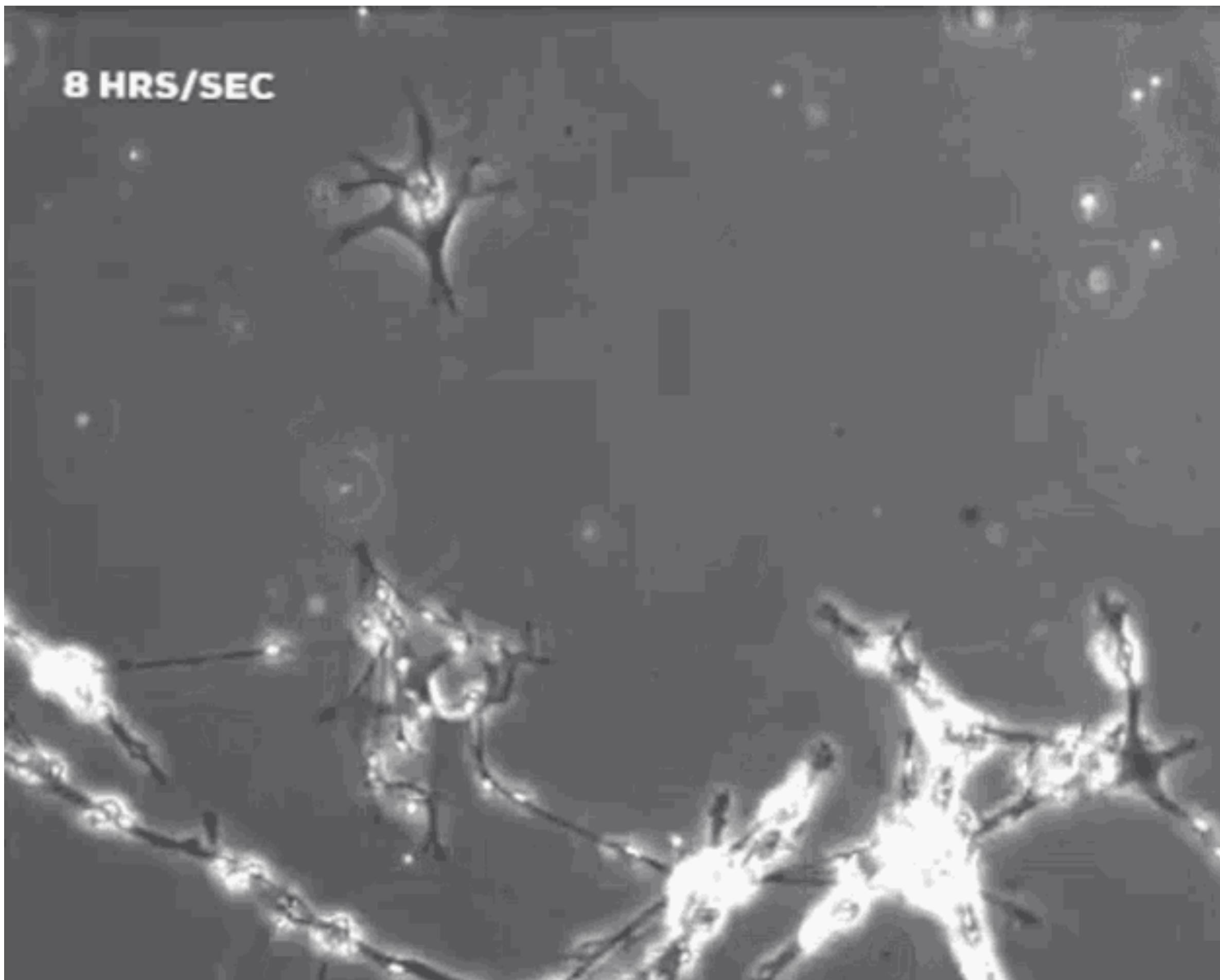


# Thinking Machines

Abstract thought as internal world model exploration

# Intelligence is a fascinating phenomenon.



= ChatGPT 4 

 You  
How do I set up my CUDA drivers on my new deep learning rig?  
 3 / 3 

 ChatGPT  
Setting up CUDA drivers on a new deep learning rig involves several steps, ensuring that your system is ready to leverage the power of NVIDIA GPUs for deep learning tasks. Here's a streamlined guide to get you started, but remember, specifics can vary based on your operating system (OS), the NVIDIA GPU model you have, and the current versions of software and drivers availab.  


 Message ChatGPT... 

ChatGPT can make mistakes. Consider checking important information.

# Building AI can help us understand intelligence.



What I cannot create, I do not understand.

— *Richard P. Feynman* —

# Building AI can help us understand intelligence.



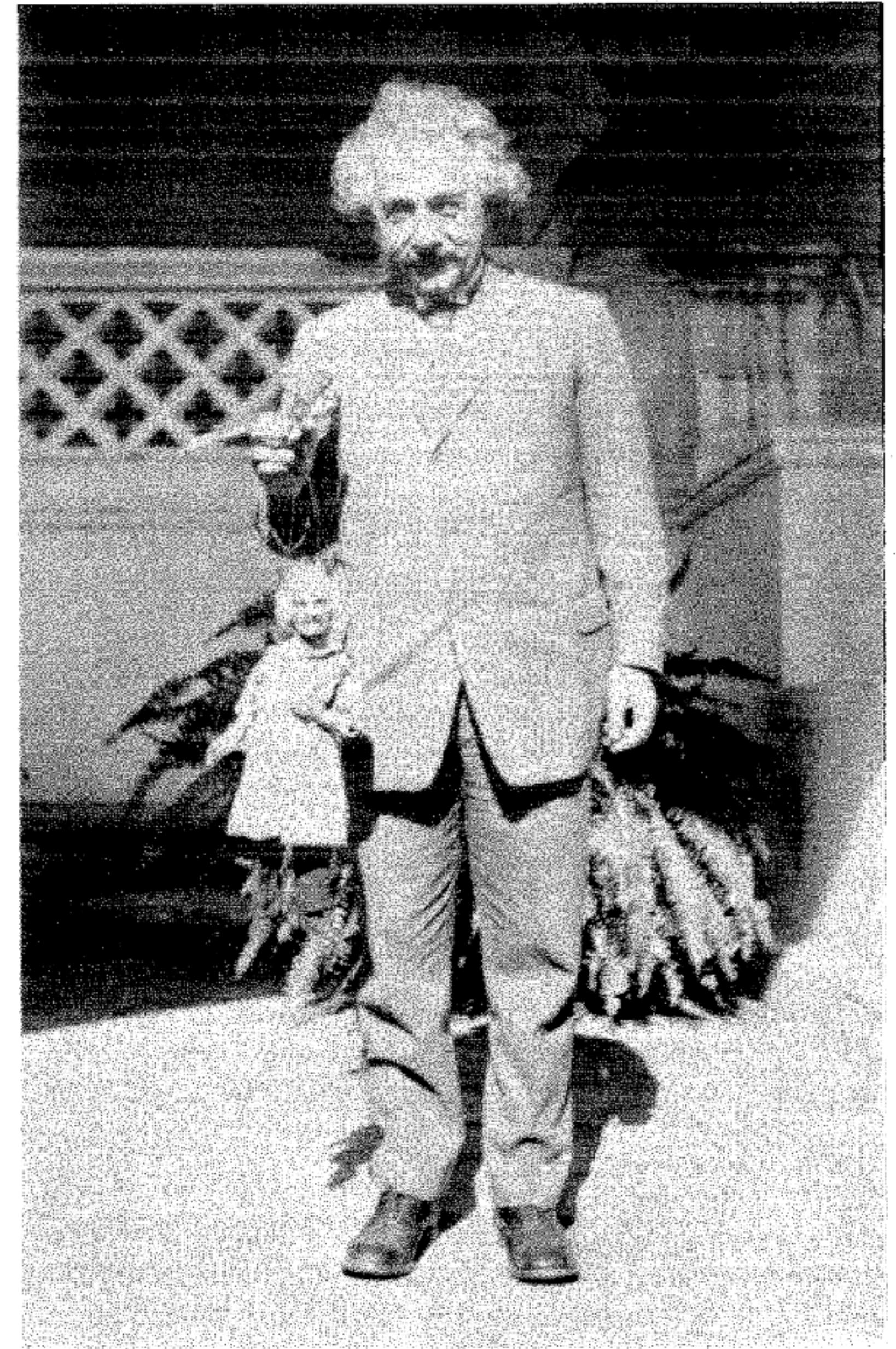
# LLMs exhibit aspects of intelligence.

## *Zero-shot learning miracle*

- **Knowledge Retrieval:** “*The Titanic sank in the year [MASK].*” (Answer: “1912”)
- **Reasoning:** “*A is taller than B. B is taller than C. Is A taller than C? Answer: [MASK]*” (Answer: “Yes”)
- **Sentiment Analysis:** “*I am sad today. The sentiment of the previous sentence was [MASK]*” (Answer: “Negative”)

**“The human mind has first to construct forms, independently, before we can find them in things.”**

**— Albert Einstein**



Two Einsteins in front of the brand-new Athenaeum in 1931 — Albert himself and a puppet from the play called *Mr. Noah*. Albert said the puppet was “good but not fat enough.”

# Thinking Machines

## Motivation

What is **thinking**?

How can we **formalize** thinking (e.g.,  
mathematically/informationally)?

How can we **build thinking machines**?

# Thinking Machines

## Motivation

What is thinking?

How can we **formalize** thinking (e.g.,  
mathematically/informationally)?

How can we **build** thinking machines?

**Why can't ChatGPT automate  
my entire PhD yet?**

# Thinking Machines

## Agenda

1. LLMs predict the next word (token).
2. World models emerge from predictive coding.
3. **BIG IDEA: Thinking = Internal World Model Exploration!**
4. Register token architecture for “Thinking LLM”.
5. Beyond LLMs – the end-game for CNS.

# LLMs predict the next token.

Try to predict the next token!

[22170, 311, 7168, 279, 1828, 4037, 0]

$x_1$

$x_2$

$x_3$

$x_4$

$x_5$

$x_6$

$x_7$

$$P_{\theta}(x_{t+1} \mid x_1, \dots, x_t)$$

# LLMs predict the next token.

Try to predict the next token!

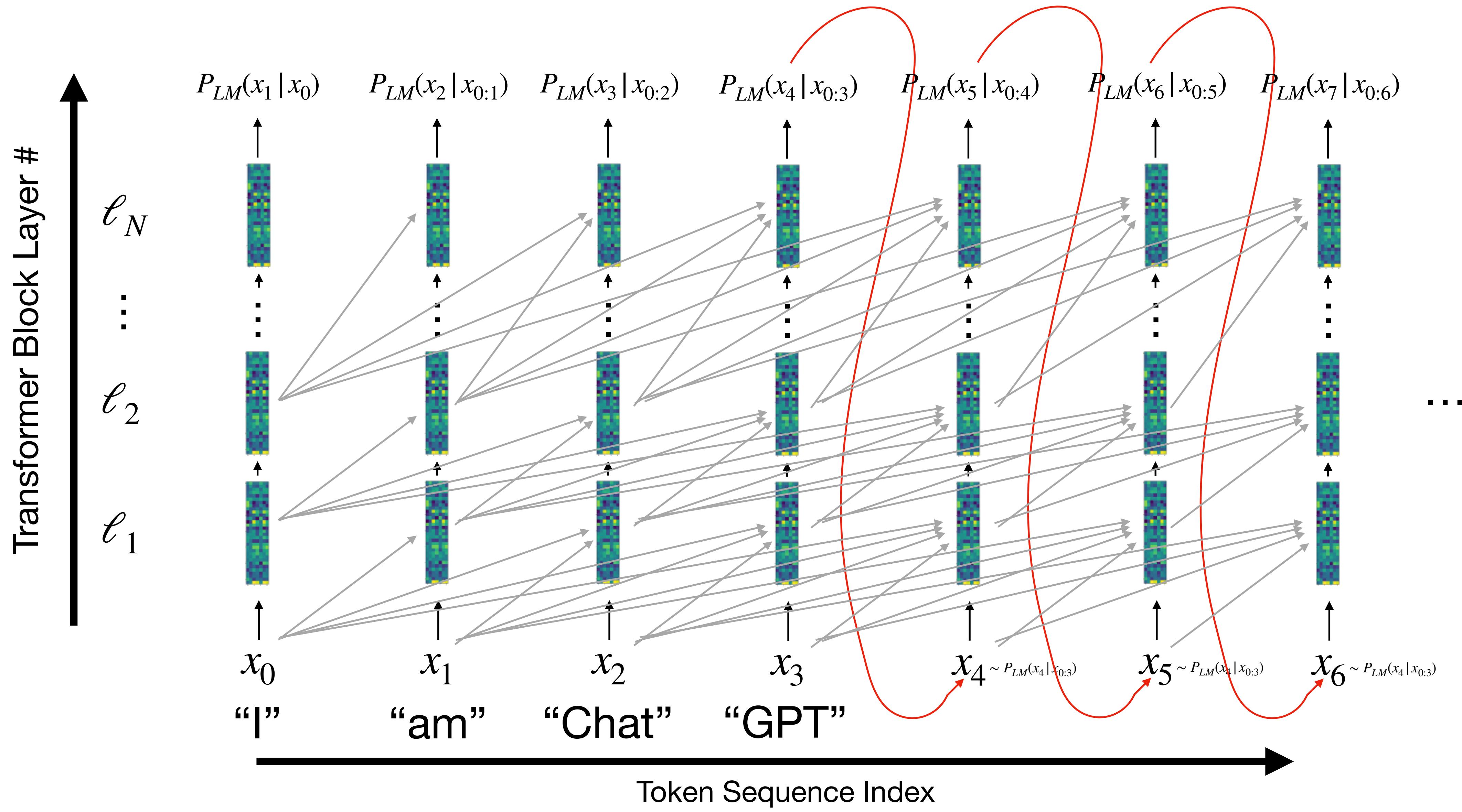
[22170, 311, 7168, 279, 1828, 4037, 0]

$x_1$        $x_2$        $x_3$        $x_4$        $x_5$        $x_6$        $x_7$

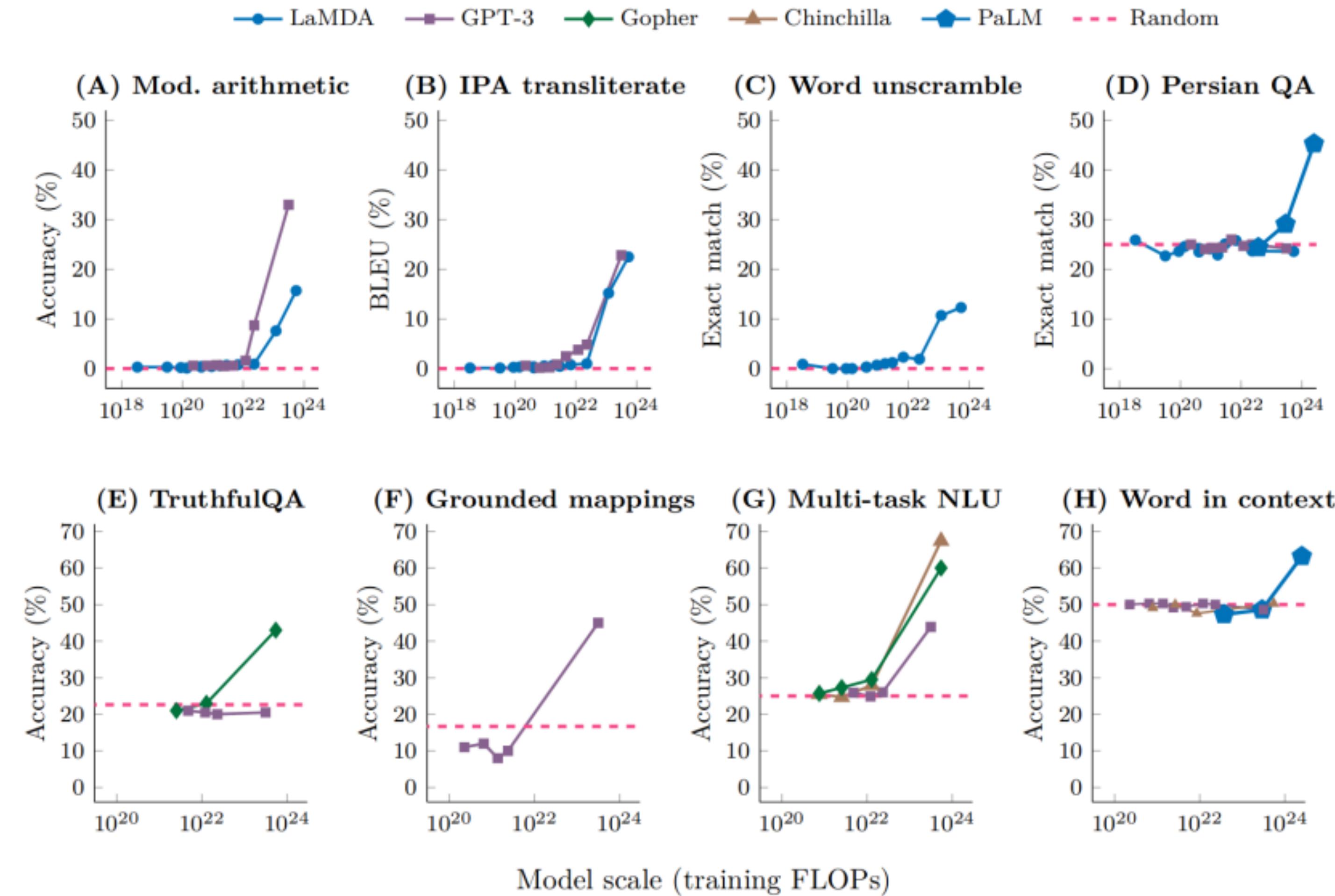
$$\theta = \arg \max_{\theta} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \left[ \sum_{i=1}^N \log P_{\theta}(x_i | x_1 \dots x_{i-1}) \right]$$

  
 $\log P_{\theta}(x_1 \dots x_N)$

# LLMs predict the next token.



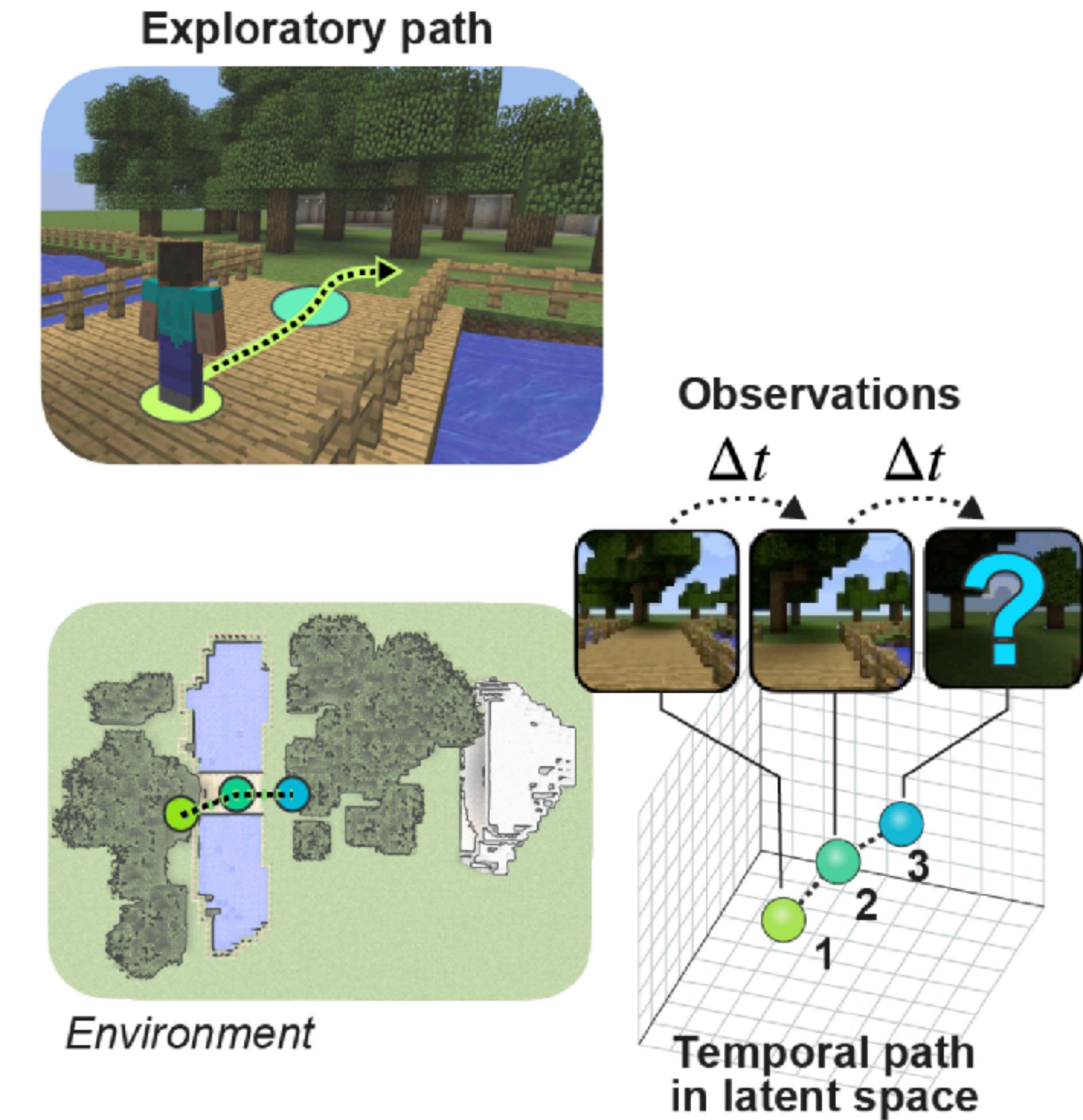
# Bigger LLM → Better Prediction → Greater Capability



From “Are Emergent Abilities of Large Language Models a Mirage?” (Schaeffer et al) — <https://arxiv.org/abs/2304.15004>

# World Models emerge from Predictive Coding.

## James Gornet (Thomson Lab, CNS)



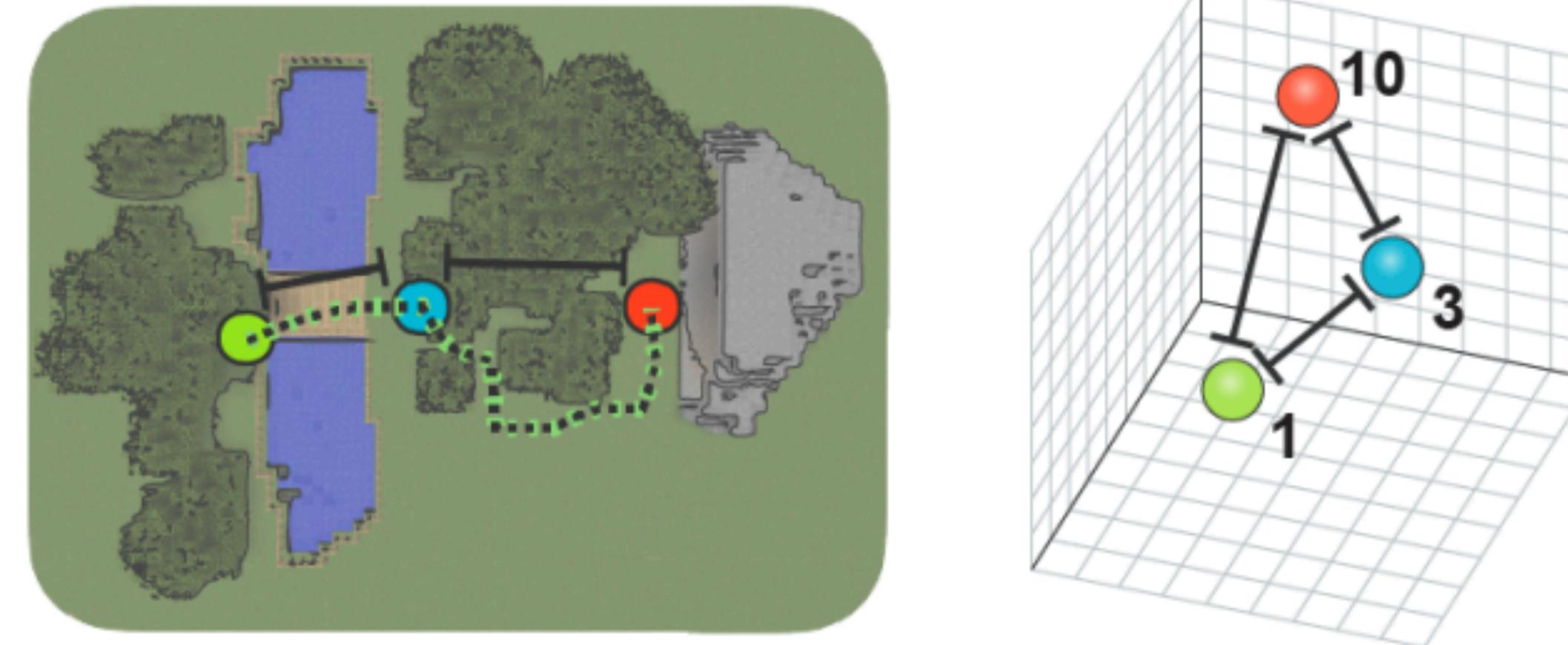
James!



