

Yale University
COGNITIVE & NEURAL COMPUTATION
LABORATORY

Email: daniel.calbick@yale.edu

A new multi-level modeling framework provides evidence for the simulation of object dynamics in the dorsomedial frontal cortex.

Daniel Calbick¹, Jason Z. Kim², Hansem Sohn³, Mehrdad Jazayeri⁴, Ilker Yildirim¹

¹ Interdepartmental Neuroscience Program, Yale University; ² Department of Physics, Cornell University;

³ Center for Neuroscience Imaging Research & Department of Biomedical Engineering, Sungkyunkwan University; ⁴ Department of Brain & Cognitive Sciences, MIT

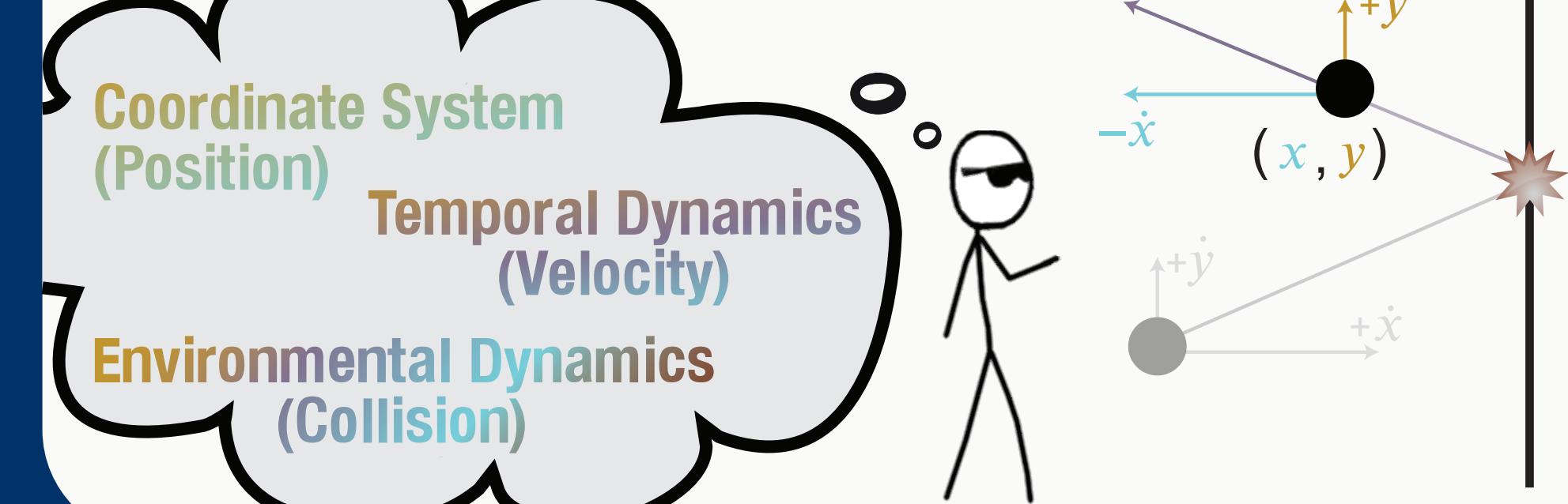
Mental Simulation & World Models

Runnable mental models of our environments enable flexible planning.



Two aspects:

- Representation of objects and how they move and interact with each other
- Planning with these object representations to allow actions

