# Mind And Consciousness Global Neural Workspace Mathematical And Computational Modeling

## J N Tavares

**Abstract**

*This paper analyzes the theory of GLobAL NeURAL WoRKspAce THeoRy (GNW), by Stanislas Dehaene and Jean Pierre Changeux, developed from the theory of GLobAL WoRKspAce THeoRy (GWT), proposed by Bernard Baars in the 1980s. GNW is a cognitive model of consciousness that seeks to explain how the brain integrates, selects, and disseminates information among dif- ferent modules or specialized systems, such as vision, hearing, memory, language, etc. Next, several proposals for mathematical and computational modeling are made, using various neural network architectures used in "Deep Learning.".*

## I.    Introduction

The Global Neuronal Workspace (GLobAL NeURonAL WoRKspAce, GNW) model has been one of the most influential theories in cognitive neuroscience for explaining the mechanisms underlying human consciousness. Initially pro- posed by Dehaene, Changeux, and Naccache ([Dehaene & Changeux 2006], [Dehaene & Changeux 1998]), GNW suggests that consciousness emerges from the integration and dissemination of information across a global network of neurons, which acts as a "workspace" for processing and sharing informa- tion between different brain regions. This model offers a robust framework for understanding how the brain selects, amplifies, and maintains relevant infor- mation, enabling conscious experience.

In this work, we explore an innovative approach to integrate the GNW model with HopfIeLD netwoRKs, a type of recurrent neural network known for its abil- ity to store and retrieve patterns stably. Hopfield networks, with their attractor- based dynamics, offer an interesting perspective for modeling the stability and resilience of the global workspace proposed by GNW. By combining these two concepts, we seek to investigate how attractor dynamics can contribute to the integration and maintenance of information in the global workspace, providing a deeper understanding of the neural mechanisms of consciousness.

The purpose of this work is, therefore, to present a hybrid model that inte- grates GNW with Hopfield networks, exploring how attractor dynamics can be applied to simulate conscious cognitive processes. Through simulations and theoretical analyses, we hope to contribute to advancing our understanding of consciousness by offering a new perspective on how the brain implements the global workspace and how the dynamic stability of neural networks can play a crucial role in this process.

It is worth mentioning that the GNW model had a precursor – the GLobAL WoRKspAce THeoRy (GWT), proposed by Bernard Baars in the 80s of the last century ([Baars 1997], [Baars 1988]). This is a cognitive model of conscious- ness that seeks to explain how the brain integrates, selects, and disseminates information among different modules or specialized systems, such as vision, hearing, memory, language, etc. It is worth briefly describing what it consists of.

The GWT presents consciousness as a "mental theater," where various sources of unconscious processing compete to "access the center stage," making certain information globally accessible to the cognitive system. The GWT argues that consciousness arises from the distribution and coordination of information across a global space accessible to multiple brain systems. 1  THe Key con- cept Is tHIs: when information is amplified and disseminated throughout the brain network (through a neuronal "ignition"), it triggers the synchronization of various processes (attention, memory, planning) into a common workspace. Thus, consciousness is the result of a global selection and diffusion of activity that makes information accessible to multiple cognitive functions. It's a sort of virtual "switchboard" in which information, after being amplified by attention, becomes widely available to all brain modulesmemory, language, planning, etc.

Most brain processing is unconscious. Consciousness emerges when cer- tain content becomes predominant and "earns the right" to be widely dissem- inated. The critical aspect of consciousness is global availability, not so much the content itself, but the fact that it can be accessed, used, and manipulated by multiple specialized systems.

## II.    Baars's Theater Of Consciousness Metaphor

Bernard Baars ([Baars 1997], [Baars 1988]) uses the analogy of a theater to explain how consciousness works. Just like in a theater, the mind has:
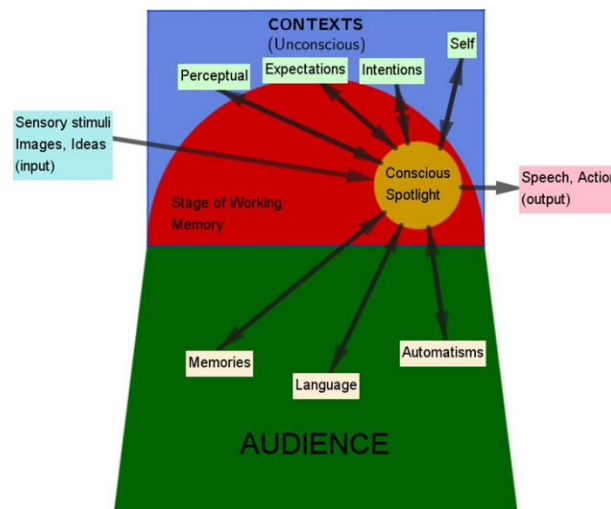


Figure 1: Metaphor from Bernard Baars's THeAteR of ConscIoUsness.

The StAge, which represents woRKIng MeMoRy – the space where you unfold the thoughts, images, and sensations of which you are conscious at any given moment. It is a limited space, just like a theater stage, which can only accom- modate a restricted number of actors and props; our consciousness can only hold a small number of items simultaneously. Events on stage are dynamic and changing: actors enter and exit the scene, lights change, and the scenery transforms, reflecting the fluid and transient nature of consciousness.

The Actors (Mental Contents) are the thoughts, sensations, images, memo- ries, and other forms of information that compete for access to the stage (con- sciousness). Only a few "actors" can be on stage at the same time, reflecting the limited capacity of conscious attention. The quality and intensity of their performance (level of neuronal excitation, salience, relevance) determine their likelihood of being selected to occupy center stage.

The Attentional Spotlight (Selective Focus) represents the mechanism of se- lective attention, which illuminates certain actors (mental contents) on stage, making them the focus of our consciousness. Attention is selectively directed to certain aspects of the experience (a speaker's voice, a captivating image, a persistent thought), while others remain in the shadows. Metaphorically, attention functions as a filter through which only certain content passes.

The AUDIence (Unconscious Processors) represents the wide range of un- conscious mental processes that operate in parallel, processing information, generating emotions, controlling the body, etc. The "audience" is not directly aware of what is happening onstage, but receives the information that is dif- fused from there.

The BAcKstAge CRew (Unconscious Context) refers to the unconscious pro- cesses that shape the conscious experience (like stage setup). Expectations, beliefs, memories, and other background information influence the way we perceive and interpret what is happening onstage (consciousness). They act as "commands" that select the direction of the performance stage.

