# A Review of Bio-Inspired Neural Arbitration: From Attractor Dynamics to Human-Centric Control

## Section 1: Foundations of Decision-Making in Recurrent Attractor Networks

The challenge of arbitrating human intent, particularly within the demanding context of Human-Machine Interaction (HMI) and Brain-Computer Interfaces (BCI), necessitates a move beyond conventional classification algorithms toward models that capture the temporal, deliberative, and often ambiguous nature of human decision-making.[1] A promising avenue lies in the adoption of bio-inspired computational models grounded in the principles of cortical dynamics. At the forefront of this approach are recurrent attractor neural networks, which frame decision-making not as a static pattern recognition problem, but as a dynamic process of evidence integration and competition. This section delves into the foundational principles of these networks, beginning with a detailed analysis of the canonical Wong & Wang (2006) model, which serves as the direct mathematical and conceptual basis for the proposed neural arbitrator. It then expands to the broader principles of Winner-Take-All (WTA) competition and the concept of attractor landscapes, which together provide a robust and biologically plausible framework for understanding how neural circuits can transform noisy, continuous inputs into stable, categorical choices.

### 1.1 The Wong & Wang (2006) Model: A Canonical Circuit for Perceptual Choice

The model proposed by Kong-Fatt Wong and Xiao-Jing Wang in their highly influential 2006 paper, "A Recurrent Network Mechanism of Time Integration in Perceptual Decisions," stands as a landmark in computational neuroscience.[2] It provides a biophysically plausible circuit

mechanism that explains how the brain can integrate noisy sensory evidence over behaviorally relevant timescales to arrive at a categorical decision. The model was specifically developed to account for key neurophysiological observations from behaving non-human primates performing two-alternative forced-choice (2AFC) tasks, such as the random-dot motion discrimination task.[2] In these experiments, neurons in cortical areas like the lateral intraparietal (LIP) cortex exhibit "ramping" activity; their firing rates gradually increase over hundreds of milliseconds, with the slope of this ramp being proportional to the strength of the sensory evidence (i.e., the coherence of motion).[2] The Wong & Wang model successfully replicates this core phenomenon, providing a bridge between neural dynamics and cognitive processes.

The architecture of the model is fundamentally that of a recurrent neural network characterized by strong competition. It consists of two selective excitatory neural populations (e.g., one tuned to "Leftward" motion and another to "Rightward" motion) that interact through a shared population of inhibitory interneurons.[5] This arrangement creates a dynamic of mutual inhibition. Crucially, each excitatory population also possesses strong recurrent self-excitation, parameterized by a weight w+, which allows for the formation of reverberatory, self-sustaining activity.[2] When presented with ambiguous sensory input, these two populations engage in a "Winner-Take-All" competition: the population receiving slightly more evidence begins to increase its activity, which in turn further excites itself and simultaneously suppresses the competing population via the inhibitory pool. This positive feedback loop continues until one population's activity dominates completely, signifying that a decision has been reached.[2]

A central and critical finding of the Wong & Wang study is the role of slow synaptic dynamics in achieving the long integration times observed in biological decision-making. The authors demonstrated that a network model relying solely on fast synaptic transmission, mediated by AMPA receptors, produces decision times that are far too short and inconsistent with experimental observations.[2] Instead, robust and slow time integration is achieved when the recurrent excitatory connections are primarily mediated by NMDA receptors. These receptors have a much longer time constant, on the order of 100 milliseconds, which endows the network with the necessary "sluggishness" to accumulate evidence over hundreds of milliseconds without reacting prematurely to transient noise.[2] This specific biological mechanism provides a strong justification for the inclusion of a time constant, $\tau$, in the rate-based differential equations used to model such systems, as proposed in the user's project.[1] The slow integration is not merely a filtering effect but a fundamental property of the circuit's dynamics, particularly in the vicinity of an unstable saddle point in the system's state space that separates the basins of attraction for the alternative choices.[2]

Perhaps the most significant contribution of the Wong & Wang model for practical implementation is its principled reduction from a complex, biophysically detailed network of thousands of spiking neurons into a simplified two-variable dynamical system.[2] This was

accomplished using a mean-field approach, a powerful technique from statistical physics where the average activity, or firing rate, of a large, homogeneous population of neurons can be represented by a single continuous variable.[2] Through a series of approximations, the dynamics of the entire complex network were captured by a pair of coupled nonlinear differential equations governing the synaptic gating variables, $S_1$ and $S_2$, for the two competing selective populations.[2] This reduction makes the model computationally tractable for simulation and amenable to dynamical systems analysis, providing a direct and well-validated precedent for implementing a "Neural Arbitrator" in a standard programming environment like Python.[1] The model's enduring impact, evidenced by its extensive citation record, establishes it not as an arbitrary choice of equations, but as a canonical and empirically grounded framework for studying the neural basis of perceptual choice.[4]

## 1.2 Winner-Take-All Dynamics: The Engine of Competition

The competitive mechanism at the heart of the Wong & Wang model belongs to a broader class of circuits known as Winner-Take-All (WTA) networks. WTA computation is considered a canonical motif in theoretical neuroscience, a fundamental building block of cortical microcircuits responsible for a wide range of cognitive functions including selective attention, pattern recognition, and categorical decision-making.[13] The core function of a WTA network is to perform selection: from a multitude of simultaneous inputs, it identifies the most salient or strongest one and amplifies its representation while actively suppressing all others.[13] This process effectively transforms a distributed, analog representation of sensory evidence into a sparse, discrete representation of a choice or category.

