

We probably need a new idea...or two?

---

**Aran Nayebi**

*K. Lisa Yang ICoN Postdoctoral Fellow  
McGovern Institute, MIT*

**CCN 2023 GAC**

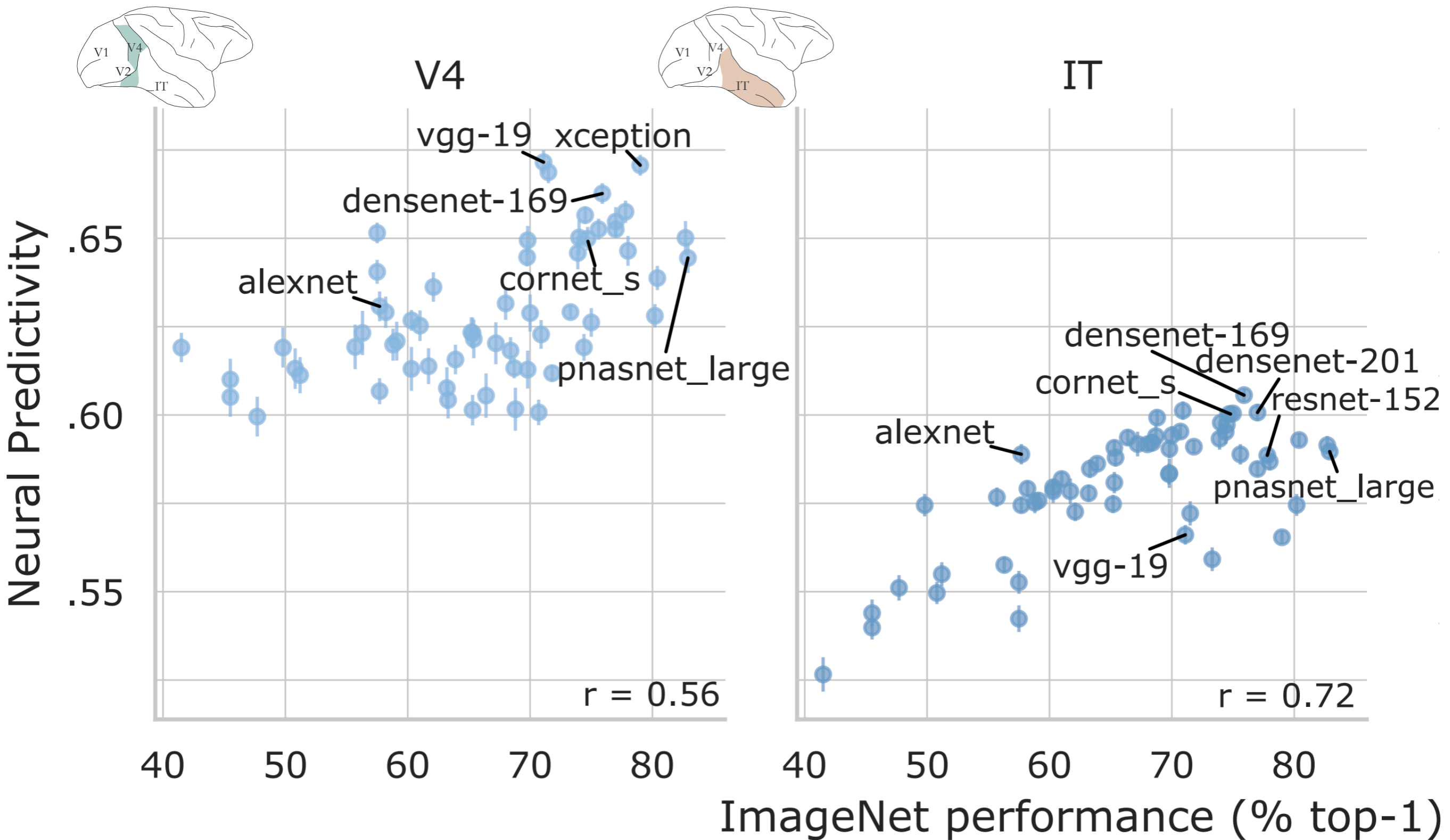
*2023.08.25*



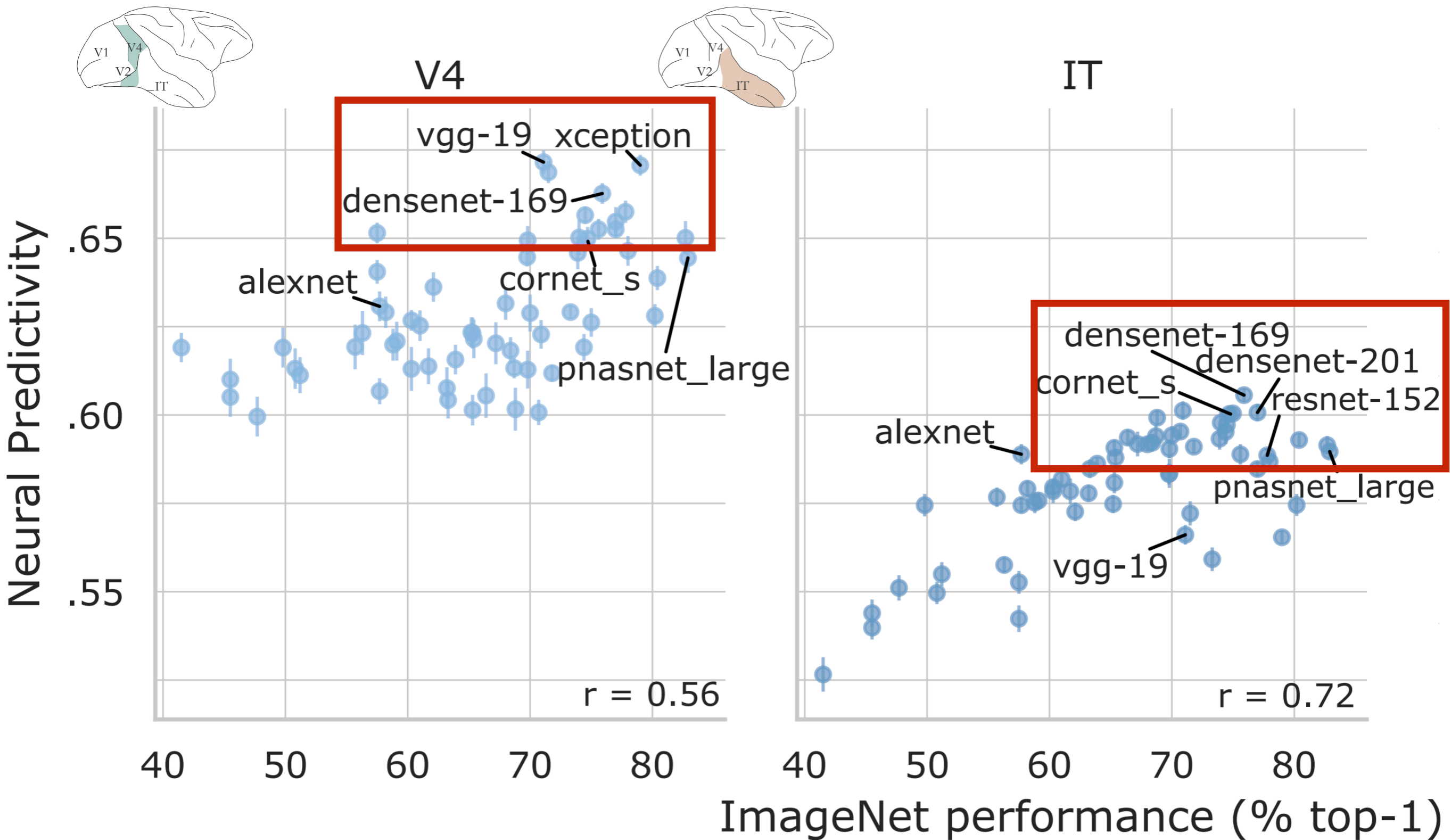
**K. LISA YANG**  
**ICoN CENTER**  
Integrative Computational Neuroscience



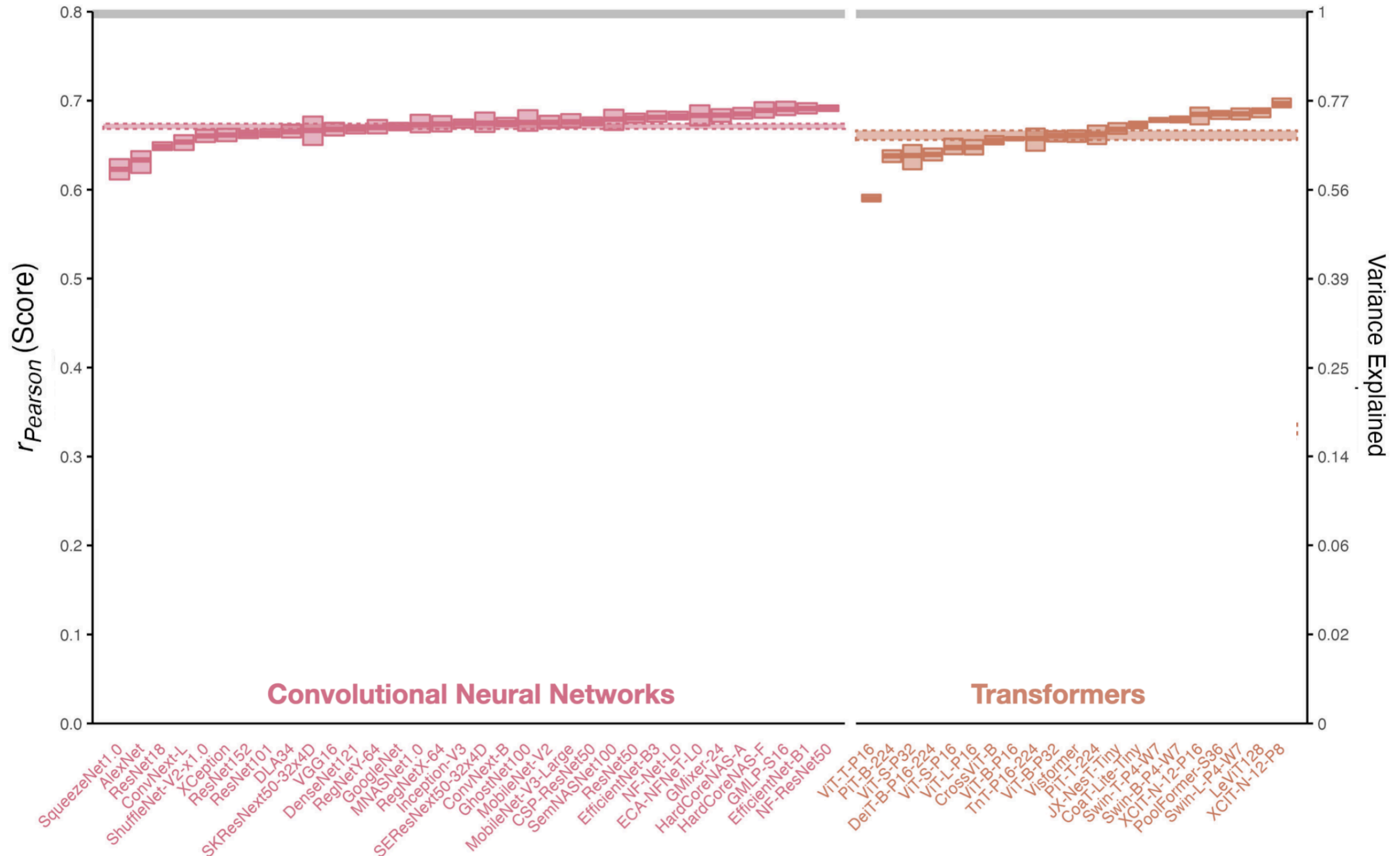
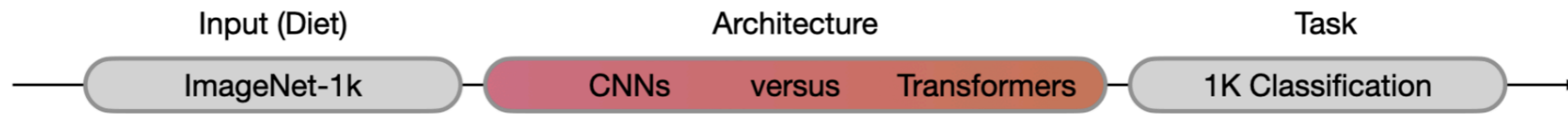
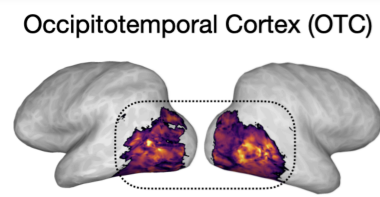
# Similar predictivities among very different CNN architectures



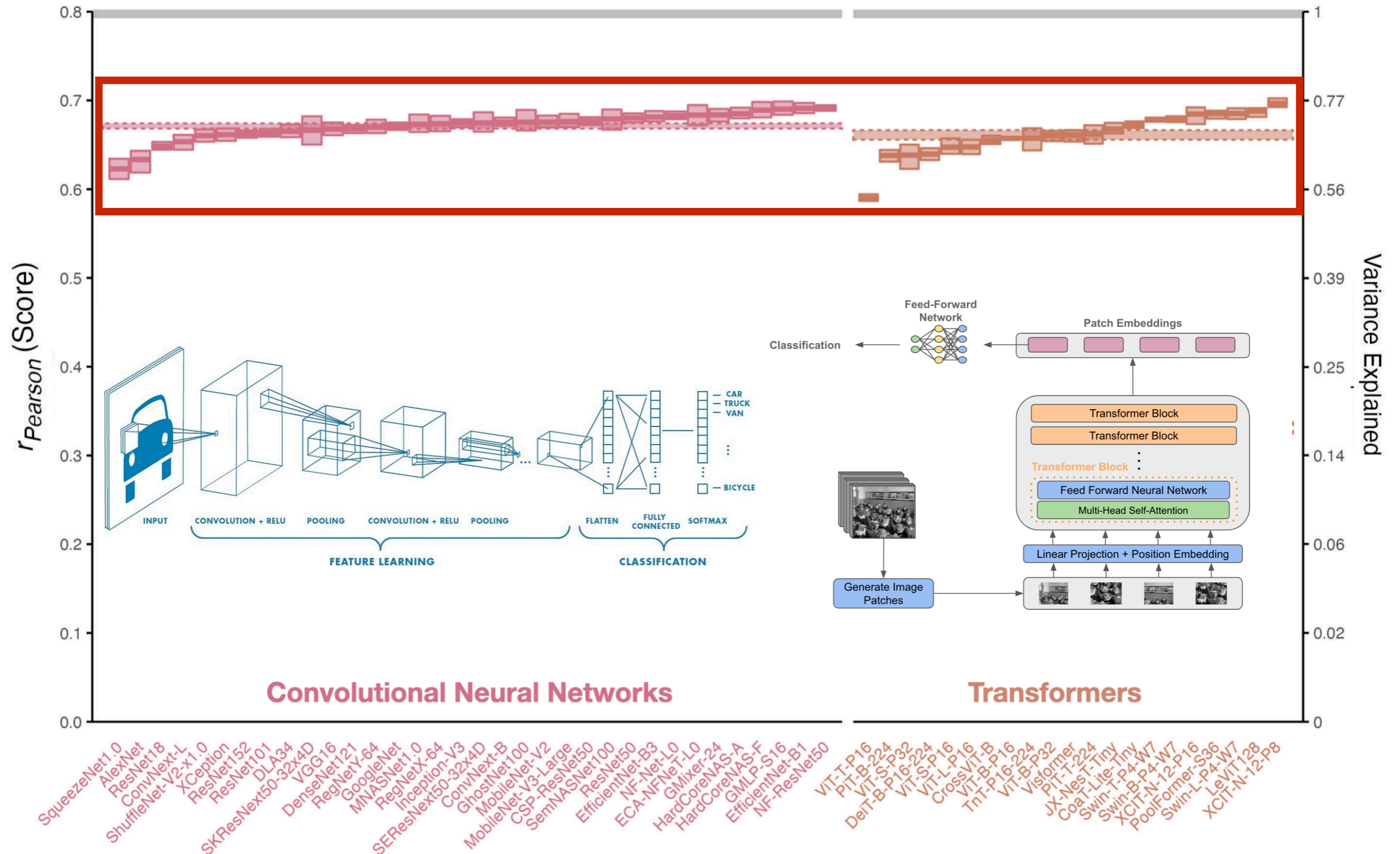
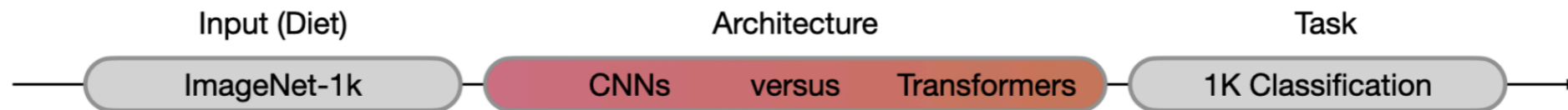
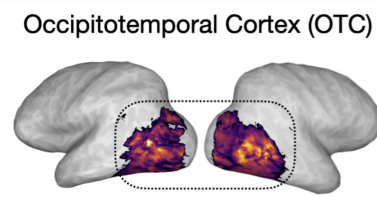
# Similar predictivities among very different CNN architectures



# Similar predictivities between CNNs vs. Transformers



# Similar predictivities between CNNs vs. Transformers



Why might this lack of distinction be?

