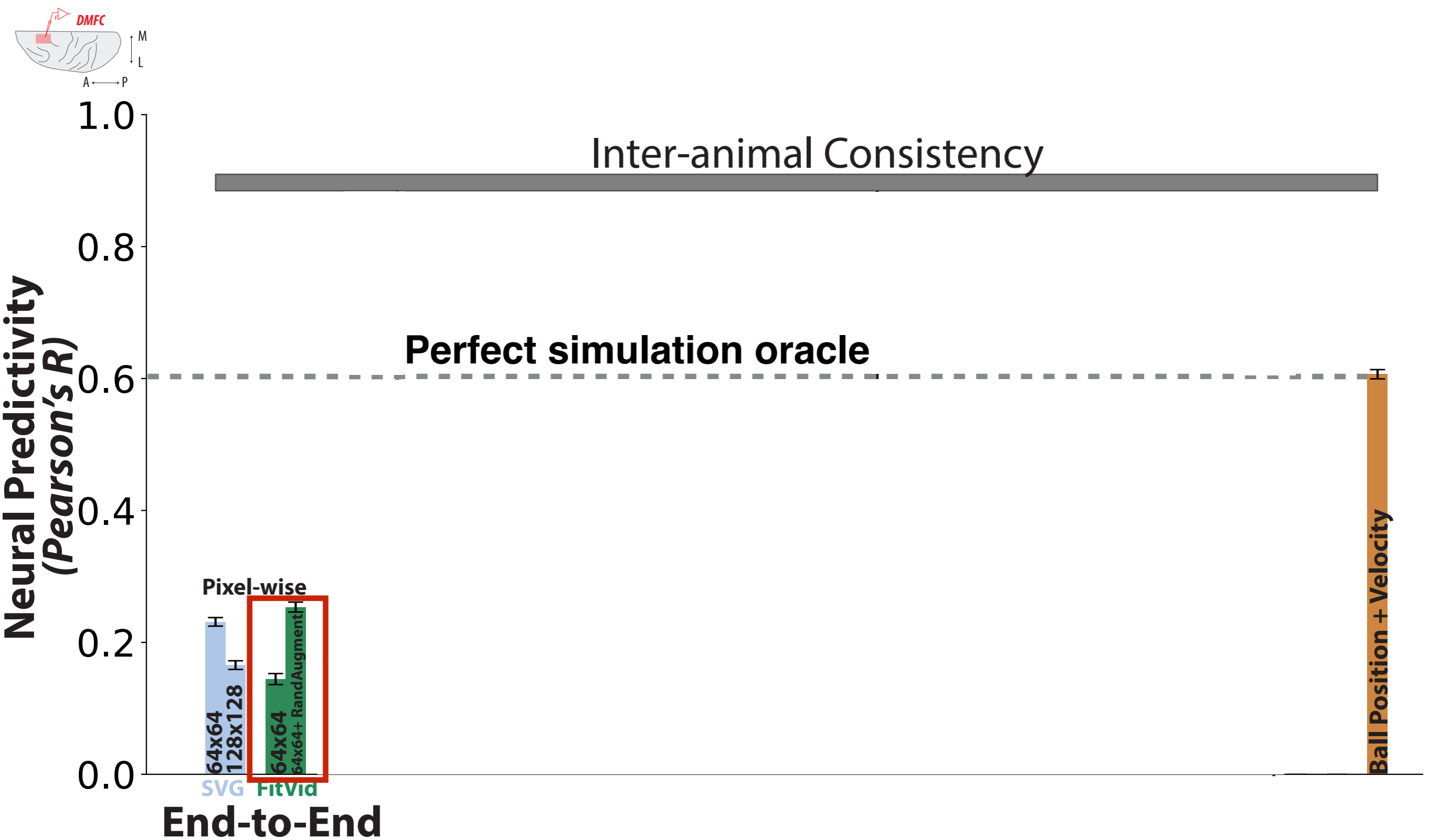
Physical Simulation Oracles Predict Neural Data Well



Pixel-wise Future Prediction Poorly Predicts Neurons



Pixel-wise Future Prediction Poorly Predicts Neurons

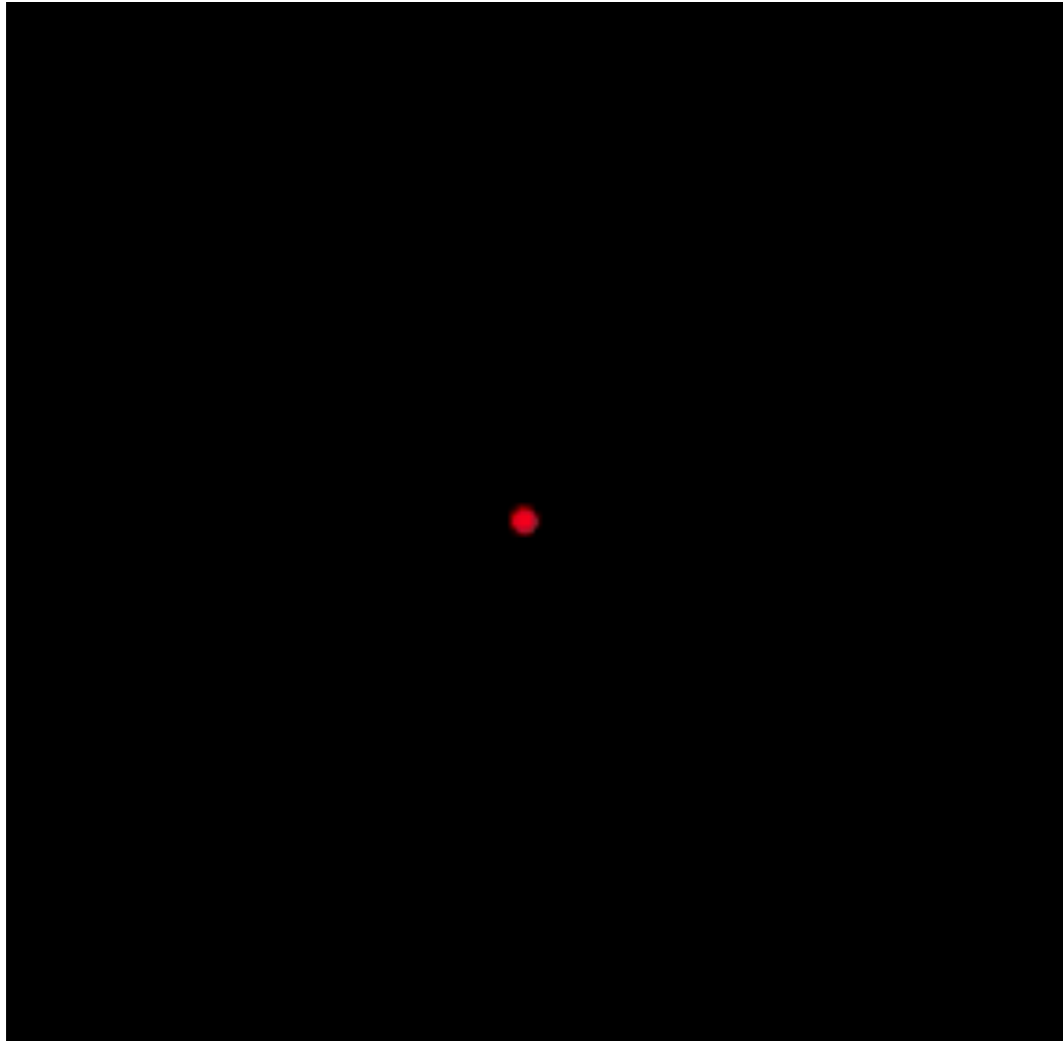


...and they struggle to generalize to Pong

Input Frames

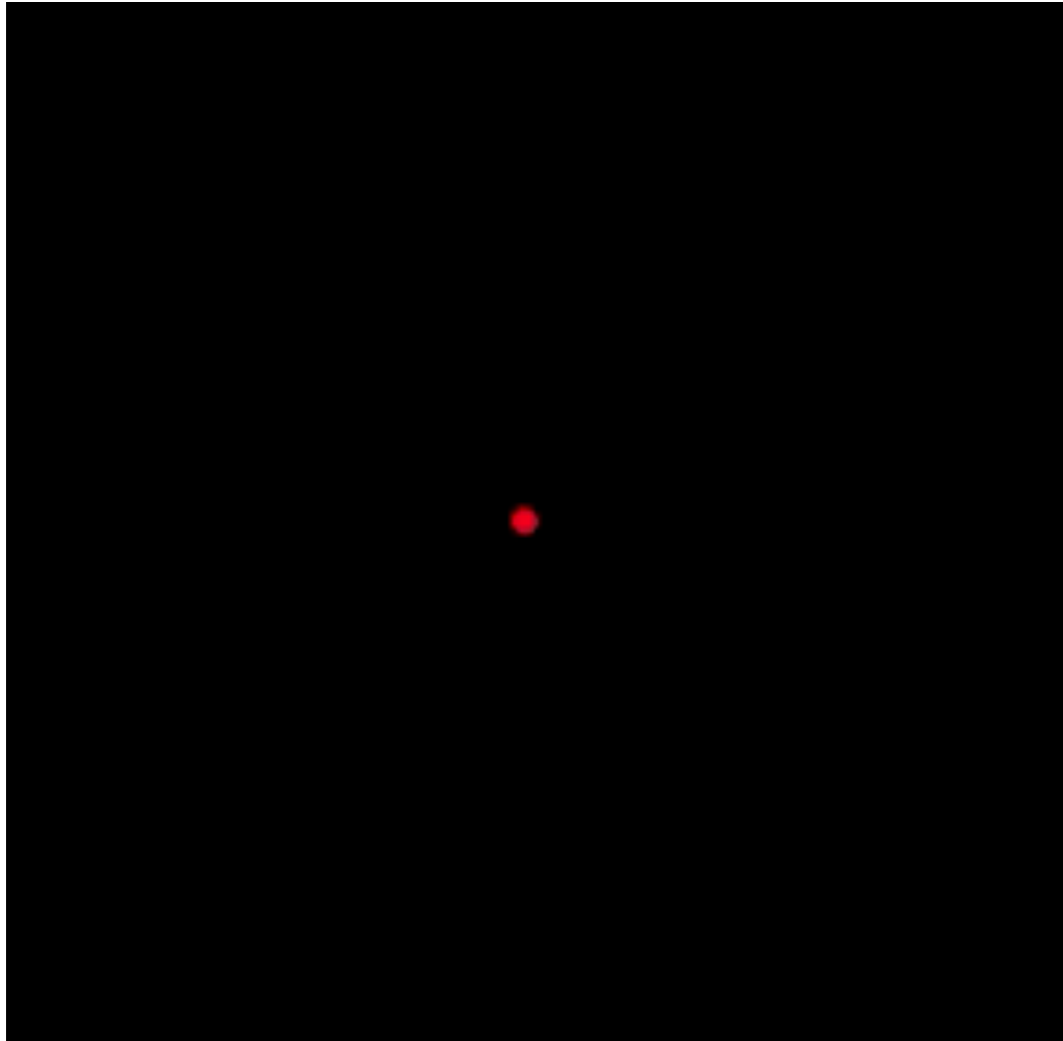
...and they struggle to generalize to Pong

Input Frames



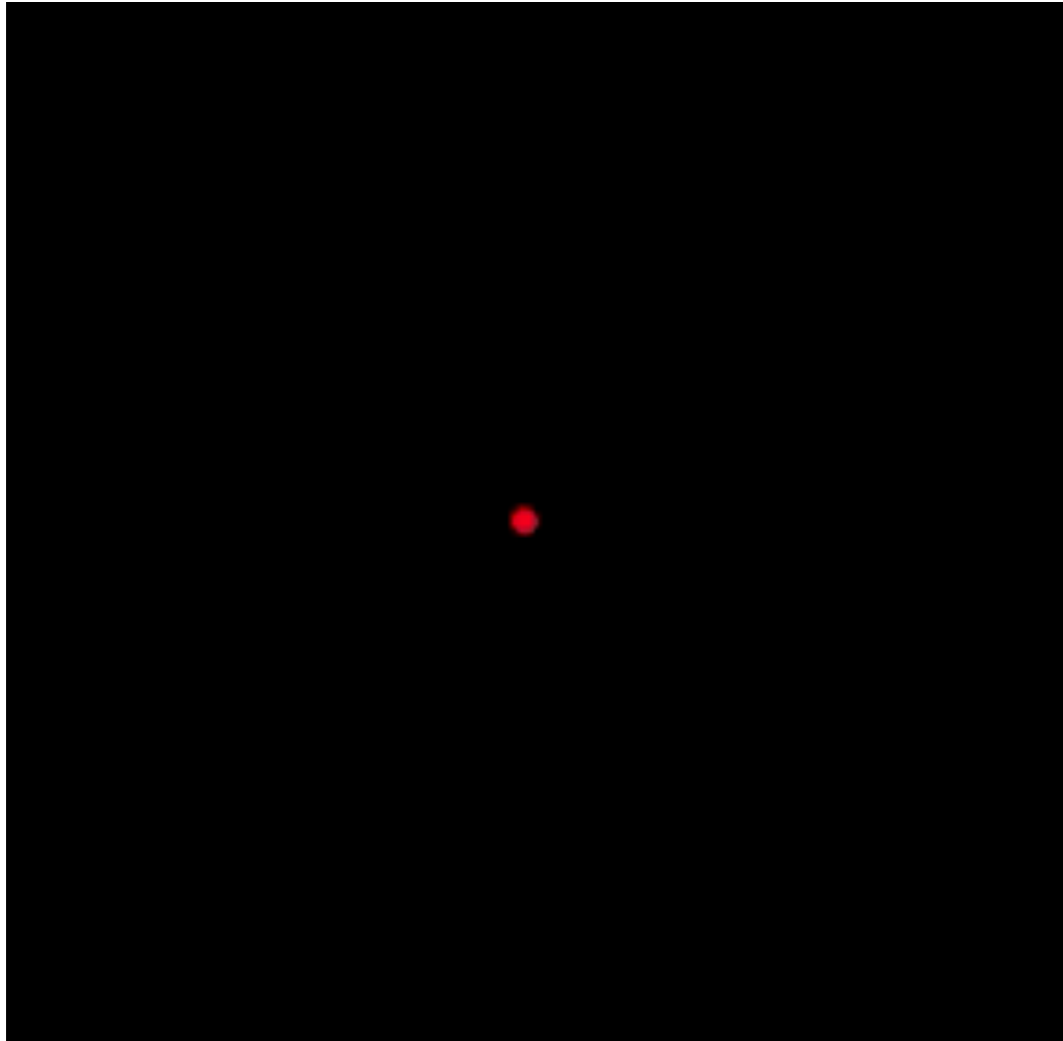
...and they struggle to generalize to Pong

