

**Irrational decisions reflect robustness constraints on value computations
implemented by orbitofrontal circuits**

Juliette Bénou¹, Mathias Pessiglione¹, Fabien Vinckier¹, Jean Daunizeau¹

¹ Paris Brain Institute - ICM, Hôpital de la Pitié Salpêtrière, Paris, France

Address for correspondence:

Jean Daunizeau

Motivation, Brain and Behavior Group

Paris Brain Institute (ICM)

47, bd de l'Hôpital, 75013, Paris, France.

Tel: +33 1 57 27 43 26

Mail: jean.daunizeau@icm-institute.org

Web: <https://sites.google.com/site/motivationbrainbehavior>

Abstract:

1 Making good decisions is essential for survival and success, yet humans and animals
2 often exhibit perplexing irrational decision-making whose biological origin remains poorly
3 understood. Recent empirical and computational work suggests that altered computations in
4 perceptual, motor and memory systems in the brain may arise from informational, metabolic or
5 robustness constraints on their internal connectivity structure. However, whether and how such
6 neurobiological constraints may have molded the architecture of decision systems (such as the
7 orbitofrontal cortex) and eventually distorted decision-relevant computations, remains largely
8 unknown. We first train cohorts of artificial neural nets to perform ten variants of rational
9 decision-relevant computations. Those variants that operate under a specific option-encoding
10 format exhibit most of the electrophysiological coding properties observed in orbitofrontal
11 neurons of monkeys making decisions under risk. We then distort these neural nets' internal
12 wiring to reproduce monkeys' irrational choices. This induces deterministic spillover
13 interferences in decision-relevant computations that generalize across individuals, at both the
14 behavioral and neural level. Importantly, although irrational nets do not seem to bring
15 informational or metabolic benefits, they display enhanced tolerance to damage and noise when
16 compared to their rational counterparts. This suggests that some forms of irrational behavior may
17 be the incidental outcome of distal evolutionary pressure on the robustness of orbitofrontal
18 circuits.

Text:

19 **Introduction**

20 People and animals arguably act, in some circumstances, against their own interest. Why
21 does irrational behavior persist, despite its potential costs to survival and fitness? Standard
22 decision theory posits that rational decisions rely on estimating and comparing the expected
23 value of each available alternative option in the choice set. Thus, irrational behavior may emerge
24 from the covert mechanisms through which the brain constructs, represents, maintains or
25 compares option values. Decades of work in human and non-human primates show that these
26 computational processes involve a specific subset of brain systems, including – but not limited
27 to – orbitofrontal (OFC), ventromedial (vmPFC) and dorsolateral prefrontal (dlPFC), as well as
28 anterior cingulate (ACC), cortices¹⁻⁹. While the relative contribution of these subsystems is not
29 well understood, a robust finding across studies is that ventromedial and/or orbitofrontal neurons
30 encode value, regardless of whether subjects are engaged in explicit decision-making or in the
31 subjective evaluation of single options¹⁰⁻¹³. Accordingly, neuropsychological studies of brain-
32 damaged patients demonstrate that lesions to the orbitofrontal cortex induce irrational value-
33 based decisions without impairing other types of high-level cognitive processes¹⁴⁻¹⁶. This means
34 that the effective rationality of decisions hinges on the integrity of OFC circuits. But even in the
35 absence of anatomical lesions, value processing in the OFC is known to exhibit systematic
36 distortions, which typically lead to overt irrational behavior. For example, value coding in the OFC
37 is modulated by its pre-stimulus baseline activity^{17,18}, adapts to the recent range of option
38 values¹⁹⁻²¹, and depends on whether a given option is the status-quo alternative²² or is currently
39 attended²³. Given the behavioral biases that ensue, these results suggest that OFC circuits are
40 organized in such a way that they process value-related information in a moderately, yet

41 pervasively, suboptimal manner. In turn, this raises the basic question of why suboptimal value
42 computations prevail in OFC circuits.

43 A common view is that seemingly irrational decisions may yield compensatory behavioral
44 benefits. For example, inattention, overconfidence or optimism biases can be shown to be
45 advantageous under the right circumstances^{24–26}. However, an alternative assumption is that
46 evolutionary pressure eventually selected for value computations that are “rational enough”,
47 given the constraints that may act at the neurobiological level^{27,28}. A prominent example is the
48 energetic budget of neural circuits, which encompasses both physiological maintenance and
49 activity-dependent firing costs^{29,30}. The former costs typically restrains the total number of
50 neurons in the brain^{31,32}, which promotes neural coding strategies that minimize redundancy
51 and/or maximize information transfer rates^{33–36}. Interestingly, variants of such mechanisms
52 explain value range adaptation effects in the OFC and the irrational behavior that ensues^{20,28,37}.
53 The latter costs induce metabolic budget constraints that are demonstrably tight. For example,
54 the mitochondrial metabolic supply of neurons is actively restricted at the expense of circuit-level
55 computational fidelity³⁸, and a scarcity of external resources (e.g., food) eventually results in
56 impaired neural processing³⁹. This supports the idea that the brain has evolved so-called energy-
57 efficient neural coding strategies that trade off computational precision for metabolic costs^{40,41}.
58 But theoretical work also emphasizes other types of tradeoffs that arise from demands on the
59 robustness or fault-tolerance of neural circuits. A widely debated notion is that neural circuits
60 must maintain their excitatory-inhibitory balance to ensure stability and/or homeostasis⁴².
61 Another possibility, which is pervasive in biological systems, is the need to minimize vulnerability
62 to localized damage^{43,44}. Although direct empirical evidence for such a constraint on neural
63 circuits is scarce, recent work indicates that frontal circuits that subtend, e.g. motor control and
64 working memory, achieve tolerance to damage through architectural redundancy^{45–47}. This is
65 important because redundant systems are notoriously inefficient, from both an informational
66 and energetic perspective^{41,48}. In other words, OFC circuits may have evolved under multiple yet

67 competing neurobiological constraints, whose impact on value-based decisions remains largely
68 unexplored.

69 In principle, the fidelity of OFC value computations, its energy budget, and its biological
70 robustness share a common determinant: its internal circuit wiring. From a neural net
71 perspective, any developmental or evolutionary pressure of the sort discussed above will
72 ultimately shape the architecture of OFC circuits in ways that distort value computations and
73 compromise decision rationality^{20,49}. This is, in fact, trivially observed in artificial neural net
74 models of the OFC trained to perform candidate value computations while complying with these
75 constraints (see Fig. 1). Critically however, the form of irrational behavior that emerges depends
76 on both the nature of the constraint and the specific value computations the OFC is assumed to
77 perform. This is because a given type of value computation requires a tailored neural net wiring,
78 whose native compliance with the above constraints is largely arbitrary. Thus, the issues of
79 identifying the nature of OFC value computations, the mechanisms through which they give rise
80 to irrational behavior, and the biological constraints that may have molded the architecture of
81 OFC circuits, are closely intertwined. In this work, we approach the problem from a
82 computational perspective.

83 We consider the paradigmatic case of binary decisions under risk, where the choice set
84 consists of two alternatives, each defined by the probability and magnitude of prospective
85 rewards. In principle, rational behavior amounts to choosing the option with the highest expected
86 reward. Extensive electrophysiological recordings from OFC neurons are available in macaque
87 monkeys performing this type of decision-making task. Here, we reanalyze an existing dataset in
88 which decision cues (i.e. option-specific reward magnitude or probability) are revealed one at a
89 time. This design provides a unique empirical estimate of the dynamics of information content in
90 the OFC as value computations unfold over within-decision time⁵⁰. In line with previous literature,
91 we distinguish between two broad types of value computations: value *synthesis* and value

92 *comparison*. The former implies that the OFC progressively integrates decision cues to compute
93 the value of both options, which can be concurrently read out on possibly orthogonal activity
94 subspaces of OFC neural ensembles^{49,51,52}. The latter reduces to directly updating the value
95 difference between the two options as a new decision cue becomes available^{53,54}. Both value
96 synthesis and value comparison can be implemented using one of five distinct neural encoding
97 formats, which vary according to how cue inputs and value outputs are framed (see Results)^{22,23}.
98 Together, this yields a total of ten candidate scenarios regarding OFC value computations.

99 We first train recurrent neural nets (RNNs) to perform each candidate value computation
100 in a rational manner, given arbitrary decision cue sequences. We then identify which, among the
101 candidate types of value computations, yield accurate RNN models of the OFC. To do so, we
102 compare the full set of recorded OFC neural responses with the activity patterns of simulated
103 RNNs exposed to the same decision trials as those experienced by the monkeys. As we will see,
104 this eventually selects two specific types of value computations, which effectively are idealized
105 models of OFC networks that would have evolved without any metabolic or robustness
106 constraint. We then distort the internal connectivity of these networks to reproduce monkeys'
107 irrational choices in the task (about 20% of all choices). As we will show, these distorted RNNs
108 make behavioral and neural predictions that generalize across monkeys. Finally, we compare
109 rational and irrational RNN models of the OFC, in terms of their compliance with energetic and
110 robustness constraints. This enables us to identify which neurobiological constraint may have
111 shaped OFC computations.

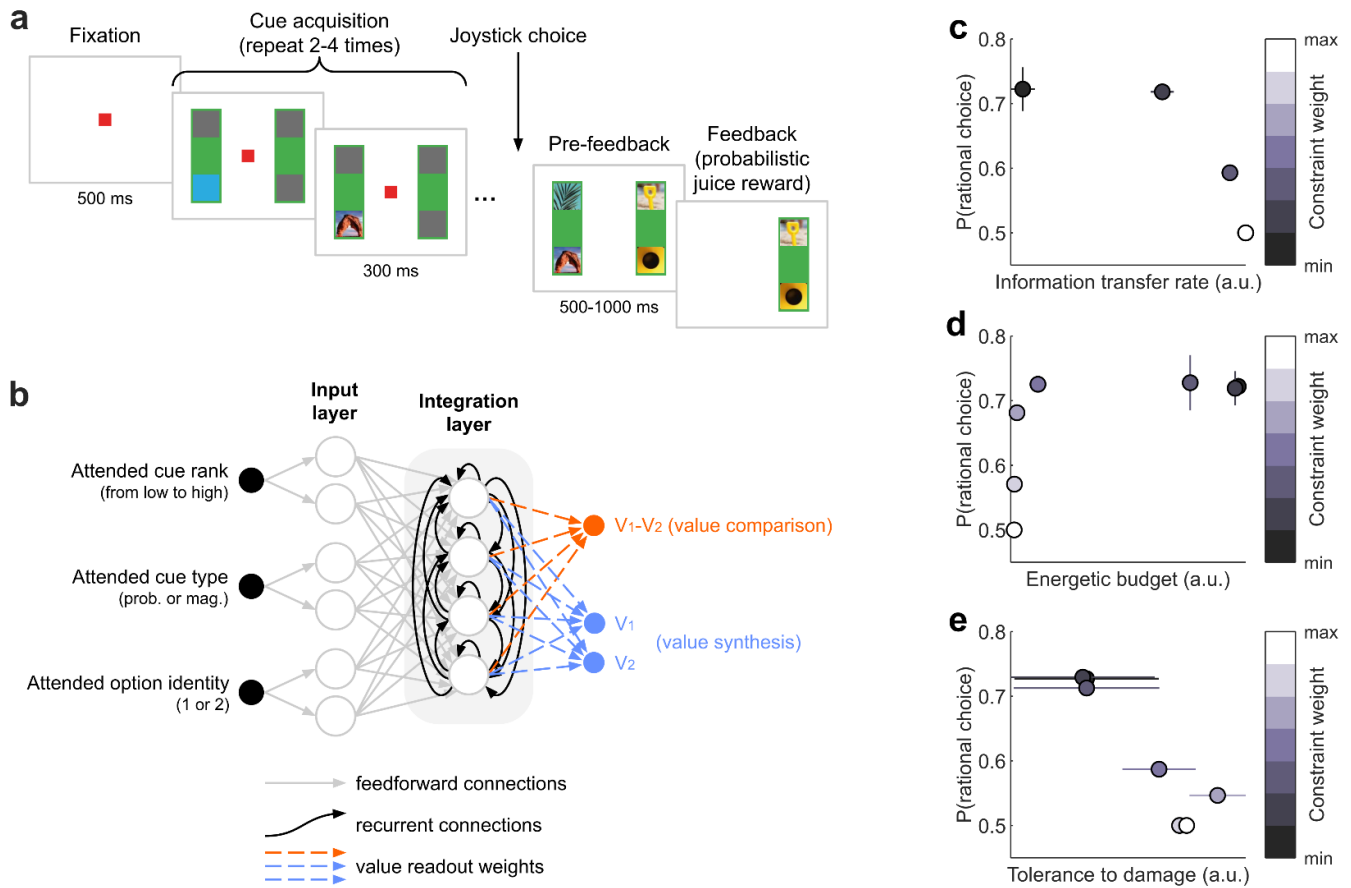
112 **Results**

113 We took advantage of an open dataset of single unit activity recordings from the OFC, the
114 dlPFC and the ACC of two macaque monkeys (n=189, 135 and 183 neurons respectively) engaged

115 in value-based decision-making (22,618 trials in total)⁵⁰. The task design is summarized on Fig.
116 1a (see also the Methods section). At each trial, monkeys chose between two options presented
117 on the left and right sides of a screen, each defined by two cues (probability and prospective
118 amount of a rewarding juice) that are revealed one at a time and could commit to a decision
119 without necessarily sampling all cues. The subjective value profiles estimated from monkey
120 choices (see Methods) indicate that they integrate both currently and previously attended cues,
121 and are highly correlated with the rational value profile (Pearson correlation, Monkey F: $\rho = 0.95$;
122 Monkey M: $\rho = 0.93$; all $p < 10^{-15}$) with no significant correlation difference between monkeys
123 (Steiger's test, $p = 0.19$).

