

Article**When Should Neural Data Inform Welfare? A Critical Framework for Policy Uses of Neuroeconomics**Louis Zhu¹¹Student (BSc, 2023-26), Dept. of Science and Technology Studies, UCL, UK;louis.zhu.23@ucl.ac.uk**Abstract**

Neuroeconomics promises to ground welfare analysis in neural and computational evidence about how people value outcomes, learn from experience and exercise self-control. At the same time, policy and commercial actors increasingly invoke neural data to justify paternalistic regulation, “brain-based” interventions and new welfare measures. This paper asks under what conditions neural data can legitimately inform welfare judgements for policy rather than merely describing behaviour. I develop a non-empirical, model-based framework that links three levels: neural signals, computational decision models and normative welfare criteria. Within an actor–critic reinforcement-learning model, I formalise the inference path from neural activity to latent values and prediction errors and then to welfare claims. I show that neural evidence constrains welfare judgements only when the neural–computational mapping is well validated, the decision model identifies “true” interests versus context-dependent mistakes, and the welfare criterion is explicitly specified and defended. Applying the framework to addiction, neuromarketing and environmental policy, I derive a Neuroeconomic Welfare Inference Checklist for regulators and for designers of NeuroAI systems. The analysis treats brains and artificial agents as value-learning systems while showing that internal reward signals, whether biological or artificial, are computational quantities and cannot be treated as welfare measures without an explicit normative model.

Keywords: *Neuroeconomics; Behavioural Welfare Economics; Reinforcement Learning; NeuroAI; Human-Inspired AI*

1 Introduction

Neuroeconomics is commonly defined as the study of the neurobiological and computational basis of value-based decision-making.¹ It aims to provide a biologically grounded account of behaviour that can be applied in both the natural and social sciences. Early work on dopaminergic reward-prediction errors and cortical value representations suggested new ways to think about intertemporal choice, risk and self-control.²

At the same time, welfare economics and public policy have begun to draw on these findings. Behavioural public economics asks how to do welfare analysis when preferences are inconsistent or context-dependent. Neuroeconomics appears to offer mechanistic evidence about which choices reflect “true” interests and which reflect mistakes or cue-triggered lapses.³ Applications now range from consumer protection and retirement saving to health and environmental policy, alongside a growing commercial ecosystem in neuromarketing and consumer neuroscience.⁴

This expansion raises a central problem. Neural data are noisy, paradigm-dependent and theory-laden, where welfare is a normative construct.⁵ It encodes judgements about when a person is better or worse off. The mere fact that a choice correlates with activation in the ventral striatum or vmPFC does not tell

¹ Antonio Rangel, Colin Camerer, and P. Read Montague, “A Framework for Studying the Neurobiology of Value-Based Decision Making,” *Nature Reviews Neuroscience* 9, no. 7 (2008): 545–556, <https://doi.org/10.1038/nrn2357>; Colin Camerer, George Loewenstein, and Drazen Prelec, “Neuroeconomics: How Neuroscience Can Inform Economics,” *Journal of Economic Literature* 43, no. 1 (2005): 9–64, <https://doi.org/10.1257/0022051053737843>.

² W. Schultz, P. Dayan, and P. R. Montague, “A Neural Substrate of Prediction and Reward,” *Science* 275, no. 5306 (1997): 1593–1599, <https://doi.org/10.1126/science.275.5306.1593>; Camerer, Loewenstein, and Prelec, “Neuroeconomics.”

³ B. Douglas Bernheim and Antonio Rangel, “Addiction and Cue-Triggered Decision Processes,” *American Economic Review* 94, no. 5 (2004): 1558–1590, <https://doi.org/10.1257/0002828043052222>;

B. Douglas Bernheim and Antonio Rangel, “Behavioral Public Economics: Welfare and Policy Analysis with Non-Standard Decision-Makers,” *NBER Working Paper* 11518 (August 2005), <https://doi.org/10.3386/w11518>;

B. Douglas Bernheim, “On the Potential of Neuroeconomics: A Critical (but Hopeful) Appraisal,” *American Economic Journal: Microeconomics* 1, no. 2 (2009): 1–41, <https://doi.org/10.1257/mic.1.2.1>;

B. Douglas Bernheim and Dmitry Taubinsky, “Behavioral Public Economics,” *SSRN Working Paper* (July 2018), https://papers.ssrn.com/sol3/papers.cfm?abstract_id=3214374.

⁴ Ming Hsu and Carolyn Yoon, “The Neuroscience of Consumer Choice,” *Current Opinion in Behavioral Sciences* 5 (October 2015): 116–121, <https://doi.org/10.1016/j.cobeha.2015.09.005>;

Ahmed H. Alsharif and Salmi Mohd Isa, “Revolutionizing Consumer Insights: The Impact of FMRI in Neuromarketing Research,” *Future Business Journal* 10, no. 1 (2024), <https://doi.org/10.1186/s43093-024-00371-z>;

Raveena Gupta, Anuj Pal Kapoor, and Harsh V. Verma, “Neuro-Insights: A Systematic Review of Neuromarketing Perspectives across Consumer Buying Stages,” *Frontiers in Neuroergonomics* 6 (July 2025), <https://doi.org/10.3389/fnrgo.2025.1542847>.

⁵ Ian Krajbich and Mark Dean, “How Can Neuroscience Inform Economics?” *Current Opinion in Behavioral Sciences* 5 (October 2015): 51–57, <https://doi.org/10.1016/j.cobeha.2015.07.005>.

us whether that choice advances the person's welfare. Without a careful framework, the authority of brain images risks being used to re-label contested welfare judgements as scientific facts.⁶

The thesis of this paper is that neural data can inform welfare only via a clearly specified computational model and a clearly defended normative standard. Brains implement algorithms. Neuroeconomics allows us to infer aspects of those algorithms from neural data and behaviour. Welfare analysis requires an additional step: deciding which outputs of those algorithms count as welfare-relevant in which contexts.

The argument proceeds in three steps:

1. I formalise the link between neural activity and computational quantities in value-based decision-making, drawing on the Rangel–Camerer–Montague five-stage framework and on reward-prediction-error theories of dopamine.⁷
2. I combine this with behavioural-welfare frameworks that distinguish welfare from observed choice, for example by treating some choices as mistakes or cue-triggered lapses.⁸
3. I derive conditions under which neural evidence can legitimately support those distinctions, and I illustrate them in three policy domains: addiction, neuromarketing and environmental policy.