Finally, the DIRectoR (Executive Functions), which represents the executive control systems located in the frontal cortex that oversee and coordinate activ- ities on stage. The director makes decisions about which actors should appear onstage, which props should be used, and how the story should be told, re- flecting the role of executive functions in regulating and organizing conscious thought. In summary, the functions of Baars's Theater Metaphor are:

- ExpLAIn tHe LIMIts of conscIoUs cApAcIty: Just as a stage can only accom- modate a limited number of actors, our consciousness can only process a limited amount of information at any given moment.
- ILLUstRAtIng tHe IMpoRtAnce of tHe UnconscIoUs: Most mental activity occurs outside our awareness, as do the actors behind the scenes and the technical crew who make the show possible.
- GIvIng sense to tHe fRee fLow of conscIoUsness: The actors enter and leave the stage, the lights go on and off, the scenery changes, reflecting the dynamic and fluid nature of our thoughts and sensations.
- FAcILItAtIng UnDeRstAnDIng of tHe fUnctIons of conscIoUsness: The metaphor of theater helps to understand how consciousness allows us to integrate information, plan actions, solve problems, and interact with the world flexibly.

Baars's metaphor has, however, been widely criticized. Some have criti- cized it for implying the existence of an internal observer (a "homunculus") that watches the spectacle of consciousness. To this, Baars responded that the "audience" in theater is not a conscious homunculus, but rather the vast array of unconscious processors in the brain, each acting according to its own competence.

Baars also argues that consciousness is not a physical location in the brain, but rather a process of information dissemination.

In sHoRt: Baars's Theater of Consciousness metaphor offers an intuitive and powerful way to understand how consciousness emerges from the inter- action between conscious and unconscious mental processes. By highlighting the importance of attention, integration, and the dissemination of information, this metaphor provides a rich framework for investigating the neural mecha- nisms underlying conscious experience.

## III.    The Global Neural Workspace (GNW)

Stanislas Dehaene and Jean-Pierre Changeux (see [Dehaene & Changeux 2006], [Dehaene & Changeux 1998], when developing the GLobAL NeURAL WoRKspAce (GNW) model, took Bernard Baars's Theater of Consciousness metaphor and "neuronalized" it, adding biological specificity, experimental testability, and computational rigor. The main changes and improvements were as follows:

The replacement of the "mental theater" with a neural architecture. Al- though Baars sees his metaphor as a cognitive architecture, with actors and spotlights representing functional processes, he does not define where this theater is located in the brain or how it is physically implemented.

Dehaene and Changeux locate the "global work" in a specific neural network involving areas of the prefrontal and parietal cortex. and cingulate2, long-range connections over long distances, pyramidal neurons, etc.

The GNW model identifies a true "neuronal work" – a distributed network, composed primarily of areas of the prefrontal, parietal, and cingulate cortex, connected by long-range axons. These "workn¨ eurons receive signals from var- ious sensory areas and can amplify and globally process the selected content.

Replacing the "attentional spotlight" with mechanisms of neuronal ignition. For Baars, the spotlight is an abstract concept of selection and amplification. Dehaene Changeux propose specific neuronal mechanisms for this amplifica- tion.

The GNW introduces the concept of neuronal ignition – when a stimulus reaches a certain threshold (by strength, novelty, attention), abruptly activates a pattern of sustained, large-scale activity (seen in EEG3, MEG4, fMRI5). Only content that triggers this sudden ignition becomes conscious; others remain subliminal or unconscious. The "ignition" is a sudden wave of synchronized activity that spreads throughout the workspace, sustained by neurons with long-range axons and specialized synapses. Feedback (reentry) between the cortical areas and the thalamus reinforces these signals, keeping the informa- tion active in working memory.

Integrated activity instead of "actors" and "audience"! For Baars, the "ac- tors" were pieces of information, and the audience the receivers of that infor- mation. For Dehaene Changeux, this distinction is replaced by a more complex system of interactions. distributed among neurons that form complex circuits.

Incorporating executive functions into feedback. The metaphor has a direc- tor, but it doesn't quite explain how this figure manages everything. Dehaene Changeux use the recurrent cycle, mediated by the thalamus, to explain how the networks of consciousness interact.

Replacing subjective reporting with objective measurements. The events of Baars' theater are highly dependent on what a person or observer sees. For Dehaene Changeux, the model allows for the identification of events that could relate to the experience in neurobiological terms, making the existence of subjective reporting unnecessary and testing the theory.

Let's delve into more detail into the main concepts that characterize the functioning of Dehaene Changeux's GNW

- ConscIoUsness As gLobAL AMpLIfIcAtIon (bRoADcAst). Consciousness emerges when information processed locally (in sensory, perceptual, etc.) zones is "amplified" and distributed throughout the brain network, becoming available to multiple systems (attention, memory, planning).
- IgnItIon (IgnItIon). There is a critical moment ("ignition") when certain patterns of neuronal activity expand rapidly and globally, for example, via gAMMA syncHRony6, correlating with conscious experience.
- ObjectIve ExpeRIences. Using techniques such as fMRI, EEG, magne- toencephalography, and paradigms such as the "visual mask," Dehaene identifies brain markers of consciousness, such as sustained activation of the parietal-frontal cortex. He advocates an empiricist stance – con- scious states can (and should) be investigated by objective methods, with- out falling into metaphysical speculations.
- MeAnIng AnD LIMIts of ConscIoUsness. Dehaene distinguishes strongly between unconscious (extensive and efficient) and conscious (limited, se- quential, but flexible) processing.

For Dehaene, the "GLobAL WoRKspAce" (GNW), consciousness allows access, coordination, and manipulation, functioning as a high-level RAM. Conscious- ness is a neurocomputational phenomenon, emerging from the selection and diffusion of information in complex brain networks, capable of being experimentally tested in the laboratory.

The GNW model is compatible with processes of attention, working memory, and conscious recognition. Sustained activation of regions of the frontopari- etal cortex is a strong neuronal marker of these states of consciousness. Each module works in parallel, unconsciously. The "GLobAL WoRKspAce" is a means of integrating and circulating information between modules. Becoming "con- scious" means being amplified into this common space. Furthermore, GNW explains many phenomena such as blocks of consciousness, attention, multi- tasking, amnesia, etc.

Dehaene and Changeux developed network models (simulated neural net- works) that reproduce experimental phenomena: conscious access, blocking, masking, "all-or-none" neural access, etc. GNW predicts clear experimental signatures: sustained fronto-parietal activation, gamma synchrony, bursts, delays in conscious reports. It has inspired many paradigms of neuroscience and experimental psychology tested in the laboratory with patients, anesthesia studies, sleep studies, etc.

In sUMMARy: The original theory (GWT, Baars) paved the way for viewing consciousness as a phenomenon of global content integration and diffusion. GNW made this model scientific, computationally simulable, and testable by identifying the actual circuitry, dynamics, and neurophysiological predictions of the "global workspace" in the human brain.
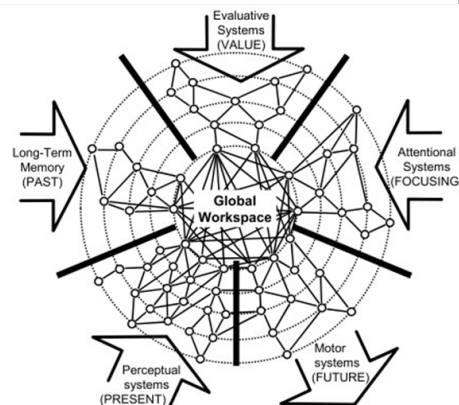
## IV.     From GWT To GNW. Formal Aspects



Figure 2: Global Neuronal Workspace. Original diagram from Dehaene et al. (1998).