Input Frames



...and they struggle to generalize to Pong

Input Frames

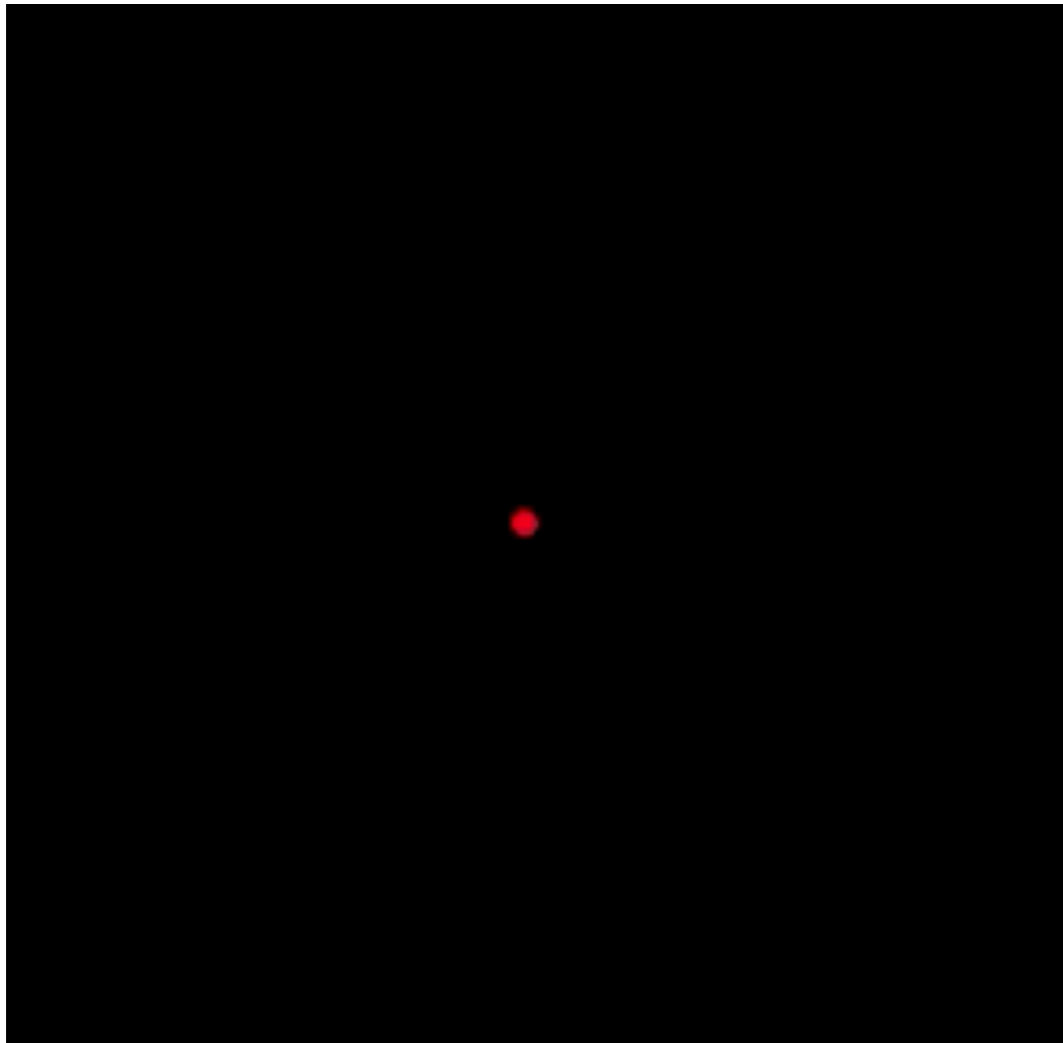


Predicted Frames

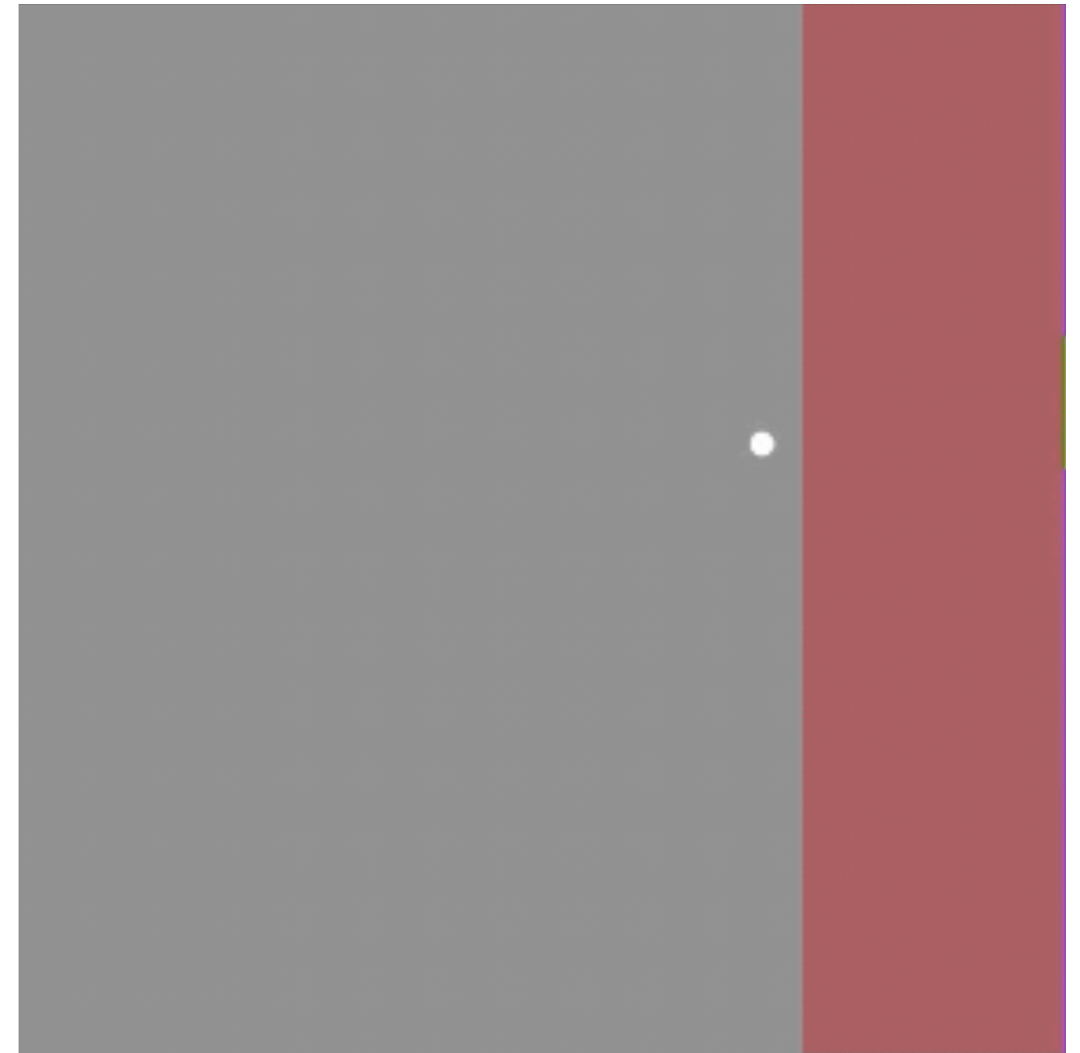


...and they struggle to generalize to Pong

Input Frames

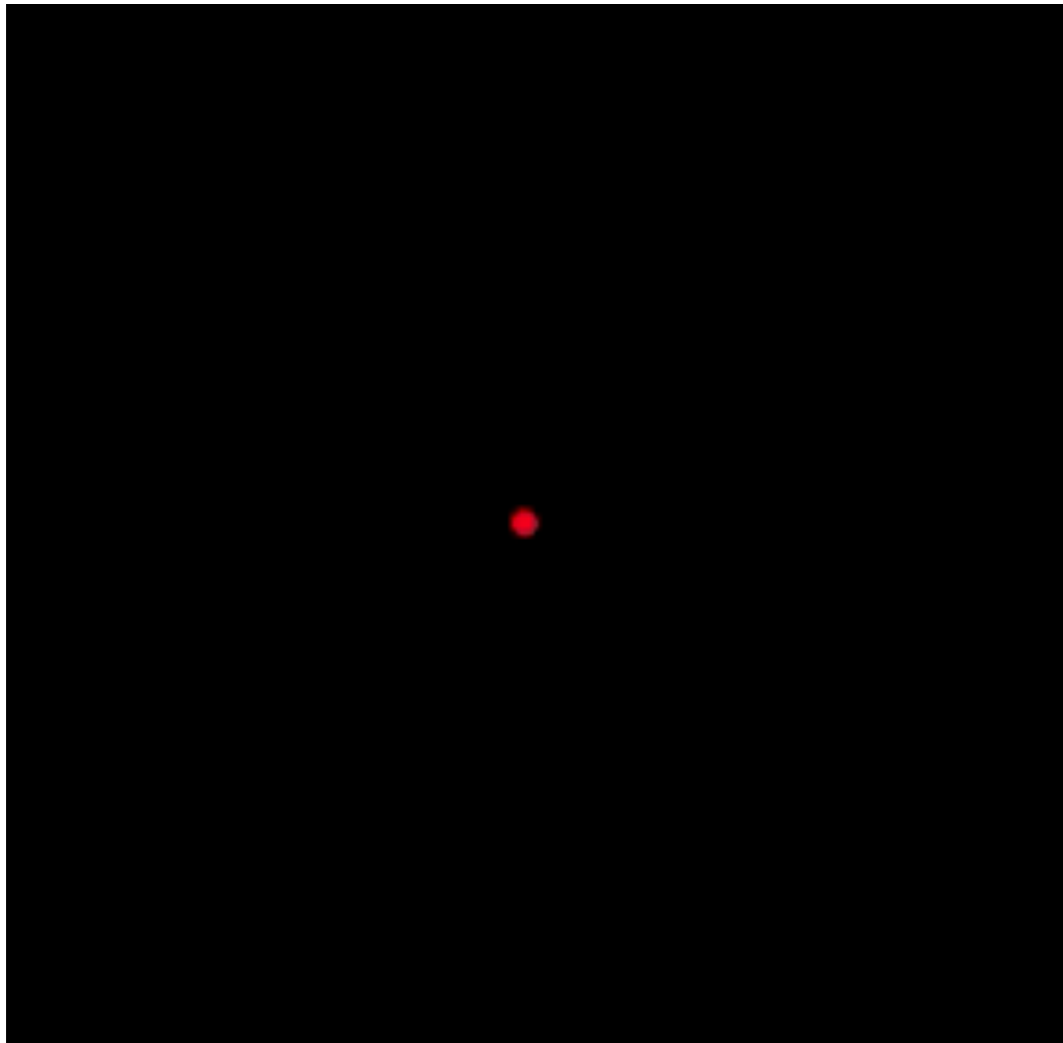


Predicted Frames

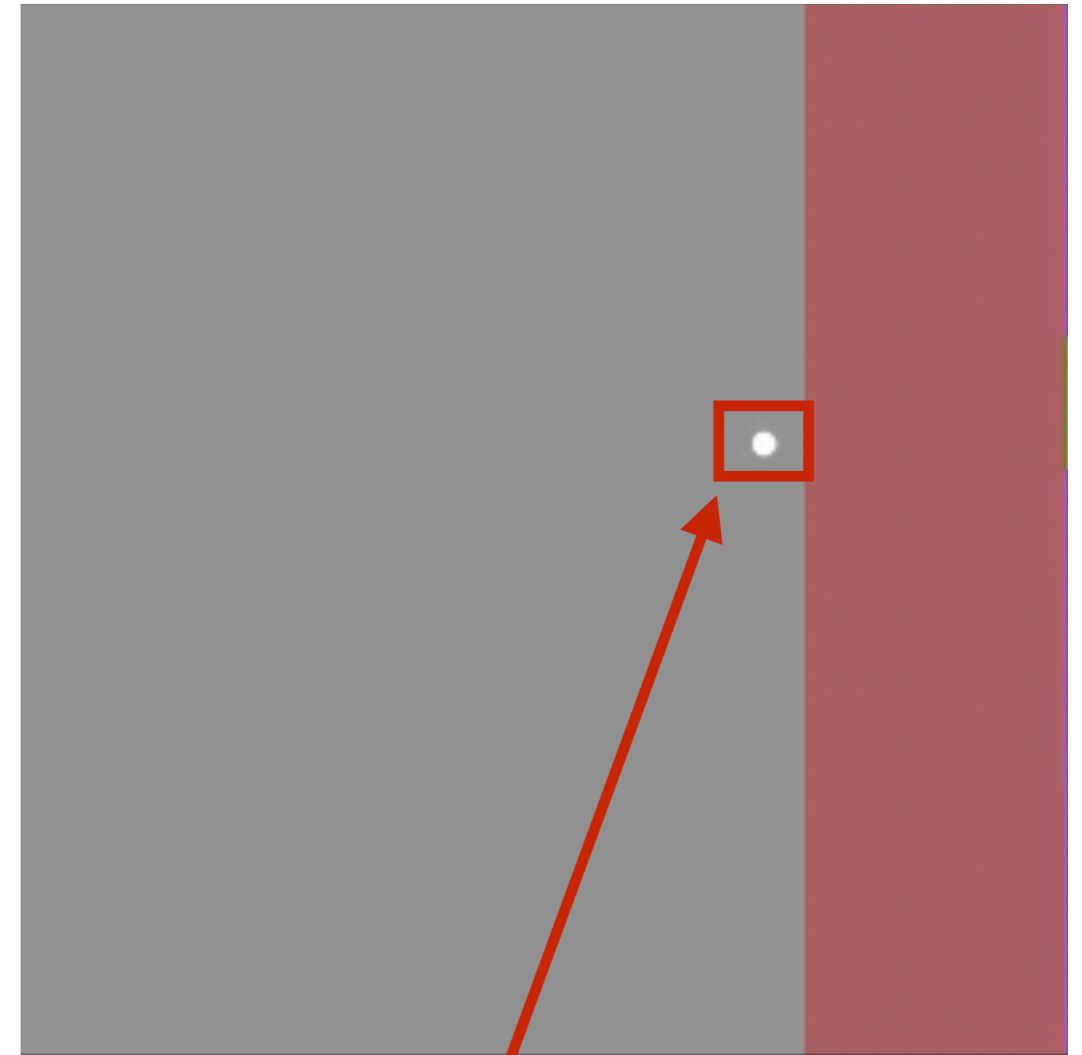


...and they struggle to generalize to Pong

Input Frames

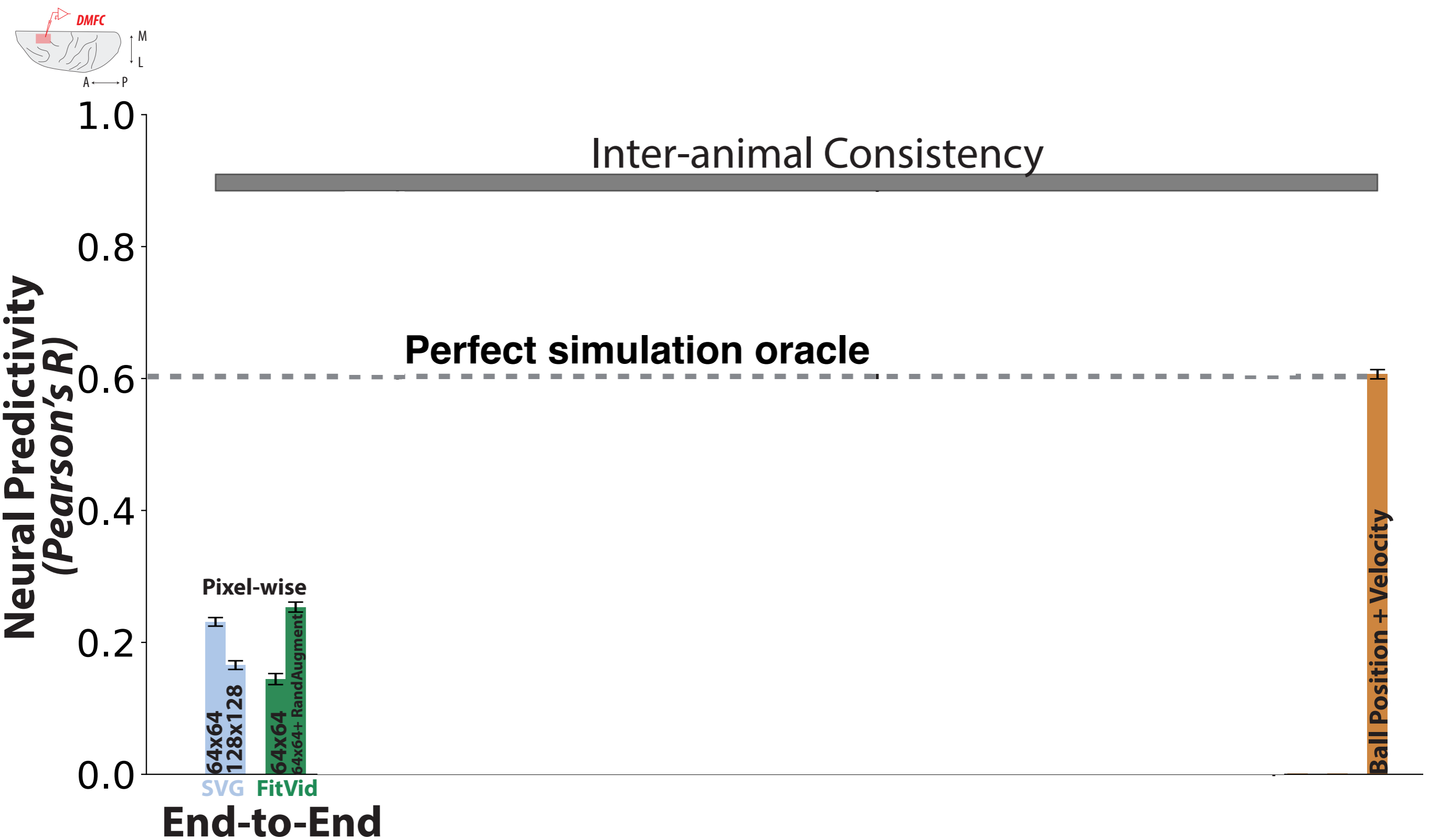


Predicted Frames



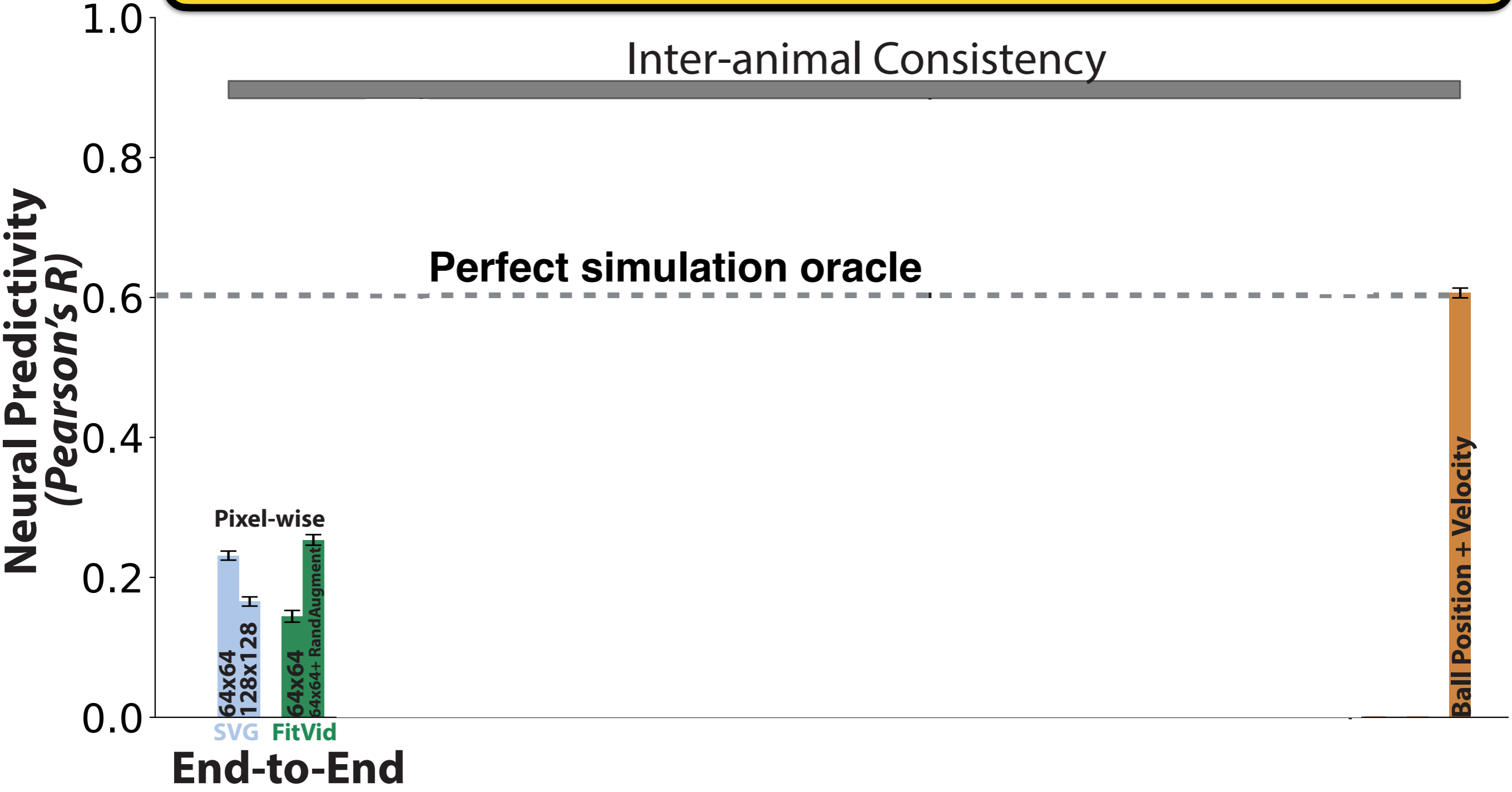
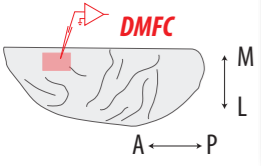
Ball stops at final input frame, in the model's "imagination"

Pixel-wise Future Prediction Poorly Predicts Neurons

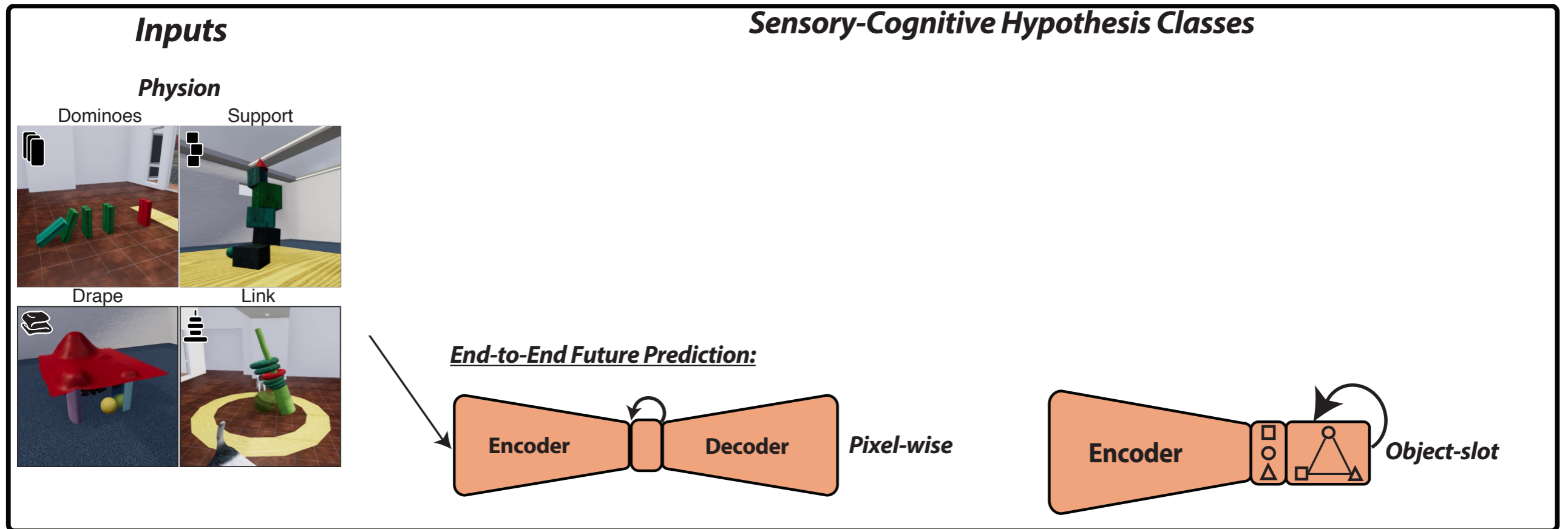


Pixel-wise Future Prediction Poorly Predicts Neurons

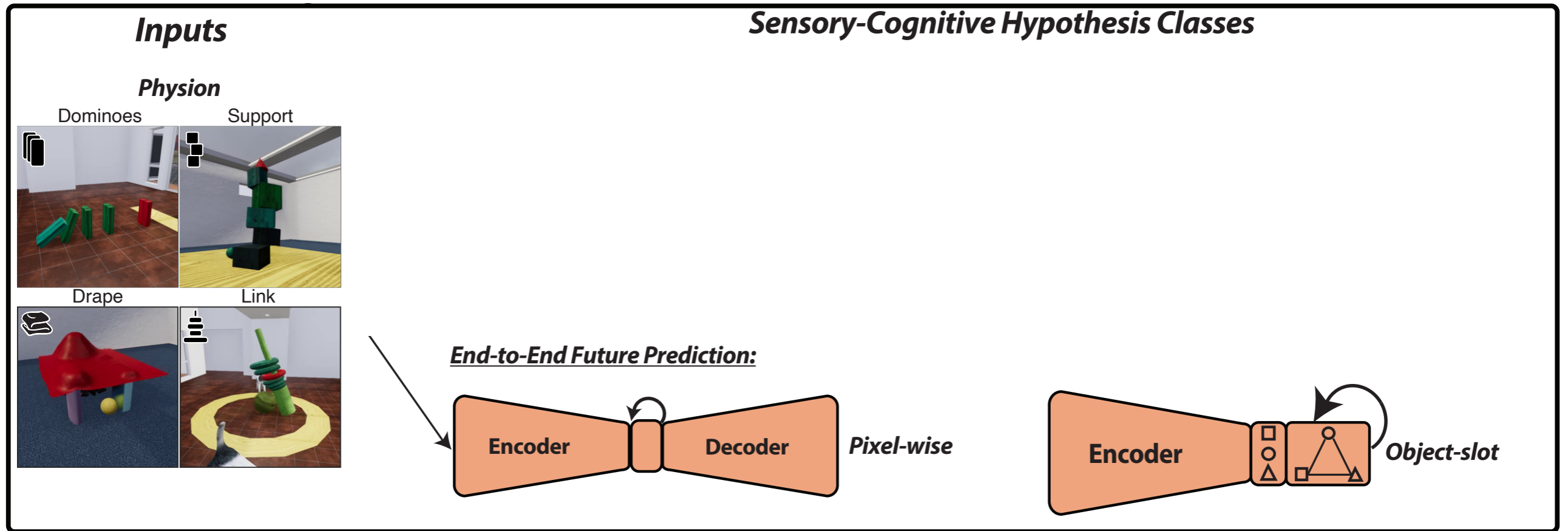
Perhaps DMFC predicts a “factorized” version of the scene?
How?