This paper makes three contributions. First, it integrates reinforcement-learning based neuroeconomic models with behavioural welfare economics by providing a unified, three-level framework that links neural signals, computational decision processes and normative welfare criteria. Second, it derives a set of necessary conditions under which neural data can genuinely constrain welfare judgements and expresses them as an operational checklist for applied work in addiction, neuromarketing and environmental policy. Third, it extends these ideas to human-inspired AI by treating actor–critic architectures and dopaminergic prediction errors as templates for artificial value-learning agents, and by showing why their internal reward and value signals cannot be equated with welfare without an explicit normative model.

⁶ Bernheim, “On the Potential of Neuroeconomics”; Cass R. Sunstein, “Behavioral Welfare Economics,” *Journal of Benefit-Cost Analysis* 11, no. 2 (2020): 196–220, <https://doi.org/10.1017/bca.2020.14>.

⁷ Rangel, Camerer, and Montague, “Neurobiology of Value-Based Decision Making”; Schultz, Dayan, and Montague, “Neural Substrate of Prediction and Reward”; Will Dabney et al., “A Distributional Code for Value in Dopamine-Based Reinforcement Learning,” *Nature* 577, no. 7792 (2020): 671–675, <https://doi.org/10.1038/s41586-019-1924-6>.

⁸ Bernheim, “On the Potential of Neuroeconomics”; Bernheim and Rangel, “Addiction and Cue-Triggered Decision Processes”; Bernheim and Rangel, “Behavioral Public Economics”; B. Douglas Bernheim and Antonio Rangel, “Beyond Revealed Preference: Choice-Theoretic Foundations for Behavioral Welfare Economics,” *Quarterly Journal of Economics* 124, no. 1 (2009): 51–104, <https://doi.org/10.1162/qjec.2009.124.1.51>; Bernheim and Taubinsky, “Behavioral Public Economics.”

Throughout, I connect to NeuroAI. Reinforcement learning and distributional reinforcement learning have become shared formalisms for modelling both dopamine systems and artificial agents.⁹ Both brains and artificial systems can be treated as value-learning processes, but welfare remains a separate, normative layer.

Conceptually, the paper distinguishes three levels: neural data (spikes, BOLD), computational models (values, prediction errors, policies) and welfare judgements (better/worse for the person or society). Neural evidence can constrain welfare analysis only by way of a computational model, and welfare claims always require an additional normative step.¹⁰ This three-level structure is made concrete in Sections 3–4 using an actor–critic reinforcement-learning architecture and its neural implementation (Figures 1–3).

2 Neuroeconomics, welfare and NeuroAI

2.1 Neuroeconomics as value-based decision science

Rangel, Camerer and Montague propose that value-based decisions can be decomposed into five processes: (1) representing the decision problem; (2) valuing candidate options; (3) selecting an action; (4) evaluating the outcome; and (5) learning to update value representations.¹¹

Empirically, orbitofrontal and ventromedial prefrontal cortex encode subjective values of options, often in a “common currency” across reward modalities, while ventral striatum and midbrain dopamine neurons exhibit activity patterns consistent with reward-prediction errors.¹² Dorsolateral prefrontal cortex and anterior cingulate cortex contribute to executive control, conflict monitoring and exploration.

A key implication is that these neural variables track subjective evaluations, contingent on expectations, reference points and affective states.¹³ High striatal activation is not a direct readout of welfare; it reflects context-dependent value signals.

⁹ Kenji Doya, “Reinforcement Learning: Computational Theory and Biological Mechanisms,” *HFSP Journal* 1, no. 1 (2007): 30–40, <https://doi.org/10.2976/1.2732246>; Dabney et al., “A Distributional Code for Value”; Sadra Sadeh and Claudia Clopath, “The Emergence of NeuroAI: Bridging Neuroscience and Artificial Intelligence,” *Nature Reviews Neuroscience* 26, no. 10 (2025): 583–584, <https://doi.org/10.1038/s41583-025-00954-x>.

¹⁰ Bernheim, “On the Potential of Neuroeconomics”; Bernheim and Taubinsky, “Behavioral Public Economics.”

¹¹ Rangel, Camerer, and Montague, “Neurobiology of Value-Based Decision Making”

¹² Schultz, Dayan, and Montague, “Neural Substrate of Prediction and Reward”; Rangel, Camerer, and Montague, “Neurobiology of Value-Based Decision Making”; Hsu and Yoon, “Neuroscience of Consumer Choice.”

¹³ Camerer, Loewenstein, and Prelec, “Neuroeconomics.”

2.2 Welfare economics and behavioural public economics

Standard welfare economics identifies welfare with the satisfaction of coherent, stable preferences revealed in choice. When preferences are complete, transitive and context-independent, a utility function can represent them and serve as a welfare criterion.

Behavioural economics documents systematic violations: present bias, framing effects, loss aversion, limited attention and self-control problems. Behavioural public economics therefore develops alternative welfare notions that can accommodate such phenomena.¹⁴

One influential strand distinguishes welfare-relevant from mistake-driven choices. For example, an “as-judged-by-themselves” criterion seeks to recover the preferences individuals would endorse under better conditions, free of errors and self-control failures.¹⁵

Neuroeconomics is attractive here because it appears to offer mechanistic evidence about which choices are genuine expressions of underlying goals and which reflect cue-triggered or pathological processes. Addiction is a central case.¹⁶

2.3 NeuroAI, human-inspired AI and value-learning systems

Reinforcement learning (RL) offers a natural bridge between neuroeconomics and AI. Temporal-difference (TD) learning updates value estimates using prediction errors, a formalism closely matching dopamine neuron activity in conditioning and choice tasks.¹⁷

Dabney et al. show that dopamine populations in mice encode a distributional value signal, with different neurons tuned to different reward quantiles, mirroring distributional RL algorithms in machine learning¹⁸. Sadeh and Clopath describe this two-way traffic as the core of “NeuroAI”: AI tools accelerate neuroscience, while neural principles inform new AI architectures.¹⁹

¹⁴ Bernheim and Rangel, “Behavioral Public Economics”; Bernheim and Rangel, “Beyond Revealed Preference”; Bernheim and Taubinsky, “Behavioral Public Economics.”; Sunstein, “Behavioral Welfare Economics.”

¹⁵ Bernheim and Rangel, “Beyond Revealed Preference”; Sunstein, “Behavioral Welfare Economics.”

¹⁶ Bernheim and Rangel, “Addiction and Cue-Triggered Decision Processes.”

¹⁷ Schultz, Dayan, and Montague, “Neural Substrate of Prediction and Reward”; Doya, “Reinforcement Learning”; Dabney et al., “A Distributional Code for Value.”