Legend Fig. 2. GLobAL NeURonAL WoRKspAce (GNW). Original diagram from Dehaene et al. (1998) illustrating the principles main aspects of the GNW hy- pothesis: local and specialized cortical processors are linked, at the central level, by a central set of highly interconnected areas containing a high density of large pyramidal neurons with long-distance axons. At any time, this archi- tecture can select information within one or several processors, amplify it, and transmit it to all other processors, thus making it consciously accessible and available for verbal reporting. Recent studies tracking global cortical connec- tivity of feedforward and feedback confirm a bowtie architecture with a central core composed primarily of parietal and prefrontal areas, forming a structural bottleneck capable of routing information between other cortical processors.

The CentRAL HypotHesIs of GNW is that consciousness arises when local information (processed by specific sensory areas) is "amplified" and rapidly diffuses into an extensive frontoparietal network of long-range neurons – the "global neuronal workspace". On the other hand, consciousness is different from unconscious processing – much information is processed unconsciously at local levels. Only when you access the global workspace does it become conscious and available for thinking, planning, or reporting.

The GNW structure consists of a network of highly connected frontal and parietal neurons (with long, branched axons). Its dynamics involve conscious ignition ("all-or-none") – a content reaches consciousness when it induces a "neuronal ignition" – a burst of sustained activity across a vast network, visible in EEG, MEG, fMRI, as gamma oscillations, sustained activation, positive delays in ERP (P3), etc. This access is abrupt: stimuli just above threshold elicit a global response, others remain subthreshold.

The GNW model confirms several predictions: conscious access occurs dis- continuously, not progressively. WoRKspAce allows content access to multiple functionsmemory, language, planning, etc. The global space allows that, al- though many processes are "concurrent", only one (or a few) can access con- sciousness, avoiding overload.

The GLobAL NeURonAL WoRKspAce (GNW) is a computational model of how consciousness can emerge in the brain. Instead of thinking of the brain as a collection of separate zones, this model imagines it as a network of inter- connected neurons, where many areas process information locally and un- consciously. However, only some information can be "diffused" throughout the network, becoming conscious. The GNW can be simulated on a computer using neural networks of various architectures.

The model attempts to match real components of the brain that make con- scious experience possible ("NeURonAL CoRReLAtes of ConscIoUsness").

Higher brain functions (such as consciousness, decision-making, attention control, etc.) are not performed by a single area, but by the joint work of many parts, as in a swarm of bees or a school of fish ("swarm behavior"), where col- lective emergent phenomena are observed. This expression ("swarm behavior") refers to the fact that in GNW, neurons work locally, but some activation pat- terns can "contagious" and involve the entire network (spreading like a "spark"), leading to emergent phenomena: conscious decisions or central control.

In sHoRt. GNW presents a robust and testable framework to explain the emer- gence of consciousness from global brain dynamics. Evidence shows that con- sciousness depends on a sudden integration (ignition) in a specific network, allowing the diffusion and manipulation of content throughout the system – and that the loss of this network corresponds to the loss of consciousness itself.

The GNW is a model that uses neural networks (real or computer-simulated) to explain how consciousness works. It shows that the brain, like a group of ants or a swarm, only makes conscious information that can be shared and amplified by a vast network, thus creating a mathematical basis for phenomena such as consciousness, decision-making, and central mental control.

## V.   GNW Model With Hopfield Networks.

In this work, we explore an innovative approach to integrate the GNW model with HopfIeLD netwoRKs, a type of recurrent neural network known for its abil- ity to store and retrieve patterns stably. Hopfield networks, with their attractor- based dynamics, offer an interesting perspective for modeling the stability and resilience of the global workspace proposed by GNW. By combining these two concepts, we seek to investigate how attractor dynamics can contribute to the integration and maintenance of information in the global workspace, providing a deeper understanding of the neural mechanisms of consciousness.

The purpose of this work is, therefore, to present a hybrid model that inte- grates GNW with Hopfield networks, exploring how attractor dynamics can be applied to simulate conscious cognitive processes. Through simulations and theoretical analyses, we hope to contribute to advancing our understanding of consciousness by offering a new perspective on how the brain implements the global workspace and how the dynamic stability of neural networks can play a crucial role in this process.

The goal of a mathematical model for GNW is to create a system that cap- tures:
1. The existence of specialized LocAL MoDULes ;
2. The "gLobAL WoRKspAce", which integrates distributed signals;
3. Ignition dynamics (all-or-none transitions and critical thresholds);
4. The notion of broadcasting and global accessibility;
5. Synaptic plasticity, attention and learning phenomena.

The GNW model we will develop in this section uses binary Hopfield networks ([Hopfield 1982], [Hertz 1991]), which offer several advantages that make them an appropriate choice for modeling both LocAL MoDULes (sensory) and WoRKspAcein the GNWmodel.

In fact, Hopfield networks are inherently associative memories. This means they can recover a complete and correct pattern from an incomplete or noisy version of it. This is extremely useful for modeling the brain's ability to rec- ognize objects or situations, even when sensory information is imperfect. Be- cause of their ability to correct errors, Hopfield networks are robust to noise. This is important in sensory modules, where the incoming information may be noisy or ambiguous.

They can be trained using Hebb's rule ([Hebb 1949]), which is a form of unsupervised learning. This means that the network can learn to recognize patterns without the need for labeled data. This is useful in environments where labeled data is scarce or nonexistent. The plasticity in weights, coupled with the ability to store memories, ensures the stability of the system.

Recurrent connections in Hopfield networks allow neurons to interact with each other, creating complex dynamics that can be used to model complex cognitive processes.

Although simplified, Hopfield networks capture some important aspects of neural computation, such as synaptic plasticity and recurrent dynamics. This can make the model biologically more plausible.

Similarly, WoRKspAce needs to integrate information from multiple sources (sensory modules). Hopfield networks are capable of combining information from different sources and generating a coherent

global state. The WoRKspAce is responsible for decision-making. Hopfield networks can be used to model decision-making, where the network's attractors represent different options or actions. The workspace needs to be flexible and adaptable to different tasks and environments. Hopfield networks can be trained to learn new patterns and adapt their behavior based on experience.

Although simple, Hopfield networks can be extended and combined with other techniques to create more complex and powerful models. The structure of Hopfield networks (neurons, connections, weights) is relatively easy to inter- pret, which can help understand how the GNW model works and what cognitive processes it simulates.

In summary, Hopfield networks offer a good balance between simplicity, computational power, and biological plausibility, making them a suitable choice for modeling both the LocAL MoDULes and the WoRKspAce in the GNW model we will develop next. They provide a solid foundation for exploring important con- cepts such as associative memory, information integration, decision-making, and consciousness.