Why might this lack of distinction be?

- **Metrics**
- **Brain Comparison Data**
- **Models**

# Why might this lack of distinction be?

- **Metrics**
  - Correlation Based
  - Causal Manipulations
  - Mapping Transform Class
- **Brain Comparison Data**
  - Organism
  - Cognitive Tasks
  - Brain Area
  - Number of Units
  - Number of Conditions
- **Models**
  - Architecture
  - Loss Function
  - Dataset



# Changing the Brain Comparison Data

- **Metrics**
  - *Correlation Based* (Regression)
  - Causal Manipulations
  - Mapping Transform Class
- **Brain Comparison Data**
  - *Organism* (Primates)
  - Cognitive Tasks
  - Brain Area
  - Number of Units
  - Number of Conditions
- **Models**
  - Architecture
  - Loss Function
  - Dataset

# Changing the Brain Comparison Data

- **Metrics**
  - *Correlation Based* (Regression)
  - Causal Manipulations
  - Mapping Transform Class

**Modify** →

- **Brain Comparison Data**
  - *Organism* (Primates)
  - Cognitive Tasks
  - Brain Area
  - Number of Units
  - Number of Conditions

- **Models**
  - Architecture
  - Loss Function
  - Dataset

# Changing the Brain Comparison Data

- **Metrics**
  - *Correlation Based* (Regression)
  - Causal Manipulations
  - Mapping Transform Class

**Modify** →

- **Brain Comparison Data**
  - *Organism* (Primates)
  - *Cognitive Tasks* (“Mental-Pong”)
  - Brain Area
  - Number of Units
  - Number of Conditions

- **Models**
  - Architecture
  - Loss Function
  - Dataset

# Changing the Brain Comparison Data

- **Metrics**
  - *Correlation Based* (Regression)
  - Causal Manipulations
  - Mapping Transform Class

**Modify** →

- **Brain Comparison Data**
  - *Organism* (Primates)
  - *Cognitive Tasks* (“Mental-Pong”)
  - *Brain Area* (frontal cortex)
  - Number of Units
  - Number of Conditions

- **Models**
  - Architecture
  - Loss Function
  - Dataset

# New Models Needed

- **Metrics**
  - *Correlation Based* (Regression)
  - Causal Manipulations
  - Mapping Transform Class

**Modify**



- **Brain Comparison Data**
  - *Organism* (Primates)
  - *Cognitive Tasks* (“Mental-Pong”)
  - *Brain Area* (frontal cortex)
  - Number of Units
  - Number of Conditions

**New ideas  
needed**

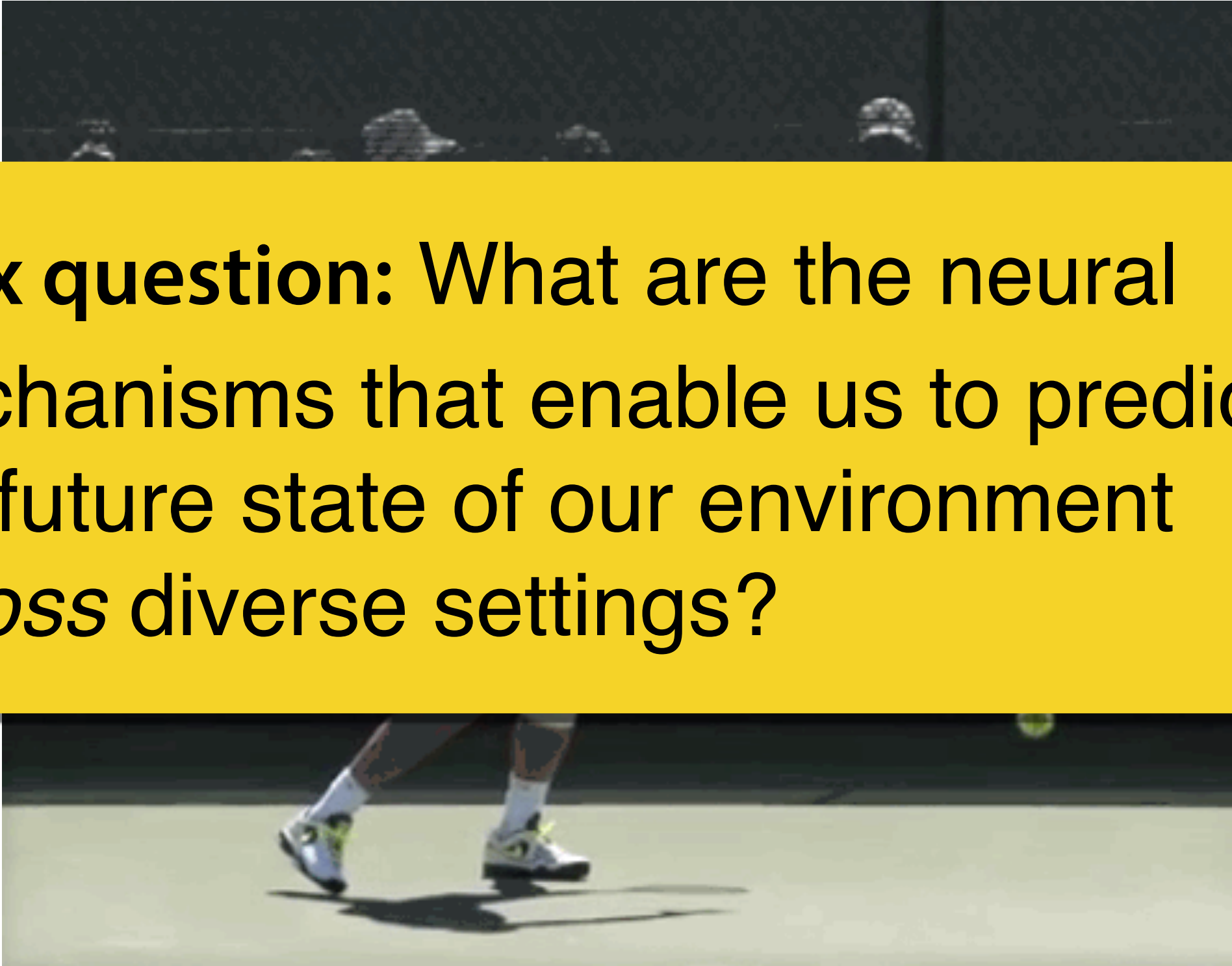


- **Models**
  - Architecture
  - Loss Function
  - Dataset

We do a lot more than passive viewing...

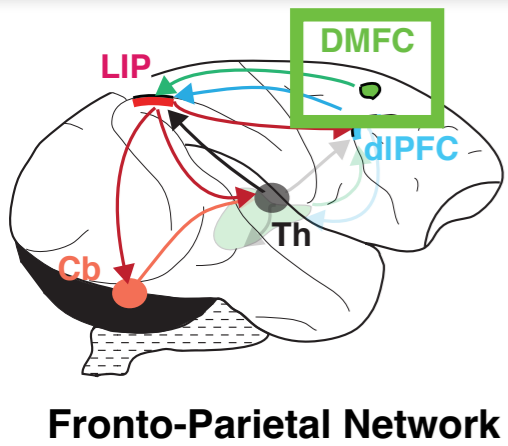


We do a lot more than passive viewing...

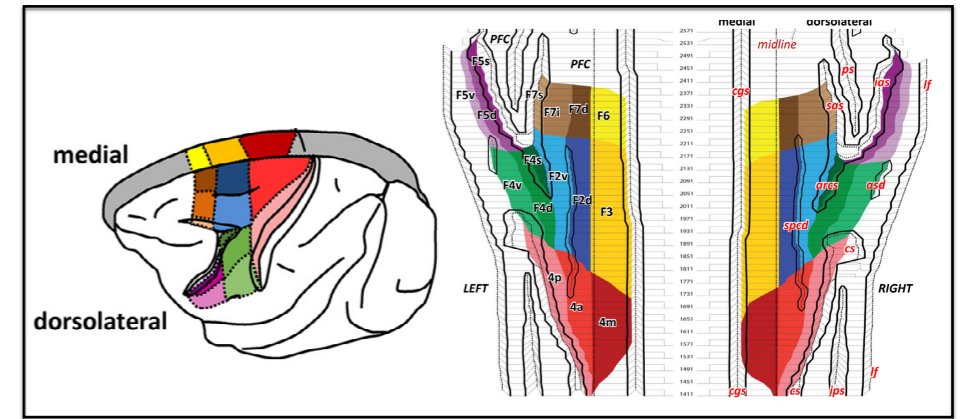
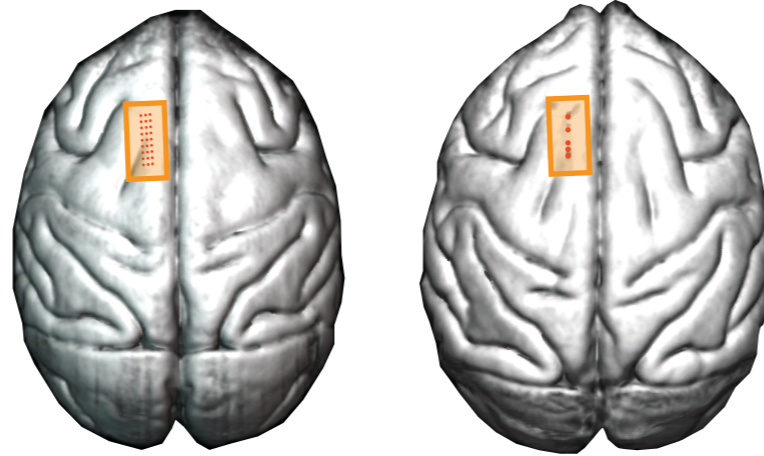
A photograph of a tennis player's legs and feet in motion on a court. The player is wearing white socks and white tennis shoes with yellow accents. The background is dark, suggesting an indoor or night-time setting. The image is partially obscured by a yellow text box.

**Crux question:** What are the neural mechanisms that enable us to predict the future state of our environment *across* diverse settings?

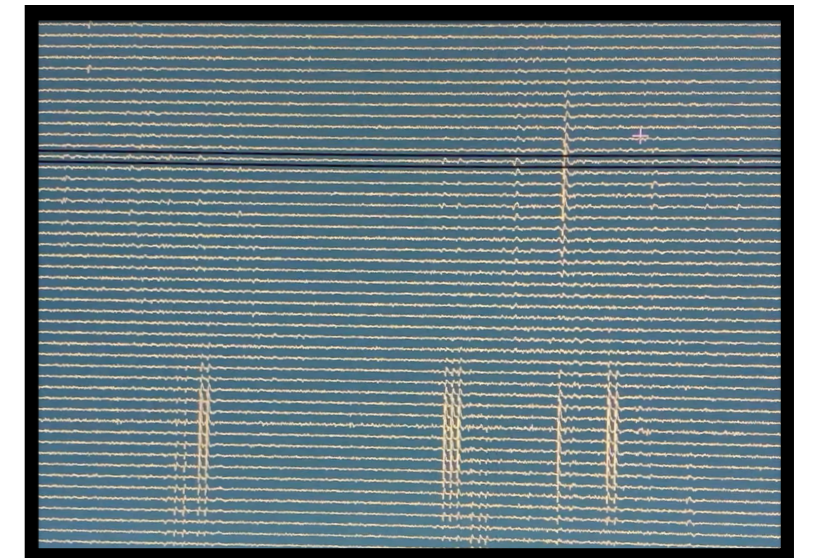
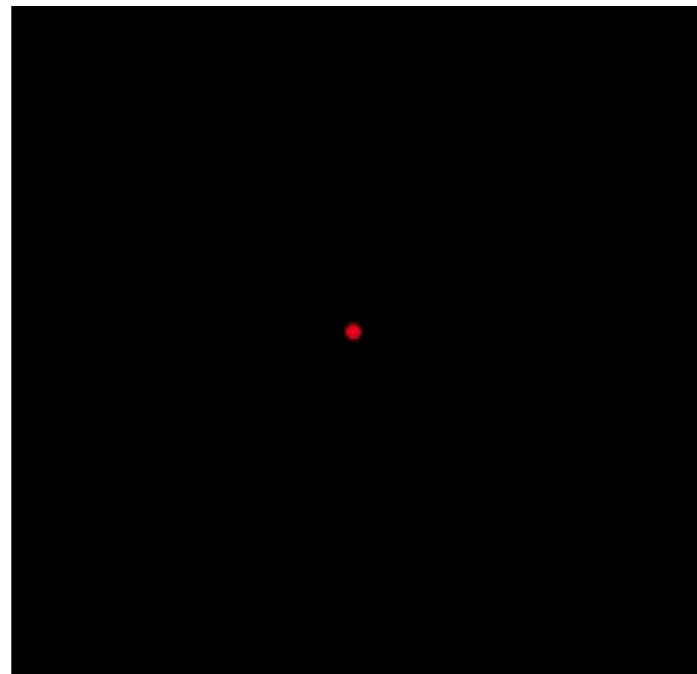
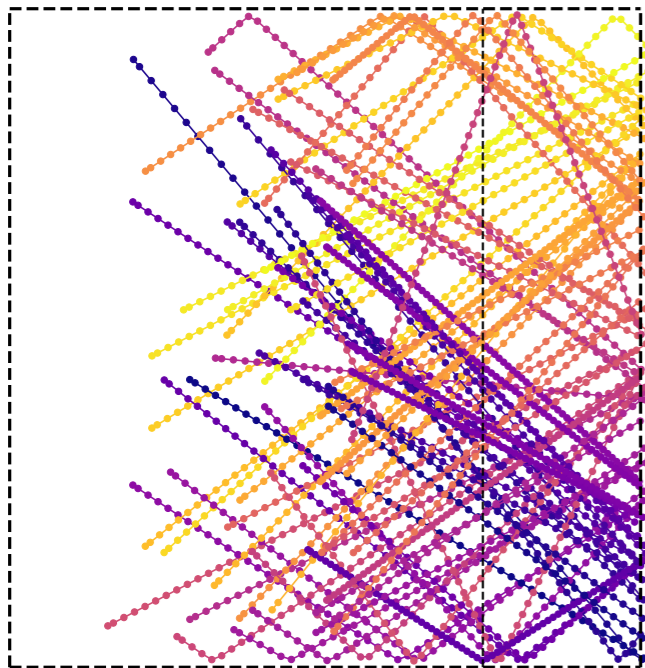
# Mental-Pong Task & Macaque Neurophysiology