¹⁸ Dabney et al., “A Distributional Code for Value.”

¹⁹ Sadeh and Clopath, “Emergence of NeuroAI.”

For present purposes, the crucial point is conceptual. Both brains and artificial agents can be treated as value-learning systems: they maintain value functions, policies and learning rules. Welfare questions arise when we ask which values these systems *ought* to learn, and which policies are better for the agents and for others.

NeuroAI can therefore sharpen our understanding of mechanisms and algorithms, but it does not, by itself, answer welfare questions. It supplies richer models of how agents update and act; welfare analysis must decide how to evaluate those behaviours.

From a human-inspired AI perspective, actor–critic TD learning, distributional value codes and neuromodulatory teaching signals are not just descriptive models of the brain but also templates for AI architectures. Many contemporary RL agents are deliberately designed to mirror these mechanisms. The framework developed in this paper therefore applies symmetrically to neural and artificial systems: it treats both as value-learning processes with internal reward and value representations, and asks under what conditions those internal quantities can justifiably be used for welfare-relevant purposes.

This is directly relevant to current practice in reinforcement learning, where agents are often trained with human-derived feedback or rewards and their internal reward signals are routinely used as proxies for user welfare; the framework explains why that identification requires additional normative argument rather than being guaranteed by the biological inspiration alone.

3 A formal mapping from neural signals to welfare claims²⁰

This section sets out a compact formal framework that makes explicit the steps from neural data to welfare-relevant judgements. The point is not to provide a full empirical model, but to show where normative assumptions enter.

3.1 Decision and learning model

Consider a sequential decision problem with discrete time t . The individual faces states $s_t \in S$, chooses actions $a_t \in A$, receives rewards r_t and transitions to new states s_{t+1} . A standard reinforcement-learning description posits a value function $V_\theta(s)$ or an action-value function $Q_\theta(s, a)$, a policy $\pi_\phi(a | s)$ and a learning rule that updates θ on the basis of prediction errors.

²⁰ *Appendix A* provides the complete mathematical formulation of the decision environment, value functions, learning rule and welfare representation.

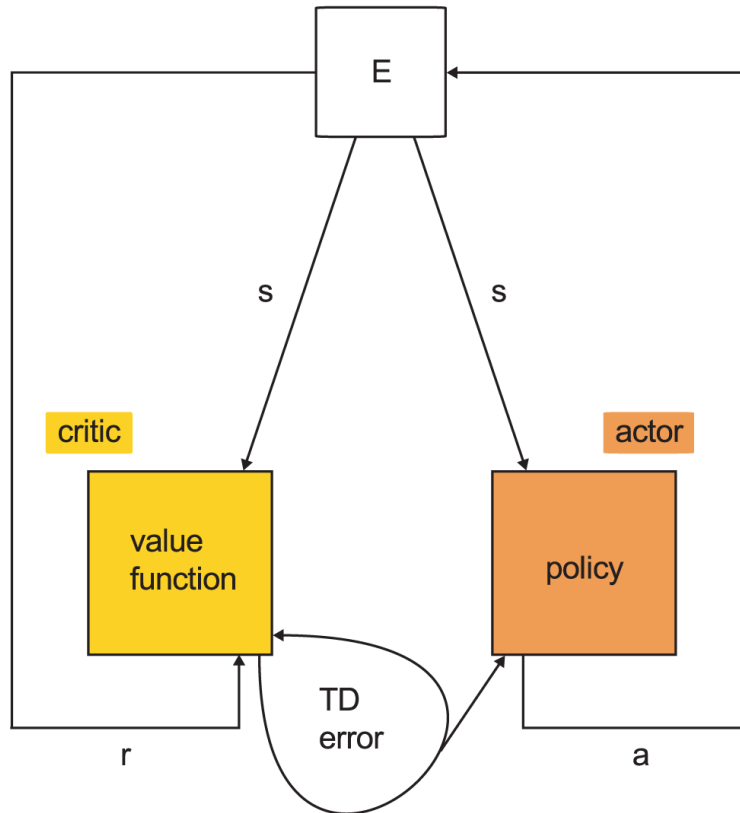


Figure 1. Actor-critic temporal-difference architecture²¹

The environment (E) informs the critic and the actor about the current state s . It also delivers the current reward r to the critic. Using the value function for the current state and the successor state together with the reward, the critic computes a temporal-difference (TD) error signal that is used to update both the value function and the policy. The actor selects an action a according to the current policy, which the environment then executes.

Source: Potjans, Diesmann and Morrison (2011, Fig. 1), PLOS Computational Biology, CC BY 4.0.

Figure 1 shows the classic actor-critic temporal-difference architecture.²² The environment E supplies the current state s_t and reward r_t to the critic (value-function module) and the current state s_t to the actor (policy module). The critic computes a temporal-difference (TD) error and uses it to update both the value function and the policy; the actor samples an action a_t from the current policy, which is fed back to the environment.

²¹ Wiebke Potjans, Markus Diesmann, and Abigail Morrison, "An Imperfect Dopaminergic Error Signal Can Drive Temporal-Difference Learning," *PLOS Computational Biology* 7, no. 5 (2011): e1001133, <https://doi.org/10.1371/journal.pcbi.1001133>.

²² Potjans, Diesmann, and Morrison, "Imperfect Dopaminergic Error Signal."

Formally, we can write a TD-style update as:

$$\delta_t = r_t + \gamma V_\theta(s_{t+1}) - V_\theta(s_t)$$

$$\theta_{t+1} = \theta_t + \alpha \delta_t \nabla_\theta V_\theta(s_t),$$

where γ is a discount factor and α is a learning rate. Choices then follow a softmax policy:

$$\pi_\phi(a \mid s_t) = \frac{\exp(\beta Q_\theta(s_t, a))}{\sum_{a'} \exp(\beta Q_\theta(s_t, a'))},$$

where β is an inverse-temperature parameter controlling sensitivity to value differences.

This actor–critic scheme is implemented both in artificial agents and in biologically plausible network models, providing a natural bridge to neuroeconomics and NeuroAI.²³

3.2 Neural signals as noisy encodings of computational quantities

Let n_t denote neural measurements at time t (for example, firing rates or BOLD activity). In the simplest case, we can treat these as noisy encodings of latent computational variables:

$$n_t = f(\delta_t) + \varepsilon_t \text{ or } n_t^{\text{vmPFC}} = g(V_\theta(s_t)) + \eta_t,$$

where f and g are monotone link functions and ε_t, η_t are noise terms. In the first case n_t tracks a temporal-difference prediction error δ_t ; in the second, n_t^{vmPFC} tracks a state value $V_\theta(s_t)$.