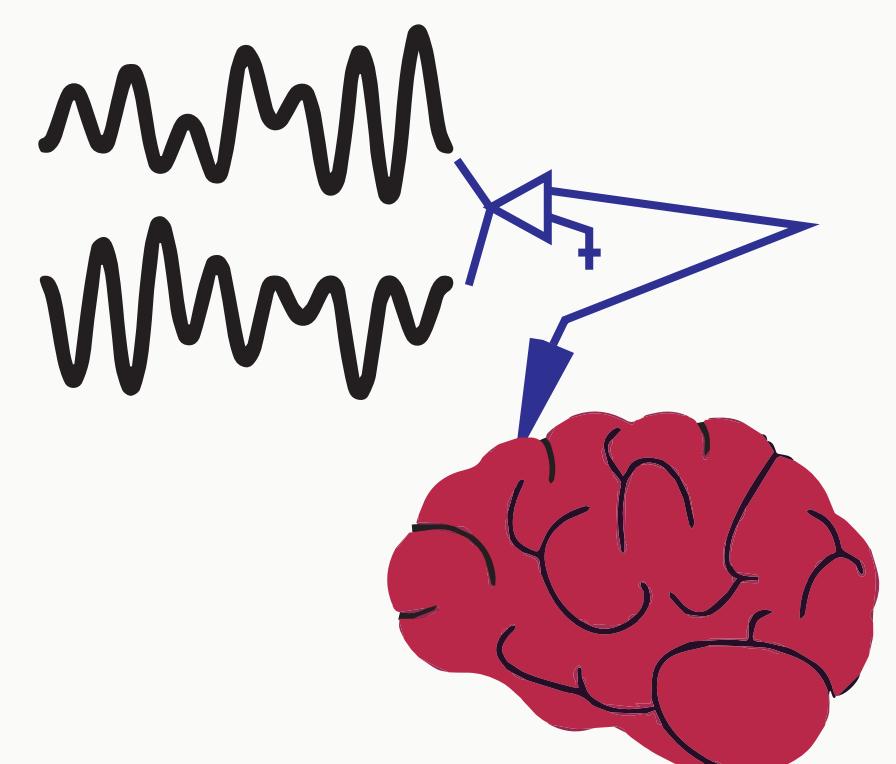
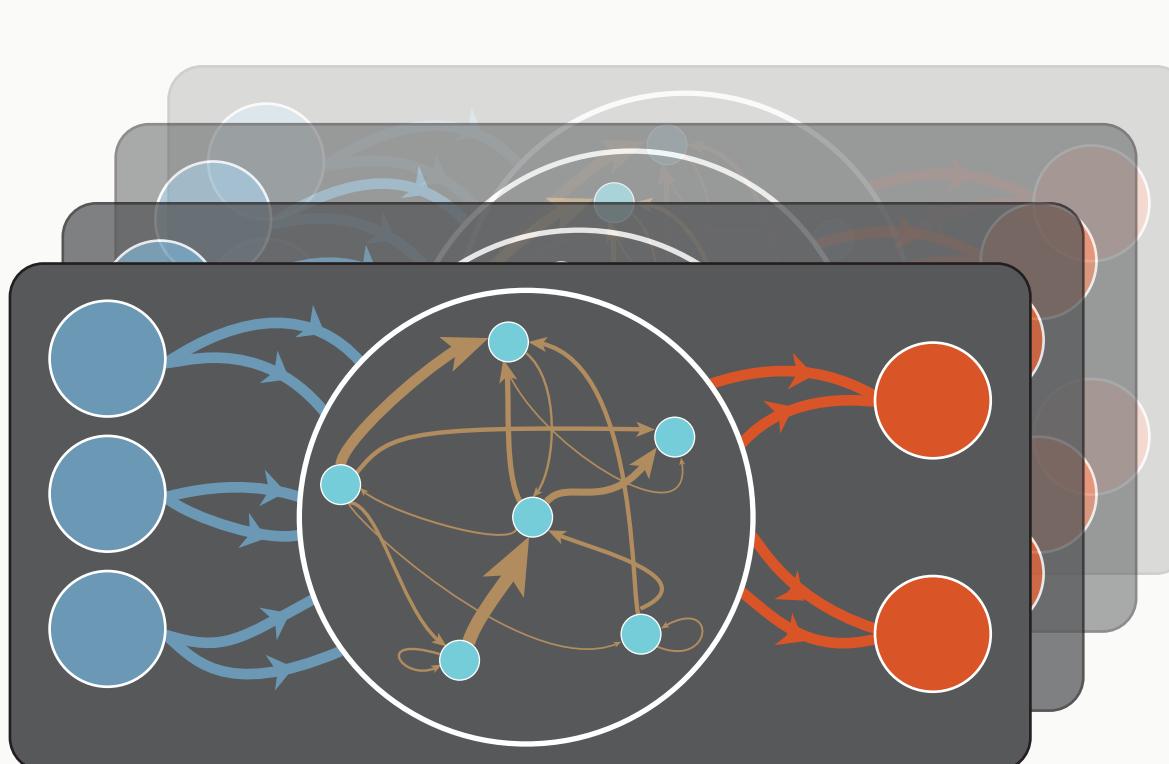


Candidate Symbolic World Model

```
x[0], dx[0] ~ NormalDist
y[0], dy[0] ~ NormalDist
x[t] = x[t-1] + dx[t]
y[t] = y[t-1] + dy[t]
if x[t] < left_wall
    dx[t] = -1 * dx[t]
if y[t] < lower_wall
    dy[t] = -1 * dy[t]
```

Key Issue: Existing frameworks do not allow for hypothesis-driven exploration of structured world models in neurobiology.

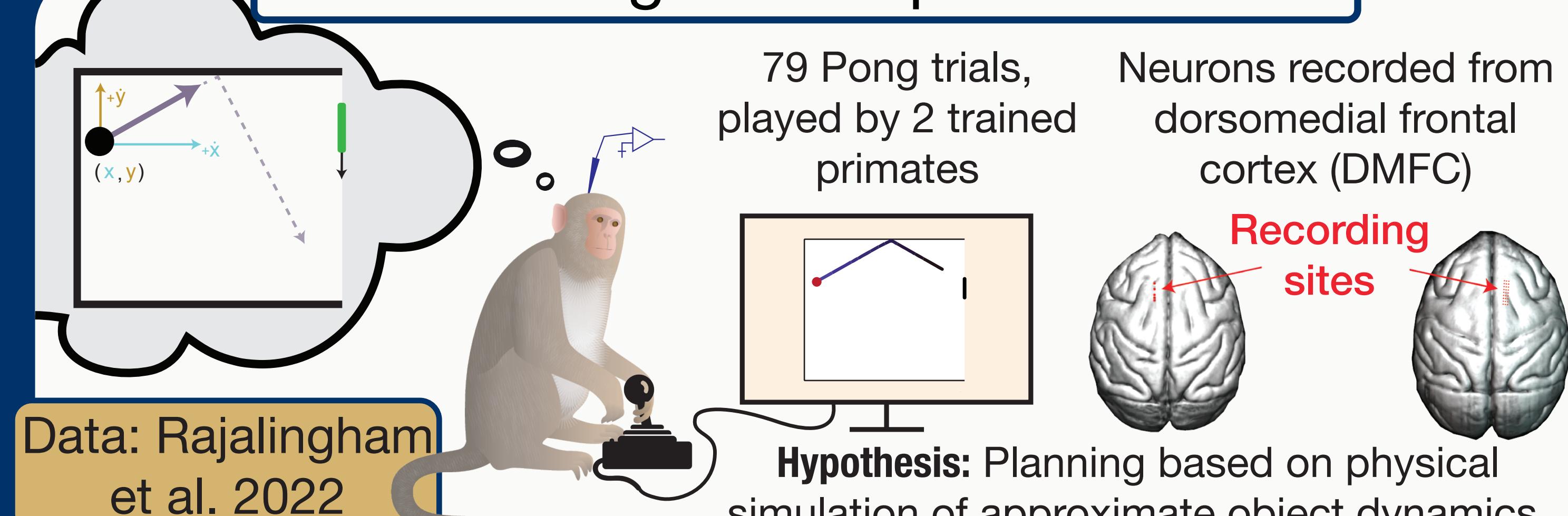
Programmable Neural Networks (Mapable to neural implementation)