While the user's proposed model is rate-based, examining the implementation of WTA dynamics in more biologically detailed spiking neural networks provides deeper insight into the underlying principles.[1] In spiking networks, WTA functionality is typically realized through the interaction of distinct populations of excitatory and inhibitory neurons.[15] A key architectural principle for achieving stable WTA behavior is the concept of "surround inhibition." In this configuration, inhibitory interneurons project broadly, suppressing excitatory neurons that are functionally or spatially distant while having weaker or no inhibitory effect on their immediate neighbors.[16] This allows a localized group of co-active excitatory neurons to form a stable "bump" of activity, which is self-sustaining due to local recurrent excitation and protected from competing inputs by the far-reaching inhibition it generates.[15] In the simplified rate-based model, this complex spatial connectivity is abstracted into a global mutual inhibition term, $w_{inhibit}$, which captures the essential competitive interaction.[1]

From a computational perspective, the WTA mechanism provides a powerful solution to the

problem of robust decision-making in the presence of noise. The competition is inherently nonlinear; a decision is not made based on a simple linear comparison of input strengths. Instead, a decision is "latched" only after one neural population has accumulated enough evidence to dynamically and decisively suppress its competitors.[1] This process creates an implicit decision threshold. Transient noise or fleeting evidence in favor of one option is insufficient to overcome the recurrent inhibition from the currently leading population. A sustained, differential input is required to push the system's state across a separatrix and into a winning configuration. This inherent robustness is a direct consequence of the circuit's competitive dynamics, a far more sophisticated mechanism than simple linear filtering. Furthermore, some theoretical work has proposed that WTA circuits may form the neural basis for probabilistic inference, where the firing rate of the winning neuron could represent the marginal probability of a particular hypothesis or state, given the sensory evidence.[18] This perspective aligns with the goal of arbitrating ambiguous human intent, where the system must infer the most likely intention from a noisy signal stream.

## 1.3 Attractor Landscapes: Encoding Decisions as Stable States

Viewing the WTA network through the lens of dynamical systems theory offers a powerful conceptual framework for understanding its function. In this view, a neural circuit is modeled as a dynamical system whose "state" at any moment is defined by the vector of firing rates of all its neurons.[19] The network's recurrent connections define a set of rules that govern how this state evolves over time. A decision is not an instantaneous output but rather the process of the system's state converging toward a stable configuration, known as an "attractor".[23]

The behavior of these systems can be intuitively visualized by imagining the system's state as a ball rolling on a multi-dimensional "energy landscape".[23] The valleys or local minima in this landscape correspond to the system's stable attractor states. For a 2AFC task, the energy landscape is configured to have two such attractor states, one corresponding to "Decision A" and the other to "Decision B".[20] In these attractor states, the neural population selective for the chosen option maintains a high, persistent firing rate, while the other population is suppressed. The user's concept of a "latched" decision is a direct reference to the system's state falling into one of these energy minima and remaining there, even if the driving stimulus is removed.[1] This property of persistent activity is a hallmark of attractor networks and is thought to be the neural substrate for working memory.[23]

The landscape is not composed solely of attractors. The set of all initial states (i.e., initial patterns of neural activity) that eventually converge to a particular attractor is defined as that attractor's "basin of attraction".[19] The boundaries between these basins are defined by unstable states (saddle points or ridges in the energy landscape). The relative size of a basin

of attraction corresponds to the a priori bias towards the associated decision.[19] A larger basin means that a wider range of initial conditions will lead to that outcome, effectively biasing the system's choice. This provides a sophisticated mechanism for modeling how factors like expectation or prior knowledge can shape a decision, potentially by modulating the network parameters to reshape the basins of attraction without altering the attractor states themselves.[19]

Finally, the stochastic nature of neural firing introduces a critical element of noise into the system's dynamics.[23] This is not merely a biological imperfection but a computationally vital component. Noise allows the system to escape from unstable equilibrium points, such as the initial, undecided state where both populations have low, symmetric activity. More importantly, it renders the decision-making process probabilistic.[23] When the sensory evidence is weak or ambiguous, the system's state will be close to a boundary between basins. In this situation, random fluctuations from neural noise can be sufficient to push the state into one basin on one trial and into the other basin on another trial. This naturally accounts for the inherent variability observed in behavioral responses to identical stimuli and forms the core principle that allows such a network to model internally generated, spontaneous choices, a key goal of the proposed project.[1]

# Section 2: A Comparative Analysis of Decision-Arbitration Frameworks

The bio-inspired neural arbitrator, grounded in the attractor network dynamics described in the previous section, represents a specific theoretical and architectural stance on how to model decision-making. To fully appreciate its unique contributions and justify its "human-centric" label, it is essential to place it within the broader landscape of computational modeling. This section provides a comparative analysis, situating the proposed model against two critical benchmarks. First, it is compared to Evidence Accumulation Models (EAMs) from cognitive science, which provide a parallel, higher-level abstraction of the same cognitive process. This comparison highlights the deep conceptual alignment between the neural implementation and established cognitive theory. Second, it is contrasted with the dominant engineering alternative for processing noisy, sequential data: standard Recurrent Neural Networks (RNNs) like LSTMs. This contrast clarifies the fundamental trade-offs between biological plausibility and intrinsic interpretability versus raw predictive power and architectural generality.

## 2.1 Evidence Accumulation Models (EAMs): The Cognitive Science Parallel

Evidence Accumulation Models (EAMs) are a highly successful class of computational models in cognitive psychology and neuroscience that provide a formal mathematical description of the latent cognitive processes underlying simple decisions.[28] Rather than focusing on the implementation details of a neural circuit, EAMs abstract the decision process as the noisy accumulation of evidence over time toward a response threshold.[29] The most prominent example for two-choice tasks is the Drift-Diffusion Model (DDM), which models the decision as a one-dimensional random walk of a particle (representing the accumulated evidence) between two boundaries (representing the two choices).[29] The first boundary to be crossed determines both the choice and the decision time.