²³ Potjans, Diesmann, and Morrison, “Imperfect Dopaminergic Error Signal”; Doya, “Reinforcement Learning.”

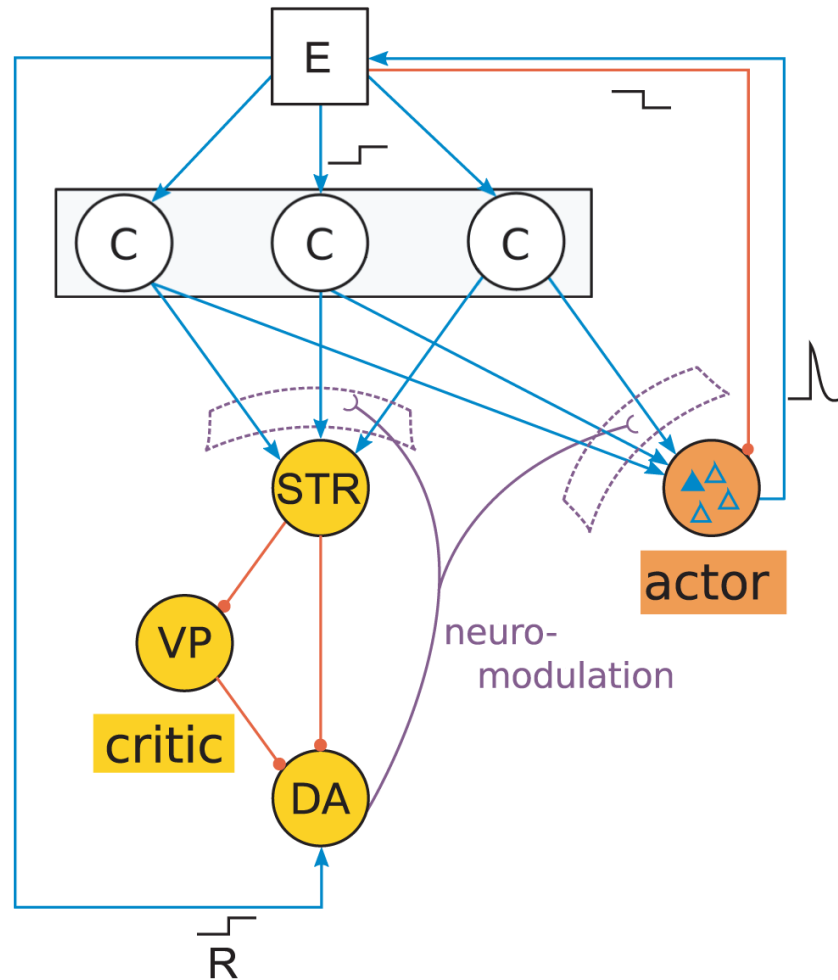


Figure 2. Neuronal actor–critic architecture generating and exploiting a dopaminergic TD error signal²⁴

The input layer consists of cortical populations (C) representing state information. The critic comprises striatal neurons (STR), ventral pallidum (VP) and dopaminergic (DA) neurons. Cortical inputs drive these populations; STR and VP project to DA with different delays. The actor consists of action-selective neurons receiving cortical input. The environment (E) stimulates state-coding cortical neurons and interprets the first-spiking actor neuron as the chosen action. Reward (R) delivered by the environment is encoded as a DC input to DA neurons. The DA signal globally modulates plasticity at cortico-striatal and cortico-actor synapses, implementing a TD-like teaching signal. Red lines indicate inhibitory connections; blue lines, excitatory connections; purple lines, dopaminergic neuromodulation.

Source: Potjans, Diesmann and Morrison (2011, Fig. 2), PLOS Computational Biology, CC BY 4.0.

²⁴ Potjans, Diesmann, and Morrison, “Imperfect Dopaminergic Error Signal.”

Figure 2 extends the abstract actor–critic of Figure 1 to a neuronal actor–critic architecture.²⁵ Cortical populations C encode state information; the critic is implemented in STR, VP and DA neurons; the actor is a pool of action-selective neurons. State-dependent cortical inputs drive STR, VP and DA; DA neurons broadcast a neuromodulatory signal that globally modulates plasticity at cortico-striatal and cortico-actor synapses. In this model, the dopaminergic signal plays the role of a TD-error-like teaching signal in the actor–critic loop.

Empirically, phasic dopamine activity does resemble a TD prediction error in classic conditioning tasks. Before learning, unpredicted rewards elicit brief dopaminergic bursts; after learning, the burst shifts to the predictive cue and disappears at reward time; omission of an expected reward produces a dip below baseline.²⁶

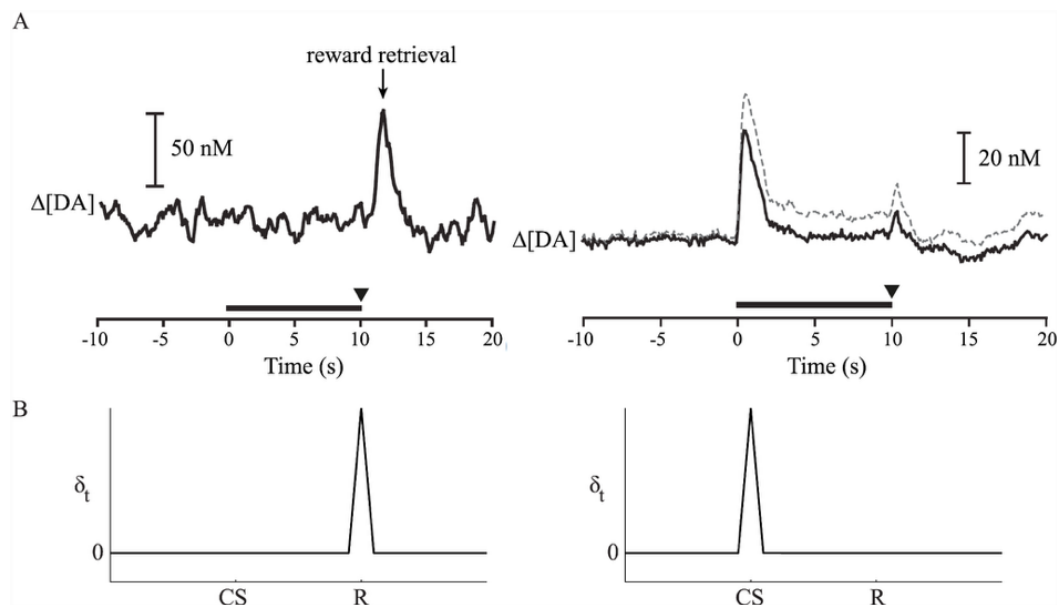


Figure 3. Phasic dopamine signals resembling a temporal-difference prediction error²⁷

Panel A: before learning, an unpredicted reward elicits a phasic increase in dopamine; after learning, the dopamine burst shifts to the predictive cue and the predicted reward elicits little or no response; omission of an expected reward produces a dip below baseline. Panel B: a temporal-difference (TD) model of reinforcement learning produces an error signal δ_t with the same qualitative pattern (positive at

²⁵ Potjans, Diesmann, and Morrison, “Imperfect Dopaminergic Error Signal.”