Probabilistic models of cognition: Allow for hypothesis-driven exploration but not relatable to neural data.

Task-optimized DNNs: Yield black-box hypothesis as a function of architectures, training sets, and objective functions

Game of Pong as Example World Model



Symbolically Programming (not training) RNNs to Implement Hypothesis

1. Coarse Network Architecture & Equations

$$\begin{aligned} \frac{1}{\gamma} \dot{\mathbf{r}} &= -\mathbf{r} + \mathbf{g} (\mathbf{A}\mathbf{r} + \mathbf{B}\mathbf{x} + \mathbf{d}) \\ \mathbf{o} &= \mathbf{W}\mathbf{r} ; \mathbf{W} = \underset{\mathbf{W}}{\operatorname{argmin}} \|\mathbf{W}\mathbf{r} - \mathbf{o}\| \end{aligned}$$

2. Decomposed Network State (\mathbf{r})

Taylor Series Decomposition

$$\begin{aligned} r_1 &\approx h^* + \frac{\partial h^*}{\partial x_1} x_1 + \frac{\partial h^*}{\partial x_2} x_2 + \dots + \frac{\partial h^*}{\partial x_k} x_k + \frac{\partial^2 h^*}{\partial x_1 \partial x_2} x_1 x_2 + \dots \\ r_2 &\approx \dots \\ r_N &\approx \dots \end{aligned}$$

3. Programmed Output Matrix (\mathbf{O}) (Example: Linear Position Update)

$$o_1 = x_1 + x_3 \quad o_3 = x_1 \\ o_2 = x_2 + x_4 \quad o_4 = x_2$$

4. Scripting & Feedback

$$\frac{1}{\gamma} \dot{\mathbf{r}} = -\mathbf{r} + \mathbf{g} (\tilde{\mathbf{B}} \tilde{\mathbf{W}} \mathbf{r} + \mathbf{B}\mathbf{x} + \mathbf{d}) \quad \mathbf{A} = \tilde{\mathbf{B}} \tilde{\mathbf{W}}$$

