

# Consciousness & AI

---

15 June 2023

Ziyuan Ye

# Overview

## □ What is consciousness?

- I. Phenomenal Consciousness
- II. Cognitive (Access) Consciousness

## □ Where does the consciousness come from?

- I. Integrated Information Theory
- II. Global Neuronal Workspace Theory
- III. Dendrites Integration Theory

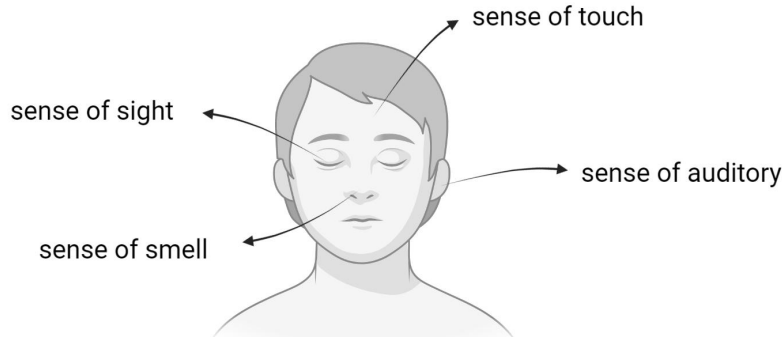
## □ Discussion about consciousness and next-generation AI models.

# Overview

- **What is consciousness?**
  - I. **Phenomenal Consciousness**
  - II. **Cognitive (Access) Consciousness**
  
- **Where does the consciousness come from?**
  - I. **Integrated Information Theory**
  - II. **Global Neuronal Workspace Theory**
  - III. **Dendrites Integration Theory**
  
- **Discussion about consciousness and next-generation AI models.**

# What is consciousness?

## Phenomenal consciousness: subjective experience



An analogue mapping between the nervous system and patterns of information in the environment or in the body.

Highly related to our perception!

## Cognitive (Access) consciousness: cognition and processing of certain information

**Cognitive Ability:** attention, memory, reasoning, problem-solving, language processing...

**Knowledge and Learning:** knowledge acquisition, processing, storage, application...

**Emotional Cognition:** understanding and recognition of our own emotions and those of others

**Self-awareness:** understanding and perception of ourselves (e.g., self-concept)

...

# Overview

## □ What is consciousness?

- I. Phenomenal Consciousness
- II. Cognitive (Access) Consciousness

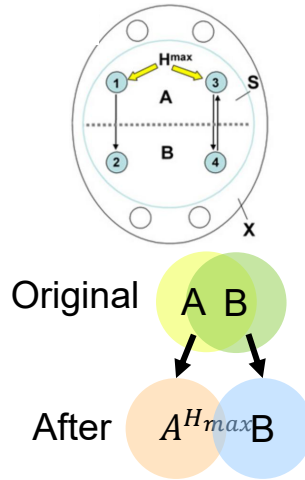
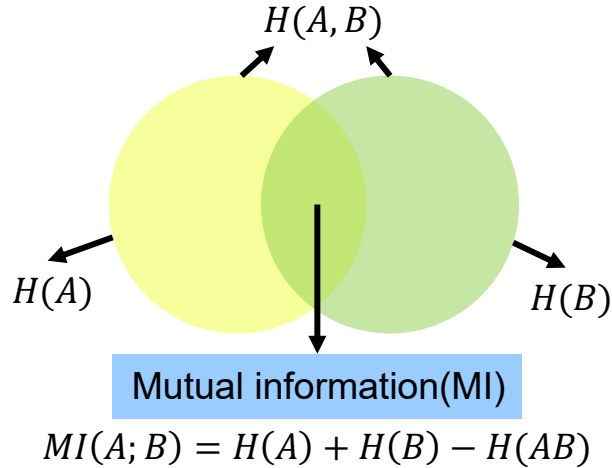
## □ Where does the consciousness come from?

- I. Integrated Information Theory
- II. Global Neuronal Workspace Theory
- III. Dendrites Integration Theory

## □ Discussion about consciousness and next-generation AI models.

# Integrated Information Theory (IIT)

Consciousness corresponds to the **capacity** of a system to **integrate information**. The **state of consciousness** can be measured as the **causally effective information (EI)  $\phi$  value** of a complex of elements.