124 ***Identifying value computations in the OFC***

125 To begin with, we aimed to identify legitimate models of an idealized OFC network, which
126 would have evolved without any metabolic or robustness constraint (and would thus perform
127 rational value computations). To do this, we adopt a normative approach that obviates the need
128 for empirical data in training models. In line with recent empirical work, we hypothesized that the
129 OFC may implement one of two candidate decision-relevant computations: (i) computing the
130 value of both options independently^{51,55} (“value synthesis”) or (ii) computing the difference
131 between option values^{9,56} (“value comparison”). Each scenario can be implemented using
132 recurrent artificial neural networks (RNNs), which operate under the exact same conditions as
133 monkeys in the task. In particular, RNNs access cues sequentially and in an encoding format that
134 specifies attribute type and rank as well as option identity (see below). At each cue onset, these
135 inputs are sent to a first hidden layer (cue-encoding), whose units feed their output forward to a
136 second hidden layer (cue-integration), from which the RNN's value outputs are readout (see Fig.
137 1b and Methods). The integration layer relies on internal recurrent connections to combine
138 currently and previously sampled cues, and progressively update its ongoing within-trial



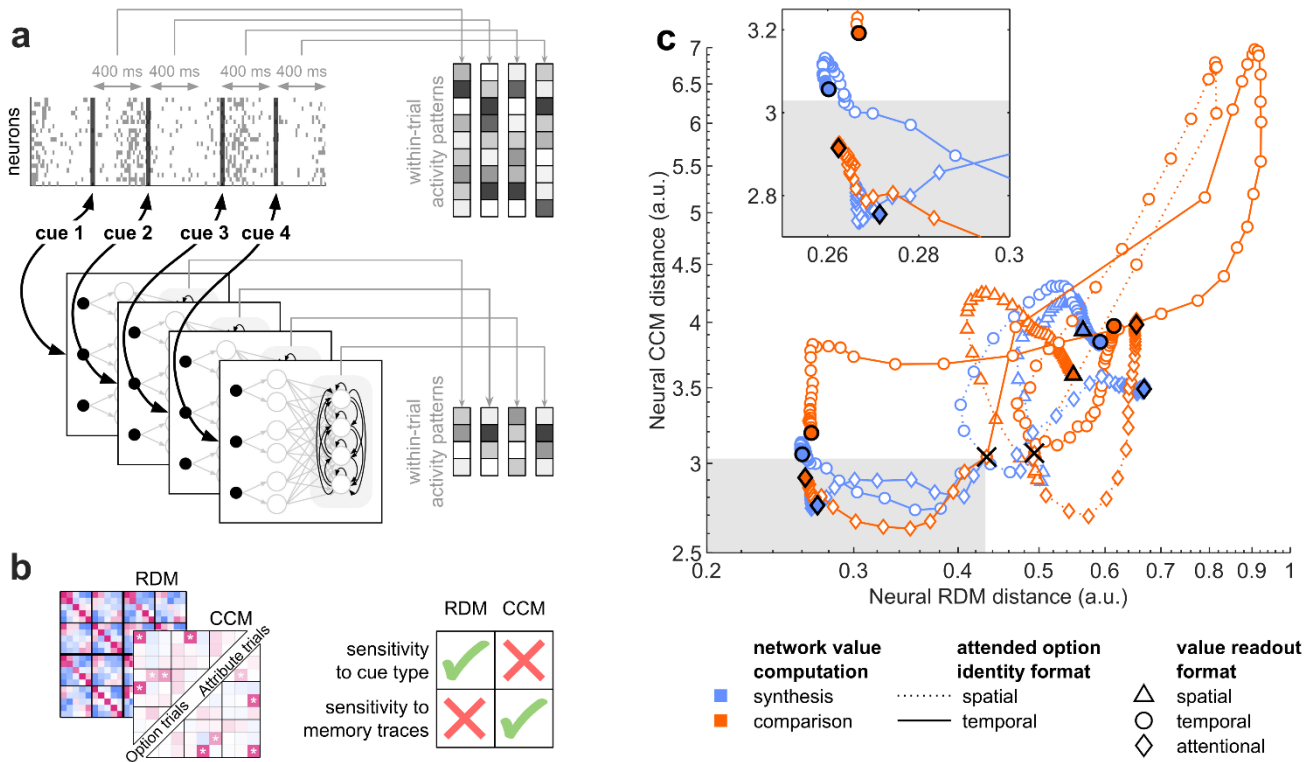
139 **Fig. 1 | Decision task, neural net design, option values and impact of constraints.** **a**, Task design. Adapted from
 140 Hunt et al., 2018⁵⁰. At each trial, monkeys choose between two options presented on the left and right sides of a screen,
 141 each defined by the probability and prospective amount of a rewarding juice. Each decision cue (representing either
 142 the probability or the magnitude of the currently attended option) appears sequentially and then disappears. Monkeys
 143 can commit to a decision at any point after the second cue without necessarily sampling the remaining cues and are
 144 free to decide which cue to sample if they decide to continue the trial. **b**, Neural net architecture (see Methods). At
 145 each cue onset, the neural net's inputs (back dots) encode the currently attended cue in terms of its three raw features
 146 (i.e., cue rank, cue type and option identity), while its outputs are the neural net's current estimate of option values
 147 (value synthesis, blue dots) or value difference (value comparison, orange dot). Feedforward connections convey input
 148 features to a layer of input-specific units, which send their outputs to an input-integration layer. Importantly, the latter
 149 is equipped with recurrent connections that carry reentrant value computations elicited from previously attended
 150 cues. Note that both the identity of the attended option and value outputs can be encoded in distinct formats (see
 151 Methods). **c**, **d** and **e**, Neural net training under multiple constraints. Neural nets can be trained to optimize a
 152 compromise between rational value computations, on the one hand, and information transfer rate (**c**), energetic budget
 153 (**d**) or tolerance to damage (**e**), on the other hand (see Methods). Each dot depicts the average rate of rational choices
 154 (y-axis) and constraint adequacy (x-axis) for a given constraint compliance weight (color scale). Error bars indicate

155 standard error. One can see that imposing stronger compliance with such constraints tends to compromise decision
156 rationality, because the underlying network wiring eventually alters value computations.

157 computations⁴⁹. Importantly, both value synthesis and comparison also require specifying how
158 options are identified, which is debated in the existing literature. The OFC may do so based on,
159 e.g., spatial location⁵⁷ (left vs. right), temporal order⁴⁹ (first vs second), or attentional focus²³
160 (attended vs. unattended). In principle, both OFC inputs (decision cues) and outputs (option
161 values or value difference) may encode option identity in a different format, irrespective of
162 whether the OFC operates value synthesis or comparison. This resulted in ten cohorts of RNNs
163 (two types of value computations combined with five input-output format variations), which we
164 train on rational option values (see Methods). This is not trivial, since it requires RNNs to maintain
165 a memory trace of previously attended cues, while remaining invariant to the order in which cues
166 are presented. Importantly, each cohort gathers 1,000 RNN instances that sample the manifold
167 of admissible wirings following random weight initializations.

168 Once trained, these RNN models can be used to simulate neural population dynamics
169 that implement rational value computations during decision trials. If triggered with the
170 sequences of cues that monkeys attend in the decision task, this provides value readouts and
171 activity patterns that can be compared to observed choices and OFC recordings (see Fig. 2).
172 When probed at monkeys' decision onset time, these rational RNNs correctly predict
173 $79\% \pm 3 \times 10^{-3}$ (mean \pm standard error) of choices (monkey F: $77\% \pm 4 \times 10^{-3}$, monkey M:
174 $80\% \pm 4 \times 10^{-3}$). Crucially, although all rational RNNs yield identical decisions, their internal
175 representations are different (see Supplementary Materials). We thus asked whether any of these
176 RNN cohorts capture key aspects of OFC neural representation geometry, despite not having
177 been exposed to neural recordings during training. To test this, we replicated the two types of
178 analysis conducted by Hunt and colleagues⁵⁰ on single units' recordings, which we also

179 performed on the RNNs' integration layer (see Fig. 2a). We first ran a representational similarity
180 analysis at first cue onset, building representational dissimilarity matrices (RDMs) by correlating
181 population activity vectors in response to all possible cues (see Methods). In brief, RDMs identify
182 which cue features elicit discriminable response patterns across neurons when only a single cue
183 is available. To track neural representation geometry at all stages of decision trials, we also
184 quantified whether and how inter-neuron differences in their sensitivity to current and past cue
185 ranks are preserved across cue onset times (cf. cross-correlation matrices or CCMs; see
186 Methods). One can think of RDMs and CCMs as two distinct summary statistics of the
187 representational geometry of distributed neural systems, with complementary properties (see
188 Fig. 2b). We then derived the two ensuing neural distance metrics by comparing OFC neurons
189 and RNN units at each stage of the training process (see Fig. 2c). Note that even untrained – i.e.
190 random – RNNs exhibit some degree of neural similarity with the OFC, because they respond to
191 value-relevant input cues. Untrained RNNs thus effectively provide the distribution of neural
192 distances under the null hypothesis. Now, when being trained to perform a specific value
193 computation, RNNs modify their representational geometry and hence their neural distance to
194 the OFC. In line with previous work on visual and language systems in the brain^{58,59}, we considered
195 that legitimate idealized RNN models of the OFC are those RNN cohorts that significantly
196 decrease both neural distance metrics as a result of training.



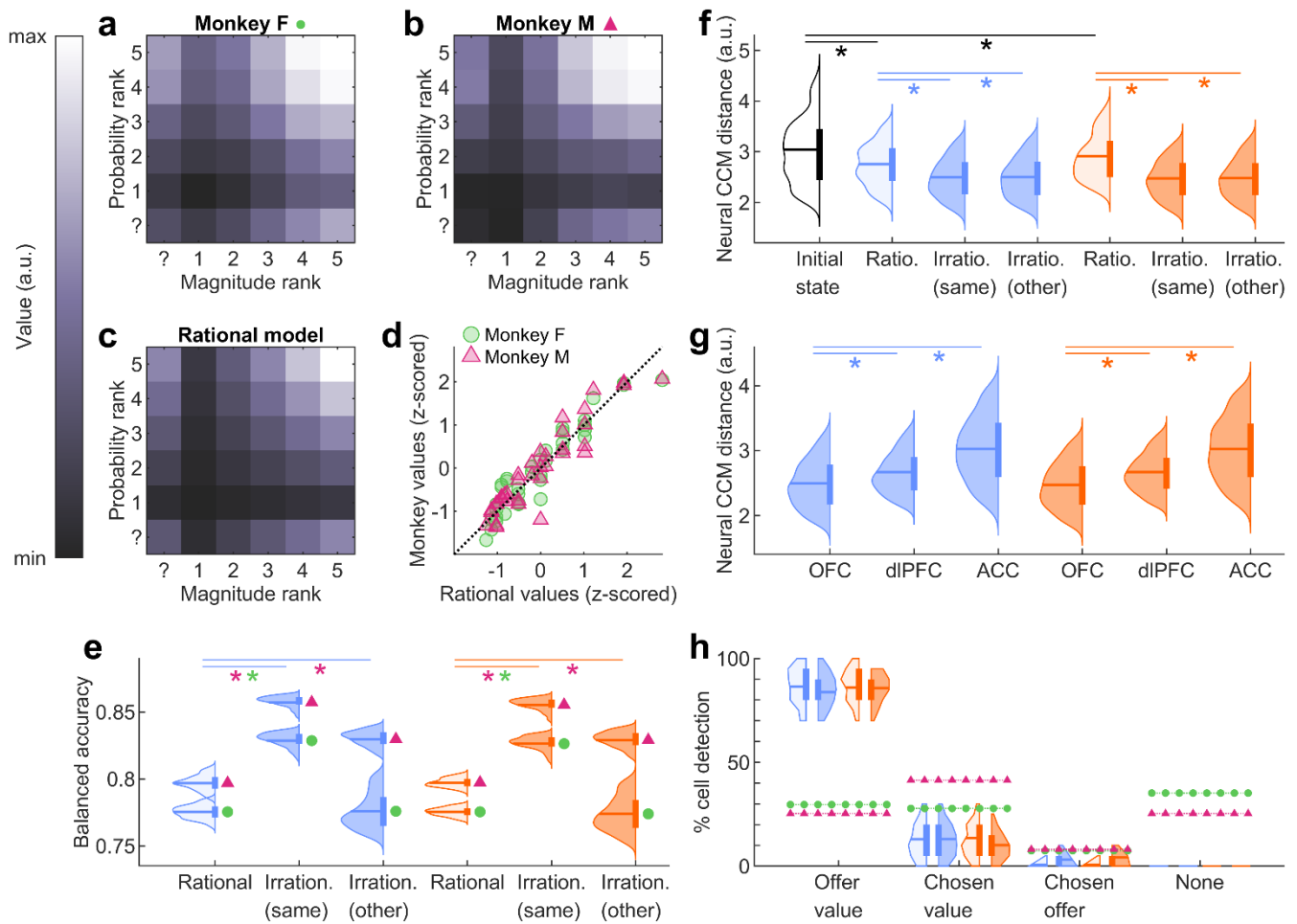
197 **Fig. 2 | Selection of candidate idealized RNN models of the OFC.** **a**, Schematic procedure for extracting within-trial
 198 activity patterns from OFC neural recordings and RNN models. At each trial, a sequence of cues is sampled until a
 199 choice is triggered. Within-trial activity patterns of OFC neurons are constructed as the average firing rate of each
 200 neuron, from 100 ms to 500 ms after each cue onset. Similarly, within-trial activity patterns of RNN models are the
 201 activation strengths of integration units in response to each cue. **b**, Summary of neural distance metrics for comparing
 202 activity patterns of OFC neurons and RNN models. RDMs quantify how dissimilar evoked activity patterns are for any
 203 pair of cues ($2 \times 2 \times 5 = 20$ possibilities at first cue onset). CCMs quantify the similarity of profiles of neural sensitivity to
 204 present and past cue ranks. Although CCM-based distance metrics are insensitive to cue type (probability or
 205 magnitude), they quantify potential internal memory traces about previously sampled cues. Full RDM and CCM
 206 summary statistics for all monkeys and brain regions can be eyeballed in the Supplementary Materials. **c**, Average
 207 neural distance trajectories between OFC and RNN cohorts (over 1,000 RNN instances), computed using either RDMs
 208 (x-axis) or CCMs (y-axis) metrics. Black crosses indicate the initial (random) state of RNN cohorts, black
 209 triangles/dots/diamonds denote their final rational state. Intermediary points show the neural distance at various
 210 stages of RNN training (from 0% to 100%, by steps of 2%), where color, line style and marker type indicate the type of
 211 computation (value synthesis vs value comparison), the identity format of the attended option (spatial vs temporal),
 212 and the value readout format (spatial vs temporal vs attentional), respectively. The gray shaded area indicates the two-
 213 dimensional region that falls below the 95% confidence interval of both pre-training neural distances.