²⁶ Schultz, Dayan, and Montague, “Neural Substrate of Prediction and Reward.”

²⁷ Kevin Lloyd and Peter Dayan, “Taming Ramping: Algorithmic, Implementational, and Computational Explanations of Phasic Dopamine Signals in the Accumbens,” *PLOS Computational Biology* 11, no. 12 (2015): e1004622, <https://doi.org/10.1371/journal.pcbi.1004622>.

unpredicted rewards, shifting to cues with learning, negative at omitted rewards). This correspondence underpins the interpretation of dopamine as a TD-like prediction-error signal.

Source: adapted from Lloyd and Dayan (2015, Fig. 1), PLOS Computational Biology, CC BY 4.0.

Figure 3 summarises this mapping graphically: panel A shows dopamine concentration (or firing rate) shifting from reward to cue over learning; panel B shows the corresponding TD error δ_t in a simple RL model. The tight qualitative match justifies, at least in this paradigm class, modelling dopaminergic activity as $n_t \approx f(\delta_t)$.²⁸

In our framework, neural signals therefore provide noisy constraints on the latent variables of an RL model: dopamine gives evidence about δ_t ; vmPFC and ventral striatum about V_θ ; prefrontal regions about control signals that govern π_ϕ . Whether and how those constrained models inform welfare analysis is the topic of Section 4.

3.3 Welfare criteria and “true interests”

Suppose we want to evaluate whether a policy p that changes the decision environment makes individuals better off. Behavioural-welfare frameworks typically distinguish between the implemented decision process and a welfare criterion.²⁹

Let $U(s, a)$ denote a welfare-relevant utility function. Let $V_\theta(s)$ and π_ϕ describe the actual value function and policy implemented by the person. These need not coincide. In particular, there may be states in which π_ϕ selects actions that reduce U .

A dual-self self-control model makes this explicit. Let $u^L(a)$ be a long-run utility and $u^S(a)$ a short-run utility. The implemented instant utility might be:

$$U^{\text{impl}}(a) = \lambda u^L(a) + (1 - \lambda) u^S(a),$$

with $\lambda \in [0,1]$ reflecting the influence of long-run considerations. In addiction, cue exposure may transiently reduce λ , shifting weight toward the short-run self.³⁰

²⁸ Lloyd and Dayan, “Tamping Ramping.”

²⁹ Bernheim and Rangel, “Beyond Revealed Preference”; Bernheim and Taubinsky, “Behavioral Public Economics.”

³⁰ Bernheim and Rangel, “Addiction and Cue-Triggered Decision Processes.”

A welfare criterion might take u^L as normative and treat episodes dominated by u^S under certain cue conditions as mistakes. The important point is that welfare is defined in terms of U , not directly in terms of V or n_t . Neural and behavioural data tell us about V , π and the environments in which they are distorted; they do not determine U .

4 When do neural data legitimately inform welfare?

Within this framework, neural data can constrain welfare analysis only under several conditions.

4.1 Measurement and representational validity

The actor–critic architecture in Figures 1–2, together with the empirical mapping between dopamine and TD error in Figure 3, motivates treating neural data as evidence about δ_t , V_θ and π_ϕ . For such evidence to bear on welfare, however, we need further conditions.

First, neural signals must be valid measures of specified computational variables.

1. Task-level validation. There should be convergent evidence that n_t tracks δ_t or V_θ across tasks of the relevant class, not just in a single paradigm. Reward-prediction-error interpretations of dopamine are supported by conditioning, decision and movement studies.³¹
2. Representational specificity. The mapping

$$n_t = f(\delta_t) + \varepsilon_t$$

should be specific enough that alternative computational models do not explain the data equally well. Otherwise, neural data do not favour any particular welfare analysis that depends on one model over another.³²

3. Robustness across individuals and contexts. Encoding must be robust, or systematic variations must be understood and modelled (for example, stress- or mood-dependent modulation of value signals).

³¹ Schultz, Dayan, and Montague, “Neural Substrate of Prediction and Reward”

³² Krjbich and Dean, “How Can Neuroscience Inform Economics?”.

Without these, neural data cannot be straightforwardly treated as measures of value or prediction error and thus have limited welfare relevance.

4.2 Model identifiability

Even with a plausible encoding, there is often a model identifiability problem. Multiple decision and learning models can reproduce the same pattern of behaviour and neural responses.³³

Formally, suppose two models M and M' with parameters (θ, ϕ) and (θ', ϕ') generate the same joint distribution of actions and neural data in the observed environment:

$$P(a_t, n_t \mid M, \theta, \phi) = P(a_t, n_t \mid M', \theta', \phi')$$

for all observed t . If the welfare analysis associated with M and M' differs because they embed different welfare criteria or mistake classifications, neural data do not discriminate between those welfare judgements.

A second condition, then, is:

Model identification condition: in the domain of application, the model class is sufficiently constrained, and alternative models with distinct welfare implications are either empirically ruled out or explicitly acknowledged.³⁴

4.3 Normative transparency

A third condition concerns the normative move from computational variables to welfare. When authors claim that one neural pattern reflects a “true” valuation and another a “biased” valuation, they rely on implicit assumptions about which brain state is authentic or welfare-relevant.

In the notation above, they are selecting a welfare function U and classifying some states as mistake states. The neuroscientific contribution lies in mapping states to mechanisms and correlations; the welfare judgement is a separate, normative step.³⁵

³³ Bernheim, “On the Potential of Neuroeconomics”; Krajbich and Dean, “How Can Neuroscience Inform Economics?”.

³⁴ Bernheim and Rangel, “Beyond Revealed Preference”; Bernheim, “On the Potential of Neuroeconomics.”