Let us now move on to the mathematical model announced above. Its goal is to model consciousness as an emergent phenomenon of distributed interaction.

From now on, this model will be referred to as the GNWMoDeL. It is based on the architecture

$$GNW = \bigcup_{m=1}^{M} M_m \cup W \qquad (1)$$

In this architecture, the LocAL MoDULes Mm are implemented as binary Hopfield networks 7, which represent specialized brain areas. The WoRKspAce W, also implemented as a Hopfield network, integrates information from local modules, interconnected by activation and feedback signals. This architecture combines the associative memory capacity of Hopfield networks with a global integration mechanism represented by the WoRKspAce.

This model will be refined in the following sections, increasing its complexity, to more realistically reproduce conscious activity. The sensory inputs for the Mm modules will be discussed in the 7 section.

We will detail the dynamic equations and how this structure influences the system's dynamics, including the concept of "ignition" and the modulation of attractors.

We will use a presentation style that facilitates computational modeling, ex- posing the theory, translating it into pseudocode that can later be programmed, for example, in Python.

## VI. Formal Architecture, Parameters, And Dynamic Variables.

The basic GNW model, which we propose, has the following formal architecture, parameters, and dynamic variables.
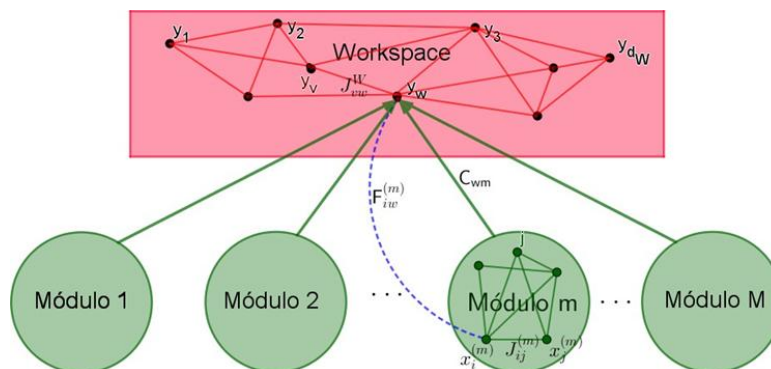


Figure 3: GNW model. Connections of the neurons of each module to the neurons of the Workspace, and their respective weights.

It has M local Hopfield modules $\{M_m\}_{m=1,2,\cdots,M}$, each with Nm neurons. Each Hopfield module represents a specialized brain area, which is composed of thousands or millions of neurons. Each module Mm has Pm memorized pat- terns (memories). The WoRKspAce W is also a Hopfield network, with NW neu- rons.

The matrix (Nm × Nm) of synaptic weights in the module Mm is

$$J^{(m)} = \left( J^{(m)}_{ij} \right)_{i,j=1,\cdots,N_m}$$

(symmetric and with zero diagonal). They are defined by Hebb's rule:

$$J_{ij}^{(m)} = \sum_{p=1}^{P_m} \xi_i^{(m,p)} \cdot \xi_j^{(m,p)} \tag{2}$$

where $\{\xi^{(m,p)}\}_{p=1,\cdots,P}$, are the $P_m$ patterns memorized in the module $M_m$.

The matrix ($N_W \times N_W$) of synaptic weights in the WORKSPACE is

$$J^W = (J^W)_{uv=1,\cdots,N_W}$$

(symmetric and with zero diagonal). They are also defined by Hebb's rule:

$$J_{uv}^W = \sum_{q=1}^{P_W} \eta_u^q \cdot \eta_v^q \tag{3}$$

where $\{\eta^q\}_{q=1,\cdots,P_W}$ are the patterns memorized in the Workspace.

The module connection matrix ($N_W \times M$) for the workspace is $C = (C_{um})$; $u = 1, \cdots, N_W$, and $m = 1, \cdots, M$. The feedback matrix ($N_m \times N_W$) from the workspace to module m is $F^{(m)} = (F_{iu}^{(m)})$, $i = 1, \cdots, N_m$, and $u = 1, \cdots, N_W$. The feedback signals from the workspace serve as inputs for modules.

$\mathbf{x}^{(m)}$ is the state vector of module m. The activation signals of the modules serve as inputs for WORKSPACE. $\mathbf{y} = \mathbf{y}^W$ is the Workspace state vector. Finally, $\vartheta$ is the activation threshold for Workspace ignition, $\alpha$ is the strength of lateral competition between modules, $\beta$ is the strength of Workspace feedback to modules, and $\sigma(x)$ is the sigmoid function.

The states of the Hopfield modules $\mathbf{x}^{(m)}(0)$, the WORKSPACE Hopfield state $\mathbf{y}(0)$, and the connection matrices C and feedback matrices F are randomly initialized.

## Dynamics and Propagation (matrix notation).

The state of module m is updated using the "update" equation:

$$\mathbf{x}^{(m)}(t + \Delta t) = \sigma \left[ J^{(m)} \mathbf{x}^{(m)}(t) + \beta F^{(m)} \mathbf{y}(t) + \boldsymbol{\varepsilon}^{(m)}(t) \right] \tag{4}$$

where $\sigma$ is the activation function Signal (for binary networks) an d, more generally, the function sigmoid, tanh, or others, for continuous networks, which we do not address in this article. See, however, [Ramsauer 2021]).

To define LATERAL COMPETITION between MODULES, we first calculate the AVERAGE ACTIVATION of EACH MODULE m:

$$A_m(t) = \frac{1}{N_m} \sum_{i=1}^{N_m} x_i^{(m)}(t) \tag{5}$$

Lateral competition is then defined by:

$$\text{Comp}_m(t) = A_m(t) - \frac{1}{M-1} \sum_{n=m} A_n(t) \tag{6}$$

that enters as "input" of the modules to the Workspace:

$$I(t) = C \cdot \text{Comp}(t); \text{ that is } I_u(t) = \sum_{m=1} C_{um}\text{Comp}_m(t) \tag{7}$$

Updating the Workspace state is done through the "update" equation:

$$\mathbf{y}(t + \Delta t) = \sigma(\mathbf{J}^W \mathbf{y}(t) + I(t)) \tag{8}$$

**Workspace Ignition.** Now let's mathematically formalize the Workspace Ignition (the critical moment of Workspace activation). First, we calculate their average activation:

$$a_W(t) = \frac{1}{N_W} \sum_{u=1}^{N_W} y_u(t) \tag{9}$$

Later, we apply an activation threshold for ignition:

$$\text{Ignition}(t) = \begin{cases} 1 & \text{if } a_W(t) > \vartheta \\ 0 & \text{otherwise} \end{cases} \tag{10}$$

## Explanatory Notes.

**1. Storing "conscious patterns" in the weights structure.** In the architecture proposed above, the weights internal to Workspace (matrix $\mathbf{J}^W$) are crucial. They determine which conscious patterns Workspace can store and recognize. By disregarding these weights, Workspace becomes a simple sum of the modules' activations, losing the capacity for associative memory and modeling of complex patterns.