The explanatory power of EAMs stems from their ability to decompose observable behavioral data—namely, distributions of choices and response times—into a small set of psychologically meaningful latent parameters.[29] These core parameters provide a quantitative description of distinct cognitive subprocesses:

- **Drift Rate ($v$):** This parameter represents the average rate of evidence accumulation. It is considered a measure of the quality of the sensory evidence and the efficiency of information processing. For a given task, a higher drift rate corresponds to an easier decision (stronger evidence), leading to faster and more accurate responses.[29]
- **Boundary Separation ($a$):** This parameter represents the distance between the two decision boundaries, signifying the amount of evidence that must be accumulated before a commitment is made. It is interpreted as a measure of response caution. A wider boundary leads to slower, more deliberate, and typically more accurate decisions, reflecting a strategic emphasis on accuracy over speed (the speed-accuracy tradeoff).[29]
- **Starting Point ($z$):** This parameter denotes the initial position of the evidence accumulator. If it is biased closer to one boundary than the other, it reflects an a priori preference or expectation for that response option.[29]
- **Non-Decision Time ($T_{er}$):** This parameter captures the time consumed by processes outside of the core decision-making, such as sensory encoding of the stimulus and the execution of the motor response.[33]

A crucial point of synthesis is that the attractor network model of Wong & Wang can be understood as a neurally plausible implementation of the cognitive process described by EAMs.[2] There is a direct and intuitive mapping between the components of the two frameworks. The ramping activity of the selective neural populations in the attractor network is the direct neural correlate of the evidence accumulation process in an EAM.[2] The strength of the external stimulus driving the network (e.g., the $\mu_{signal}$ term in the user's proposal) corresponds directly to the drift rate $v$ in the DDM.[1] The decision threshold in an EAM is

analogous to the effective firing rate threshold in the attractor network, which, when crossed, leads to the winner-take-all dynamics "latching" into a stable state.[5] This deep connection is not merely an interesting parallel; it provides a powerful method for validating the "human-centricity" of the neural arbitrator. By demonstrating that the model can replicate canonical findings from the EAM literature—such as the inverse relationship between task difficulty (lower drift rate) and reaction time—one can provide strong evidence that the neural model captures fundamental properties of human cognition.[1] This positions the project not just as an engineering endeavor, but as a contribution to the field of computational cognitive neuroscience by bridging the gap between a high-level cognitive abstraction (the EAM) and a low-level neural implementation (the attractor network).

## 2.2 Standard Recurrent Neural Networks (RNNs) for Noisy Signal Processing

While the attractor network offers a principled, bio-inspired approach, the dominant paradigm in modern machine learning for handling sequential and noisy data is the use of generic, data-driven Recurrent Neural Network (RNN) architectures, most notably Long Short-Term Memory (LSTM) and Gated Recurrent Units (GRU) networks.[35] These models have achieved state-of-the-art performance in a vast array of domains that require processing information over time, making them the primary "black-box" alternative to the proposed neural arbitrator.

Standard RNNs are widely and successfully applied to problems analogous to arbitrating human intent from noisy signals. They are used for automatic speech recognition from continuous acoustic spectrograms [37], for denoising audio features for robust ASR [38], for noise reduction in signals from physical sensors like MEMS gyroscopes [39], and for analyzing complex physiological time-series data, including EEG.[36] Within the BCI domain specifically, various deep learning models, including RNNs and Convolutional Neural Networks (CNNs), are the standard tools for decoding user intent, such as classifying different types of motor imagery from multi-channel EEG signals.[40] The power of these models lies in their ability to learn complex, high-dimensional temporal dependencies directly from large datasets without being constrained by a predefined, biologically inspired architecture.[35]

The architectural contrast between a standard LSTM and the proposed WTA attractor network is fundamental and highlights the core trade-offs of the user's approach. The WTA model has a highly specific, low-dimensional structure: two (or more) competing populations, with explicit parameters for self-excitation ($w_{self}$), mutual inhibition ($w_{inhibit}$), and a time constant ($\tau$) that map directly onto hypothesized circuit properties.[1] Its internal state can be easily visualized as the competing firing rates $R_1$ and $R_2$. In stark contrast, an LSTM network is composed of a chain of generic "cells," each containing a set of input,

output, and forget gates whose behavior is controlled by large matrices of weights learned via backpropagation through time.[35] The internal state of an LSTM is a high-dimensional vector that evolves in a complex, nonlinear way. While this high dimensionality grants it immense expressive power, it does not typically map onto an easily understandable, low-dimensional "decision space" like the one defined by the attractor model's dynamics.[2] This fundamental architectural difference is the origin of the interpretability challenge for standard RNNs, making it difficult to understand *how* or *why* they arrive at a particular decision. This comparison thus sharpens the user's value proposition: the neural arbitrator trades the universal approximation power of a generic RNN for a constrained, but therefore more interpretable and biologically grounded, model of a specific cognitive function.