**Effective information from A to B**

$$EI(A \rightarrow B) = MI(A^{H_{max}}; B)$$

**Effective information between A and B**

$$EI(A \leftrightarrow B) = MI(A^{H_{max}}; B) + MI(A; B^{H_{max}})$$

**Normalized effective information between A and B**

$$EI(\widetilde{A \leftrightarrow B}) = EI(A \leftrightarrow B) / H^{max}(A \leftrightarrow B)$$

$$H^{max}(A \leftrightarrow B) = \min(A^{H_{max}}; B^{H_{max}})$$

**Minimum information bipartition (MIB)**

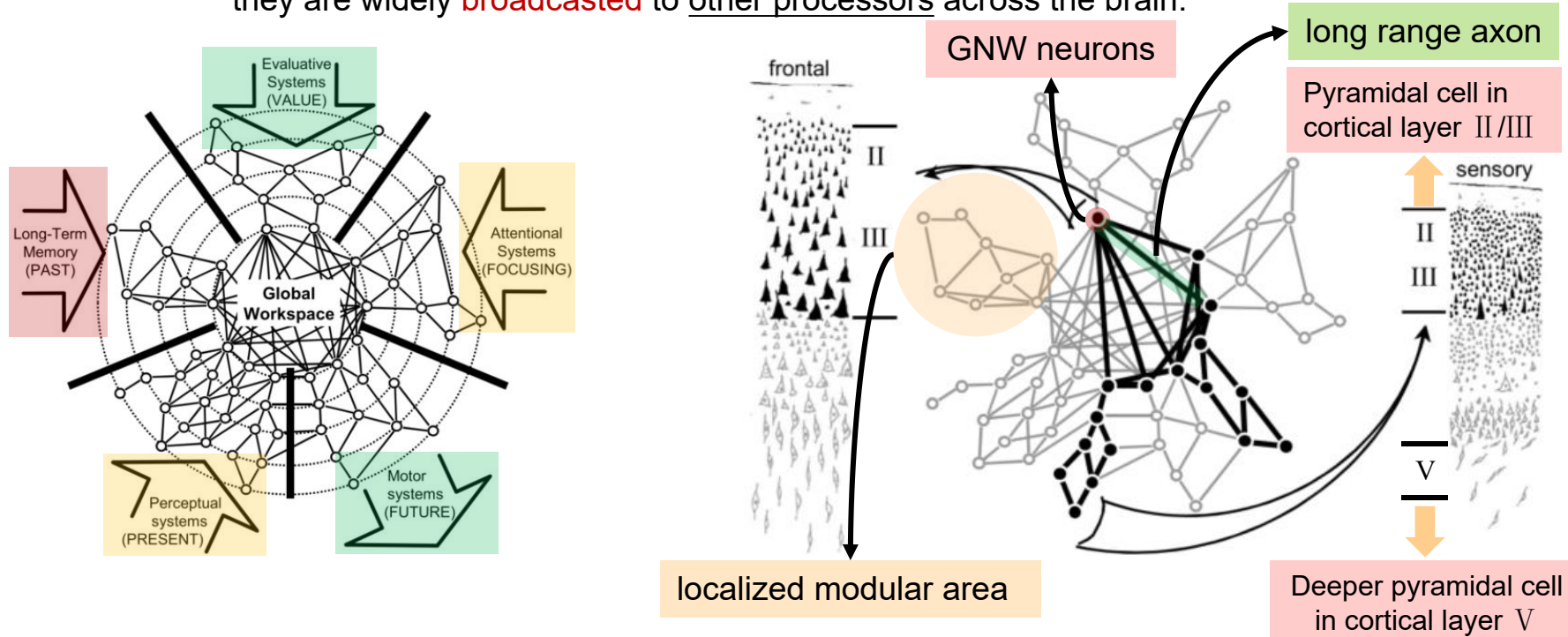
$$MIB A \leftrightarrow B = \operatorname{argmin}\{EI(\widetilde{A \leftrightarrow B})\}$$