To define $\mathbf{J}^W$, Hebb's rule is used, as in local modules, with specific patterns that represent stored memories induced by sensory inputs (section 7).

Patterns capable of accessing the Workspace (i.e., becoming conscious) are encoded in the memories of the Workspace Hopfield network, becoming available for global propagation ("broadcast"). "Conscious access" to the Workspace can be seen as the activation of one of these "global patterns."

**2. Lateral Competition Between Modules.** Lateral competition is a mechanism in the GNW model that simulates competition between different brain areas (represented by Hopfield modules $M_m$) for access to the global WorkspaceThe

idea is that only a subset of brain areas can be "conscious," or have access to the WORKSPACE at a given time.

In the previous formulas (see (5) and (6)), lateral competition is implemented by calculating the average activation of each module, $A_m(t)$, and subtracting a fraction ($\alpha$) from the sum of all the activations of the other modules. This means that if the sum of the activations of the other modules $n$ ($n' = m$) is high, the competition for module m will be lower (more negative), inhibiting its ability to access WORKSPACE; that is, modules with relatively high activation inhibit modules with lower activation.

A high value for $\alpha$ means stronger competition, while a low value means weaker competition.

In essence, this formulation implements a form of lateral inhibition, where more active modules tend to suppress less active modules, simulating competition for cognitive resources.

### 3. Module Connections to the Workspace.
In the model we implemented, modules are indirectly linked to the WORKSPACE through the inputs $I_u(t) = \sum_{m=1}^{M} C_{um}Comp_m(t)$, defined in 7. The competition between modules influences the amount of activation each module can "send" to the Workspace.

$$q_W(t + \Delta t) = \frac{1}{N_W} \sum_{u=1}^{\infty} w\, y_u(t + \Delta t) = \frac{1}{N_W} \sum_{u=1}^{N} w\, \sigma\left( J_{uv}^W y_v(t) + \sum_{m=1}^{\infty} C_{um}Comp_m(t)\right)$$

This means that the activation of each module, adjusted for lateral competition, contributes to the total activation of WORKSPACE.

Furthermore, in the model, there is feedback from WORKSPACE to the modules, which is implemented in the Hopfield module update:

$$\mathbf{x}^{(m)}(t + \Delta t) = \sigma\left( J^{(m)}\mathbf{x}^{(m)}(t) + \beta\, F^{(m)}\mathbf{y}(t) + \xi^{(m)}(t)\right)$$

Here, the term $\beta\, F^{(m)}\mathbf{y}(t)$ represents the feedback from WORKSPACE to the module $M_m$. The parameter $\beta$ controls the strength of this feedback. When WORKSPACE is active ($q_W$ high), it influences the local fields of the modules, which, in turn, affect the activation of the modules in the next iteration.

### 4. Ignition and Feedback.
In the proposed architecture, ignition can be seen as the transition of WORKSPACE to an active state that can trigger changes in local modules, and, in particular, alter the memories (patterns) stored in the Hopfield networks of these modules.

Indeed, workspace ignition is not limited to being a passive event; it triggers a feedback signal that is sent to the local modules. This feedback modulates the activity of the neurons in the modules, influencing their states and, consequently, their ability to influence WORKSPACE in future iterations.

More specifically, feedback from WORKSPACE alters the local fields of neurons in the local modules. The local field $h_i^{(m)}$ of a neuron $i$ in module m is the weighted sum of the inputs it receives from other neurons in that module, and also from any external input or feedback:

$$h_i^{(m)} = \sum_{j=1}^{N_m} J_{ij}^{(m)} x_j + \beta \sum_{u=1}^{N_W} F_{iu}^{(m)} y_u^W$$

The feedback, therefore, adds an additional term to this local field, which could cause the state of neuron $i$ to change.

If the feedback strength ($\beta$) is high enough, module attractors can be significantly altered or even eliminated, allowing the network to explore different local states (memories), which can be seen as a form of "Attention" or "cognitive Reconfiguration", where WORKSPACE directs the activity of modules to explore new areas of the Hopfield network's state space.

The most interesting effect is how workspace feedback interacts with synaptic plasticity (the ability of connection weights to change over time. See the 6 section). By applying Hebb's rule (2), feedback can cause connection weights to adjust to stabilize or amplify the resulting state. In other words, feedback can "imprint" a new state in the network's memory.

The incorporation of stochastic dynamics (which allows the exploration of different states) with feedback from WORKSPACE and synaptic plasticity allows the system to learn and adapt to the environment in a flexible and robust manner.

**5. Ignition can be viewed as a nonlinear phase transition.** A "local" pattern refers to a specific activation state, $\mathbf{x}^{(m)} \in \{-1, +1\}^{N_m}$, of a sensory module, $M_m$, of the GNW network. This pattern can be triggered by a sensory stimulus (e.g., the presentation of an image, a sound, or a tactile sensation). If this sensory stimulus is strong enough, the corresponding pattern in the sensory module will become more relevant, competing with the activity of other modules to access WORKSPACE.

Lateral competition ensures that only the most relevant and informative patterns can influence WORKSPACE. This is represented by the average activation, $A_m(t)$, of the module $M_m$, and by the lateral competition $\text{Comp}_m(t)$, defined, respectively, by equations (5) and (6). A high value of $\text{Comp}_m(t)$ indicates that module m is "winning" the competition, as explained earlier.

The propagation of a local pattern to WORKSPACE depends on the strength of the connections between the local module and WORKSPACE, represented by the connection matrix $C = (C_{um})_{u=1,\cdots,N_W; m=1,\cdots,M}$ (see Figure 3).

If the weights $C_{um}: u = 1, \cdots, N_W$, corresponding to a given module m, are high, then the activity of that module will have a greater impact on WORKSPACE.

This is represented by the input to the workspace (see equation (7)):

$$I_u(t) = \sum_{m=1}^{M} C_{um} Comp_m(t)$$

For a local pattern to be "aware" or integrated into WORKSPACE, it needs to match a memory pattern $\eta^q$, already stored in WORKSPACE. This match can be measured by the "distance" between the input to the workspace $I(t)$ and the Hopfield network attractors in the workspace. If the distance is small, the Hopfield network in the workspace will converge to a nearby attractor, representing the "awareness" or integration of the pattern.

Workspace ignition is modeled as an abrupt ("all-or-none") transition to an attractor in the global state space. This occurs when the average workspace activation $a_w$ exceeds a threshold $\vartheta$:

$$\text{If } a_w > \vartheta, \text{ then the pattern becomes "conscious."}$$

This behavior represents a nonlinear phase transition, where below the threshold $\vartheta$ the pattern remains "latent" (not conscious) and above the threshold, it spreads throughout the network through long-range connectivity.

The critical temperature can be related to the difficulty with which this transition occurs. Intuitively, lateral competition selects the most relevant input pattern, the connection matrix amplifies this signal to the workspace, and if this signal is strong enough, the work is ignited.