**Table 1: Comparative Analysis of Decision-Making Model Architectures**

| Criterion | WTA/Attractor Network (User's Model) | Evidence Accumulation Model (DDM) | Black-Box RNN (LSTM) |
|---|---|---|---|
| **Biological Plausibility** | High (Models cortical microcircuit dynamics, synaptic types) | Medium (Cognitive abstraction of neural processes) | Low (Engineering solution, not based on neural architecture) |
| **Mechanism of Action** | Competitive dynamics of recurrently coupled neural populations converging to a stable state. | Stochastic accumulation of a 1D evidence variable to a decision boundary. | Learning complex, high-dimensional patterns in sequential data via gated memory cells. |
| **Interpretability** | **Intrinsic (Process-level):** Internal state (firing rates of competing populations) is directly interpretable as the decision process. | **Parametric (Cognitive-level):** Parameters (drift, boundary) map to latent cognitive constructs. | **Post-hoc:** Requires external explanation tools (e.g., LIME, SHAP) to approximate internal logic. |
| **Handling of Signal Noise** | Inherent robustness via slow | Explicitly modeled as variance (noise) | Learned implicitly from large, noisy |

| | time integration and competitive inhibition dynamics. | in the accumulation process. | training datasets; can be brittle to out-of-distribution noise. |
|---|---|---|---|
| **Primary Application** | Modeling cortical circuits for decision, working memory, and attention. | Explaining choice and response time distributions in cognitive psychology/neuroscience. | General-purpose sequence modeling (e.g., NLP, speech recognition, time-series forecasting). |

# Section 3: Human-Centricity in Brain-Computer Interfaces

The concept of a "human-centric" system goes beyond mere accuracy; it encompasses qualities like trustworthiness, interpretability, and robustness, which are paramount in high-stakes HMI applications such as BCI. The proposed neural arbitrator is founded on the premise that its bio-inspired design endows it with these properties, setting it apart from conventional black-box approaches.[1] This section directly substantiates these claims by linking the model's architectural features to established challenges in the BCI field. It first outlines the critical dilemma of achieving both interpretability and robustness with noisy EEG signals. It then argues that the WTA attractor model provides a form of intrinsic interpretability that circumvents the limitations of post-hoc explanation methods. Finally, it draws on findings from the application of EAMs in other safety-critical domains to demonstrate how such models can provide deep diagnostic insights into the user's cognitive state, paving the way for truly adaptive and robust interfaces.

## 3.1 The Interpretability-Robustness Dilemma in BCI

Brain-Computer Interfaces, particularly non-invasive systems relying on electroencephalography (EEG), face a formidable signal processing challenge. EEG signals are notoriously difficult to work with, characterized by a low signal-to-noise ratio (SNR), high dimensionality, and significant non-stationarity.[41] The desired neural signatures of intent are often buried in noise from various sources, including environmental electromagnetic

interference, muscle artifacts (electromyography), eye movements (electrooculography), and inherent biological variability both within and across individuals.[40]

This challenging data landscape has driven the adoption of powerful deep learning models, which can often achieve high classification accuracy by learning complex, subtle patterns from the raw data.[42] However, this performance comes at a significant cost: a lack of transparency. The complex, nonlinear transformations learned by these "black-box" models are opaque to human understanding, creating a major barrier to their adoption in clinical and safety-critical settings.[43] A BCI that correctly classifies a user's intent 95% of the time is not trustworthy if there is no way to verify that its decisions are based on genuine brain activity rather than spurious correlations with artifacts like jaw clenching or eye blinks.[43] This lack of interpretability raises profound issues of safety, accountability, and user trust.[43]

In response, a vibrant field of "Explainable AI" (XAI) has developed methods to peer inside these black boxes. These methods are often categorized as backpropagation-based (e.g., Class Activation Maps, which highlight important input features by tracking gradients), perturbation-based (e.g., LIME and SHAP, which build simple, local surrogate models to explain individual predictions), or rule-based.[43] While useful, these techniques are fundamentally *post-hoc*—they attempt to explain a model that was not designed to be explainable. Critics argue that such explanations can be fragile, misleading, and may not faithfully represent the model's true internal logic, making them potentially unsuitable for high-stakes decisions where the reason for a failure must be understood completely.[47] This creates a dilemma: developers are often forced to choose between the high performance of opaque models and the lower performance of simpler, more transparent ones.

## 3.2 Bio-Inspired Architectures as Inherently Interpretable Systems

The proposed neural arbitrator offers a compelling path out of this dilemma by belonging to a class of models that are *inherently interpretable*.[47] Its interpretability is not an afterthought applied via a post-hoc tool; it is a direct consequence of its bio-inspired design, where the model's structure and state variables have a clear, pre-defined meaning rooted in neuroscience.[1] This represents a paradigm shift from explaining a black box to building a "white box" or "glass box" from the outset.

The central claim of "Process Interpretability" in the project proposal is realized through the model's low-dimensional and meaningful internal state.[1] The state of the system is not an abstract high-dimensional vector but is simply the firing rates of the competing neural populations (e.g., $R_{Left}$ and $R_{Right}$). This allows for a direct and intuitive visualization of the entire decision-making process as it unfolds in time. One can literally plot

the "struggle" between the competing options, watch as evidence accumulates for one over the other, and observe the moment of "latching" when the winner-take-all dynamic suppresses the losing population.[1] This dynamic trace provides a complete and faithful explanation for every decision the model makes. If the model makes an error, the reason is immediately apparent from the trajectory of its internal state—for example, the evidence for the incorrect choice may have been transiently stronger, or noise may have prematurely pushed the system into the wrong attractor basin.

