

# Transformer as a hippocampal memory consolidation model based on NMDAR-inspired nonlinearity

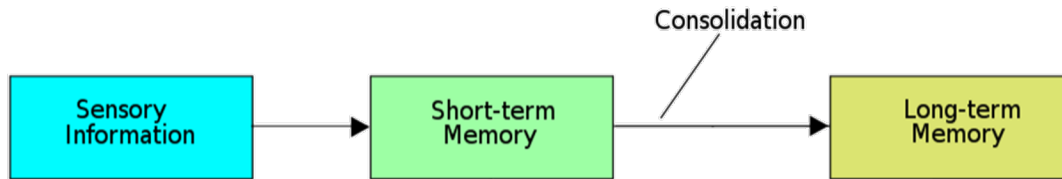
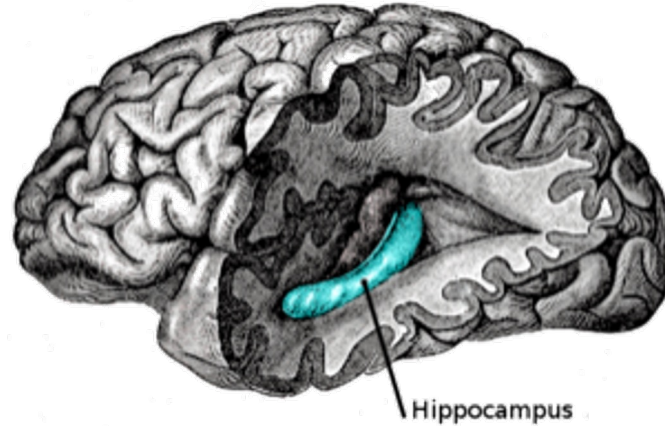
Dong-Kyum Kim<sup>1\*</sup> Jea Kwon<sup>2\*</sup> Meeyoung Cha<sup>1,3†</sup> C. Justin Lee<sup>2†</sup>



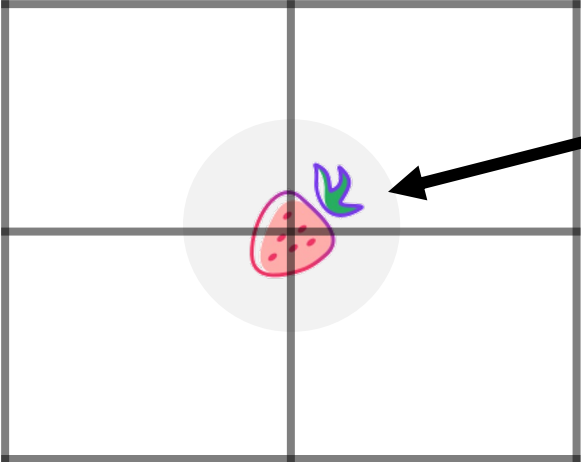
<sup>1</sup>Data Science Group, IBS    <sup>2</sup>Center for Cognition and Sociality, IBS    <sup>3</sup>School of Computing, KAIST

\*Equal contributions; †Corresponding authors.