Hypothesis Class 2: Object Slots

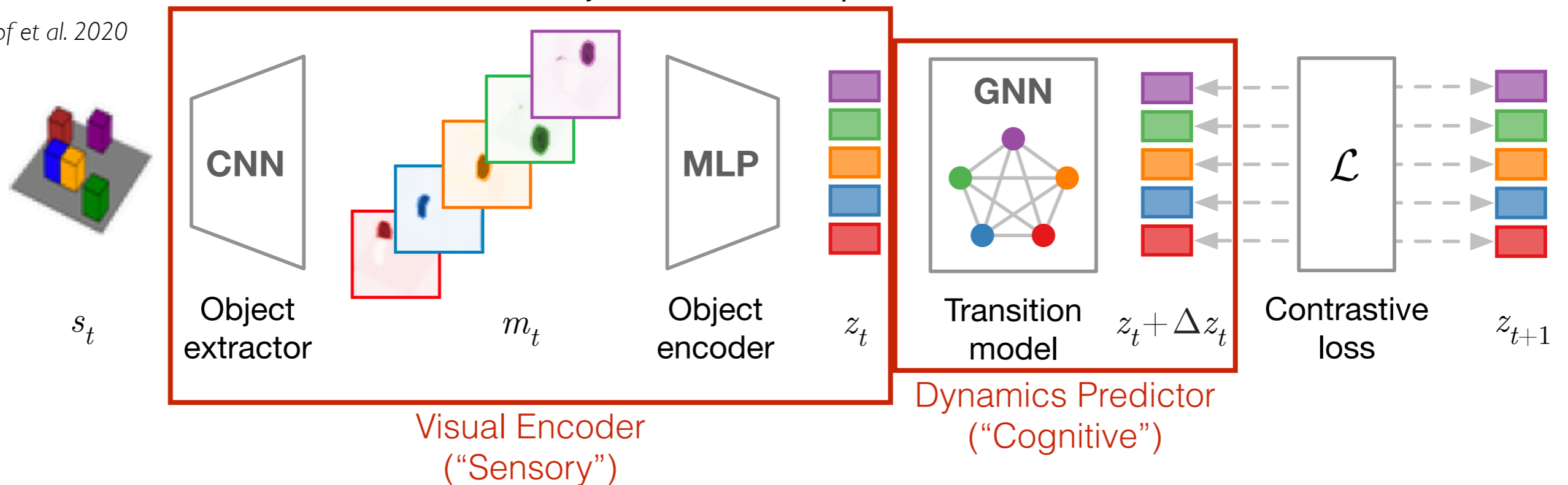


Hypothesis Class 2: Object Slots

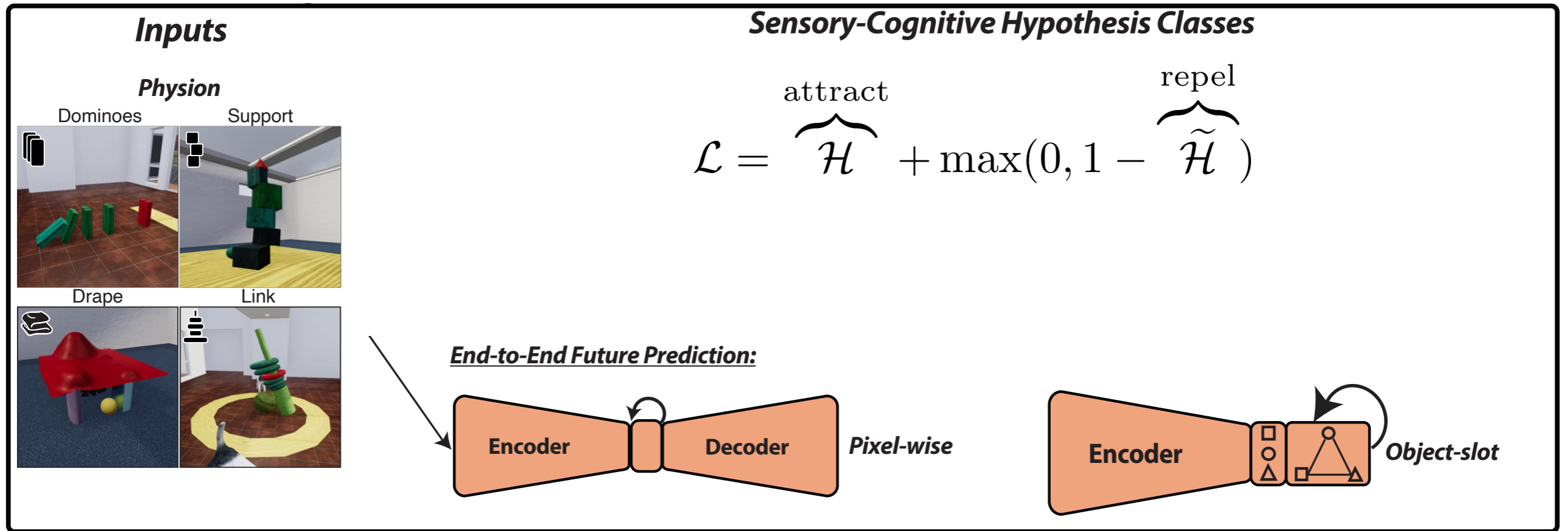


Predicts at the level of object slot representations and their relations

Kipf et al. 2020

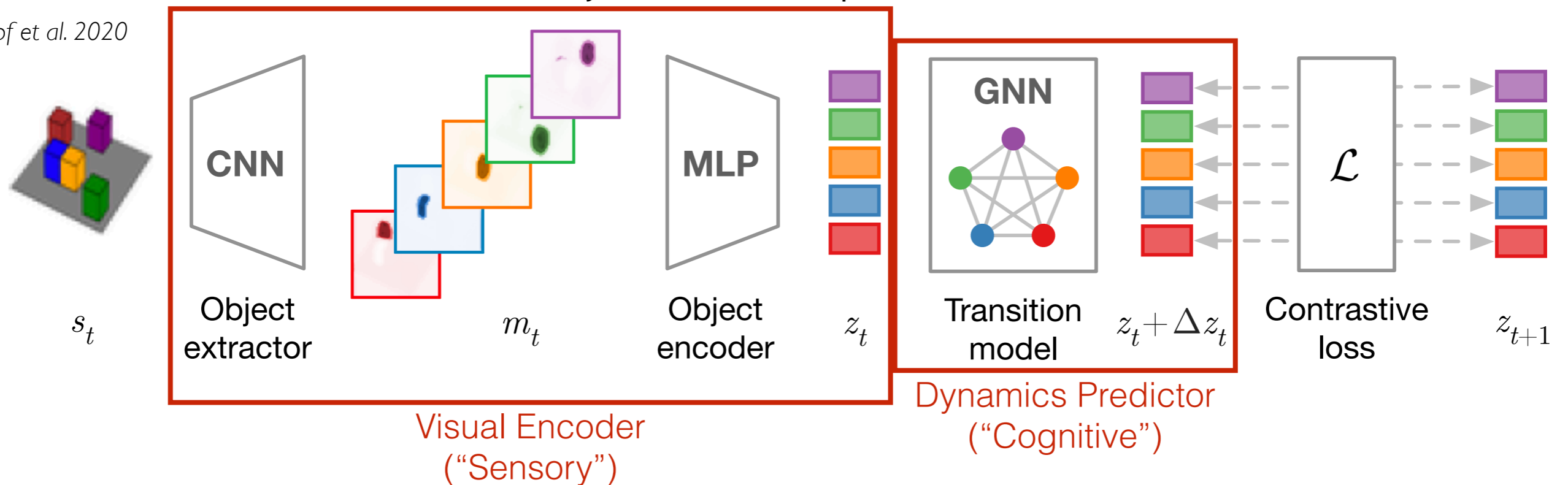


Hypothesis Class 2: Object Slots

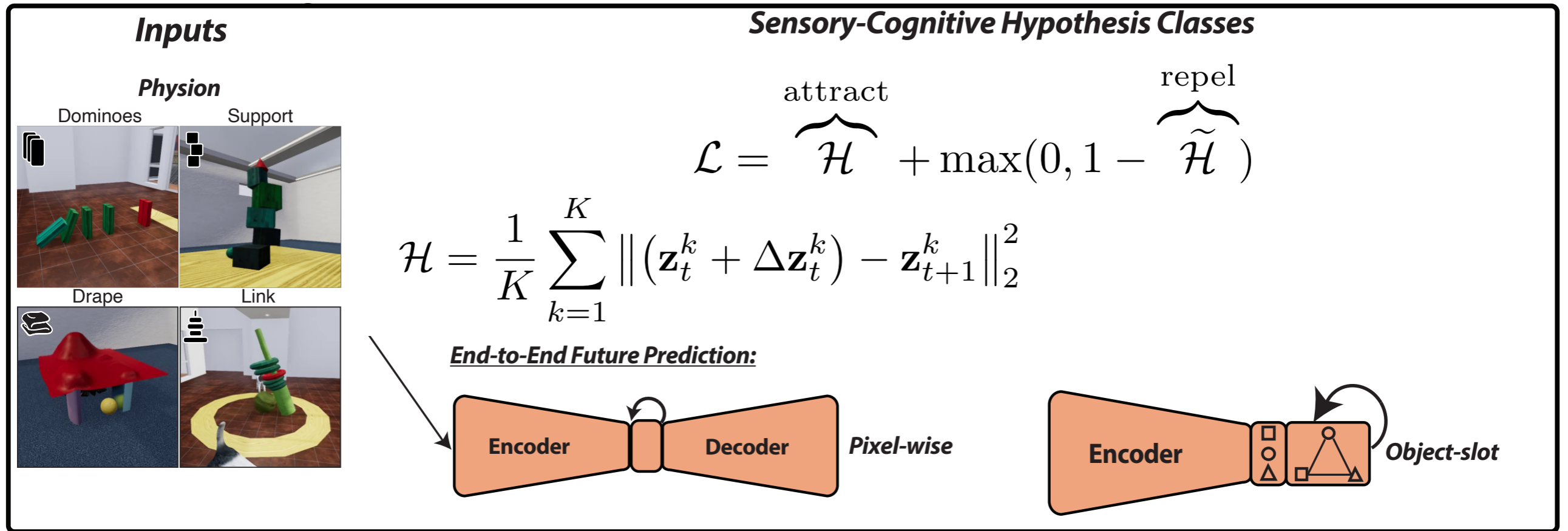


Predicts at the level of object slot representations and their relations

Kipf et al. 2020

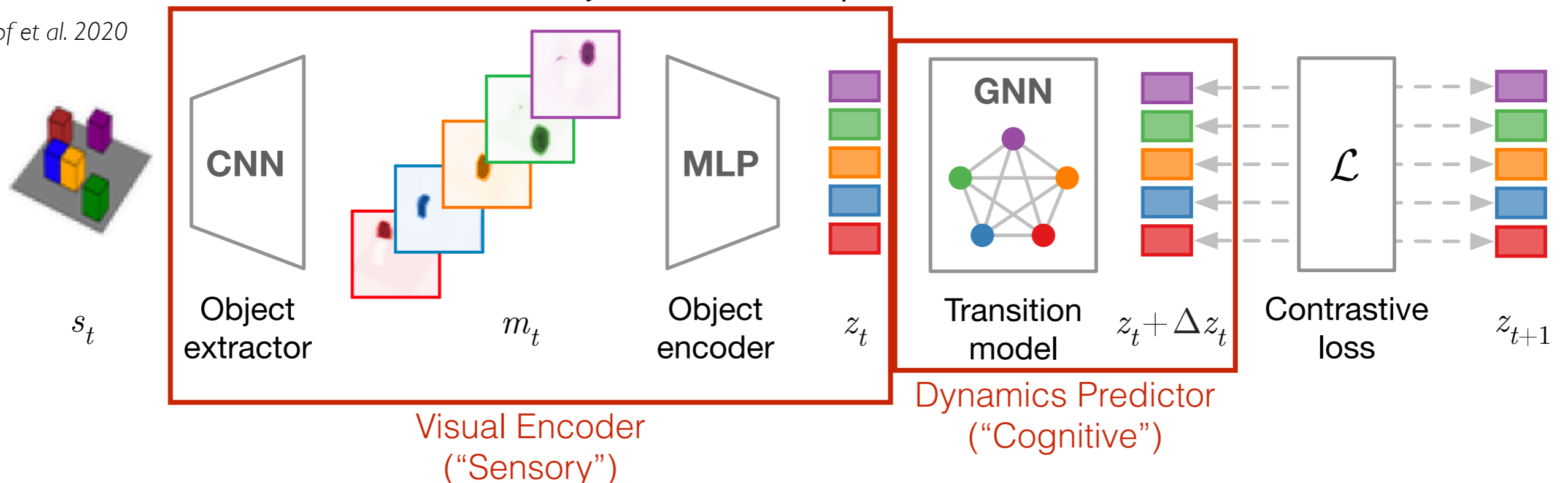


Hypothesis Class 2: Object Slots

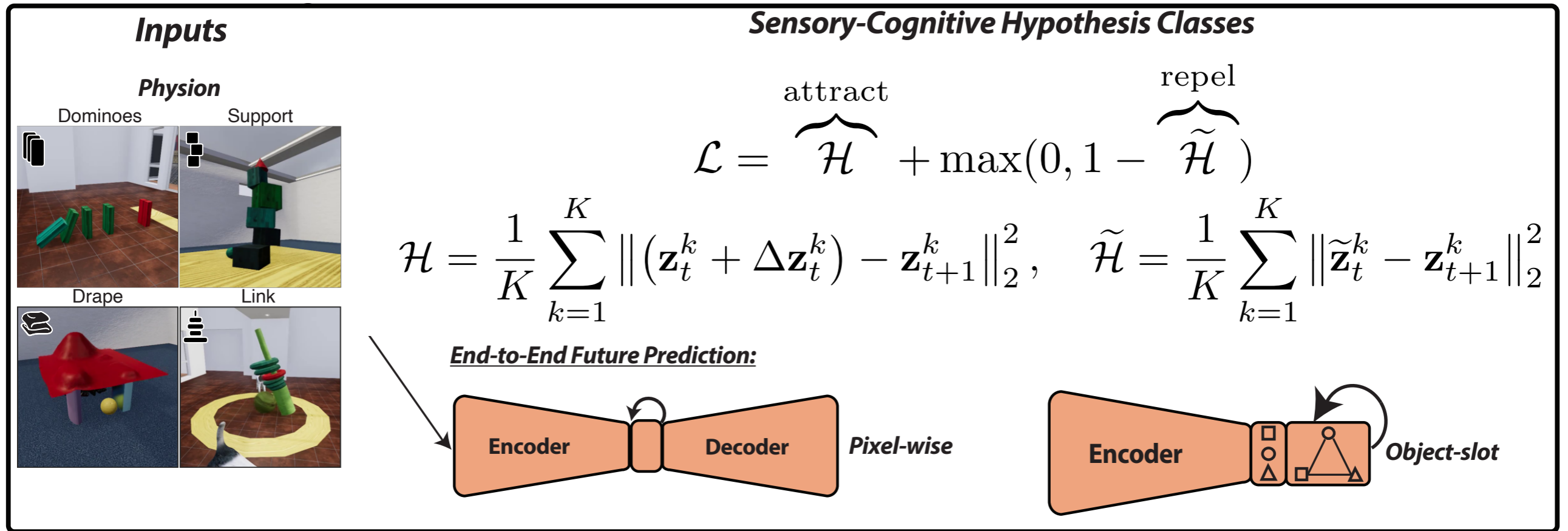


Predicts at the level of object slot representations and their relations

Kipf et al. 2020

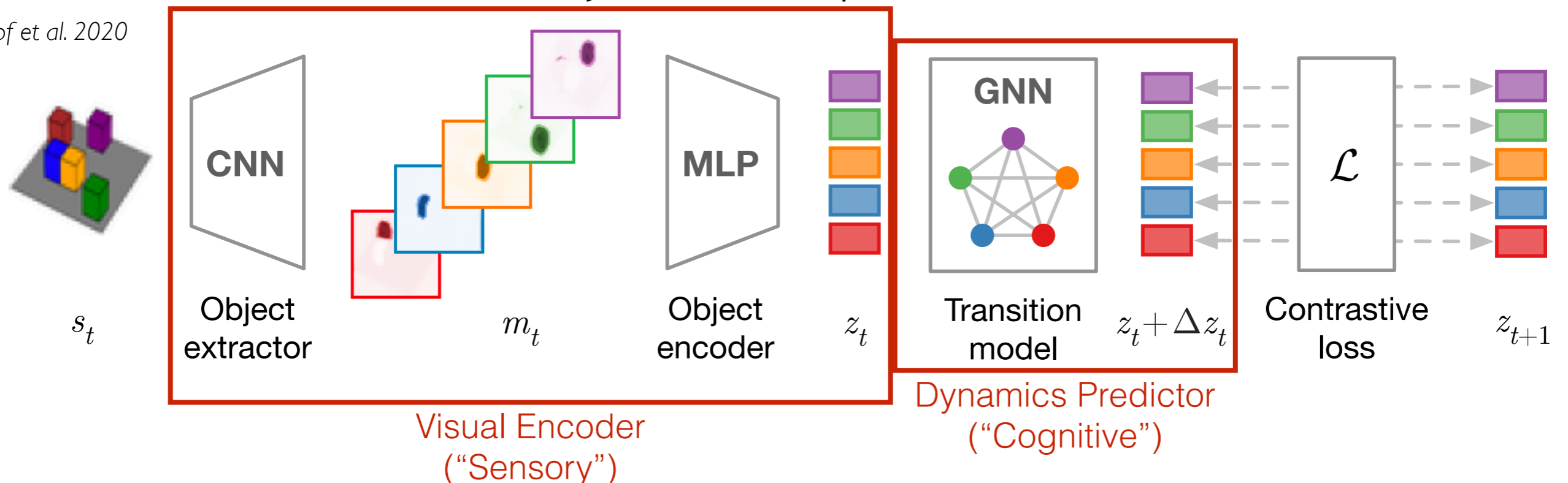


Hypothesis Class 2: Object Slots

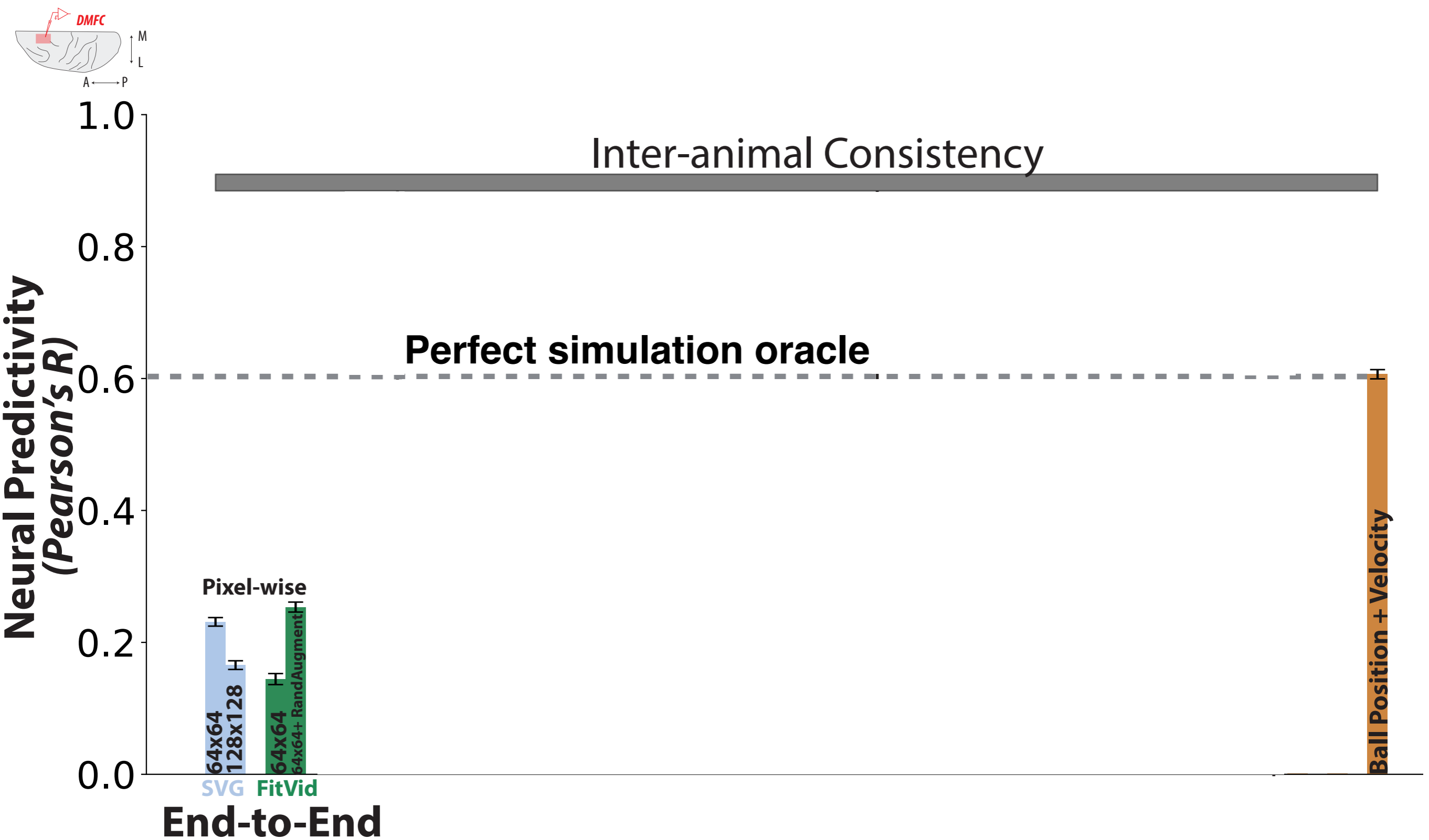


Predicts at the level of object slot representations and their relations

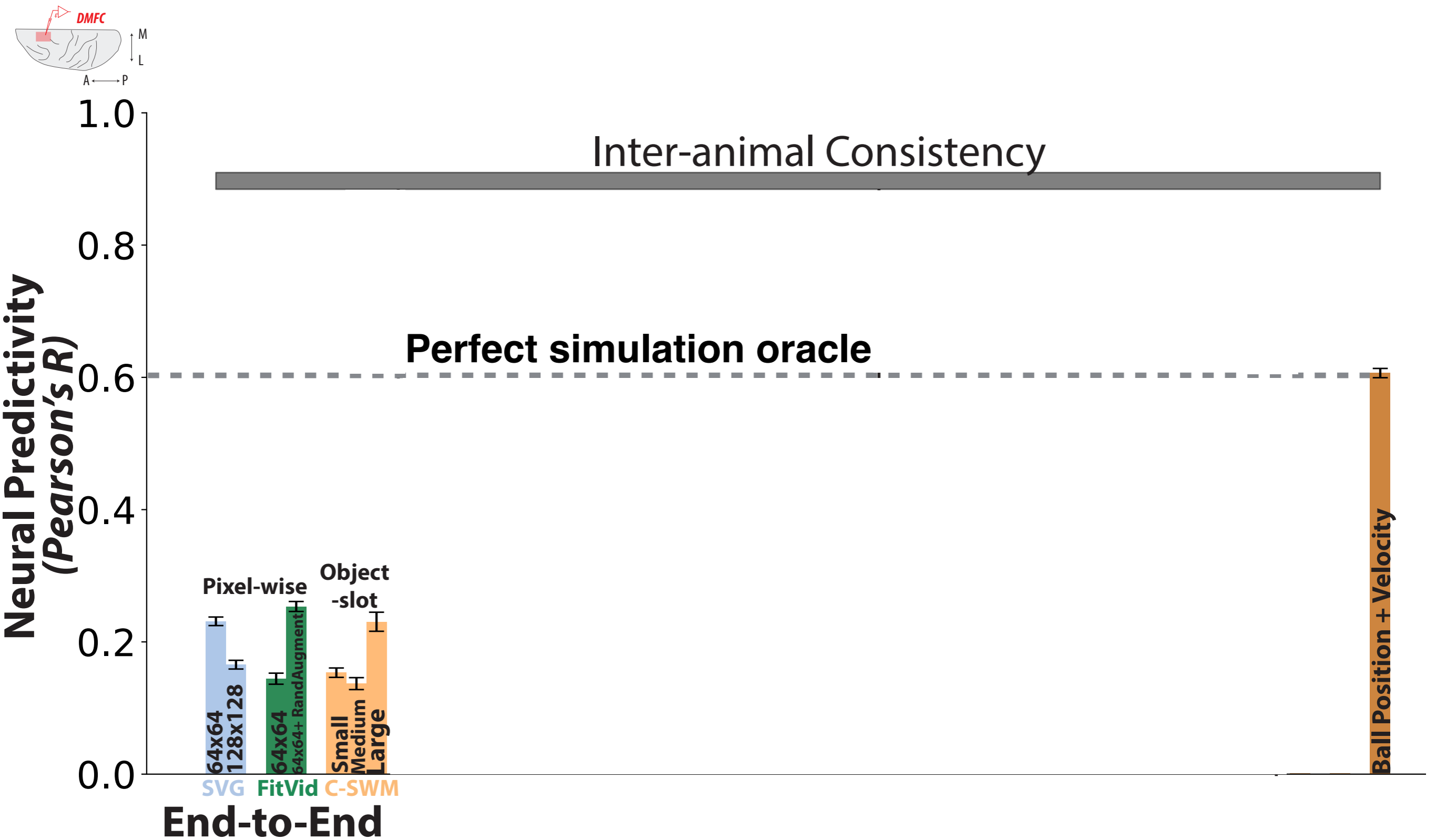
Kipf et al. 2020



Pixel-wise Future Prediction Poorly Predicts Neurons

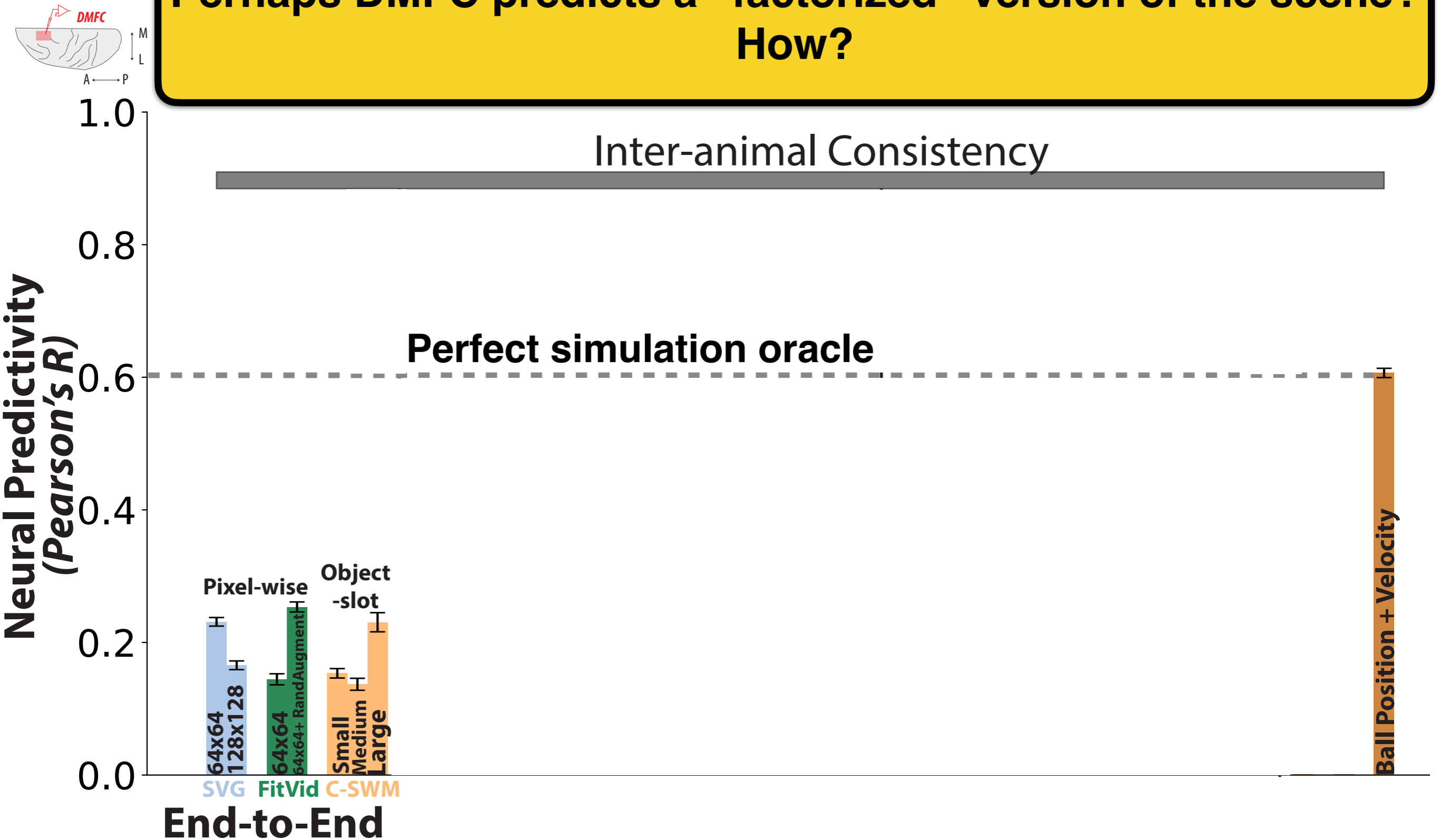


Object Slot Future Prediction Poorly Predicts Neurons



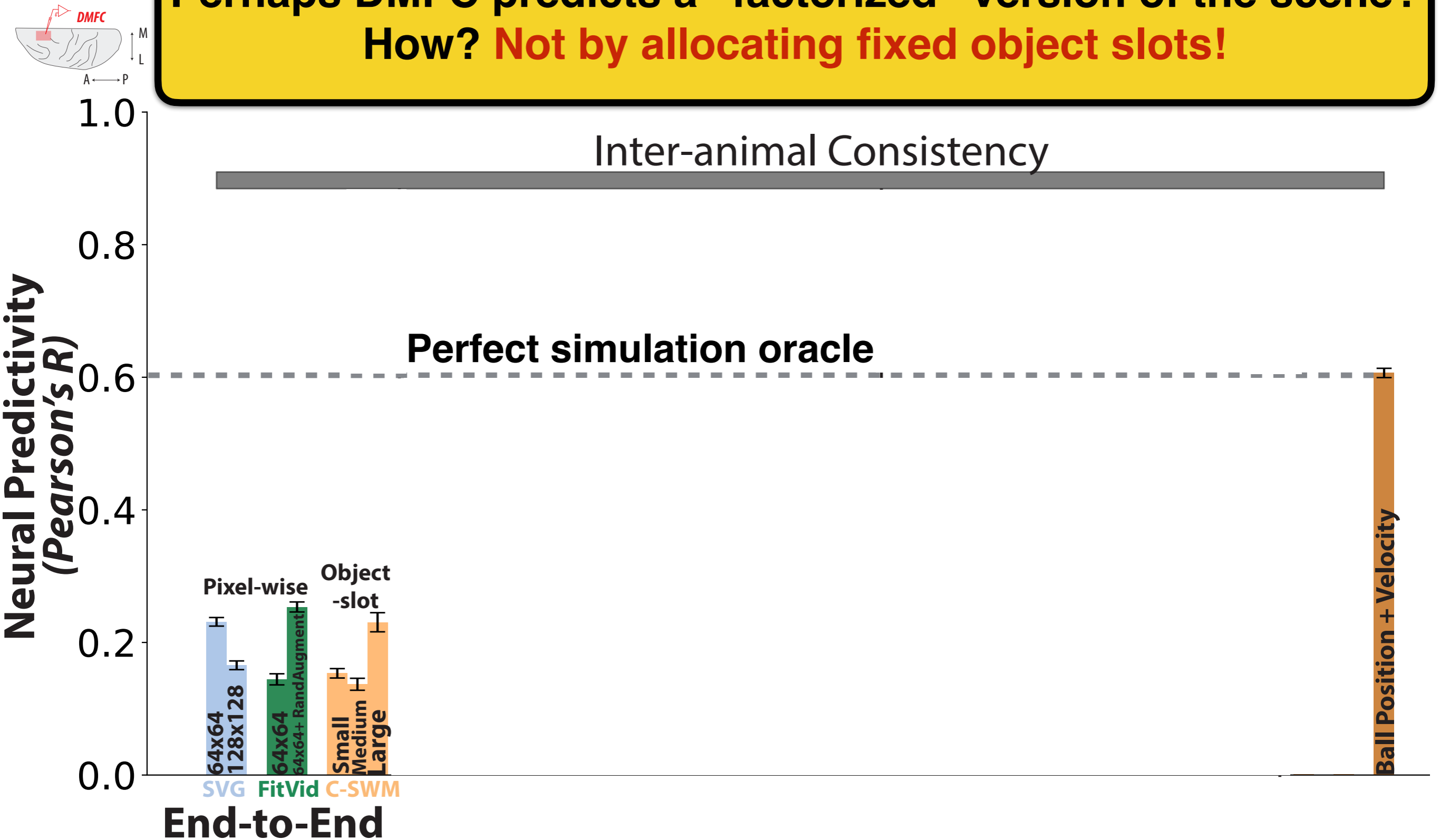
Object Slot Future Prediction Poorly Predicts Neurons

Perhaps DMFC predicts a “factorized” version of the scene?
How?

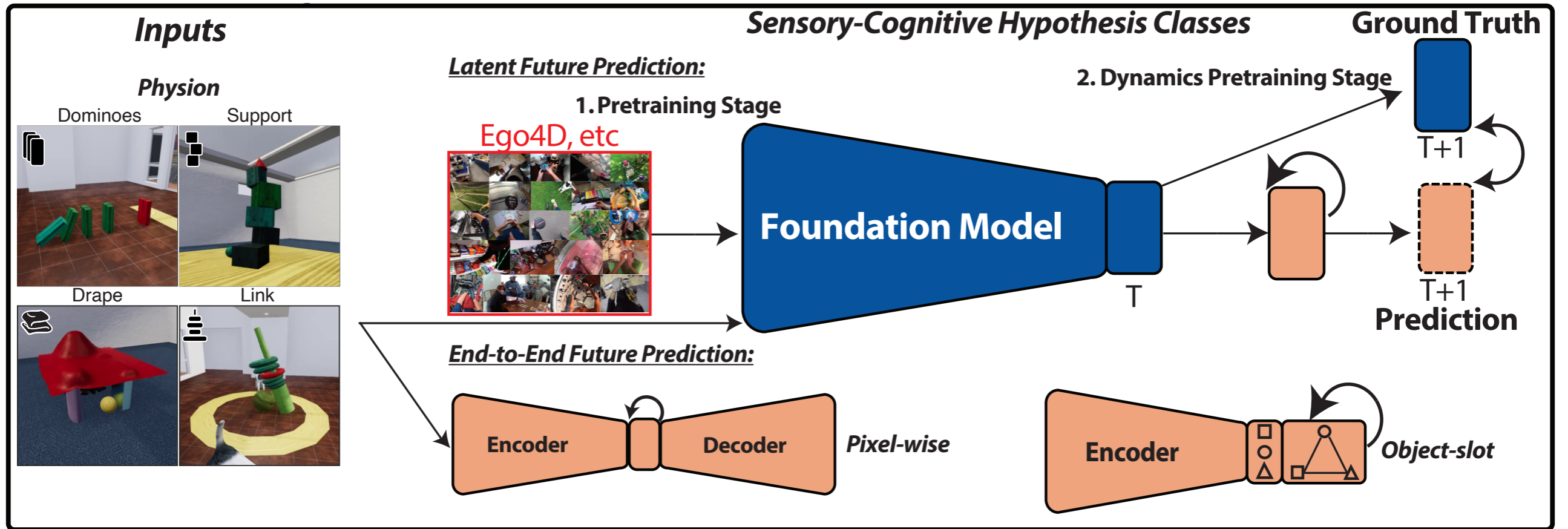


Object Slot Future Prediction Poorly Predicts Neurons

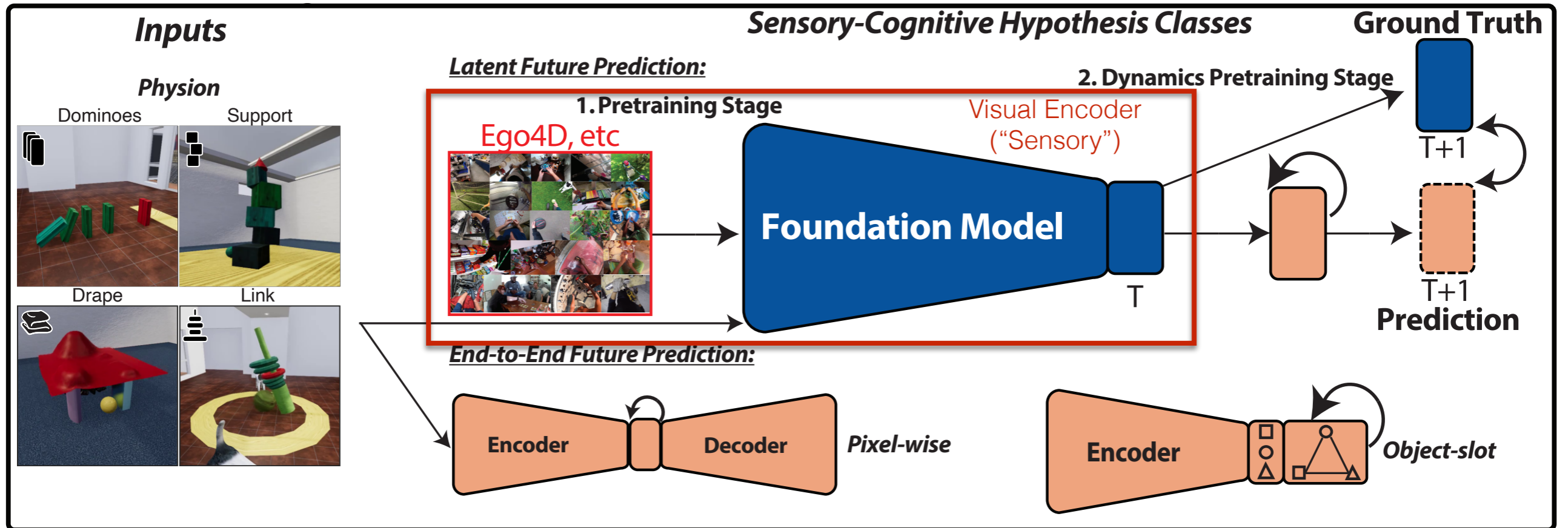
Perhaps DMFC predicts a “factorized” version of the scene?
How? **Not by allocating fixed object slots!**



Hypothesis Class 3: Latent Future Prediction

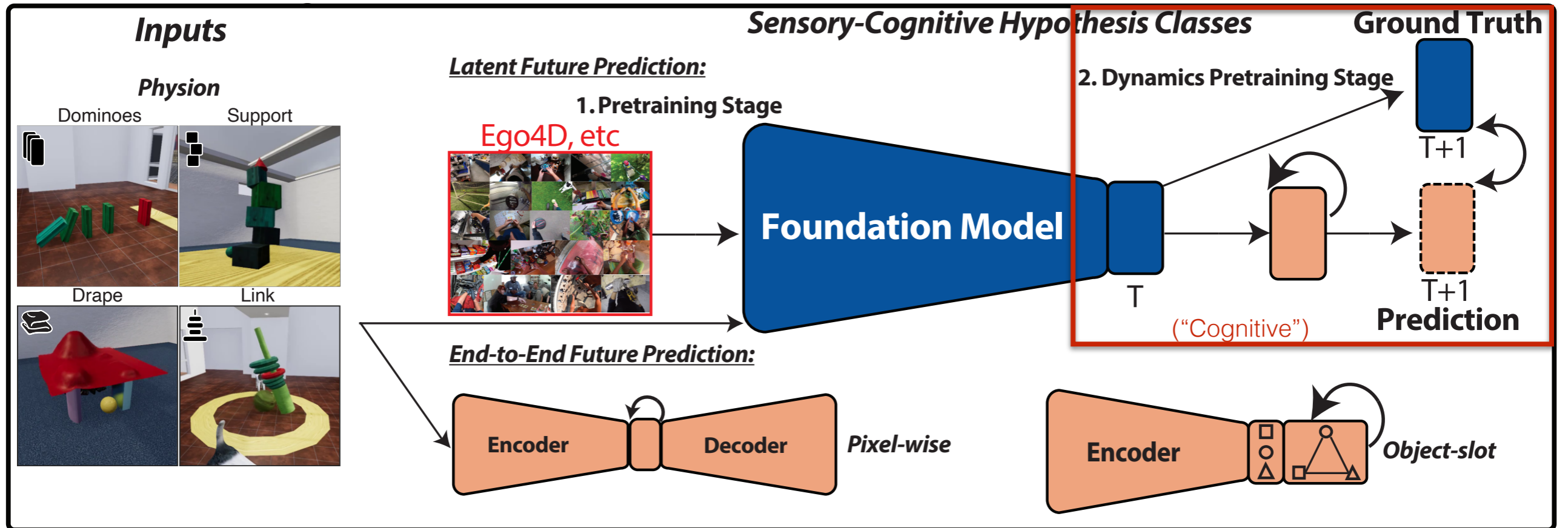


Hypothesis Class 3: Latent Future Prediction



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

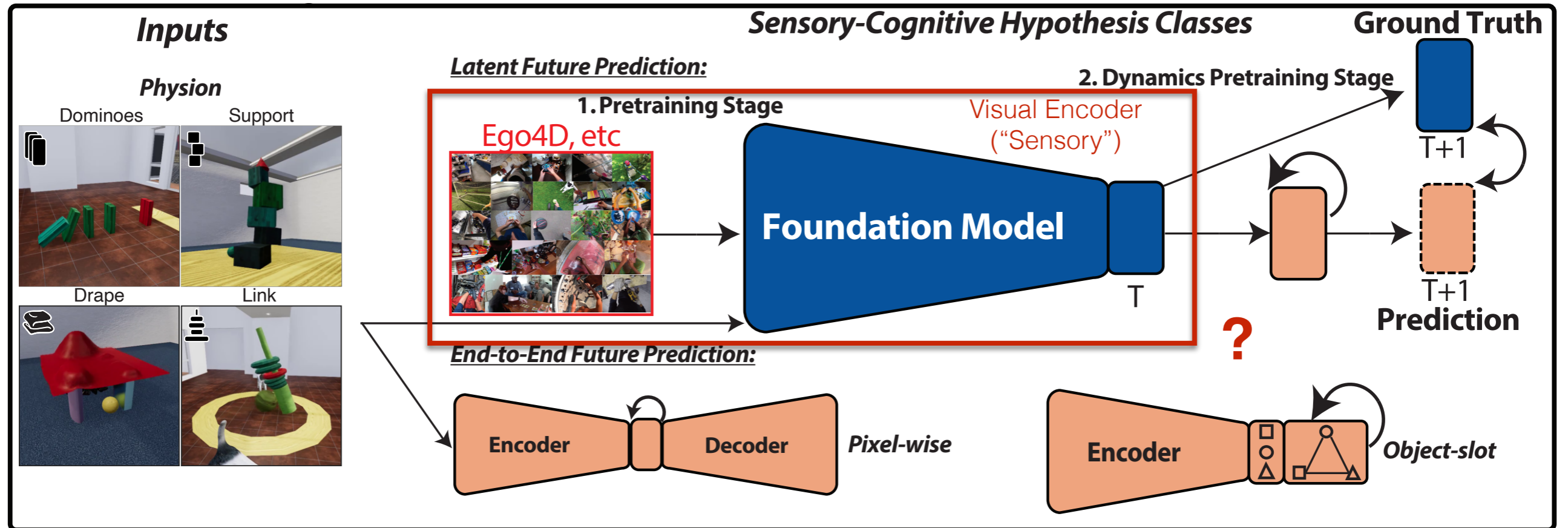
Hypothesis Class 3: Latent Future Prediction



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

Hypothesis Class 3: Foundation Models



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

What vision task?