³⁵ Bernheim and Rangel, “Behavioral Public Economics”; Bernheim and Taubinsky, “Behavioral Public Economics”; Sunstein, “Behavioral Welfare Economics.”

Normative transparency condition: any paper or policy that draws welfare conclusions from neural data should explicitly state the welfare criterion, justify it, and explain how neural and behavioural data identify contexts in which welfare and implemented choice diverge.

4.4 Policy relevance and institutional context

Finally, translating individual-level neural evidence into policy involves further steps. Policies operate at scale, under constraints of fairness, privacy and feasibility. Using neural data to justify, for example, targeted mandates or tax treatment of particular groups raises distributional and political concerns.³⁶

Policy context condition: any proposed use of neural evidence in welfare analysis should specify how such evidence would be operationalised in practice, whether classification based on neural traits is legitimate, and how uncertainty and heterogeneity will be handled.

These four conditions form the core of the Neuroeconomic Welfare Inference Checklist (Section 6).

5 Case domains

To make the framework concrete, I sketch three domains where neural data have been proposed as inputs to welfare analysis.

5.1 Addiction and self-control

Neuroscience has transformed theories of addiction. Drug exposure alters dopaminergic signalling, cue reactivity and plasticity; drug-related cues acquire strong motivational salience; and many studies report reduced prefrontal control in addicted individuals.³⁷

A standard normative argument runs as follows: at least some consumption episodes reflect cue-triggered mistakes rather than genuine preferences; therefore, paternalistic policies such as sin taxes, mandatory cooling-off periods or access to commitment devices, can be justified as tools to restore an individual's own long-run interests.

³⁶ Bernheim and Rangel, “Behavioral Public Economics”; Bernheim and Taubinsky, “Behavioral Public Economics.”

³⁷ Bernheim and Rangel, “Addiction and Cue-Triggered Decision Processes.”

In the formal framework, addiction can be represented as a regime in which, in “cue” states s^{cue} , the effective value function and policy are distorted:

$$V_{\theta}^{\text{cue}}(s, a) = V_{\theta}^{\text{base}}(s, a) + \kappa C(s, a),$$

where $C(s, a)$ captures learned cue reactivity and κ is large in cue states, while λ in the dual-self utility shrinks. Neural evidence in favour of this picture includes enhanced striatal and amygdala responses to drug cues and attenuated dorsolateral prefrontal activation during drug-seeking.

Within the checklist:

- Measurement validity: imaging and electrophysiology support distinct cue-driven and baseline regimes.
- Model identifiability: alternative interpretations (for example, adaptive re-valuation given changed internal states) must be considered.
- Normative transparency: privileging baseline valuations (higher λ , lower κ) as welfare-relevant is a normative choice.

Neural data can strengthen the case that some episodes are mistakes in the sense of Bernheim and Rangel’s cue-triggered model, but they do not by themselves determine whether, for example, abstinence is universally welfare-superior to managed use. Those conclusions require ethical and economic argument.

5.2 Neuromarketing and behavioural public policy

Neuroeconomics and behavioural insights now inform both public policy (“nudges”) and commercial practice (neuromarketing, personalised advertising).³⁸

Public-sector applications include using fMRI and EEG to identify message framings that better engage brain networks linked to self-control or pro-social motivation, for example in retirement saving or health campaigns. Private-sector applications include optimising adverts to maximise reward-related activation and approach tendencies.

³⁸ Hsu and Yoon, “Neuroscience of Consumer Choice”; Nikki Leeuwis, Tom van Bommel, and Maryam Alimardani, “A Framework for Application of Consumer Neuroscience in Pro-Environmental Behavior Change Interventions,” *Frontiers in Human Neuroscience* 16 (September 2022), <https://doi.org/10.3389/fnhum.2022.886600>; Alsharif and Isa, “Revolutionizing Consumer Insights”; Gupta, Kapoor, and Verma, “Neuro-Insights.”

In welfare terms, there is an asymmetry. The same neural insights that support welfare-enhancing nudges can be used to exploit biases. Recommendation algorithms that maximise click-through or engagement effectively reshape the reward landscape in which human neural RL systems operate. Over time, they may lead to over-consumption of immediately rewarding but harmful content, potentially narrowing preferences and weakening self-regulation.³⁹

In the formal vocabulary, platforms choose the reward function r and state-to-cue mappings; AI systems learn policies that maximise platform objectives; human value functions V_θ adapt. Neural evidence, for example that a given design induces heightened ventral striatal activation and more habitual responding, shows that the platform is effectively tuning the human learning system.

Neural data support welfare claims when combined with:

- an explicit welfare criterion (e.g. long-run financial security, mental health);
- a model of how design changes shift behaviour away from that criterion; and
- normative transparency about which neural responses are taken as welfare-relevant.

Without this, there is a risk that high reward-system activation is mis-interpreted as high welfare.

5.3 Environmental and energy policy

Environmental and energy policy is another emerging domain for “environmental neuroeconomics”.⁴⁰ Neuroimaging studies examine how people respond to different framings of climate risks or energy savings messages, with the aim of designing communications that close the energy-efficiency gap.

Neural responses, for example, strong activation in regions linked to affective value or social cognition when viewing climate impacts, are sometimes interpreted as evidence of latent concern that is under-expressed in choices. Policy advocates may argue that such findings justify stronger environmental protection than stated preferences alone would suggest.

The framework suggests caution:

³⁹ Hsu and Yoon, “Neuroscience of Consumer Choice”; Alsharif and Isa, “Revolutionizing Consumer Insights.”

⁴⁰ Nik Sawe, “Adapting Neuroeconomics for Environmental and Energy Policy,” *Behavioural Public Policy* 3, no. 1 (2018): 17–36, <https://doi.org/10.1017/bpp.2018.2>;

Nik Sawe and Kiran Chawla, “Environmental Neuroeconomics: How Neuroscience Can Inform Our Understanding of Human Responses to Climate Change,” *Current Opinion in Behavioral Sciences* 42 (December 2021): 147–154, <https://doi.org/10.1016/j.cobeha.2021.08.002>;

Leeuwis, van Bommel, and Alimardani, “Framework for Application of Consumer Neuroscience”

- Neural data can show that certain messages more strongly engage circuits associated with valuation, social norm processing or moral concern.
- They can reveal heterogeneity: some individuals show little neural engagement even when they express pro-environmental attitudes, or vice versa.⁴¹

But moving from “people show neural responses to environmental loss” to “policies should weight environmental quality more heavily than current choices imply” requires a normative step. The welfare criterion must treat certain neural responses, e.g. guilt, awe, moral outrage, as part of what makes outcomes better or worse for individuals.