# Memory formation is key in intelligent systems

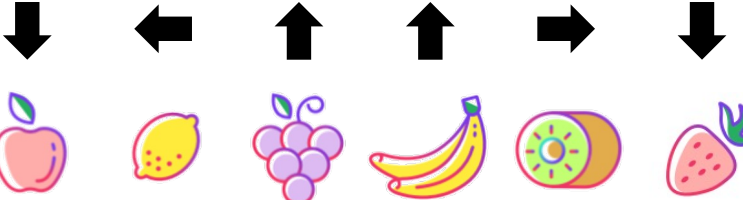


# Short-term working memory use transient information



**Trial 1**

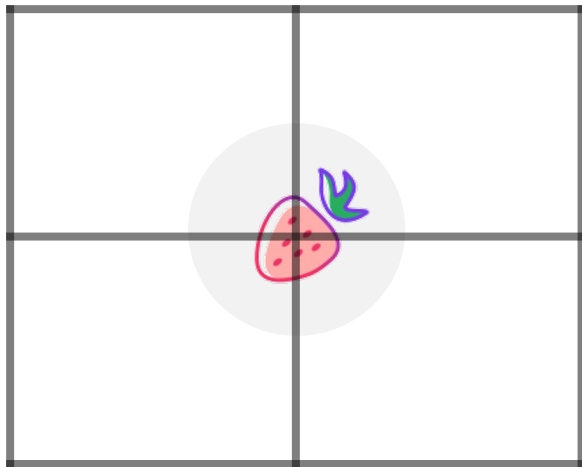
Action  
Observation



Q. What comes next?

→ Inferred by within-trial observations

# Long-term reference memory retrieve previous information



**Trial 2**

Action



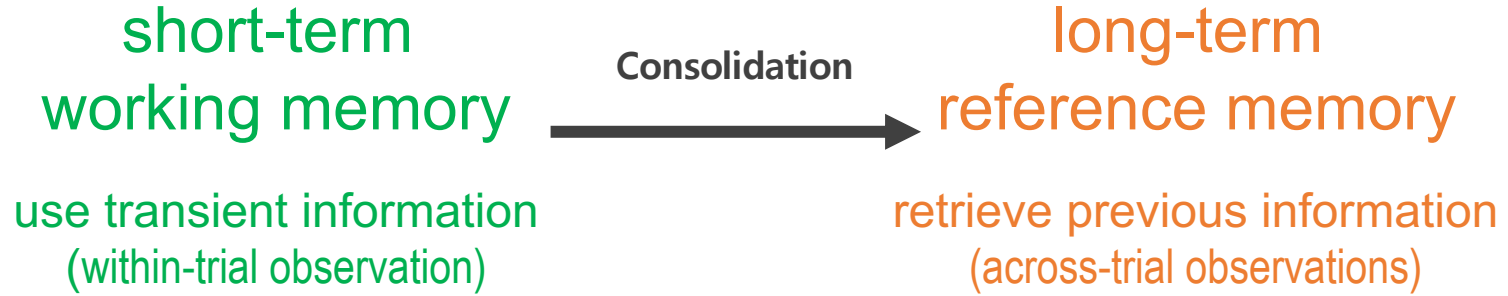
**Q. What comes next?**

Observation

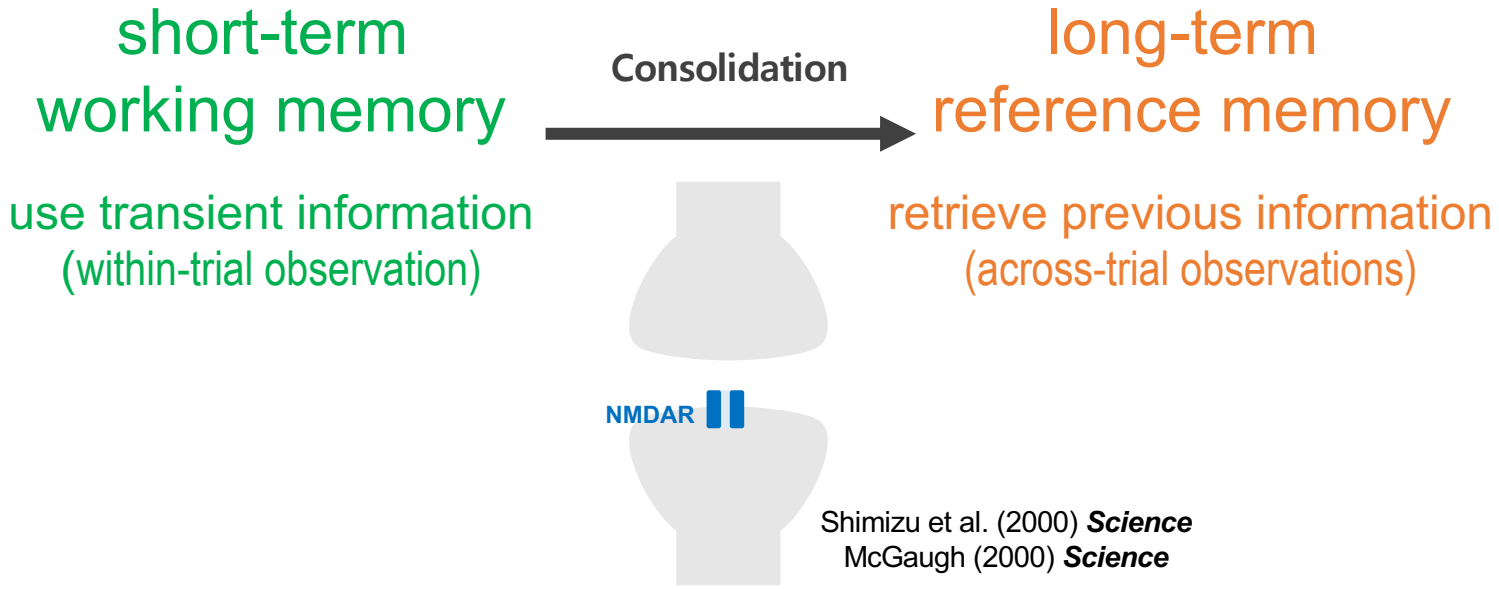


→ Inferred by across-trial observations

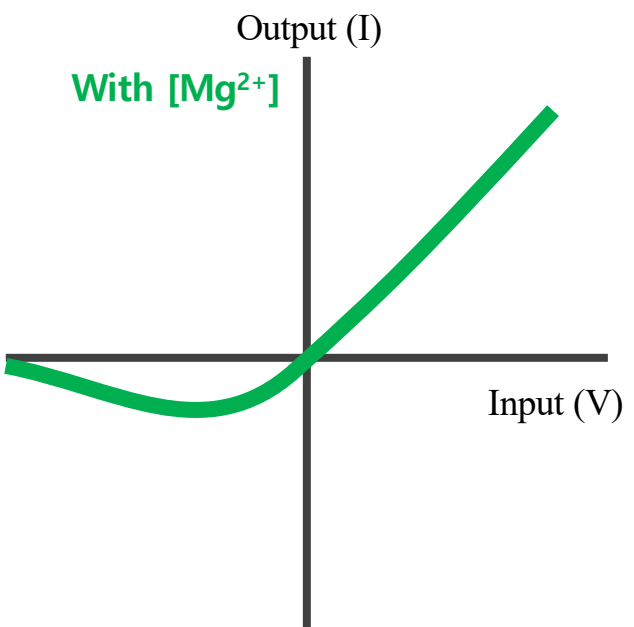
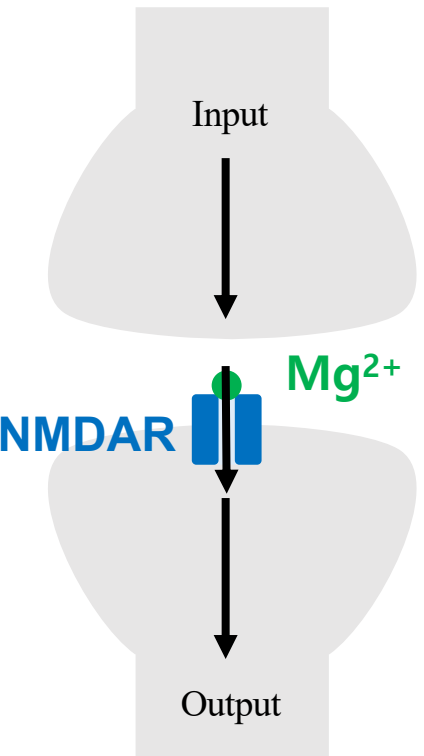
# Memory consolidation is a process that transforms STWM into LTRM