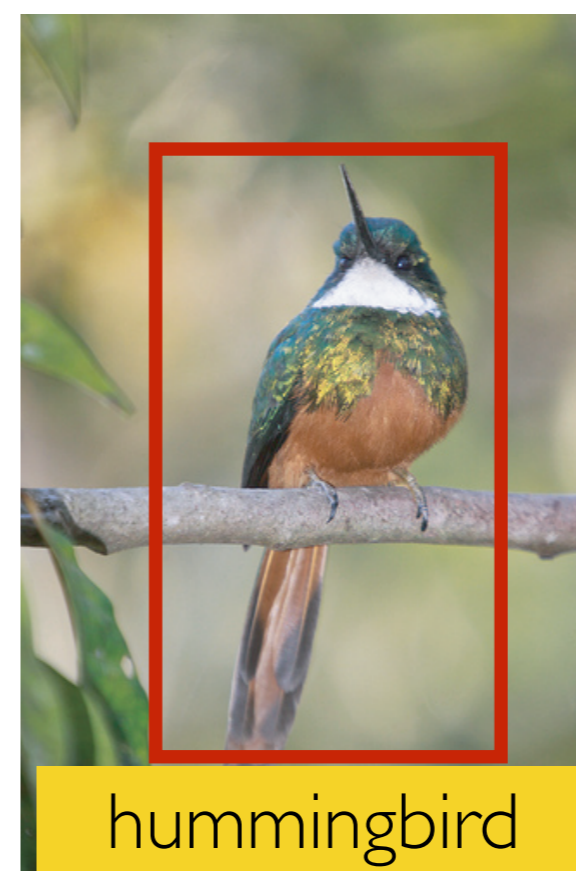
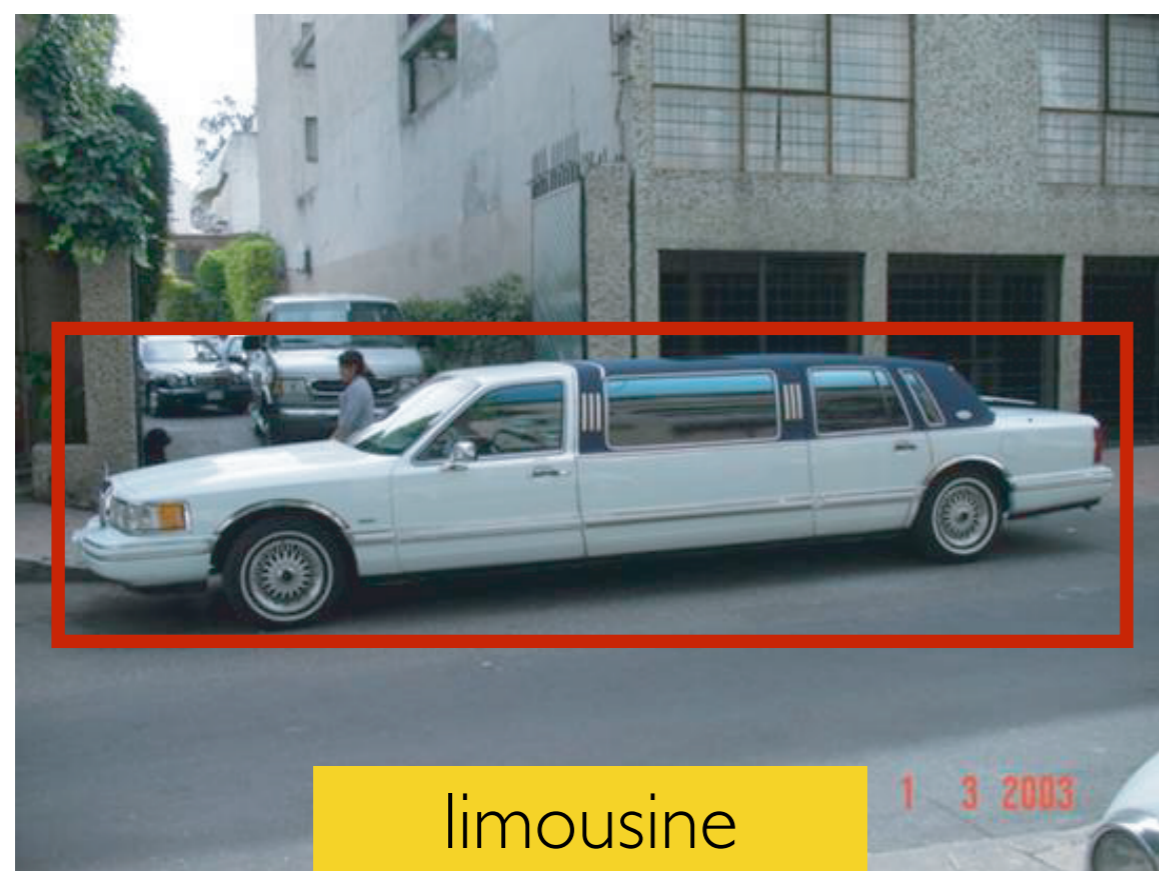
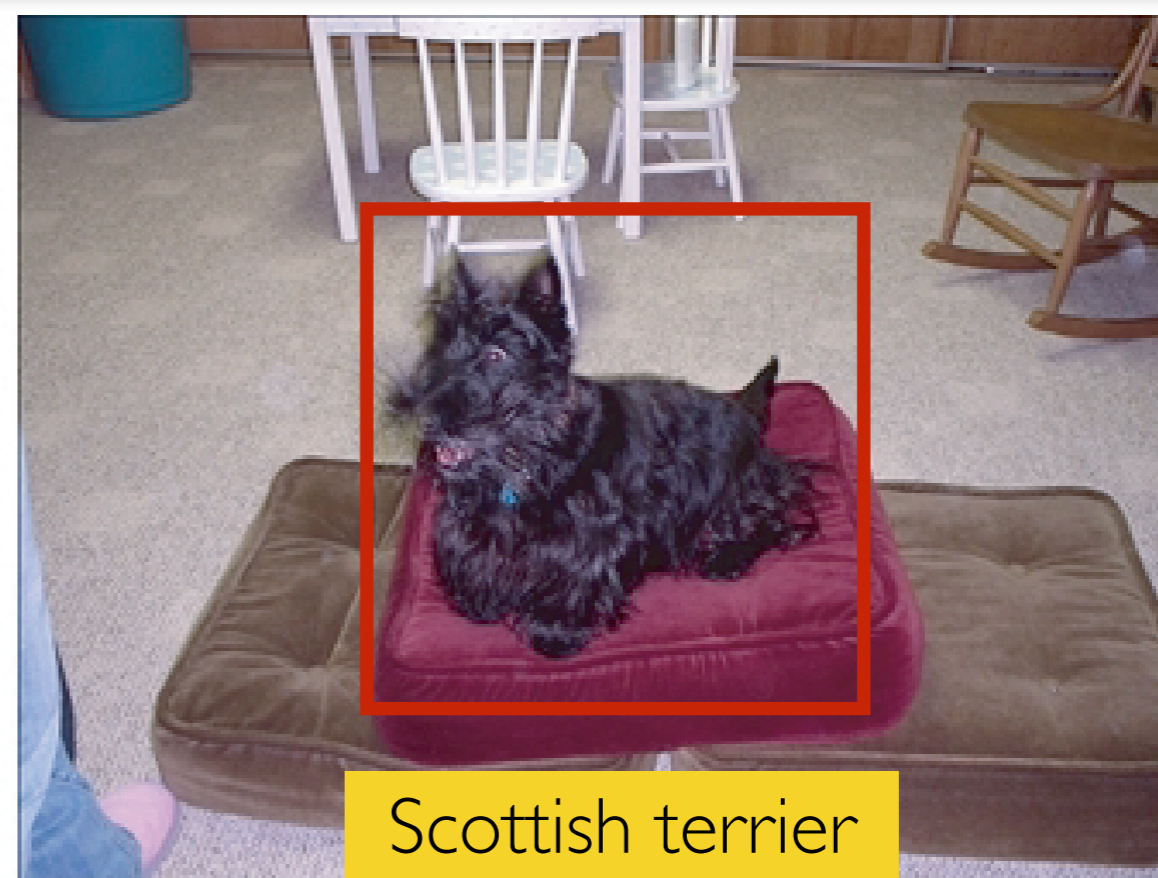
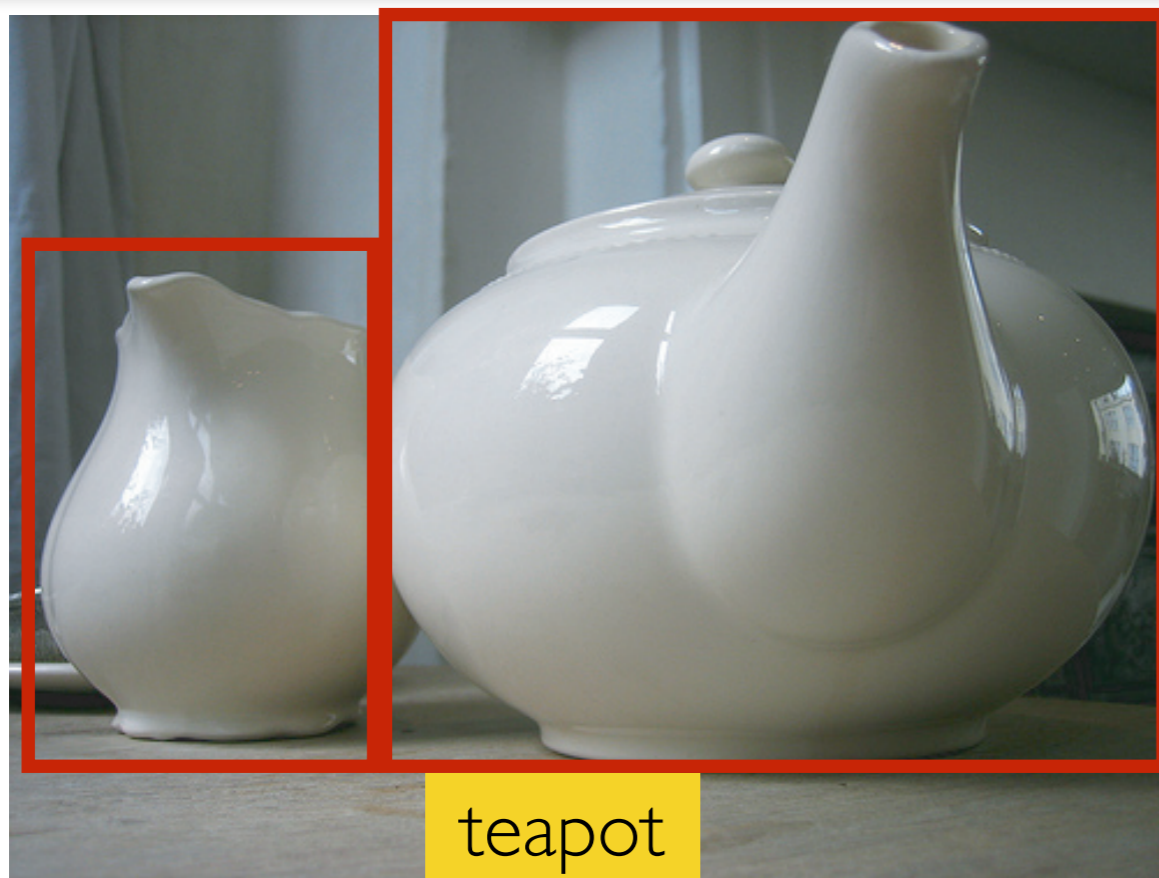
Leverage these dynamics to do explicit future prediction

Hypothesis Class 3: Static Image Foundation Models

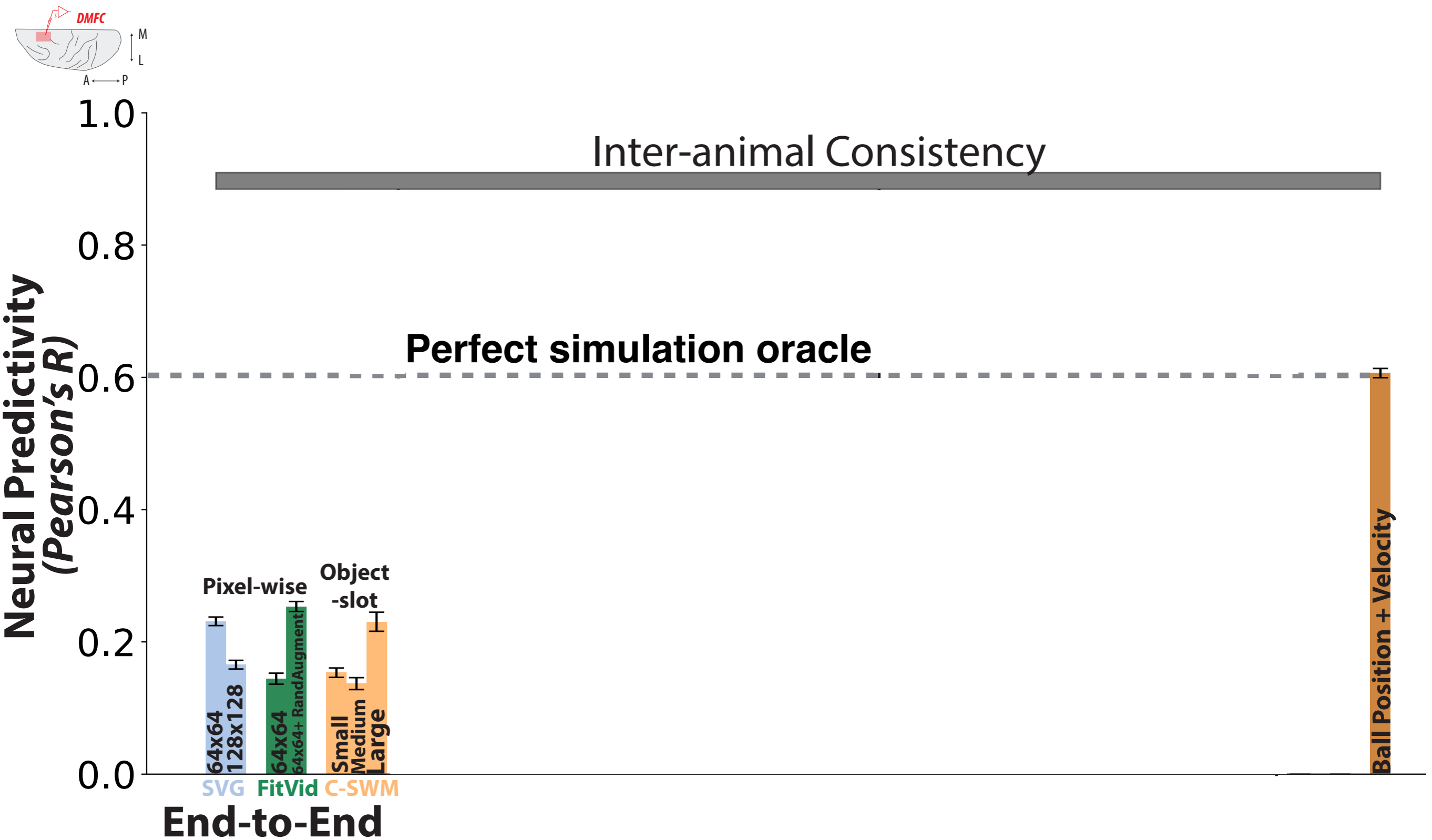
Hypothesis Class 3: Static Image Foundation Models



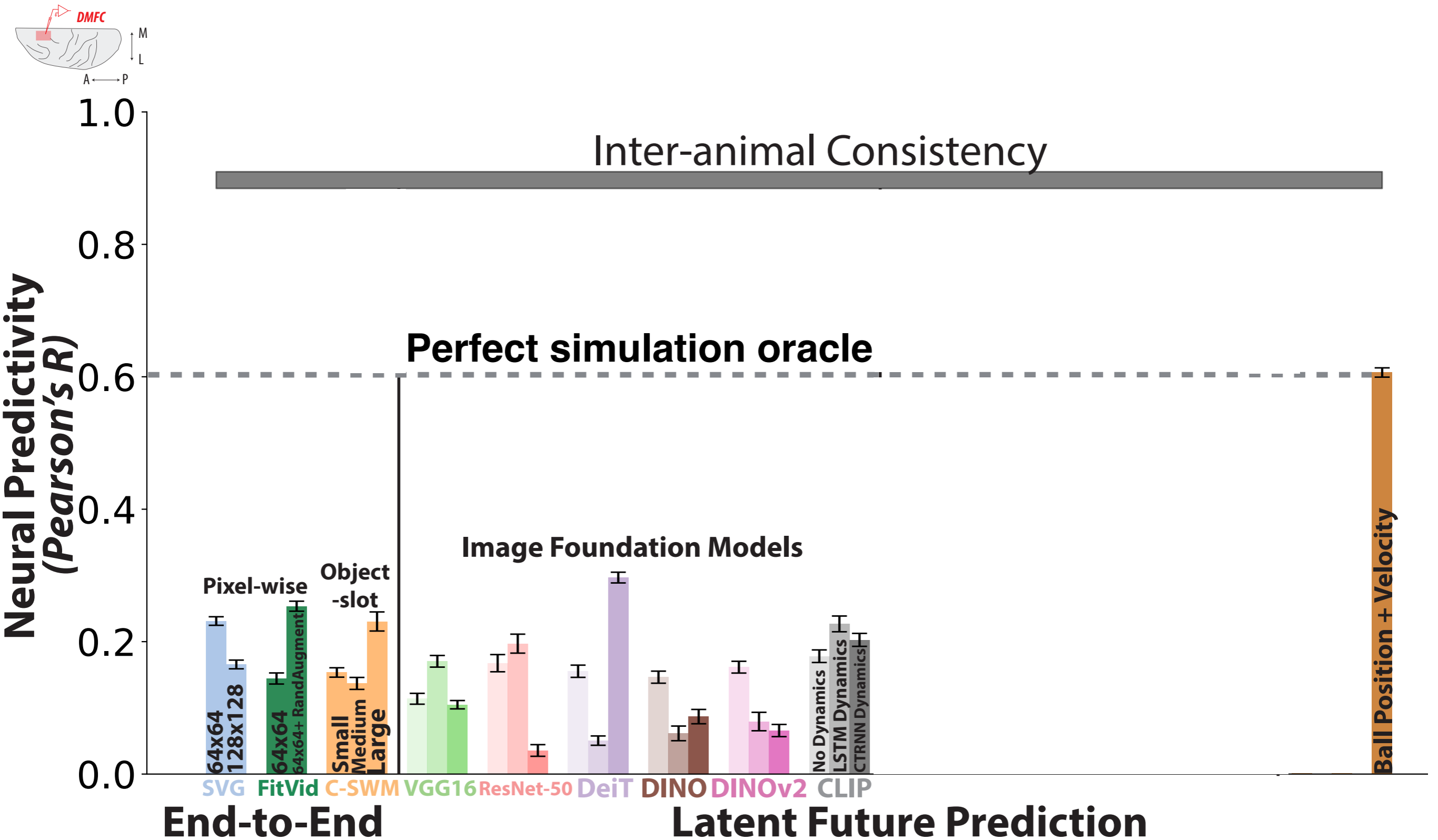
Hypothesis Class 3: Static Image Foundation Models



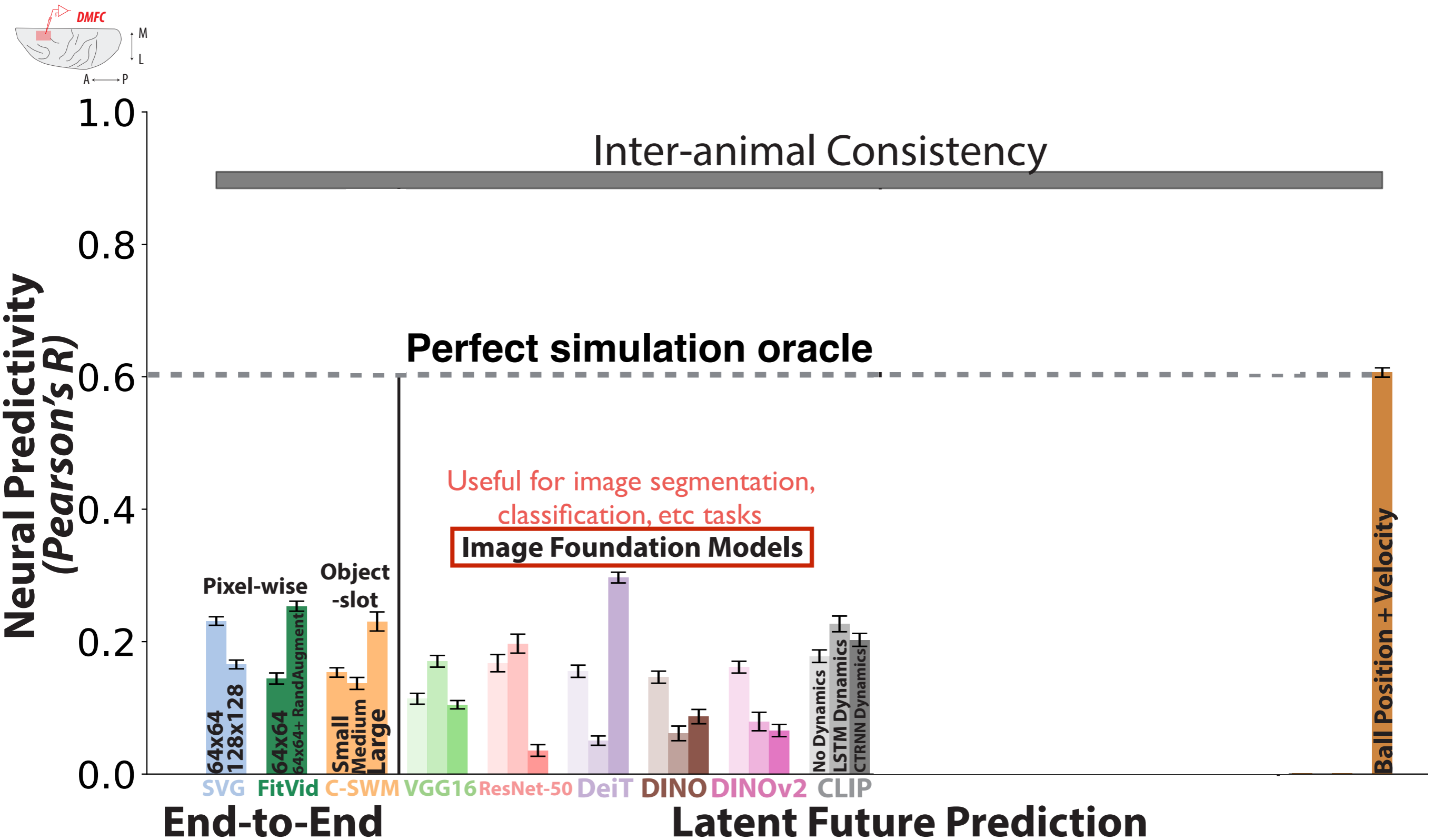
Object Slot Future Prediction Poorly Predicts Neurons



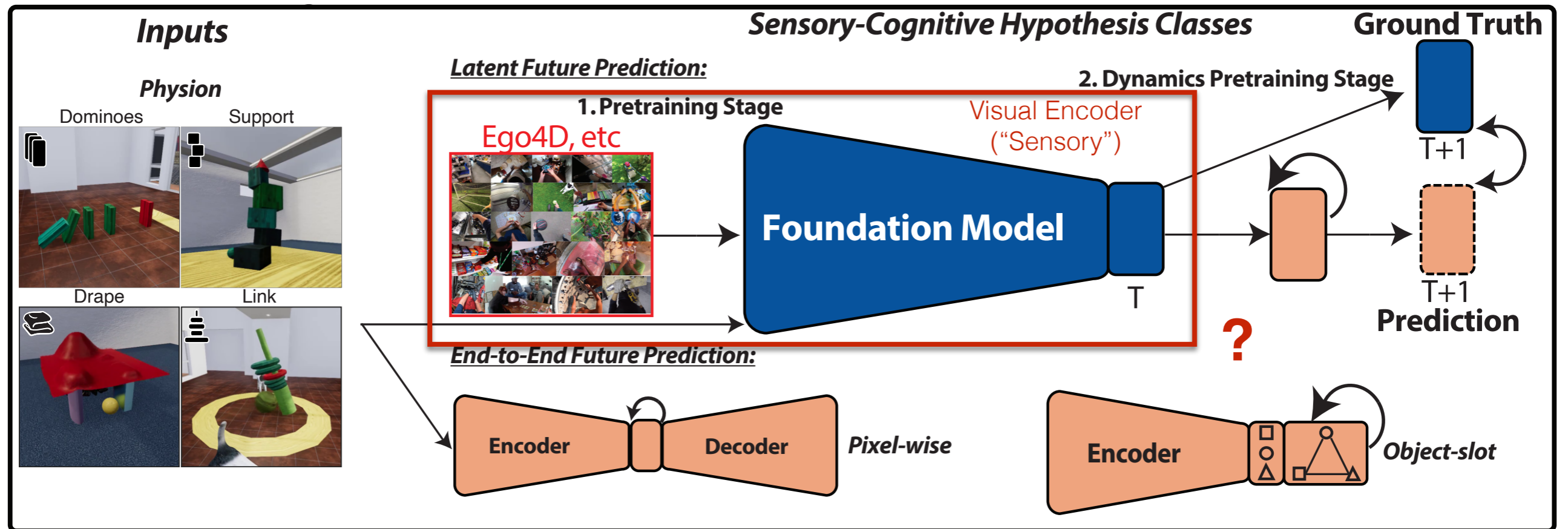
Static Image Foundation Future Prediction Poorly Predicts Neurons



Static Image Foundation Future Prediction Poorly Predicts Neurons



Hypothesis Class 3: Foundation Models

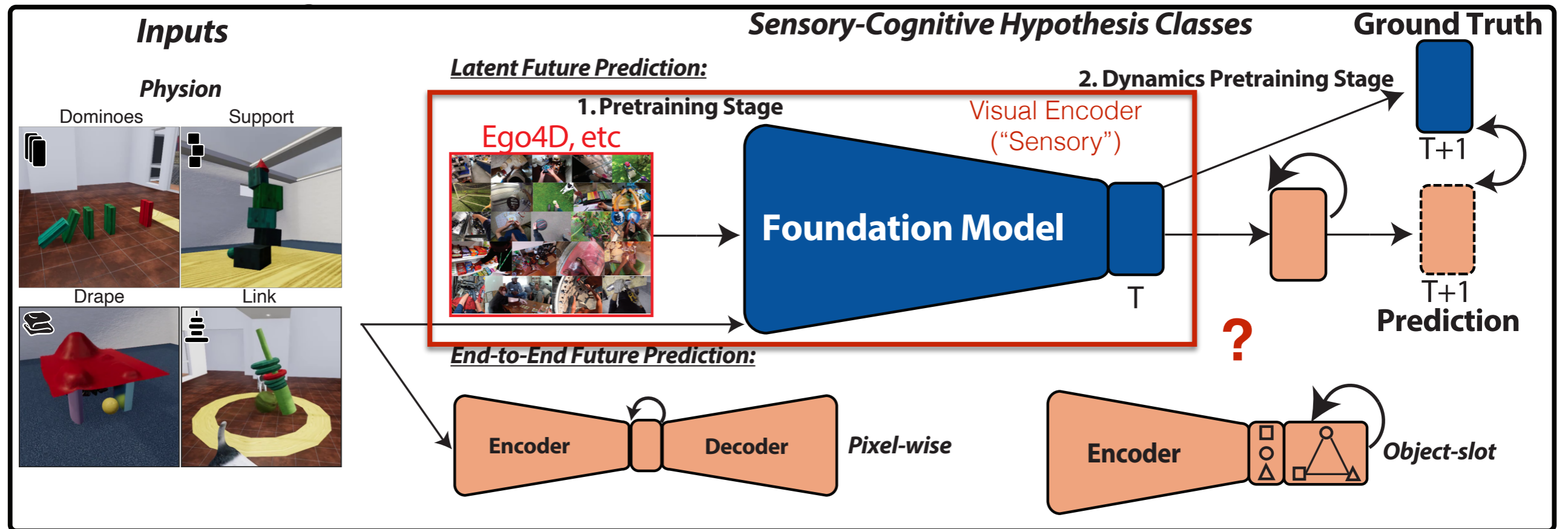


Learn a partial, *implicit* representation of the physical world by performing a challenging vision task ("foundation model")

What vision task?