Again, the checklist applies: task validity (do the lab stimuli correspond to real policy trade-offs), model identification (are alternative explanations for neural activity considered) and normative transparency (why are those responses taken as welfare relevant).

6 A Neuroeconomic Welfare Inference Checklist and implications for NeuroAI

Drawing the argument together, *Table 1* summarises a practical checklist for using neural data in welfare analysis.

Step	Task	Details
1	Define the welfare criterion U	Specify whether welfare is identified with long-run preferences, experienced utility, or another standard. Justify this choice in the policy context.
2	Specify the computational model	State the decision and learning model (for example, reinforcement learning with parameters θ, ϕ). Explain which behavioural patterns it captures.
3	Validate neural encodings	Provide evidence that neural variables encode specific computational quantities (for example, δ_t, V_θ), and assess the robustness and specificity of these encodings across tasks and individuals.
4	Assess model identifiability	Consider alternative models that fit the same behavioural and neural data. Discuss how, if at all, neural evidence distinguishes between these models and their associated welfare analyses.
5	Locate welfare-relevant divergences	Identify states or contexts where implemented choice (given V_θ and policy π_ϕ) diverges from U . Use neural evidence to characterise mechanisms behind such divergences, not to define U itself.

⁴¹ Sawe and Chawla, “Environmental Neuroeconomics”; Leeuwis, van Bommel, and Alimardani, “Framework for Application of Consumer Neuroscience.”

6	Analyse policy implementation	Explain how neural-based classifications or interventions would be implemented in practice. Address fairness, privacy, legal and broader political constraints on using neural data in policy.
---	-------------------------------	--

Table 1. Neuroeconomic Welfare Inference Checklist

For NeuroAI, the checklist has two main implications.

First, it supports modelling both brains and artificial agents as value-learning systems whose internal variables can be compared. When AI systems are designed to assist human decision-makers, designers can adopt an explicit welfare criterion U and treat human neural and behavioural data as evidence about V and π . The checklist disciplines how that evidence is used.

Second, it guards against conflating optimisation of neural activation with optimisation of welfare. An AI system that maximises expected striatal activation might push users toward behaviours with high short-run reward but low long-run welfare. Welfare-aligned NeuroAI should instead optimise a defensible U , with neural data serving only as one class of inputs and constraints.

For human-inspired AI, the checklist can be read as a set of design principles for brain-inspired value-learning agents. Architectures that borrow from dopaminergic TD learning or cortical–striatal actor–critic loops (Figures 1–2) should treat their internal reward signals and value functions as computational variables, not welfare criteria. A designer who wants such an agent to act in the interests of a human user must specify and defend a separate welfare model U and justify how the agent’s learning objective approximates it. The checklist then governs how neural data (for example, user-specific prediction-error patterns) may or may not be used to adapt that objective. In this way, the paper contributes not only to welfare analysis of human behaviour but also to the conceptual foundations of human-inspired, welfare-aware AI, bridging brain-inspired RL architectures and welfare-aware AI design and filling a gap between technical work on value learning and the normative questions that arise when such systems are deployed in policy-relevant domains.

7 Conclusion

Neuroeconomics is sometimes presented as if it could directly deliver welfare-relevant quantities by “reading” the brain. This paper argues that such expectations are misplaced. Neural data are valuable, but they inform welfare only through computational and normative lenses.

By formalising the path from neural signals, through reinforcement-learning and dual-self models, to welfare criteria, I have derived conditions under which neural data can legitimately inform policy. These conditions emphasise representational validity, model identification, normative transparency and institutional context. Case sketches in addiction, neuromarketing and environmental policy show both the promise and limits of current practice.

The alignment with NeuroAI is natural: RL, prediction-error coding and distributional value representation provide a shared language for brains and machines.⁴² Yet welfare remains a separate layer that requires philosophical and economic argument. NeuroAI systems that increasingly shape human decision environments should incorporate explicit welfare models and adhere to something like the Neuroeconomic Welfare Inference Checklist, rather than treating neural correlates of reward as welfare endpoints.

Neural data can refine our models of how people choose and illuminate mechanisms behind departures from rationality. Embedded in transparent welfare frameworks, such models can support more precise and humane policy design. What neural data cannot do is relieve us of the need to argue about what counts as a good life, and whose judgement should prevail when brain signals and choices diverge, that is, whether the decision-maker is a human or a human-inspired AI system trained on those signals.⁴³

Future work can proceed in three directions. Empirically, the checklist can be used to audit concrete policy uses of neural evidence in domains such as addiction treatment, personalised advertising or climate communication, testing which of the conditions are actually met. Methodologically, the formal framework can be extended to richer value learning models, including inverse reinforcement learning and multi agent settings, to analyse how welfare claims arise when AI systems infer or influence human reward functions. Normatively, the account can be integrated with broader debates in political philosophy and AI ethics about legitimacy and authority in welfare judgements, especially when human decisions are mediated by human inspired AI systems.

⁴² Doya, “Reinforcement Learning”; Dabney et al., “A Distributional Code for Value”; Sadeh and Clopath, “Emergence of NeuroAI.”

⁴³ Bernheim, “On the Potential of Neuroeconomics”; Krajbich and Dean, “How Can Neuroscience Inform Economics?”.