# NMDA receptor in brain makes this happen



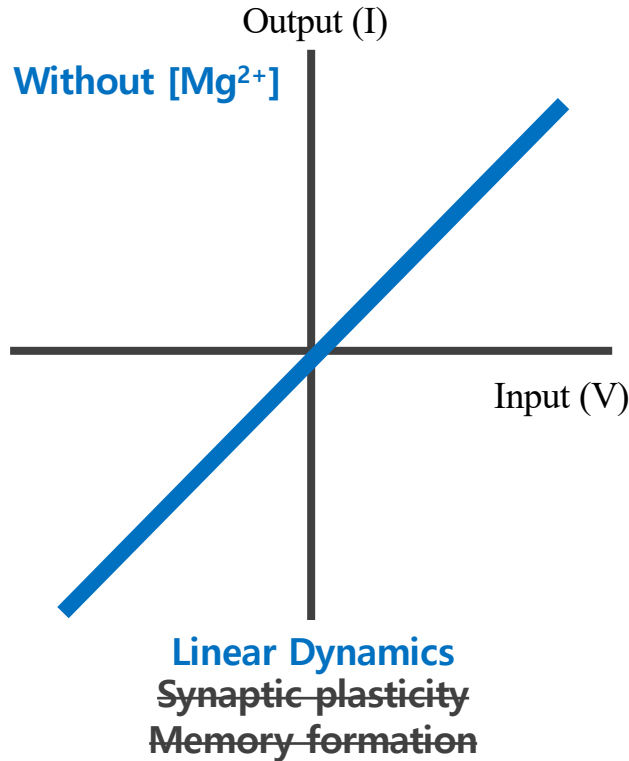
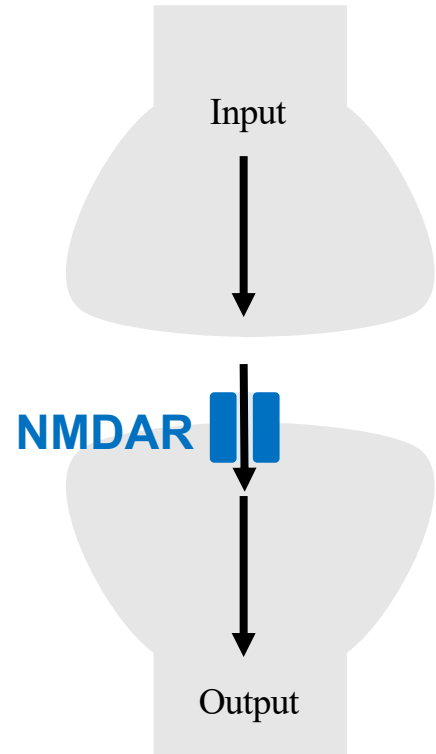
# NMDA receptor has nonlinearity



Nonlinear Dynamics  
Synaptic plasticity  
Memory formation

$$V \cdot \frac{1}{1 + \frac{[\text{Mg}^{2+}]}{K_{\text{Mg}^{2+}}} \cdot \exp(-\beta V)}$$