# World Models emerge from Predictive Coding.

## James Gornet (Thomson Lab, CNS)

$$P_{\theta}(x_{t+1} \mid x_1, \dots, x_t)$$



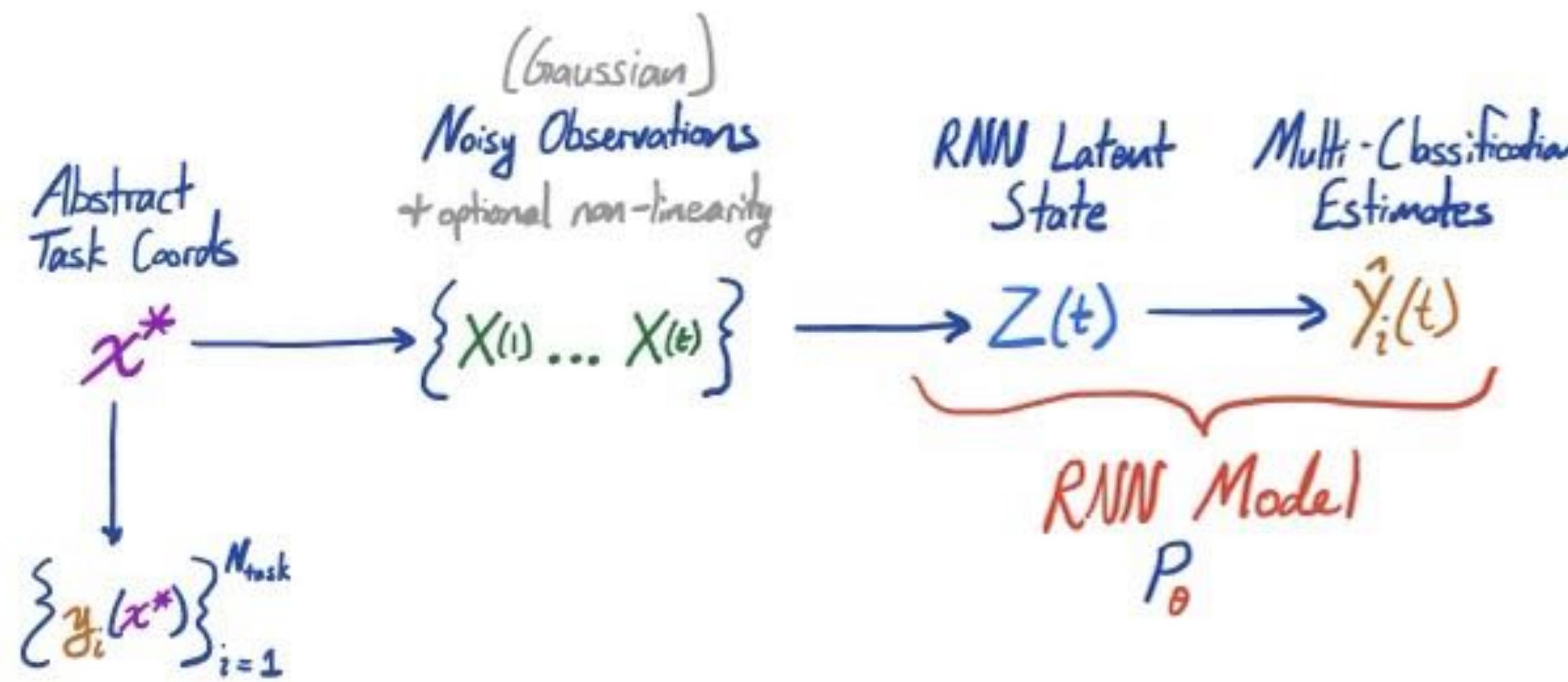
Distance in latent space maps environment

James!



# World Models emerge from Evidence Aggregation.

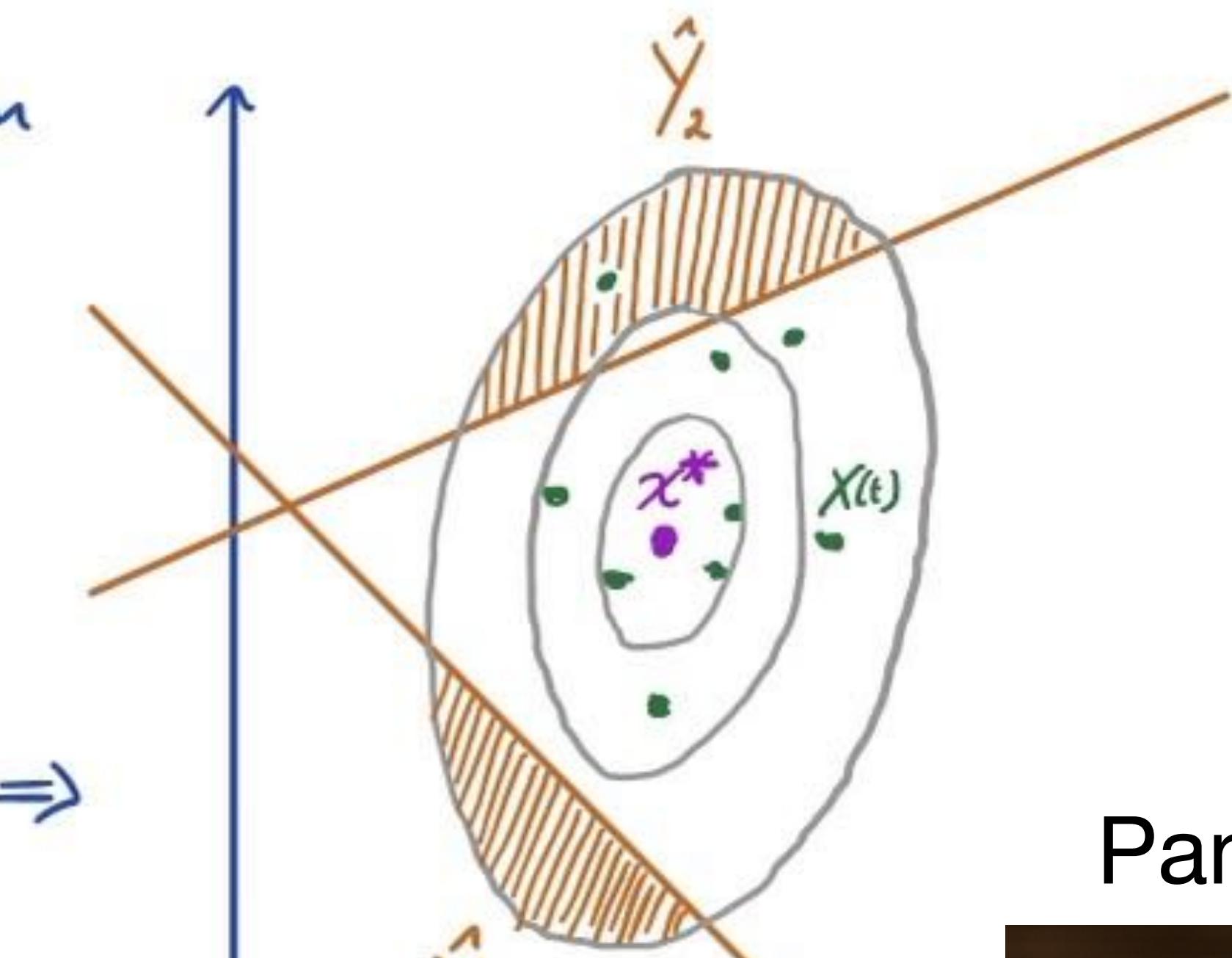
## With Pantelis Vafidis (Rangel Lab, CNS)



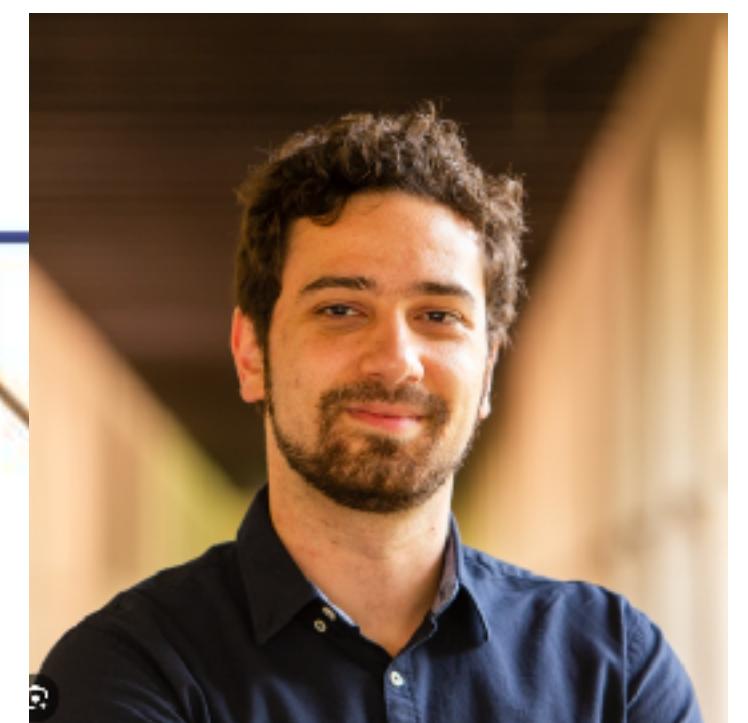
Generalization Theorem: Optimal estimation of  $\hat{y} \Rightarrow$

$Z(t)$  contains optimal  $x^*$  estimate based on  $X(1) \dots X(t)$   
if  $N_{task} \geq \dim(x^*)$ .

[Bonus] Gaussian noise + tanh map from  $Z(t) \rightarrow \hat{y}_i$   
 $\Rightarrow \hat{x}^* = \text{lin func}(Z(t))$



Pantelis!



# World Models emerge from Predictive Coding.

- Latent  $\mathbf{Z}(t)$  has provably high mutual information with world  $I(\mathbf{Z}(t); \mathbf{x}^*(t))$  for a variety of prediction tasks.

Abstraction/thought/language:

Latent/Brain State:

Noisy/Partial Observations:

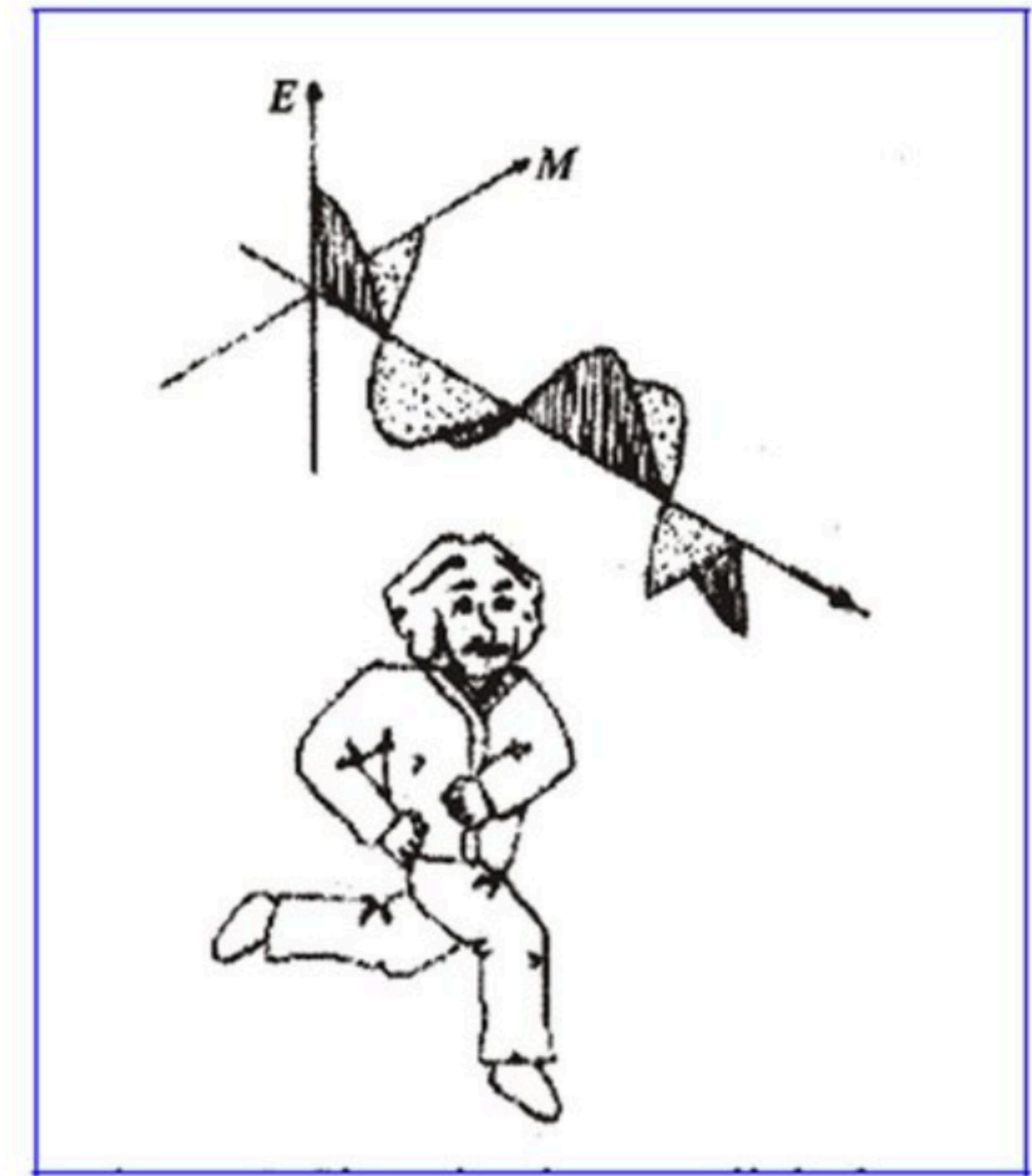
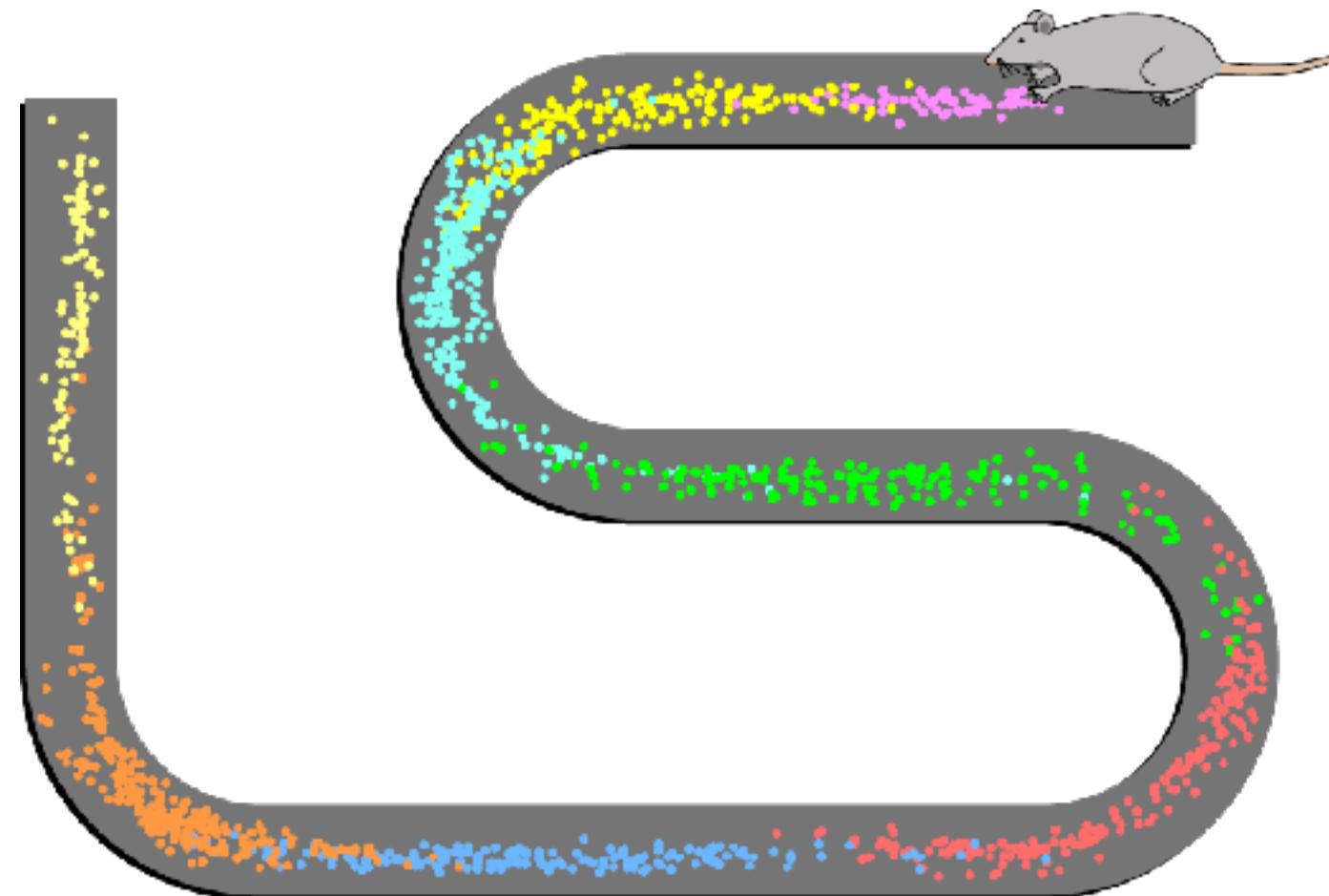
World State:



# Thinking is internal world model exploration.

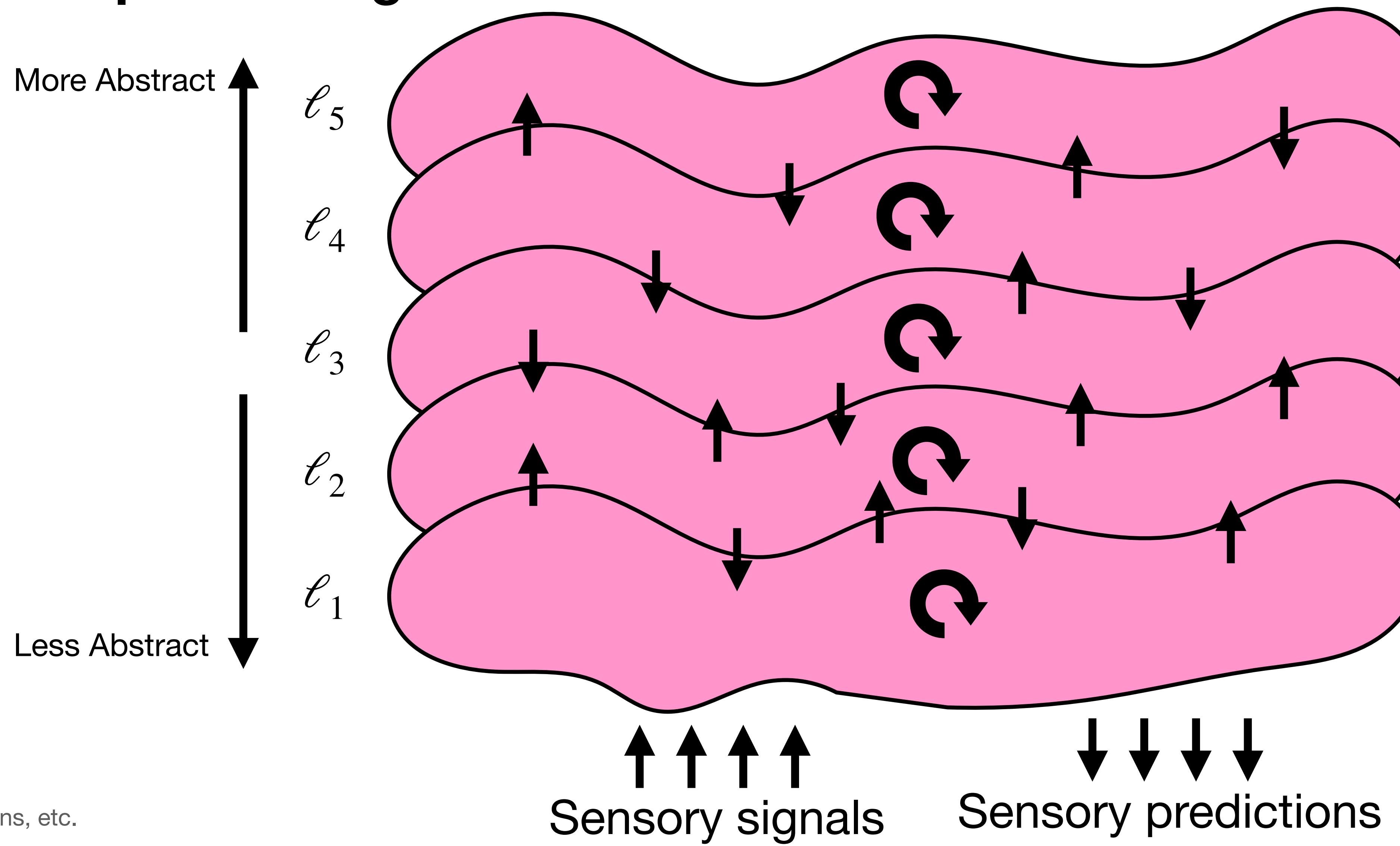
## Leveraging *implicit understanding*.