Leverage these dynamics to do explicit future prediction

Hypothesis Class 3: Foundation Models



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

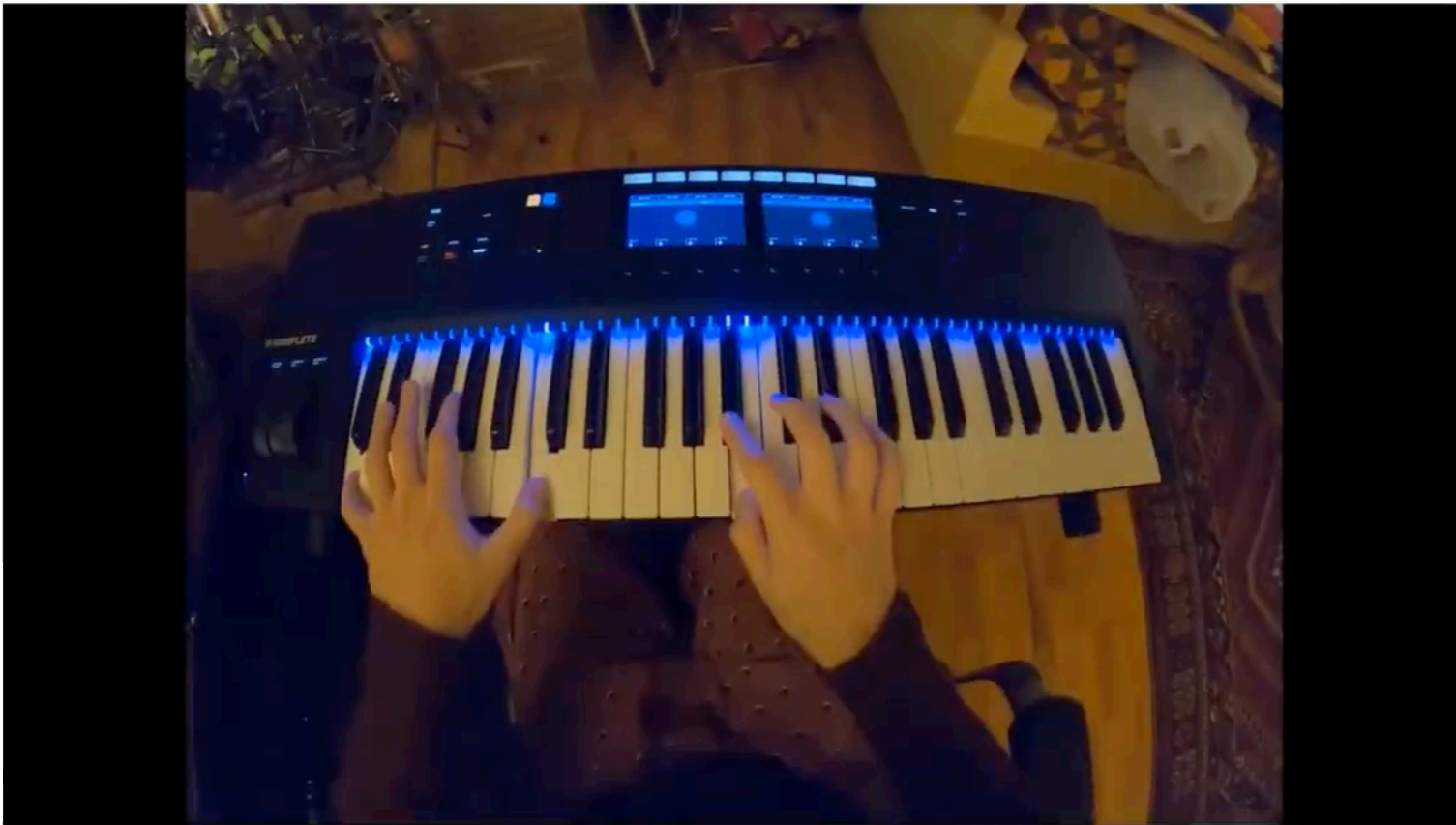
What vision task?

We do far more than engage with static images!

Leverage these dynamics to do explicit future prediction

Hypothesis Class 3: Video Foundation Models

Ego4D: everyday activity around the world



Ego4D: A massive-scale egocentric dataset

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

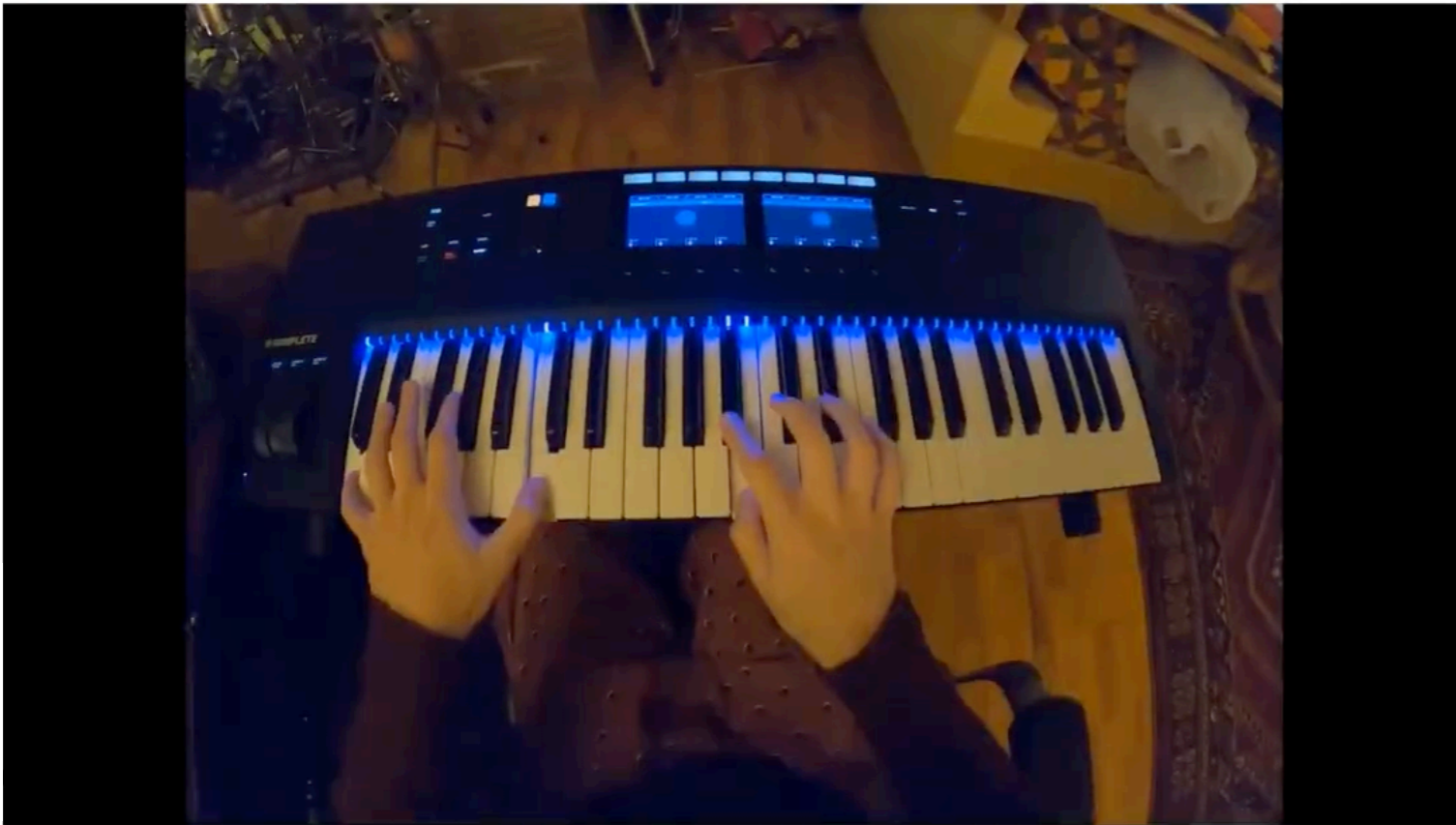
931 participants from 74 worldwide locations

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



Hypothesis Class 3: Video Foundation Models

Ego4D: everyday activity around the world



Ego4D: A massive-scale egocentric dataset

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

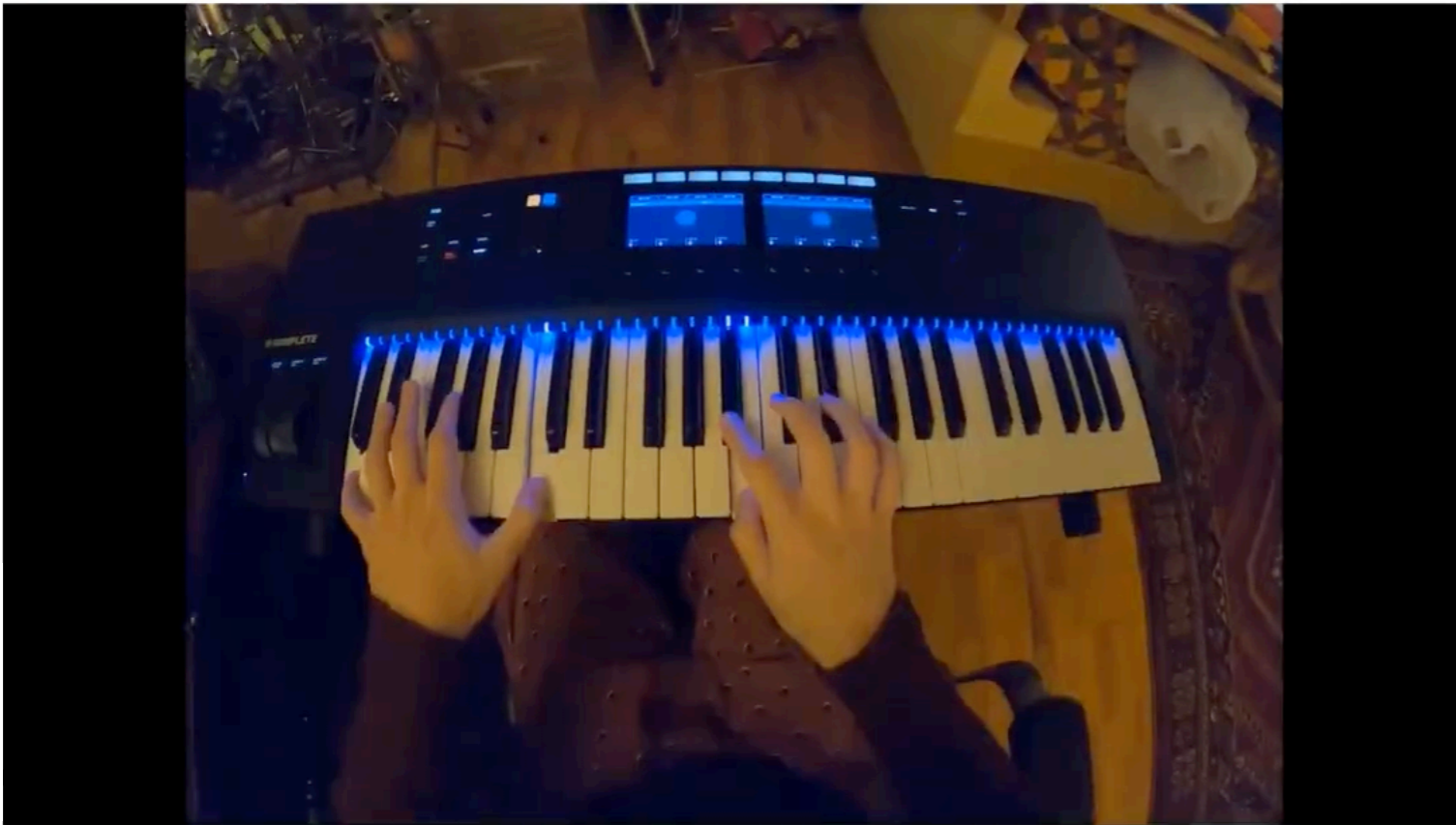
931 participants from 74 worldwide locations

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



Hypothesis Class 3: Video Foundation Models

Ego4D: everyday activity around the world



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

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

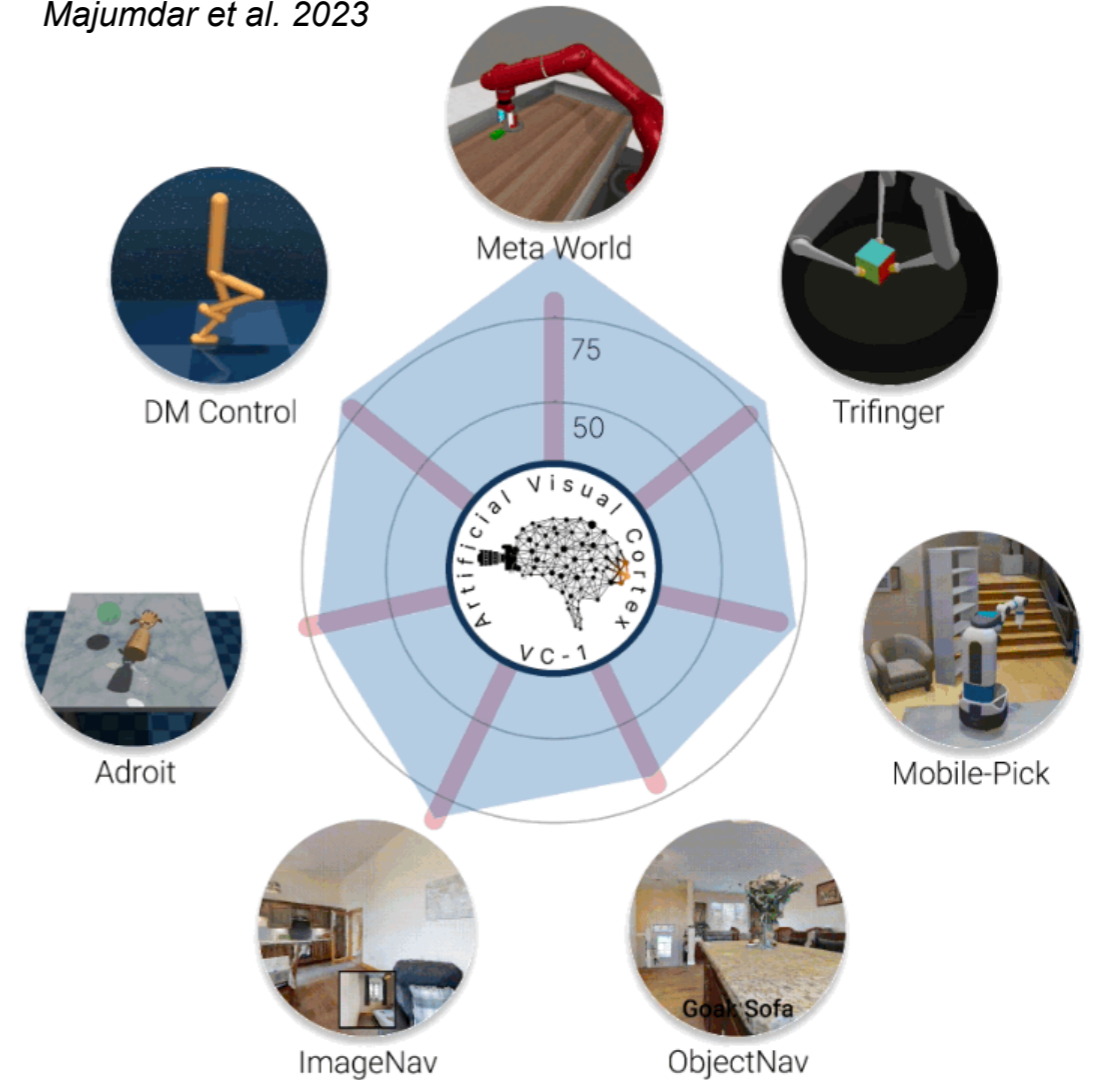


Hypothesis Class 3: Video Foundation Models

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

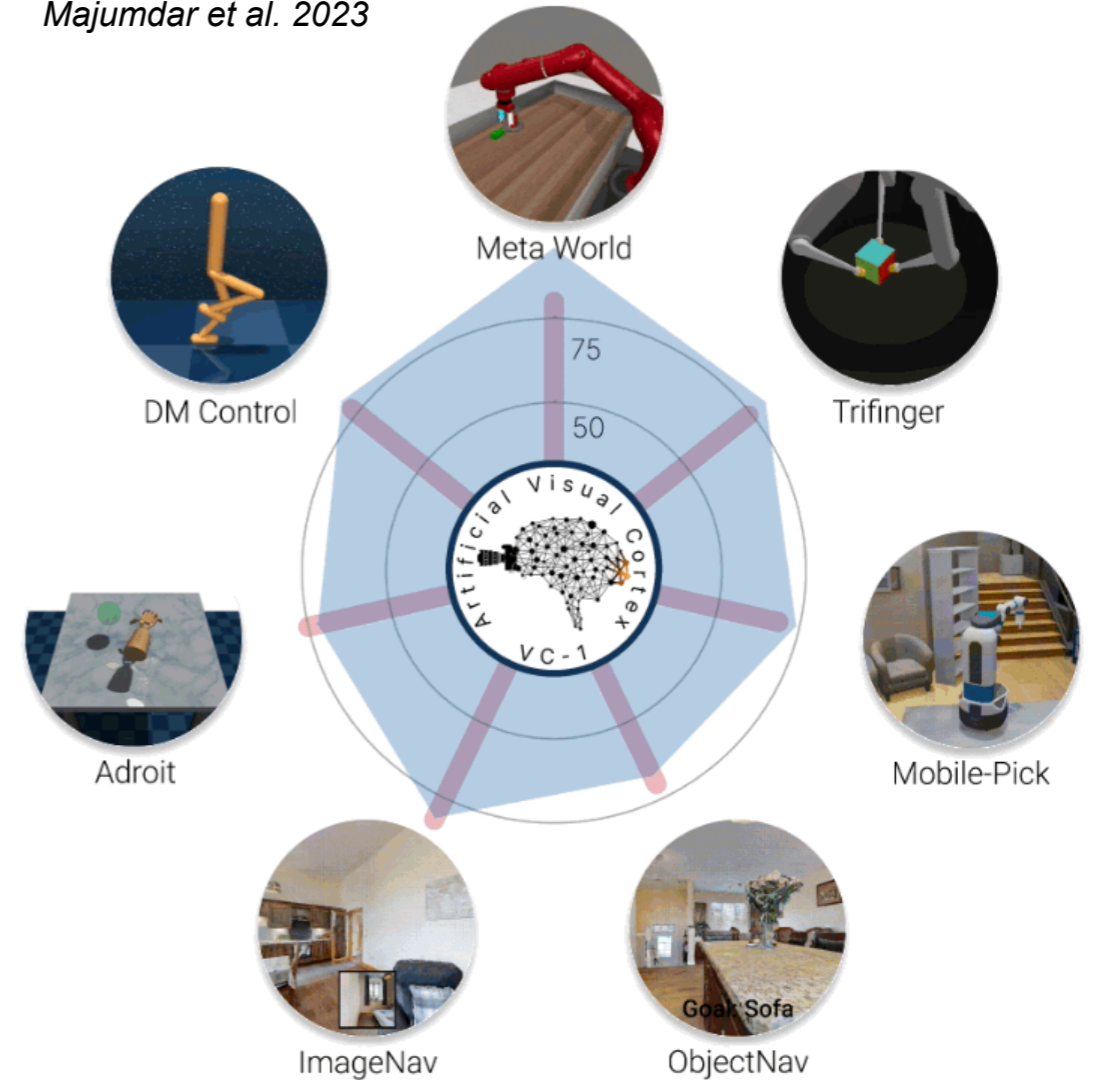


Hypothesis Class 3: Video Foundation Models

Ego4D: everyday activity around the world



Majumdar et al. 2023



Ego4D: A massive-scale egocentric dataset

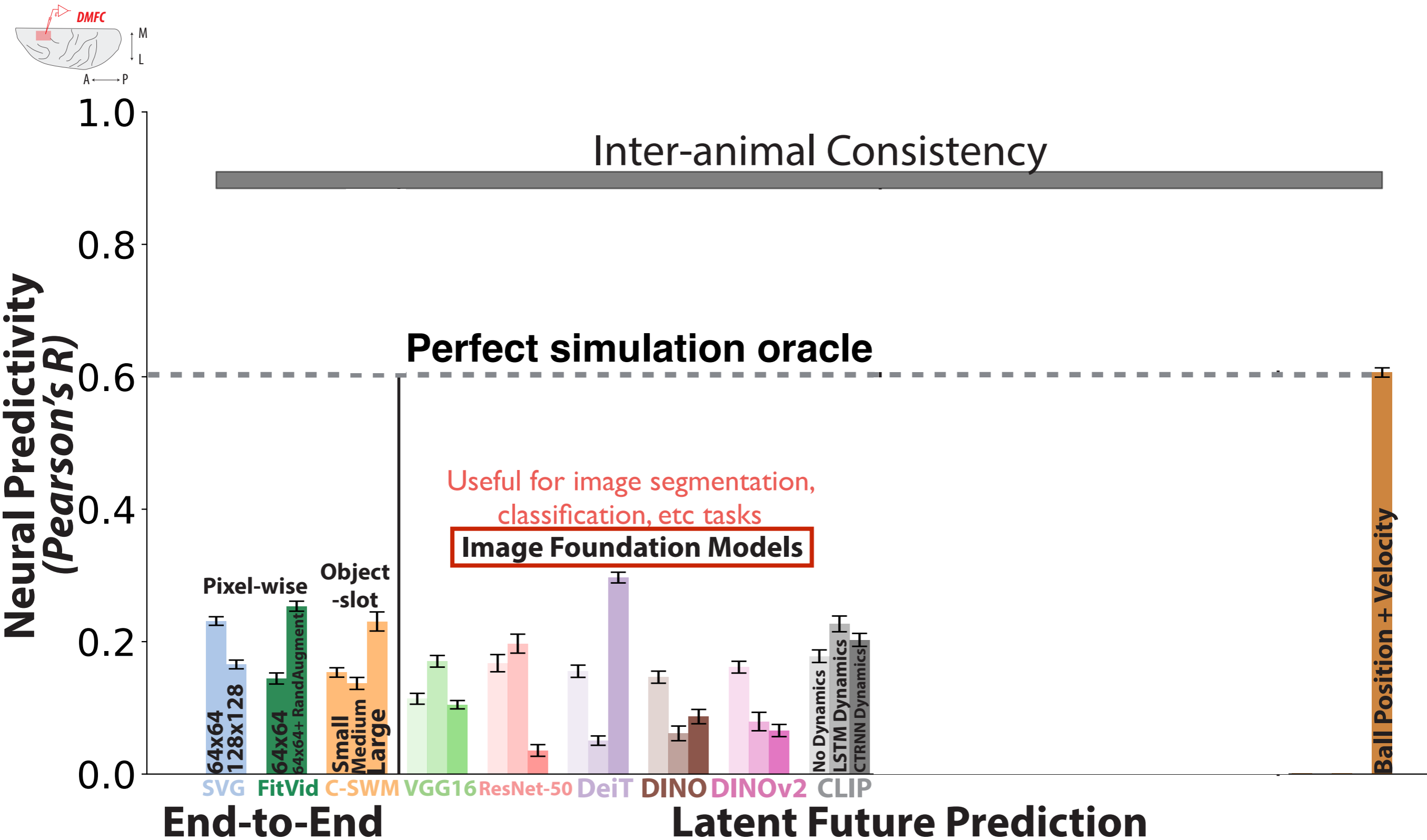
3,670 hours of in-the-wild daily life activity

931 participants from 74 worldwide locations

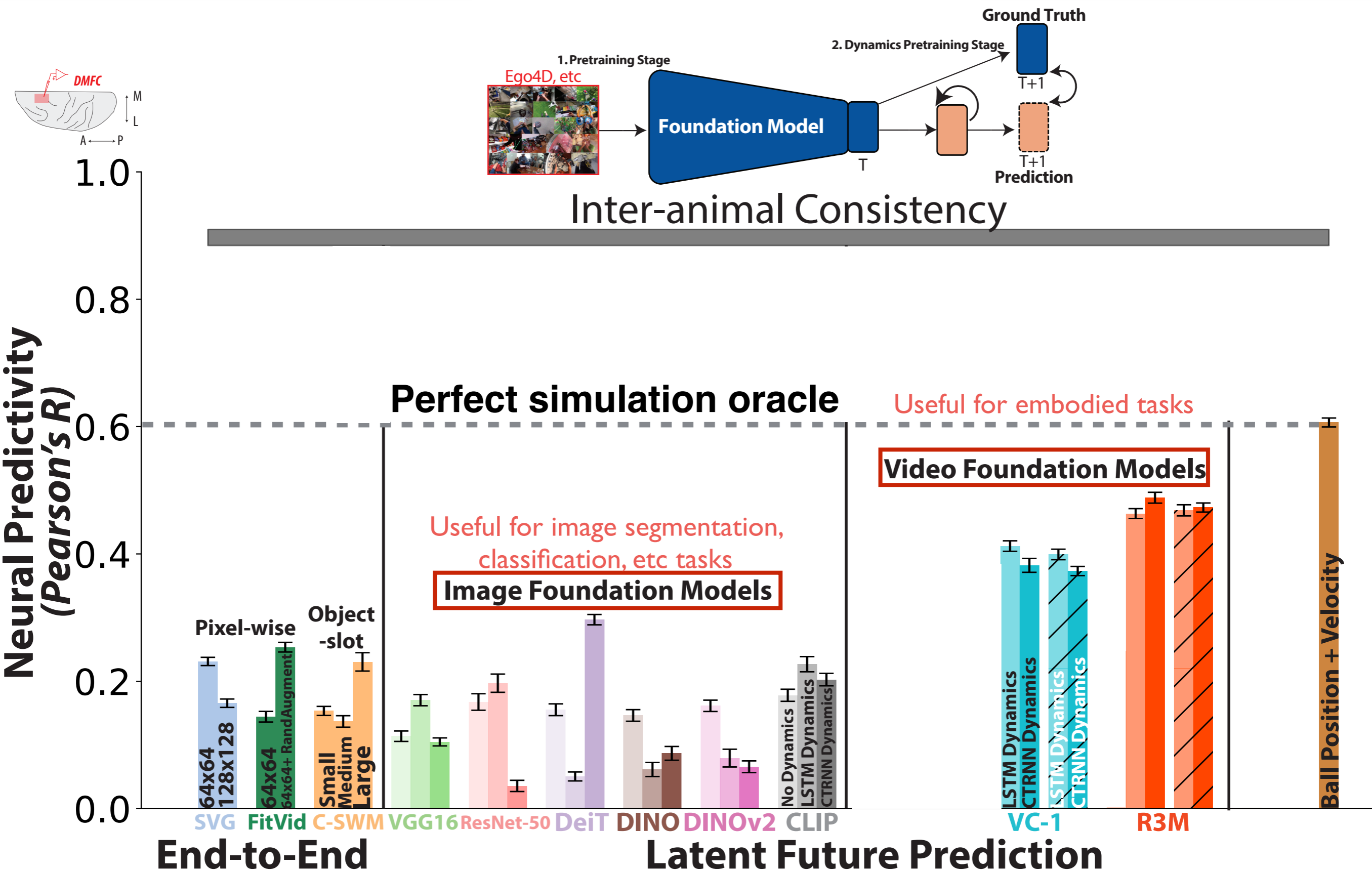
Multimodal: audio, 3D scans, IMU, stereo, multi-camera