$$\tilde{\mathbf{A}} = \tilde{\mathbf{B}} \tilde{\mathbf{W}}$$

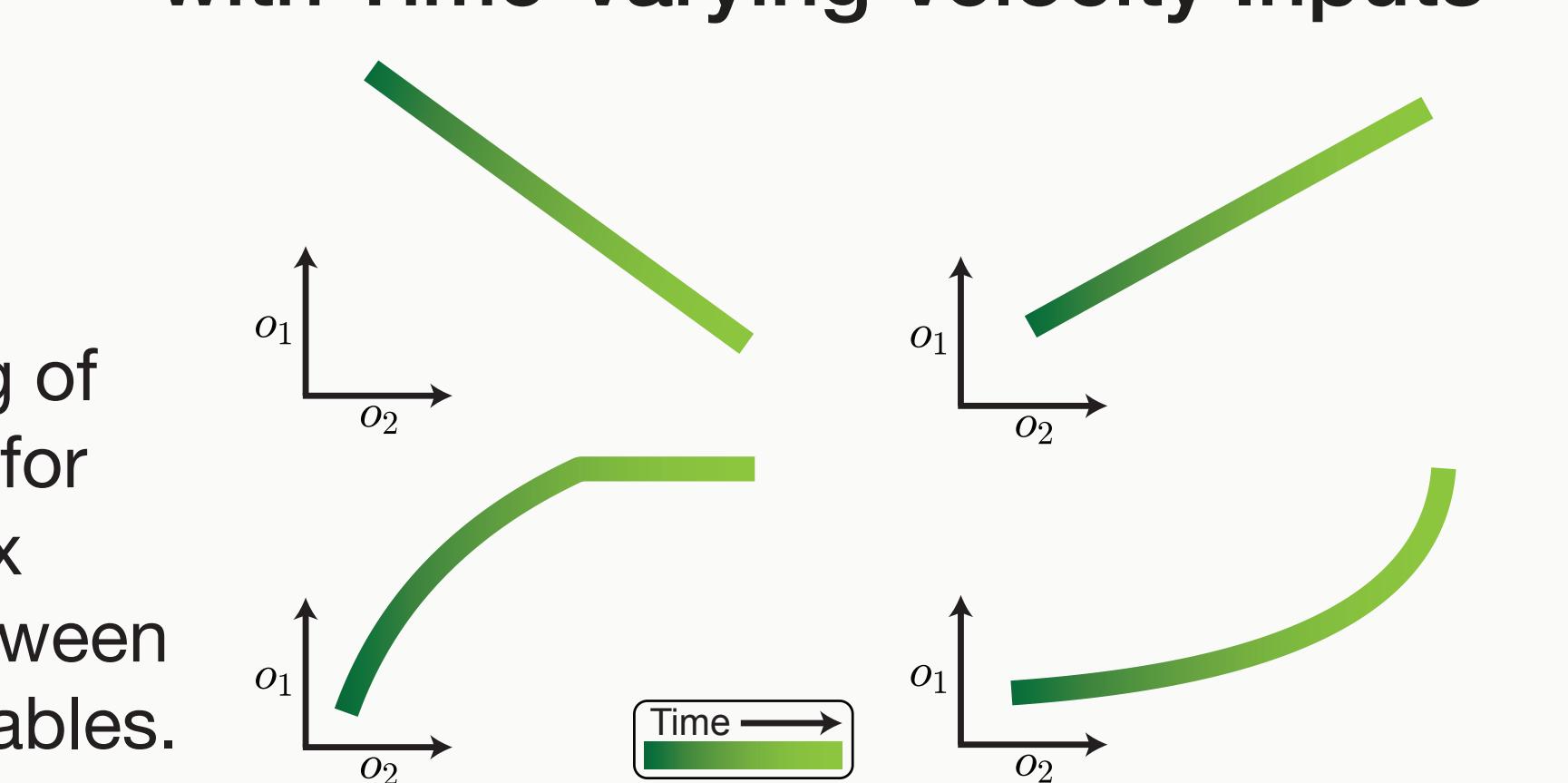
$$\tilde{\mathbf{x}}_1 = \tilde{\mathbf{x}}_1 + \tilde{\mathbf{A}} \tilde{\mathbf{r}} \quad \tilde{\mathbf{x}}_2 = \tilde{\mathbf{x}}_2 + \tilde{\mathbf{A}} \tilde{\mathbf{r}}$$

$$\tilde{\mathbf{x}}_3 = \tilde{\mathbf{x}}_3 + \tilde{\mathbf{A}} \tilde{\mathbf{r}} \quad \tilde{\mathbf{x}}_4 = \tilde{\mathbf{x}}_4 + \tilde{\mathbf{A}} \tilde{\mathbf{r}}$$

$$o_1 = \begin{matrix} \text{---} & \text{---} & \dots & \text{---} \\ \text{---} & \text{---} & \dots & \text{---} \\ \text{---} & \text{---} & \dots & \text{---} \end{matrix} \quad o_2 = \begin{matrix} \text{---} & \text{---} & \dots & \text{---} \\ \text{---} & \text{---} & \dots & \text{---} \\ \text{---} & \text{---} & \dots & \text{---} \end{matrix} \quad \dots$$

$$o_3 = \begin{matrix} \text{---} & \text{---} & \dots & \text{---} \\ \text{---} & \text{---} & \dots & \text{---} \\ \text{---} & \text{---} & \dots & \text{---} \end{matrix} \quad o_4 = \begin{matrix} \text{---} & \text{---} & \dots & \text{---} \\ \text{---} & \text{---} & \dots & \text{---} \\ \text{---} & \text{---} & \dots & \text{---} \end{matrix}$$

5. Simulated Trajectories from RNN with Time-varying Velocity Inputs



Feedback allows for the linking of simple programmed functions for more complex dynamics between symbolic variables.