- **Neuro:** We observe rodents “mentally exploring” a maze in place cell activations.
- **Phil of Mind:** Thought experiments, hypotheticals, imagined movement/environments.



# Thought as internal world model exploration

## Cortical processing cartoon sketch



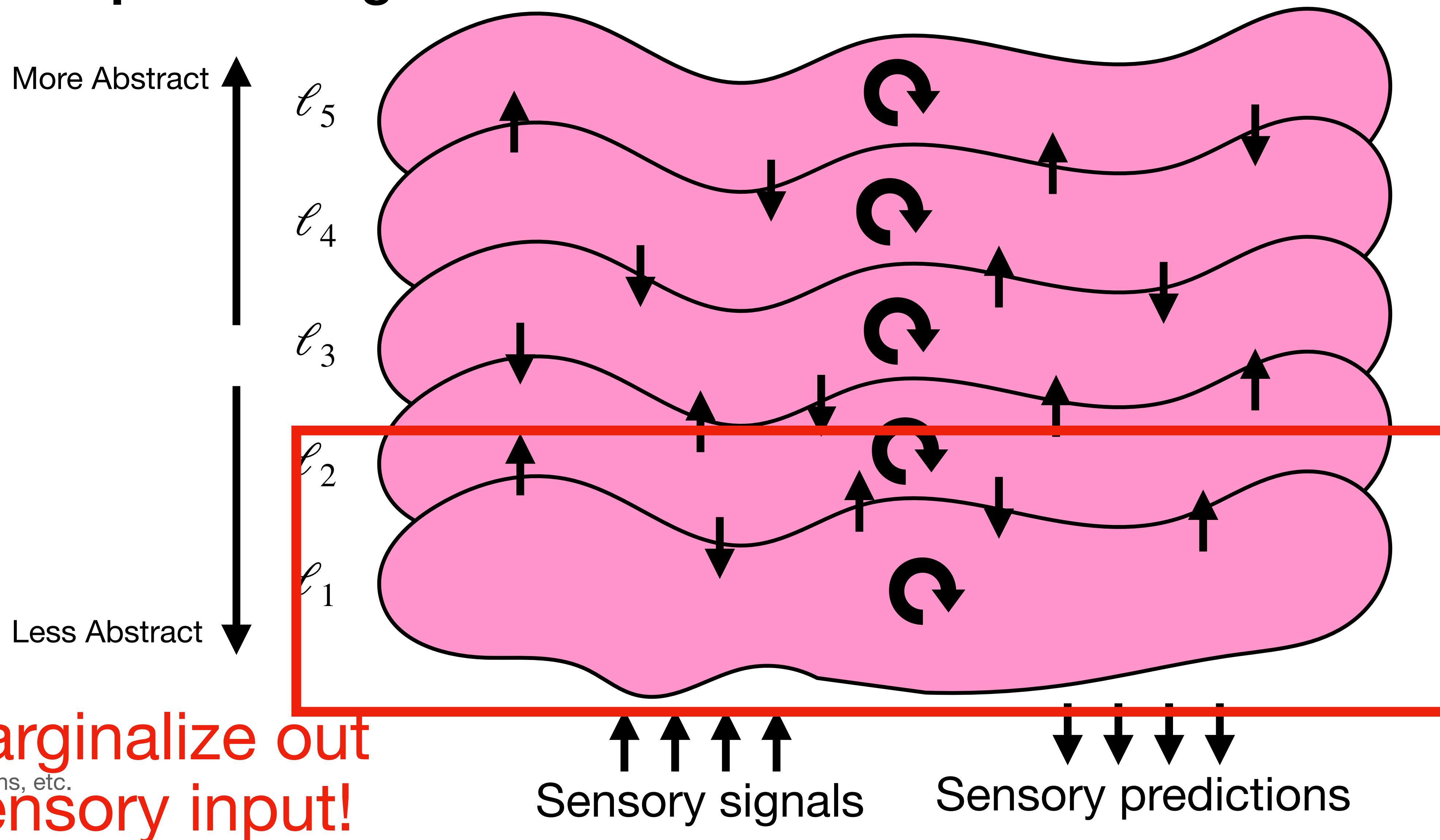






# Thought as internal world model exploration

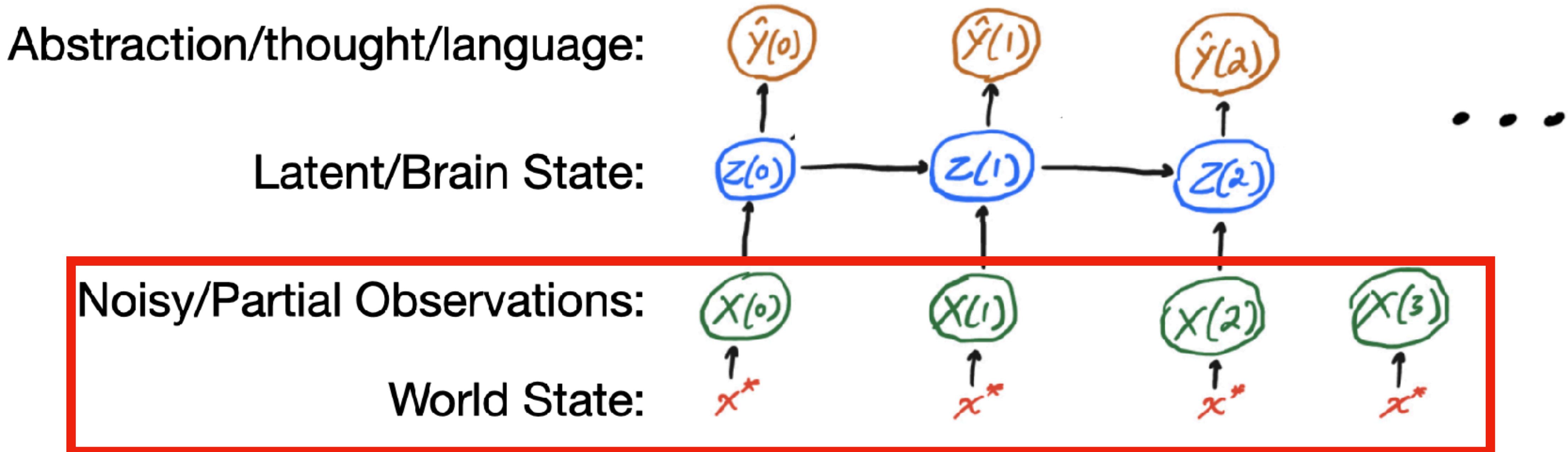
## Cortical processing cartoon sketch



# World Models emerge from Predictive Coding.

- Latent  $\mathbf{Z}(t)$  has provably high mutual information with world  $I(\mathbf{Z}(t); \mathbf{x}^*(t))$  for a variety of prediction tasks.

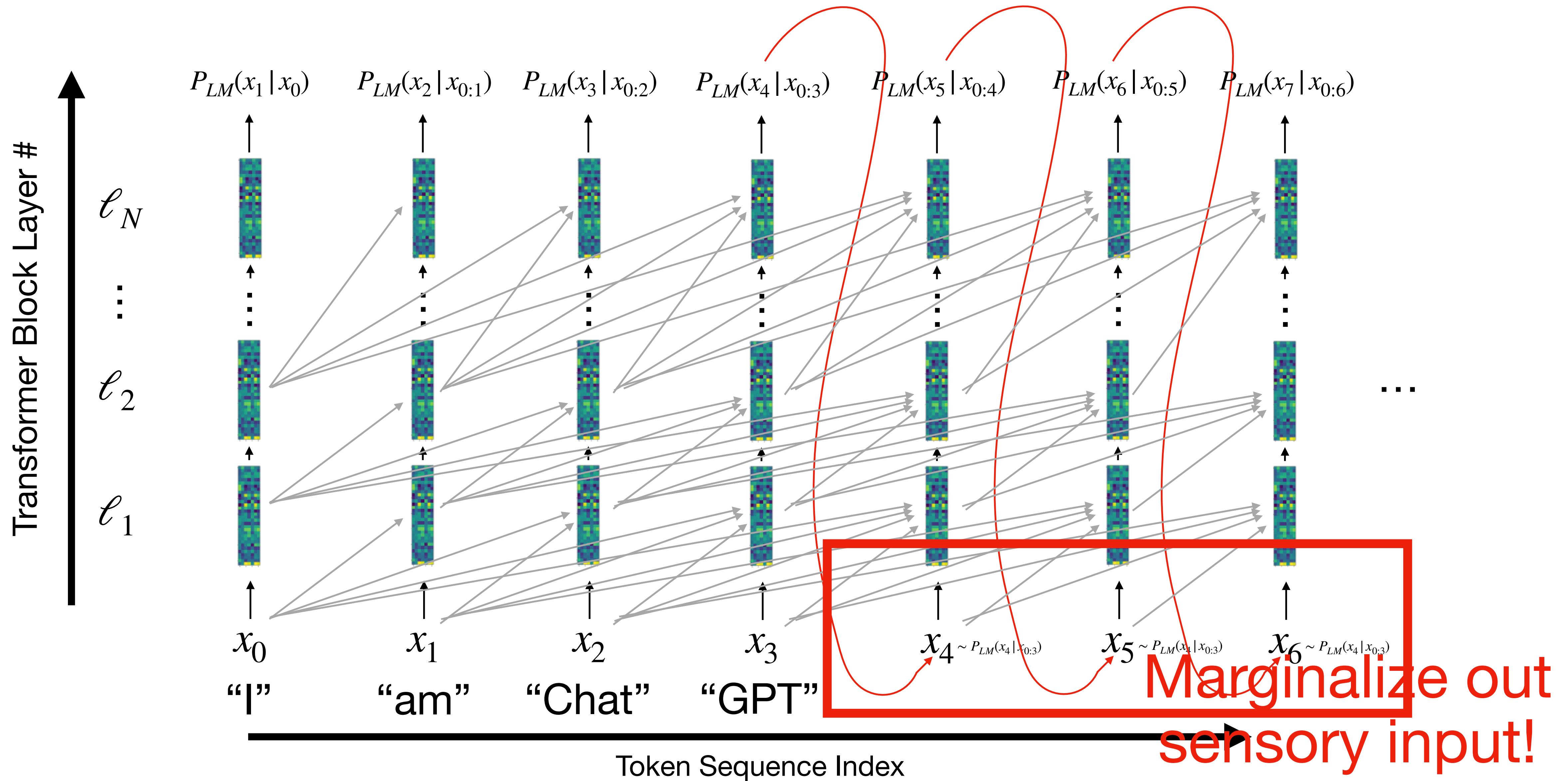
Abstraction/thought/language:



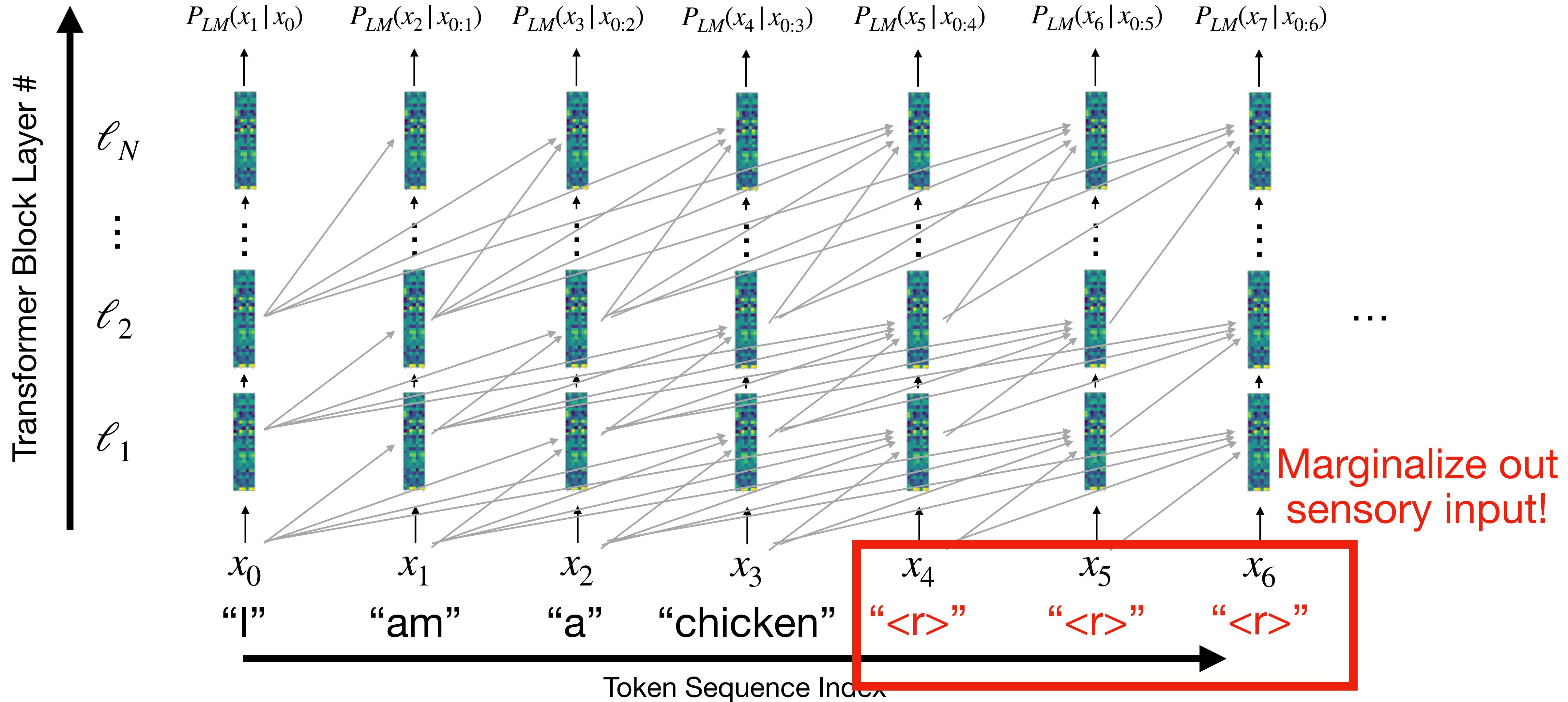
Marginalize out sensory input!

# Thinking LLMs via Register Tokens

# Register Tokens enable LLM world model exploration.



# Register Tokens enable LLM world model exploration.



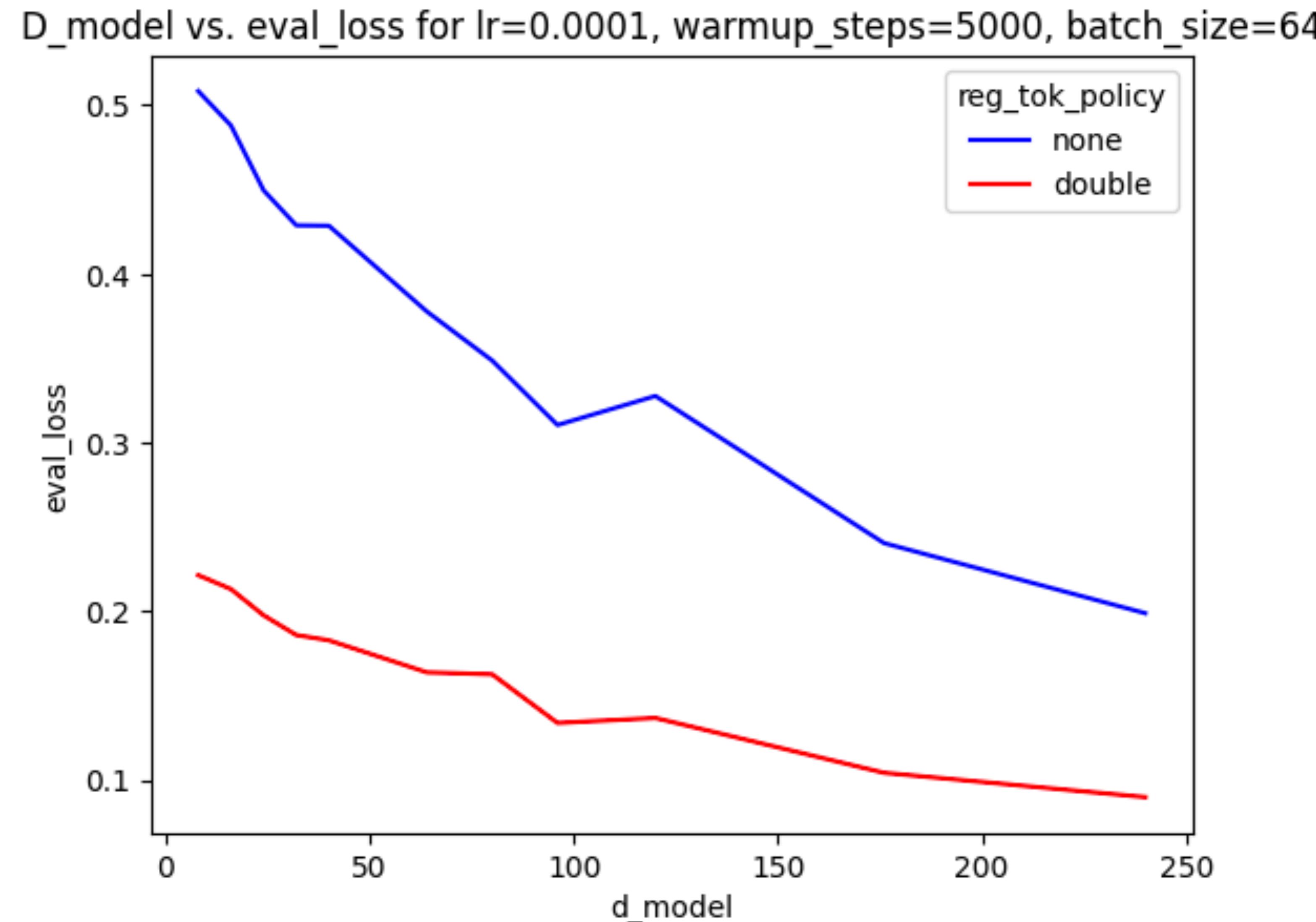
# Register tokens improve arithmetic performance.

```
1  Dataset = [  
2      {  
3          "Question": "8 * 98 = <r><r>...<r>",  
4          "Answer": "784"  
5      },  
6      {  
7          "Question": "32 * 98891 = <r><r>...<r>",  
8          "Answer": "3164512"  
9      },  
10 ]
```

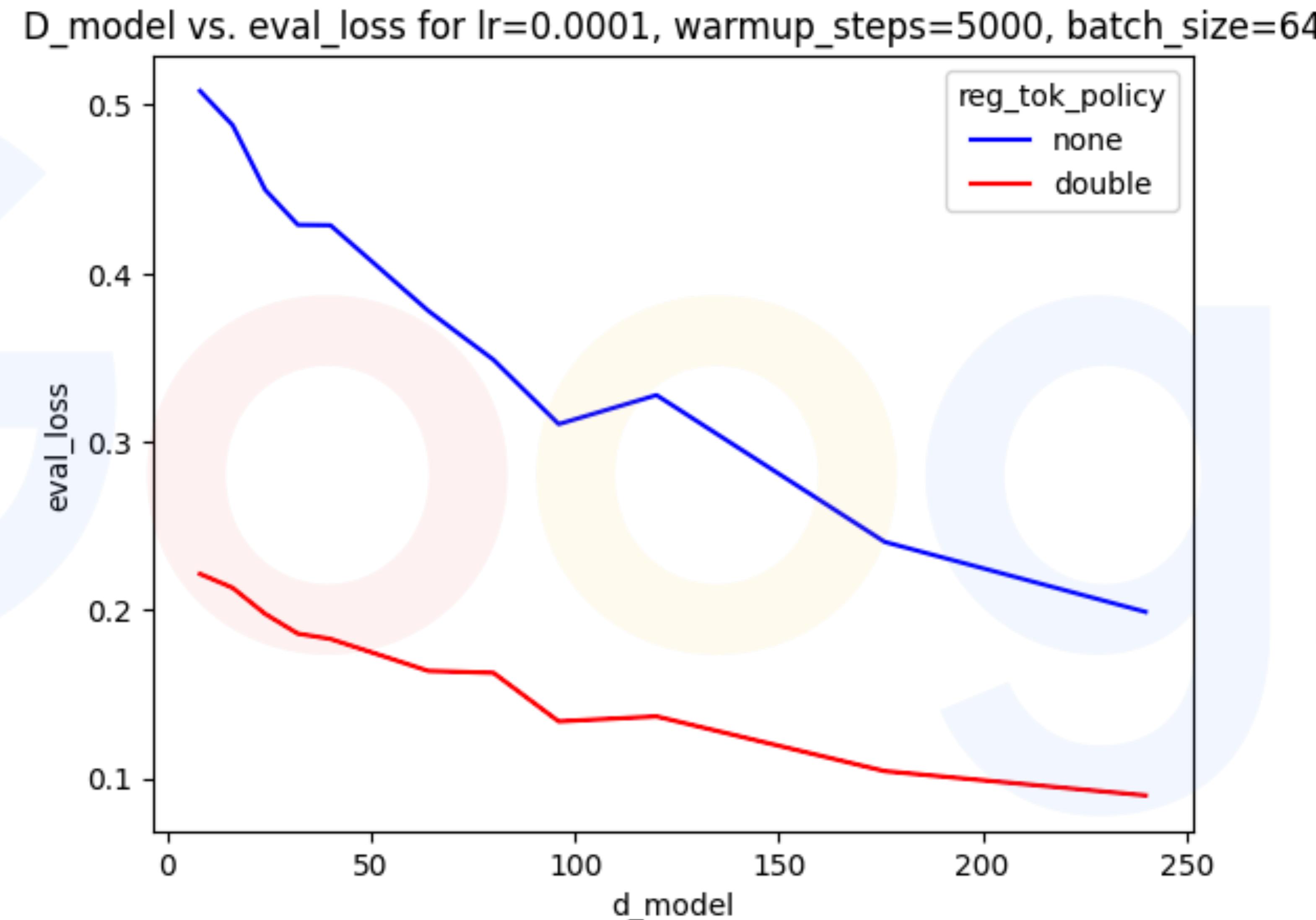
# Register tokens improve arithmetic performance.

```
1  Dataset = [  
2      {  
3          "Question": "8 * 98 = ",  
4          "Answer": "784"  
5      },  
6      {  
7          "Question": "32 * 98891 = ",  
8          "Answer": "3164512"  
9      },  
10 ]
```

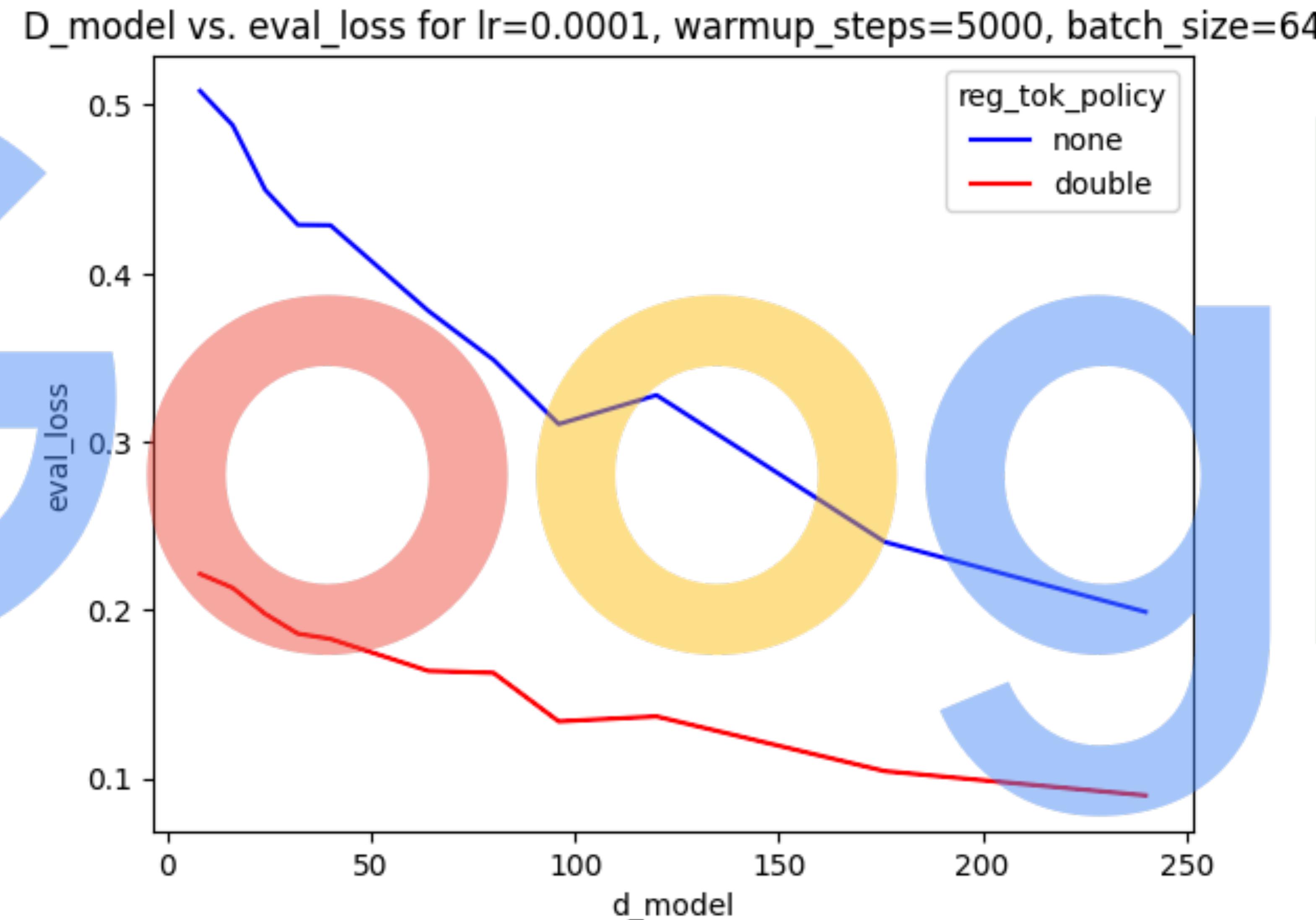
# Register tokens improve arithmetic performance.