**Conclusion.** The GNW can be formalized as an extended Hopfield network, with local modules, a work, and specialized connectivities. The activation of the modules, influenced by lateral competition, contributes to the activation of the work. "Consciousness" emerges as the synchronized activation ("ignition") of a global pattern/attractor of the work – a phase transition in the state space. The work sends feedback to the modules, influencing their local fields and, consequently, their activations.

# 6 Plasticity in the GNW Base Model

In the context of the GNWbase model, from the 5 section, plasticity refers to the ability of connections and model parameters to adapt and change over time in response to experience or learning. In biological terms, synaptic plasticity is the basis of learning and memory.

Incorporating plasticity is an important improvement to the GNWmodel, which can make it more adaptive, robust, and capable of modeling complex cognitive processes. Plasticity can be incorporated at several levels:

**Synaptic Plasticity in Hopfield Modules.** Adjust the weight matrices ($J^{(m)}$) of Hopfield modules to memorize new patterns or strengthen existing ones. Adjust the weight matrix ($J^W$) of WORKSPACE to learn new global patterns or refine existing ones. Use Hebbian learning rules (e.g., Oja's rule) to update the weights based on neuron activity.

**Plasticity in Module-Workspace Connections.** Adjust the connection matrices (C and F) to optimize the flow of information between modules and WORKSPACE. Allow connections to strengthen or weaken based on the relevance of the information.

**Plasticity in Self-Attention Modules (if used).** Adjust the parameters of the self-attention modules to improve feature extraction and information weighting. To be covered in a later article.

We will now describe how to incorporate synaptic plasticity into the base GNW model, allowing the network to learn and adapt dynamically.

The parameters, dynamic variables, connectivities, and initialization are the same as those in the base model in the 5 section. Let's move on to the dynamics of the model with synaptic plasticity. For each $t$, we update the state of module m:

$$x^{(m)}(t + \Delta t) = \sigma \left( \sum_i J^{(m)}(t)x^{(m)}(t) + \sum_u F^{(m)}(t)y_u(t) \right)$$

and we apply synaptic plasticity to update the weight matrix $J^{(m)}(t)$, using Hebb's rule:

$$J^{(m)}(t + \Delta t) = J^{(m)}(t) + \eta \left( \mathbf{x}^{(m)}(t + \Delta t)\mathbf{x}^{(m)}(t + \Delta t)^\top - J^{(m)}(t) \right) \tag{11}$$

The Workspace state is updated similarly.

$$y_u(t + \Delta t) = \sigma \left( J^W_{uv}(t)y_v(t) + \sum_{m=1} C_{um}(t)L_m(t) \right) \tag{12}$$

We then apply synaptic plasticity to update the matrix of weights $J^W(t)$ using Hebb's rule:

$$J^W(t + \Delta t) = J^W(t) + \eta \left( \mathbf{y}(t + \Delta t)\mathbf{y}(t + \Delta t)^\top - J^W(t) \right) \tag{13}$$

Here, $\eta$ is the learning rate.

To apply plasticity to Module Connections⇌Workspace, we update the matrices C and $F$ using a rule based on the correlation between module and WORKSPACE activity:

$$C_{um}(t + \Delta t) = C_{um}(t) + \eta_w \left( y_u(t + \Delta t) \cdot A_m(t) - C_{um}(t) \right) \tag{14}$$

$$F_{iu}^{(m)}(t + \Delta t) = F_{iu}^{(m)}(t) + \eta_t \ x_i^{(m)}(t + \Delta t) \cdot y_u(t) - F_{iu}^{(m)}(t) \tag{15}$$

where $\eta_w$ and $\eta_f$ are the learning rates for the connections. Lateral Competition between Modules and Workspace Ignition are calculated as before.

**Important note.** The reason for using $\mathbf{y}(t + \Delta t)$ instead of $\mathbf{y}(t)$ in the right-hand side of the Hebb update rule equation, equation (16), is to ensure that the weight update is based on the resulting system activity after applying the input and network dynamics. For the other plasticity updates, the argument is the same.

Hebb's basic rule states that "*neurons that ftre together, they wire.*" This means that if two neurons $i$ and $j$ are active simultaneously, the strength of the connection between them ($J_{ij}$) should be increased. Synaptic plasticity refers to the ability of synapses (connections between neurons) to change their strength over time, in response to neuronal activity. This is fundamental for learning and memory.

The update equations implement Hebb's rule with synaptic plasticity. To be more concrete, let's take as an example:

$$\Delta J^W(t) = \eta \ \mathbf{y}(t + \Delta t)\mathbf{y}(t + \Delta t)^\top - J^W(t) \tag{16}$$

where $\Delta J^W(t)$ is the change in the weight matrix $J^W$ at time $t$; $\mathbf{y}(t+\Delta t)\mathbf{y}(t+\Delta t)^\top$ is the correlation between the neurons at time $t+dt$, representing the resulting *network activity*; $J^W(t)$ is the weight matrix at the current time $t$ (subtracting $J^W(t)$ implements a "forgetting" or "regularization" mechanism); $\eta$ is the learning rate).

The main reason for using $\mathbf{y}(t + \Delta t)$ is to capture the activity caused by the input and the network dynamics. Hebb's rule adjusts the network's memory ($J^W(t)$) so that it better remembers this new state ($\mathbf{y}(t + \Delta t)$). If we used $\mathbf{y}(t)$, we would be adjusting the network's memory based on its previous state, which would not be as effective for learning.

In this way, the network learns to associate the inputs with the states resulting from the network dynamics. The "forgetting" term helps stabilize the network, and the equation models the adaptability of the connections over time. In conclusion, the use of $\mathbf{y}(t + \Delta t)$ in Hebb's rule with synaptic plasticity ensures that weight updates are based on the resulting network activity, leading to more effective and stable learning.

# 7 Sensory Inputs and Feedforward Networks in the GNW Model with Hopfield Networks

The local Hopfield modules, $M_m$, are sensory: one for vision, one for hearing, one for touch, etc. For these modules, in the GNW model, to work effectively with

sensory inputs (images, sounds, sensations, etc.), adequate preprocessing is required to extract relevant features from these inputs and encode them in a format compatible with Hopfield networks.

In fact, Hopfield modules alone are not designed to handle raw sensory data directly. They work best with binary representations, ($-1$ and $+1$), that encode specific patterns or features. Therefore, it is crucial to have (multilayer) feed-forward networks, or other preprocessing mechanisms that extract meaningful features from sensory inputs and transform them into suitable representations for Hopfield networks.

It is also convenient to use convolutional networks (CNNs) and generative networks (GANs and VAEs) in the preprocessing of sensory inputs. CNNs are especially effective in extracting features from data such as images or sounds. They use convolutional filters (whose weights are calculated by progressive learning), which detect local and hierarchical patterns. If the inputs are visual (e.g., images), CNNs can be used to identify edges, textures, objects, and other relevant visual features. For auditory inputs (e.g., sounds), CNNs can learn to extract spectral features, such as frequencies and time-frequency patterns.