*Dorsomedial frontal cortex (DMFC)*



**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*

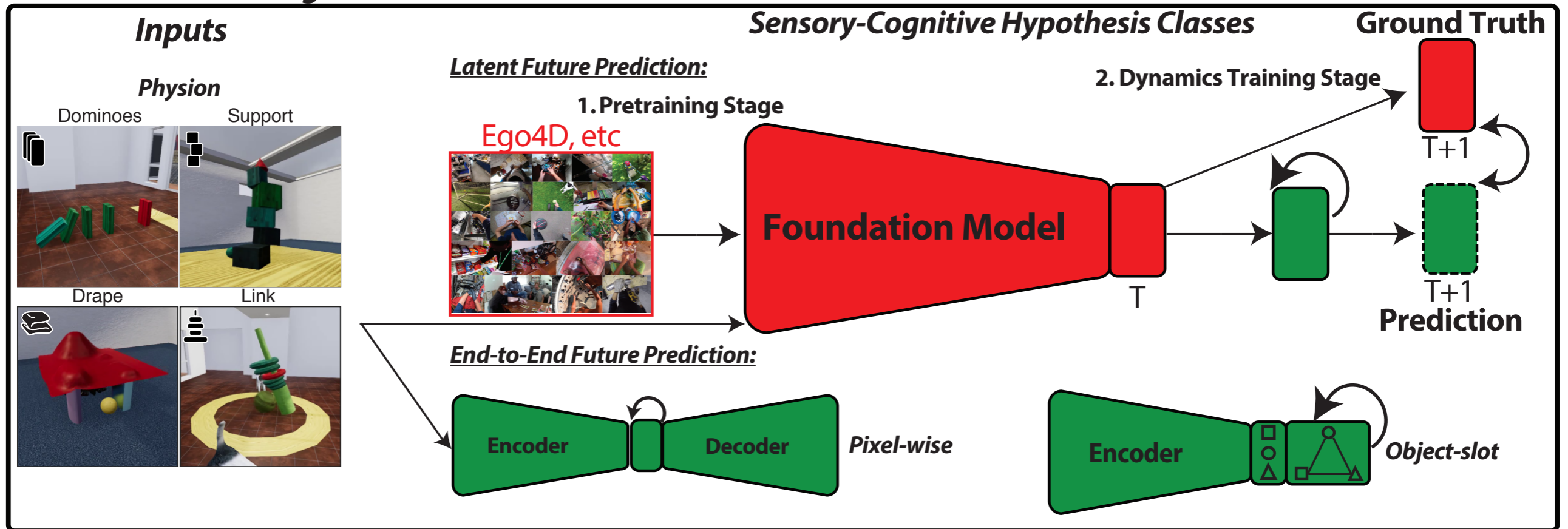


Rishi Rajalingham



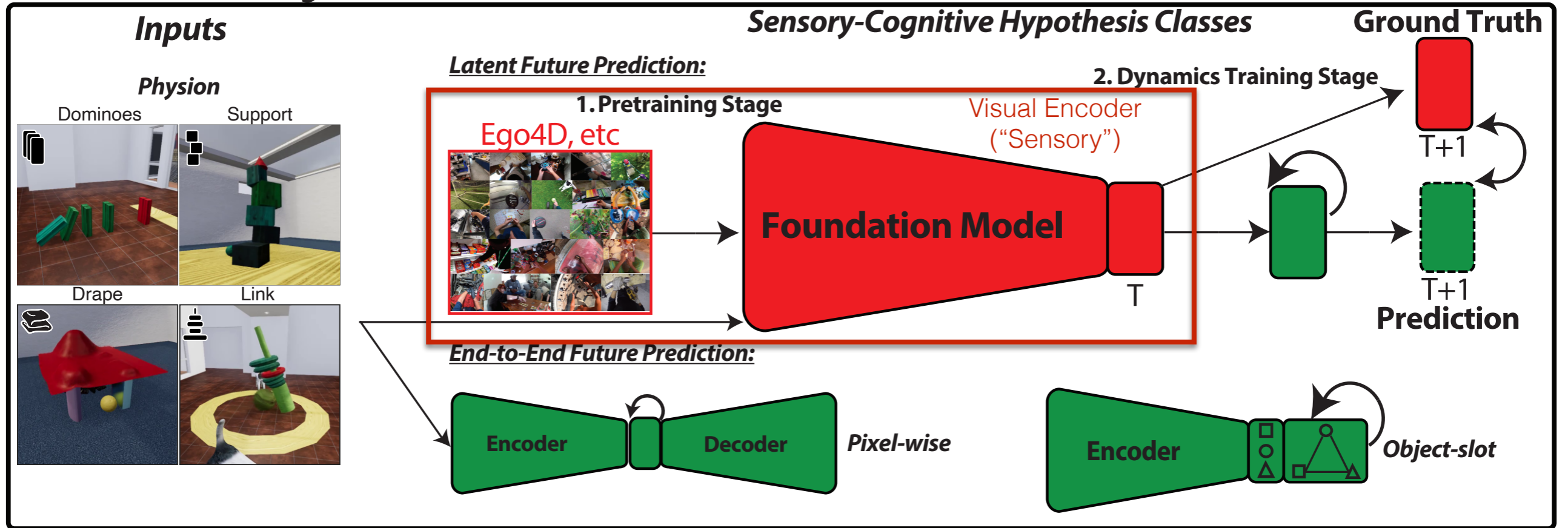
# Overall Approach

## (A) Model Pretraining



# Overall Approach: Foundation Models

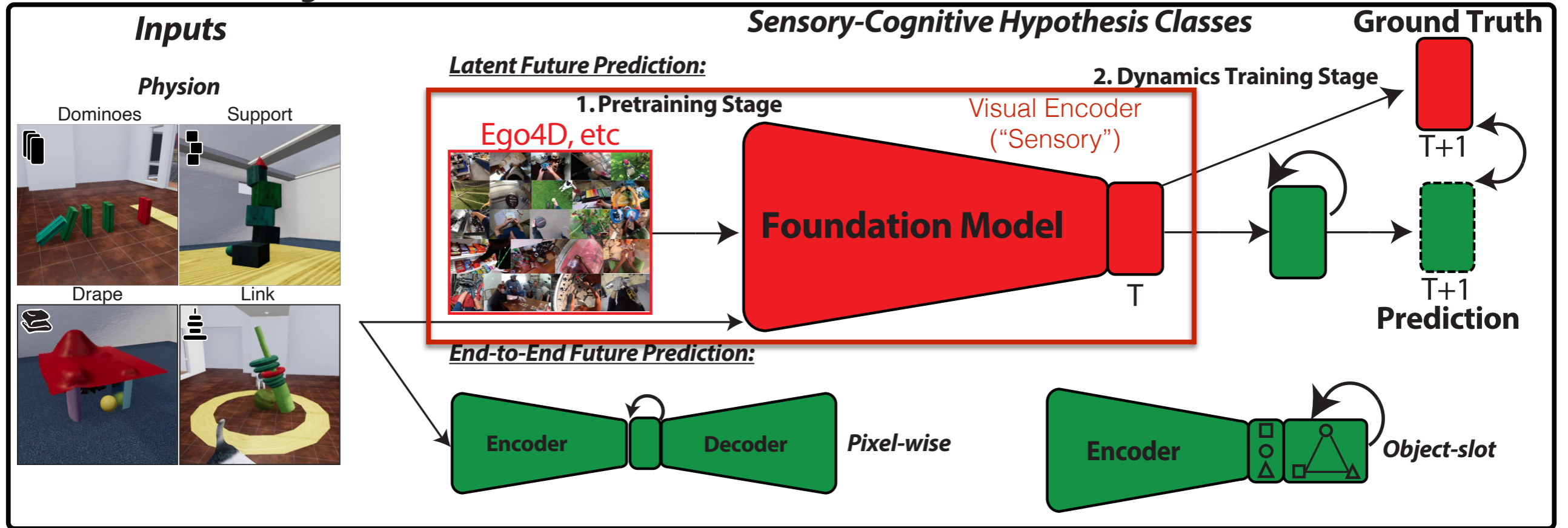
## (A) Model Pretraining



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

# Overall Approach: Foundation Models

## (A) Model Pretraining

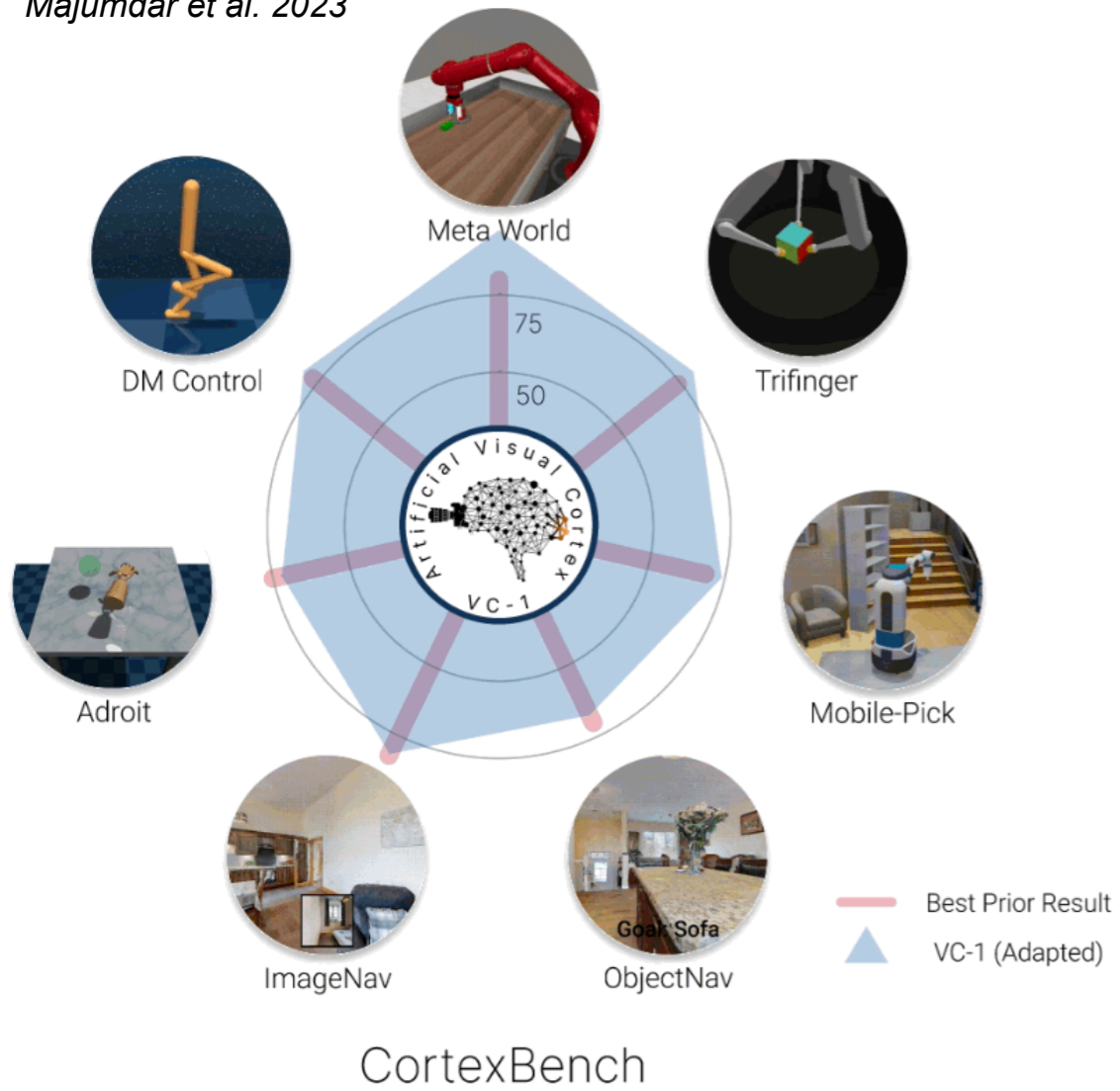


Learn a partial, *implicit* representation of the physical world by performing a challenging vision task ("foundation model")

Emphasis on *reusability!*

# Overall Approach: Foundation Models

Majumdar et al. 2023



# Overall Approach: Foundation Models

Majumdar et al. 2023



CortexBench

## 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

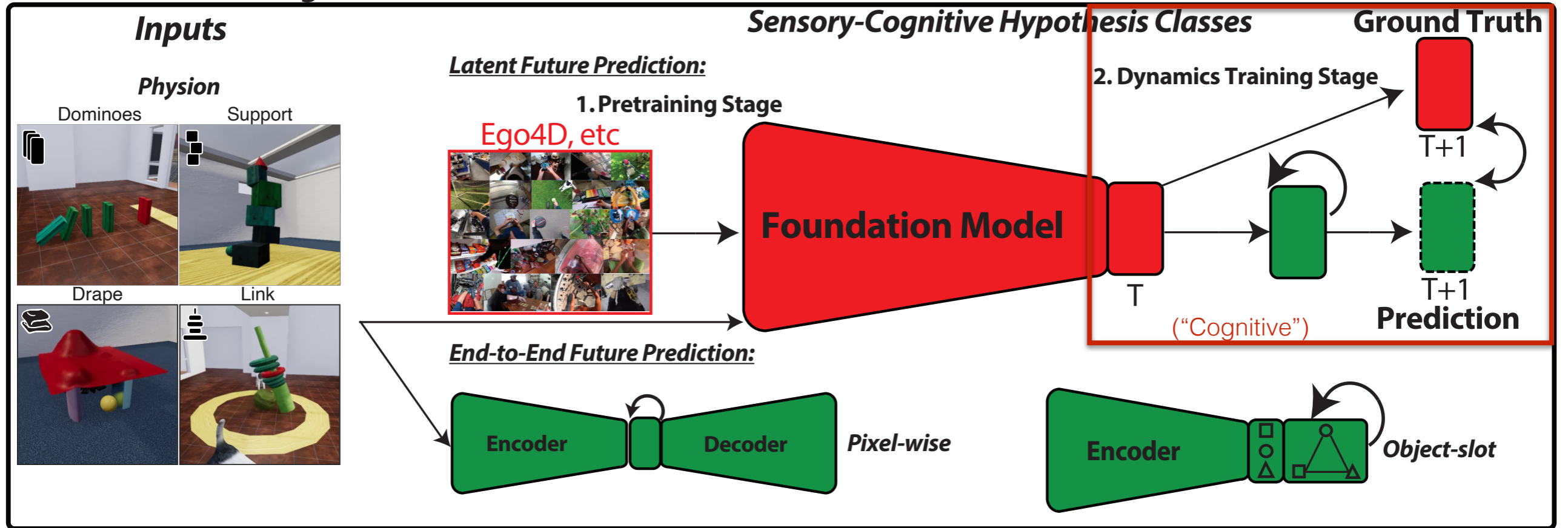
## Ego4D: everyday activity around the world



Grauman et al. 2022

# Overall Approach: Foundation Models + Dynamics

## (A) Model Pretraining



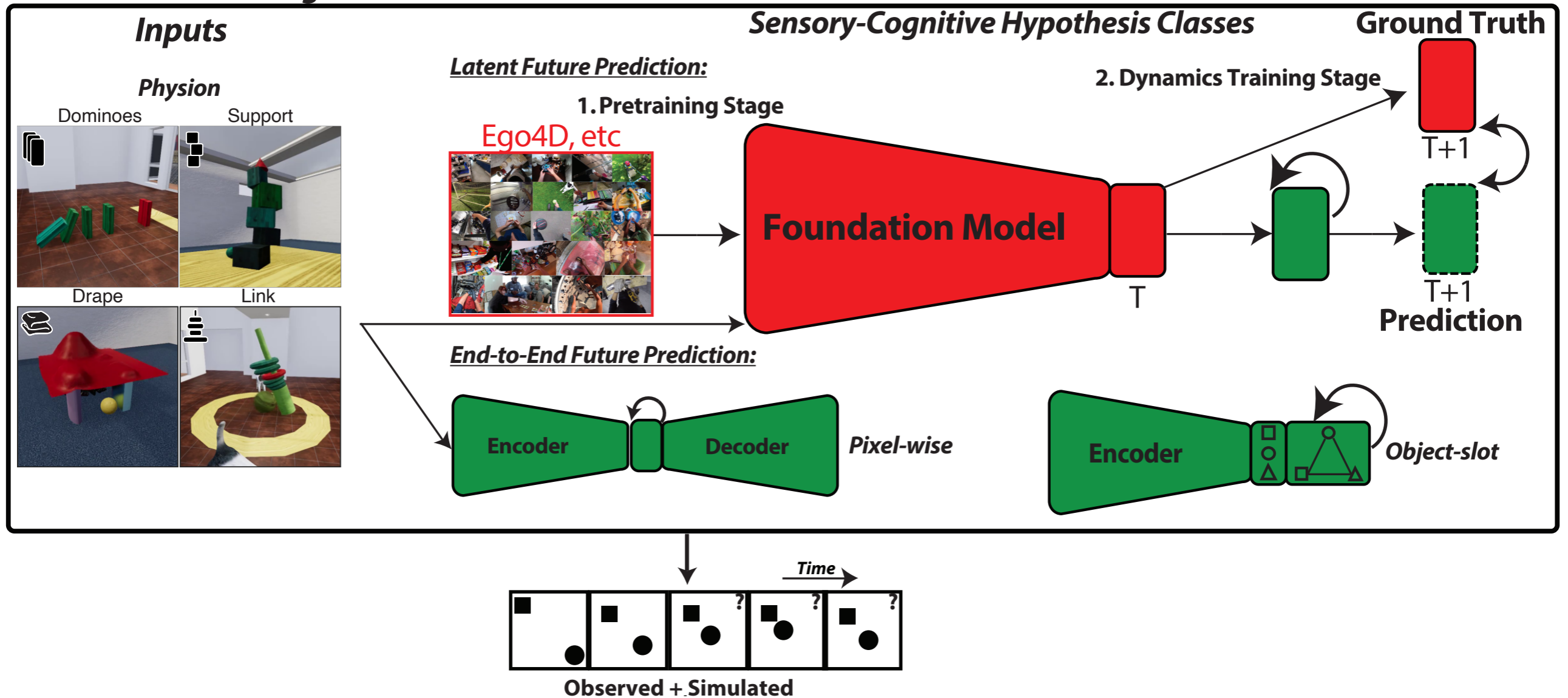
Learn a partial, *implicit* representation of the physical world by performing a challenging vision task (“foundation model”)