# Register tokens improve arithmetic performance.



# Register tokens improve arithmetic performance.



# Register tokens improve arithmetic performance.

Published as a conference paper at ICLR 2024

---

## THINK BEFORE YOU SPEAK: TRAINING LANGUAGE MODELS WITH PAUSE TOKENS

**Sachin Goyal\***<sup>†</sup>

Machine Learning Department  
Carnegie Mellon University  
[sachingo@andrew.cmu.edu](mailto:sachingo@andrew.cmu.edu)

**Ziwei Ji**

Google Research, NY  
[ziwei.ji@google.com](mailto:ziwei.ji@google.com)

**Ankit Singh Rawat**

Google Research, NY  
[ankitsrawat@google.com](mailto:ankitsrawat@google.com)

**Aditya Krishna Menon**

Google Research, NY  
[adityakmenon@google.com](mailto:adityakmenon@google.com)

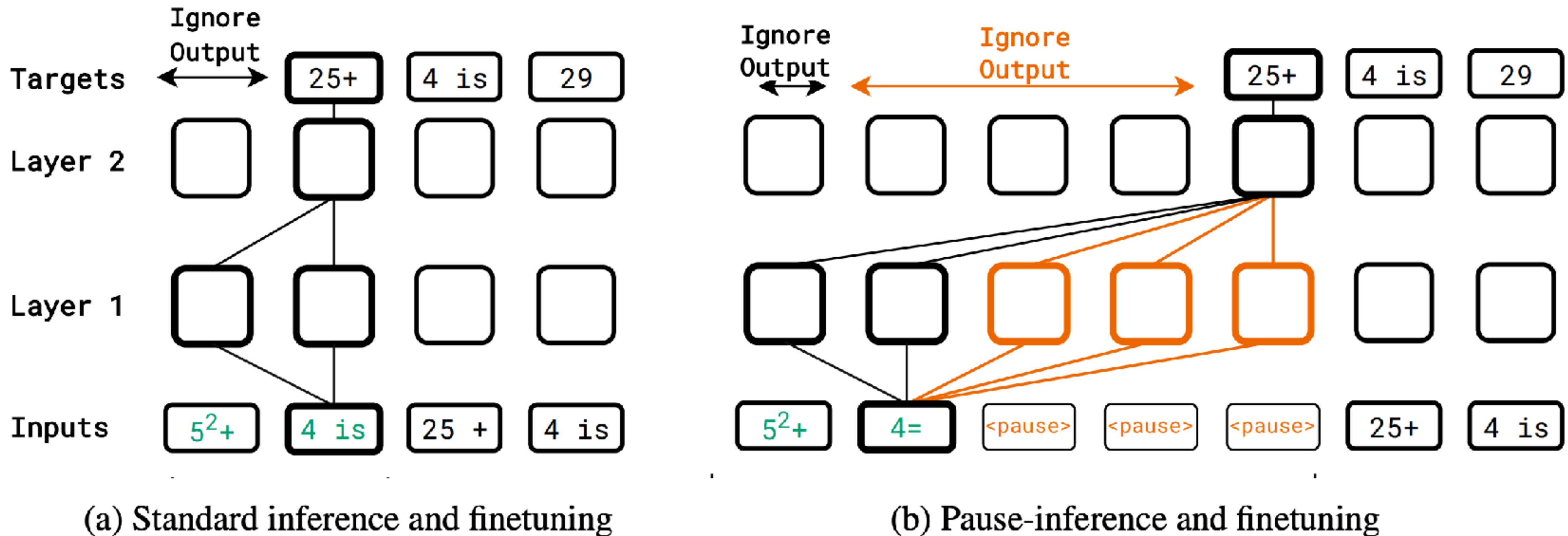
**Sanjiv Kumar**

Google Research, NY  
[sanjivk@google.com](mailto:sanjivk@google.com)

**Vaishnavh Nagarajan<sup>†</sup>**

Google Research, NY  
[vaishnavh@google.com](mailto:vaishnavh@google.com)

# Register tokens improve arithmetic performance.



# Register tokens improve arithmetic performance.



# Register tokens improve arithmetic performance.

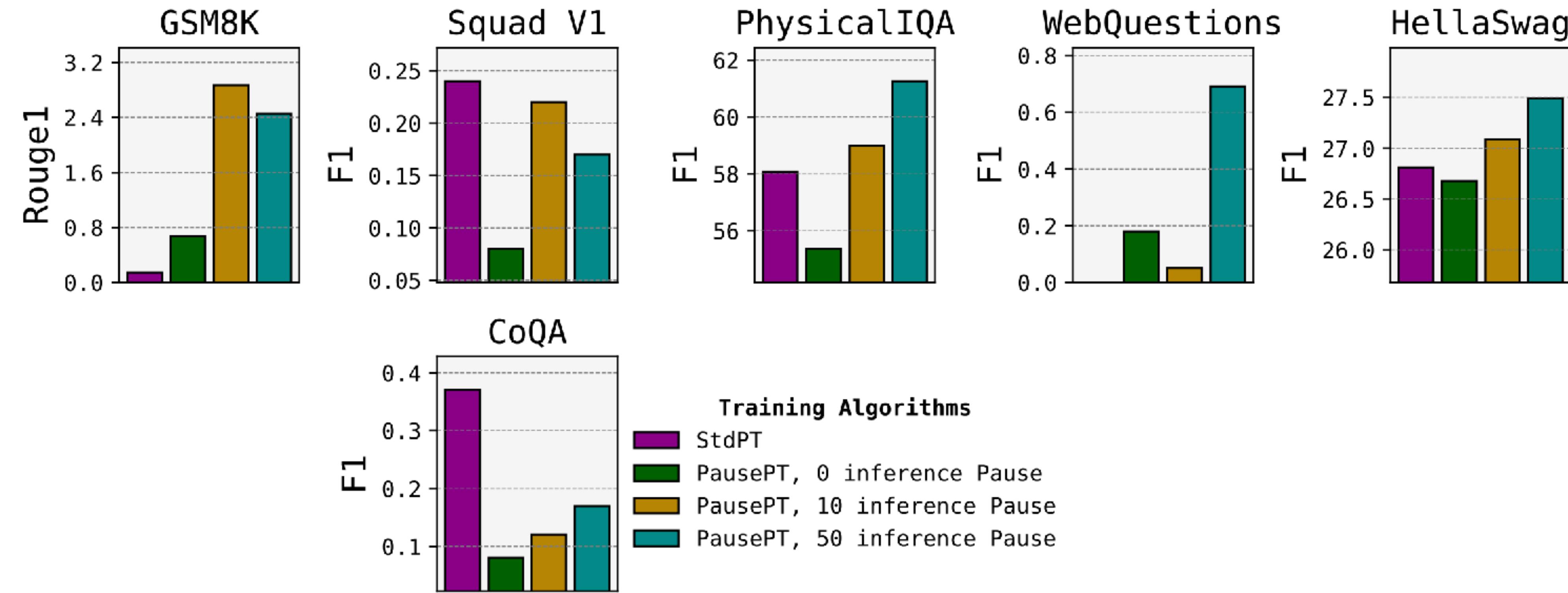


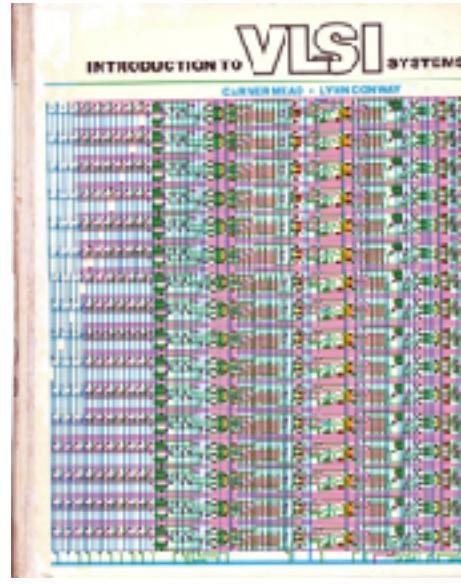
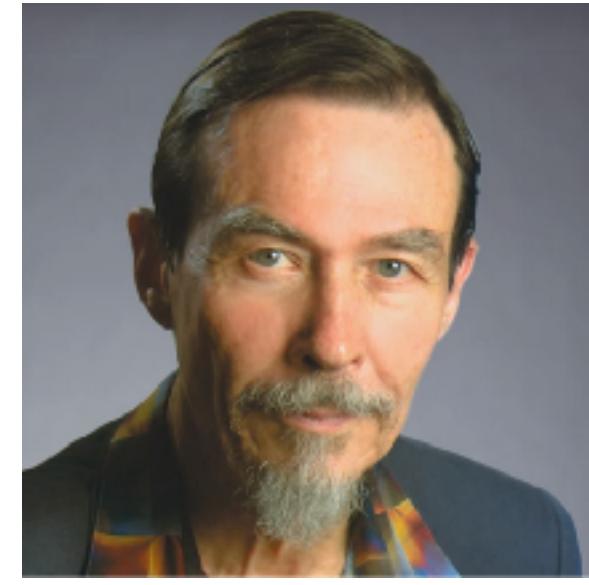
Figure 6: **Zero-shot evaluation of pause-pretrained models.** Zero-shot inference with pause tokens on a pause-pretrained model gives gains on tasks like GSM8k and HellaSwag. However, we note that our zero-shot accuracies are quite low, as we experiment with a small 1B parameter model.

# Register tokens improve arithmetic performance.

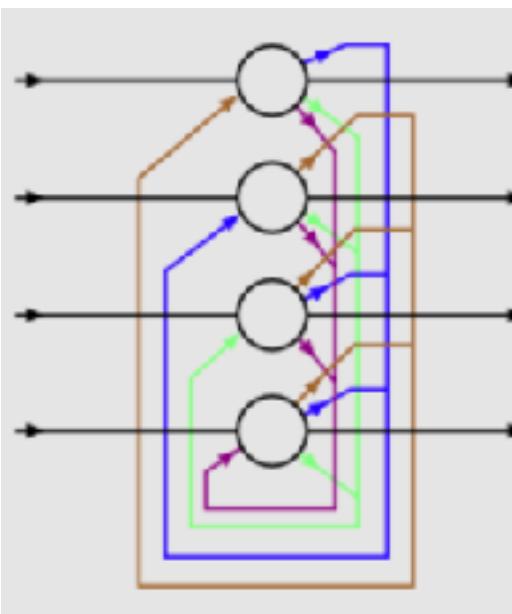
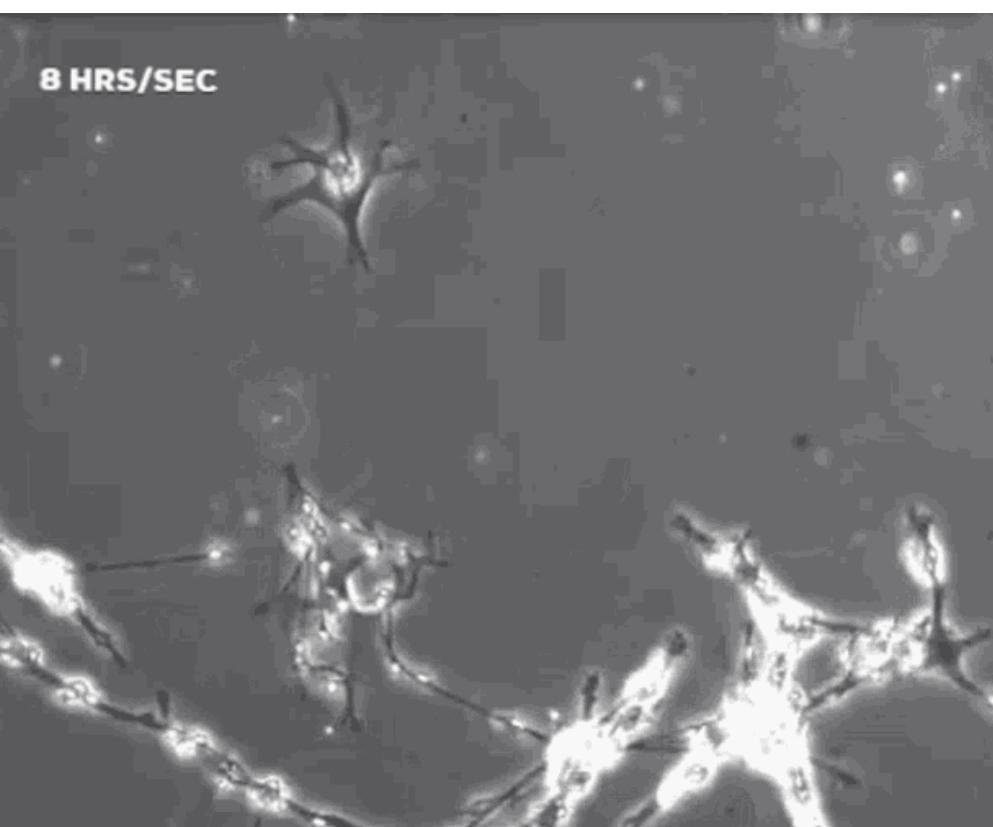
- They don't know about **world models!**
- Their experiments are **limited**.
  - Pre-training objective/dataset/etc.
  - Register token clumping.
- Can I add **way more registers -> way better performance?**

# CNS Heroes

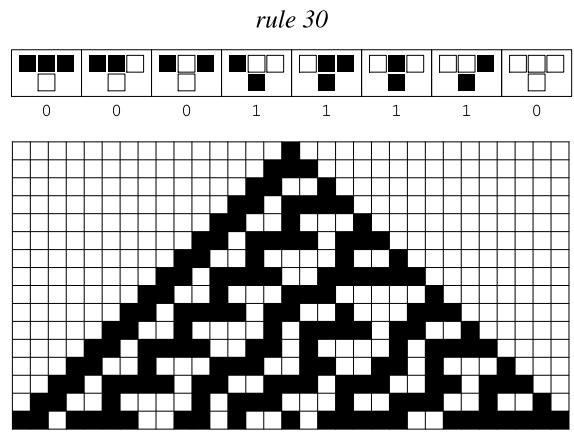
## Toward Cellular Computing



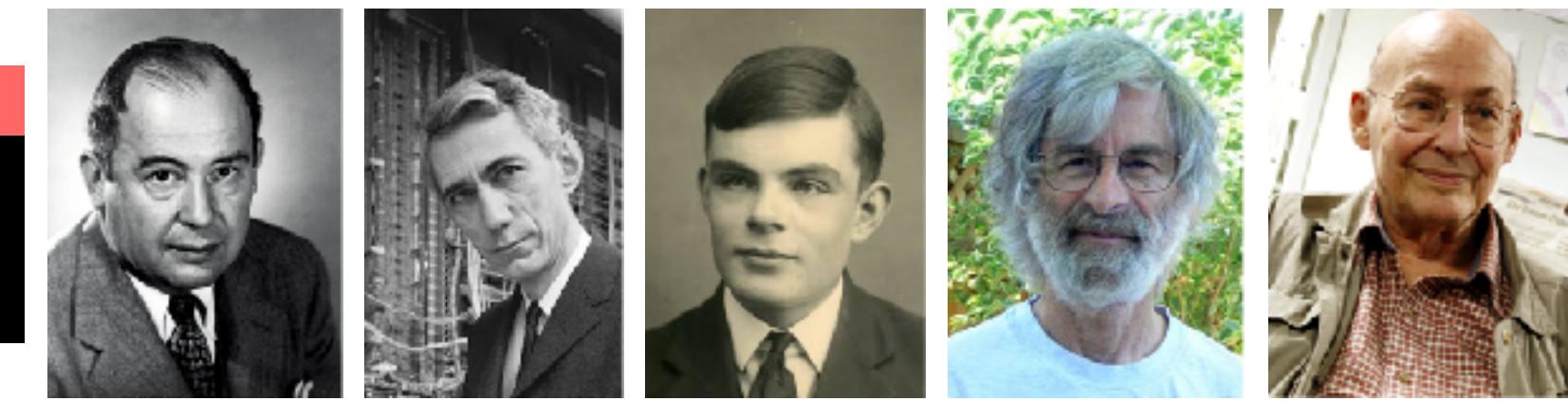
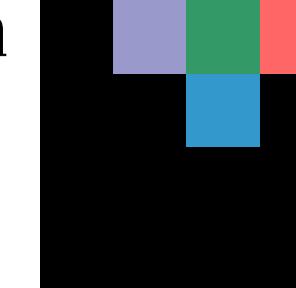
Mead/Conway Neuromorphics