# NMDA receptor nonlinearity disappears without Magnesium ion (Mg<sup>2+</sup>)

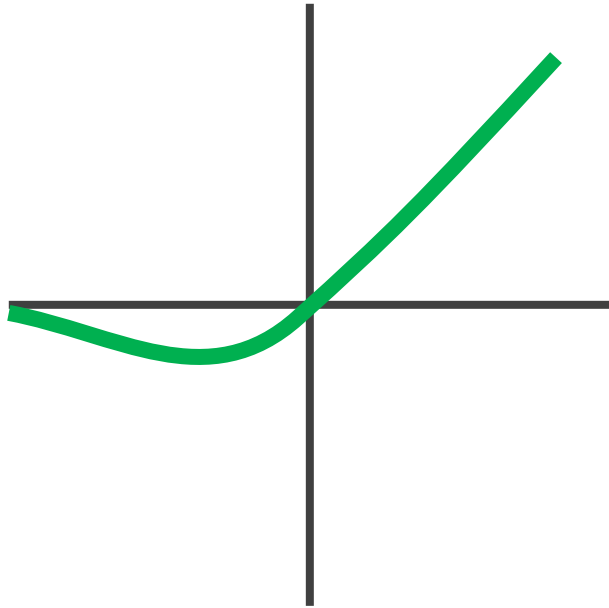


$$V \cdot \frac{1}{1 + \frac{0 \cdot \exp(-\beta V)}{K_{\text{Mg}^{2+}}}}$$

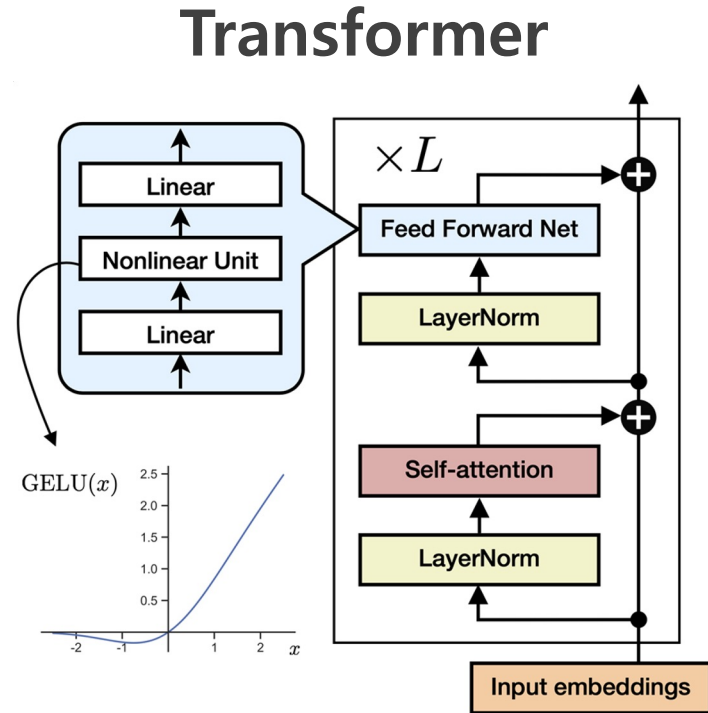
The equation shows the effect of zero magnesium concentration. A blue box highlights the term  $\frac{0 \cdot \exp(-\beta V)}{K_{\text{Mg}^{2+}}}$ , with a blue arrow pointing from the '0' in the numerator to the '0' in the denominator of the fraction in the equation above.



# We observe NMDAR nonlinearity in brain resembles the nonlinear activation function in Transformer

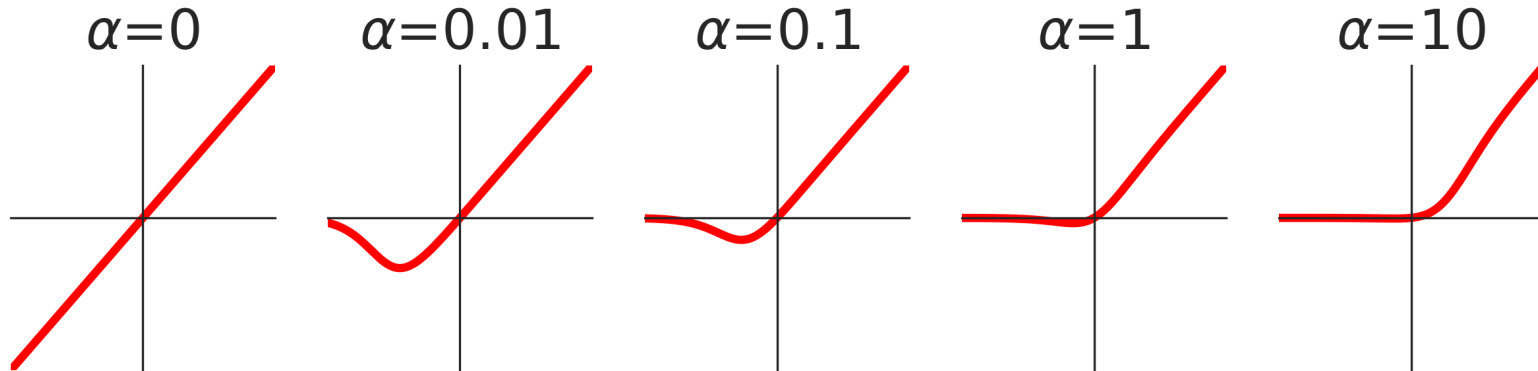


**NMDAR Nonlinearity**



We propose a new activation function for transformers that mimics brain's nonlinearity ( $\text{Mg}^{2+}$  of NMDAR)