*Emphasis on reusability!*

Leverage these dynamics to do explicit physical simulation

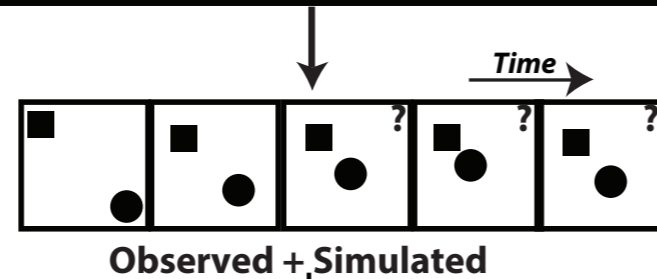
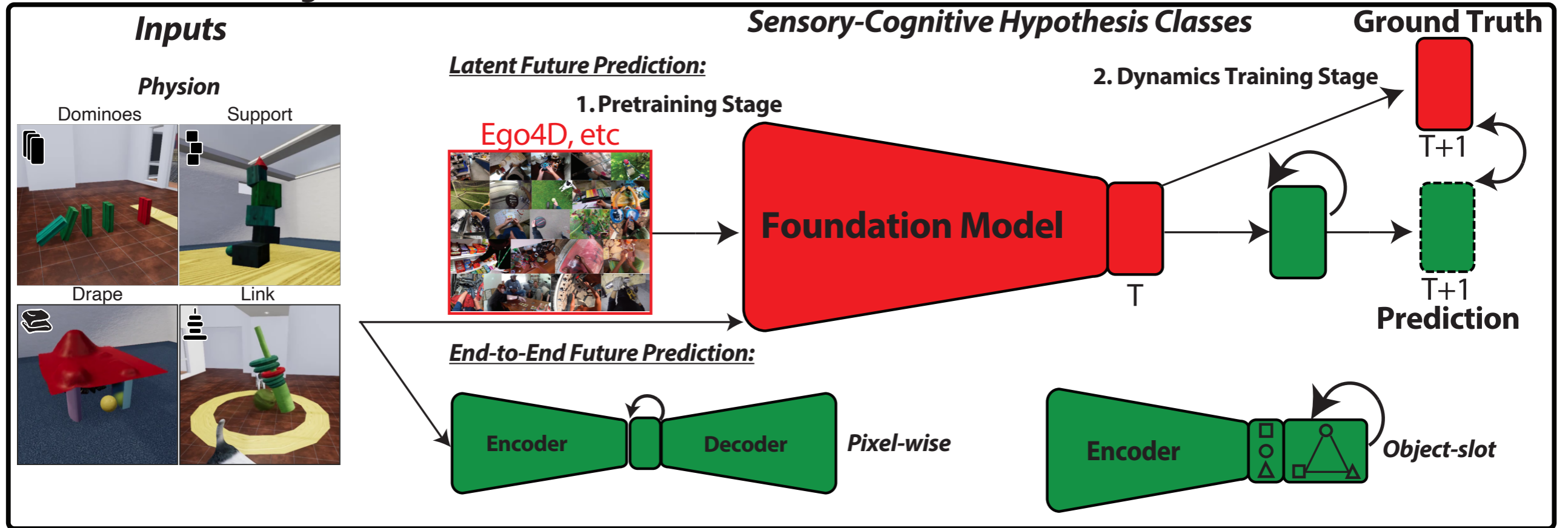
# Overall Approach

## (A) Model Pretraining



# Overall Approach: Model Evaluations

## (A) Model Pretraining



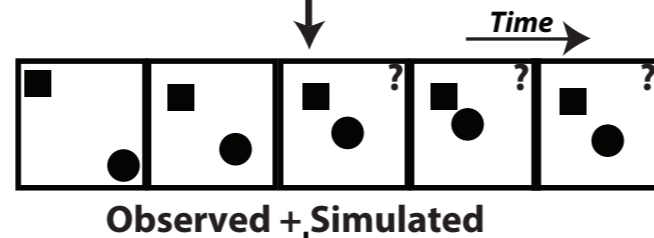
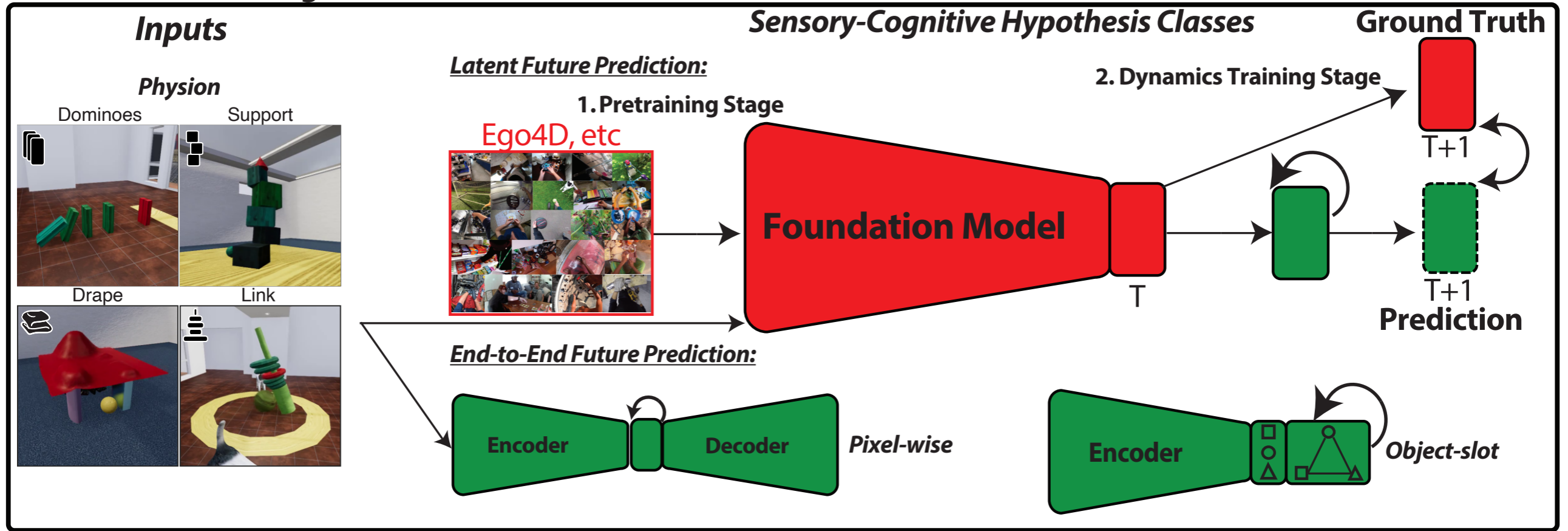
## (B) Model Evaluations



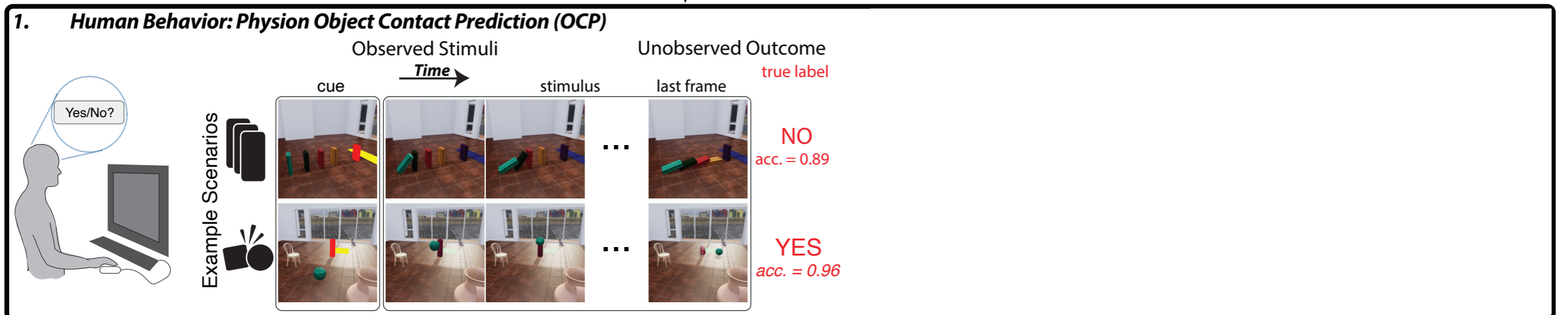


# Overall Approach: Model Evaluations (Human Behavior)

## (A) Model Pretraining

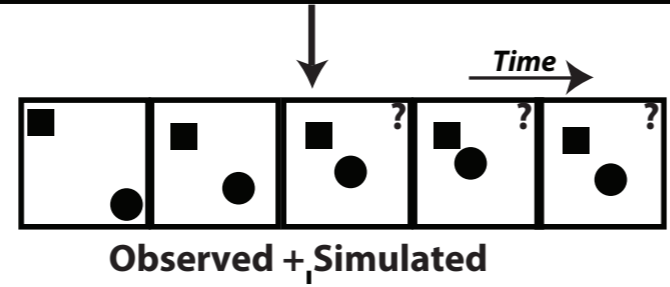
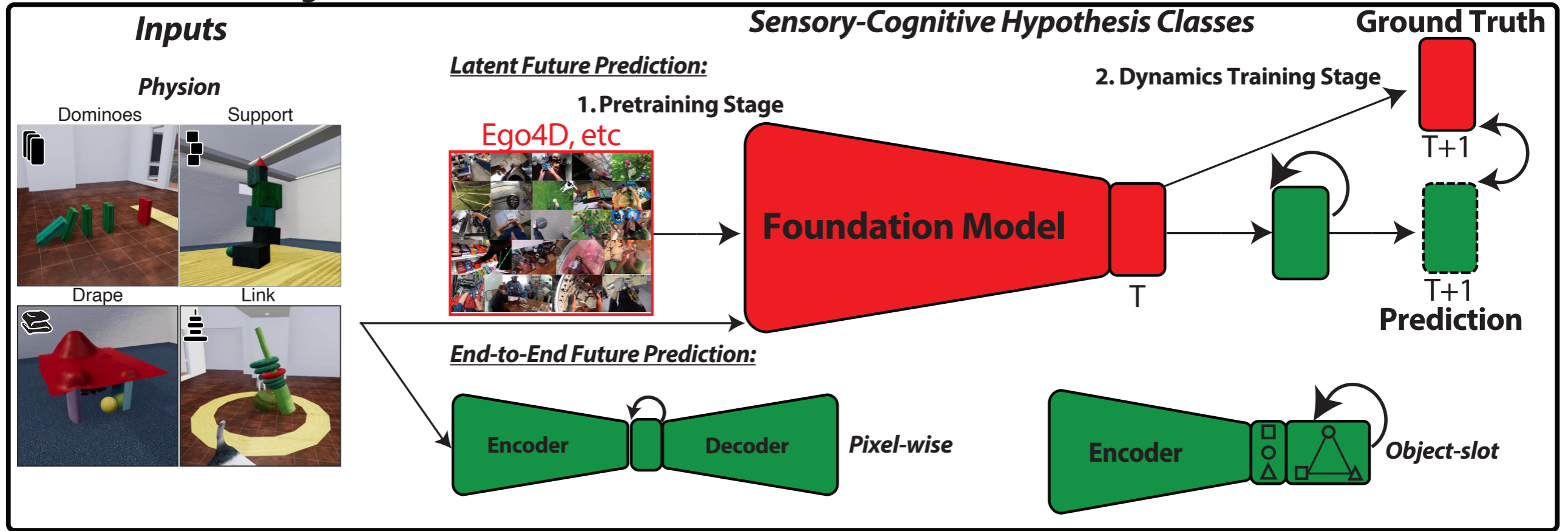


## (B) Model Evaluations

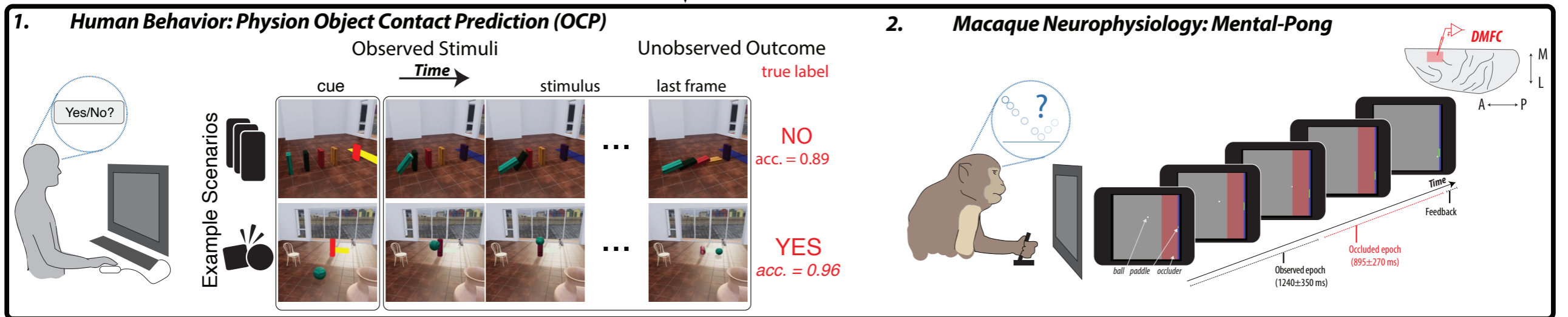


# Overall Approach: Model Evaluations (Macaque Physiology)

## (A) Model Pretraining

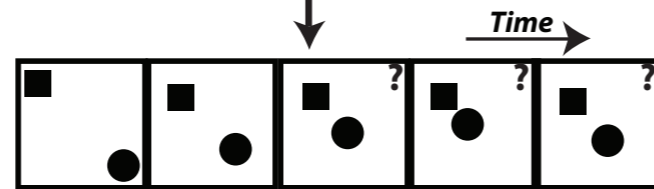
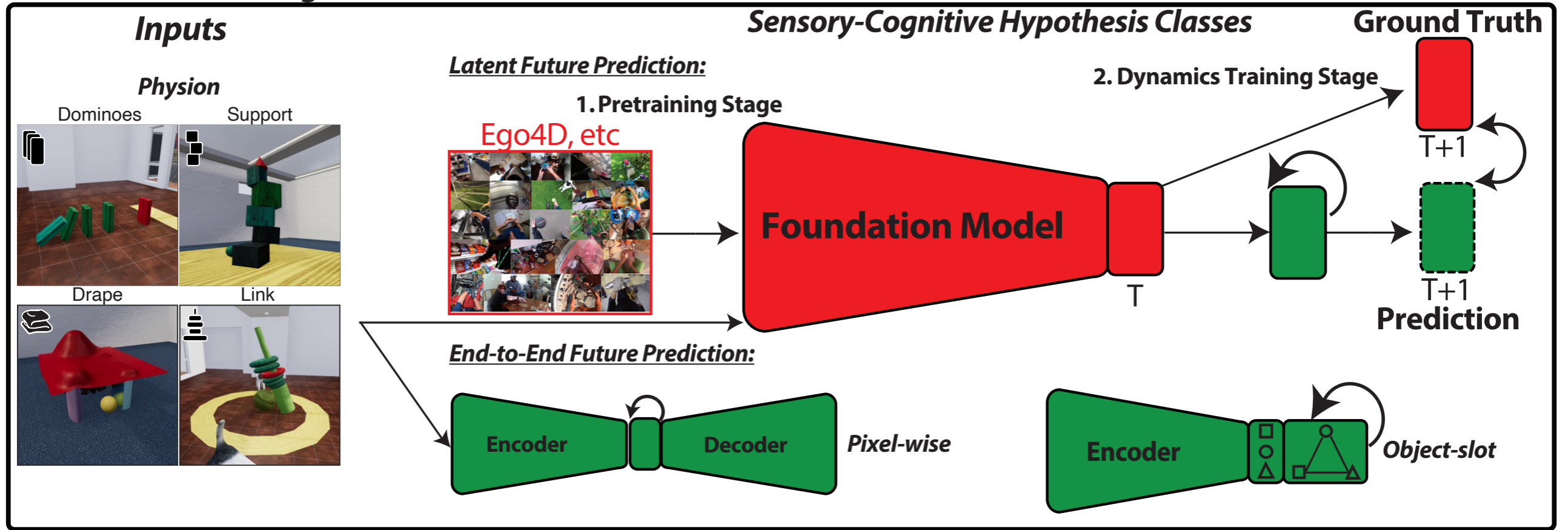