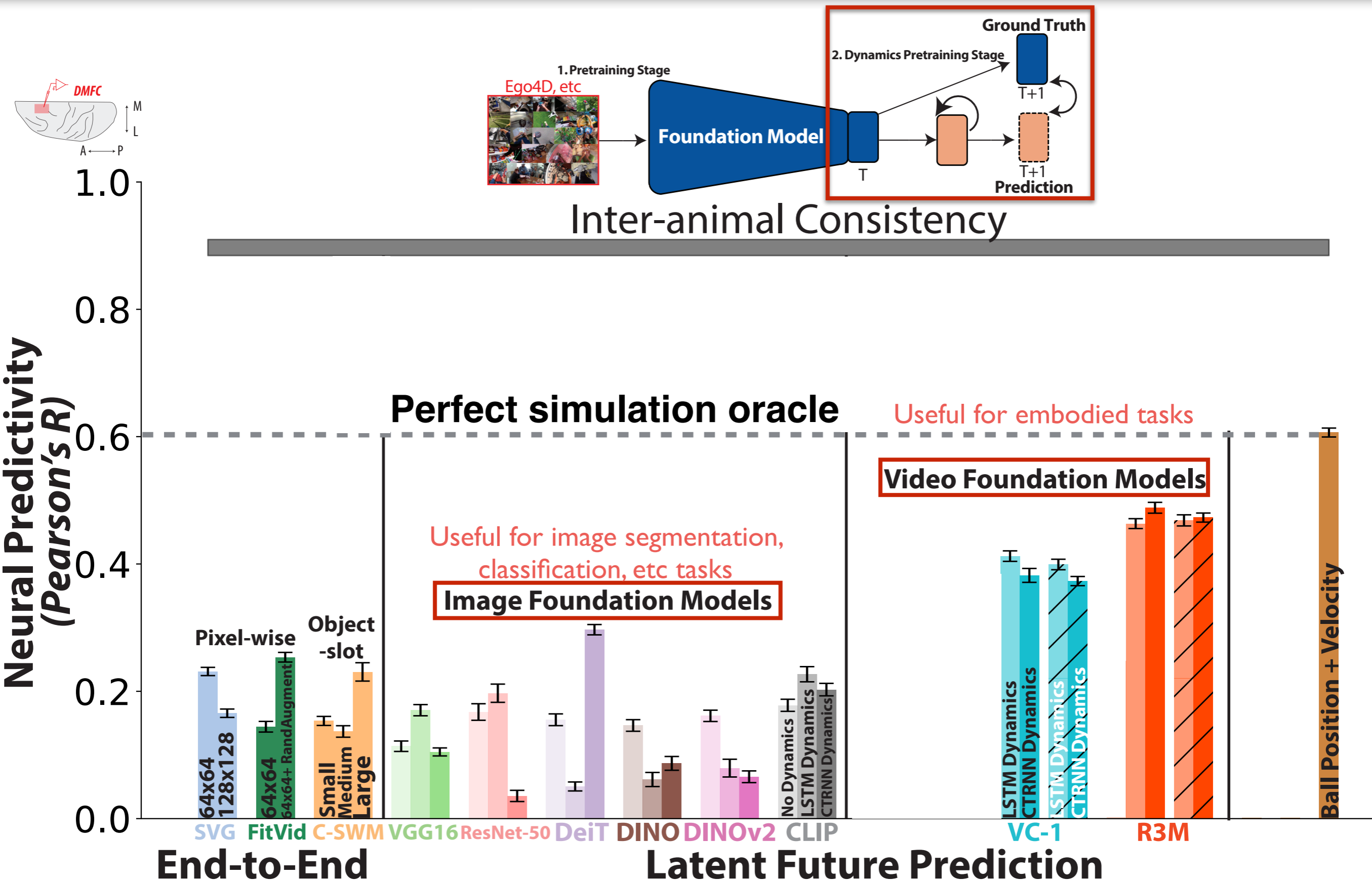
Static Image Foundation Future Prediction Poorly Predicts Neurons



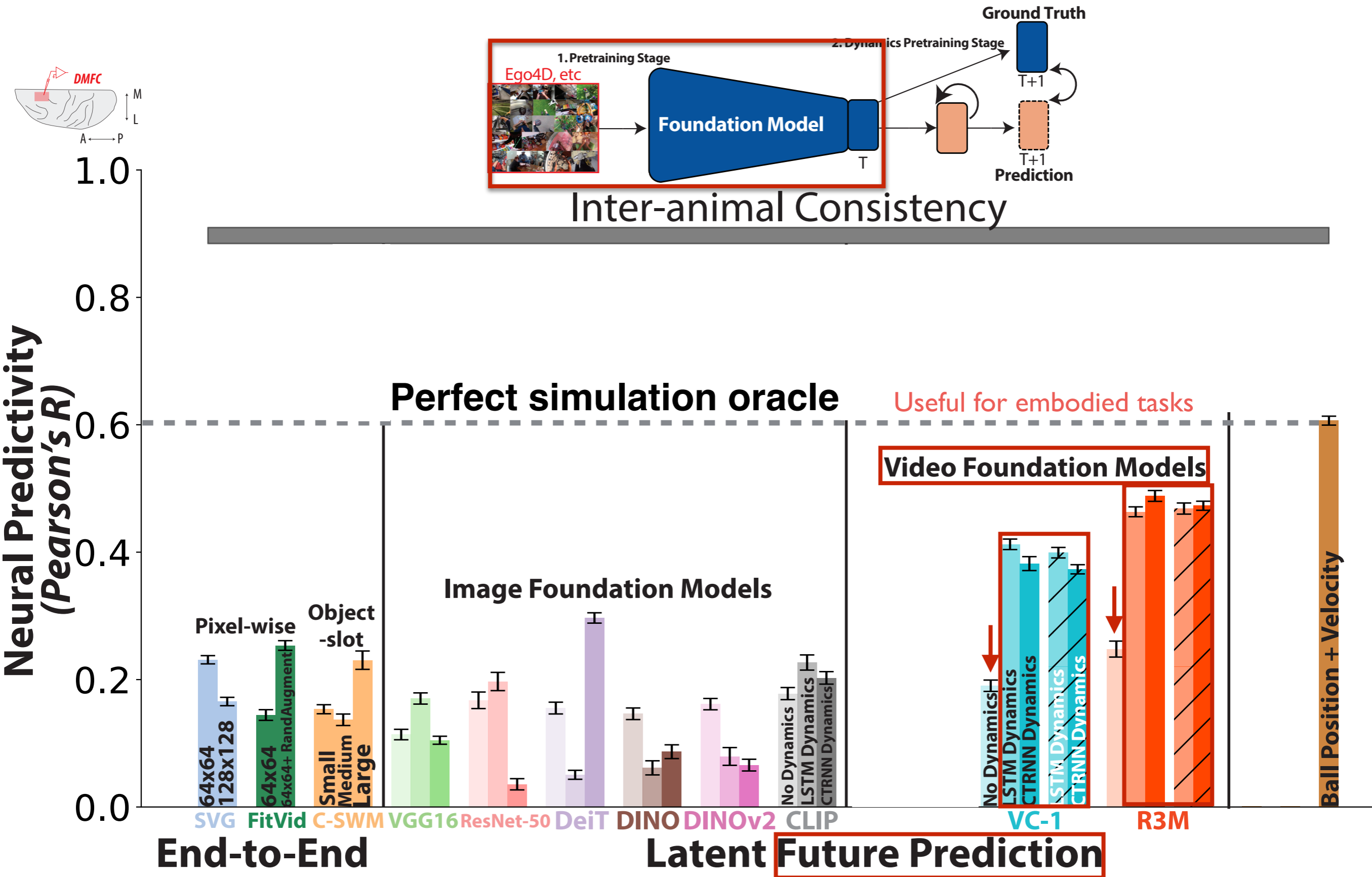
Video Foundation Future Prediction Best Predict Neurons



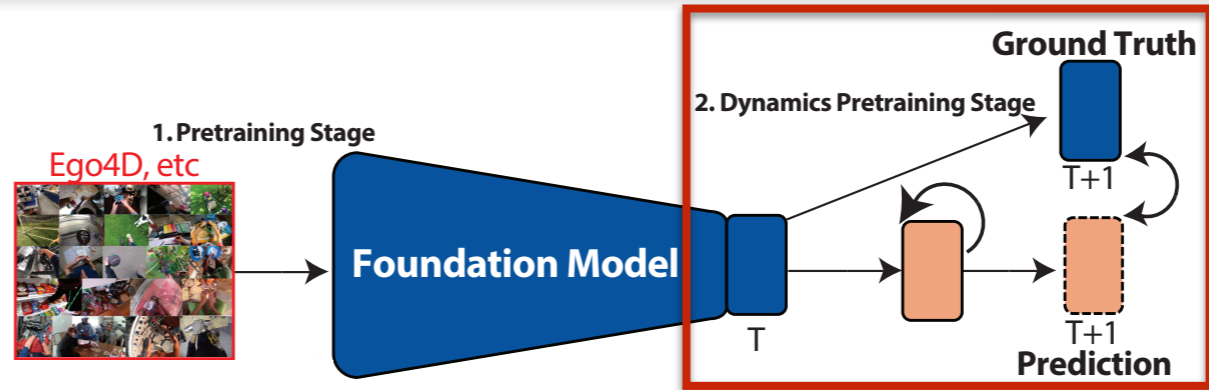
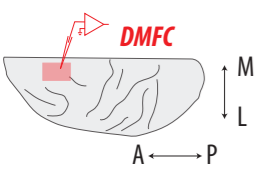
Video Foundation Future Prediction Best Predict Neurons



Video Foundation Future Prediction Best Predict Neurons



Video Foundation Future Prediction Best Predict Neurons



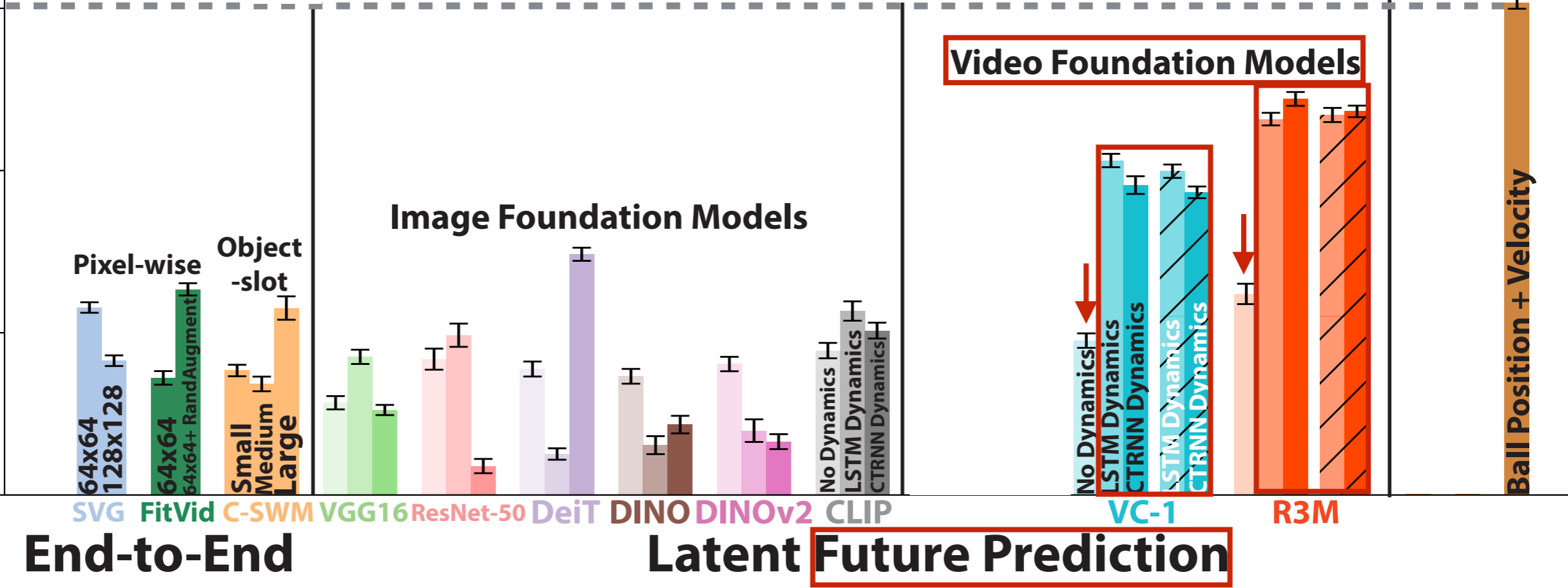
**Pretraining on Ego4D is not enough on its own:
Need explicit future prediction!**

Neural Predictivity
(Pearson's R)

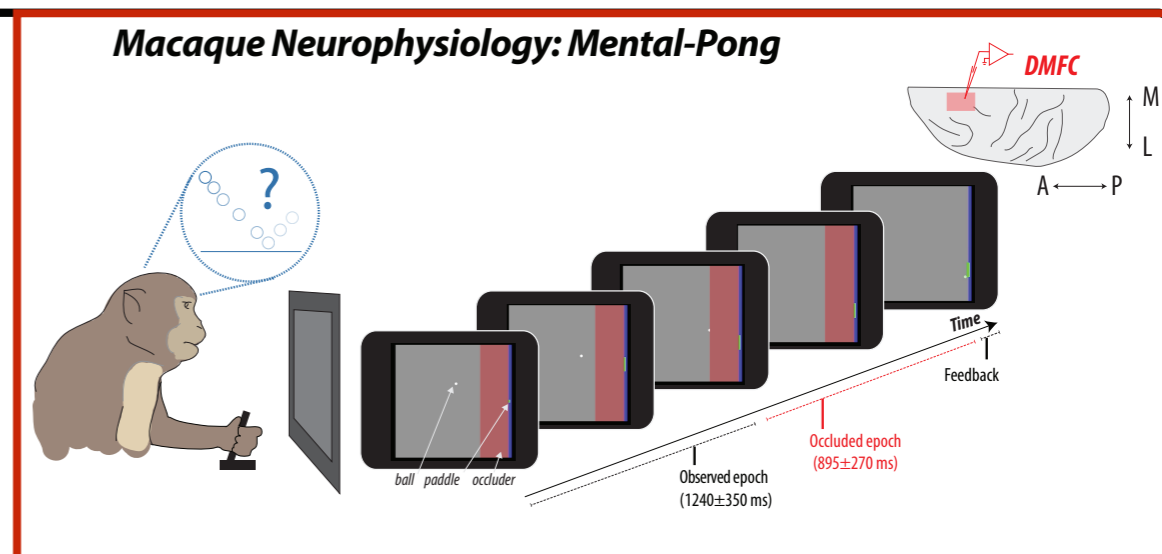
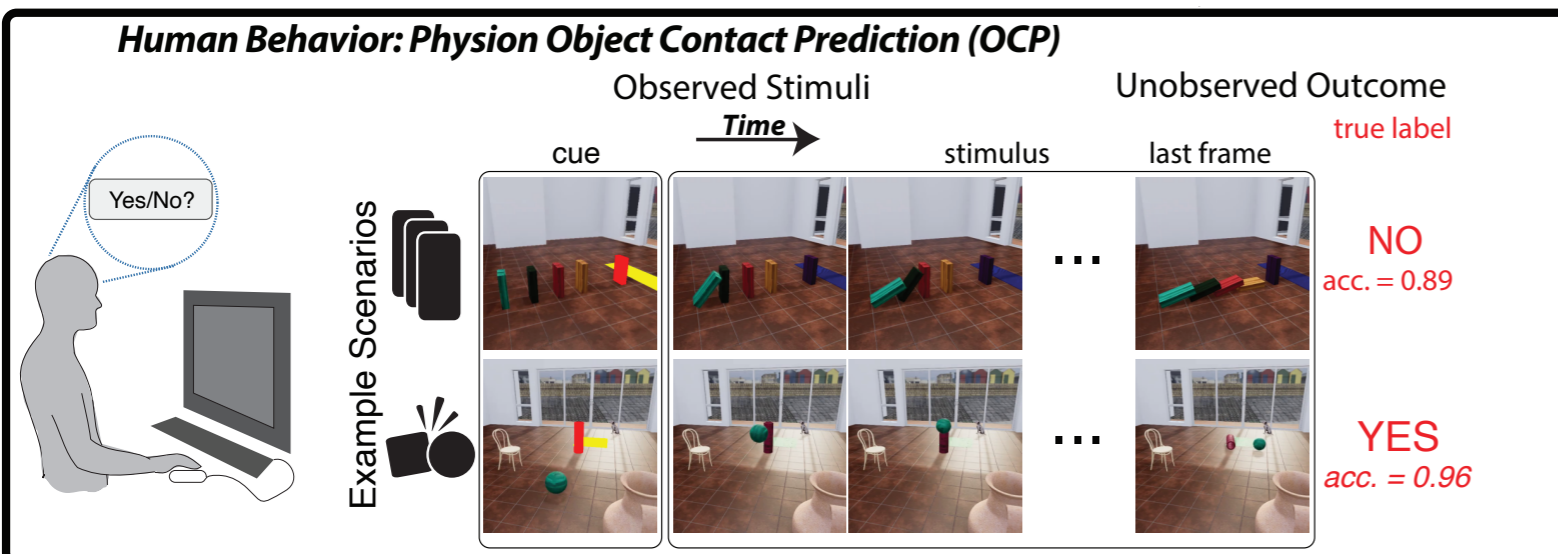
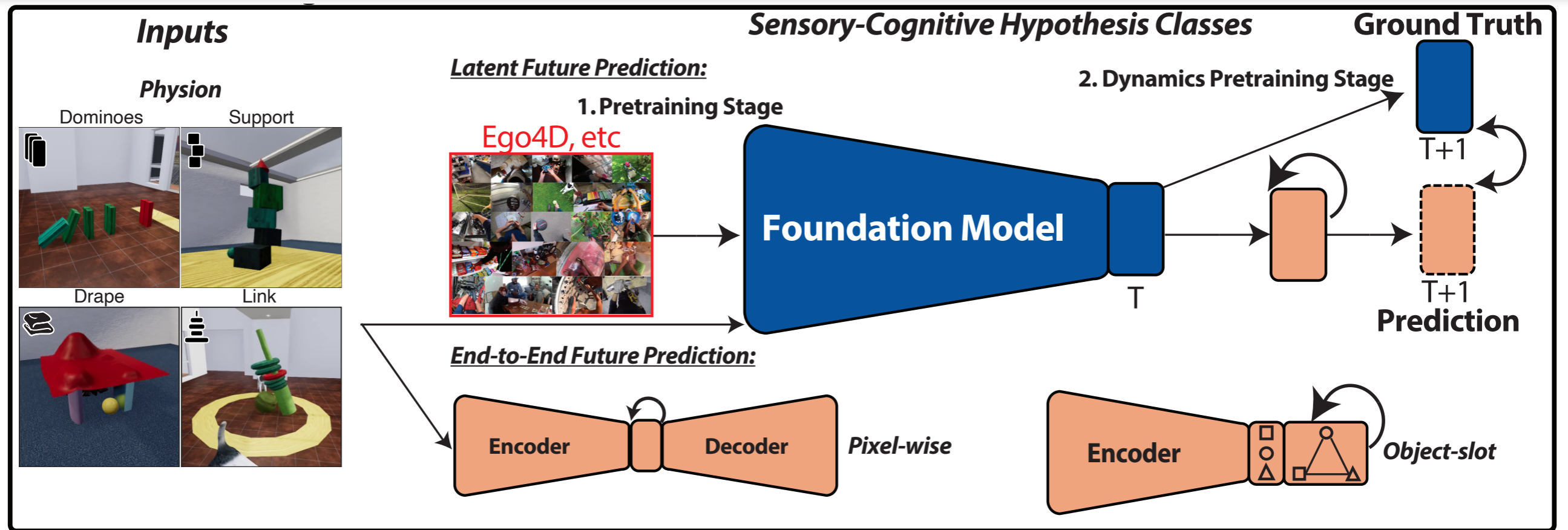
1.0
0.8
0.6
0.4
0.2
0.0

Perfect simulation oracle

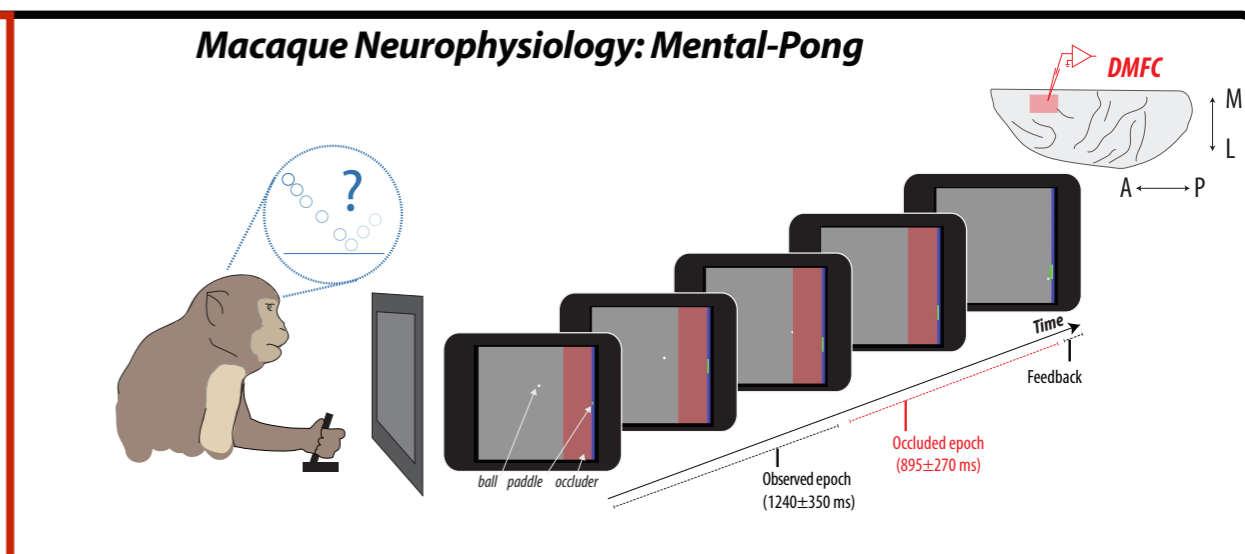
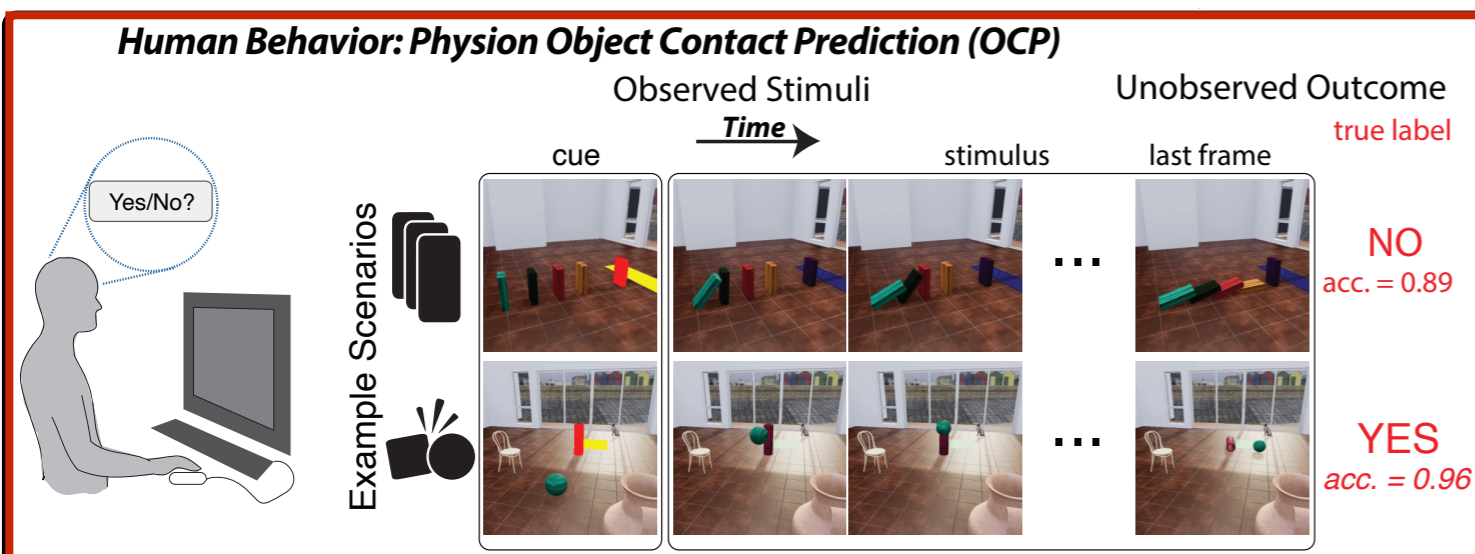
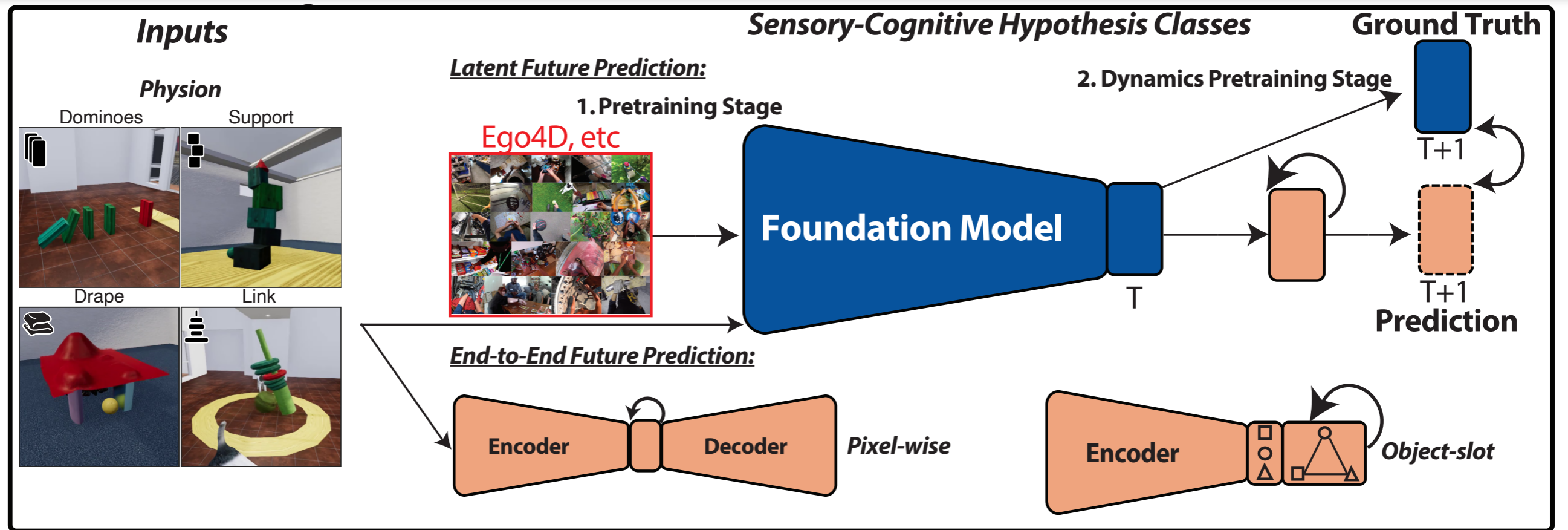
Useful for embodied tasks