$$\text{NMDA}_{\alpha}(x) = \frac{x}{1 + \alpha e^{-x}}$$



# RQ 1: How would NMDAR-like nonlinearity perform in Transformer?

Memory Consolidation  
STWM  $\rightarrow$  LTRM



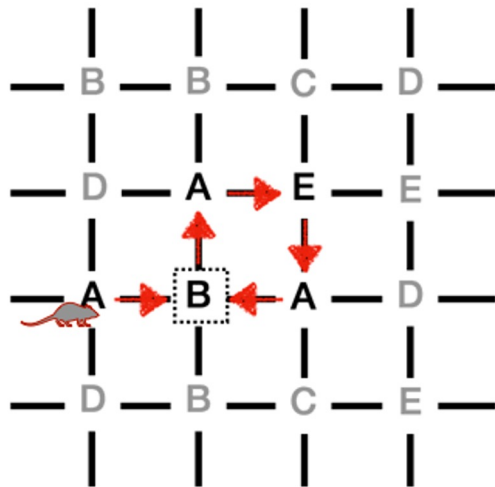
NMDAR  
Nonlinearity

Transformer

# We will test two memory functions:

## Short-term working memory

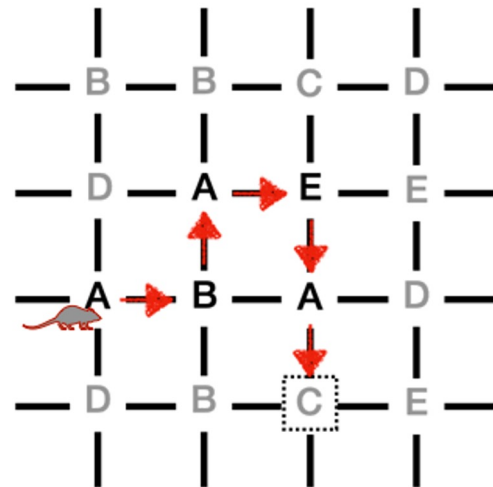
Visited place prediction



Action (a)                    → ↑ → ↓ ←  
Observation (x) A B A E A

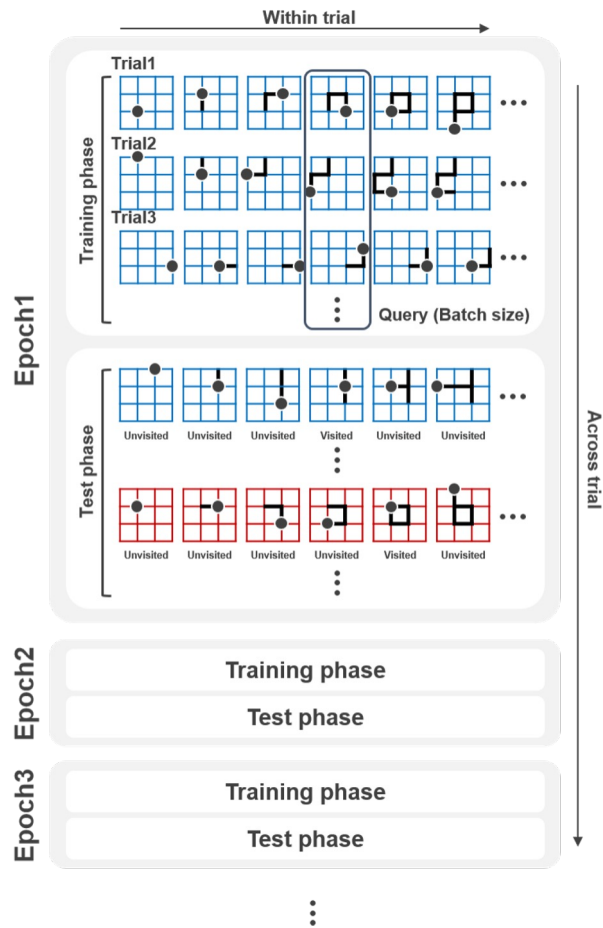
## Long-term reference memory

Unvisited place prediction



   → ↑ → ↓ ↓  
A B A E A

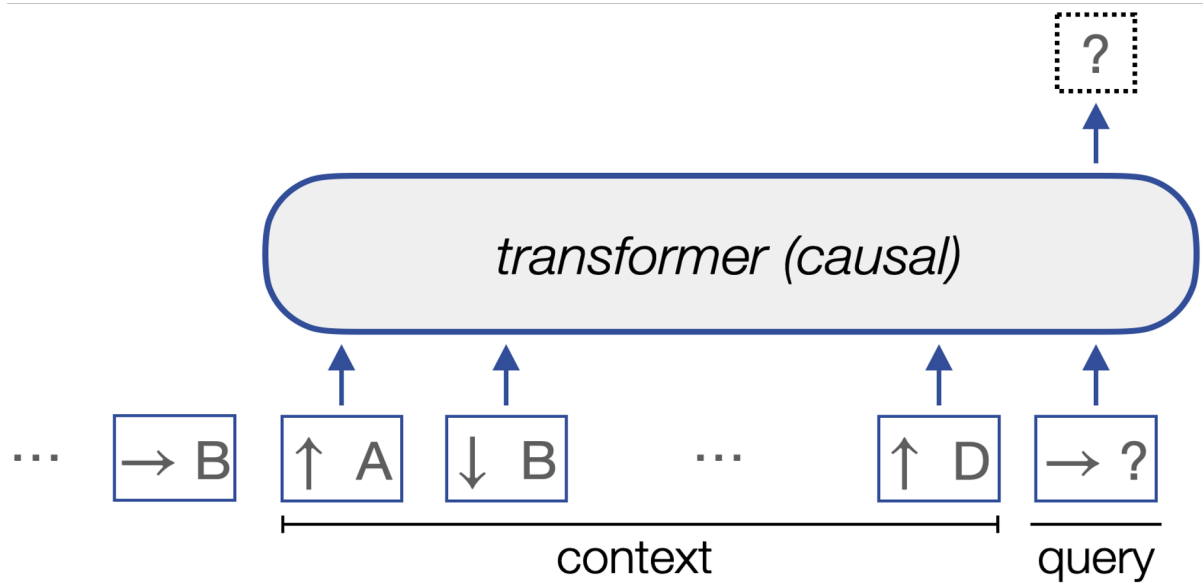
# by designing a 2D navigation task



## Experimental details

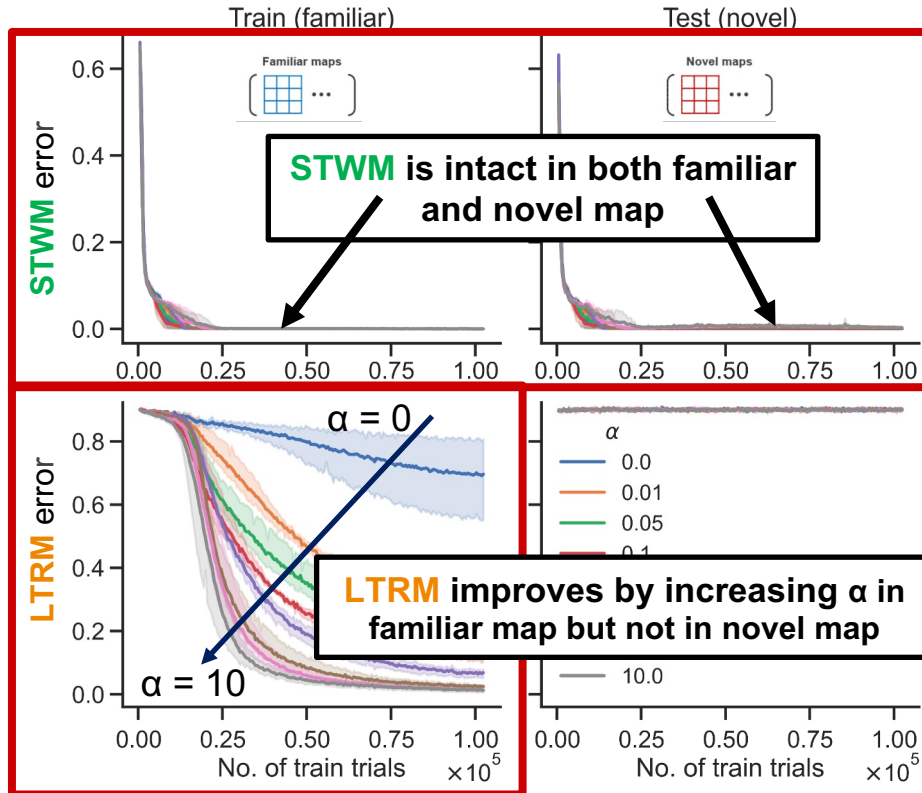
- We prepare  $N$  number of training maps.
- For each map, a random sensory observation among 10 letters are randomly placed at each node on each map.
- In a single trial, a randomly selected map from training maps is given to the agent.
- The agent starts at a random position and initiates a random walk on the selected map.