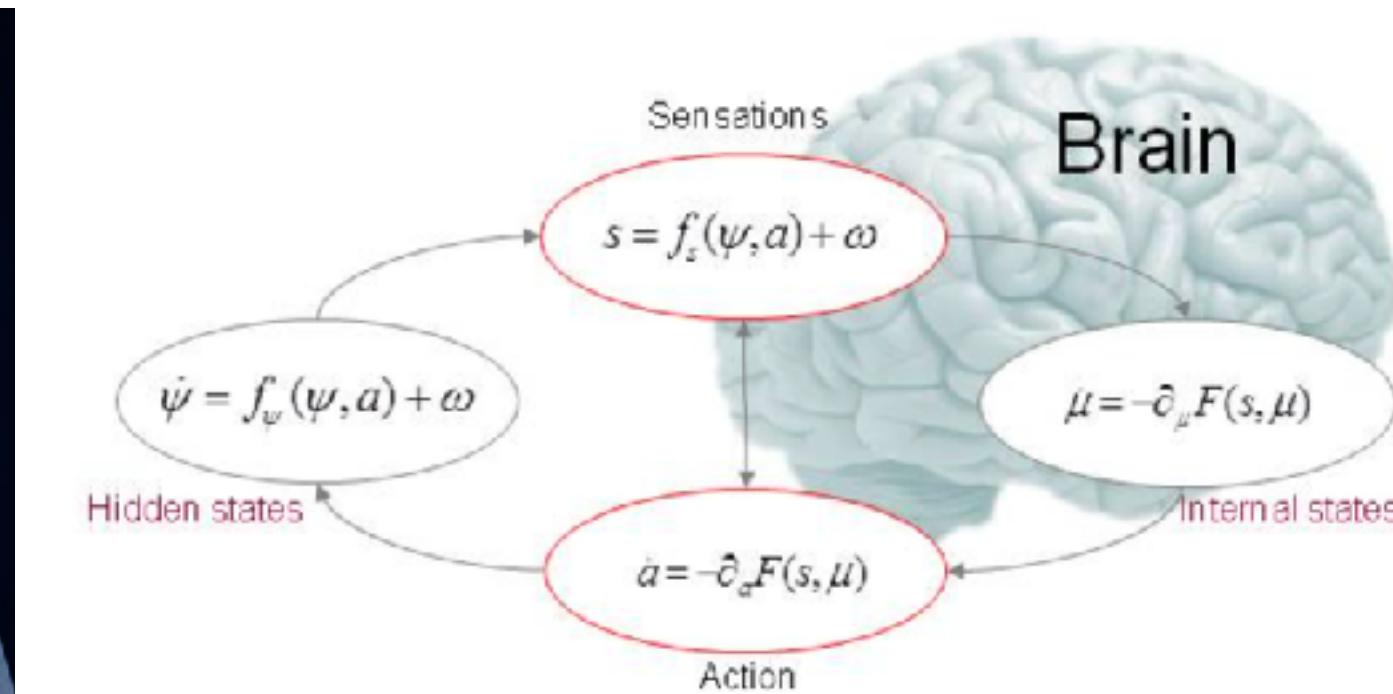
Hopfield/Comp Neuro



Thinking Machines Corporation



Feynman/Wolfram/Winfree/Levin/Thomson  
Cellular/bio/chemical computation



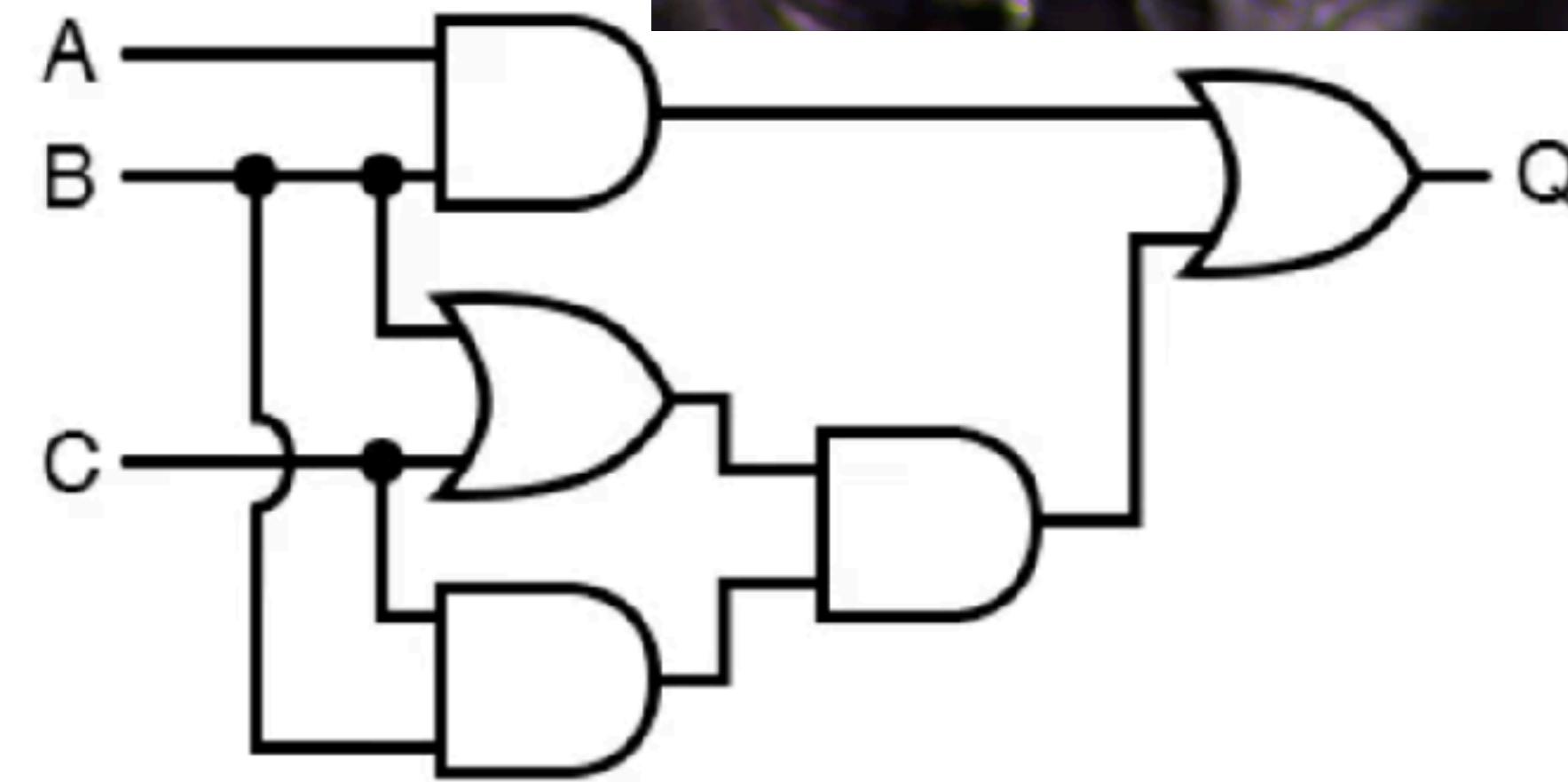
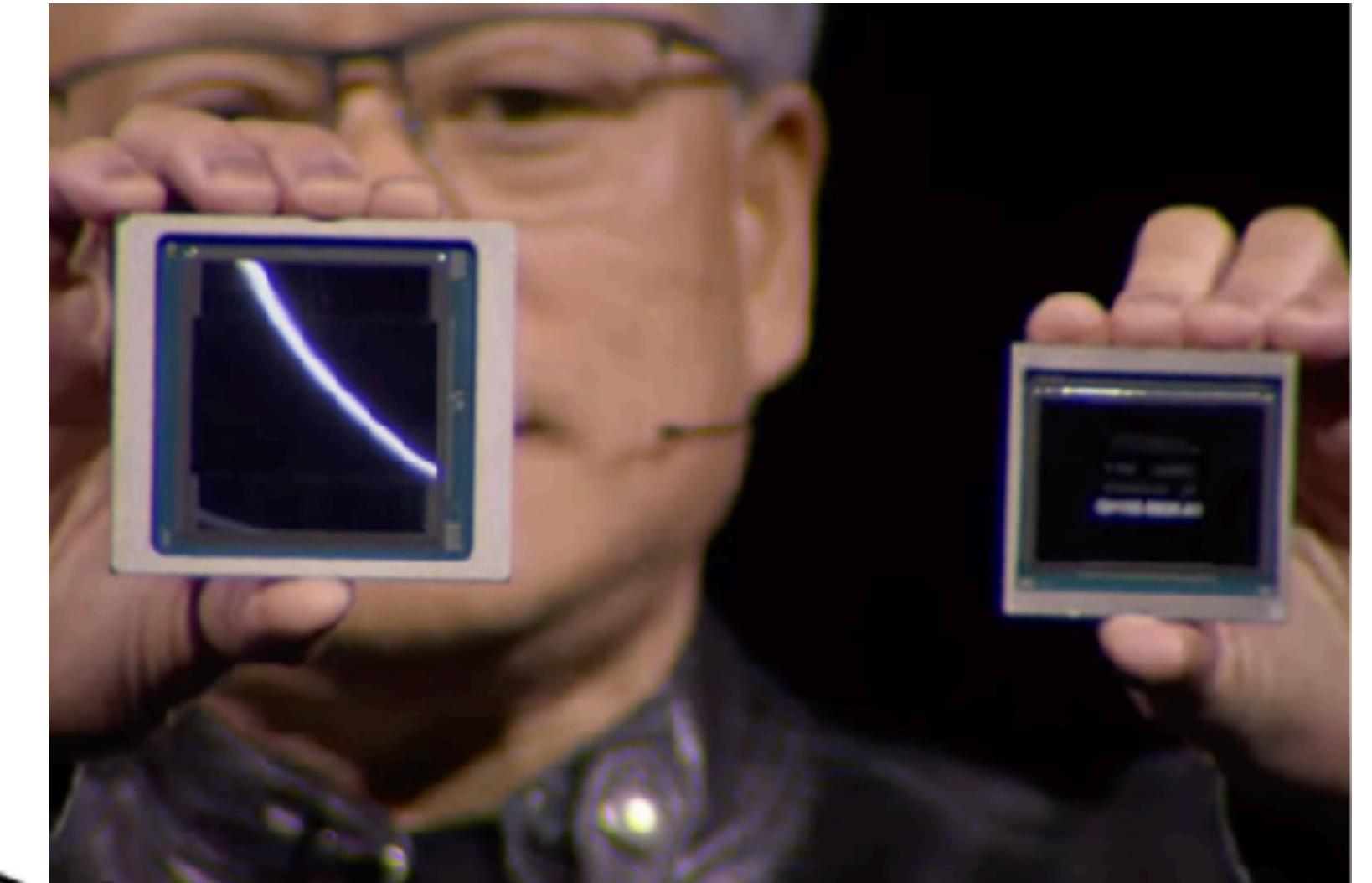
Friston/Free Energy Principle

# Beyond LLMs

## Toward Cellular Computing



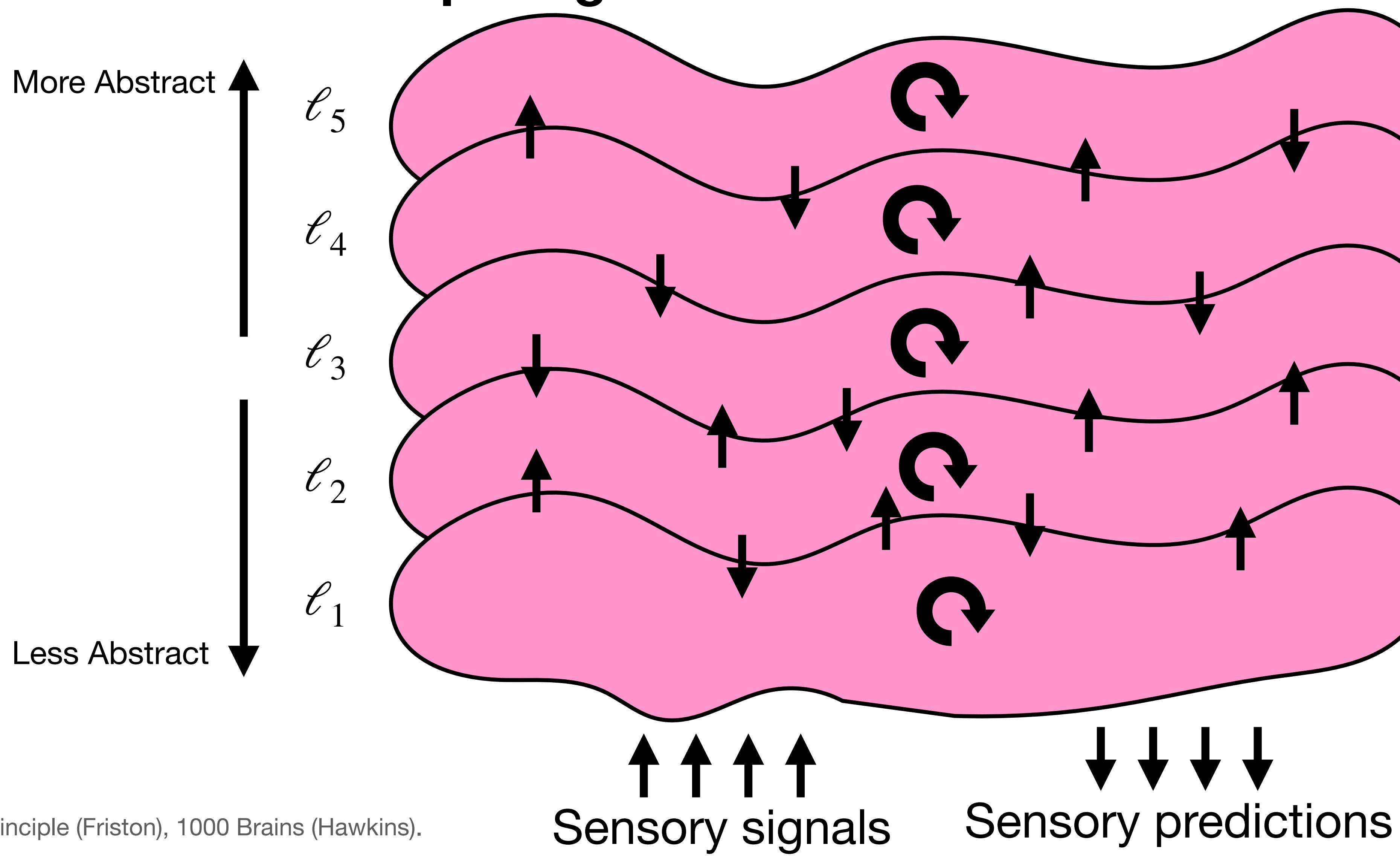
Sand



Thinking Sand

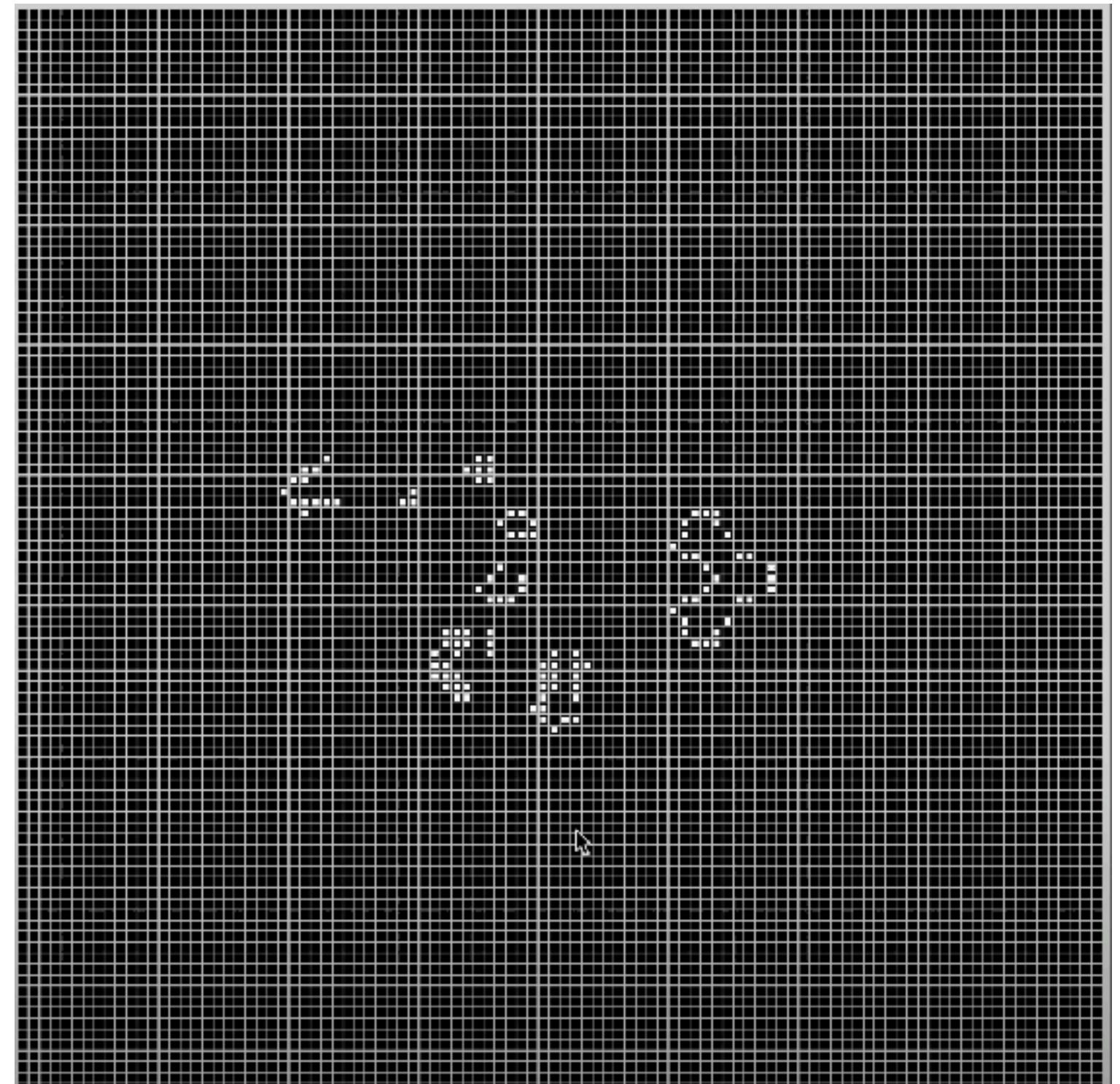
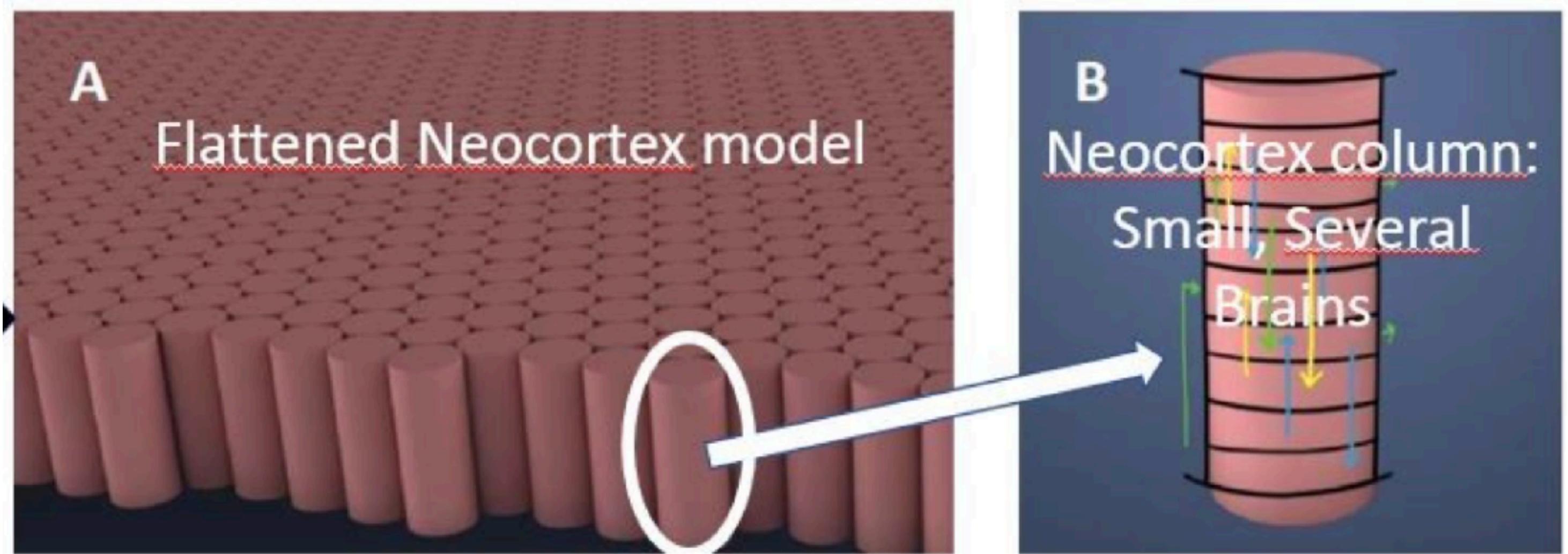
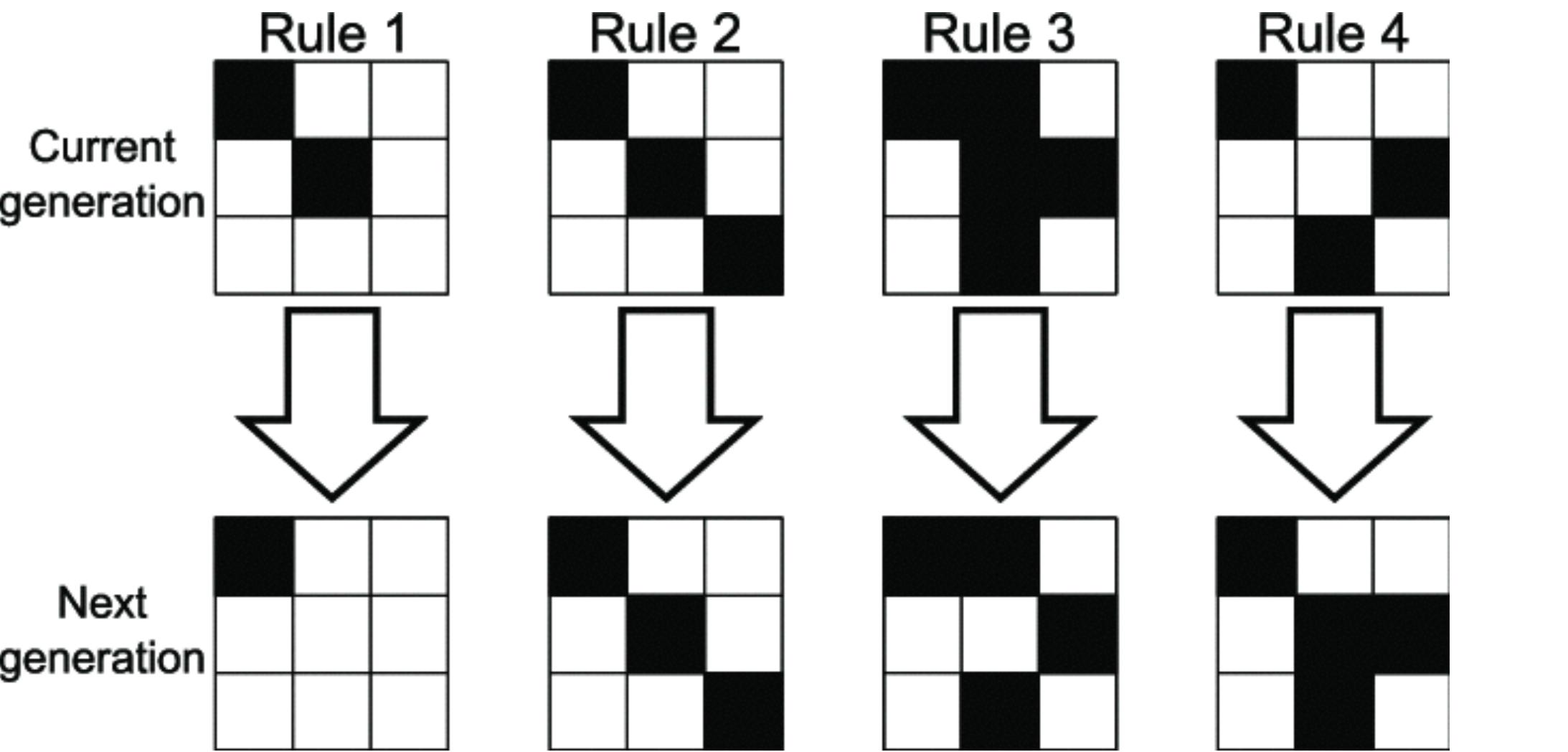
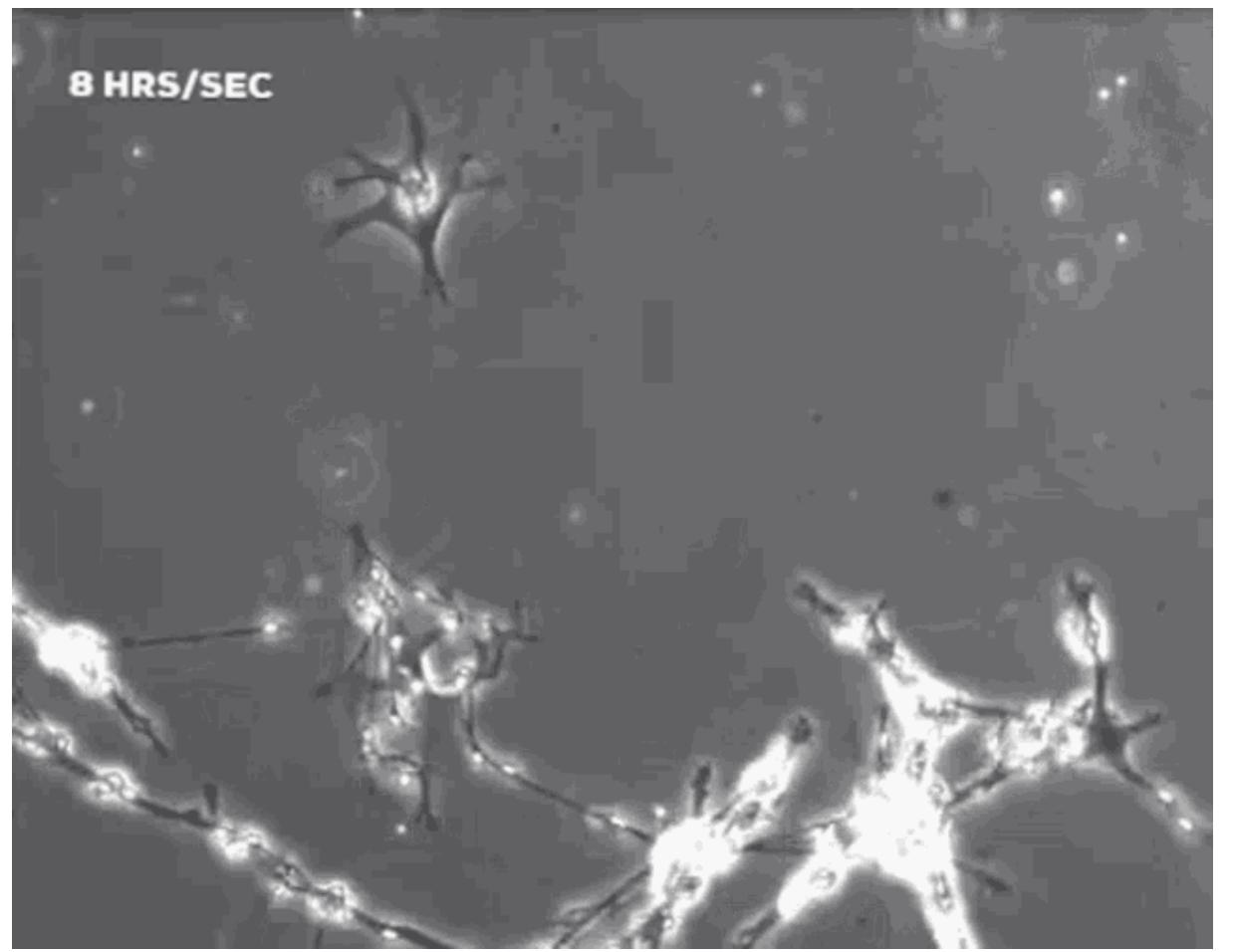
# Beyond LLMs