Generative networks, such as GANs or VAEs ("Variational Autoencoders"), can learn latent representations of sensory inputs. These "latent representations" capture the essential features of the data in a lower-dimensional space.

Generative networks can be used to generate new examples of data that are similar to the original inputs. This can be useful for augmenting the training dataset or exploring different variations of the inputs. Generative networks can be used to simulate "imagination", generating internal representations of possible sensory inputs that the model has not yet experienced.

It is easy to extend the base GNW model, from the 5 section, by incorporating feeDfoRWARD netwoRKs (multilayer), with convolutional layers, and even generative networks in the local Hopfield modules of the GNW model (Global Neuronal Workspace), allowing the model to process and interpret sensory data in a more advanced way.

The subsequent steps (updating the Hopfield Modules, Lateral Competition, Workspace, and Feedback) follow the same dynamics described in the previous sections, with the Hopfield modules receiving the inputs processed by CNNs and generative networks. CNNs and generative networks can be trained separately (pre-training) or together with Hopfield networks, using labeled or unlabeled training data.

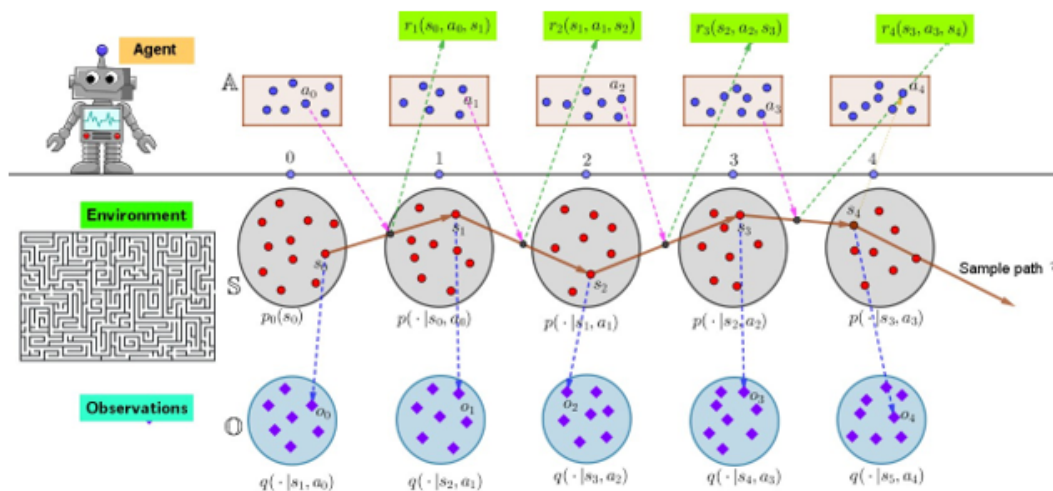# 8 Role of Reinforcement Learning (RL) in the GNWModel



**Figure 4:** Reinforcement Learning

In the base GNWmodel, from the 5 section, Reinforcement Learning (RL) [Sutton 2018] acts as a mechanism to optimize system behavior in a given environment or task. Instead of manually fixing the model's connections and parameters, the RL agent is allowed to learn to adjust them to maximize total reward.

The Reinforcement Learning Components in the GNW Model are as follows:

Agent. The agent is an algorithm that interacts with the environment (the GNW model) and learns to make decisions to maximize reward. In this case, the agent can be a neural network (e.g., a deep Q-learning model) that receives the state of the environment (the GNW state) as input and produces an action as output. The agent's goal is to learn a policy (a function that associates an action with each state) that maximizes the expected reward over time.

**Environment (World).** The environment is the GNW model itself: GNW = $\cup_{m=1}^{M} M_m \cup W$. It consists of the local Hopfield modules, the Workspace Hopfield, and the connections between them. The environment receives actions from the agent and evolves according to the dynamic equations of the GNWmodel.

**World States.** The world state is a representation of the current state of the environment. In the GNWmodel, the state of the world can be defined as the combination of the states of the $M$ local Hopfield and Workspace Hopfield mod-

ules:

$$\mathbf{s}(t) = \mathbf{x}^{(1)}(t), \mathbf{x}^{(2)}(t), \cdots, \mathbf{x}^{(M)}(t), \mathbf{y}^{W}(t) \tag{17}$$

This state can be represented as a concatenated vector of the states of all nodes (neurons) in the Hopfield networks.

**Actions Available to the Agent.** Actions are the decisions the agent can take to influence the environment. In the GNW model, actions can include: (i). AD-justing the Strength of Lateral Competition ($\alpha$). The agent can increase or decrease the strength of lateral competition between modules; (ii). Modifying Connections Modules/Workspace (C, F). The agent can adjust the connection matrices that control the flow of information between local modules and Workspace; (iii). Controlling the Activation Threshold ($\vartheta$). The agent can adjust the activation threshold that determines when Workspace "ignites," when a given piece of information becomes conscious.

Actions can be discrete (e.g., increase, decrease, maintain) or continuous (e.g., a numerical value between 0 and 1).

**Rewards.** The immediate reward is a signal, $r(\mathbf{s}, a)$, that evaluates the quality of the world transition when, in state $\mathbf{s}$, it receives an action $a$ from the agent. The reward is defined based on the goal we intend the GNW model to achieve. Examples of rewards in the GNW model are:

(i). Workspace Activation. Positive reward if Workspace is active (e.g., average activation exceeds a certain threshold);

(ii). Pattern Recognition. Positive reward if the model recognizes a specific pattern (e.g., a pattern memorized in one of the modules or in Workspace);

(iii). Classification. Positive reward if the model correctly classifies an input;

(iv). Stability. Negative reward if the system becomes unstable (e.g., the module states oscillate too rapidly).

By using RL, the GNW model can learn to:

(i). Allocating Cognitive Resources. Adjust lateral competition to prioritize relevant modules and suppress irrelevant modules;

(ii). Integrating Information. Adjust the Modules/Workspace connections, to optimize information flow and integration in Workspace;

(iii). Controlling Awareness. Adjust the activation threshold to control when Workspace becomes aware of a given information;

(iv). ADApting to DynAMic EnviRonMents. Adjust your parameters in real-time to adapt to changes in the environment or task.

We will now describe pseudocode for a Q-Learning episode, applied to the base GNW model, to optimize a policy.

In addition to those already defined in the **??** section, we have the following new additional hyperparameters: $\epsilon$ = exploration probability (exploration rate) and $N_{episodes}$ = number of training episodes. We also have the function Q, defined on the set {states} × {actions}, whose value $Q(\mathbf{s}, a)$ estimates the expected reward for performing action $a$ in state $\mathbf{s}$.

After randomly initializing the parameters and dynamic variables:

$$\mathbf{x}^{(m)}(0),\ \mathbf{y}^{W}(0),\ J^{(m)},\ J^{W}(0),\ C(0),\ F^{(m)}(0)$$

we randomly initialize the function $Q(\mathbf{s}, a)$, for all possible states and actions.