# Two types of memory error in Transformer

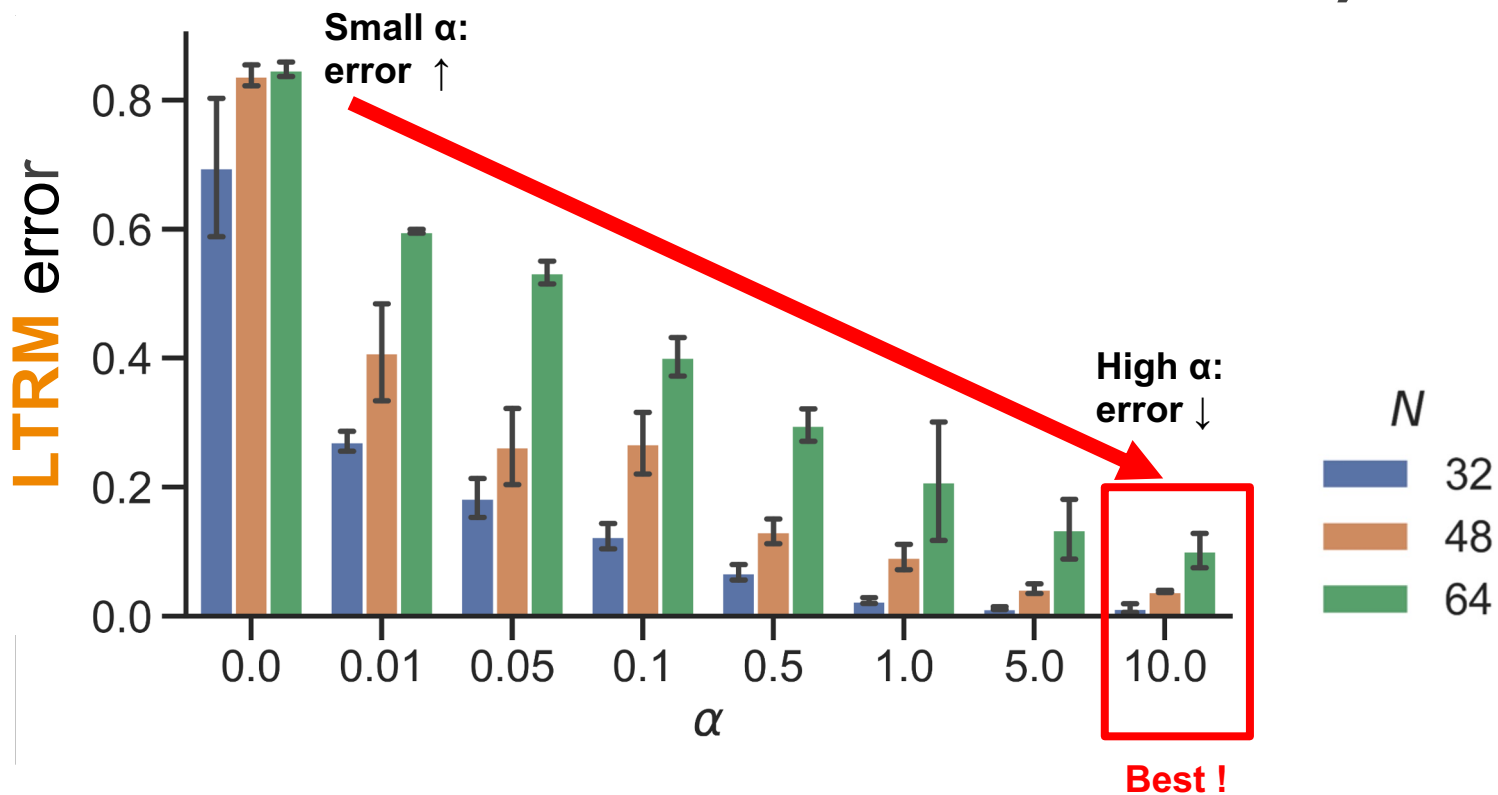


When answer is in the *context* - **STWM** error  
When answer is not in the *context* - **LTRM** error