## (B) Model Evaluations

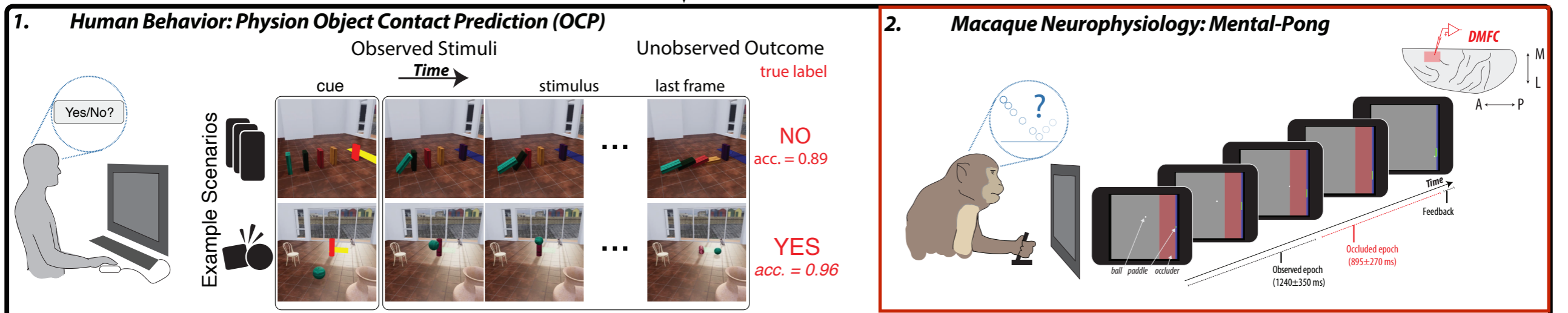


# Model Evaluations: Macaque Neurophysiology

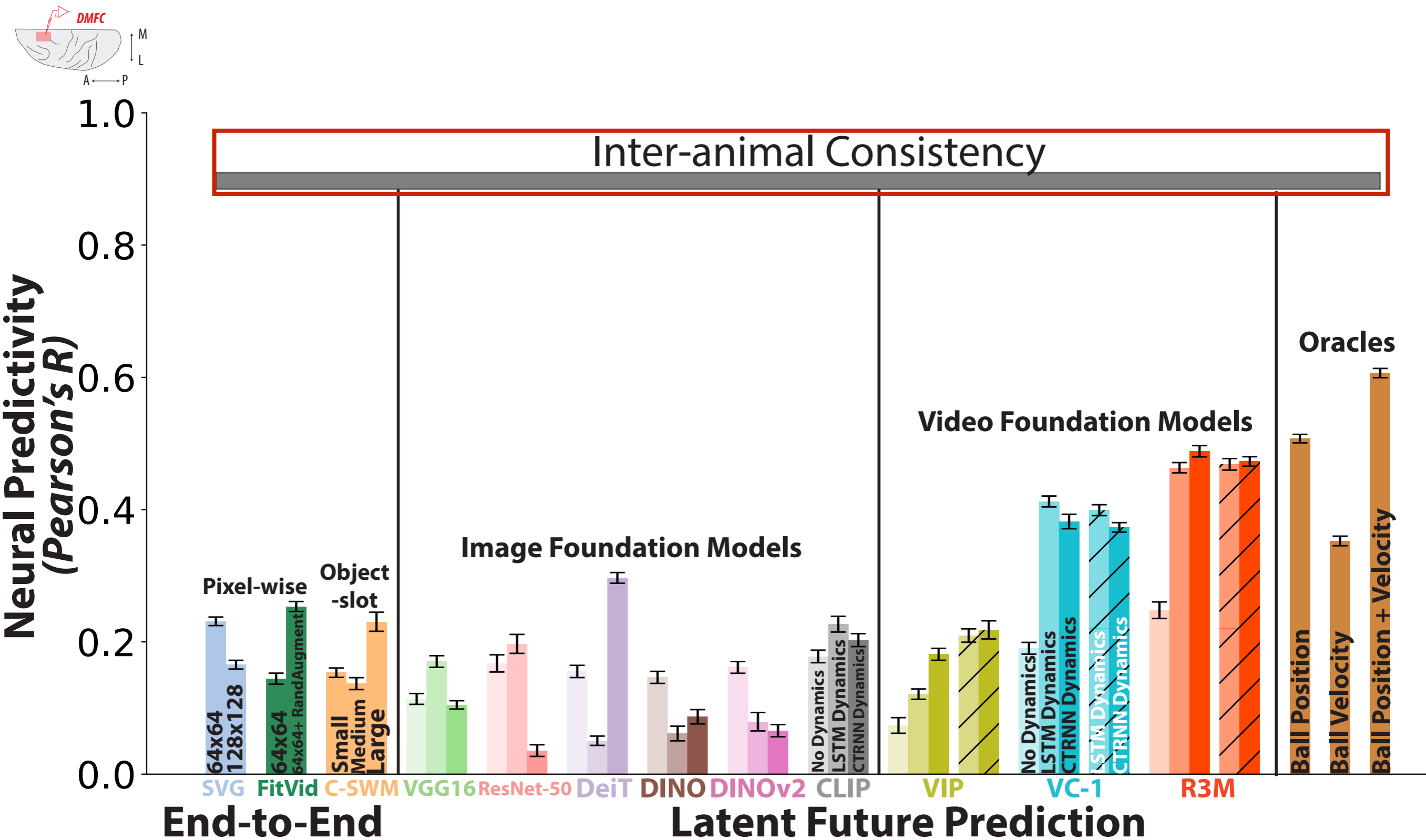
## (A) Model Pretraining



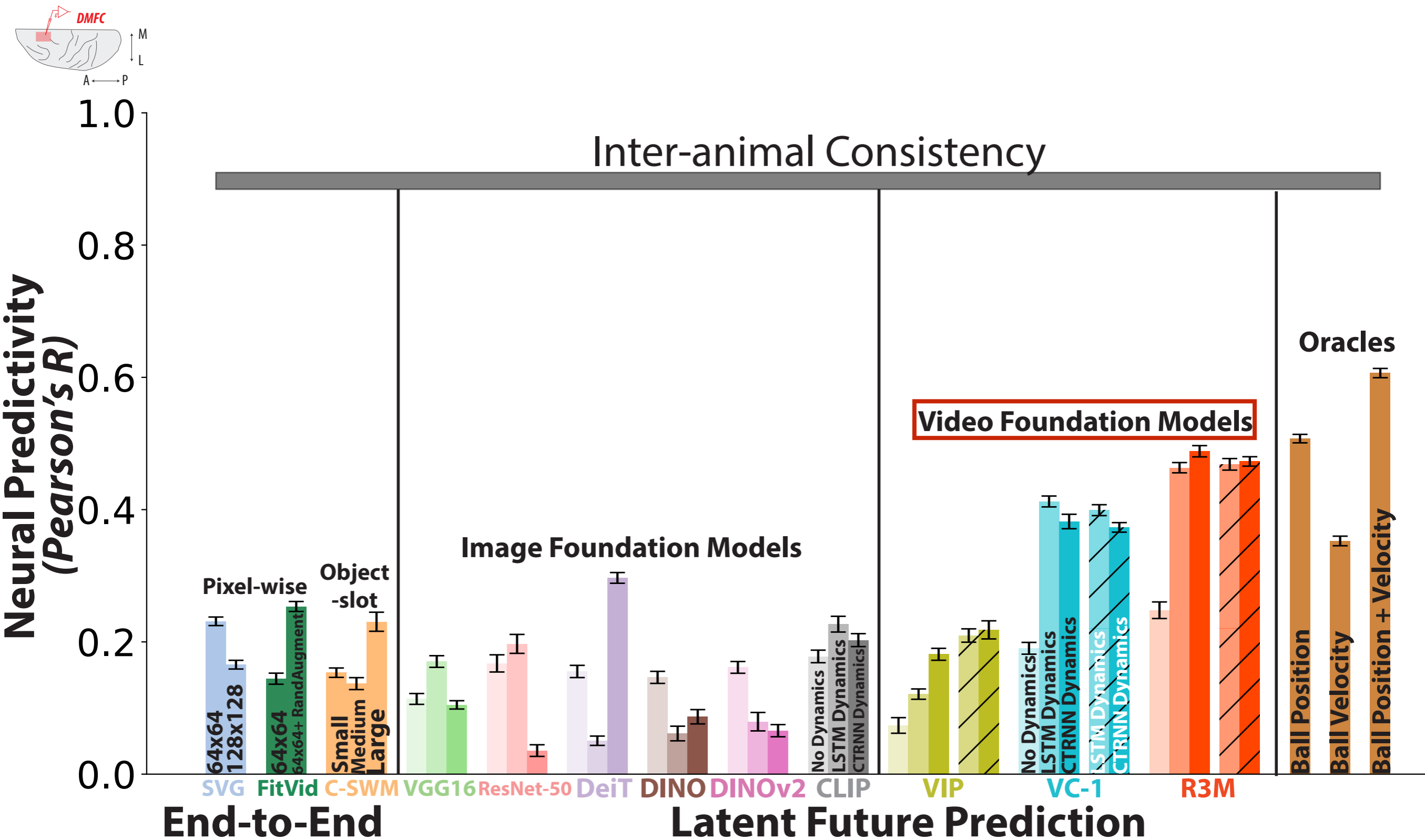
## (B) Model Evaluations



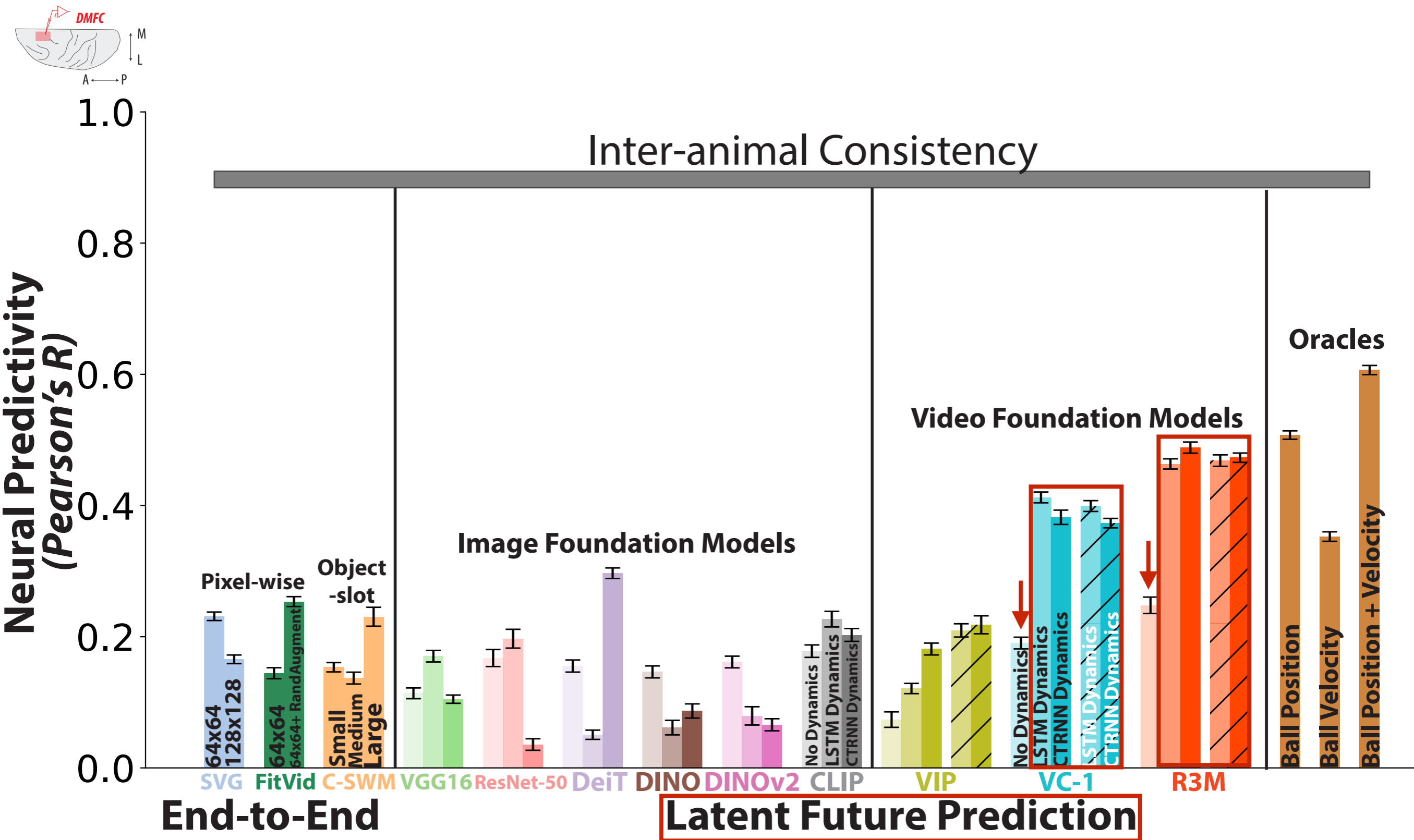
# Neural response predictivity strongly separates models



# Neural response predictivity strongly separates models



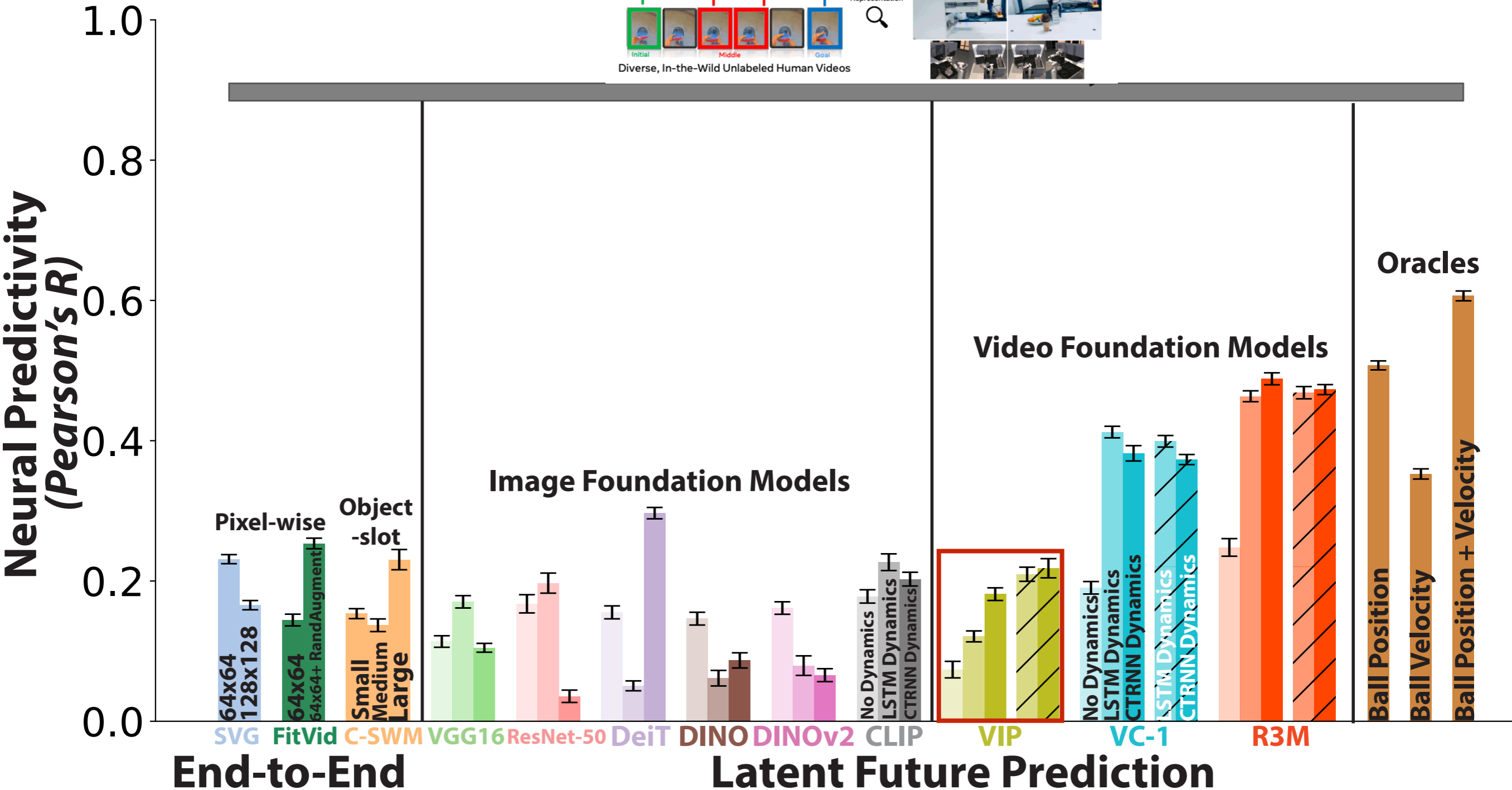
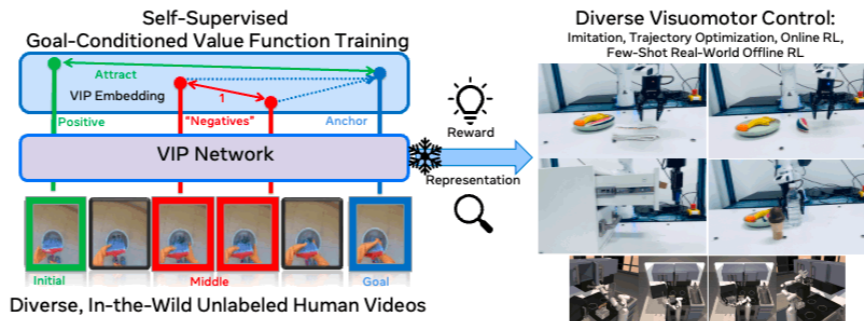
# Neural response predictivity strongly separates models