Bibliography

- Alsharif, Ahmed H, and Salmi Mohd Isa. 2024. "Revolutionizing Consumer Insights: The Impact of FMRI in Neuromarketing Research." *Future Business Journal* 10 (1). <https://doi.org/10.1186/s43093-024-00371-z>.
- Bernheim, B. Douglas. 2009. "On the Potential of Neuroeconomics: A Critical (but Hopeful) Appraisal." *American Economic Journal: Microeconomics* 1 (2): 1–41. <https://doi.org/10.1257/mic.1.2.1>.
- Bernheim, B. Douglas, and Antonio Rangel. 2004. "Addiction and Cue-Triggered Decision Processes." *American Economic Review* 94 (5): 1558–90. <https://doi.org/10.1257/0002828043052222>.
- . 2005. "Behavioral Public Economics: Welfare and Policy Analysis with Non-Standard Decision-Makers." *National Bureau of Economic Research*, August. <https://doi.org/10.3386/w11518>.
- . 2009. "Beyond Revealed Preference: Choice-Theoretic Foundations for Behavioral Welfare Economics." *The Quarterly Journal of Economics* 124 (1): 51–104. <https://doi.org/10.1162/qjec.2009.124.1.51>.
- Bernheim, B. Douglas, and Dmitry Taubinsky. 2018. "Behavioral Public Economics." Ssrn.com. July 2018. https://papers.ssrn.com/sol3/papers.cfm?abstract_id=3214374.
- Camerer, Colin, George Loewenstein, and Drazen Prelec. 2005. "Neuroeconomics: How Neuroscience Can Inform Economics." *Journal of Economic Literature* 43 (1): 9–64. <https://doi.org/10.1257/0022051053737843>.
- Dabney, Will, Zeb Kurth-Nelson, Naoshige Uchida, Clara Kwon Starkweather, Demis Hassabis, Rémi Munos, and Matthew Botvinick. 2020. "A Distributional Code for Value in Dopamine-Based Reinforcement Learning." *Nature* 577 (7792): 671–75. <https://doi.org/10.1038/s41586-019-1924-6>.
- Doya, Kenji. 2007. "Reinforcement Learning: Computational Theory and Biological Mechanisms." *HFSP Journal* 1 (1): 30–40. <https://doi.org/10.2976/1.2732246>.
- Eckstein, Maria K., Sarah L. Master, Ronald E. Dahl, Linda Wilbrecht, and Anne G.E. Collins. 2022. "Reinforcement Learning and Bayesian Inference Provide Complementary Models for the Unique Advantage of Adolescents in Stochastic Reversal." *Developmental Cognitive Neuroscience* 55 (June): 101106. <https://doi.org/10.1016/j.dcn.2022.101106>.
- Gupta, Raveena, Anuj Pal Kapoor, and Harsh V Verma. 2025. "Neuro-Insights: A Systematic Review of Neuromarketing Perspectives across Consumer Buying Stages." *Frontiers in Neuroergonomics* 6 (July). <https://doi.org/10.3389/fnrgo.2025.1542847>.
- Hsu, Ming, and Carolyn Yoon. 2015. "The Neuroscience of Consumer Choice." *Current Opinion in Behavioral Sciences* 5 (October): 116–21. <https://doi.org/10.1016/j.cobeha.2015.09.005>.

- Krajovich, Ian, and Mark Dean. 2015. "How Can Neuroscience Inform Economics?" *Current Opinion in Behavioral Sciences* 5 (October): 51–57. <https://doi.org/10.1016/j.cobeha.2015.07.005>.
- Leeuwis, Nikki, Tom van Bommel, and Maryam Alimardani. 2022. "A Framework for Application of Consumer Neuroscience in Pro-Environmental Behavior Change Interventions." *Frontiers in Human Neuroscience* 16 (September). <https://doi.org/10.3389/fnhum.2022.886600>.
- Lloyd, Kevin, and Peter Dayan. 2015. "Tamping Ramping: Algorithmic, Implementational, and Computational Explanations of Phasic Dopamine Signals in the Accumbens." Edited by Olaf Sporns. *PLOS Computational Biology* 11 (12): e1004622. <https://doi.org/10.1371/journal.pcbi.1004622>.
- Potjans, Wiebke, Markus Diesmann, and Abigail Morrison. 2011. "An Imperfect Dopaminergic Error Signal Can Drive Temporal-Difference Learning." Edited by Tim Behrens. *PLoS Computational Biology* 7 (5): e1001133. <https://doi.org/10.1371/journal.pcbi.1001133>.
- Rangel, Antonio, Colin Camerer, and P. Read Montague. 2008. "A Framework for Studying the Neurobiology of Value-Based Decision Making." *Nature Reviews Neuroscience* 9 (7): 545–56. <https://doi.org/10.1038/nrn2357>.
- Sadeh, Sadra, and Claudia Clopath. 2025. "The Emergence of NeuroAI: Bridging Neuroscience and Artificial Intelligence." *Nature Reviews. Neuroscience* 26 (10): 583–84. <https://doi.org/10.1038/s41583-025-00954-x>.
- Sawe, Nik. 2018. "Adapting Neuroeconomics for Environmental and Energy Policy." *Behavioural Public Policy* 3 (1): 17–36. <https://doi.org/10.1017/bpp.2018.2>.
- Sawe, Nik, and Kiran Chawla. 2021. "Environmental Neuroeconomics: How Neuroscience Can Inform Our Understanding of Human Responses to Climate Change." *Current Opinion in Behavioral Sciences* 42 (December): 147–54. <https://doi.org/10.1016/j.cobeha.2021.08.002>.
- Schultz, W., P. Dayan, and P. R. Montague. 1997. "A Neural Substrate of Prediction and Reward." *Science* 275 (5306): 1593–99. <https://doi.org/10.1126/science.275.5306.1593>.
- Sunstein, Cass R. 2020. "Behavioral Welfare Economics." *Journal of Benefit-Cost Analysis* 11 (2): 196–220. <https://doi.org/10.1017/bca.2020.14>.

Appendix A. Formal model (collected)

For convenience, the main formal elements are collected here.

1. Decision environment

- States: $s_t \in S$
- Actions: $a_t \in A$
- Rewards: $r_t \in \mathbb{R}$
- Dynamics: $P(s_{t+1} \mid s_t, a_t)$

2. Value and policy

- State value: $V_\theta(s)$
- Action value: $Q_\theta(s, a)$
- Policy:

$$\pi_\phi(a \mid s_t) = \frac{\exp(\beta Q_\theta(s_t, a))}{\sum_{a'} \exp(\beta Q_\theta(s_t, a'))}$$

3. Learning rule (TD)

$$\delta_t = r_t + \gamma V_\theta(s_{t+1}) - V_\theta(s_t)$$

$$\theta_{t+1} = \theta_t + \alpha \delta_t \nabla_\theta V_\theta(s_t)$$

4. Neural encoding

- Prediction-error encoding:

$$n_t = f(\delta_t) + \varepsilon_t$$

- Value encoding:

$$n_t^{\text{vmPFC}} = g(V_\theta(s_t)) + \eta_t$$

5. Dual-self welfare representation

- Long-run utility: $u^L(a)$
- Short-run utility: $u^S(a)$
- Implemented instant utility:

$$U^{\text{impl}}(a) = \lambda u^L(a) + (1 - \lambda) u^S(a),$$

with λ state-dependent (for example, reduced in cue states for addiction).

6. Addiction-style cue distortion

$$V_{\theta}^{\text{cue}}(s, a) = V_{\theta}^{\text{base}}(s, a) + \kappa C(s, a),$$

where $C(s, a)$ captures cue-associated value and κ the strength of cue distortion.

These components can be instantiated as computational models, fit to behaviour and neural data, and then combined with explicit welfare criteria U to assess policy proposals.