Macaque Neurophysiology: Mental Pong



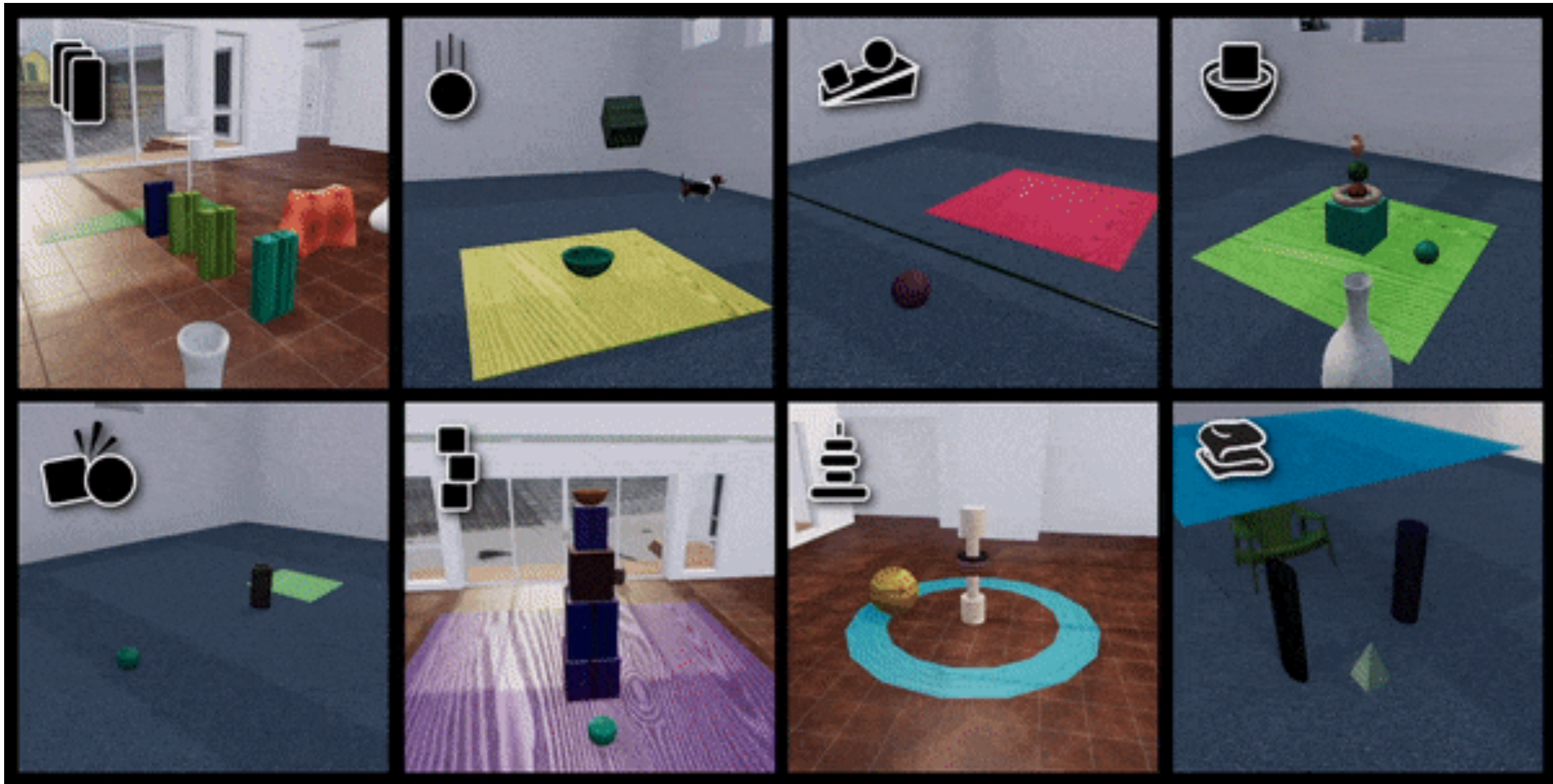
Human Behavior: Object Contact Prediction



Object Contact Prediction Environment

Physion/ThreeD World (TDW)

Bear et al. 2021



Focus on everyday physical understanding



Daniel Bear



Joshua Tenenbaum



Daniel Yamins

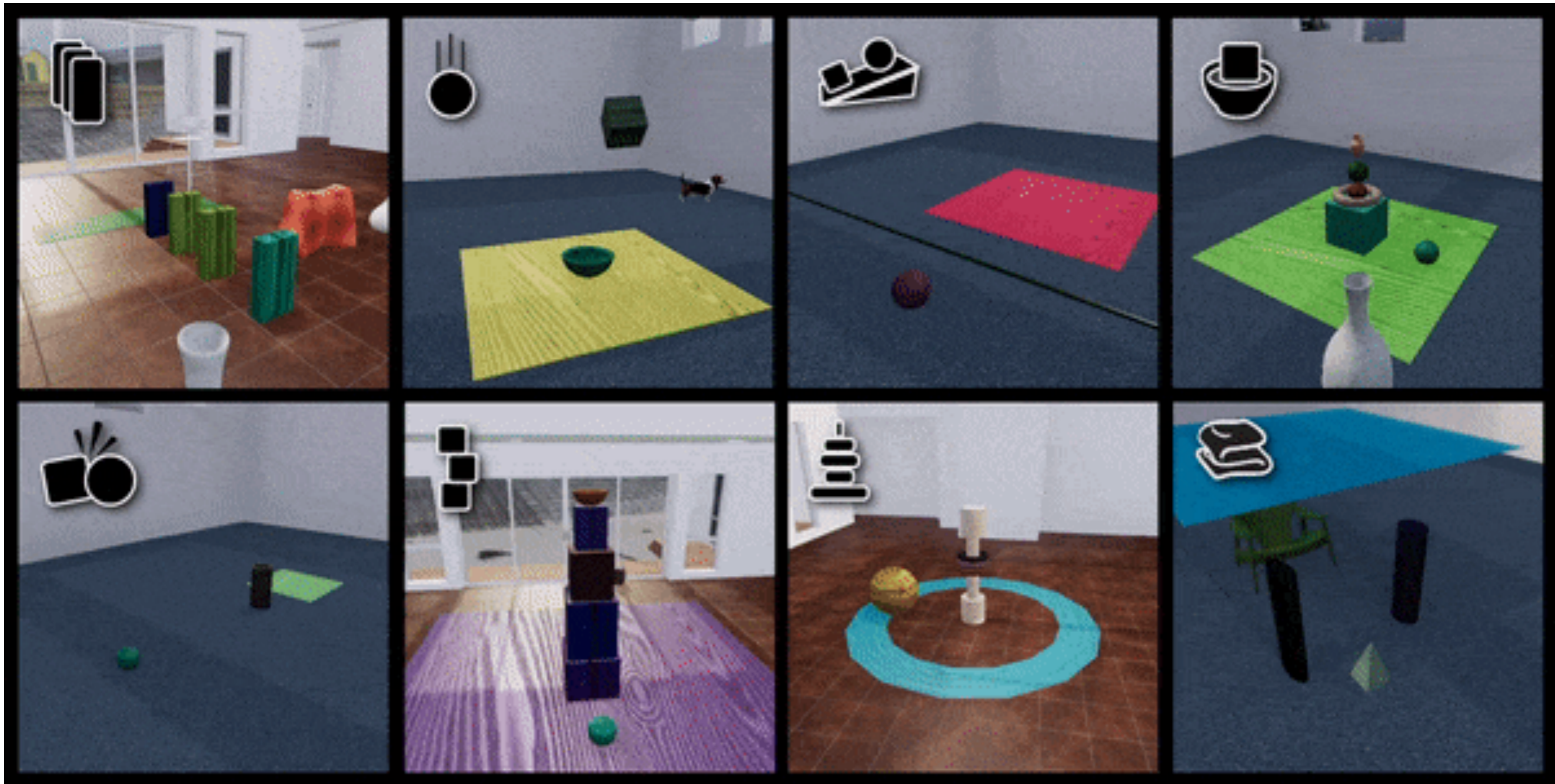


Judith Fan

Object Contact Prediction Environment

Physion/ThreeD World (TDW)

Bear et al. 2021



Focus on everyday physical understanding



Daniel Bear



Joshua Tenenbaum



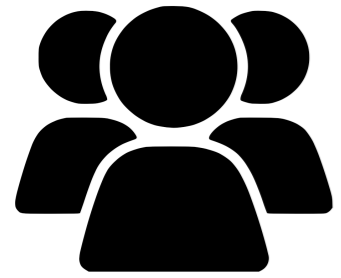
Daniel Yamins



Judith Fan

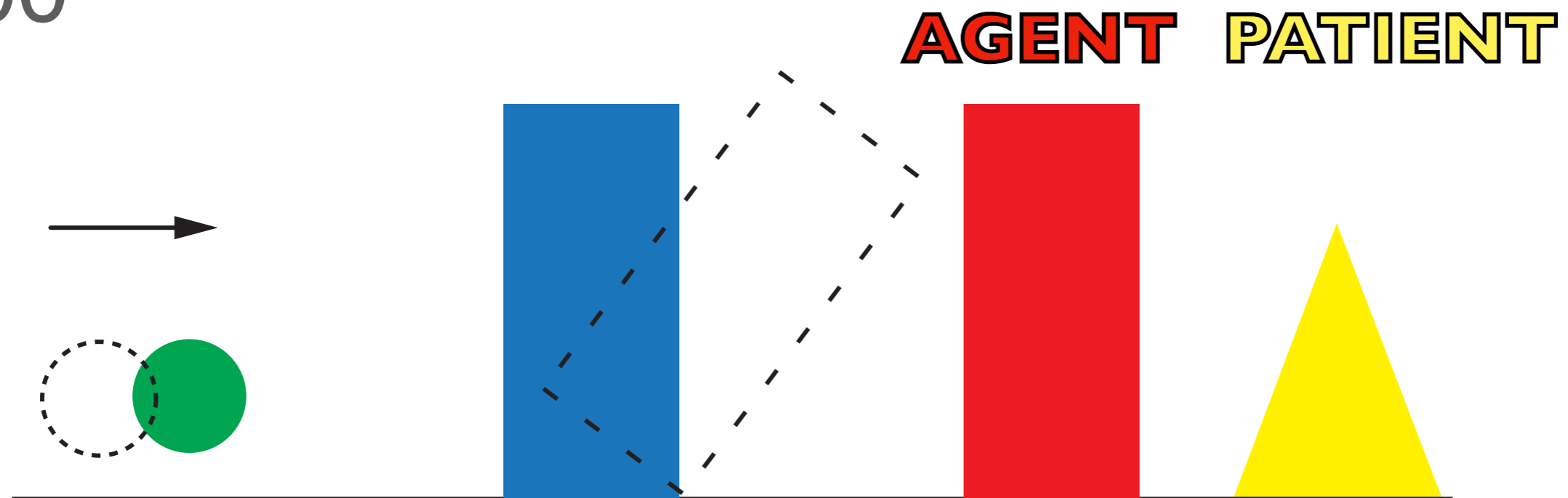
Human Behavior: Object Contact Prediction

Bear et al. 2021



“Will the *agent* object contact the *patient* object?”

n=100



Daniel Bear



Joshua Tenenbaum



Daniel Yamins



Judith Fan

Bear et al. 2021



YES

NO

Is the red object going to hit the yellow area?

Bear et al. 2021



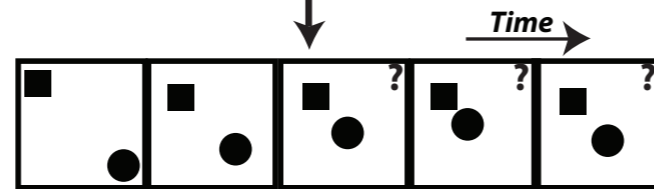
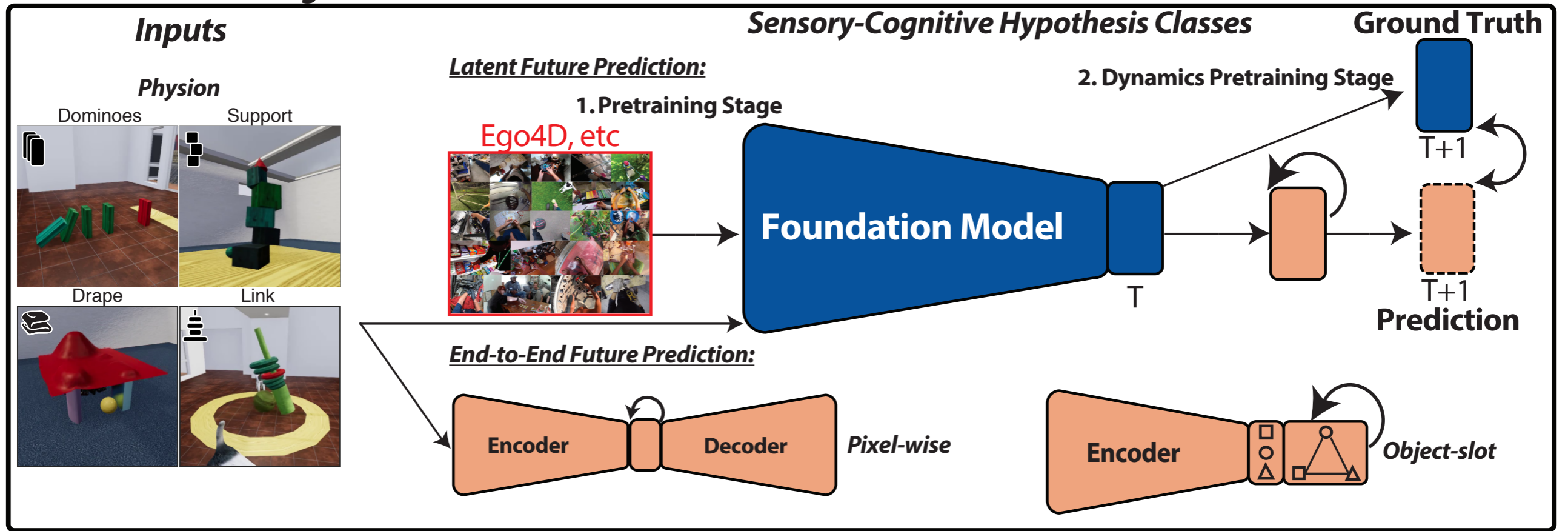
YES

NO

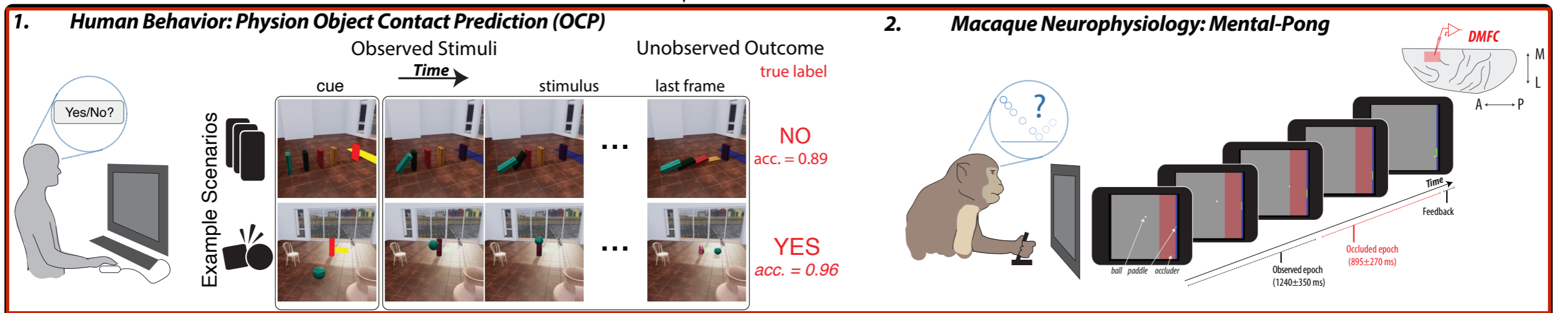
Is the red object going to hit the yellow area?

Model Evaluations: Both Metrics

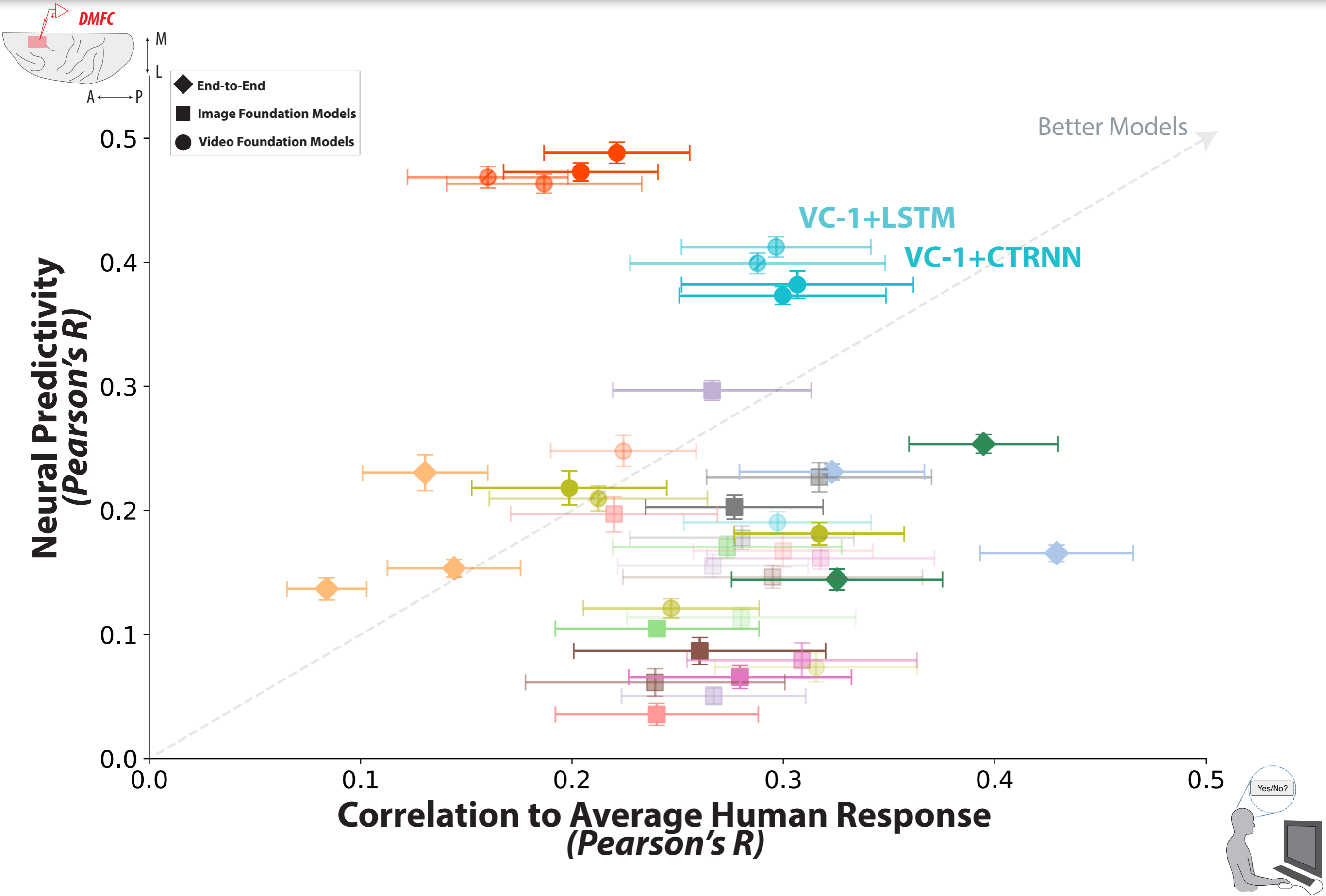
(A) Model Pretraining



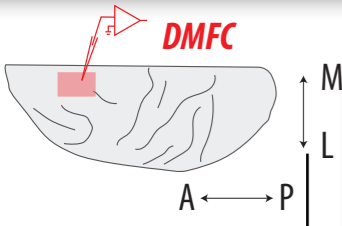
(B) Model Evaluations



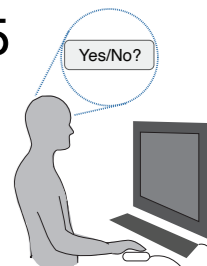
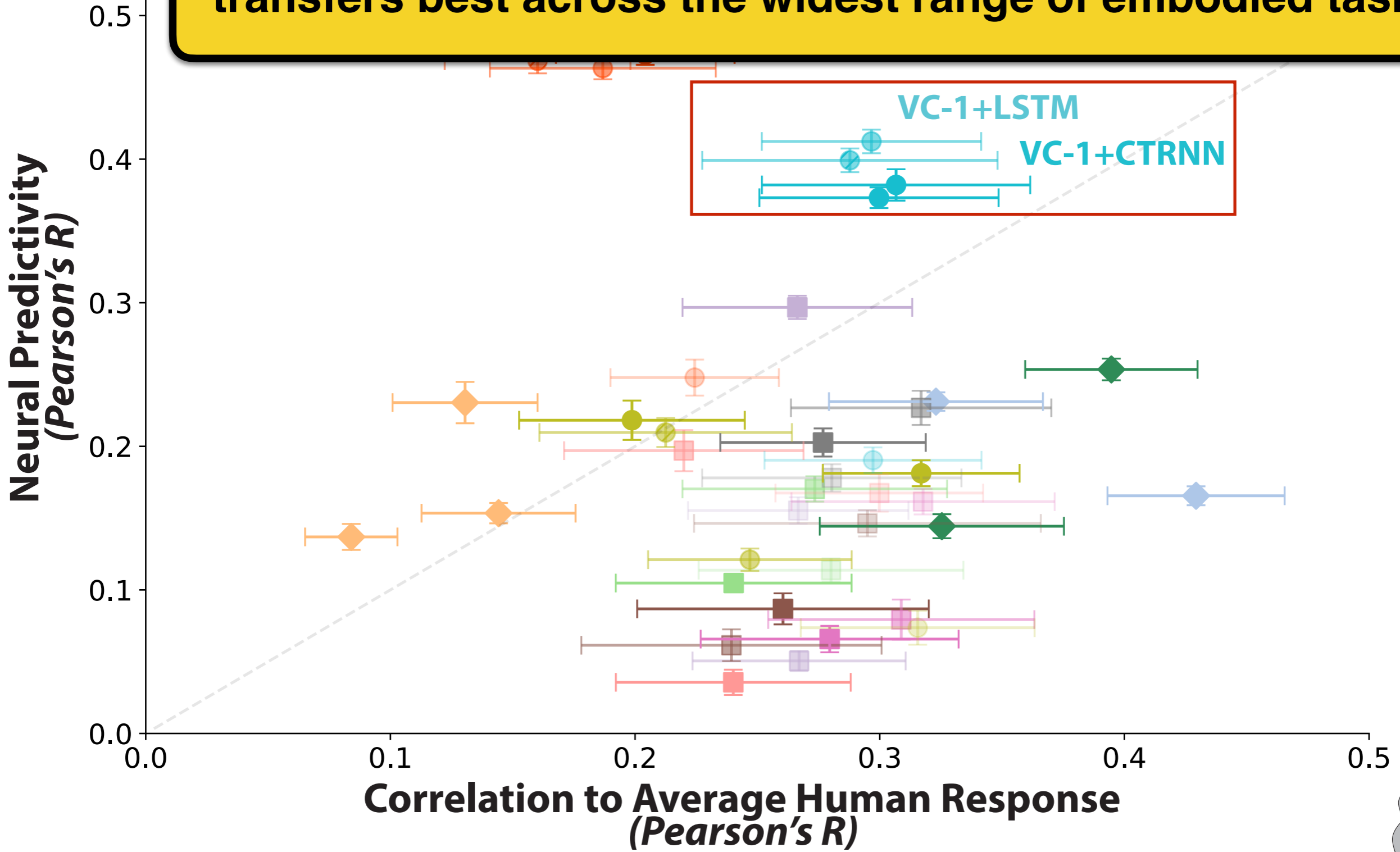
Dynamically-Equipped Video Foundation Models Can Match Both



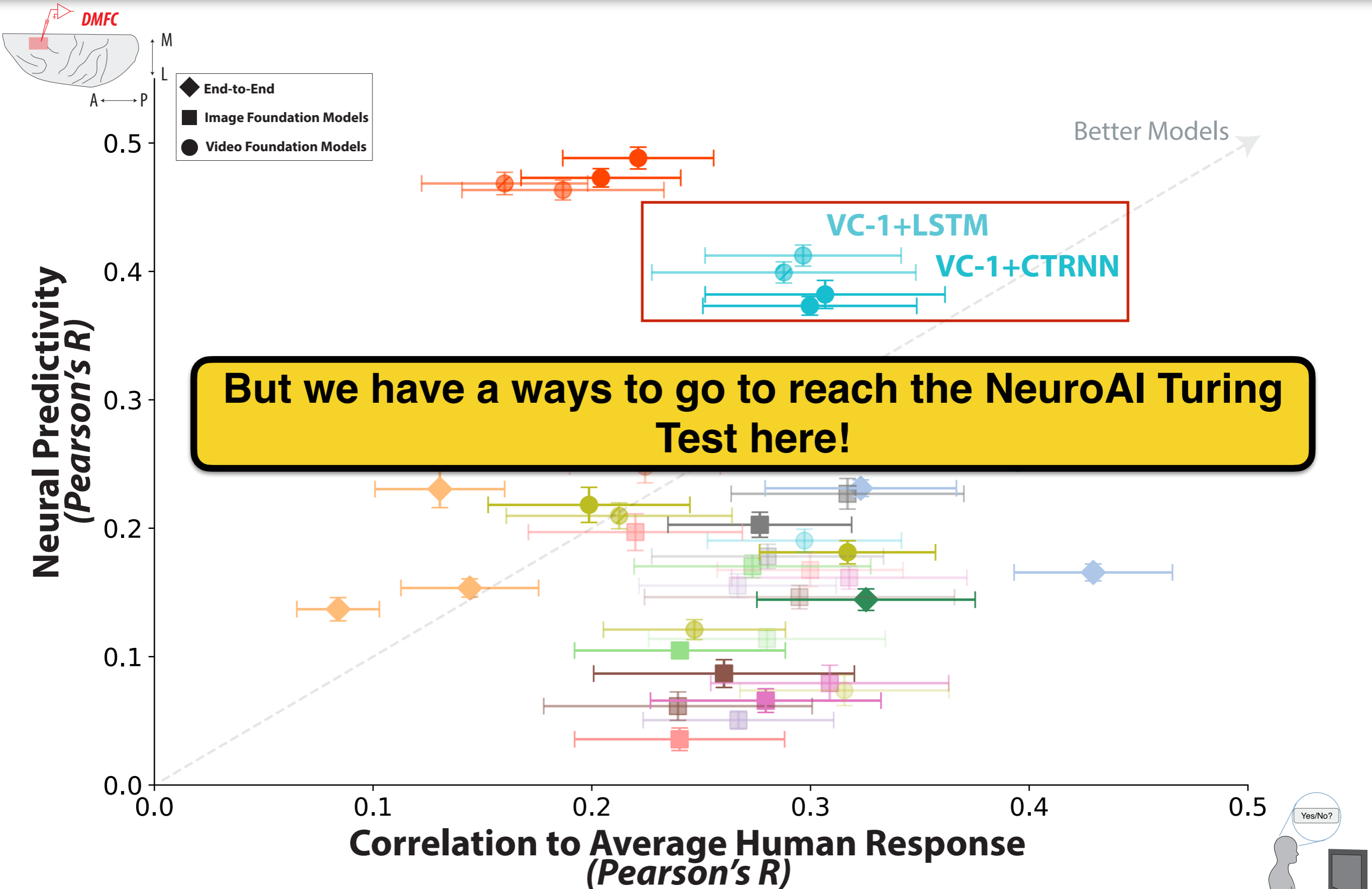
Dynamically-Equipped Video Foundation Models Can Match Both