This form of intrinsic interpretability is qualitatively different from and arguably superior to post-hoc explanations for high-stakes applications. It is grounded in physiology, allowing for a crucial sanity check: one can verify that the model is making its decisions based on neurologically plausible activity patterns (e.g., signals originating from the motor cortex during a motor imagery task) rather than by keying in on confounding artifacts.[43] This directly addresses a primary concern in BCI development. In essence, the model's bio-inspired constraints, which might limit its ability to learn arbitrary patterns compared to a generic LSTM, are the very source of its trustworthiness. By building a model that computes in a way that is analogous to the brain, its reasoning becomes understandable in human and neuroscientific terms. This aligns with a growing movement advocating for the use of inherently interpretable models in domains where decisions have a profound impact on human lives.[47]

## 3.3 EAMs in Practice: Lessons from Safety-Critical Systems

The value of using models that provide insight into latent cognitive processes is not merely theoretical. Research applying Evidence Accumulation Models—the cognitive science counterpart to the neural arbitrator—to complex, real-world, safety-critical domains has demonstrated their immense practical utility for improving human-machine systems.[28] Studies in simulated air-traffic control (ATC), driving, and medical image analysis have shown that EAMs can move beyond coarse performance metrics like accuracy to provide a fine-grained diagnosis of the *cognitive sources* of operator error or success.[29]

By fitting an EAM to an operator's behavioral data, researchers can quantitatively determine whether a performance decrement under high workload is caused by:

1. **Inefficient Information Processing:** A reduction in the EAM's drift rate ($v$), suggesting that the operator is struggling to extract a clear signal from the display.
2. **A Shift in Strategy:** A lowering of the decision boundary ($a$), indicating a strategic shift towards speed at the expense of accuracy, perhaps due to time pressure.
3. **An Increase in Bias:** A shift in the starting point ($z$), revealing that the operator has developed a bias to respond one way over another.

4. **Perceptual or Motor Delays:** An increase in non-decision time ($T_{er}$), pointing to issues with initial stimulus encoding or response execution.

This diagnostic power has direct and actionable implications for HMI design. If analysis reveals that errors are consistently linked to a low drift rate, it suggests a flaw in the interface itself—the information is not being presented in a salient or easily processable way, and the display should be redesigned.[31] Conversely, if errors are linked to a suboptimal decision threshold, it points to a strategic issue that could be addressed through improved operator training or feedback.[31]

The proposed neural arbitrator, as a neural implementation of an EAM, inherits this diagnostic capability. The dynamics of its internal state—such as the average slope of the ramping activity (proxy for drift rate) or the time taken to latch (proxy for response time, influenced by the boundary)—can serve as a real-time window into the user's cognitive state. This opens up a powerful new paradigm for adaptive BCIs. An interface could monitor the internal parameters of its own arbitration model. If it detects a sustained drop in the accumulation rate, suggesting the user is confused or the task is too difficult, it could proactively simplify the display, provide more salient cues, or slow down the interaction pace. This represents a shift from a simple command-and-control BCI to a truly symbiotic HMI, where the machine intelligently adapts to the user's cognitive capacity and state, thereby enhancing both robustness and performance.

# Section 4: Modeling Higher-Order Cognition: Deliberation and Spontaneity

The most scientifically ambitious and novel contribution of the proposed neural arbitrator lies in its capacity to model not only perceptual decisions driven by external stimuli but also spontaneous, internally generated actions.[1] This endeavor connects the engineering goal of building a robust HMI to a deep and long-standing inquiry in neuroscience and philosophy concerning the nature of volition and free will. This section explores this connection by first reviewing the seminal Libet experiment and the debate it ignited. It then details the modern computational resolution to this debate offered by the Stochastic Decision Model (SDM). Finally, it synthesizes these concepts to demonstrate how the proposed WTA attractor network provides a single, unified mechanism capable of arbitrating the full spectrum of human intent, from stimulus-driven reaction to endogenous volition.

## 4.1 The Neuroscience of Spontaneous Action: The Libet Experiment

## and Its Legacy

In a series of landmark experiments in the early 1980s, neurophysiologist Benjamin Libet investigated the temporal relationship between brain activity, conscious intention, and voluntary action.[51] Participants were asked to perform a simple, spontaneous motor act (e.g., flexing their wrist) at a time of their own choosing, while their brain activity was recorded using EEG. They were also asked to report the precise moment they first became aware of the "urge" or "intention" to move by noting the position of a rapidly rotating spot on a clock-like face.[51]

The results were striking and controversial. Libet discovered that a specific pattern of brain activity, a slow negative-going electrical potential over the motor cortex known as the "Readiness Potential" (RP), consistently began to build up long before the action itself. This was expected. The unexpected finding was that the RP began, on average, about 550 milliseconds before the movement, whereas the participants' reported time of conscious intention (dubbed "W" for will) occurred only about 200 milliseconds before the movement.[51] This left a gap of roughly 350 milliseconds during which the brain appeared to be unconsciously preparing the action before the individual was consciously aware of having decided to perform it.[52]

Libet's interpretation of this finding sent shockwaves through the fields of neuroscience, philosophy, and law.[51] He concluded that voluntary actions are initiated unconsciously by the brain. Conscious will, in this view, does not initiate the action but appears later in the process, perhaps serving as a final gatekeeper with the power to "veto" the action before it is executed.[52] This interpretation was seen by many as a direct scientific challenge to the traditional notion of free will, suggesting that our subjective experience of consciously causing our actions might be a form of illusion.[51] For decades, this result has fueled a debate about the causal role of consciousness in human behavior.

## 4.2 Computational Models of "Free Will": The Stochastic Decision Model (SDM)

For many years, the interpretation of the RP as a specific, deterministic plan for action remained the dominant view. However, a more recent and computationally grounded alternative has emerged that fundamentally reframes the meaning of the RP. The Stochastic Decision Model (SDM), most prominently advanced by Schurger, Sitt, and Dehaene in 2012, proposes a new interpretation based on the principles of evidence accumulation models.[56]

The core insight of the SDM is that the RP may not be a representation of an unconscious decision at all. Instead, it could be an artifact that emerges from averaging noisy neural signals under specific experimental constraints. The SDM posits that in a state of waiting to make a spontaneous movement, the brain's motor system exhibits ongoing, random (stochastic) fluctuations in neural activity.[56] The decision to move is not made at the beginning of the RP ramp; rather, it is triggered at the very end, at the moment these random fluctuations happen to cross a fixed motor threshold.[56] The characteristic slow, ramping shape of the RP arises from the experimental procedure of averaging many trials, time-locked to the movement onset and viewed *backwards in time*. Because only trials where the noise happened to drift upwards to the threshold are selected, the average of these traces will inevitably show a gradual ramp leading up to the threshold crossing.