214 Strikingly, all RNN cohorts that rely on the spatial (left versus right) encoding of the
215 attended option's identity tend to increase both neural distances to OFC neurons as training
216 unfolds. Among the remaining four cohorts, those that compute option values using a temporal
217 format (first versus second) eventually shorten their RDM distance, but at the cost of
218 compromising their CCM distance. Ultimately, only those two RNN cohorts that encode the
219 attended option identity using the temporal format (first versus second) while computing option
220 values in the attentional format (attended versus unattended) significantly shorten both neural
221 distances during training. These differ in terms of the type of value computation: one RNN cohort
222 performs value synthesis (CCM distance, paired two-sided t-test: $t(999) = 9.5$, $p < 10^{-15}$, Cohen's
223 $d = 0.30$; RDM distance: $t(999) = 44$, $p < 10^{-15}$, Cohen's $d = 1.39$), whereas the other performs value
224 comparison (CCM: $t(999) = 4.6$, $p = 5 \times 10^{-6}$, Cohen's $d = 0.14$; RDM: $t(999) = 48$, $p < 10^{-15}$, Cohen's
225 $d = 1.52$). Although we cannot yet arbitrate between these two scenarios, we have clearly
226 narrowed the set of plausible models of the idealized OFC network. In the remainder of the paper,
227 we focus on these two RNN cohorts (some extended results for all model variants are shown in
228 the Supplementary Materials).

229 At this point, we asked whether and how models of the idealized OFC network can be
230 modified to explain monkeys' irrational behavior. We thus retrained the selected rational RNNs to
231 predict monkeys' choices in the task, of which about 20% are irrational (Monkey F: $21\% \pm 3 \times 10^{-2}$;
232 Monkey M: $19\% \pm 3 \times 10^{-2}$). To preserve the interpretability of RNN value computations while
233 allowing for, e.g., nonlinear interferences and spill-over effects between attended cues, RNNs
234 were initialized with their trained rational weights, and retraining was restricted to recurrent
235 connections within the integration layer. When probed at monkeys' decision onset time, retrained
236 RNNs achieve an average choice-prediction accuracy of about $83\% \pm 1 \times 10^{-2}$ (monkey F: $80\% \pm$
237 5×10^{-2} , monkey M: $84\% \pm 2 \times 10^{-2}$), significantly outperforming rational models (synthesis models,
238 paired two-sided t-test: $t(1999) = 154$, $p < 10^{-15}$, Cohen's $d = 3.44$; comparison models: $t(1999) =$
239 174 , $p < 10^{-15}$, Cohen's $d = 3.88$; see Fig. 3e). Moreover, models trained on one monkey make

240 choice predictions on the other monkey that still significantly outperform rational RNNs
241 (synthesis models, paired two-sided t-test: $t(1999) = 25$, $p < 10^{-15}$, Cohen's $d = 0.55$; comparison
242 models: $t(1999) = 22$, $p < 10^{-15}$, Cohen's $d = 0.49$; see Fig. 3e). This means that retrained RNNs
243 capture hidden deterministic mechanisms underlying irrational behavior that generalize across
244 trials and individuals.

245 Although we have leveraged the flexibility of RNNs to capture irrational choices, we have
246 not yet demonstrated that irrational RNNs are accurate models of actual OFC computations.
247 Remarkably, their neural CCM distance to the OFC decreases even further compared to their
248 rational counterparts (synthesis models, paired two-sided t-test: $t(1999) = 17$, $p < 10^{-15}$, Cohen's
249 $d = 0.38$; comparison models: $t(1999) = 29$, $p < 10^{-15}$, Cohen's $d = 0.64$). Furthermore, this
250 improvement generalizes across monkeys, as shown when evaluating the neural distance of
251 retrained RNNs to the other monkey (synthesis models, paired two-sided t-test: $t(1999) = 16$,
252 $p < 10^{-15}$, Cohen's $d = 0.36$; comparison models: $t(1999) = 27$, $p < 10^{-15}$, Cohen's $d = 0.61$; see
253 Fig. 3f). This is despite yielding activity patterns that accurately predict inter-individual
254 differences in CCM matrices (see Supplementary Materials). However, one may argue that
255 informing RNN models about monkeys' actual choices may have facilitated the resemblance to
256 any brain system that contributes to behavioral control in the task, thus challenging the
257 anatomical specificity of our results. To address this point, we also computed the neural distance
258 of retrained RNNs to dlPFC and ACC neurons. We first checked that empirical summary statistics
259 of neural information geometry vary more across brain regions than across monkeys (see
260 Supplementary Materials). We then compared neural distances across brain regions; we found
261 that retrained RNNs were significantly closer to the OFC than to the dlPFC and the ACC (synthesis
262 and comparison models, paired two-sided t-test: $t(1999) > 29$, $p < 10^{-15}$ and Cohen's $d > 0.65$ for
263 all comparisons between areas, see Fig. 3g).



264 **Fig. 3 | Behavioral and neural realism of candidate RNN models of the OFC.** **a** and **b**, Estimated option values (see
 265 Methods) for monkey F (**a**) and M (**b**) are shown as a function of reward's magnitude rank (x-axis) and probability rank
 266 (y-axis). **c**, Rational options values are defined as the expected reward, i.e. the product of reward magnitude and
 267 probability. Note that when an attribute is unknown (cf. question mark), the rational model replaces it with its a priori
 268 expected rank. **d**, Monkeys' estimated values (y-axis) are plotted against rational values (x-axis). Each point represents
 269 a specific combination of probability and magnitude ranks, including cases where one or both attributes are unknown.
 270 **e**, Balanced accuracy for predicting monkey choices (green circle: monkey F, pink triangle: monkey M), under each
 271 candidate model (blue: value synthesis, orange: value comparison). Lighter distributions correspond to rational
 272 models, darker distributions to irrational models, and right-most distributions represent irrational models trained on
 273 one monkey and tested on the other. Within each violin plot, the horizontal line denotes the mean, and the thicker
 274 vertical line represents the interquartile range (25th – 75th percentile). Asterisks indicate significant differences, with p-
 275 value < 0.005. **f**, Neural CCM distance between models and the OFC. The white distribution corresponds to random
 276 RNN initializations (identical for both RNN cohorts). **g**, Neural CCM distance (averaged across monkey-specific
 277 distances) between irrational models and the OFC, the dIPFC and the ACC. **h**, Percentage of units classified as *offer*

278 *value*, *chosen value* or *chosen option* cells, in rational and irrational RNN models and in recorded OFC neurons (green
279 circle: monkey F, pink triangle: monkey M), at the time of choice.

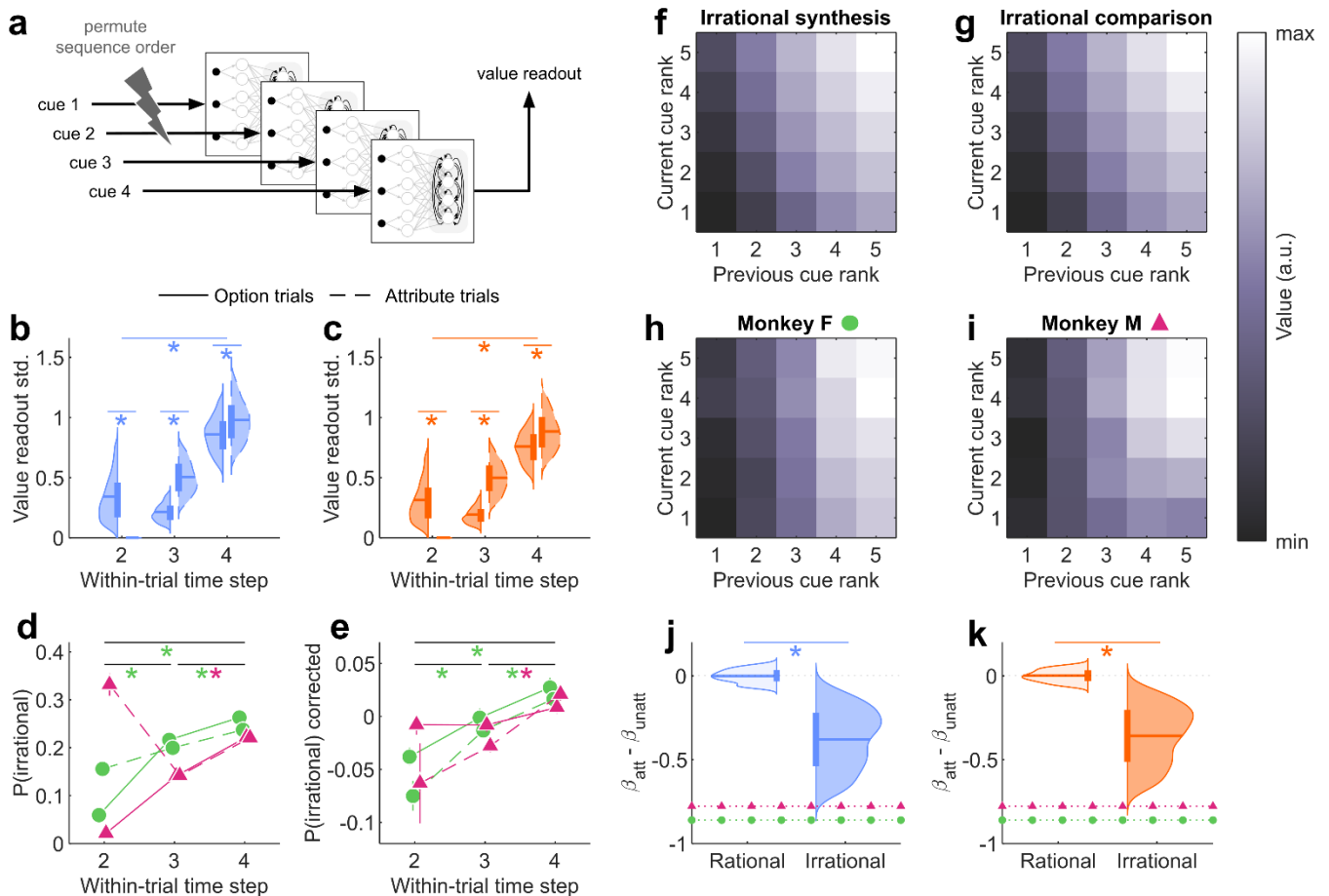
280 One may also ask whether selected RNNs exhibit the established mixed selectivity of OFC
281 neurons. In line with the existing literature^{57,60}, we attempted to classify units according to three
282 distinct response profiles (see Methods): *option value cells*, which encode the value of a single
283 option; *chosen option cells*, which encode the binary identity of the chosen option; and *chosen*
284 *value cells*, which encode the value of the chosen option. As expected, we found that the trial-
285 by-trial variations of OFC neurons' firing rate at the time of choice can be matched to one of the
286 three response profiles. Importantly, this is also the case for integration units of selected RNNs,
287 albeit with a slight over-representation of *offer value* units (see Fig. 3h). Note that this gap is much
288 reduced if we remove the OFC neurons that do not match any of those response profiles (about
289 30% of OFC neurons).

290 Together, these findings suggest that the selected RNNs perform value computations that
291 are realistic, from both a behavioral and neural standpoint. In particular, adjusting the RNNs'
292 internal wiring to account for irrational choices yields both accurate and specific predictions of
293 the OFC's representational geometry. This completes the identification of idealized (rational) and
294 actual (irrational) neural net models of value computations in the OFC.

295 ***Characterizing irrational interferences in value computations***

296 We still lack a computational explanation for why monkeys exhibit irrational behavior in the task.
297 Thus, we now seek to characterize the systematic distortions in cue processing that give rise to
298 irrational choices. First, we quantified potential nonlinear interference effects across decision
299 cues. Recall that, by assumption, rational choices should be solely driven by the informational
300 content of decision cues and thus remain invariant with regard to cue presentation order. In

301 contrast, irrational interference effects would manifest as variability in RNNs' value outputs
 302 across random permutations of cue presentation order, all else being equal. We thus performed
 303 systematic simulations of selected RNNs, quantifying the standard deviation of value outputs
 304 across permuted cue presentation orders, for all possible combinations of two, three or four cues
 305 (see Methods).