# Our data show NMDA receptor nonlinearity controls **long-term reference memory (LTRM)**

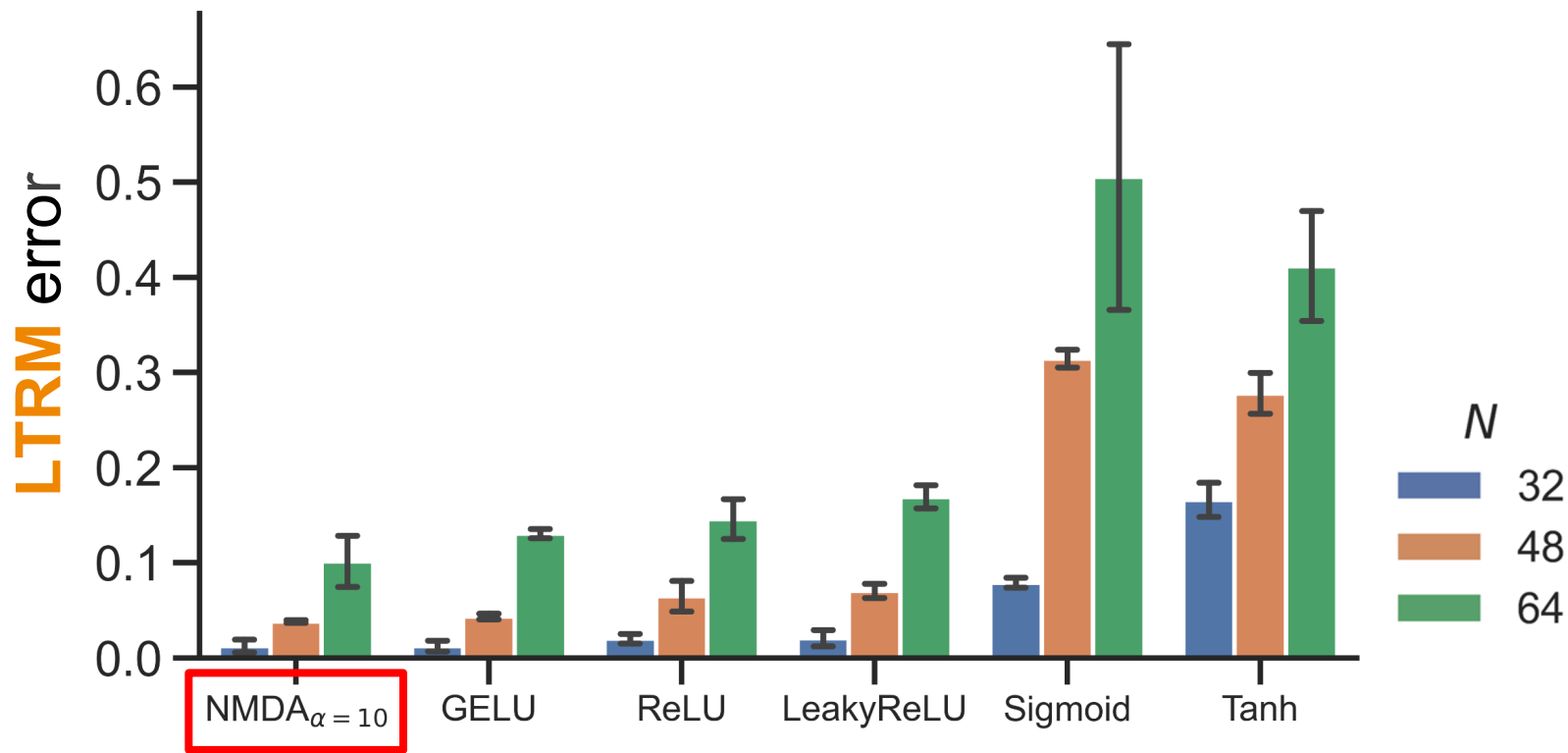


# LTRM improves with increasing $\alpha$ (similar to how animal brain works)

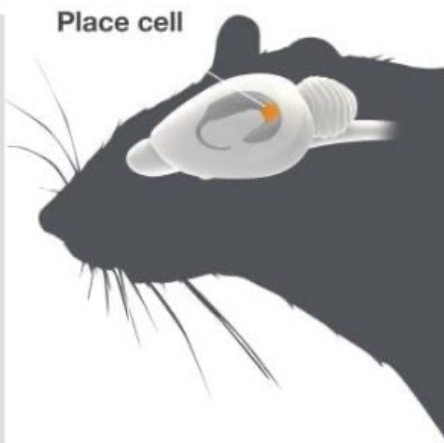




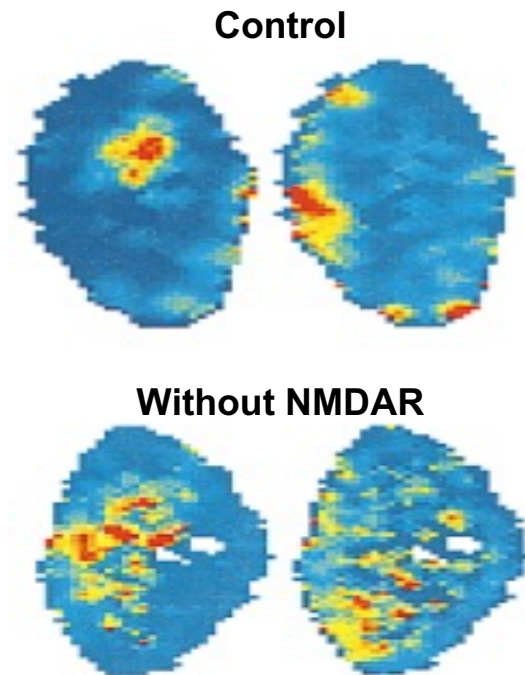
# The improvement saturates after certain $\alpha$ (similar to how animal brain works)