Method: Kim & Bassett 2023

Conclusions

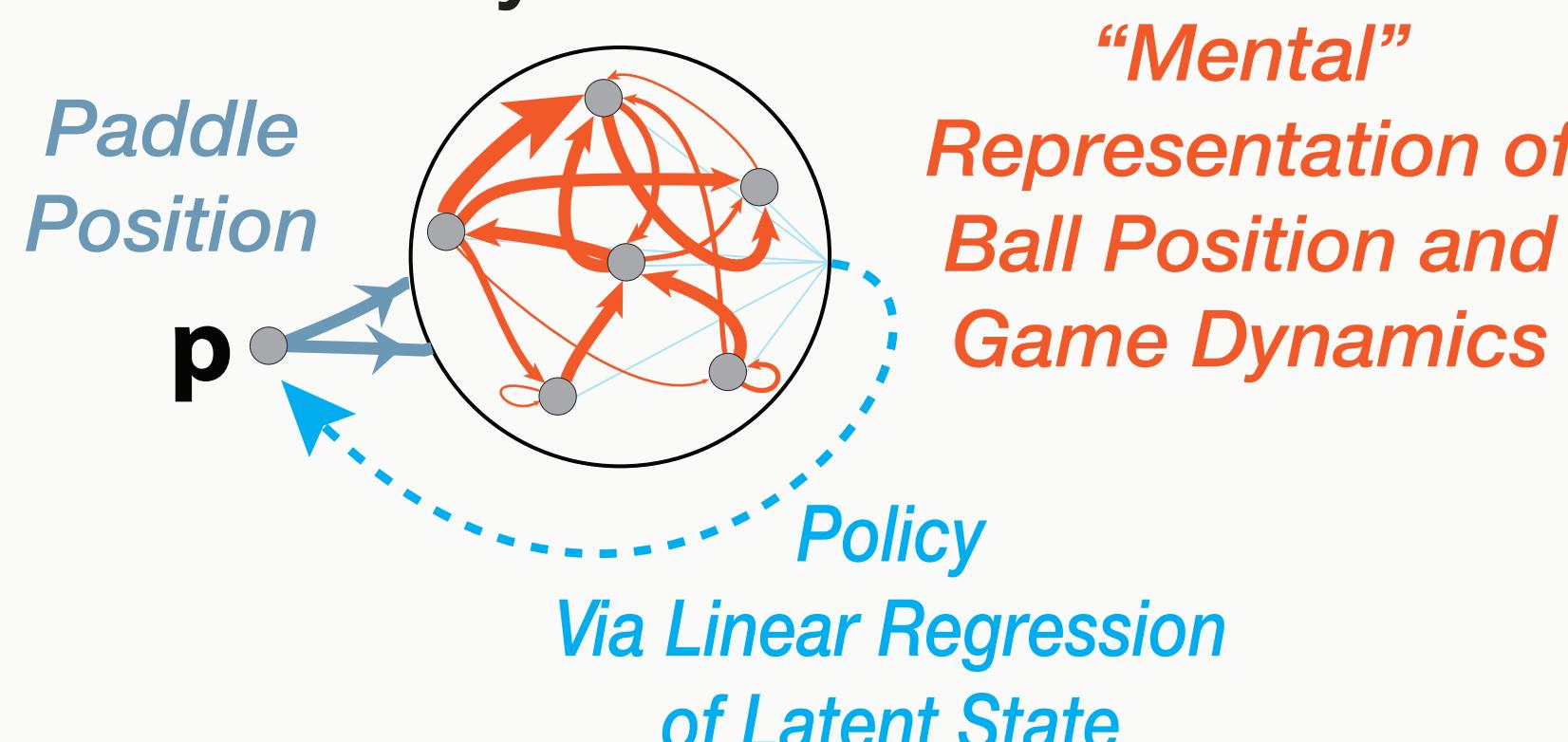
- Evidence DMFC neurons may be simulating object dynamics
- Our findings highlight a new multi-level modeling framework and its potential for exploring world models in the brain.

Programmed RNN World Model of Pong

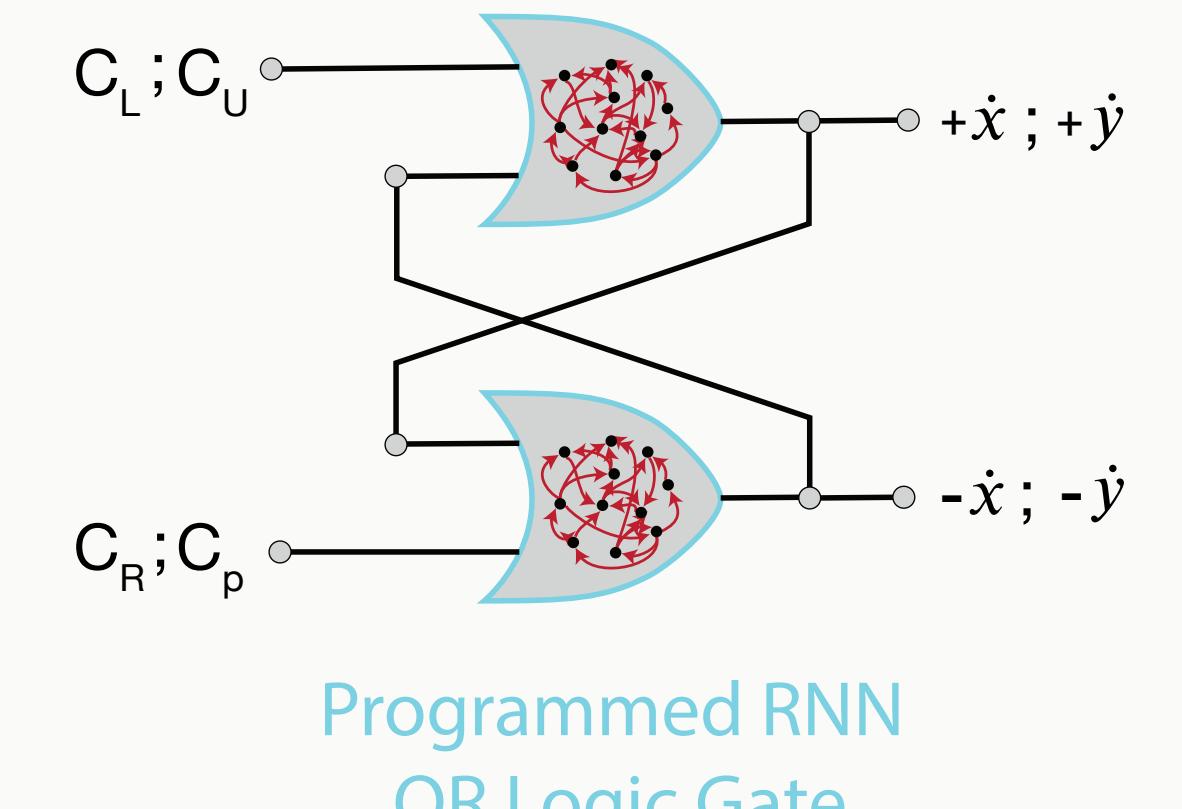
Hypothesized World Model Codifying Pong Dynamics