306 **Fig. 4 | Interference mechanisms in irrational models and monkeys.** **a**, Schematic procedure for evaluating the
 307 RNNs' sensitivity to cue presentation order. **b**, Standard deviation of the irrational value synthesis RNNs' outputs in
 308 response to random permutations of cue sequence orders (y-axis), as a function of cue onset times (x-axis) during
 309 option trials only (solid line) or attribute trials only (dashed line). Asterisks indicate p-value < 0.005. **c**, Same format as
 310 panel **b**, but for irrational value comparison RNNs. **d**, Rate of monkeys' irrational choices (y-axis; green circle: monkey
 311 F, pink triangle: monkey M), as a function of cue onset time, for both option (solid) and attribute (dashed) trials. Vertical
 312 error bars show the standard error. Asterisks indicate that the difference between time steps (averaged over both trial
 313 types) are significant, with p-value < 0.02. **e**, Average residual irrational choice rate, once decision difficulty has been
 314 regressed away (same format as panel **d**). **f**, Average value output (color scale) of irrational value synthesis RNNs, as a

315 function of the rank of both previously (x-axis) and currently (y-axis) attended cues (see Methods). **g**, Same format as
316 panel **f**, but for irrational value comparison RNNs. **h** and **i**, Same format as panel **f**, but for both monkeys (h: monkey F,
317 i: monkey M). **j**, Difference between the contributions of cue ranks to the attended pseudo-value (attended cue minus
318 unattended cue; see Methods), for both rational (light) and irrational (dark) variants of value synthesis RNNs. The
319 asterisk denotes a significant difference between rational and irrational RNNs, with p -value < 0.01 . Colored dotted
320 lines indicate the attended/unattended cue contribution difference for both monkeys. **k**, Same format as panel **j**, but
321 for value comparison RNNs.

322 By construction, rational RNN models exhibit almost no variability. Thus, irrational RNNs
323 exhibit significantly stronger interference effects than their rational counterparts (paired two-
324 sided t-test at each time step: all $t(1999) > 64$, $p < 10^{-15}$ and Cohen's $d > 1.44$). Importantly,
325 interference effects increase as within-trial decision time unfolds (paired two-sided t-test within
326 each cohort between step 2 and step 4: both $t(1999) > 227$, $p < 10^{-15}$, Cohen's $d > 5.08$; see Fig. 4b
327 and Fig. 4c). This suggests that systematic perturbations in sequential cue processing
328 accumulate over time. Accordingly, monkeys' choices become more irrational – i.e. less
329 consistent with their average preferences (cf. Fig. 3a and Fig. 3b) – as decision time unfolds
330 (Monkey F, paired two-sided t-test between pairs of steps: all $t(24) > 3.3$, $p < 3 \times 10^{-3}$, Cohen's $d >$
331 0.68 ; Monkey M, step 3 vs. step 4: $t(29) = 13$, $p = 2 \times 10^{-13}$, Cohen's $d = 2.45$; see Fig. 4d). One may
332 argue that this accumulating interference effect may only be apparent, because decisions that
333 monkeys trigger later in time tend to be more difficult (see Supplementary Materials). To control
334 for the effect of decision difficulty, we regressed irrational choice rates onto the absolute
335 subjective value difference, across trials. Reassuringly, the residuals of this regression still
336 increase as decision time unfolds (Monkey F, two-sample two-sided t-test between pairs of
337 steps: all $t > 3.4$, $p < 7 \times 10^{-4}$, Cohen's $d > 0.07$; Monkey M, step 3 vs. step 4: $t(12830) = 5.2$,
338 $p = 2 \times 10^{-7}$, Cohen's $d = 0.09$; see Fig. 4e). This means that monkeys' rationality deteriorates with
339 decision time beyond what can be expected from decision difficulty alone.

340 A possibility is that cue traces within the RNNs' integration layer may leak into one
341 another, either across options or across attributes. To investigate this, we separated *option trials*
342 – where the second cue reveals the missing attribute of the same option as the first cue – from
343 *attribute trials* – where the second cue reveals the same attribute as the first cue, but for the other
344 option. At the second cue onset, interference effects are significantly stronger in option trials
345 than in attribute trials, for both RNN types (paired two-sided t-test within each cohort: both
346 $t(1999) > 64$, $p < 10^{-15}$, Cohen's $d > 1.45$). This is also the case for monkey F, based on residual
347 irrational choice rates (Monkey F, two-sample two-sided t-test: $t(1169) = 2.4$, $p = 2 \times 10^{-2}$, Cohen's
348 $d = 0.14$; Monkey M: $t(321) = 1.6$, $p = 0.11$, Cohen's $d = 0.18$; see Fig. 4e). This suggests that cue
349 interference effects are more pronounced within options – i.e. across attributes – than across
350 options. Thus, when attending a given option, we expect the integration of previously and
351 currently attended cues to be asymmetrical, above and beyond differences induced by the type
352 of information that they convey – i.e. reward probability vs magnitude. To test this, we quantified
353 a pseudo-value profile as a function of both previously and currently attended cue ranks,
354 irrespective of cue types (see Methods). In contrast to rational RNNs, irrational RNNs output
355 pseudo-values that are mostly influenced by the previously attended cue (see Fig. 4f and Fig. 4g).
356 When quantified in terms of the relative contribution of the unattended versus attended cue ranks
357 (see Methods), we find that the asymmetry is significantly stronger in retrained RNNs than in
358 rational RNNs (paired two-sided t-test within each cohort: both $t(1999) > 220$, $p < 10^{-15}$, Cohen's d
359 > 4.96 ; see Fig. 4j and Fig. 4k). This asymmetry is also significantly present in monkeys' choices
360 (Monkey F, one-sample two-sided t-test: $t(23) = -18$, $p < 10^{-14}$, Cohen's $d = 3.72$; Monkey M:
361 $t(28) = -17$, $p < 10^{-15}$, Cohen's $d = 3.24$; see Fig. 4j and 4k). Together, these results suggest that
362 previously attended cues leave a persisting value trace that partly resists novel value-relevant
363 information, yielding accumulating perturbations in value computations.

364

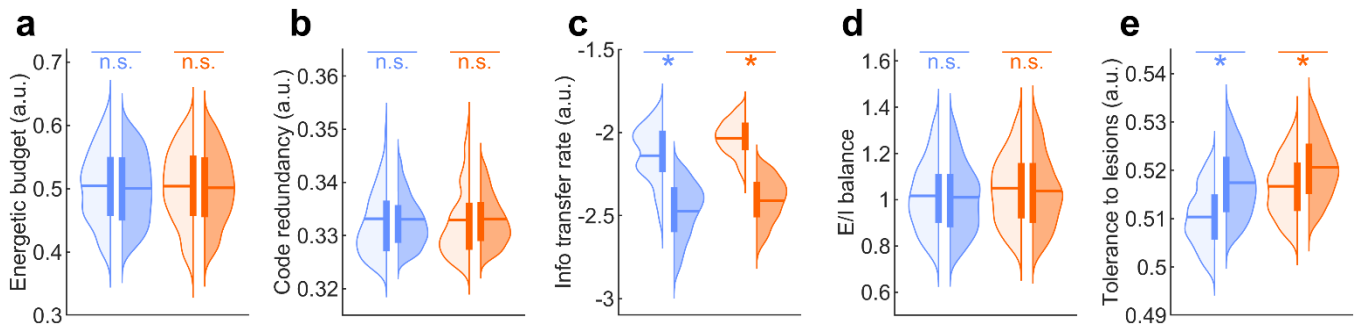
Comparing biological constraint compliance in rational and irrational RNNs

365 In contrast to irrational RNNs, rational RNNs are idealized neural net models of OFC
366 circuits that would have evolved in the absence of energetic or robustness constraints. We now
367 ask whether the wiring peculiarities of irrational RNNs may bring some form of biological
368 advantage that may have overcompensated the behavioral suboptimality that they induce.

369 First, we compared rational and irrational RNNs in terms of the metabolic cost of
370 sustaining their respective structures. Since action potentials are a major source of energetic
371 consumption in the brain⁶¹, we quantified metabolic cost in terms of the average network activity
372 (see Methods). However, we found no significant difference in energetic budget between rational
373 and irrational RNNs (synthesis models, paired two-sided t-test: $t(1999) = 2.0$, $p = 0.04$, Cohen's
374 $d = 0.05$; comparison models: $t(1999) = 1.2$, $p = 0.25$, Cohen's $d = 0.03$; see Fig. 5a).

375 Second, we took inspiration from other variants of efficient coding models, which rather
376 suggests that brain circuits with limited neural resources self-organize to either minimize code
377 redundancy or maximize information transfer rate. We quantify these in terms of the average rate
378 of units' co-activation across all possible units pairs^{62,63} and the entropy of units' response
379 outputs^{20,64}, respectively (see Methods). We found no significant difference in code redundancy
380 (synthesis models, paired two-sided t-test: $t(1999) = 0.4$, $p = 0.68$, Cohen's $d = 0.01$; comparison
381 models: $t(1999) = -0.8$, $p = 0.42$, Cohen's $d = 0.02$; see Fig. 5b). Interestingly however, we found
382 that irrational RNNs exhibit significantly lower information transfer rate than their rational
383 counterparts (synthesis models, paired two-sided t-test: $t(1999) = 51$, $p < 10^{-15}$, Cohen's $d = 1.13$;
384 comparison models: $t(1999) = 82$, $p < 10^{-15}$, Cohen's $d = 1.82$; see Fig. 5c). This implies that the
385 wiring peculiarities that yield irrational behavior also incur an additional cost in terms of
386 information transfer rate.

387 Third, we reasoned that irrational circuits may benefit from better robustness. We start
 388 with excitatory-inhibitory balance, which would ensure stability and/or homeostasis⁴². However,
 389 we found no significant difference in the relative proportion of negative and positive connection
 390 weights (see Methods) between rational and irrational RNNs (synthesis models, paired two-sided
 391 t-test: $t(1999) = 1.2$, $p = 0.23$, Cohen's $d = 0.03$; comparison models: $t(1999) = 2.3$, $p = 0.02$,
 392 Cohen's $d = 0.05$; see Fig. 5d).



393 **Fig. 5 | Potential biological benefits of irrational circuits.** For all panels, asterisks indicate a significant difference
 394 between rational (light) and irrational (dark) RNNs (blue: value synthesis, orange: value comparison), with
 395 p -value < 0.005 . **a**, Metabolic cost, measured as the average network activity, over all trials, trial steps, and units.
 396 **b**, Neural code redundancy, measured as the average co-activation probability over all integration units pairs.
 397 **c**, Information transfer rate, measured as the entropy of units response outputs. **d**, Excitatory-inhibitory balance,
 398 measured as the relative proportion of positive and negative connection weights. **e**, Tolerance to damage, measured
 399 as the retained rate of rational choice from 10% to 50% of lesioned units.

400 We then reasoned that irrational circuits may simply be more tolerant to external
 401 perturbations such as damage or noise. To test this, we simulated random virtual lesions of RNN
 402 integration units and measured the retained rate of rational choices. As expected, rational choice
 403 rate monotonically decreases when the fraction of lesioned units increases, for all types of
 404 models. Thus, we quantify tolerance to damage in terms of the rational choice rate averaged
 405 within a range of lesion sizes (from 10% to 50% of integration units; see Methods). We find that
 406 irrational RNNs exhibit significantly stronger tolerance to damage than their rational
 407 counterparts, irrespective of value computations (synthesis models, paired two-sided t-test:

408 $t(1999) = -27, p < 10^{-15}$, Cohen's $d = 0.61$; comparison models: $t(1999) = -16, p < 10^{-15}$, Cohen's
409 $d = 0.36$; see Fig. 5e). Interestingly, perturbing OFC computations with neural noise or imposing
410 virtual disconnections within the integration layer yield qualitatively identical results (see
411 Supplementary Materials).

412 **Discussion**

413 In this work, we asked whether irrational behavior may be explained by distal constraints
414 that act on the neurobiology of brain decision-making systems. First, we adopted a normative
415 approach to identify models of an idealized OFC network, which would have evolved without any
416 metabolic or robustness constraint. We found that only a specific subset of candidate RNNs
417 reproduces the representational geometry of the OFC – specifically, those that receive inputs
418 encoding option identity in a temporal format (first vs. second option), while computing option
419 values or value difference in an attentional format (attended vs. unattended option). We discuss
420 this finding below. Second, we retrained the selected RNNs to account for monkeys' irrational
421 choices when making decisions under risk. Importantly, these retrained RNNs eventually make
422 out-of-sample behavioral and neural predictions that generalize across trials and individuals. We
423 also showed that their peculiar wiring induces deterministic interferences in value computations
424 that explain the irrational variability of monkeys' choices across within-trial attentional
425 trajectories. Finally, we compared the potential biological benefits of rational and irrational
426 variants of OFC circuits and show that the latter exhibits greater tolerance to damage or noise.
427 Irrational interferences in value computation may thus be understood as an incidental byproduct
428 of selective pressure favoring the robustness of OFC circuits to anatomical damage.

429 That irrational behavior is the incidental outcome of neurobiological constraints is not a
430 novel idea. However, to our knowledge, this work is the first attempt to compare distinct types of
431 constraints and eventually demonstrate the importance of tolerance to circuit damage in this
432 context. We contend that this demonstration is theoretical in essence, at least when compared
433 to empirical work that employ causal – e.g., optogenetic – manipulations to disclose proximal
434 neurobiological constraints^{38,39}. Arguably however, it would have been difficult to provide direct
435 empirical evidence for our claim, at least in mammals. This is inherent to the distal nature of the

436 constraint, which is more readily addressed from a computational perspective. In turn, our
437 conclusions rely on a set of modeling assumptions: we will now discuss these.