## Toward Cellular Computing



# Beyond LLMs

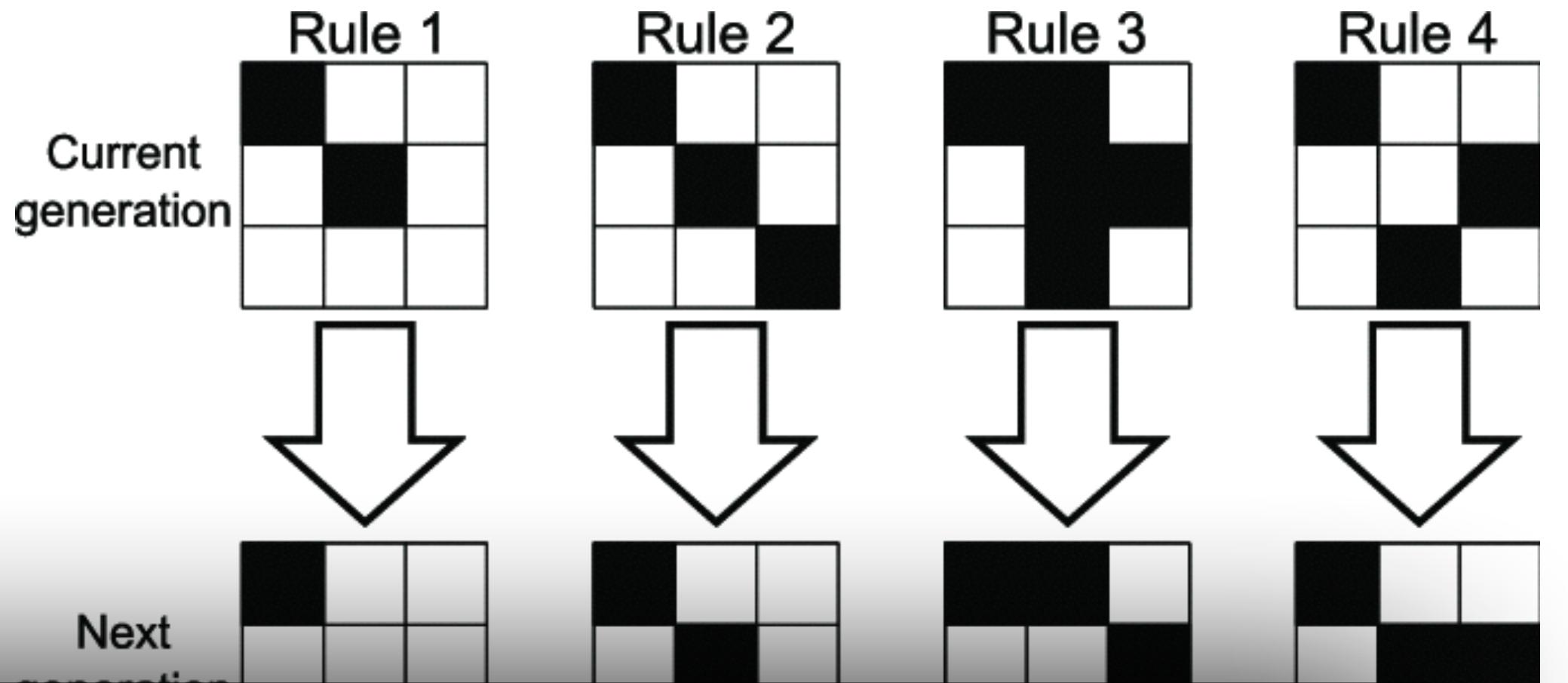
## Cortical Universal Update Rule



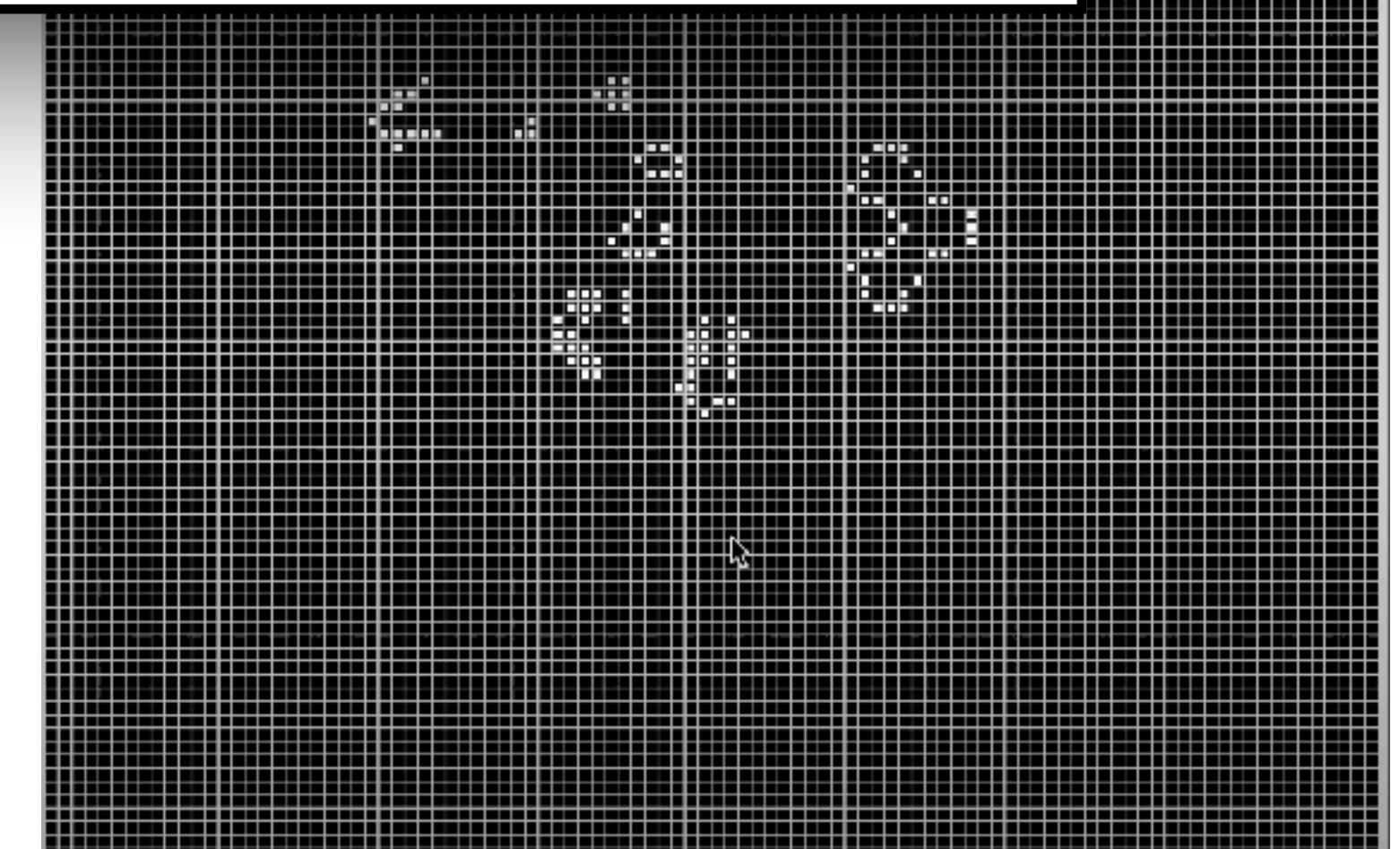
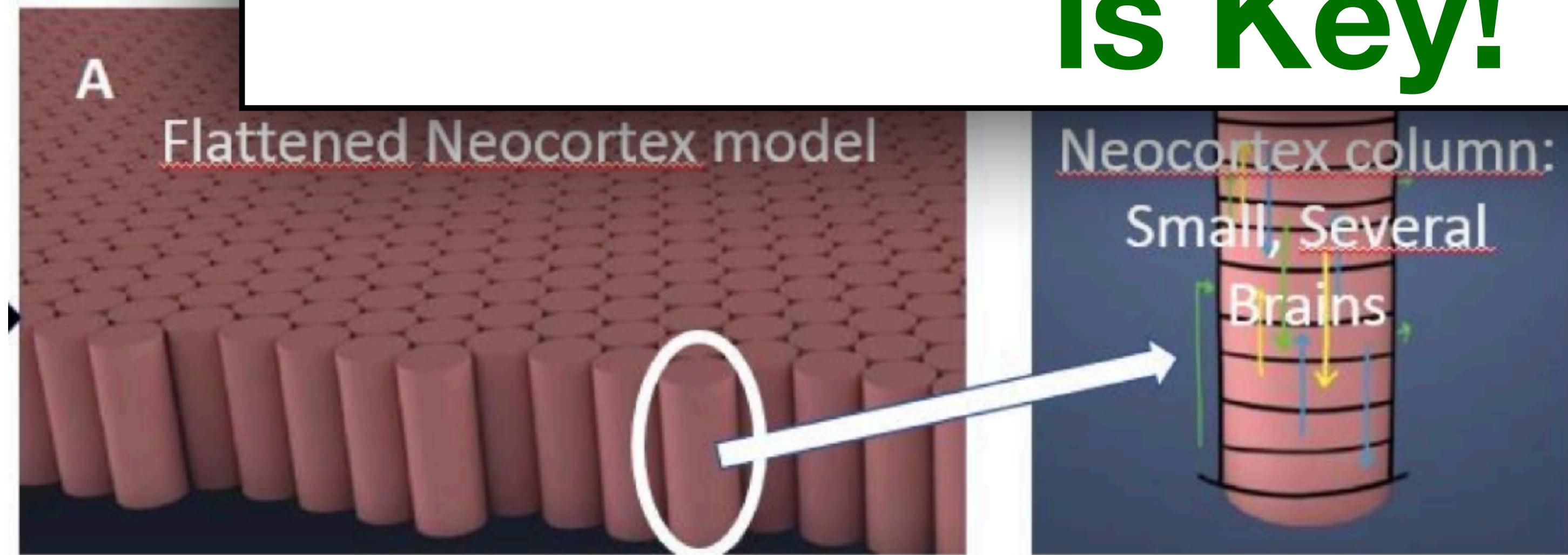
FEP (Friston), 1000 Brains (Hawkins), New Kind of Science (Wolfram), etc.

# Beyond LLMs

## Cortical Universal Update Rule



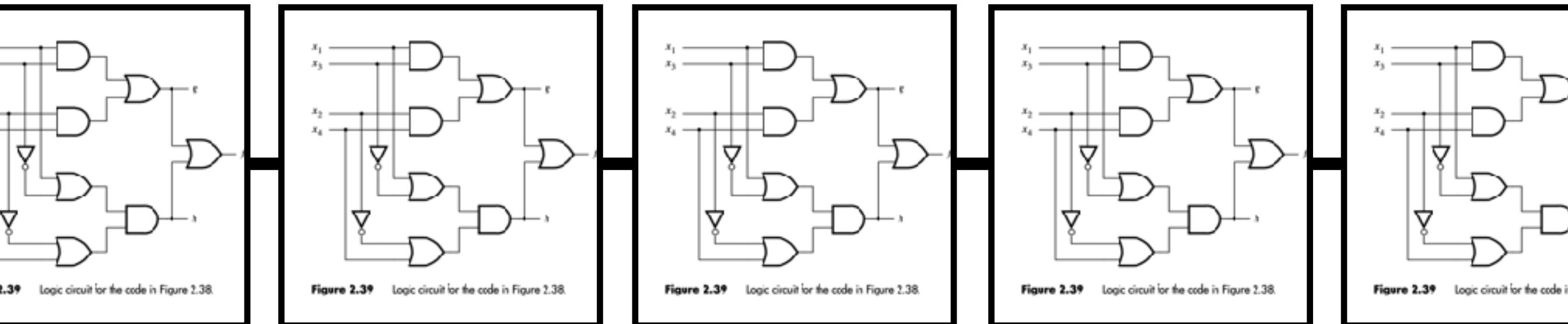
# Local Information Processing is Key!



# Beyond LLMs

## Cortical Universal Update Rule

*Third Edition*

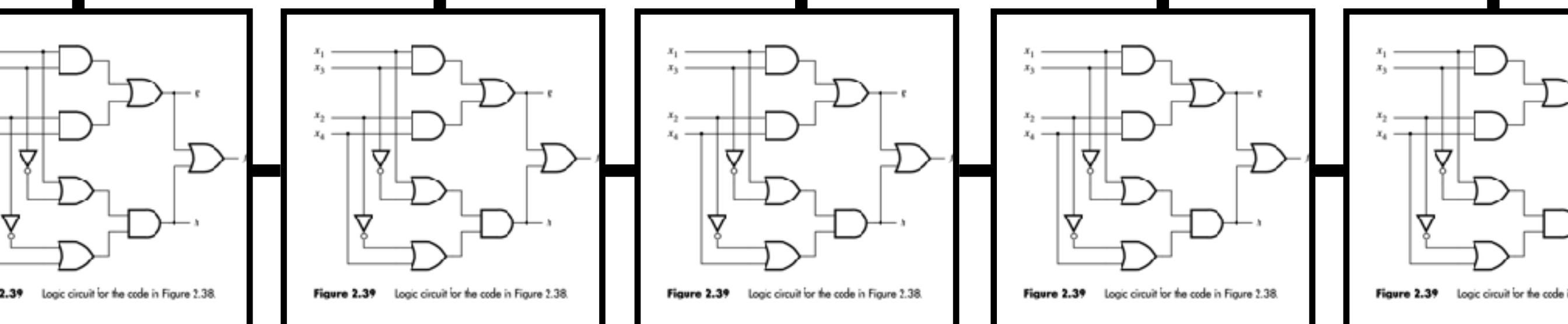


**Figure 2.39** Logic circuit for the code in Figure 2.38.

**Figure 2.39** Logic circuit for the code in Figure 2.38.

**Figure 2.39** Logic circuit for the code in Figure 2.38.

**Figure 2.39** Logic circuit for the code in Fig.

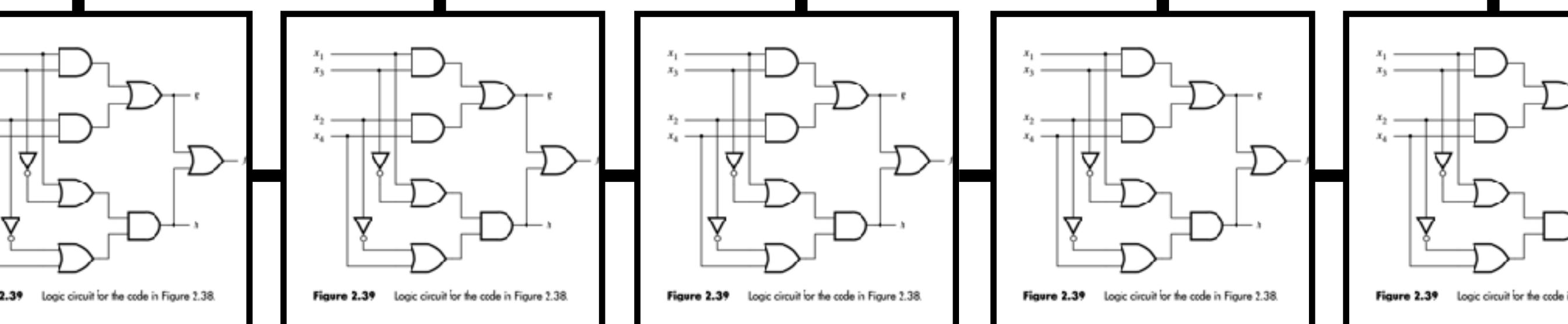


**Figure 2.39** Logic circuit for the code in Figure 2.3.

**Figure 2.39** Logic circuit for the code in Figure

**Figure 2.39** Logic circuit for the code in Fig.

**Figure 2.39** Logic circuit for the code in Fig.

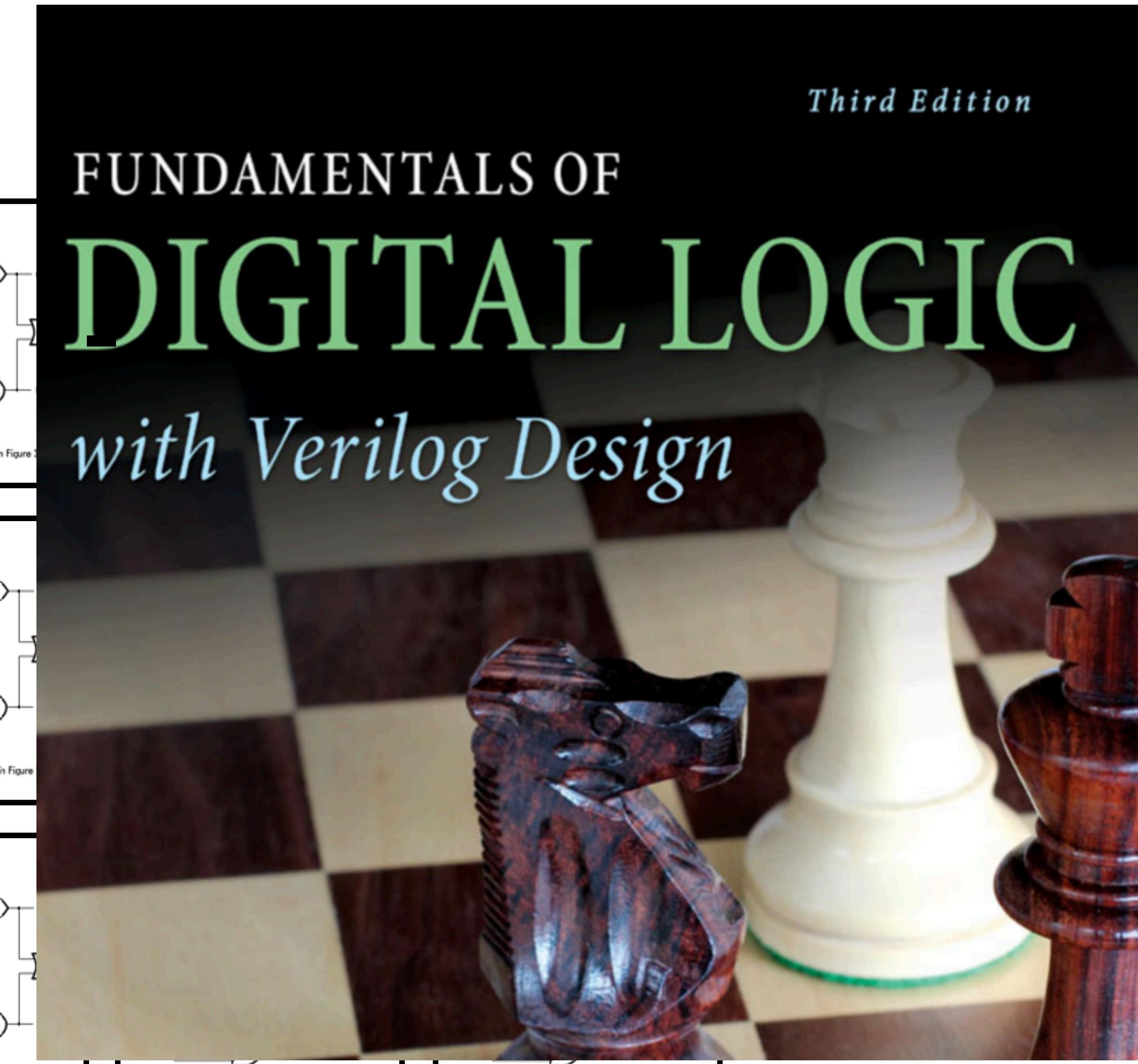


**Figure 2-39** Logic circuit for the code in Figure 2-38.

**Figure 2-39** Logic circuit for the code in Figure 1.

**Figure 2-39** Logic circuit for the code in Fig.

**Figure 2-39** Logic circuit for the code in Table 2-1.



FEP (Friston). 1000 Brains (Hawkins). New Kind of Science (Wolfram). etc.