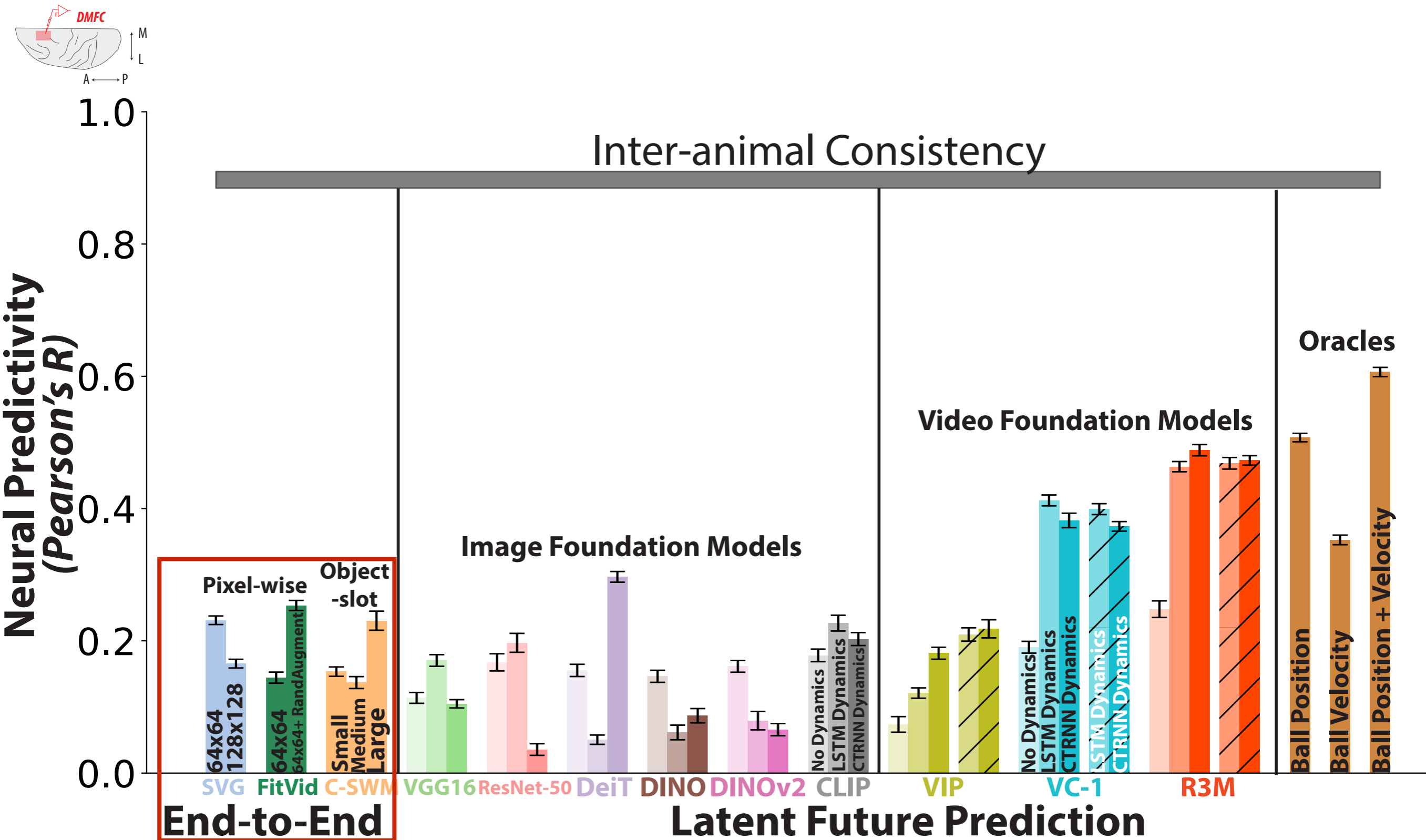
# Neural response predictivity strongly separates models

Ma et al. 2023

## VIP: Towards Universal Visual Reward and Representation Via Value-Implicit Pre-Training

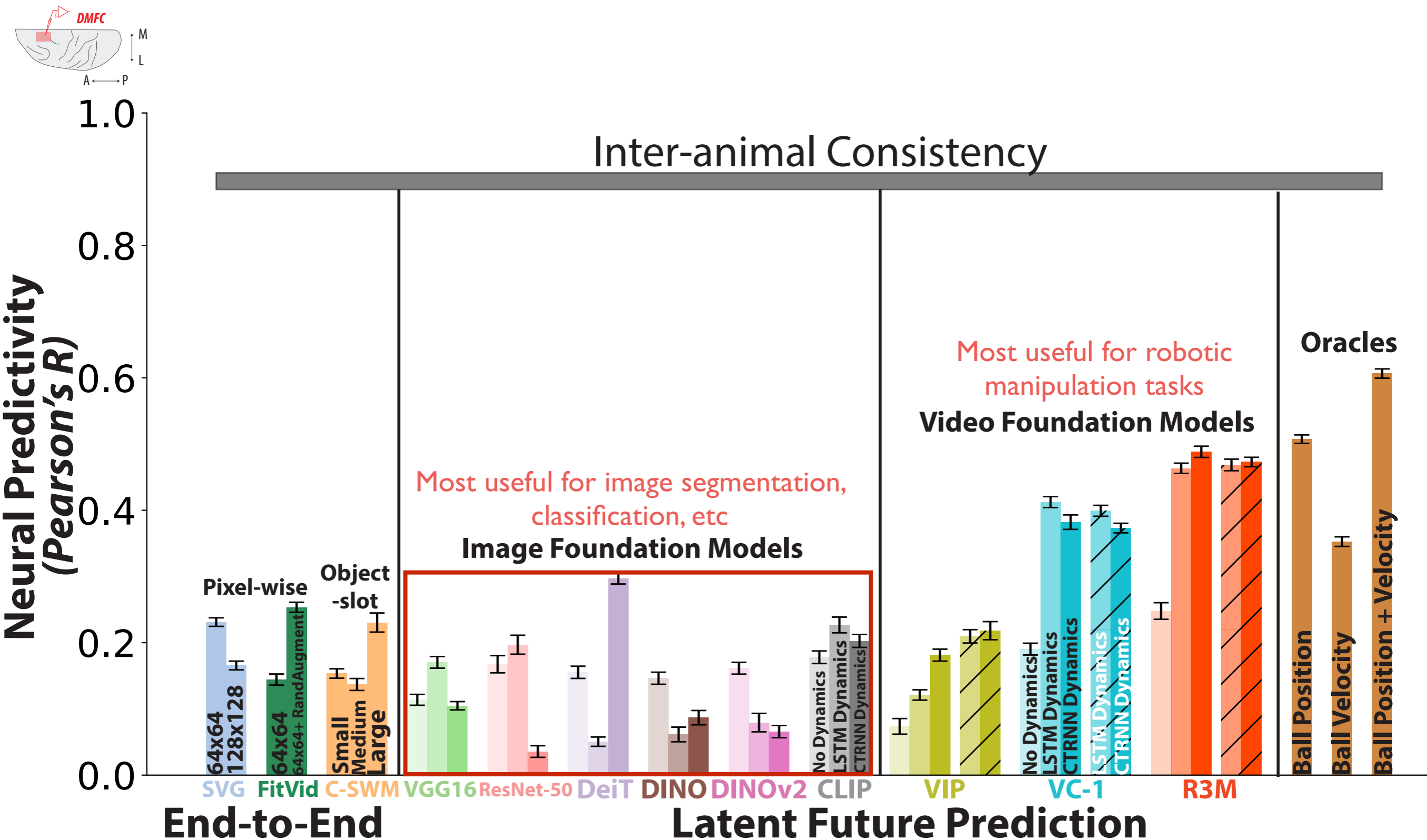


# Neural response predictivity strongly separates models



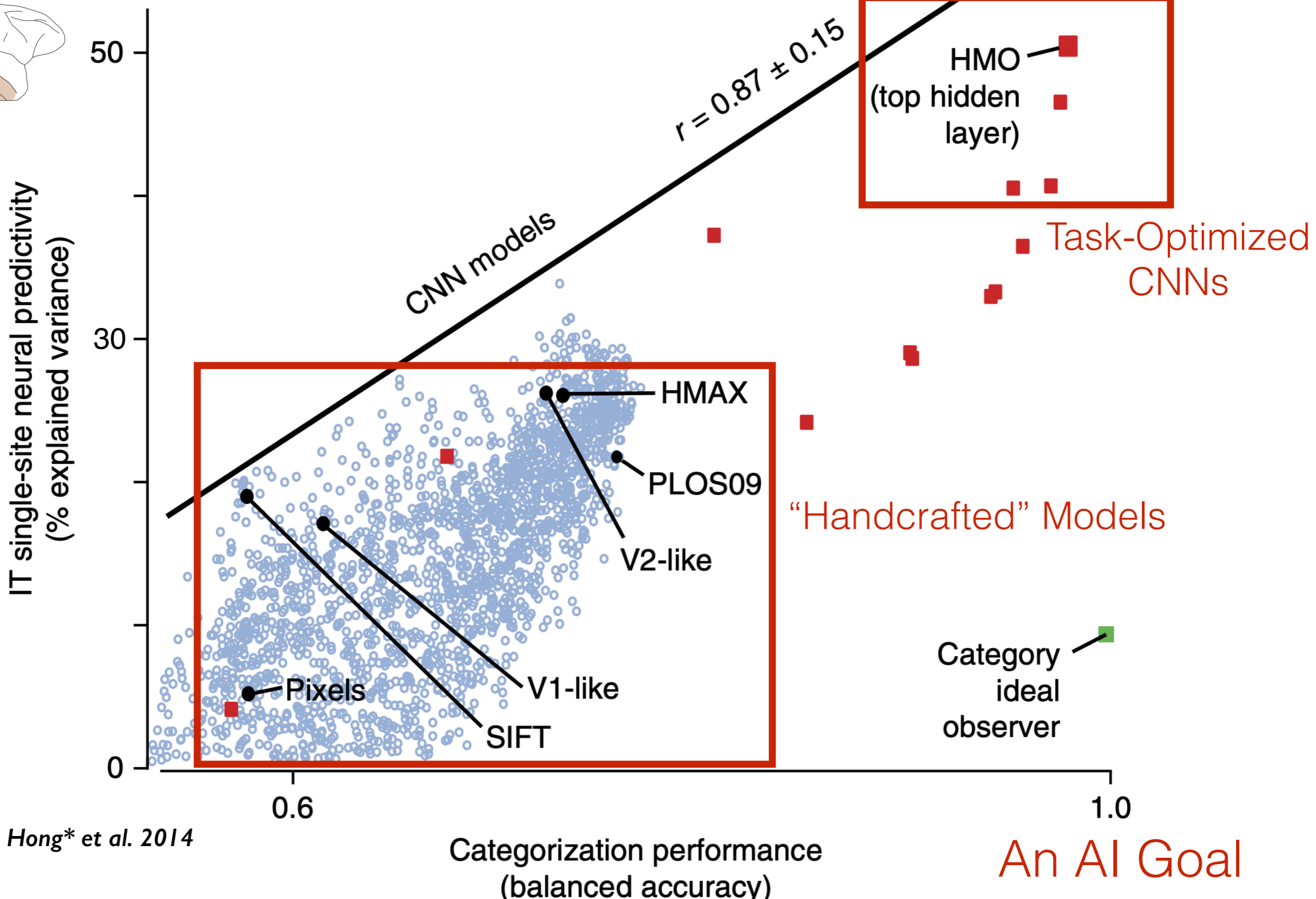
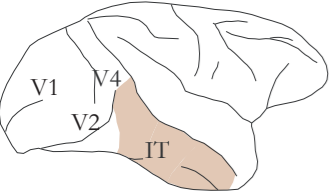


# Neural response predictivity strongly separates models



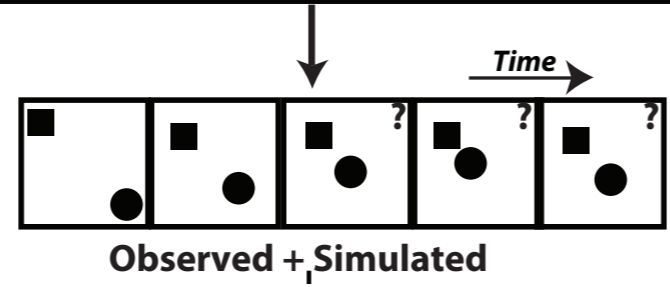
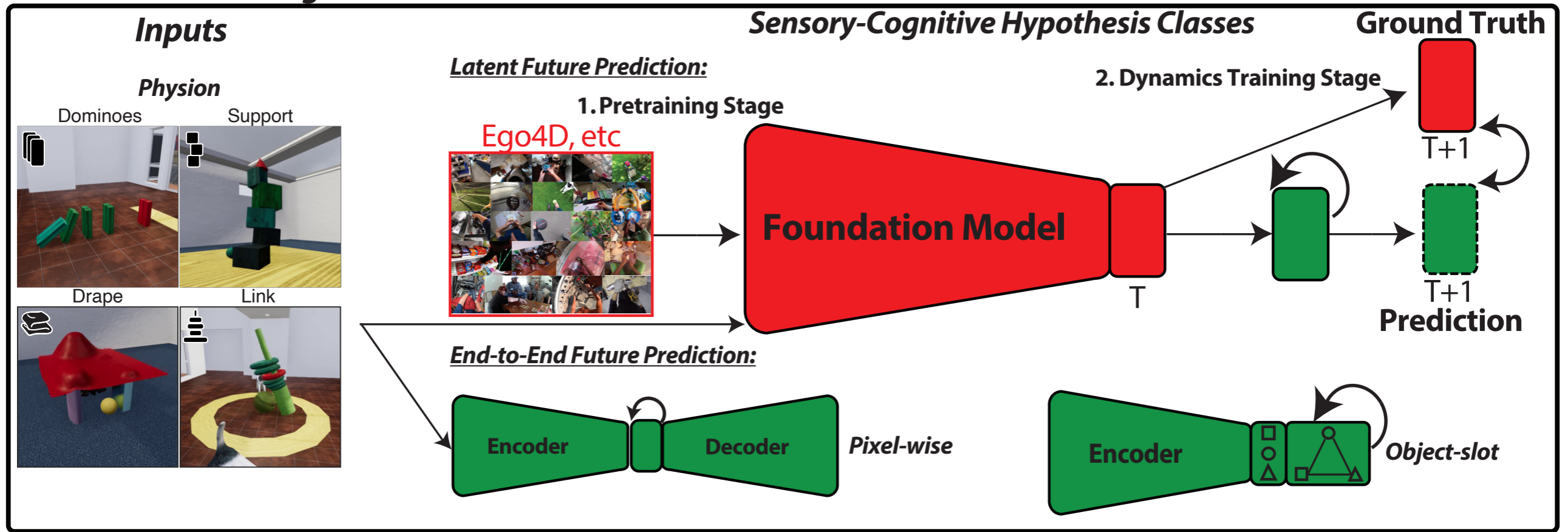
# Prior Results in Inferior Temporal (IT) Cortex