438 To begin with, we restricted the set of candidate OFC computations to variants of value
439 synthesis and value comparison. Although a few recent empirical studies consider other types of
440 OFC computations⁶⁵, this prior selection is representative of current debates regarding OFC's
441 contributions to decision making¹. Importantly, we showed that some of these variants reproduce
442 complex features of the OFC's representational geometry, even without being informed with
443 behavioral and/or neural data. This includes established results regarding the mixed selectivity of
444 OFC neural populations (cf. *option value cells*, *chosen value cells* and *choice cells*)^{56,57}.
445 Moreover, we showed that these computational scenarios are anatomically specific, in that their
446 neural predictions do not resemble electrophysiological recordings in either dlPFC or ACC.
447 Retrospectively, this assumption may thus not be so restrictive. Note that the specific RNN
448 variants that we validated using OFC single unit recordings in monkeys are consistent with
449 landmark fMRI studies of value coding in the human vmPFC. In particular, our results directly
450 confirm fMRI studies promoting the attentional format of value coding²³. They are also consistent
451 with other results; for example, if a default option can be identified prior to decision onset (e.g.,
452 in terms of a prior preference over superordinate categories), then pre-stimulus activity in the
453 vmPFC varies with its subjective value, and the magnitude of these variations predicts peoples'
454 irrational attachment to their default preference²². In other words, the vmPFC may use a value
455 coding format that rather distinguishes default versus alternative options. Interestingly, this also
456 aligns with our neural and behavioral results, under the assumption that early preferences – e.g.,
457 based upon the first attended cue – set a default option. The reason is twofold. First, as long as
458 attention remains focused on the first option, attentional and default/alternative value-coding
459 formats are formally indistinguishable. Second, the persisting value trace of the firstly attended
460 cue will, on average, appear as a bias towards the default option. In summary, although the
461 statistical resemblance to the default/alternative hypothesis may be stronger in trials where

462 decisions are triggered prematurely – i.e., before all relevant cues have been processed – we
463 argue that our findings remain compatible with known representational frameworks of value
464 coding in the human vmPFC.

465 Beyond value-coding format issues, one may find it disappointing that we could not
466 disambiguate computational scenarios of value comparison and value synthesis. The underlying
467 question here is whether the OFC directly implements choice, or whether its role is limited to
468 assigning values to available options^{49,66}. When implemented in the form of winner-take-all
469 networks, the former scenario explains established findings in electrophysiological and
470 neuroimaging studies, in particular: the observed mixed selectivity of OFC cells^{56,60}, as well as the
471 apparent encoding of the value difference between chosen and unchosen options – at least close
472 to choice onset (not shown)⁶⁷. Interestingly, we have shown that such findings can be equally well
473 reproduced by RNNs performing either value synthesis or value comparison, which calls for
474 experiments that are designed to distinguish these kinds of computational scenarios, as opposed
475 to testing one of them.

476 Also, we did not vary the global architecture of our artificial neural nets, which consisted
477 of a layer of attribute-encoding units sending their outputs to a layer of recurrently connected
478 integration units. In line with recent neural net approaches to value computations in the OFC^{20,49},
479 we adopted the minimal architecture that ensures universal approximation capabilities while
480 using a limited number of artificial units^{68,69}. Note that a major computational bottleneck of both
481 value synthesis and value comparison scenarios is OFC circuits' capacity for combining value-
482 relevant attributes of arbitrary number and type⁵⁵. Now, the above two-layer architecture provides
483 a flexible and simple solution to this problem that rests on the second layer's trained ability to
484 integrate arbitrary sequences of attributes, whose type and rank are encoded in separate pools
485 of the first layer units. This circumvents the need for otherwise unrealistic, context-dependent
486 changes in connectivity with upstream brain systems involved in recognizing or storing value-

487 relevant information. Nevertheless, the relative simplicity of our design contrasts with previous
488 studies that favored off-the-shelf deep neural nets to approximate the hierarchical organization
489 of, e.g., primates' visual ventral stream⁵⁸ or humans' language networks⁷⁰. From a machine
490 learning perspective, tasks such as visual perception and speech comprehension are inherently
491 difficult problems, which remained unsolved until the advent of deep neural nets trained on
492 massive, labeled datasets. In these domains, objective task performance reliably predicts
493 statistical similarity with neural data. This relationship, however, does not readily generalize to
494 our findings: RNNs resemble more closely OFC data when they permit systematic, error-inducing
495 interferences. In retrospect, it is remarkable that our value synthesis/comparison RNNs exhibit
496 such realistic features, at both the behavioral and neural levels. This is despite the degeneracy of
497 RNN wiring profiles with regard to each type of value computation, which we systematically
498 explored by repeating the training process across many random initializations of RNN
499 parameters. Arguably, the ensuing marginalization process renders our results robust to local
500 minima issues. This statistical benefit would have been prohibitively costly to match using deeper
501 neural net architectures.

502 One might also argue that rational and irrational RNNs may have been compared in an
503 unfair manner. For example, we chose to train rational RNNs under a normative approach, which
504 precludes idiosyncratic variations in risk attitudes. The rationale here was twofold. First, we
505 aimed at selecting neural nets that could serve as neutral and fully interpretable reference points,
506 in that their computational objective was under our control – i.e. computing expected values, as
507 prescribed by rational decision theory. Second, we wanted to assess the neural realism of
508 idealized OFC network models that were derived from first principles alone, without any access
509 to data from the decision task. When evaluating the resemblance of candidate RNNs to OFC
510 recordings, this leaves little room for methodological objections to the ensuing model selection.
511 Now, we acknowledge that, when it comes to measuring statistical similarity to neural recordings,
512 irrational RNNs may benefit from being trained on individual behavioral datasets. However, the

513 fact that irrational RNNs make out-of-sample predictions that generalize across trials and
514 individuals rather suggests that they have captured hidden, yet shared, decision mechanisms. In
515 any case, there is no reason to think that this training difference would favor irrational RNNs in
516 terms of, e.g., tolerance to circuit damage. A related concern is whether the latter may be the
517 artefactual byproduct of re-training *per se*, which may provide an additional opportunity for
518 improving efficiency or robustness. This is the reason why we also explored another training
519 strategy for irrational RNNs, which starts from the same randomly initialized parameter sets as
520 rational RNNs. Reassuringly, our conclusions remain unchanged under this alternative training
521 strategy (see Supplementary Materials).

522 In conclusion, we believe our modeling assumptions are tenable, at least when compared
523 to state-of-the-art computational studies in the field. Importantly, they have enabled us to
524 reverse the usual approach to disclosing distal neurobiological constraints on rationality, which
525 typically rests on highlighting conflicts with the demands of behavioral performance (cf. Fig. 1c,
526 1d and 1e). In contrast, we identify realistic mechanisms that explain observed deviations to
527 rationality and explore their potential neurobiological advantages. We believe that this may be a
528 fruitful method for investigating related evolutionary or developmental issues in cognitive
529 neuroscience.

530 **Methods**

531 ***Task design***

532 The decision task is summarized in Fig. 1a. Monkeys were seated in a behavioral chair
533 with their heads restrained. Each trial began when the monkey fixated on a central fixation cue for
534 500 ms. At the start of the trial, two options were presented, each consisting of two hidden cues
535 initially masked by grey squares. One of these squares then turned blue, indicating the first cue
536 available for sampling. When the subject fixated on the blue square, the corresponding picture
537 cue was revealed and had to be continuously fixated for 300 ms before it was re-masked. All
538 picture cues had been previously learned and were associated with a specific rank of a specific
539 reward attribute (either probability or magnitude). Probability cues indicated reward probabilities
540 of 10%, 30%, 50%, 70%, or 90%, while magnitude cues represented reward magnitudes of 0.15,
541 0.35, 0.55, 0.75, or 0.95 (arbitrary units of an appetitive juice). Following the initial cue, a second
542 blue square highlighted the next available cue, which had to be sampled using the same
543 procedure. This second cue was either the other attribute of the same option (*option trial*) or the
544 same attribute but for the other option (*attribute trial*). After the second cue, the two remaining
545 cues were simultaneously highlighted with blue squares, allowing the subject to freely choose
546 which one to sample next, or to select one of the two options using a joystick. If a third cue was
547 sampled, the subject could then either sample the final cue or make a choice. Once the fourth
548 cue was revealed, the subject was required to make a choice.

549 ***Value profile estimation***

550 We estimated the subjective value profile of each monkey (and each model) using
551 standard statistical procedures, based on the observed choices (cf. Figures 3a, 3b and 3c). More

552 precisely, we fitted the underlying value function, under the assumption that choices followed a
553 simple softmax mapping of the difference in option values:

$$554 \quad P(\text{choose option 1}) = \frac{1}{1 + \exp\left(-\left(V(p_1, m_1) - V(p_2, m_2)\right)\right)} \quad (1)$$

555 where p_i and m_i denote the reward probability and magnitude of option i , as known by the agent
556 at the time of choice, and $V(p, m)$ is the corresponding subjective value. Equation 1 provides a
557 binomial likelihood function for observed choices, given the unknown monkeys' value function.
558 Parameterizing the value function then enables us to regress trial-by-trial choices against option
559 attributes.

560 To maximize modelling flexibility, we employed a semi-parametric approach, whereby
561 each possible combination of probability and magnitude - including cases in which one or both
562 attributes were unknown at the time of choice - is captured using a specific model parameter
563 that quantifies the corresponding value. Given that each attribute has five possible ranks (plus
564 the unknown attribute case), this means that we estimate $6 \times 6 = 36$ parameters from the choices.
565 This semi-parametric approach enables us to capture entirely arbitrary value profiles over the bi-
566 dimensional space spanned by reward probability and magnitude, under the constraints that (i)
567 the same value function applies to all options, and (ii) the value function does not depend upon
568 the cue presentation order. The ensuing probabilistic logistic regression was performed using the
569 VBA academic freeware⁷¹, which relies on the variational Laplace approach to approximate
570 Bayesian inference^{72,73}.

571 **RNN architecture**

572 The basic RNN architecture is summarized in Fig. 1b. Let $t \in \{1,2,3,4\}$ denote the time
 573 step at which cues are revealed or attended within a decision trial. The RNN component variables
 574 are defined as follows:

- 575 - $\mathbf{U}(t) \in \mathbb{R}^3$: input vector at time t , whose entries include the attribute rank (from 1 to 5)
 576 and type (probability or magnitude) of the currently attended cue, as well as the identity
 577 of the currently attended option (see below).
- 578 - $\mathbf{X}_1(t) \in \mathbb{R}^9$: unit activation vector in the first hidden layer at time t . This *input layer* is com-
 579 posed of units that are selective to input entries (units 1 to 5, 6 to 7, and 8 to 9 receive
 580 $U_1(t)$, $U_2(t)$ and $U_3(t)$, respectively). It forms a standard population code of the currently
 581 attended cue.
- 582 - $\mathbf{X}_2(t) \in \mathbb{R}^{10}$: unit activation vector in the second hidden layer at time t . Feedforward con-
 583 nections convey the current activity of the first layer ($\mathbf{X}_1(t)$) to the second layer's units.
 584 But these units are also connected with each other through recurrent connections that
 585 carry reentrant dynamics ($\mathbf{X}_2(t-1)$). If properly wired (see below), this *integration layer*
 586 can thus integrate currently and previously attended cues.
- 587 - $\mathbf{Y}(t) \in \mathbb{R}^1$ (for value comparison models) or $\mathbf{Y}(t) \in \mathbb{R}^2$ (for value synthesis models): value
 588 readout at time t .

589 At any time t within a decision trial, cue-related information propagates through the
 590 network according to the following equations:

$$591 \quad \mathbf{X}_1(t) = f(\mathbf{W}_{encode} \cdot \mathbf{U}(t) - \mathbf{B}_1) \quad (2)$$

$$592 \quad \mathbf{X}_2(t) = f(\mathbf{W}_{forward} \cdot \mathbf{X}_1(t) + \mathbf{W}_{recurrent} \cdot \mathbf{X}_2(t-1) - \mathbf{B}_2) \quad (3)$$

$$593 \quad \mathbf{Y}(t) = \mathbf{W}_{readout} \cdot \mathbf{X}_2(t) \quad (4)$$

594 under the constraint that activity within the integration layer is reset at each decision trial, i.e.
595 $X_2(0) = 0$ by convention.

596 Here, W_{\blacksquare} refers to matrices of connection weights, and $B_{1:2}$ are bias vectors applied to
597 the corresponding hidden layers. The weights W_{encode} and biases B_1 are set such that each
598 admissible cue rank (U_1) preferentially activates a dedicated unit in a rank-specific pool of first
599 hidden layer units. Similarly, each admissible cue type (U_2) and option identity (U_3) preferentially
600 activates one out of two units each (again in secluded pools of first hidden layer units). To ensure
601 distributed encoding within each pool, the activation profiles of first layer units were set to tile
602 the domain of their specific input uniformly: whenever one unit's activity reached 75% of its
603 maximum, the next "adjacent" unit in the pool was 25% active.

604 To impose bounds on units' firing rates, we use a standard sigmoid activation function f
605 for all hidden units:

$$606 \quad f: x \mapsto \frac{1}{1 + \exp(-x)} \quad (5)$$

607 Importantly, when structurally organized into two such hidden layers, neural nets with a
608 limited number of sigmoidal units possess universal approximation capabilities^{68,69}. This means
609 that RNNs of this type can be trained to store memory traces of all previously attended cues,
610 albeit in a format that may not be directly accessible. More generally, it should be possible to train
611 such RNNs (see below) to compute any arbitrary mapping of the sequence of attended cues.

612 The RNN receives inputs one at a time, in a sequential manner, as monkeys did in the task.
613 The sequence order is determined by the exogenous control of attention, which samples cues in
614 an arbitrary fashion within a decision trial. Let $U_1(t)$, $U_2(t)$ and $U_3(t)$ denote the entries of the
615 input vector $\mathbf{U}(t) \in \mathbb{R}^3$:

616 - $U_1(t)$ encodes the normalized rank of the attended cue, with the following mapping:

Magnitude cue	Probability cue	Cue rank	U_1
0.15 AU	10 %	1	0.1
0.35 AU	30 %	2	0.3
0.55 AU	50 %	3	0.5
0.75 AU	70 %	4	0.7
0.95 AU	90 %	5	0.9

- 617 - $U_2(t)$ encodes the attribute type, i.e. probability: $U_2 = 0$ and magnitude: $U_2 = 1$.
- 618 - $U_3(t)$ encodes the identity of the attended option, i.e. option 1: $U_3 = 0$, option 2: $U_3 = 1$.