# Toward Cellular Computing

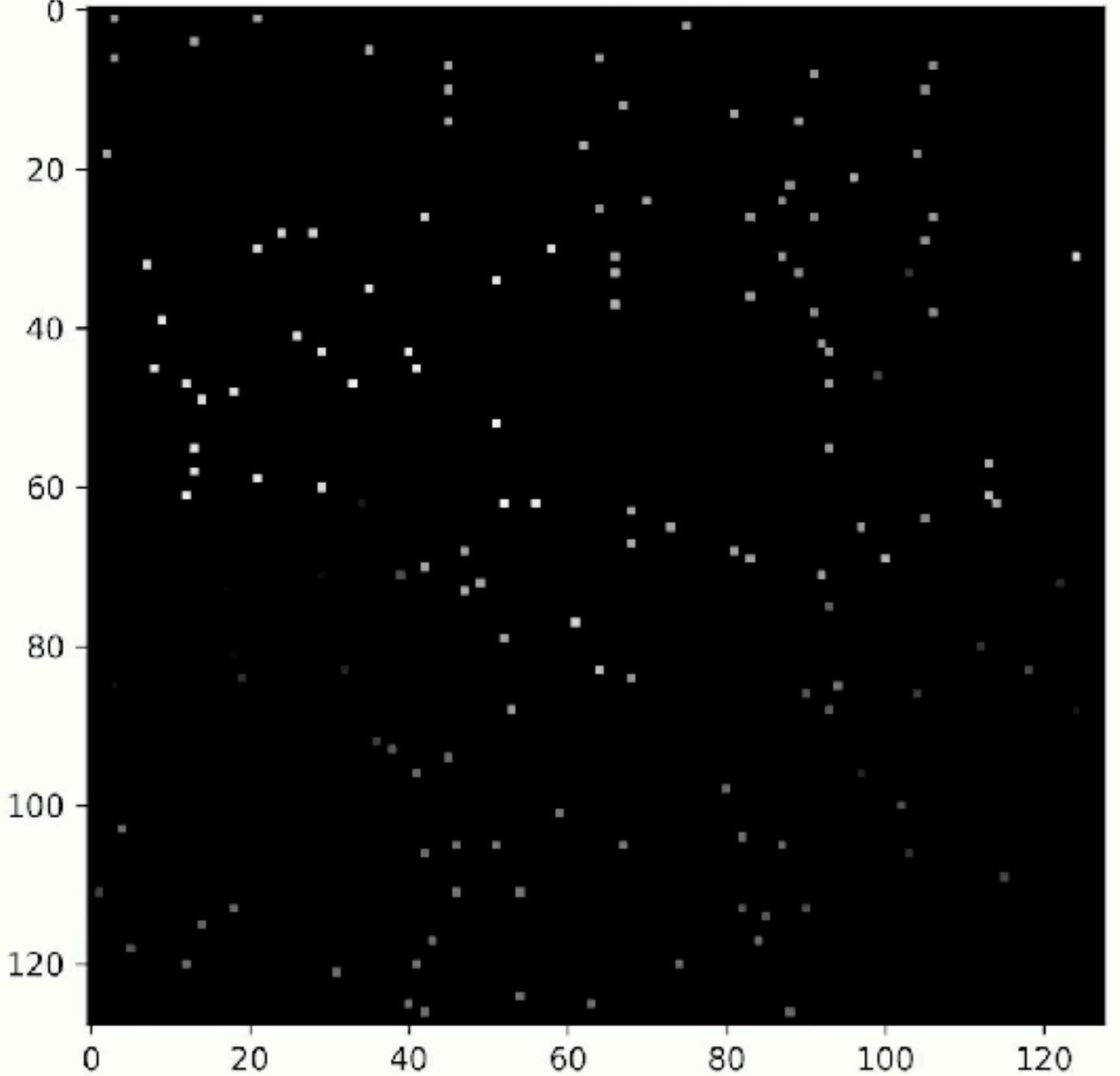
## Limit case of physically-realizable computation

Neural Cellular Automata do Active Inference

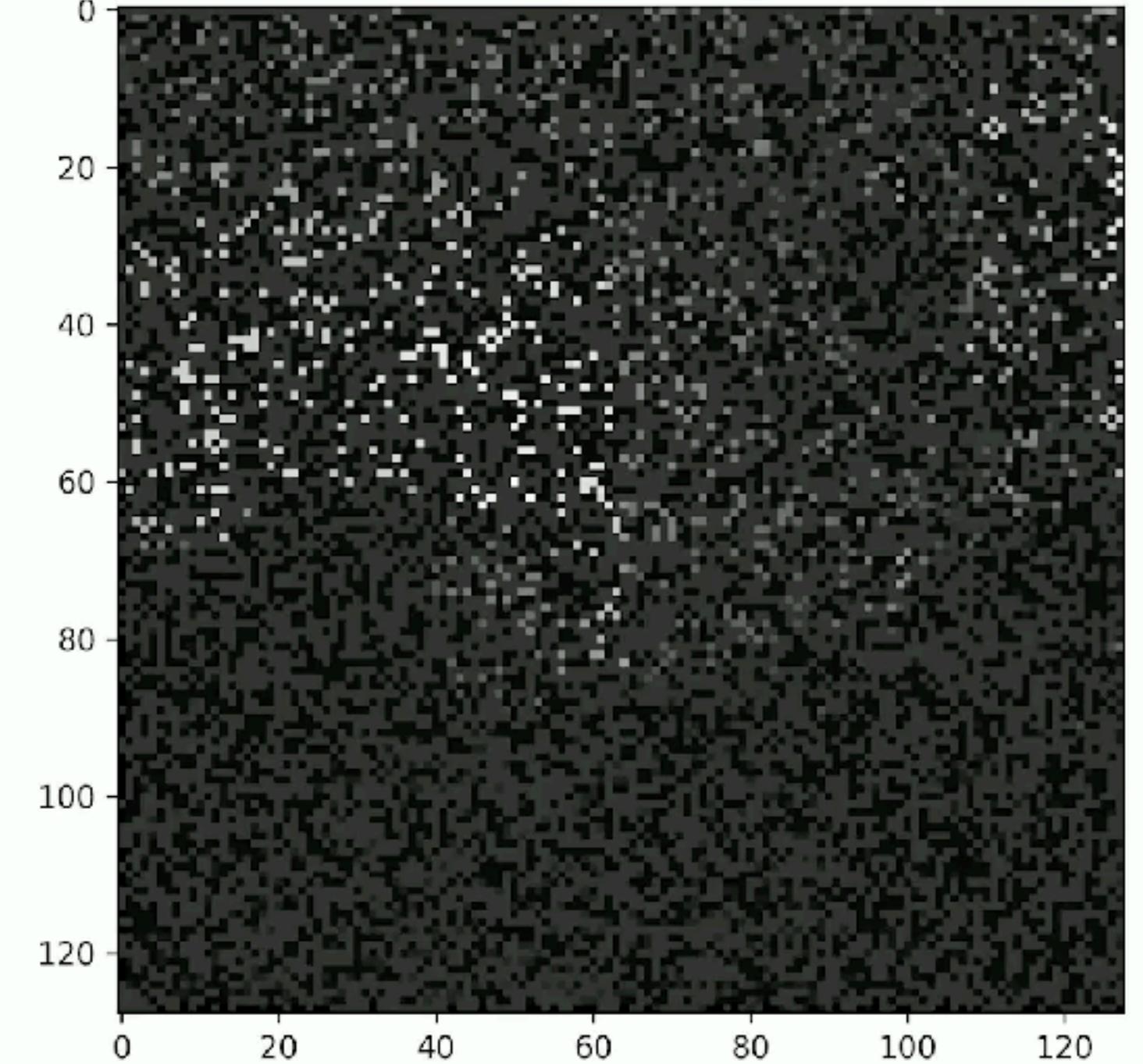
Source Video



Raw Info Stream -- 04\_0.5\_update frame 0

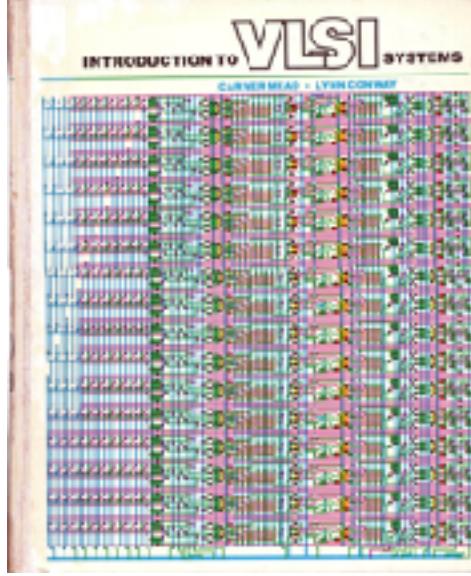
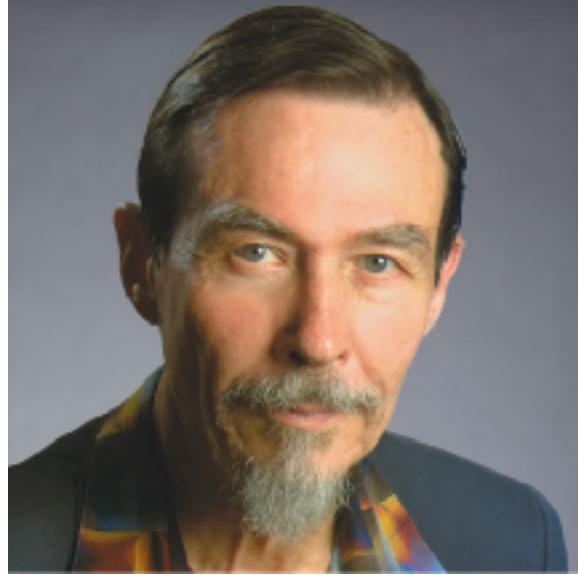


Trained Model State Estimate -- 02.5\_diffusion\_test frame 0



# Toward Cellular Computing

## Limit case of physically-realizable computation



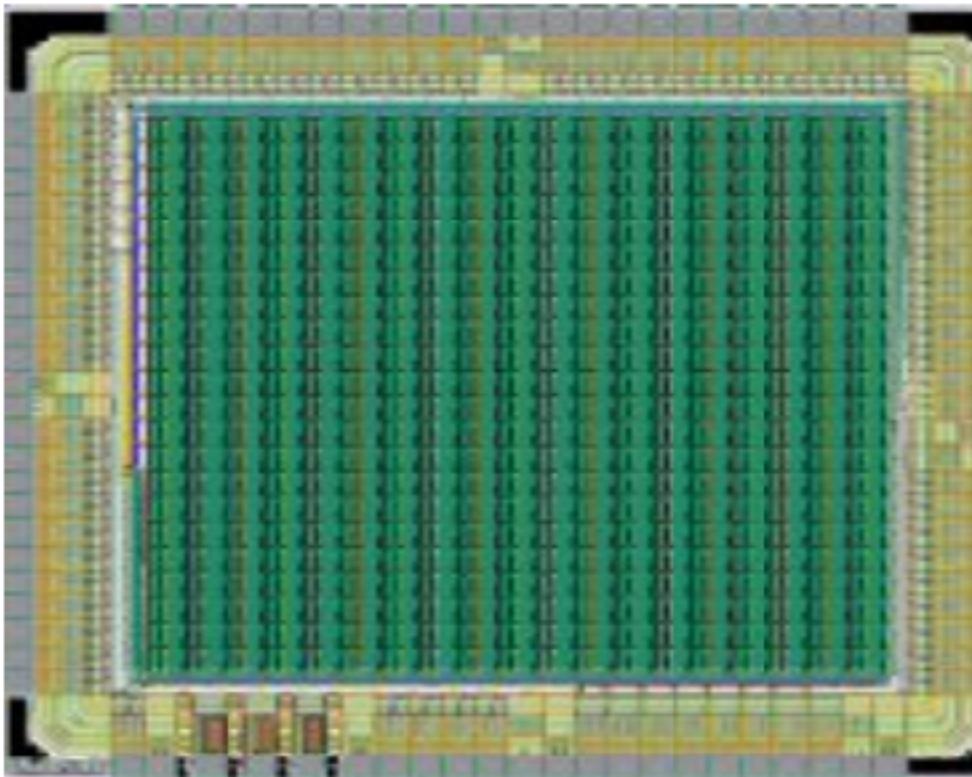
Piotr Dudek

The University of Manchester

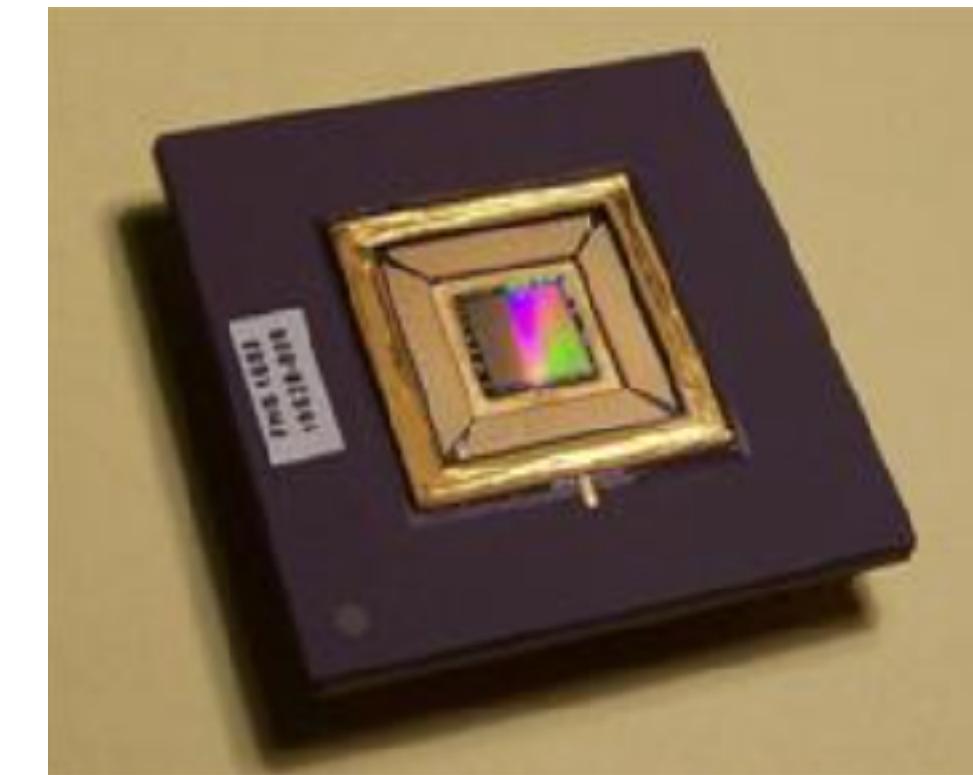
Verified email at manchester.ac.uk - Homepage

VLSI Design Neuromorphic Engineering Imag

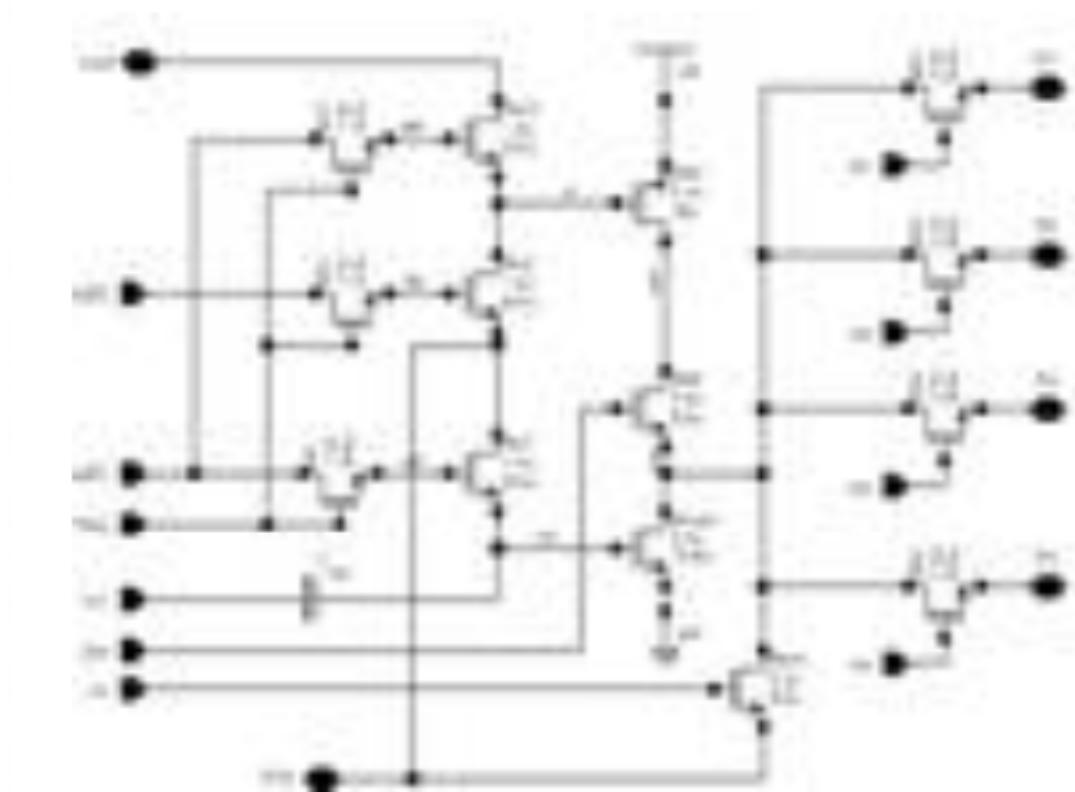
Mead/Conway Neuromorphics



Cellular processors  
arrays



Vision Chips



Analog Sig Proc

# Thank You!

## Questions?

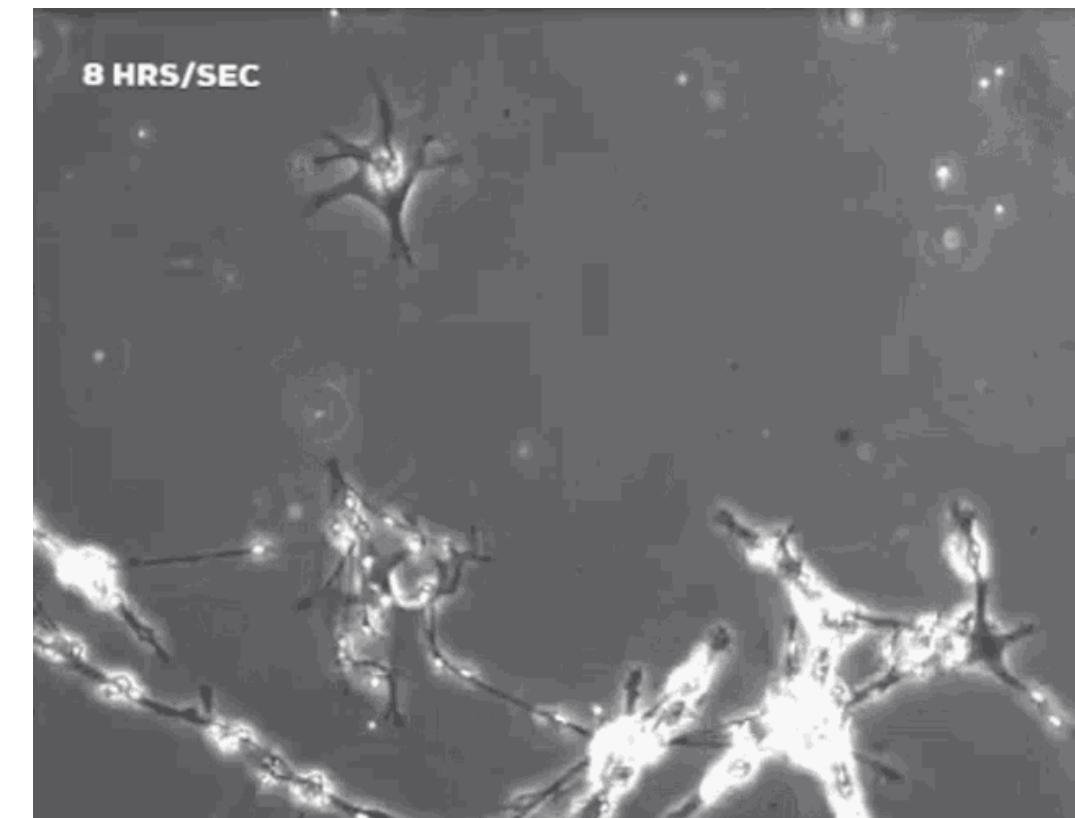
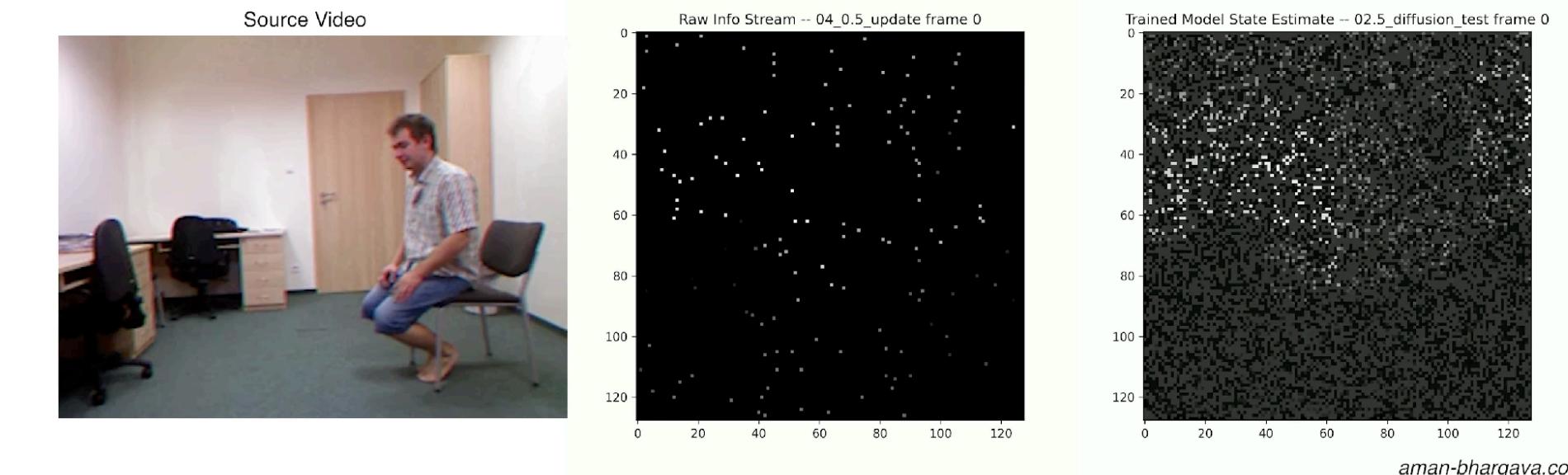
### Mentors

- Matt Thomson
- Erik Winfree
- Ralph Adolphs
- Frederick Eberhart
- Rob Phillips
- Steve Mann (UToronto)
- Milad Lankarany (UToronto)
- Michael Levin (Tufts)

### Friends + Collaborators

- Pantelis Vafidis
- Cameron Witkowski (UToronto)
- Salvador Buse
- Cayden Pierce (MIT)
- James Gornet
- Meera Prasad
- Michael Zellinger
- Hersh Bhargava (UCSF)
- Mango Weng
- Mingshi Chi (UToronto/York)

Neural Cellular Automata do Active Inference



**[EOS]**

# Transformer theory supports register tokens.

## Circuit Complexity of Chain-of-Thought

step towards theoretically answering these questions. Specifically, we examine the expressivity of LLMs with CoT in solving fundamental mathematical and decision-making problems. By using circuit complexity theory, we first give impossibility results showing that bounded-depth Transformers are unable to directly produce correct answers for basic arithmetic/equation tasks unless the model size grows *super-polynomially* with respect to the input length. In contrast, we then prove by construction that autoregressive Transformers of *constant size* suffice to solve both tasks by generating CoT derivations using a commonly used math language format. Moreover, we show LLMs with CoT can handle a general class of decision-making problems known as Dynamic Programming, thus justifying their power in tackling