## A Neuroscience Goal



# Model Evaluations: Object Contact Prediction (OCP)

## (A) Model Pretraining



## (B) Model Evaluations

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

Observed Stimuli      Unobserved Outcome

cue      stimulus      last frame      true label

Time

Example Scenarios

NO  
acc. = 0.89

YES  
acc. = 0.96

### 2. Macaque Neurophysiology: Mental-Pong

DMFC

M

L

A ← P

ball paddle occluder

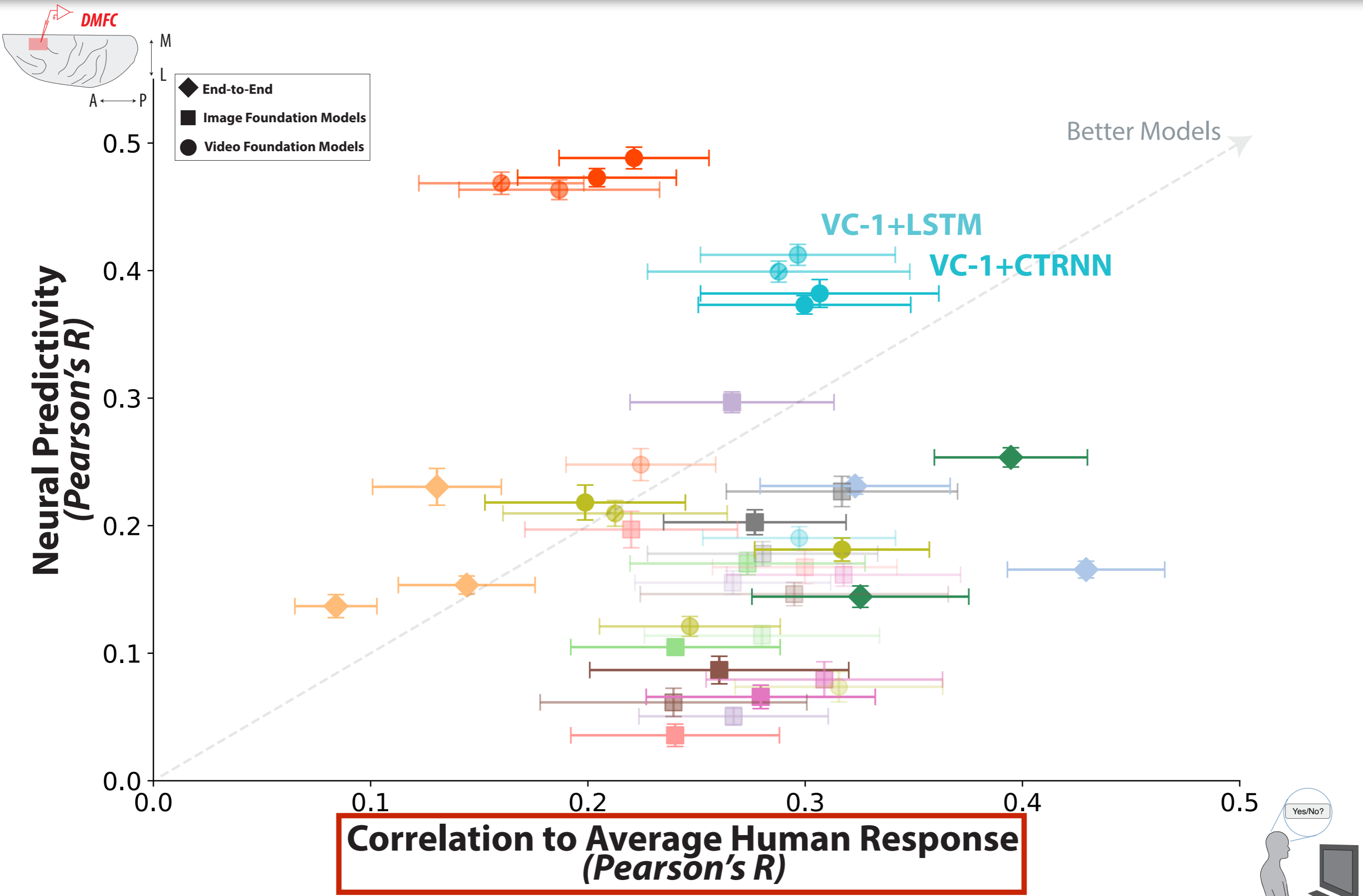
Observed epoch (1240 ± 350 ms)

Occluded epoch (895 ± 270 ms)

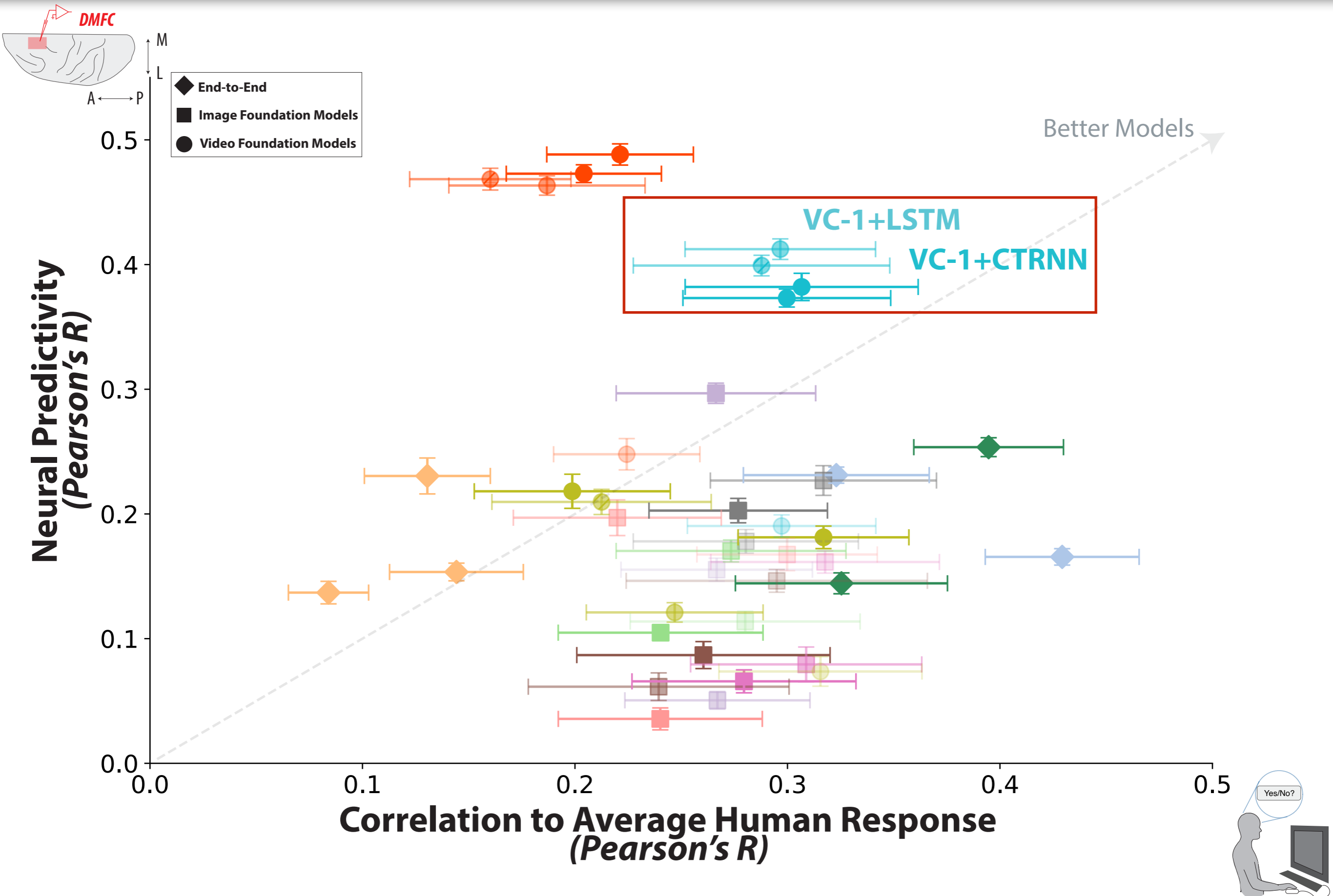
Feedback

Time

# Video Foundation Models Can Better Match Both



# Video Foundation Models Can Better Match Both



# Takeaways so far...

1. Mental simulation appears to be primarily relevant to predicting the **future** state of the environment in a suitable **latent** space.

# Takeaways so far...

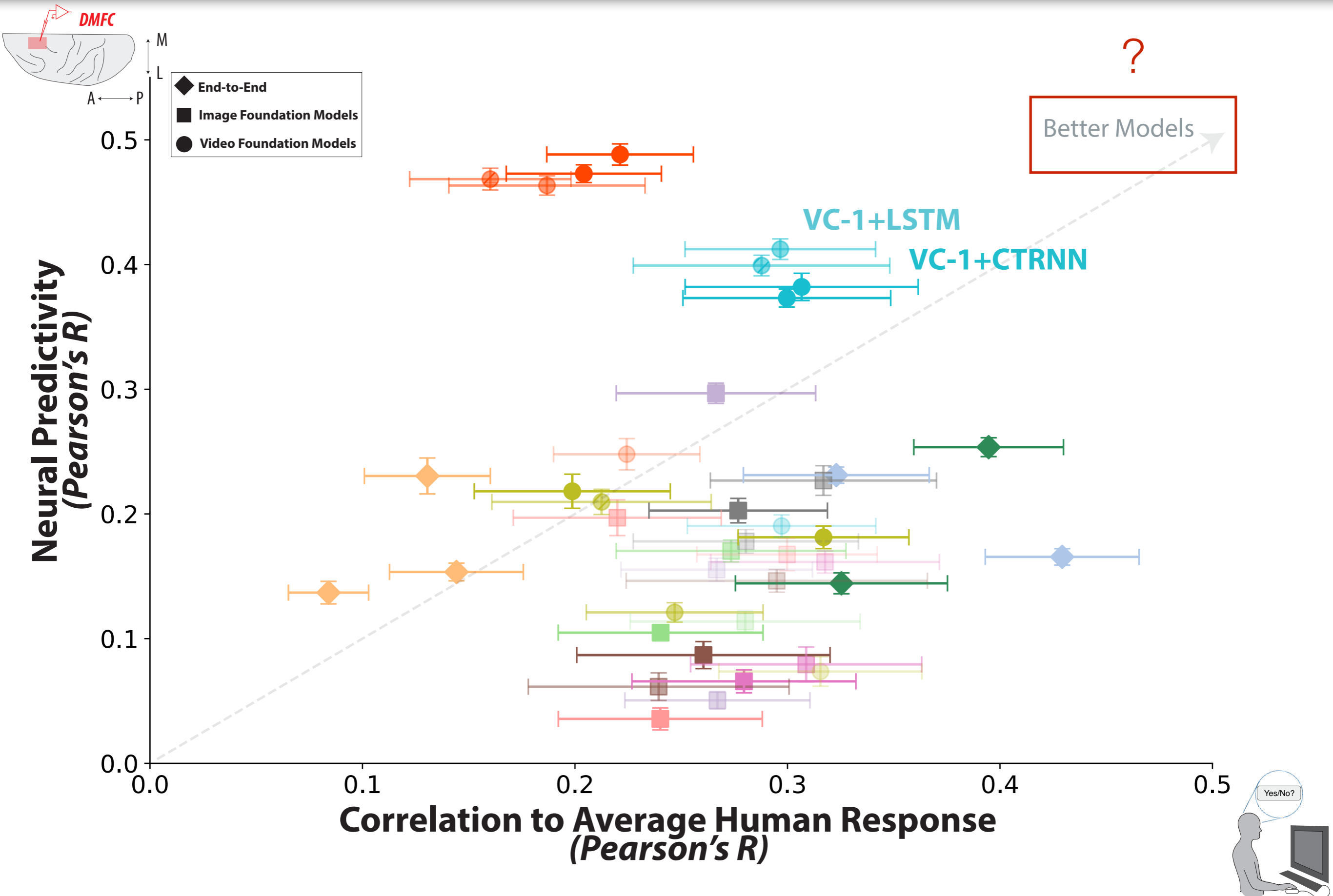
1. Mental simulation appears to be primarily relevant to predicting the **future** state of the environment in a suitable **latent** space.
2. In particular, this latent space is highly constrained -- it doesn't appear to consist of bespoke object slots or prioritize fine-grained details (e.g. at the level of pixels), but rather mainly has to be **reusable** across *dynamic* scenes.

# Takeaways so far...

1. Mental simulation appears to be primarily relevant to predicting the **future** state of the environment in a suitable **latent** space.
2. In particular, this latent space is highly constrained -- it doesn't appear to consist of bespoke object slots or prioritize fine-grained details (e.g. at the level of pixels), but rather mainly has to be **reusable** across *dynamic* scenes.
3. So far a correspondence between the ability to predict neural & behavioral responses, and developing useful representations for Embodied AI more generally (rather than classic computer vision tasks e.g. classification, segmentation, etc).



# Future Directions

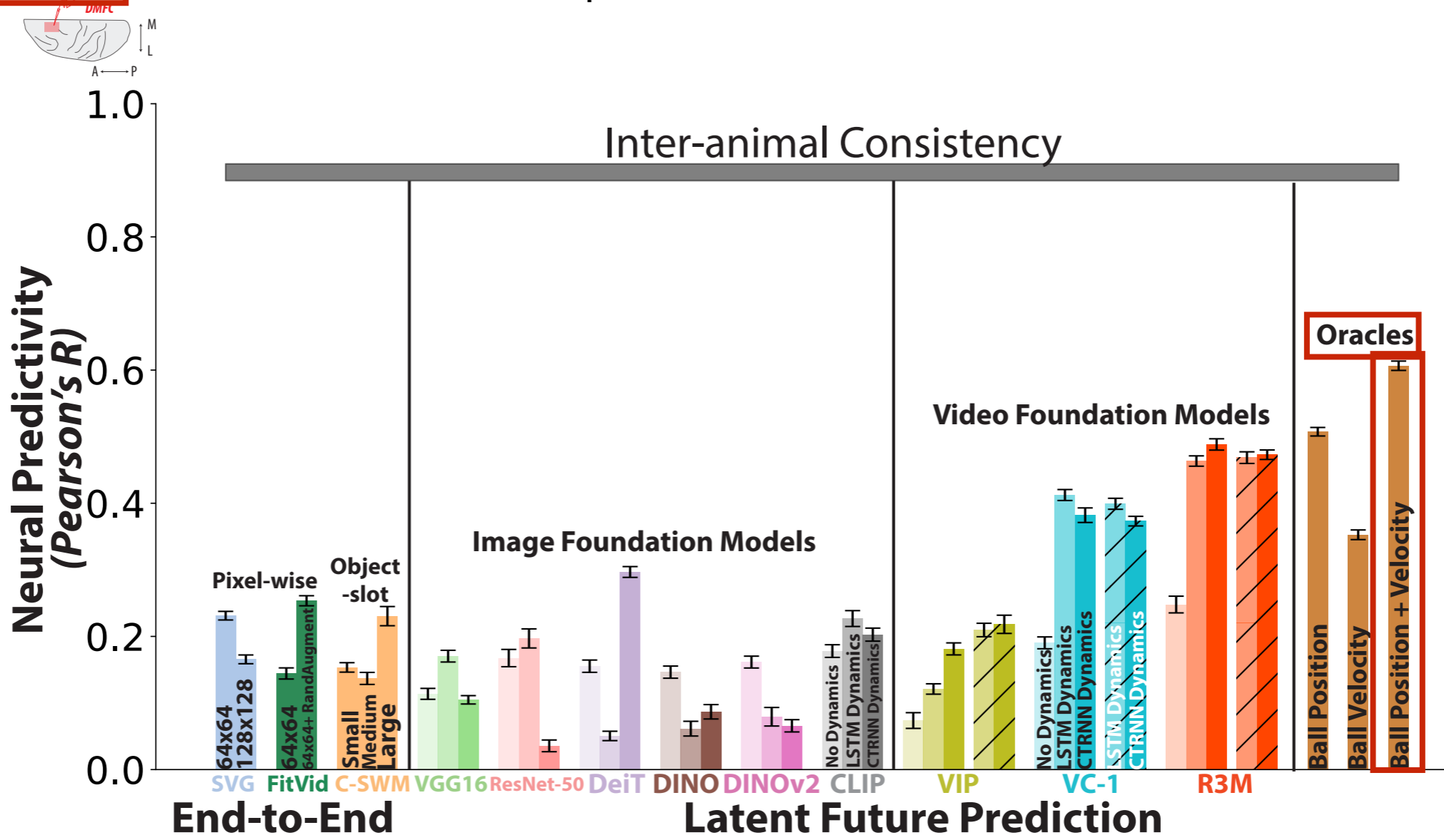


# Future Directions

1. **Sensory**: Better leverage temporal relationships to learn a more “factorized” *and* reusable representation:

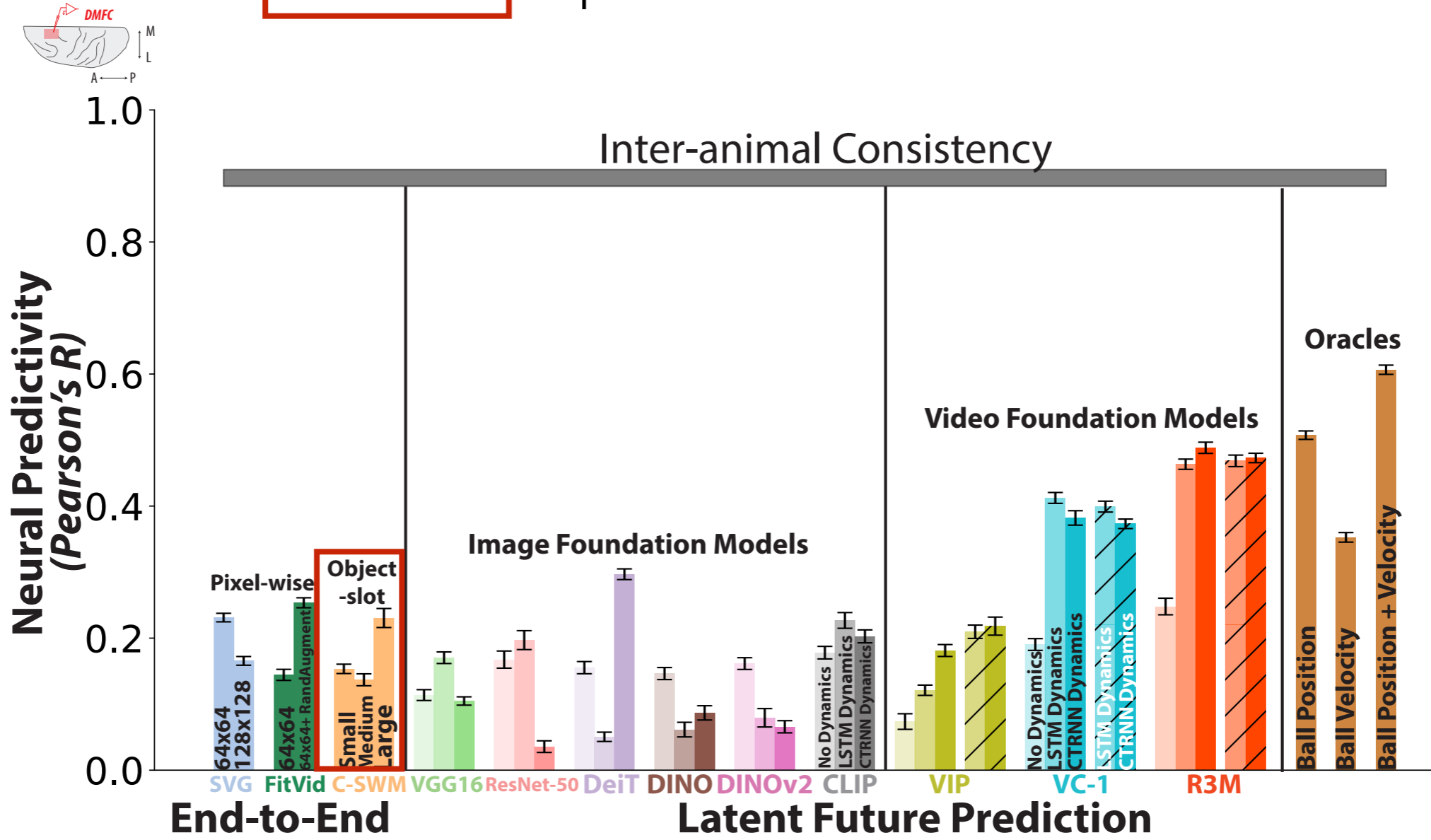
# Future Directions

1. **Sensory:** Better leverage temporal relationships to learn a more “factorized” and reusable representation:



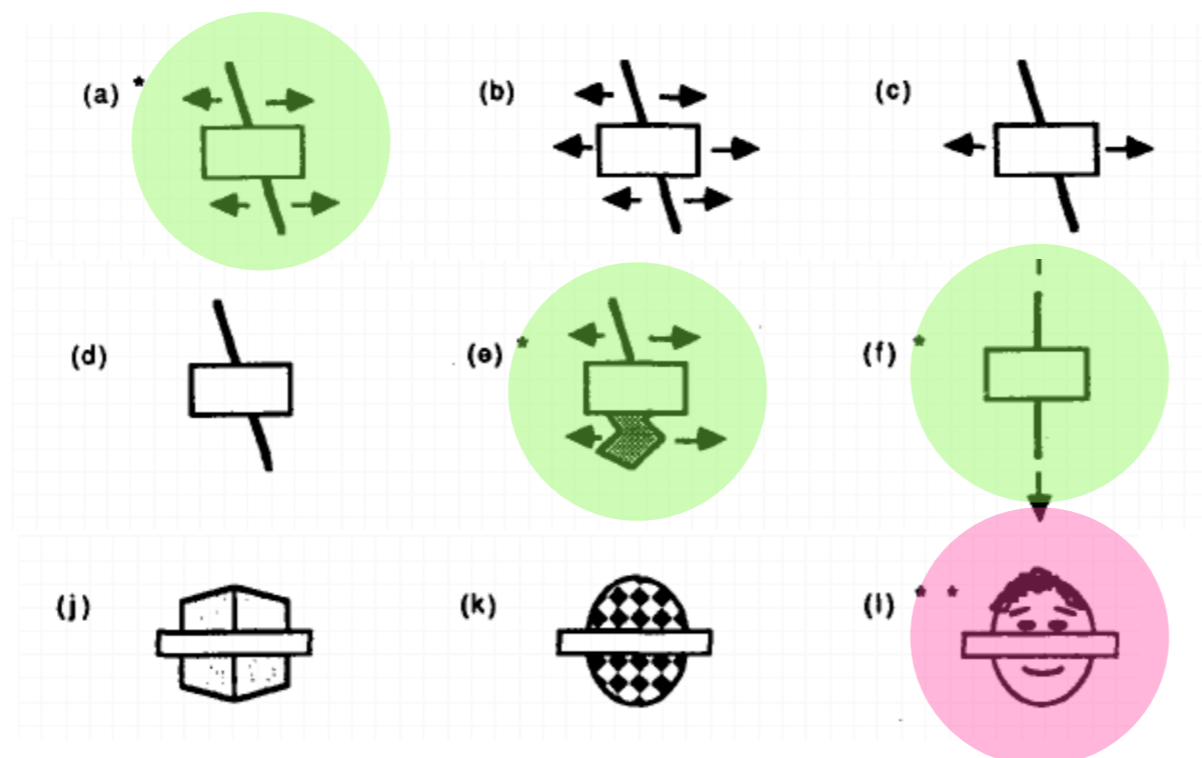
# Future Directions

1. **Sensory:** Better leverage temporal relationships to learn a more “factorized” *and* reusable representation:



# Future Directions

1. **Sensory:** Better leverage temporal relationships to learn a more “factorized” *and* reusable representation: **object-centric, video foundation model?**



*Principles of Object Perception* Elizabeth Spelke, 1990



Elizabeth Spelke

# Future Directions

1. **Sensory**: Better leverage temporal relationships to learn a more “factorized” *and* reusable representation: **object-centric, video foundation model?**
2. **Cognitive**: Hierarchy/modularization of timescales in dynamics?

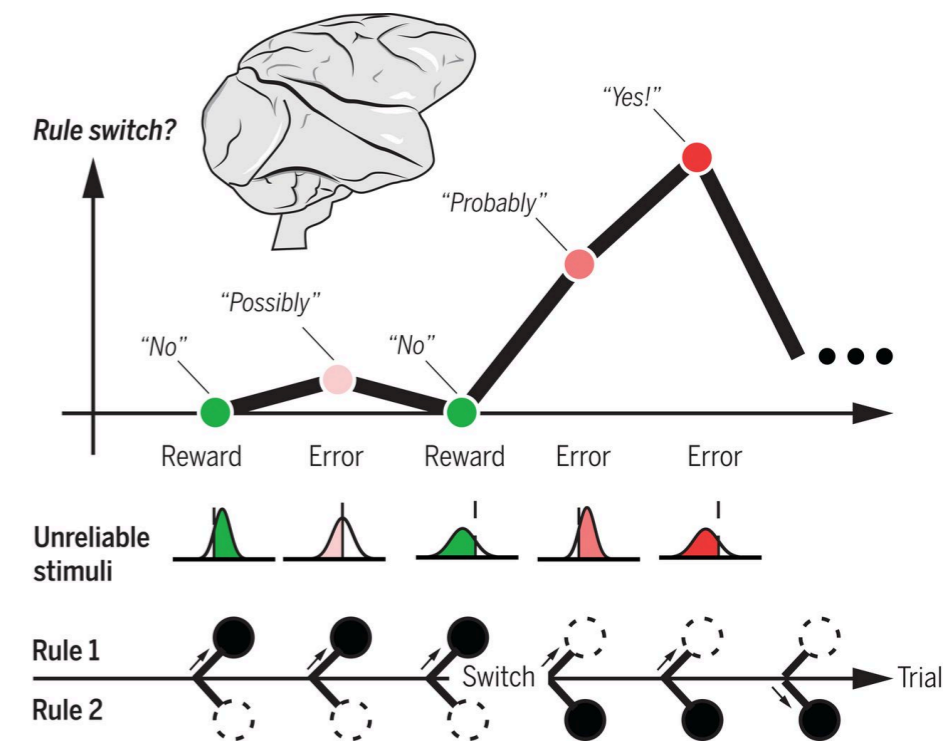
# Future Directions

1. **Sensory:** Better leverage temporal relationships to learn a more “factorized” *and* reusable representation: **object-centric, video foundation model?**
2. **Cognitive:** Hierarchy/modularization of timescales in dynamics?

## Hierarchical reasoning by neural circuits in the frontal cortex

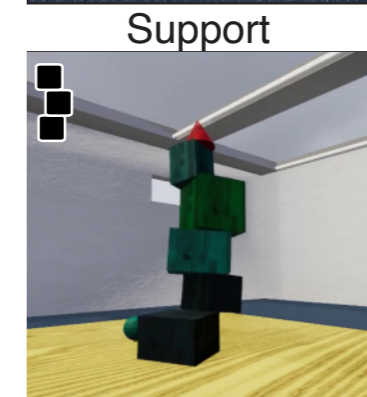
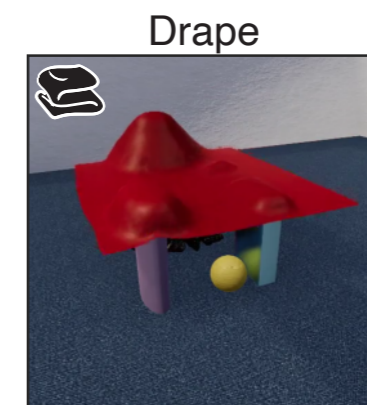
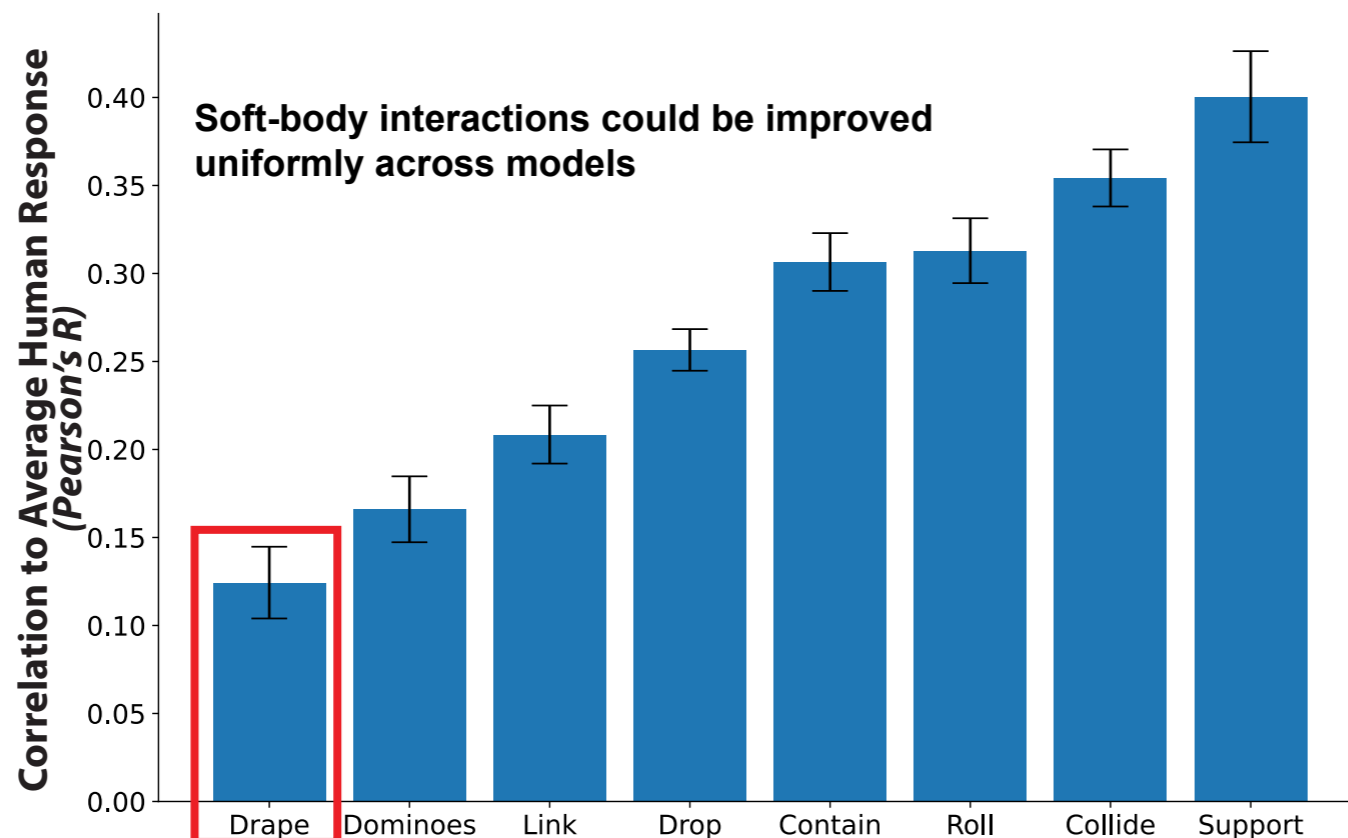
MORTEZA SARAFYAZD  AND MEHRDAD JAZAYERI  [Authors Info & Affiliations](#)

SCIENCE • 17 May 2019 • Vol 364, Issue 6441 • DOI: 10.1126/science.aav8911



# Future Directions: Learning Diverse Material Properties

1. **Sensory**: Better leverage temporal relationships to learn a more “factorized” *and* reusable representation: **object-centric, video foundation model?**
2. **Cognitive**: Hierarchy/modularization of timescales in dynamics?
3. **Data**: More complex 2D and 3D scenes/real world objects





# Acknowledgements




Rishi Rajalingham



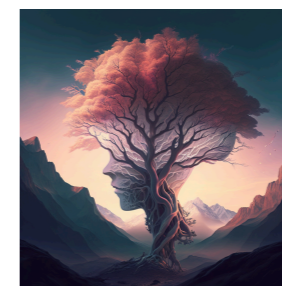
Mehrdad Jazayeri



Guangyu Robert Yang

Contact:  
[anayebi@mit.edu](mailto:anayebi@mit.edu)  
 [@aran\\_nayebi](https://twitter.com/aran_nayebi)

**Preprint:** <https://arxiv.org/abs/2305.11772>



YangLab