619 Note that the identity of the attended option can be encoded in two different
620 representation formats: spatial (left vs. right) or temporal (first vs. second). This distinction
621 affects the encoding of U_3 , as illustrated in the following example trials:

Trial ID	Time step	Attended option side	U_3 in the <i>spatial</i> frame (right = 0, left = 1)	U_3 in the <i>temporal</i> frame (first = 0, second = 1)
1	1	Right	0	0
1	2	Left	1	1
1	3	Left	1	1
1	4	Right	0	0
2	1	Left	1	0
2	2	Left	1	0
2	3	Right	0	1
2	4	Right	0	1

622 Although the encoding frame of the attended option's identity does not modify the raw
623 information that is conveyed to the network, it nonetheless determines the pattern of activity in
624 the input layer. In turn, the encoding format is likely to modify the sequence of activity patterns of

625 the whole RNN over time steps within a decision trial. In other terms, the way the RNN responds
 626 to a sequence of cues will depend upon the encoding format of the attended option's identity.

627 Similarly, the value outputs $Y(t)$ of the network can be expressed in different
 628 representation formats: spatial, temporal, or attentional (attended vs. unattended). The example
 629 trials below illustrate how the encoding format of option values varies across these frames. Let
 630 V_{left} and V_{right} denote the values of the left and right options as estimated by the network at each
 631 cue onset. The statistical similarity between representation formats depends on the actual
 632 sequence order of cue attendance:

Trial ID	Time step	Attended option side	Y in the <i>spa-</i> <i>tial</i> frame	Y in the <i>tem-</i> <i>poral</i> frame	Y in the <i>atten-</i> <i>tional</i> frame
1	1	Right	$V_{right} \& V_{left}$	$V_{right} \& V_{left}$	$V_{right} \& V_{left}$
1	2	Left	$V_{right} \& V_{left}$	$V_{right} \& V_{left}$	$V_{left} \& V_{right}$
1	3	Left	$V_{right} \& V_{left}$	$V_{right} \& V_{left}$	$V_{left} \& V_{right}$
1	4	Right	$V_{right} \& V_{left}$	$V_{right} \& V_{left}$	$V_{right} \& V_{left}$
2	1	Left	$V_{right} \& V_{left}$	$V_{left} \& V_{right}$	$V_{left} \& V_{right}$
2	2	Left	$V_{right} \& V_{left}$	$V_{left} \& V_{right}$	$V_{left} \& V_{right}$
2	3	Right	$V_{right} \& V_{left}$	$V_{left} \& V_{right}$	$V_{right} \& V_{left}$
2	4	Right	$V_{right} \& V_{left}$	$V_{left} \& V_{right}$	$V_{right} \& V_{left}$

633 Note that not all combinations of input/output formats are admissible. More precisely,
 634 when the currently attended option's identity is encoded using the spatial format, then value
 635 outputs can be encoded in all representation formats (3 possibilities). However, when the
 636 currently attended option's identity is encoded using the temporal format, then the spatial
 637 information is lost, which leaves only 2 possible value encoding formats (temporal and
 638 attentional frames). This also means that there are only 5 combinations of input/output
 639 representation formats in total. Training the RNN to encode value computations within a given

640 format will require specific settings of recurrent connection weights. This implies that the
641 encoding format of value readouts is also likely to modify the way the RNN responds to a
642 sequence of cues.

643 ***RNN training***

644 *Rational training*

645 Under standard decision theory, rational choices maximize the *expected value*, where the
646 expected value of an option is defined as the product between reward probability and magnitude:
647 $EV = p \times m$. Formally, options' expected value can be defined at each time step within a decision
648 trial, under the assumption that option attributes that are unknown (i.e., yet to be sampled) can
649 be replaced by their expectation under the task distribution. In principle, *value synthesis* and
650 *value comparison* models can thus be trained to output the expected value of both options or,
651 respectively, the difference in expected values, in response to each cue presentation (at each
652 time step within a decision trial). It turns out that RNNs of the sort we describe above can reliably
653 learn to solve the problem of integrating previously and currently attended cues to output options'
654 expected values from almost any random initialization of their connectivity.

655 All RNN models were implemented and trained using the VBA academic freeware⁷¹. The
656 RNN parameters subject to training ($\mathbf{W}_{forward}$, $\mathbf{W}_{recurrent}$, $\mathbf{W}_{readout}$ and \mathbf{B}_2) were initialized as
657 samples from an i.i.d. Gaussian distribution with mean 0 and variance 0.5. For each RNN model,
658 the training procedure was repeated with a different initial random sample, until 1,000 trained
659 models reached 95% accuracy on held-out test data. This enables us to sample the unknown
660 manifold of RNN wirings that can perform a given type of value computation. In the main text, we
661 refer to the ensemble of trained RNNs as a *cohort*, each of which corresponds to a given type of

662 value computation (value synthesis versus value comparison) and a given combination of
663 input/output representation format (see above).

664 For each model instance in a given RNN cohort, a training set and a testing set consisting
665 of 500 trials each were generated. Every trial consisted of a sequence of four cues, randomly
666 chosen among the set of different option pairings, and presented in a random sequence order.
667 Importantly, trained models carry no memory trace of preceding decision trials (their internal
668 state is reset to baseline levels at each trial onset). Also, we did not endow RNNs with the capacity
669 to decide which cue to attend to or when to commit to a decision; rather, we trained them to
670 operate value computations independently of such processes, which are treated as exogenous
671 and arbitrary. Note that training and testing trials could be classified post-hoc as either *attribute*
672 *trials* or *option trials*, depending on whether attention switched to the second option at the
673 second cue onset, or not.

674 Training was terminated when the absolute change in variational free energy between VBA
675 successive iterations fell below 10. However, the trained RNN was kept in its corresponding
676 model cohort only if it reached at least 95% of explained variance (in trial-by-trial expected
677 values) on its testing set. Each RNN cohort consisted of 1,000 independently trained model
678 instances, each with a unique training set, testing set, and parameter initialization. Importantly,
679 random seeds were shared across cohorts, which allowed for paired comparisons across
680 cohorts.

681 Importantly, we stored the updated model parameters at each step of the training
682 procedure, for each model instance of each RNN cohort. This enabled us to test RNNs' neural
683 predictions against OFC neural data (see below) as training unfolds (cf. Fig. 2c).

684 Rational training with biological constraints

685 Decision systems in the brain may have been shaped under the multiple imperatives of
686 producing rational choices and satisfying biological (energetic, informational or robustness)
687 constraints. Under a neural net perspective, all behavioral and biological properties of such
688 systems are determined by their internal wiring. Therefore, training RNNs to perform rational
689 value computations while simultaneously complying with these constraints should enable us to
690 reveal possible conflicts between behavioral and biological imperatives.

691 Here, we focused on three main biological constraints: namely, maximal information
692 transfer rate, minimal energy budget, and maximal tolerance to damage (see section *Biological*
693 *constraints* below). In brief, RNN training followed the exact same procedure as above, except
694 that the objective function included an additional biological constraint adequacy term. To control
695 the balance between behavioral and biological imperatives, we introduced a constraint
696 compliance weight hyperparameter that varied exponentially from 10^{-3} to 10^3 (the higher the
697 constraint weight, the tighter the biological constraint on the RNN wiring) and trained 10 RNN
698 instances per constraint weight. This enabled us to compare RNN cohorts trained under
699 systematically varied levels of biological pressure.

700 For example, Fig. 1c, 1d and 1e plot the mean rate of rational choices against the achieved
701 biological constraint adequacy level, for each compliance weight (and each biological
702 constraint). This reveals the shape of the ensuing Pareto front, which tells how conflicting
703 behavioral and biological imperatives are. If the biological constraint spans the null space of
704 value computations, then increasing the compliance weight has no bearing on choice rationality.
705 However, if behavioral and biological imperatives conflict with each other, tightening the
706 biological constraint results in larger rationality losses, which is what we see here.

707 Re-training RNNs to explain monkeys' irrational choices

708 To begin with, recall that monkeys often commit to a choice before having sampled all
709 cues. However, this does not imply that such choices are irrational, in that monkeys may still
710 select the option with the highest expected value (where unknown attributes are replaced by their
711 expectation under the task distribution). In fact, according to this criterion, only about 20% of
712 monkeys' choices are irrational. To preserve the interpretability of value computations and
713 input/output representation formats of rational RNNs, the re-training of RNN models was
714 restricted to $W_{recurrent}$ (i.e. all other parameters were frozen). This effectively restricts the
715 admissible sources of irrational choices to within-trial nonlinear interferences and spill-over
716 effects between attended cues.

717 In contrast to the rational training phase, where value outputs can be evaluated at each
718 cue onset within decision trials, irrational re-training relies solely on observed monkey choices.
719 The latter only constrain the value readouts of retrained RNNs at the time of choice (i.e. after the
720 last sampled cue). Moreover, the only available information is the choice itself, which we
721 compare to RNN value readouts via the same softmax mapping as in Equation 1. From this
722 perspective, the retraining of RNNs may be seen as a way to relax the two constraints of the
723 logistic regression in Equation 1 (value invariance across options and cue presentation orders).

724 Each RNN instance within each cohort was then re-trained to fit the choices of each
725 individual monkey, using a training dataset of 2,000 trials randomly selected from monkeys'
726 recorded sessions. This procedure produces two twin versions of each retrained irrational model
727 - one for each monkey. We then test their respective behavioral (Figure 3e) and neural (Figure 3f,
728 3g and 3h) predictions within and across monkeys. This enables us to evaluate their inter-trial and
729 inter-individual generalization ability.

730 Note that neither rational training, nor irrational retraining is informed by neural
731 recordings. This greatly simplifies the validation of quantitative model predictions regarding the
732 representational geometry of OFC neural populations.

733 ***Analysis of representational geometry within neural populations***

734 *Representational similarity analysis*

735 Let $\mathbf{X}_2^U(t)$ denote the vector of activations in the RNN's integration layer at time t in
736 response to a sequence of input $\mathbf{U}(1), \dots, \mathbf{U}(t)$. At this first cue onset ($t = 1$), this vector can be
737 computed for each possible input \mathbf{U}_k , which yields 20 distinct activation patterns (i.e., 5 cue
738 ranks \times 2 cue types \times 2 options). The ensuing representational dissimilarity matrix (*RDM*) is
739 constructed element by element by computing pairwise similarities between these activation
740 patterns⁷⁴:

$$741 \quad RDM_{k,l} = r(\mathbf{X}_2^{U_k}(1), \mathbf{X}_2^{U_l}(1)) \quad (6)$$

742 where r denotes Pearson's correlation. If $RDM_{k,l}$ is strongly positive, then activity patterns are
743 mostly invariant to differences between inputs \mathbf{U}_k and \mathbf{U}_l , i.e. the neural representation of these
744 inputs are similar. In brief, RDMs enables us to identify what input features need to change to
745 elicit distinct neural responses.

746 The same procedure is applied to recordings of OFC neurons (as well as to neural
747 recordings within the dlPFC and the ACC), where $\mathbf{X}_2^U(1)$ is replaced with the average firing rate in
748 a 100-400 ms window after the first cue onset. This yields two RDMs: one for the model
749 (RDM^{model}) and one for the OFC data (RDM^{OFC}). Full RDM summary statistics for all monkeys
750 and brain regions can be eyeballed in Figure S2.

751 Finally, the similarity between these matrices is quantified using a rank-based distance
752 metric:

$$753 \quad d_{RDM} = 1 - \rho(RDM_{upper}^{OFC}, RDM_{upper}^{model}) \quad (7)$$

754 Here, ρ denotes Spearman's correlation and RDM_{upper} refers to the upper triangular half
755 of the matrix, excluding the diagonal. We used a rank-based metric because experimental neural
756 data is typically much noisier than RNN activations, resulting in compressed correlation ranges
757 that are more appropriately captured by rank correlations.

758 Cross-correlation matrices

759 The benefit of the above representational similarity analysis is to account for all attributes
760 of a given cue, i.e. attended option identity, attended cue type and attended cue rank.
761 Unfortunately, it does not scale well with the number of cue combinations. In our context, its
762 statistical cost is prohibitive for later phases of decision trials, when more than one cue has been
763 attended. For example, at the second cue onset, there are $20 \times 20 = 400$ possible cue
764 combinations, which would induce RDMs with almost 79,800 elements. This is incompatible with
765 the number of acquired decision trials in the task. Thus, we resort to another type of summary
766 statistics, which was proposed by Hunt and colleagues⁵⁰. In brief, this analysis enables us to
767 quantify and compare the multiple traces that cue sequences leave on units' activity, at the cost
768 of neglecting differences induced by cue types. This simplifying assumption exploits the
769 observed quasi-symmetrical impact of reward probability and magnitude on monkeys' subjective
770 value profiles (see Fig. 2a).

771 Let $X_2^U(i, t)$ denote the activation of unit i in the RNN's integration layer at time t in
772 response to a sequence of input $\mathbf{U}(1), \dots, \mathbf{U}(t)$. We regress each unit's trial-by-trial activity
773 variations at cue onset t concurrently onto trial-by-trial variations of normalized attribute rank in
774 all cues, while identifying cues by their appearance order in the sequence:

775
$$X_2^U(i, t) = \beta_{i,1}(t) \times U_1(1) + \dots + \beta_{i,k}(t) \times U_1(k) + \dots + \beta_{i,t}(t) \times U_1(t) \quad (8)$$

776 Note that we also include two additional regressors, which encode how consistent the 2nd
 777 and 3rd cues are with regard to the currently preferred option, as well as an intercept term (not
 778 shown in Equation 6). This approach aims at detecting nontrivial memory traces of previously
 779 attended cues, while ruling out mere confirmation effects in value coding neurons. Importantly,
 780 we separate *option trials* (where the first two cues belong to the same option) from *attribute trials*
 781 (where the first two cues describe the same attribute – i.e. probability or magnitude – but for both
 782 options) prior to performing the regression analyses. This yields one set of regression coefficient
 783 estimates $\beta_{i,k}(t)$ per trial type, for each unit and each cue in the sequence.