Mathematically, the SDM is modeled as a leaky stochastic accumulator, often described by an Ornstein-Uhlenbeck process, which is a type of random walk with a tendency to drift back towards a baseline.[56] The key components of the model are:

1. **A Noise Term:** This represents the ongoing, spontaneous, and random fluctuations in neural activity. In the context of a spontaneous decision, this is the primary driver of the process.
2. **A Leak Term:** This causes the accumulated activity to decay back to baseline over time, preventing it from drifting away indefinitely.
3. **A Decision Threshold:** A fixed level of neural activity which, when crossed, triggers the commitment to act.
4. **An Imperative Term (Optional):** A small, constant drift or bias towards the threshold, which can represent the subject's general intention to comply with the experimental instruction to move *at some point*.[56]

Crucially, the SDM demonstrates that the entire phenomenon—the shape of the RP and the distribution of waiting times before movement—can be quantitatively reproduced by a simple accumulator model driven largely or entirely by noise, without needing to posit a specific, unconscious decision process that begins hundreds ofmilliseconds in advance.

## 4.3 Synthesis: The Neural Arbitrator as a Unified Model of Intent

The development of the Stochastic Decision Model provides a powerful and direct theoretical validation for the most innovative aspect of the proposed neural arbitrator project.[1] The user's proposal to simulate a "Spontaneous Decision" by setting the mean of the external input streams to zero ($\mu\_1 = \mu\_2 = 0$) and allowing the system's dynamics to be driven solely by the internal stochastic noise term ($I_{noise}$) is a precise implementation of the core hypothesis of the SDM.[1] This experiment is therefore not merely an interesting curiosity but a

computational replication of the leading modern theory for the neural basis of spontaneous action.

This synthesis reveals the profound elegance and power of the WTA attractor network architecture. It demonstrates that a single, unified neural mechanism can account for two apparently distinct modes of decision-making, distinguished only by the nature of their driving inputs:

- **Perceptual Decisions:** When the system is driven by external sensory evidence that contains a meaningful signal (i.e., a non-zero mean difference between the inputs, $\mu_{signal} - \mu_{baseline} > 0$), the network functions as a classic evidence accumulator, integrating that signal over time to make a stimulus-based choice. In this mode, it behaves like the neural implementation of a Drift-Diffusion Model.
- **Spontaneous Decisions:** When the system receives no differential external evidence (i.e., the inputs are balanced with a zero-mean difference, $\mu_1 = \mu_2 = 0$), the network's dynamics are governed by the integration of its own internal noise. A decision is made when these random fluctuations accumulate and, through the network's recurrent dynamics, are amplified until they cross the threshold for one of the attractor states. In this mode, it behaves as a neural implementation of the Stochastic Decision Model.

Therefore, the proposed "Neural Arbitrator" is far more than a simple signal processing tool for BCI. It represents a unified computational model that synthesizes core principles from attractor network theory, evidence accumulation models of perception, and stochastic models of volition. It offers a single, biologically plausible circuit that can compellingly arbitrate the full spectrum of human intent, from reacting to the faintest external cues to generating the most endogenous, "freely-willed" actions. This unified perspective also yields a powerful, non-obvious prediction: because the same underlying network parameters (e.g., the strength of inhibition $w_{inhibit}$, the time constant $\tau$) govern the dynamics in both input regimes, these parameters should have predictable, correlated effects on both perceptual and spontaneous tasks. For example, a change that makes the network more cautious in a perceptual task (e.g., by effectively raising the decision threshold) should also increase the average waiting time for a spontaneous action, providing a clear avenue for rigorous experimental validation of the model's unifying claims.

### Alıntılanan çalışmalar

1. CMPE 58I - Project Proposal.docx
2. A Recurrent Network Mechanism of Time Integration in Perceptual ..., erişim tarihi Ekim 28, 2025, https://www.jneurosci.org/content/26/4/1314
3. A recurrent network mechanism of time integration in perceptual decisions - PubMed - NIH, erişim tarihi Ekim 28, 2025, https://pubmed.ncbi.nlm.nih.gov/16436619/
4. Xiao-Jing Wang - Google Scholar, erişim tarihi Ekim 28, 2025,

https://scholar.google.com/citations?user=cv-YgL0AAAAJ&hl=en

5. A Recurrent Network Mechanism of Time Integration in Perceptual Decisions, erişim tarihi Ekim 28, 2025, https://www.jneurosci.org/content/26/4/1314/tab-figures-data

6. A recurrent network mechanism of time integration in perceptual decisions - Pure - Ulster University's Research Portal, erişim tarihi Ekim 28, 2025, https://pure.ulster.ac.uk/files/11301562/JNS06_supp.pdf

7. Changes of Mind in an Attractor Network of Decision-Making | PLOS ..., erişim tarihi Ekim 28, 2025, https://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1002086

8. (PDF) A Recurrent Network Mechanism of Time Integration in Perceptual Decisions (2006) | Kong Fatt Wong | 1058 Citations - SciSpace, erişim tarihi Ekim 28, 2025, https://scispace.com/papers/a-recurrent-network-mechanism-of-time-integration-in-36pozk7nim?citations_page=5