**Information integration of subset S**

$$\phi(S) = EI(MIB A \leftrightarrow B)$$

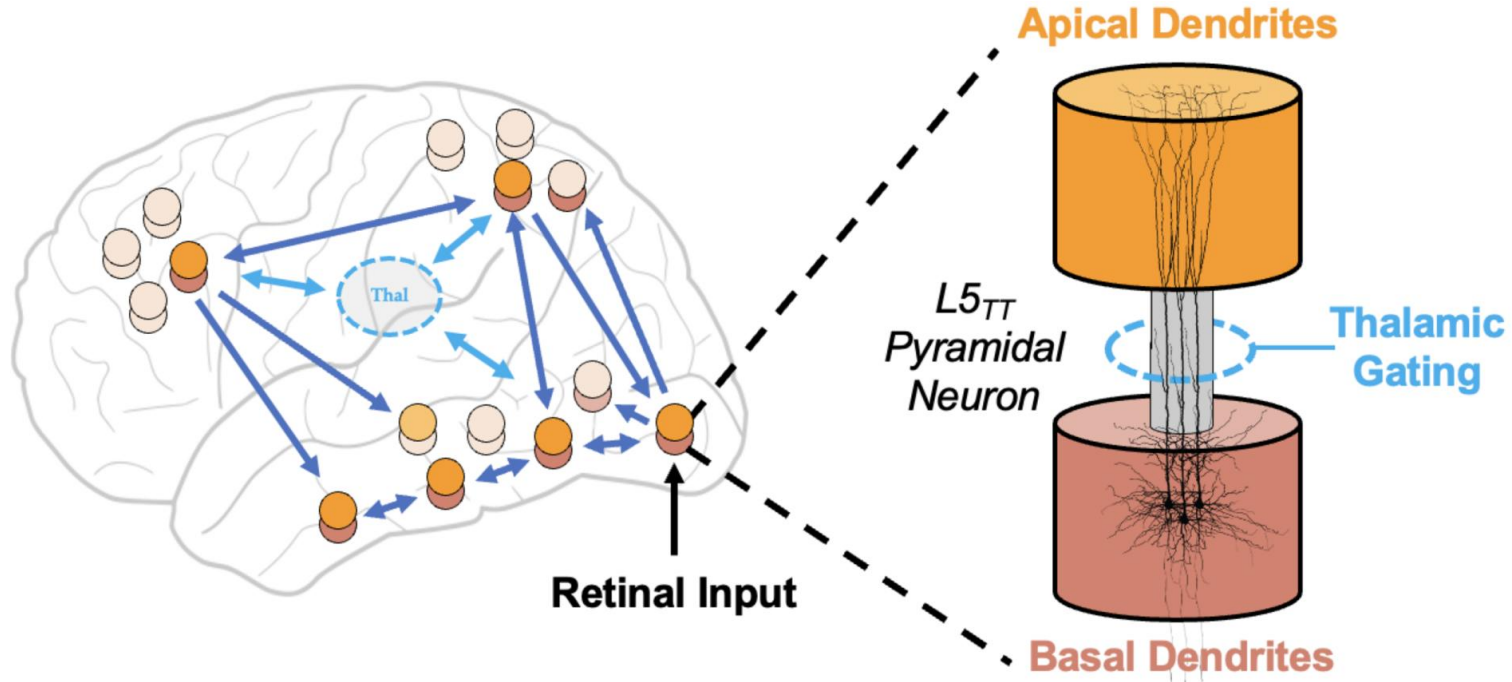
# Global Neuronal Workspace Theory (GNWT)

Perceptual contents are acted upon by localized processors, only become **conscious** when they are widely **broadcasted** to other processors across the brain.



# Dendrites Integration Theory (DIT)

Conscious states and content emerge in **deep pyramidal neurons**.