784 Let $\tilde{\beta}_k(t) \in \mathbb{R}^{n_{units}}$ denote the vector of t-statistics associated with regression coefficient
 785 estimates for the k^{th} attended cue ($k \in \{1, 2, 3\}$), across integration units. This vector measures
 786 how sensitive to the k^{th} attended cue second layer units are (at time t) in normalized signal-to-
 787 noise ratio units. This enables a direct quantitative comparison across units, cue presentation
 788 orders and decision times. Note that $\tilde{\beta}_k(t)$ vectors that involve cues presented after the time
 789 when units' activity is sampled (i.e. when $k > t$) are statistically meaningless.

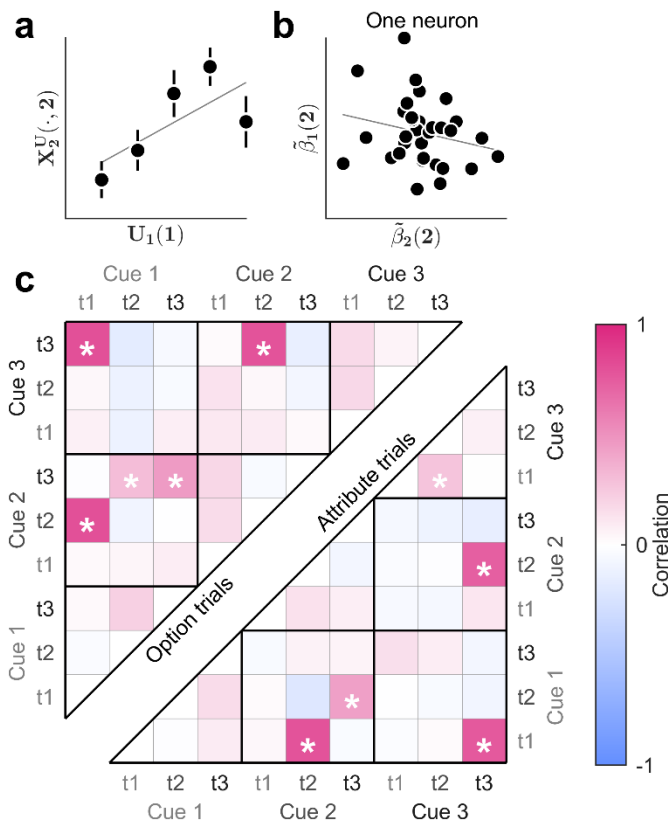
790 We then define the cross-correlation matrix (*CCM*) as follows:

791
$$CCM_{k,k',t,t'} = r(\tilde{\beta}_k(t), \tilde{\beta}_{k'}(t')) \quad (9)$$

792 where r denotes Pearson's correlation. A strongly positive CCM cell indicates that those units
 793 that are most sensitive to the rank of the k^{th} attended cue at time t are also those most sensitive
 794 to the rank of the k'^{th} cue at time t' . Importantly, the sign of regression coefficients matters. This
 795 means that CCM analyses can detect when a unit respond positively to some cue while
 796 responding negatively to another one (as may be the case when cues belong to distinct options).

797 We obtain full CCMs by systematically varying cue presentation orders (k and k') as well
 798 as activity sampling times (t and t'), yielding a 9 by 9 symmetrical matrix. We then remove CCM

799 cells that are meaningless ($t < k$ or $t' < k'$) to avoid statistical illusions possibly induced by
 800 imperfections in trial randomizations. We repeat this process for both *option trials* and *attribute*
 801 *trials*, yielding two CCM types. Differences between the two types of CCM cells that involve the
 802 first and second cue onset times (i.e. $CCM_{1,2,\blacksquare,\blacksquare}$) signal a shift in how the attended option affects
 803 the network's distributed computations. In particular, if neurons respond to the value difference
 804 between options, then one expects $CCM_{1,2,2,2}$ to be positive for option trials, and negative for
 805 attribute trials⁵⁰.



806 **Fig. M1 | Derivation of CCM statistics in OFC neurons.** **a**, First, we regress the mean firing rate of each neuron at each
 807 cue onset against the rank of all previously attended cues (across trials). Here, we show the apparent statistical
 808 relationship between the activity $X_2^U(., t)$ of one example OFC neuron sampled at time $t = 2$ (y-axis), plotted against
 809 the rank $U_1(k)$ of the first presented cue ($k = 1$). **b**, Second, we measure the correlation (across neurons), between the
 810 ensuing regression coefficients for different activity sampling and cue presentation times. Here, we plot the sensitivity
 811 $\hat{\beta}_k(t)$ of OFC neurons (one dot is one neuron) sampled at time $t = 2$ to the first cue ($k = 1$, y-axis) against their
 812 sensitivity to the second cue ($k = 2$, x-axis). **c**, Each cell in the CCM matrix shows the correlation across neurons for a
 813 given pair of regression coefficients (neurons pooled across monkeys). The upper half of the matrix shows the results

814 computed on *option trials* (where the two first cues characterize the same option), while the lower half corresponds to
815 *attribute trials* (where the two first cues characterize the same attribute, but different options). Asterisks indicate
816 significant correlations, with p-value < 0.0007 (correction for multiple comparisons across CCM cells).

817 We apply the exact same analysis on recorded data from OFC neurons (as well as neurons
818 in the dlPFC and ACC), where $X_2^U(i, t)$ is replaced with the average firing rate in a 100-400 ms
819 window after each cue onset. This provides summary statistics whose temporal resolution
820 matches that of RNN models. The procedure is summarized on Figure M1 above.

821 To compare the informational geometry of RNNs and OFC neural populations, we simply
822 compute the Euclidian distance between the meaningful CCM cells:

$$823 \quad d_{CCM} = \sqrt{\sum_{k \leq t, k' \leq t'} (CCM_{k,k',t,t'}^{OFC} - CCM_{k,k',t,t'}^{model})^2} \quad (10)$$

824 Note that, although this approach scales well with the number of cue combinations, it
825 typically is insensitive to the type of attribute (probability vs magnitude) of attended cues.

826 *Mixed selectivity: offer value cells, chosen value cells and choice cells*

827 To identify possible “offer value”, “chosen value”, and “choice” cells, we replicated the
828 analysis previously introduced by Padoa-Schioppa and colleagues⁵⁷. For each unit, we performed
829 four separate regressions of activity sampled at the time of choice, across all trials, against the
830 value of option 1, the value of option 2, the value of the chosen option (defined as the monkey’s
831 recorded choice), or the identity of the chosen option, respectively. Each unit was assigned to the
832 category that yielded the highest percentage of explained variance, provided the regression was
833 significant (p-value < 0.05). Otherwise, no category was assigned. When applied to neural
834 recordings in the OFC, we relied on subjective value profiles, as estimated from monkeys’ choices
835 in the task (see subsection Value profile estimation above). To maximize the match between

836 analyses, we also use model-specific value profiles for RNNs (using the encoding formats that
837 corresponds to each RNN model).

838 ***Analysis of computational interferences in irrational RNNs***

839 *Dependency on cue sequence order*

840 In principle, rational behavior in the task only depends upon the content of value-relevant
841 information, but not on its presentation sequence order. Under this view, any observed
842 dependency on cue sequence order violates rationality.

843 Let $\Delta V^U(t)$ denote the value difference between options, as can be readout from the
844 RNN's activity pattern at time t , in response to a given sequence of input $U(1), \dots, U(t)$. Note that,
845 for value synthesis models, we compute $\Delta V^U(t)$ by simply subtracting the readouts of both
846 option values. To quantify the dependency on cue presentation order, we first simulate the RNN
847 response to all admissible permutations of cue orderings while keeping the combination of t
848 attended cues constant, measure the standard deviation of $\Delta V^U(t)$, and then average the results
849 over all possible cue combinations. Importantly, we repeat this process separately for option
850 trials and attribute trials, meaning that we only consider cue order permutations that are
851 admissible for each trial type.

852 Let \check{U} be the set of all possible combinations of t cues, and for each member set $u \in \check{U}$,
853 let $S(u)$ denote the set of admissible permuted sequence orderings of those cues (restricted to
854 the relevant trial type). Then, we define the dependency on sequence order at time t , denoted
855 $d_S(t)$, as follows:

$$856 \quad d_S(t) = \frac{1}{|\check{U}|} \sum_{u \in \check{U}} \sqrt{\text{Var}(\Delta V^U(t) | U \in S(u))} \quad (11)$$

857 Note that $d_S(t)$ is defined for all decision times starting from the second cue onset up
858 until the last cue onset ($2 \leq t \leq 4$). As we show in the Results section, eyeballing $d_S(t)$ as a
859 function of decision time enables us to track the possible accumulation of interferences in RNN
860 computations. Reassuringly, rational RNNs show no detectable dependency on sequence order,
861 i.e. $d_S(t) \approx 0$, irrespective of cue onset time (not shown).

862 Note that this analysis cannot be directly applied to monkeys' choices, as we cannot have
863 access to the monkeys' internal value estimates for each cue sequence order. This is because
864 the total number of admissible cue sequence orders is prohibitive, at least for late cue inset
865 times. This means that decision trials in the task do not sample the set of admissible cue
866 sequence order in a dense manner. Therefore, we resort to simpler measures of apparent
867 deviations to rational choice, which effectively reduce to detecting trials that are incongruent with
868 monkeys' value profile estimates (see subsection Value profile estimation above). This effectively
869 yields a trial-by-trial binary rationality flag: its sample average provides the apparent rate of
870 rational choices for each decision time ($2 \leq t \leq 4$). Now, monkeys' choice stochasticity may
871 render difficult decisions more likely to violate strict preference orderings. We thus need to
872 correct the apparent rate of rotational choices for decision difficulty. To do this, we perform a
873 logistic regression of the trial-by-trial rationality flag variable onto decision difficulty, as given by
874 the absolute difference in estimated subjective option values (separately for each monkey). The
875 residuals of this regression quantify the amount of trial-by-trial rationality that cannot be
876 explained by variations in decision difficulty. The sample average of these residuals, for each
877 decision time, is what we call the corrected rate of rational choices.

878 Persisting value traces

879 The above dependency on sequence order may be partly driven by a directional bias,
880 whereby the effective weight of each cue is determined by its onset time. For example, previously
881 attended cues may weigh more on value outputs than currently attended cues, all else being

882 equal. We developed a specific method for detecting such persisting value traces, which can be
 883 equally applied to both RNN simulations and monkeys' behavior in the task.

884 We start by re-estimating value profiles, while allowing for value differences between
 885 options that are currently or previously attended (at the time of choice), and having separated
 886 trials by the type of attended cue (reward probability vs magnitude). Let V_{att}^{prob} denote the pseudo-
 887 value function of the attended option when a probability cue is attended at the time of choice,
 888 and V_{unatt}^{prob} that of the other (unattended) option. Let p_{att} and m_{att} be the ranks of the attended
 889 option's probability and magnitude, and p_{unatt} and m_{unatt} those of the unattended option. The
 890 probability of choosing the attended option is given by:

$$891 \quad P(\text{choose attended option}) = \frac{1}{1 + \exp\left(-\left(V_{att}^{prob}(p_{att}, m_{att}) - V_{unatt}^{prob}(p_{unatt}, m_{unatt})\right)\right)} \quad (12)$$

892 This provides a binomial likelihood function for observed choices that are triggered when
 893 a probability cue is attended. To estimate the pseudo-value profiles V_{att}^{prob} and V_{unatt}^{prob} , we use the
 894 same semi-parametric approach as before (see subsection Value profile estimation above). The
 895 pseudo-value profiles V_{att}^{mag} and V_{unatt}^{mag} can be estimated similarly, given observed choices that
 896 are triggered when a magnitude cue is attended.

897 Recall that V_{att}^{prob} (resp., V_{att}^{mag}) is the pseudo-value that ensues from currently attending
 898 a probability (resp., a magnitude) cue, while the magnitude (resp., probability) cue was previously
 899 attended (if ever). To quantify the relative impact of currently and previously attended cues while
 900 marginalizing over cue types, we then combine V_{att}^{prob} and V_{att}^{mag} to form the following average
 901 pseudo-value profile V_{att} :

$$902 \quad V_{att} = \frac{1}{2} \left(V_{att}^{prob} + V_{att}^{mag\top} \right) \quad (13)$$

903 Importantly, V_{att} is a 6×6 pseudo-value profile whose first dimension (columns) spans the
904 rank of the currently attended cue, while its second dimension (rows) spans the rank of the
905 previously attended cue – including the case where it is unknown at the time of choice. Note that
906 rational agents exhibit a strictly symmetric average pseudo-value profile, irrespective of their
907 subjective value profile over the bidimensional space spanned by native option attributes.
908 However, irrational value computations that induce persisting value traces do exhibit
909 asymmetrical V_{att} profiles.