# Register tokens help vision transformers.

## Vision Transformers Need Registers

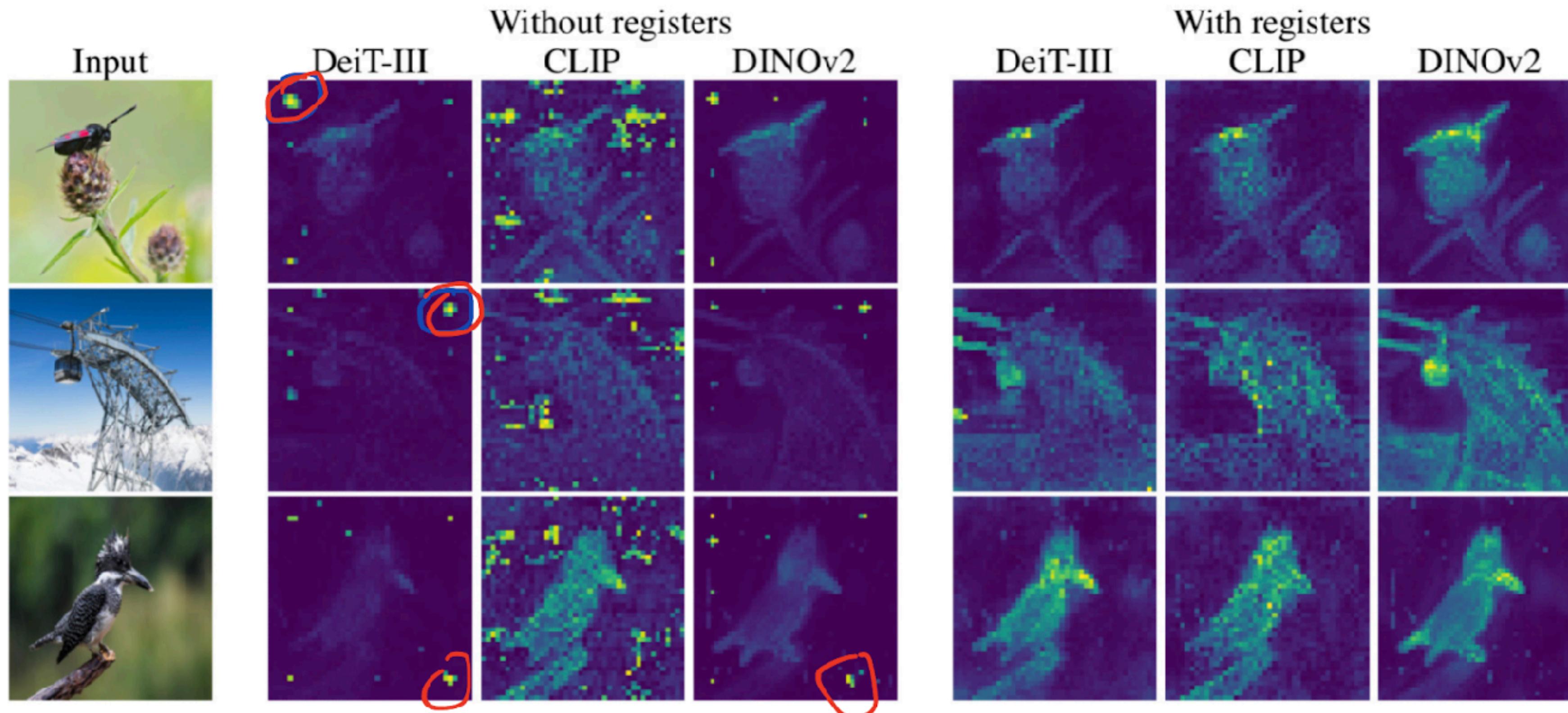
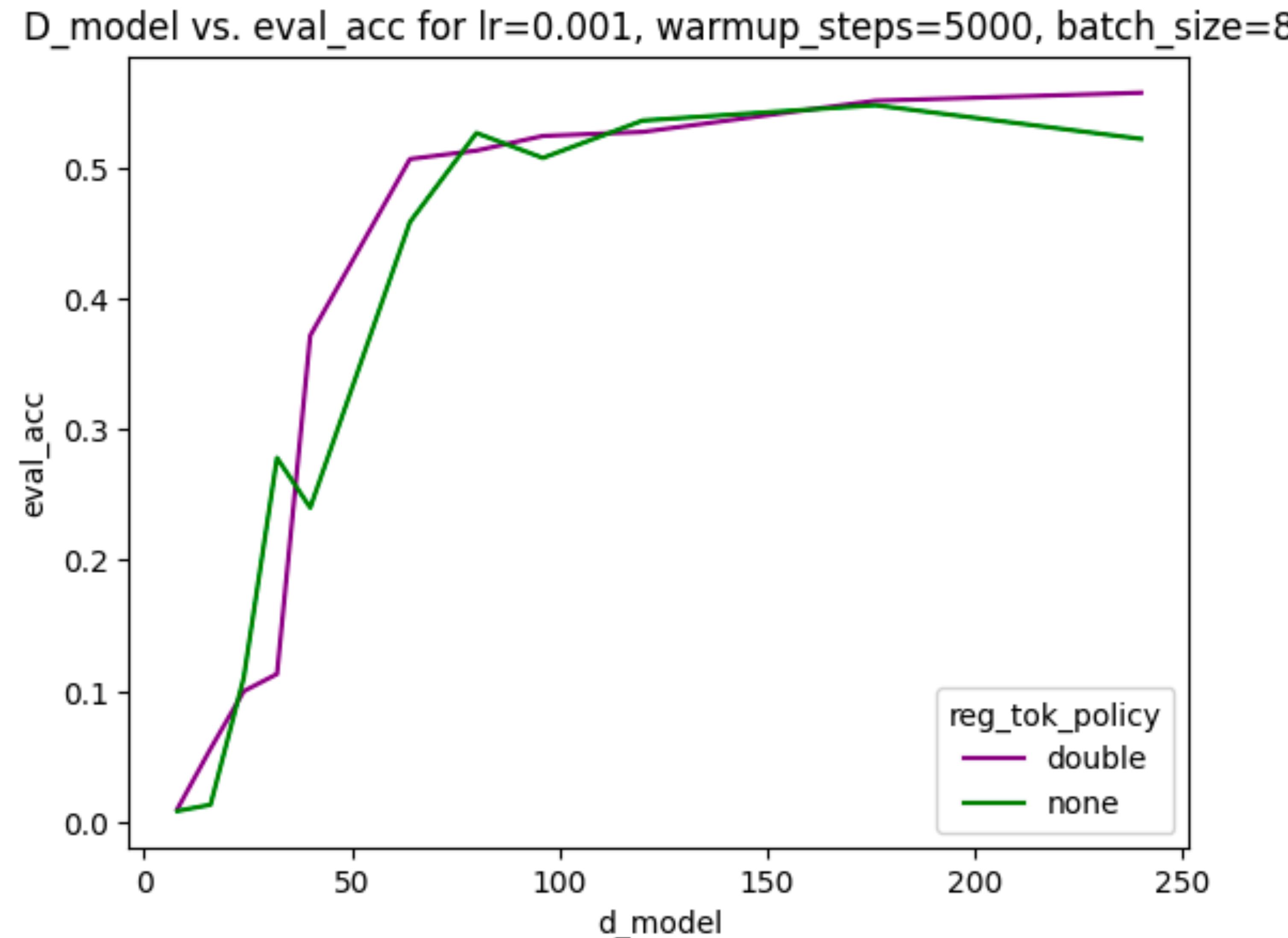


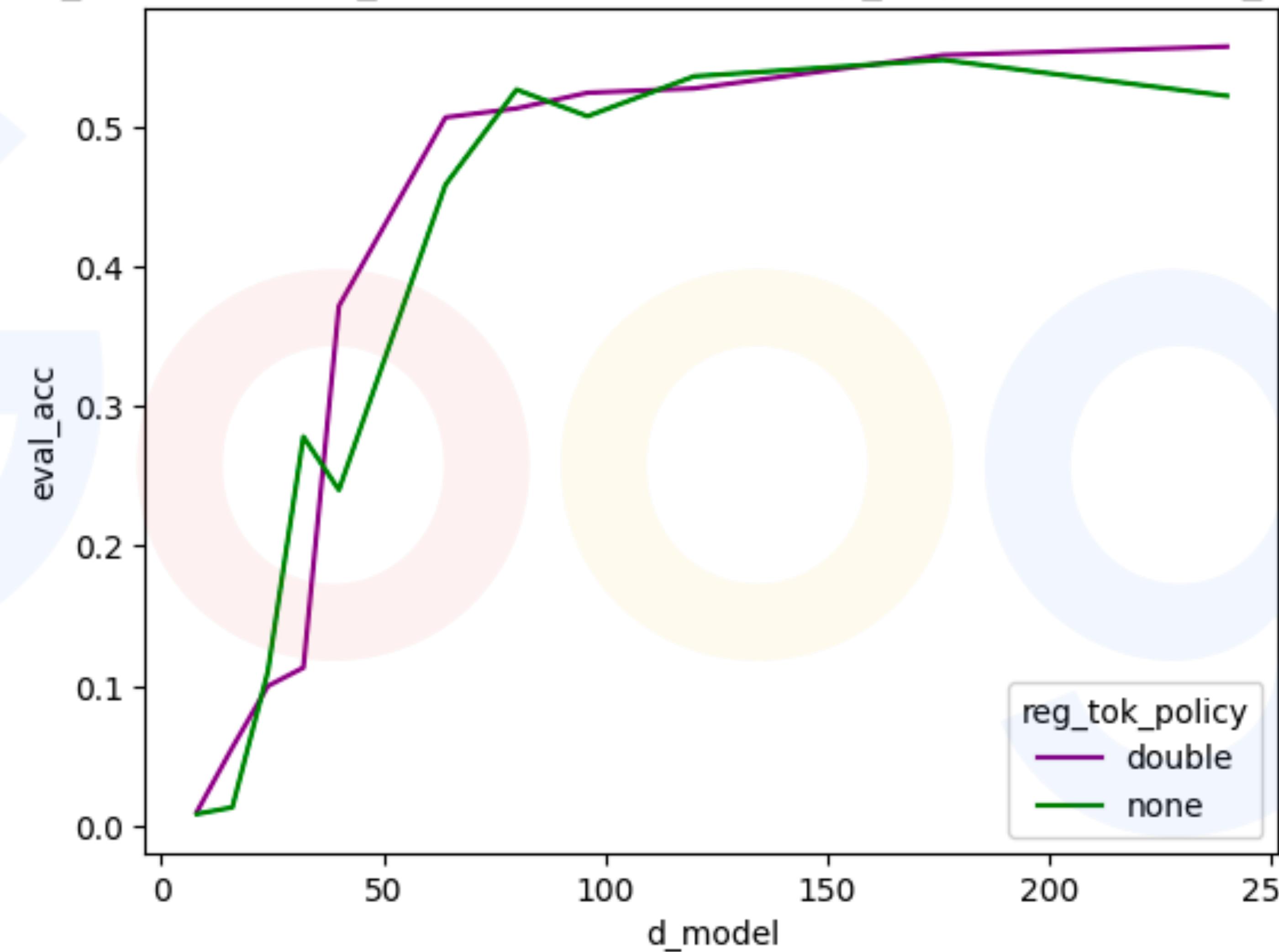
Figure 1: Register tokens enable interpretable attention maps in all vision transformers, similar to the original DINO method (Caron et al., 2021). Attention maps are calculated in high resolution for better visualisation. More qualitative results are available in appendix D.

# Register tokens improve arithmetic performance.

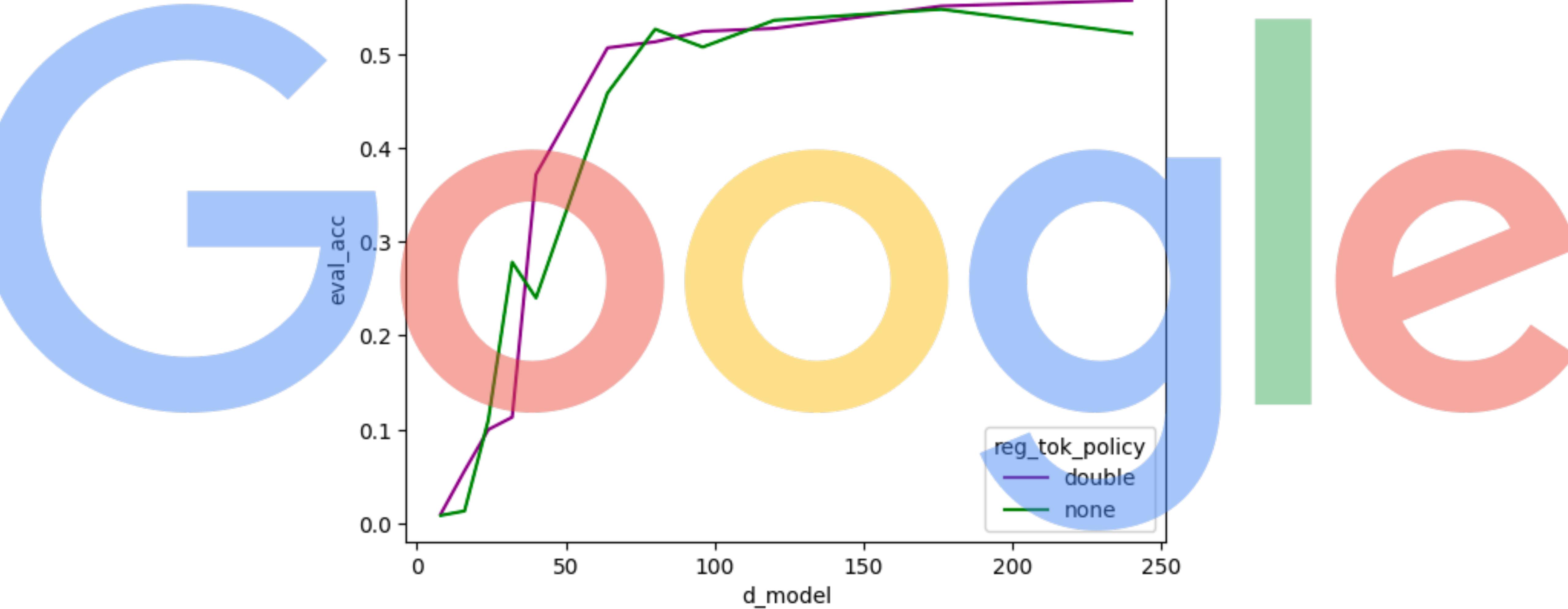


# Register tokens improve arithmetic performance.

D\_model vs. eval\_acc for lr=0.001, warmup\_steps=5000, batch\_size=8



# Register tokens improve arithmetic performance.



# Register Tokens enable LLM world model exploration.

## Big-Picture

-  Should be able to “**stop and think**”, leveraging the language “world model”.
-  Thinking should be optimized **during pre-training**.
-  **Inference-time scalability** would be nice.

# Thought as internal world model exploration

## Cortical processing cartoon sketch



*AppliedMath*



---

*Article*

## Gradient-Free Neural Network Training via Synaptic-Level Reinforcement Learning

Aman Bhargava <sup>1,2</sup>, Mohammad R. Rezaei <sup>2,3,4</sup> and Milad Lankarany <sup>2,3,4,5,\*</sup>

<sup>1</sup> Division of Engineering Science, University of Toronto, Toronto, ON M5S 2E4, Canada;  
aman.bhargava@mail.utoronto.ca

<sup>2</sup> Division of Clinical and Computational Neuroscience, Krembil Brain Institute, University Health Network,  
Toronto, ON M5G 1L7, Canada; mr.rezaei@mail.utoronto.ca

<sup>3</sup> Institute of Biomedical Engineering, University of Toronto, Toronto, ON M5S 2E4, Canada

<sup>4</sup> KITE, Toronto Rehabilitation Institute, University Health Network, Toronto, ON M5G 1L7, Canada

<sup>5</sup> Department of Physiology, University of Toronto, Toronto, ON M5S 2E4, Canada

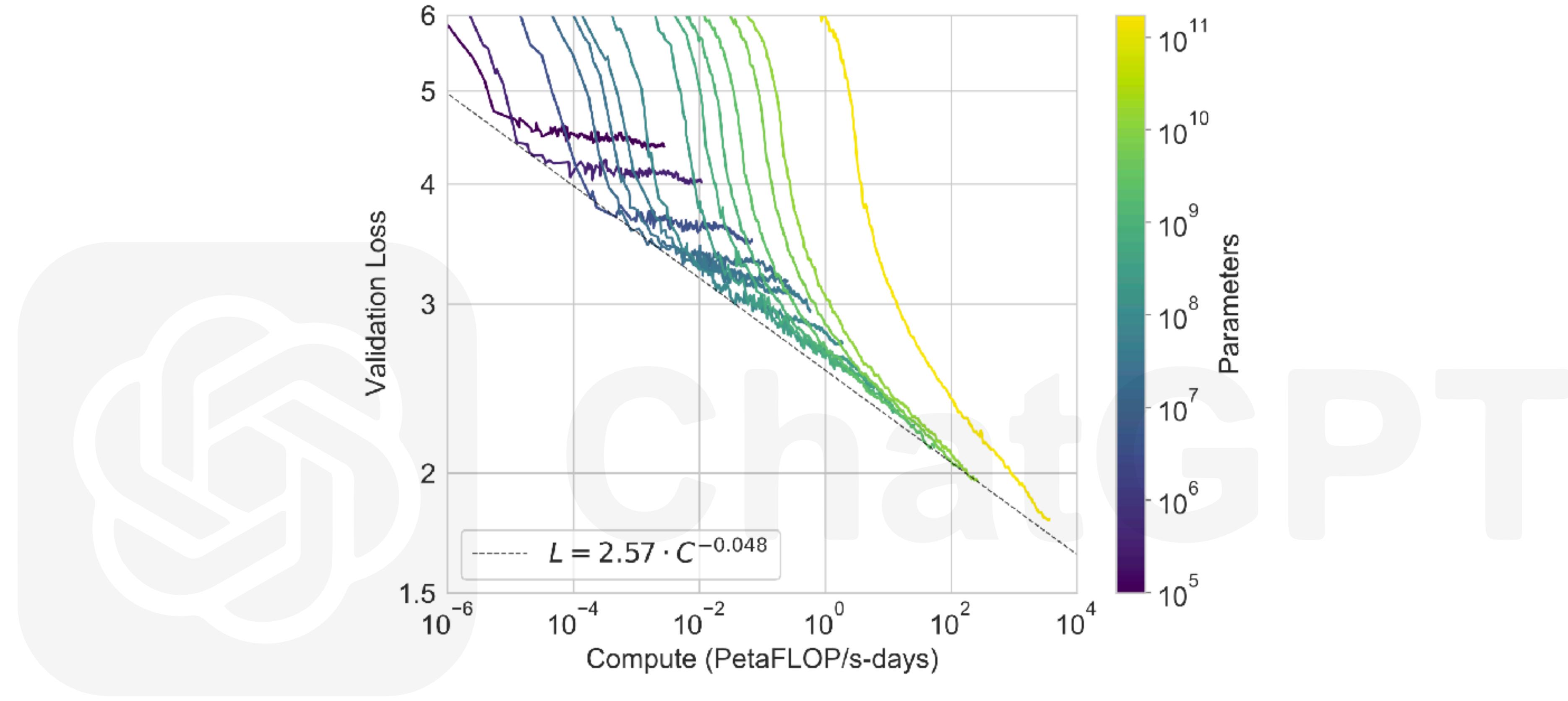
\* Correspondence: milad.lankarany@uhnresearch.ca

# World Models emerge from Predictive Coding.

Ilya Sutskever (Chief Scientist & Co-Founder, OpenAI)



# Bigger LLM → Better Prediction → Greater Capability



**Figure 3.1: Smooth scaling of performance with compute.** Performance (measured in terms of cross-entropy validation loss) follows a power-law trend with the amount of compute used for training. The power-law behavior observed in [KMH<sup>+</sup>20] continues for an additional two orders of magnitude with only small deviations from the predicted curve. For this figure, we exclude embedding parameters from compute and parameter counts.

# March 2024 LLMs Suck at Reasoning.

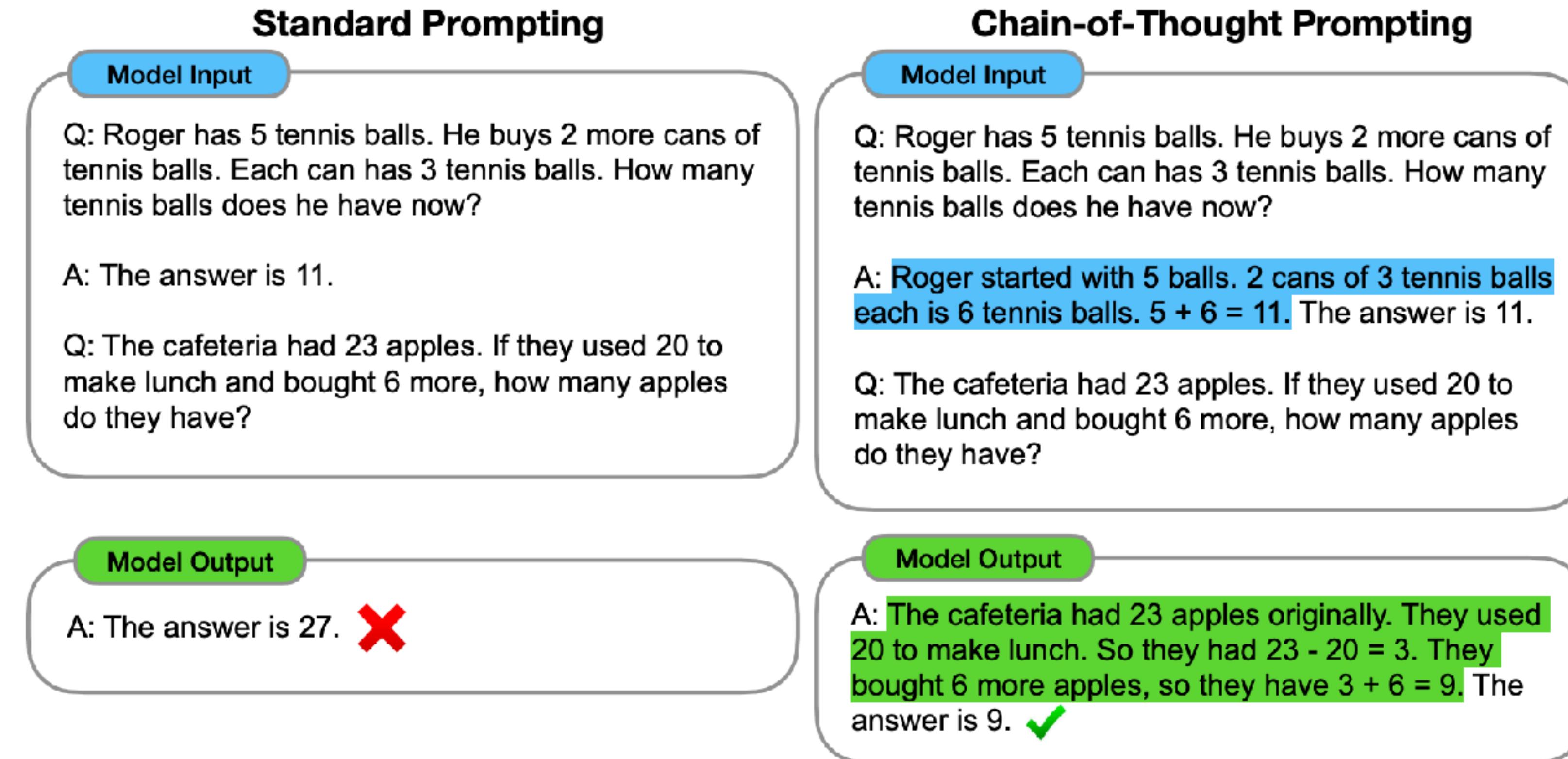


Figure 1: Chain-of-thought prompting enables large language models to tackle complex arithmetic, commonsense, and symbolic reasoning tasks. Chain-of-thought reasoning processes are highlighted.

# Register Tokens enable LLM world model exploration.

Text — Question: What is  $(3-4)/7$ ? <r><r><r>...<r> Answer:  $\frac{1}{7}$

Tokens — 74 157 323 99807 366 92 42 42 42 ... 42 844 323

# Register tokens improve arithmetic performance.