$$\begin{aligned} x &= c_u \\ +\dot{x} &= c_l \\ -\dot{x} &= c_r \\ y &= c_p \\ +\dot{y} &= m_x(x + \dot{x})/\gamma \\ -\dot{y} &= m_y(y + \dot{y})/\gamma \\ C_u &= m_x(x + \dot{x})/\gamma \\ C_l &= m_x(-\dot{x})/\gamma \\ C_r &= m_y(-\dot{y})/\gamma \\ C_p &= m_y(y + \dot{y})/\gamma \\ m_x &= \dots \\ m_y &= \dots \\ p &= \dots \end{aligned}$$

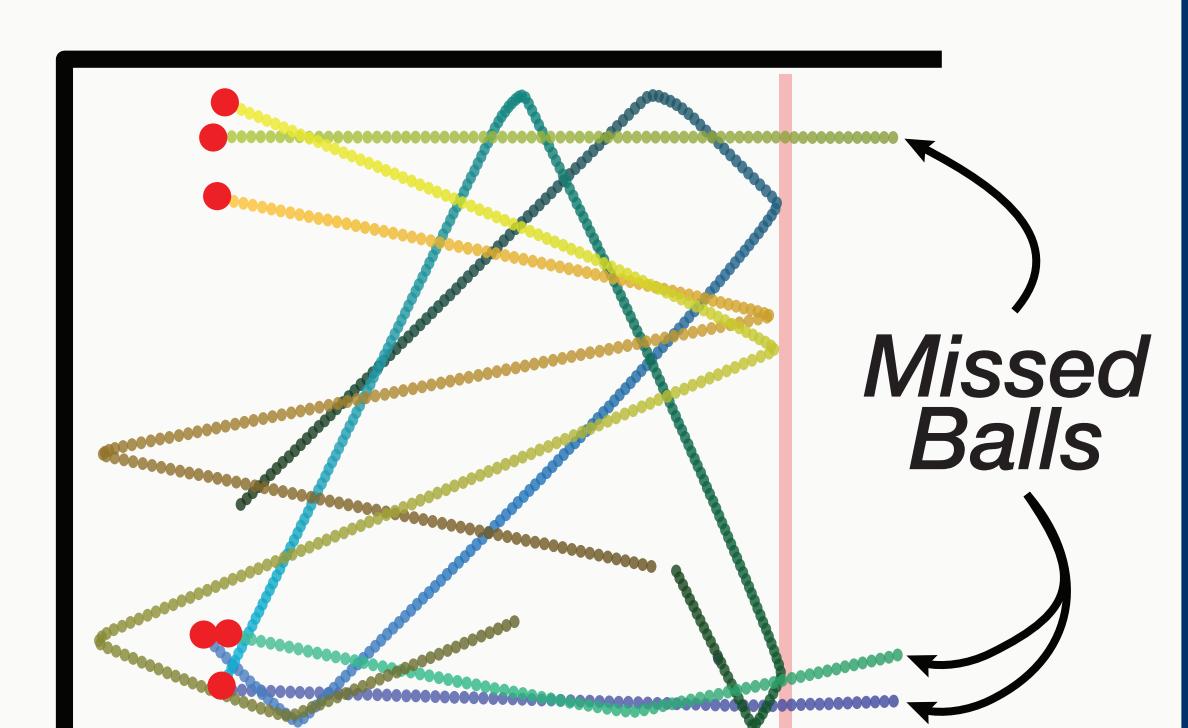
Simple Model-based Policy with Linear



Collision Detection & Resolution via Programmed Neural S-R Latch Circuit

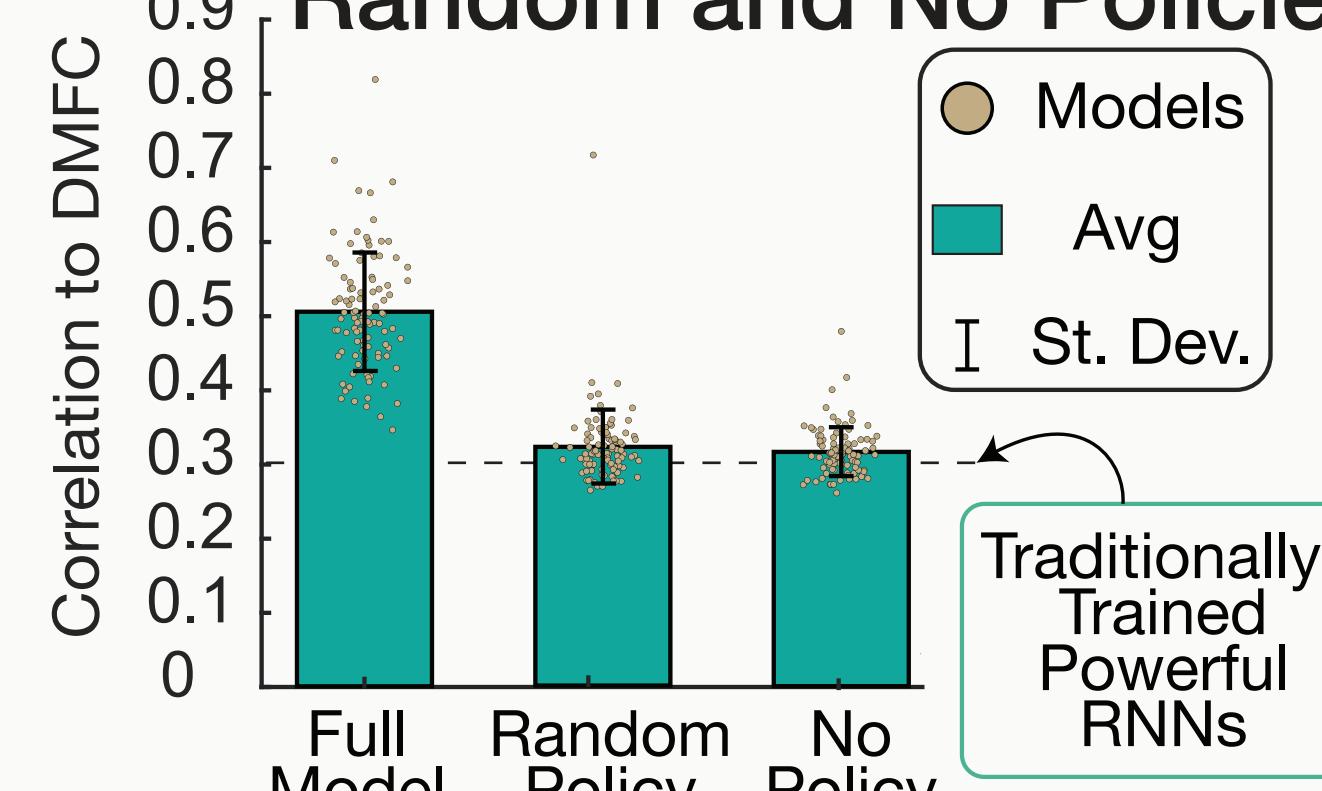


Example Gameplay

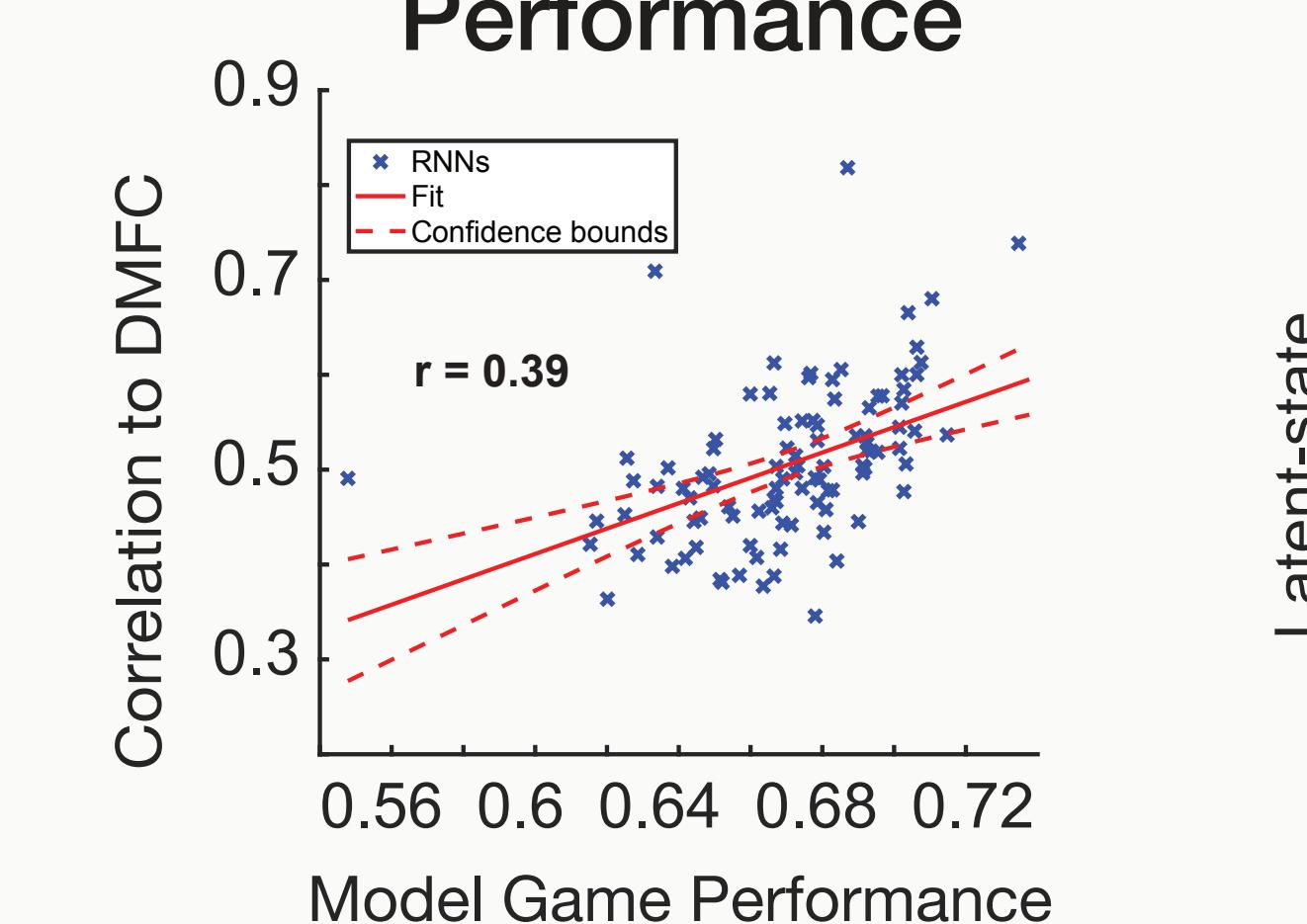


Results