9. Temporal dynamics underlying perceptual decision making: Insights from the interplay between an attractor model and parietal neurophysiology - Frontiers, erişim tarihi Ekim 28, 2025, https://www.frontiersin.org/journals/neuroscience/articles/10.3389/neuro.01.028.2008/full

10. A Recurrent Network Mechanism of Time Integration in Perceptual Decisions, erişim tarihi Ekim 28, 2025, https://www.semanticscholar.org/paper/A-Recurrent-Network-Mechanism-of-Time-Integration-Wong-Lin-Wang/3cacf7985827c0f17457f7c93ceba9855910c3bd

11. Neural circuit dynamics underlying accumulation of time-varying evidence during perceptual decision making - Frontiers, erişim tarihi Ekim 28, 2025, https://www.frontiersin.org/journals/computational-neuroscience/articles/10.3389/neuro.10.006.2007/full

12. 12. Perceptual Decision Making (Wong & Wang) - Neuronal Dynamics: Python Exercises, erişim tarihi Ekim 28, 2025, https://neuronaldynamics-exercises.readthedocs.io/en/latest/exercises/perceptual-decision-making.html

13. Spiking Inputs to a Winner-take-all Network, erişim tarihi Ekim 28, 2025, https://proceedings.neurips.cc/paper/2005/file/881c6efa917cff1c97a74e03e15f43e8-Paper.pdf

14. Learning and stabilization of winner-take-all dynamics through interacting excitatory and inhibitory plasticity - PMC - NIH, erişim tarihi Ekim 28, 2025, https://pmc.ncbi.nlm.nih.gov/articles/PMC4086298/

15. Mechanisms of Winner-Take-All and Group Selection in Neuronal Spiking Networks - PMC, erişim tarihi Ekim 28, 2025, https://pmc.ncbi.nlm.nih.gov/articles/PMC5399521/

16. Mechanisms of Winner-Take-All and Group Selection in ... - Frontiers, erişim tarihi Ekim 28, 2025, https://www.frontiersin.org/journals/computational-neuroscience/articles/10.3389/fncom.2017.00020/full

17. (PDF) Mechanisms of Winner-Take-All and Group Selection in Neuronal Spiking Networks, erişim tarihi Ekim 28, 2025, https://www.researchgate.net/publication/316433304_Mechanisms_of_Winner-Take-All_and_Group_Selection_in_Neuronal_Spiking_Networks

18. [1808.00675] Winner-Take-All as Basic Probabilistic Inference Unit of Neuronal Circuits, erişim tarihi Ekim 28, 2025, https://arxiv.org/abs/1808.00675

19. [2508.07471] Modeling bias in decision-making attractor networks - arXiv, erişim tarihi Ekim 28, 2025, https://arxiv.org/abs/2508.07471

20. Modeling Bias in Decision-Making Attractor Networks - arXiv, erişim tarihi Ekim 28, 2025, https://arxiv.org/html/2508.07471v1

21. MIT Open Access Articles Attractor and integrator networks in the brain, erişim tarihi Ekim 28, 2025, https://dspace.mit.edu/bitstream/handle/1721.1/148794/2112.03978.pdf?sequence=2&isAllowed=y

22. Attractor and integrator networks in the brain, erişim tarihi Ekim 28, 2025, https://mcgovern.mit.edu/wp-content/uploads/2024/05/s41583-022-00642-0.pdf

23. Attractor networks - Oxford Centre for Computational Neuroscience, erişim tarihi Ekim 28, 2025, https://www.oxcns.org/papers/476%20Rolls%20Attractor%20Networks%202010.pdf

24. Perceptual Decision-Making: Biases in Post-Error Reaction Times Explained by Attractor Network Dynamics | Journal of Neuroscience, erişim tarihi Ekim 28, 2025, https://www.jneurosci.org/content/39/5/833

25. (PDF) Modeling bias in decision-making attractor networks - ResearchGate, erişim tarihi Ekim 28, 2025, https://www.researchgate.net/publication/394440147_Modeling_bias_in_decision-making_attractor_networks

26. Attractor and integrator networks in the brain | Request PDF - ResearchGate, erişim tarihi Ekim 28, 2025, https://www.researchgate.net/publication/365092043_Attractor_and_integrator_networks_in_the_brain

27. Noise in Attractor Networks in the Brain Produced by Graded Firing Rate Representations - Research journals - PLOS, erişim tarihi Ekim 28, 2025, https://journals.plos.org/plosone/article/file?id=10.1371/journal.pone.0023630&type=printable

28. Evidence accumulation modelling in the wild: understanding safety, erişim tarihi Ekim 28, 2025, https://research-repository.uwa.edu.au/en/publications/evidence-accumulation-modelling-in-the-wild-understanding-safety-

29. Evidence accumulation modelling in the wild: understanding safety-critical decisions, erişim tarihi Ekim 28, 2025, https://pure.uva.nl/ws/files/109604890/1_s2.0_S1364661322002947_main.pdf

30. Evidence accumulation modelling in the wild: understanding safety ..., erişim tarihi Ekim 28, 2025, https://pubmed.ncbi.nlm.nih.gov/36473764/

31. Evidence accumulation modelling in the wild: Understanding safety ..., erişim tarihi

Ekim 28, 2025,
https://www.ampl-psych.com/wp-content/uploads/2022/11/Boag_et_al_2022_TIC
S_Evidence_accumulation_in_the_wild_preprint.pdf

32. Joint Modelling of Latent Cognitive Mechanisms Shared Across Decision-Making Domains, erişim tarihi Ekim 28, 2025,
https://pmc.ncbi.nlm.nih.gov/articles/PMC10899373/