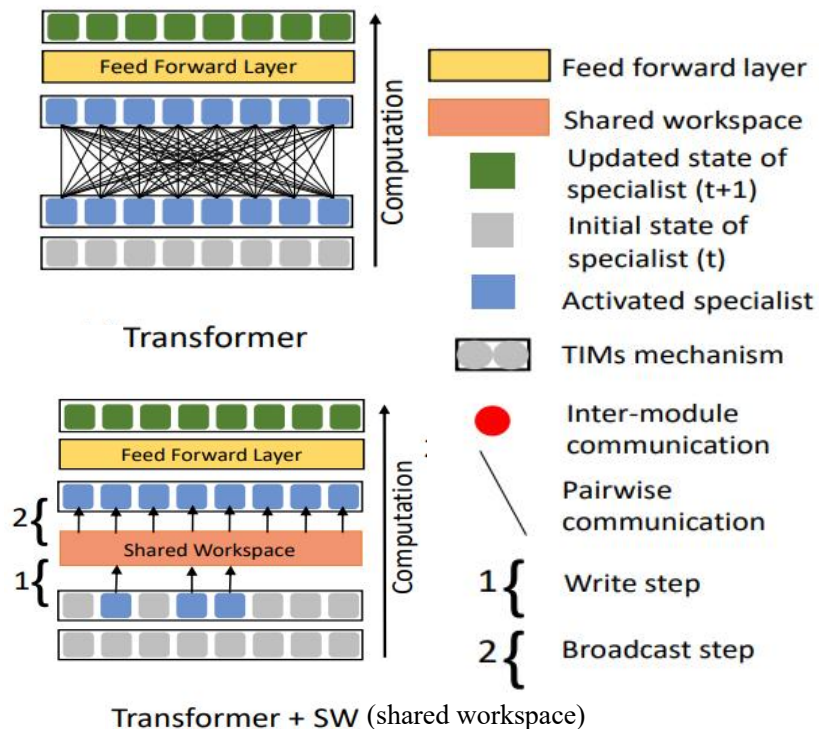




# Overview

- **What is consciousness?**
  - I. **Phenomenal Consciousness**
  - II. **Cognitive (Access) Consciousness**
  
- **Where does the consciousness come from?**
  - I. **Integrated Information Theory**
  - II. **Global Neuronal Workspace Theory**
  - III. **Dendrites Integration Theory**
  
- **Discussion about consciousness and next-generation AI models.**

# Transformer + GNWT



## Two step process

**Write step: write information into workspace**

**Embedding:**  $R \in \mathbb{R}^{n_s \times n_h}$  (specialist, hidden state)

Each row of  $R$ :  $h_t^k \in \mathbb{R}^{1 \times n_h}$ ,  $k \in \{1, \dots, n_s\}$ ,  $t$ : stage number

**Shared workspace:**  $M \in \mathbb{R}^{1 \times n_h}$

**Query:**  $\tilde{Q} = M\tilde{W}^q$

**Key/Value:**  $R$

**Update Shared Workspace:**  $M \leftarrow \text{softmax}\left(\frac{\tilde{Q}(R\tilde{W}^e)^T}{\sqrt{d_e}}\right) R\tilde{W}^v$

Soft competition

Hard competition: Select Top-k specialists

**Broadcast step: broadcast information from workspace**

**Query:**  $\hat{q}_k = h_t^k \hat{W}^q$

**Key:**  $\hat{\kappa}_j = (m_j \hat{W}^e)^T$

**Value:**  $\hat{v}_j = m_j \hat{W}^v$

$$h_t^k \leftarrow h_t^k + \sum_j \text{softmax}\left(\frac{\hat{q}_k \hat{\kappa}_j}{\sqrt{d_e}}\right) \hat{v}_j$$

# Discussion

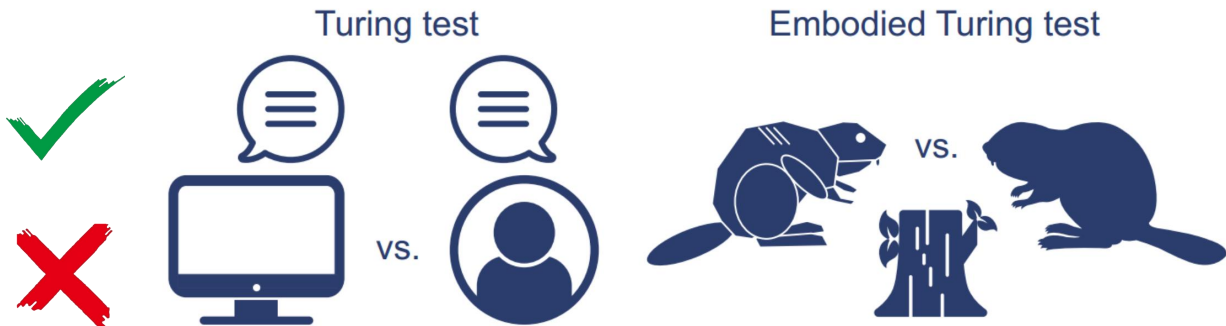
1. Measure a **AI model's consciousness** by using Integrated Information Theory. Improve model consciousness by optimizing model architecture design.
2. Design novel AI model architecture with Global Neuronal Workspace Theory / Dendrites Integration Theory.

How can we go beyond LLM?

**Cognitive (Access)**

**consciousness:** cognition and processing of certain information

**Phenomenal consciousness:**  
subjective experience



**Thanks for your attention!**

---