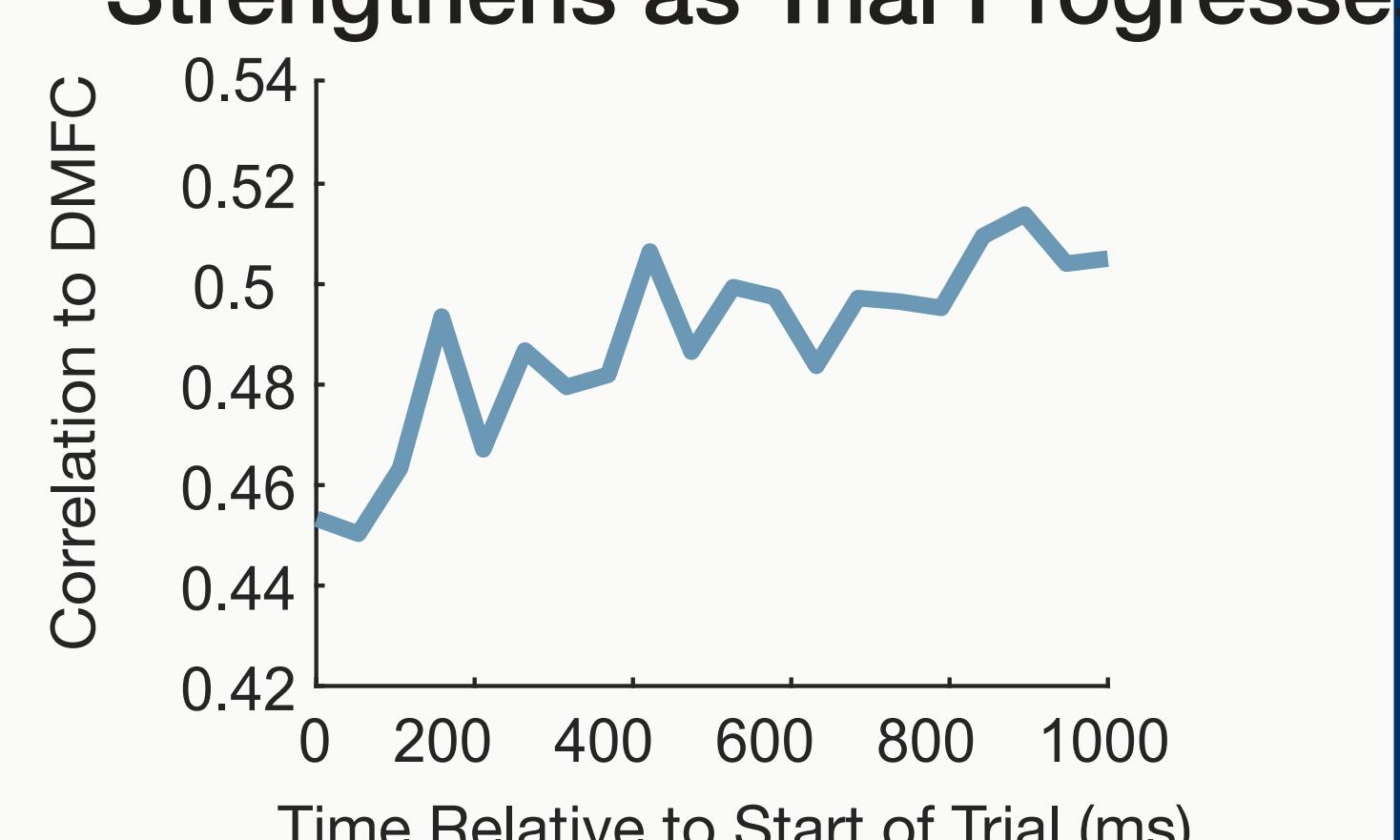
Programmed RNN explains DMFC Activity Better Than Random and No Policies



Model-Brain Similarity is Correlated with Model's Game Performance



Similarity to DMFC Strengthens as Trial Progresses



Non-linearity Seen in Prediction Using Programmed RNN and DMFC but Not in Traditional RNNs