# RQ 2: Can NMDAR nonlinearity attribute to place cell emergence in Transformer?

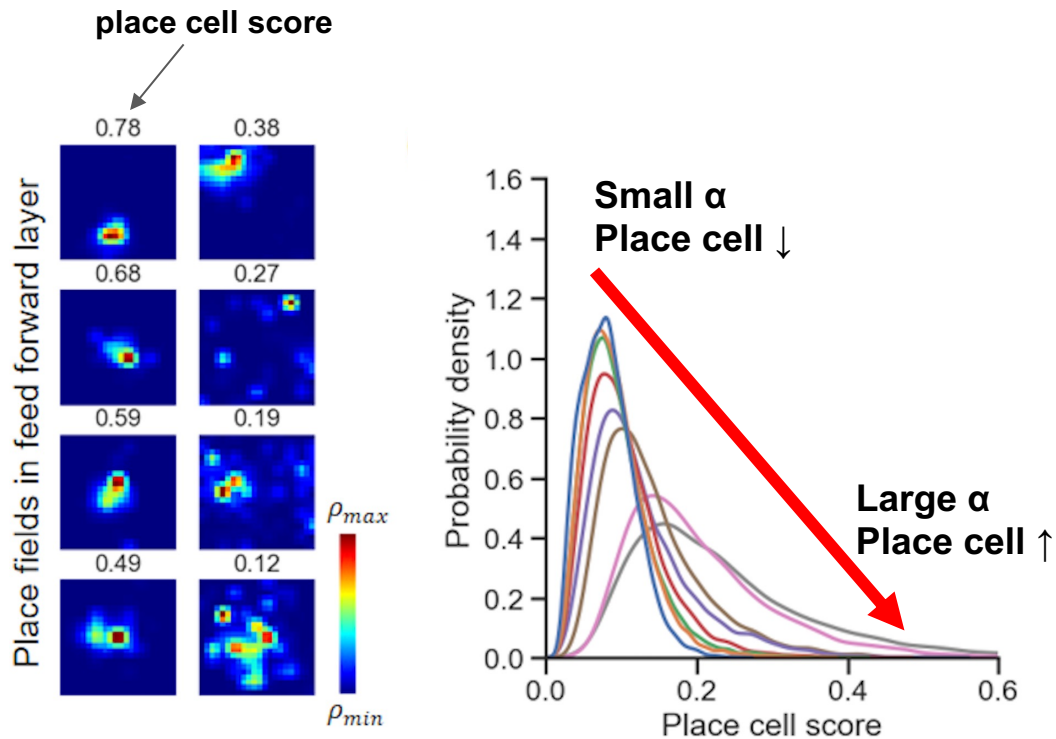


Creator: Ulrika Royen

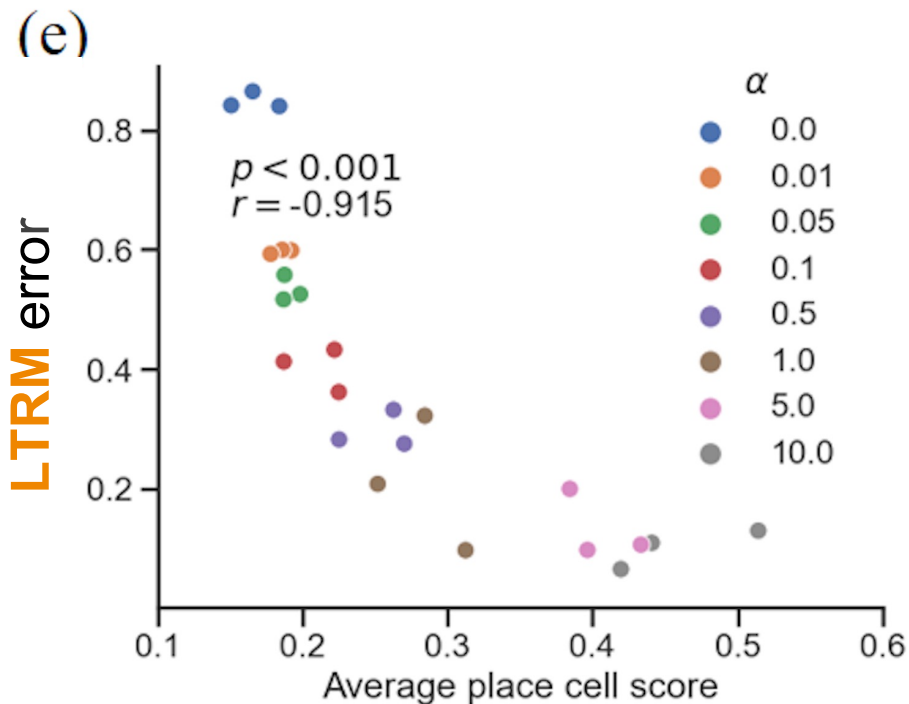


(McHugh et al., *Cell* 1996; Kentros et al., *Science* 1998)

# Our data show place cells emerge in the Transformer's FFN (feed-forward network)



# Our data show transformers with high long-term memory forms more place cells (similar to brain)



# Summary

1. We proposed a new activation function that is inspired by the nonlinear dynamics of NMDA receptors in brain (**NMDA $\alpha$** ) where  $\alpha$  mimics the  $\text{Mg}^{2+}$  concentration level.

$$\text{NMDA}_\alpha(x) = \frac{x}{1 + \alpha e^{-x}}$$

2. We developed a method for accessing the **reference memory**.
3. We evaluated the reference memory errors of **transformer** models with  $\text{NMDA}_\alpha$ . The results shows that *reference memory can be controlled by  $\alpha$* .
4.  $\text{NMDA}_\alpha$  with  $\alpha=10$  shows the **best reference memory performance** when compared to other widely used nonlinear activation functions.
5. We demonstrated the emergence of **place cells** in feed-forward networks of transformer for the first time.
6. Reference memory is impaired when the value of  $\alpha$  in  $\text{NMDA}_\alpha$  is low and this **resembles long-term memory loss in brain**; low  $\text{Mg}^{2+}$  concentration in brain causes long-term memory loss.