33. A computational framework for integrating Predictive processes with evidence Accumulation Models (PAM) - ResearchGate, erişim tarihi Ekim 28, 2025,
https://www.researchgate.net/publication/386013677_A_computational_framewo
rk_for_integrating_Predictive_processes_with_evidence_Accumulation_Models_P
AM

34. Neural underpinnings of the evidence accumulator - Hanks Lab, erişim tarihi Ekim 28, 2025,
https://hankslab.faculty.ucdavis.edu/wp-content/uploads/sites/305/2016/03/Brody
_Hanks_2016.pdf

35. Recurrent Neural Networks: A Comprehensive Review of Architectures, Variants, and Applications - MDPI, erişim tarihi Ekim 28, 2025,
https://www.mdpi.com/2078-2489/15/9/517

36. A review of recurrent neural network-based methods in computational physiology - PMC, erişim tarihi Ekim 28, 2025,
https://pmc.ncbi.nlm.nih.gov/articles/PMC10589904/

37. Recurrent neural networks as neuro-computational models of human speech recognition, erişim tarihi Ekim 28, 2025,
https://www.biorxiv.org/content/10.1101/2024.02.20.580731v1.full-text

38. Recurrent Neural Networks for Noise Reduction in Robust ASR - Google Research, erişim tarihi Ekim 28, 2025,
https://research.google/pubs/recurrent-neural-networks-for-noise-reduction-in-
robust-asr/

39. Hybrid Deep Recurrent Neural Networks for Noise Reduction of MEMS-IMU with Static and Dynamic Conditions - MDPI, erişim tarihi Ekim 28, 2025,
https://www.mdpi.com/2072-666X/12/2/214

40. A composite improved attention convolutional network for motor imagery EEG classification - PMC - NIH, erişim tarihi Ekim 28, 2025,
https://pmc.ncbi.nlm.nih.gov/articles/PMC11841462/

41. Bidirectional feature pyramid attention-based temporal convolutional network model for motor imagery electroencephalogram classification - Frontiers, erişim tarihi Ekim 28, 2025,
https://www.frontiersin.org/journals/neurorobotics/articles/10.3389/fnbot.2024.13
43249/full

42. Decoding Brain Signals: A Convolutional Neural Network Approach for Motor Imagery Classification - ResearchGate, erişim tarihi Ekim 28, 2025,
https://www.researchgate.net/publication/378062630_Decoding_Brain_Signals_A
_Convolutional_Neural_Network_Approach_for_Motor_Imagery_Classification

43. Interpretable and Robust AI in EEG Systems: A Survey - arXiv, erişim tarihi Ekim 28, 2025, https://arxiv.org/html/2304.10755v4

44. Uncertainty-Aware Deep Learning for Robust and Interpretable MI EEG Using Channel Dropout and LayerCAM Integration - MDPI, erişim tarihi Ekim 28, 2025, https://www.mdpi.com/2076-3417/15/14/8036

45. Interpretable and Robust AI in EEG Systems: A Survey - SciSpace, erişim tarihi Ekim 28, 2025, https://scispace.com/pdf/interpretable-and-robust-ai-in-eeg-systems-a-survey-3ebeqo00.pdf

46. Interpretable and Robust AI in EEG Systems: A Survey - arXiv, erişim tarihi Ekim 28, 2025, https://arxiv.org/html/2304.10755v3

47. Stop Explaining Black Box Machine Learning Models for High Stakes Decisions and Use Interpretable Models Instead - PMC - NIH, erişim tarihi Ekim 28, 2025, https://pmc.ncbi.nlm.nih.gov/articles/PMC9122117/

48. Interpretable and Robust AI in Electroencephalogram Systems - DR-NTU, erişim tarihi Ekim 28, 2025, https://dr.ntu.edu.sg/bitstreams/478efd43-900d-4879-9ced-b1c413e02c66/download

49. Demystifying the Black Box: The Importance of Interpretability of Predictive Models in Neurocritical Care - PMC - NIH, erişim tarihi Ekim 28, 2025, https://pmc.ncbi.nlm.nih.gov/articles/PMC9343258/

50. Interpretable and Robust AI in EEG Systems: A Survey | Request PDF - ResearchGate, erişim tarihi Ekim 28, 2025, https://www.researchgate.net/publication/370213217_Interpretable_and_Robust_AI_in_EEG_Systems_A_Survey

51. Volition and the Brain – Revisiting a Classic Experimental Study - PMC - NIH, erişim tarihi Ekim 28, 2025, https://pmc.ncbi.nlm.nih.gov/articles/PMC6024487/

52. Readiness Potential and Neuronal Determinism: New Insights on Libet Experiment, erişim tarihi Ekim 28, 2025, https://www.jneurosci.org/content/38/4/784

53. [PDF] Benjamin Libet Do We Have Free Will - Semantic Scholar, erişim tarihi Ekim 28, 2025, https://www.semanticscholar.org/paper/Benjamin-Libet-Do-We-Have-Free-Will-Libet/89db6d0272f2f086309ca778ac52b5e1c28c2e3d

54. Free Will and Neuroscience: From Explaining Freedom Away to New Ways of Operationalizing and Measuring It - Frontiers, erişim tarihi Ekim 28, 2025, https://www.frontiersin.org/journals/human-neuroscience/articles/10.3389/fnhum.2016.00262/full

55. Deliberation period during easy and difficult decisions: re-examining Libet's "veto" window in a more ecologically valid framework | Neuroscience of Consciousness | Oxford Academic, erişim tarihi Ekim 28, 2025, https://academic.oup.com/nc/article/2017/1/nix002/3066355

56. Clarifying the nature of stochastic fluctuations and ... - Frontiers, erişim tarihi Ekim 28, 2025, https://www.frontiersin.org/journals/psychology/articles/10.3389/fpsyg.2023.1271180/full