The TRAining Loop (Q-Learning) is as follows. For each episode $e$ from 1 to $N_{episodes}$:

- "Reset" the environment (module states and WoRKspAce)

- Initialize the environment state $\mathbf{s}(0)$.

- For each time step $t$, from 1 to $T$:

- with probability $\epsilon$, select a random action $a(t)$ (farm);

- with probability $1 - \epsilon$, select the action that maximizes the function Q:

- $a(t) = \arg\max_a Q(\mathbf{s}(t), a)$

  GNW DynAMics:

- Adjusting the strength of lateral competition $\alpha$ based on the action $a(t)$.

- Updating Hopfield Modules. For each module m:

$$\mathbf{x}^{(m)}(t + \Delta t) = \sigma\ J^{(m)}\mathbf{x}^{(m)}(t) + \beta\, F^{(m)}\mathbf{y}(t) + \xi^{(m)}(t) \qquad (18)$$

- Entry to Workspace:

$$I_u(t) = \sum_m C_{um}\, A_m(t) - \alpha \sum_{n=m} A_n(t) \qquad (19)$$

- Workspace Update:
$$\mathbf{y}(t + \Delta t) = \sigma(\mathbf{J}^W \mathbf{y}(t) + \mathbf{I}(t)) \tag{20}$$

- Workspace Ignition:
$$\text{Ignition}(t) = \begin{cases} 1 & \text{if } a_W(t) > \vartheta \\ 0 & \text{otherwise} \end{cases} \tag{21}$$

where $a_W(t) = \dfrac{1}{N_W} \sum_{u=1}^{N_W} y_u^W(t)$ is the average workspace activation.

Reward And Next Status:

- Calculate the reward $r(t)$ based on the Workspace state (e.g., average activation).

- Observe the next state $\mathbf{s}(t + \Delta t)$, based on $\mathbf{x}^{(m)}(t + \Delta t)$ and $\mathbf{y}(t + \Delta t)$.

  Update QFunction:

- Update the Qfunction using the Bellman equation:
$$Q(\mathbf{s}(t), a(t)) \leftarrow Q(\mathbf{s}(t), a(t)) + \eta \left[ r(t) + \gamma \max_a Q(\mathbf{s}(t + dt), a) - Q(\mathbf{s}(t), a(t)) \right] \tag{22}$$

where $\eta$ is the learning rate.

- Output: Optimized Q function $Q(\mathbf{s}, a)$.

In summary: For each episode, the environment is "reset," and for each time step within the episode:

- The agent selects an action using a $\epsilon$-greedy policy.

- The GNW dynamics are updated using the model equations.

- The reward is calculated, and the next state is observed.

- The Q function is updated using the Bellman equation.

# 9 Conclusion

In this work, we present a hybrid model of the Global Neuronal Workspace (GNW) that combines Hopfield networks with synaptic plasticity mechanisms to simulate complex cognitive processes. The proposed model seeks to integrate information from various brain areas (represented by local Hopfield modules) into a global workspace (global Workspace), enabling conscious decision-making and adaptation to different environments and tasks.

The simulations demonstrated that:

- The dynamics of the Hopfield network in local modules allows the recog- nition of sensory patterns and the retrieval of associative memories.
- Synaptic plasticity, implemented through Hebb's rule, allows the model to learn and adapt to new inputs and tasks. Feedforward networks and CNNs allow the preprocessing of complex sensory inputs, extracting features relevant for recognition and classification. Generative networks (VAEs) offer the ability to learn latent representations of data, allowing the generation of new inputs (imagination) and a better understanding of the data structure. Reinforcement Learning (RL) allows you to optimize model behavior by adjusting parameters and connections to maximize re- ward in specific tasks.

This model offers an interesting perspective on the architecture and dy- namics of the GNW, suggesting that the combination of associative memory mechanisms (Hopfield networks) with global integration processes and synap- tic plasticity, hierarchical processing of sensory inputs (feedforward networks and CNNs), latent representation and data generation (generative networks), and behavior optimization through RL, may be fundamental for the emergence of consciousness and other higher cognitive functions.

While this model captures some important aspects of the GNW, it also has limitations. In particular, the simulation of natural language is still simpli- fied and the model does not take into account the complexity of interactions between different brain areas.

In future work, we intend to expand the model to include:
- More detailed modeling of natural language using natural language pro- cessing techniques.
- Integration of multimodal information (vision, audition, touch, etc.).
- Incorporation of more sophisticated control and attention mechanisms.
- Validation of the model with experimental data from neuroimaging and cognitive psychology.

In conclusion, this hybrid GNW model represents an interesting step to- ward a more complete understanding of the neural mechanisms of conscious- ness and cognition. By combining different theoretical and computational ap- proaches, we hope this work can inspire further research and lead to significant advances in our understanding of the human brain.

## References

[1]. [Dehaene & Changeux 2006] Dehaene, S., Changeux, J.P., Naccache, L., Sackur, J., Sergent, C. (2006). Conscious, Preconscious, And Subliminal Processing: A Testable Taxonomy. Trends In Cognitive Sciences, 10(5), 204- 211.

[2]. [Dehaene & Changeux 1998] Dehaene, S., Kerszberg, M., Changeux, J.P. (1998). A Neuronal Model Of A Global Workspace In Effortful Cognitive Tasks. Proceedings Of The National Academy Of Sciences, 95(24), 14529-14534.

[3]. [Coolen 2005] Coolen, A. C. C., Ku¨ Hn, R., & Sollich, P., ”Theory Of Neuronal Information Processing Systems”. Oxford University Press 2005.

[4]. [Hopfield 1982] Hopfield, J.J. (1982). Neural Networks And Physical Systems With Emergent Collective Computational Abilities. Proceedings Of The Na- Tional Academy Of Sciences, 79(8), 2554-2558.

[5]. [Hertz 1991] , J., Krogh, A., Palmer, R.G. (1991). Introduction To The Theory Of Neural Computation. Addison-Wesley Publishing Company.

[6]. [Huang 2022] Haiping Huang ”Statistical Mechanics Of Neuronal Networks”. Springer 2022.

[7]. [Oja 1982] Oja, E. (1982). Simplified Neuron Model As A Principal Component Analyzer. Journal Of Mathematical Biology, 15(3), 267-273.

[8]. [Hebb 1949] Hebb, D.O. (1949). The Organization Of Behavior. A Neuropsycho- Logical Theory. New York: Wiley Sons.

[9]. [Baars 1997] Bernard J. Baars, ”In The Theater Of Consciousness: The Workspace Of The Mind”, OUP USA, ISBN-10: 0195102657.

[10]. [Baars 1988] Baars, B. J. (1988). A Cognitive Theory Of Consciousness. Cam- Bridge University Press.

[11]. [Ramsauer 2021] Hubert Ramsauer And Others: ”Hopfield Networks Is All You Need”, Arxiv:2008.02217.

[12]. [Sutton 2018] Sutton, R.S., Barto, A.G. (2018). Reinforcement Learning: An Introduction. MIT Press.