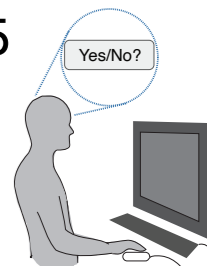
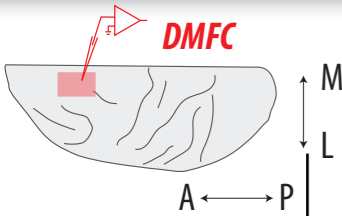
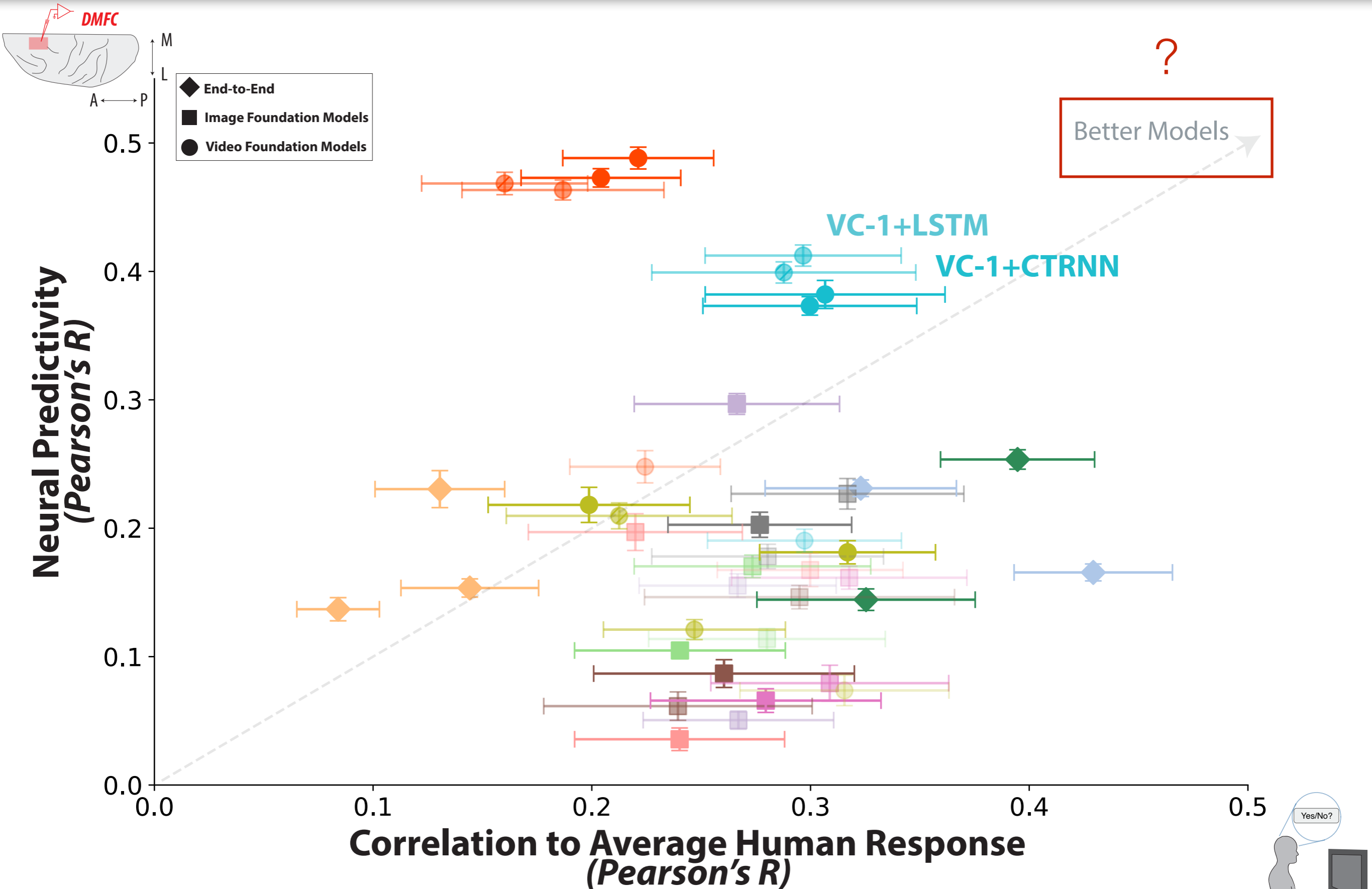
Exposed to the largest variety of egocentric video sources & transfers best across the widest range of embodied tasks.



Dynamically-Equipped Video Foundation Models Can Match Both



Future Directions



Outline

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

Mouse visual cortex (so far) is a low-acuity, shallow network that makes best use of the mouse's limited resources to create a general-purpose visual system, that can be deployed in novel environments and embodied contexts.

- ▶ Reusable Latent Representations for Primate Mental Simulation

- ▶ Heuristics for Interrogating Natural Intelligence

Outline

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

Mouse visual cortex (so far) is a low-acuity, shallow network that makes best use of the mouse's limited resources to create a general-purpose visual system, that can be deployed in novel environments and embodied contexts.

▶ Reusable Latent Representations for Primate Mental Simulation

Mental simulation crucially relies on explicit future prediction of a “factorized description” of visual scenes, where this “factorized description” is strongly constrained and must enable a wide range of dynamic embodied abilities.

▶ Heuristics for Interrogating Natural Intelligence

Outline

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

Mouse visual cortex (so far) is a low-acuity, shallow network that makes best use of the mouse's limited resources to create a general-purpose visual system, that can be deployed in novel environments and embodied contexts.

▶ Reusable Latent Representations for Primate Mental Simulation

Mental simulation crucially relies on explicit future prediction of a “factorized description” of visual scenes, where this “factorized description” is strongly constrained and must enable a wide range of dynamic embodied abilities.

▶ Heuristics for Interrogating Natural Intelligence

Heuristics for Interrogating Natural Intelligence

- **Incorporating Neuroscience Insights:**

- **Incorporating AI Insights:**

Heuristics for Interrogating Natural Intelligence

- **Incorporating Neuroscience Insights:**
- **Connectomics:**
- **Ethology:**
- **Incorporating AI Insights:**

Heuristics for Interrogating Natural Intelligence

- **Incorporating Neuroscience Insights:**

- **Connectomics:** Not usually a 1-1 mapping from a connectome to a functional model, and easy to get wrong. Rather, the best model often requires an *iterative balance* of functional optimization with macroscale structural constraints (e.g. shallow vs. deep cortex).

- **Ethology:**

- **Incorporating AI Insights:**

Heuristics for Interrogating Natural Intelligence

- **Incorporating Neuroscience Insights:**

- **Connectomics:** Not usually a 1-1 mapping from a connectome to a functional model, and easy to get wrong. Rather, the best model often requires an *iterative balance* of functional optimization with macroscale structural constraints (e.g. shallow vs. deep cortex).
- **Ethology:** Ethology alone does not always give us the “correct” functional optimization target (e.g. for intermediate brain areas). Can be helpful to consider that the brain area is embodied in a larger cognitive agent.

- **Incorporating AI Insights:**

Heuristics for Interrogating Natural Intelligence

- **Incorporating Neuroscience Insights:**

- **Connectomics:** Not usually a 1-1 mapping from a connectome to a functional model, and easy to get wrong. Rather, the best model often requires an *iterative balance* of functional optimization with macroscale structural constraints (e.g. shallow vs. deep cortex).
- **Ethology:** Ethology alone does not always give us the “correct” functional optimization target (e.g. for intermediate brain areas). Can be helpful to consider that the brain area is embodied in a larger cognitive agent.

- **Incorporating AI Insights:**

- End-to-end reinforcement learning (RL) does *not* seem to give us neurally-aligned visual systems in both rodents and primates.

Heuristics for Interrogating Natural Intelligence

- **Incorporating Neuroscience Insights:**

- **Connectomics:** Not usually a 1-1 mapping from a connectome to a functional model, and easy to get wrong. Rather, the best model often requires an *iterative balance* of functional optimization with macroscale structural constraints (e.g. shallow vs. deep cortex).
- **Ethology:** Ethology alone does not always give us the “correct” functional optimization target (e.g. for intermediate brain areas). Can be helpful to consider that the brain area is embodied in a larger cognitive agent.

- **Incorporating AI Insights:**

- End-to-end reinforcement learning (RL) does *not* seem to give us neurally-aligned visual systems in both rodents and primates.
- Suggests a possible functional modularization of optimization targets, with reusable SSL representations best matching visual areas overall.

Heuristics for Interrogating Natural Intelligence

- **Incorporating Neuroscience Insights:**

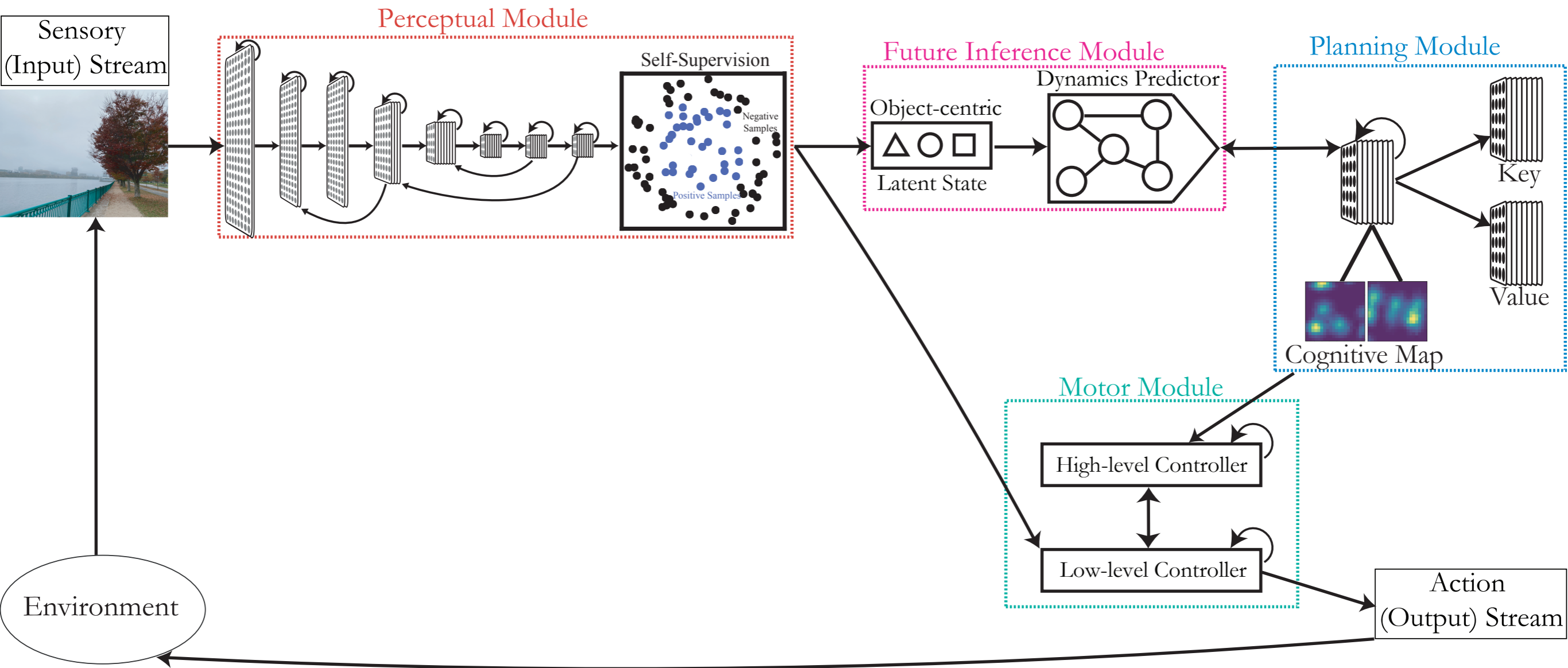
- **Connectomics:** Not usually a 1-1 mapping from a connectome to a functional model, and easy to get wrong. Rather, the best model often requires an *iterative balance* of functional optimization with macroscale structural constraints (e.g. shallow vs. deep cortex).
- **Ethology:** Ethology alone does not always give us the “correct” functional optimization target (e.g. for intermediate brain areas). Can be helpful to consider that the brain area is **embodied** in a **larger cognitive agent**.

- **Incorporating AI Insights:**

- End-to-end reinforcement learning (RL) does *not* seem to give us neurally-aligned visual systems in both rodents and primates.
- Suggests a possible functional **modularization** of optimization targets, with reusable SSL representations best matching visual areas overall.

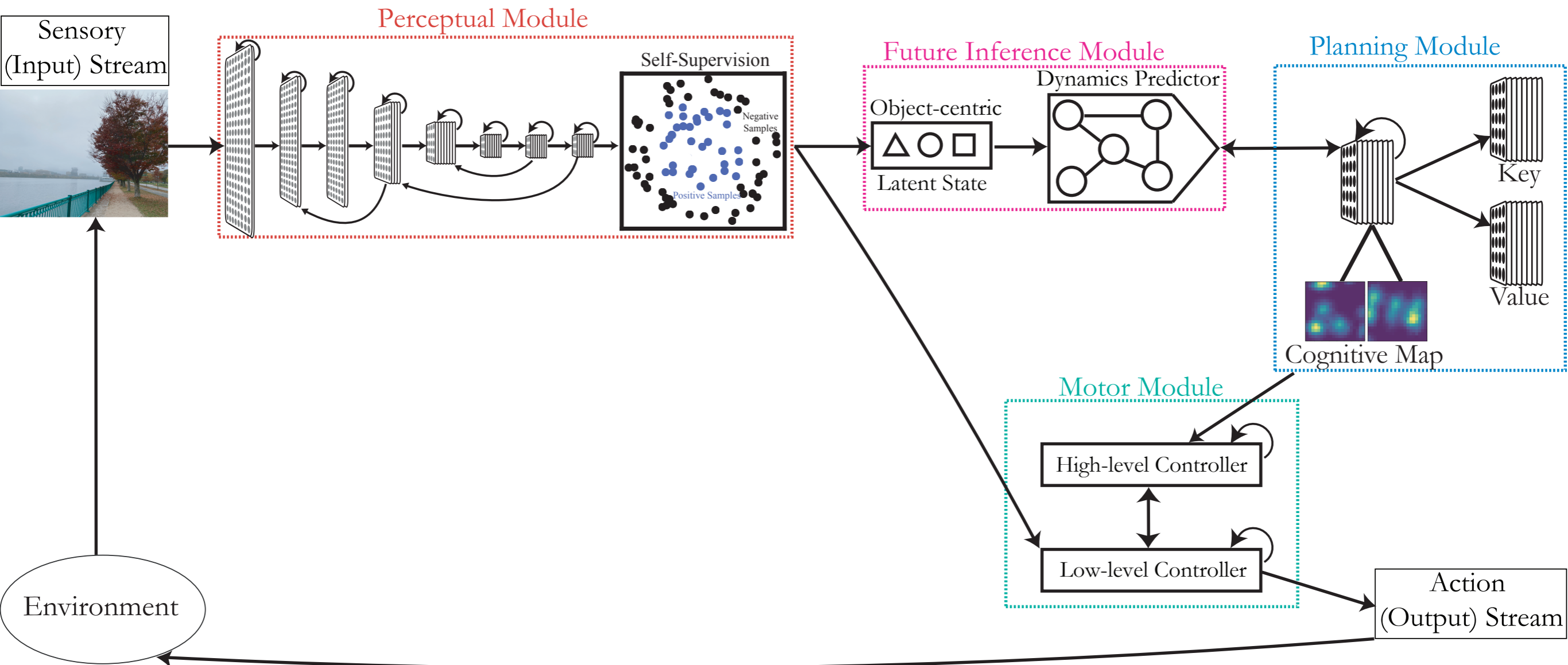
Next Steps: Modularized, Embodied Agents?

Next Steps: Modularized, Embodied Agents?



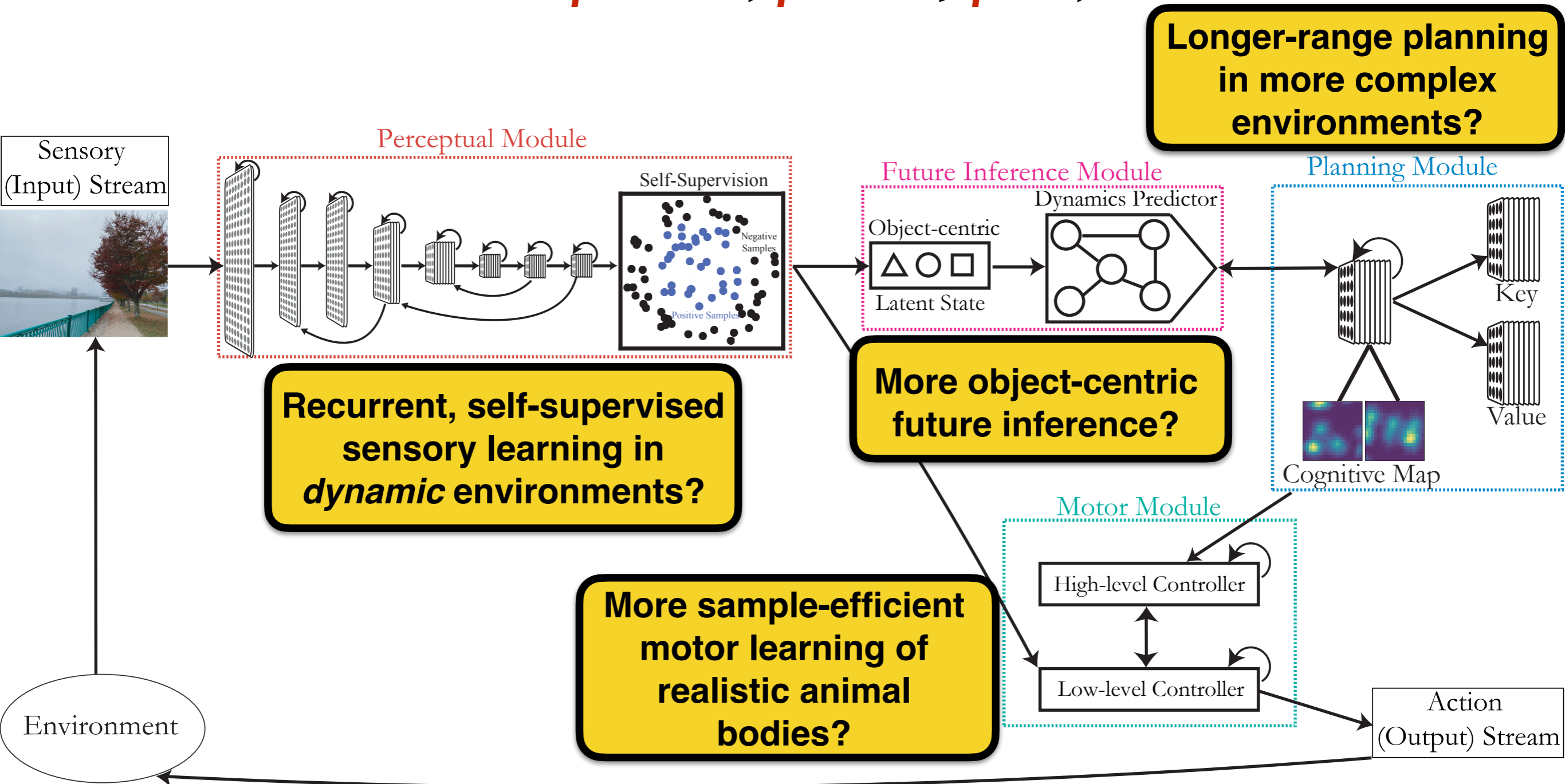
Next Steps: Modularized, Embodied Agents?

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



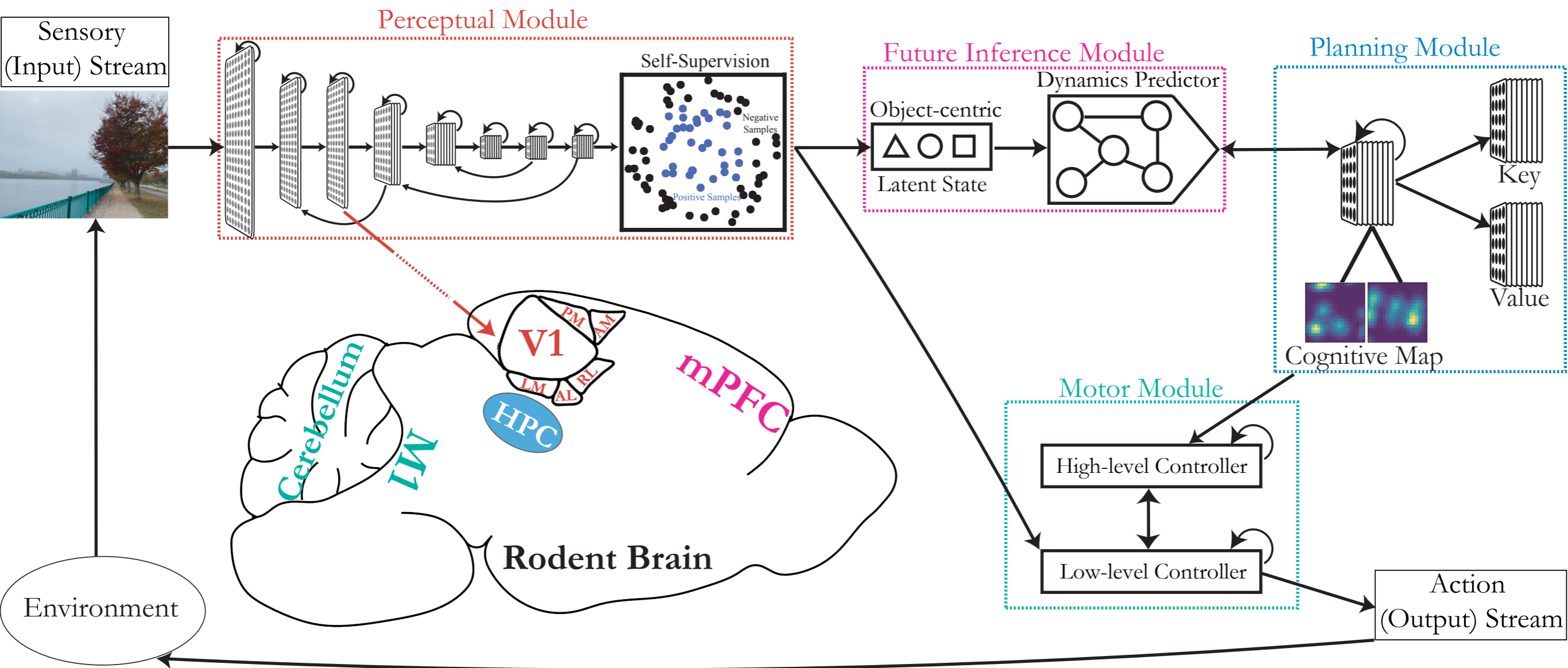
Next Steps: Modularized, Embodied Agents?

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



Next Steps: Modularized, Embodied Agents?

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



What Happens When We Get There?

What Happens When We Get There?

Barriers and Pathways to Human-AI Alignment: A Game-Theoretic Approach

ARAN NAYEBI, Carnegie Mellon University, USA

<https://arxiv.org/abs/2502.05934>

What Happens When We Get There?

Barriers and Pathways to Human-AI Alignment: A Game-Theoretic Approach

ARAN NAYEBI, Carnegie Mellon University, USA

<https://arxiv.org/abs/2502.05934>

PROPOSITION 2.6 (LOWER BOUND). *There exist functions f_j , input sets S_j , and prior distributions $\{\mathbb{P}_j^i\}^{i \in [N]}$ for all $j \in [M]$, such that any protocol among N agents needs to exchange $\Omega(M N^2 \log(1/\varepsilon))$ bits⁶ to achieve $\langle M, N, \varepsilon, \delta \rangle$ -agreement on $\{f_j\}_{j \in [M]}$, for ε bounded below by $\min_{j \in [M]} \varepsilon_j$.*

What Happens When We Get There?

Barriers and Pathways to Human-AI Alignment: A Game-Theoretic Approach

ARAN NAYEBI, Carnegie Mellon University, USA

<https://arxiv.org/abs/2502.05934>

PROPOSITION 2.6 (LOWER BOUND). *There exist functions f_j , input sets S_j , and prior distributions $\{\mathbb{P}_j^i\}_{i \in [N]}$ for all $j \in [M]$, such that any protocol among N agents needs to exchange $\Omega(M N^2 \log(1/\varepsilon))$ bits⁶ to achieve $\langle M, N, \varepsilon, \delta \rangle$ -agreement on $\{f_j\}_{j \in [M]}$, for ε bounded below by $\min_{j \in [M]} \varepsilon_j$.*

There will always be a task that is *impossible* to align efficiently over, even if the agents are computationally unbounded

What Happens When We Get There?

Barriers and Pathways to Human-AI Alignment: A Game-Theoretic Approach

ARAN NAYEBI, Carnegie Mellon University, USA

<https://arxiv.org/abs/2502.05934>

PROPOSITION 2.6 (LOWER BOUND). *There exist functions f_j , input sets S_j , and prior distributions $\{\mathbb{P}_j^i\}^{i \in [N]}$ for all $j \in [M]$, such that any protocol among N agents needs to exchange $\Omega(MN^2 \log(1/\varepsilon))$ bits⁶ to achieve $\langle M, N, \varepsilon, \delta \rangle$ -agreement on $\{f_j\}_{j \in [M]}$, for ε bounded below by $\min_{j \in [M]} \varepsilon_j$.*

There will always be a task that is *impossible* to align efficiently over, even if the agents are computationally unbounded

$$O\left(\frac{(1100)^{\frac{2097152}{(1/4)^6}}}{(1/2)^{\frac{256}{(1/4)^2}}}\right) = O(1.31 \times 10^{26125365467})$$

If the agents are computationally *bounded*, this can currently take more time than the number of atoms in the universe! ($\sim 4.8 \times 10^{79}$)

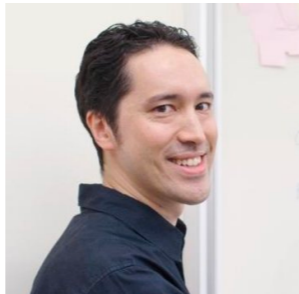
Acknowledgements



Nathan C.L. Kong



Chengxu Zhuang



Justin L. Gardner



Anthony M. Norcia



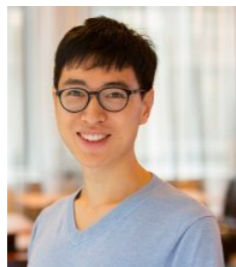
Daniel Yamins



Rishi Rajalingham

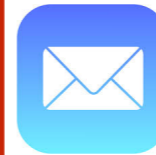


Mehrdad Jazayeri



Guangyu Robert Yang

Contact:



anayebi@cs.cmu.edu



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



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



Funding:

Burroughs Wellcome Fund CASI Award

K. Lisa Yang ICoN Postdoctoral Fellowship,
McGovern Institute, MIT

Stanford Neurosciences PhD Program

Stanford Mind, Brain, Computation and
Technology Training Program,
Wu Tsai Neurosciences Institute