910 To quantify potential asymmetries in V_{att} , we regressed the ranks of the attended and
911 unattended cues onto V_{att} using a generalized linear model:

$$912 \quad V_{att} = \beta_0 + \beta_{att} \times R_{att} + \beta_{unatt} \times R_{unatt} + \varepsilon \quad (14)$$

913 where R_{att} and R_{unatt} denote the ranks of the attended and unattended cue, respectively, each
914 ranging from 1 to 5, whereas β_{att} and β_{unatt} quantify the contribution of attended and unattended
915 cues to the pseudo-value profile. In the Results section, we report the difference $\Delta\beta = \beta_{att} -$
916 β_{unatt} between estimates of attended and unattended cue contributions. If $\Delta\beta < 0$, then the
917 pseudo-value profile V_{att} is asymmetric and more sensitive to the unattended cue. This is the
918 hallmark of a persisting value trace that resists novel (currently attended) information.

919 **Biological constraints**

920 We now derive measures of biological constraint compliance, as can be derived from
921 extensive numerical simulations of trained (rational or irrational) RNN models of the OFC. In what
922 follows, $N = 10$, $|U| = 10,000$ and $T = 4$ are the number of units in the integration layer, the
923 number of simulated input sequences and the number of time steps within a given decision trial.

924 Energetic budget: average network firing rate

925 The brain’s energetic budget constraints are tight. Since action potentials are a major
 926 source of energetic consumption in neurons, we quantified the pseudo-energetic budget \bar{E} of an
 927 RNN in terms of the average activation of RNNs’ integration units, across all units, time steps, and
 928 admissible cue sequences:

$$929 \quad \bar{E} = \frac{1}{N \times |U| \times T} \sum_{i=1}^N \sum_U \sum_{t=1}^T X_2^U(i, t) \quad (15)$$

930 where $X_2^U(i, t)$ is the activity level of unit i at time t in response to the input sequence \mathbf{U} . Since
 931 units’ activation functions are standard sigmoid mappings, $0 \leq \bar{E} \leq 1$.

932 Efficient coding: code redundancy and information transfer rate

933 Efficient coding models suggest that brain circuits with limited neural resources self-
 934 organize to either minimize code redundancy or maximize information transfer rate.

935 First, recall that population codes are redundant if neural elements tend to be active at
 936 the same time, across the range of stimuli that drive their responses. Here, a unit i is deemed
 937 “active” if its output $X_2(i)$ exceeds the a^{th} percentile of its marginal activity distribution (see
 938 below). Let $N_{active}(a, \mathbf{U}, t)$ denote the number of active units at decision time t , for the input
 939 sequence \mathbf{U} , under the exceedance threshold a . The probability that two randomly selected units
 940 are simultaneously active is given by:

$$941 \quad P(\text{co-activation}) = \frac{N_{active}(a, \mathbf{U}, t) \times (N_{active}(a, \mathbf{U}, t) - 1)}{N \times (N - 1)} \quad (16)$$

942 The higher the co-activation probability, the more redundant the neural code. However,
 943 changing the exceedance thresholds modifies the co-activation probability estimate. In analogy
 944 to Receiver Operating Characteristic (ROC) analyses, we thus systematically vary the
 945 exceedance threshold $a \in \{0, 1, \dots, 100\}$, and derive the ensuing probability of co-activation. We

946 then define the code redundancy \bar{R} as the average probability of co-activation over exceedance
 947 thresholds:

$$948 \quad \bar{R} = \frac{1}{101 \times |U| \times T} \sum_{a=0}^{100} \sum_{\mathbf{U}} \sum_{t=1}^T \frac{N_{active}(a, \mathbf{U}, t) \times (N_{active}(a, \mathbf{U}, t) - 1)}{N \times (N - 1)} \quad (17)$$

949 By construction, $0 \leq \bar{R} \leq 1$. When $\bar{R} \approx 0$, code redundancy is minimal, i.e. units almost
 950 never co-activate across trials and decision time steps.

951 Second, we measure information transfer rate in terms of the average entropy of the unit's
 952 output (across integration units). Let $f: x \mapsto y$ be the input-output activation function of neural
 953 net units. At the low noise limit, information transfer rate \bar{I} reduces to the expected, log-
 954 transformed, absolute gradient of the activation function⁶⁴:

$$955 \quad \bar{I} = \mathbb{E} \left[\ln \left| \frac{\partial f}{\partial x}(x) \right| \right] \approx \frac{1}{N \times |U| \times T} \sum_{i=1}^N \sum_{\mathbf{U}} \sum_{t=1}^T \ln \left(X_2^{\mathbf{U}}(i, t) \times (1 - X_2^{\mathbf{U}}(i, t)) \right) \quad (18)$$

956 where the expectation is taken under the distribution of admissible cue sequences \mathbf{U} and we have
 957 used standard results of sigmoidal activation functions. Note that the maximum information
 958 transfer rate is achieved when the distribution of inputs to each integration units exactly matches
 959 the gradient of their activation function (here: $\bar{I} \lesssim -1.38$ nats).

960 *Robustness: excitatory/inhibitory balance and tolerance to damage*

961 In this work, we consider two distinct types of robustness.

962 First, the excitatory/inhibitory balance of a circuit is critical for maintaining its dynamical
 963 stability. Formally, it refers to the relative contribution of excitatory and inhibitory inputs on
 964 features of the circuit's evoked responses (e.g., selective tuning). In electrophysiological studies,
 965 E/I balance is typically evaluated using intracellular conductance estimates across a wide range
 966 of conditions and contexts. Here, we quantify a simple form of structural balance \bar{B} , which we

967 define as the ratio between the number of positive (excitatory) and negative (negative) connection
 968 weights:

$$969 \quad \bar{B} = \frac{\sum_j 1_{\{W(j)>0\}}}{\sum_j 1_{\{W(j)<0\}}} \quad (19)$$

970 where $1_{\{W(j)\leq 0\}}$ is binary indicator that flags whenever the j^{th} entry of the augmented connection
 971 weight matrix $\mathbf{W} = \mathbf{W}_{forward} \cup \mathbf{W}_{recurrent}$ is positive or negative. By construction, $\bar{B} \geq 0$. Note
 972 that we expect rational RNNs to be relatively well balanced ($\bar{B} \approx 1$), because excessive excitatory
 973 or inhibitory connections may induce runaway or saturating activity dynamics, which precludes
 974 accurate value computations (at least in late phases of decision trials).

975 Second, biological systems need to maintain function despite compromised structural
 976 integrity. Here, we define the RNN's tolerance to damage in terms of their achieved rate of rational
 977 choices under varying levels of artificial lesions on the integration layer. Let $n \in \{1, 2, \dots, 9\}$ and
 978 \mathbf{C}_n denote the number of lesioned units in the integration layer, and the $N \times 1$ binary lesion map
 979 vector that indicates the combination of lesioned units, respectively. Artificial lesions are
 980 performed by forcing lesioned units to stay silent across all time steps and trials. Let
 981 $z_{model}(\mathbf{U}, t, \mathbf{C}_n) \in \{0, 1\}$ denote the RNN's simulated choice at time t in response to an input
 982 sequence \mathbf{U} , under a lesion \mathbf{C}_n of its integration layer. Let $z_{rational}(\mathbf{U}, t)$ denote the rational
 983 choice (i.e., the preferred option based upon options' expected value) given the same input
 984 sequence. We define the tolerance to damage \bar{T} as the retained rate of rational choice, averaged
 985 over all possible lesion maps involving 10% to 50% of all units in the integration layer:

$$986 \quad \bar{T} = \frac{5}{2N \times |\mathbf{U}| \times T} \sum_{n=N/10}^{N/2} \frac{1}{\binom{N}{n}} \sum_{\mathbf{C}_n} \sum_{\mathbf{U}} \sum_{t=1}^T 1_{\{z_{model}(\mathbf{U}, t, \mathbf{C}_n) = z_{rational}(\mathbf{U}, t)\}} \quad (20)$$

987 where $1_{\{z_{model}(\mathbf{U}, t, \mathbf{C}_n) = z_{rational}(\mathbf{U}, t)\}}$ is a binary indicator that flags whenever the lesioned RNN
 988 emits a rational choice, and $\binom{N}{n}$ is the number of possible combinations (lesion maps) when

989 lesioning n among N . By construction, $0 \leq \bar{T} \leq 1$. If $\bar{T} \approx 1$, then RNNs' value computations are
990 mostly unaffected by virtual lesions. We note that the results we report in the article are mostly
991 invariant to the chosen range of lesion extent (here: from 10% to 50%).

References

- 992 1. Knudsen, E. B. & Wallis, J. D. Taking stock of value in the orbitofrontal cortex. *Nat. Rev.*
993 *Neurosci.* **23**, 428–438 (2022).
- 994 2. Shidara, M. & Richmond, B. J. Anterior cingulate: single neuronal signals related to degree of
995 reward expectancy. *Science* **296**, 1709–1711 (2002).
- 996 3. Stoll, F. M. & Rudebeck, P. H. Preferences reveal dissociable encoding across prefrontal-
997 limbic circuits. *Neuron* **112**, 2241–2256.e8 (2024).
- 998 4. Kable, J. W. & Glimcher, P. W. The neural correlates of subjective value during intertemporal
999 choice. *Nat. Neurosci.* **10**, 1625–1633 (2007).
- 1000 5. Hare, T. A., Camerer, C. F. & Rangel, A. Self-control in decision-making involves modulation
1001 of the vmPFC valuation system. *Science* **324**, 646–648 (2009).
- 1002 6. Christopoulos, G. I., Tobler, P. N., Bossaerts, P., Dolan, R. J. & Schultz, W. Neural Correlates of
1003 Value, Risk, and Risk Aversion Contributing to Decision Making under Risk. *J. Neurosci.* **29**,
1004 12574–12583 (2009).
- 1005 7. Boorman, E. D., Rushworth, M. F. & Behrens, T. E. Ventromedial Prefrontal and Anterior
1006 Cingulate Cortex Adopt Choice and Default Reference Frames during Sequential Multi-
1007 Alternative Choice. *J. Neurosci.* **33**, 2242–2253 (2013).
- 1008 8. Kahnt, T., Park, S. Q., Haynes, J.-D. & Tobler, P. N. Disentangling neural representations of
1009 value and salience in the human brain. *Proc. Natl. Acad. Sci.* **111**, 5000–5005 (2014).
- 1010 9. Rushworth, M. F., Kolling, N., Sallet, J. & Mars, R. B. Valuation and decision-making in frontal
1011 cortex: one or many serial or parallel systems? *Curr. Opin. Neurobiol.* **22**, 946–955 (2012).
- 1012 10. Lopez-Persem, A. *et al.* Four core properties of the human brain valuation system
1013 demonstrated in intracranial signals. *Nat. Neurosci.* **23**, 664–675 (2020).

- 1014 11. Lebreton, M., Abitbol, R., Daunizeau, J. & Pessiglione, M. Automatic integration of
1015 confidence in the brain valuation signal. *Nat. Neurosci.* **18**, 1159–1167 (2015).
- 1016 12. Ballesta, S., Shi, W., Conen, K. E. & Padoa-Schioppa, C. Values encoded in orbitofrontal
1017 cortex are causally related to economic choices. *Nature* **588**, 450–453 (2020).
- 1018 13. Bartra, O., McGuire, J. T. & Kable, J. W. The valuation system: A coordinate-based meta-
1019 analysis of BOLD fMRI experiments examining neural correlates of subjective value.
1020 *NeuroImage* **76**, 412–427 (2013).
- 1021 14. Camille, N., Griffiths, C. A., Vo, K., Fellows, L. K. & Kable, J. W. Ventromedial Frontal Lobe
1022 Damage Disrupts Value Maximization in Humans. *J. Neurosci.* **31**, 7527–7532 (2011).
- 1023 15. Fellows, L. K. Deciding how to decide: ventromedial frontal lobe damage affects
1024 information acquisition in multi-attribute decision making. *Brain* **129**, 944–952 (2006).
- 1025 16. Fellows, L. K. The role of orbitofrontal cortex in decision making: a component process
1026 account. *Ann. N. Y. Acad. Sci.* **1121**, 421–430 (2007).
- 1027 17. Abitbol, R. *et al.* Neural mechanisms underlying contextual dependency of subjective
1028 values: converging evidence from monkeys and humans. *J. Neurosci. Off. J. Soc. Neurosci.*
1029 **35**, 2308–2320 (2015).
- 1030 18. Padoa-Schioppa, C. Neuronal Origins of Choice Variability in Economic Decisions.
1031 *Neuron* **80**, 1322–1336 (2013).
- 1032 19. Padoa-Schioppa, C. Range-Adapting Representation of Economic Value in the
1033 Orbitofrontal Cortex. *J. Neurosci.* **29**, 14004–14014 (2009).
- 1034 20. Brochard, J. & Daunizeau, J. Efficient value synthesis in the orbitofrontal cortex explains
1035 how loss aversion adapts to the ranges of gain and loss prospects. *eLife* **13**, e80979 (2024).
- 1036 21. Louie, K., Khaw, M. W. & Glimcher, P. W. Normalization is a general neural mechanism
1037 for context-dependent decision making. *Proc. Natl. Acad. Sci.* **110**, 6139–6144 (2013).
- 1038 22. Lopez-Persem, A., Domenech, P. & Pessiglione, M. How prior preferences determine
1039 decision-making frames and biases in the human brain. *eLife* **5**, e20317 (2016).

- 1040 23. Lim, S.-L., O'Doherty, J. P. & Rangel, A. The Decision Value Computations in the vmPFC
1041 and Striatum Use a Relative Value Code That is Guided by Visual Attention. *J. Neurosci.* **31**,
1042 13214–13223 (2011).
- 1043 24. Johnson, D. D. P. & Fowler, J. H. The evolution of overconfidence. *Nature* **477**, 317–320
1044 (2011).
- 1045 25. Sharot, T. The optimism bias. *Curr. Biol.* **21**, R941–R945 (2011).
- 1046 26. Dewan, A. & Neligh, N. Estimating information cost functions in models of rational
1047 inattention. *J. Econ. Theory* **187**, 105011 (2020).
- 1048 27. Gershman, S. J., Horvitz, E. J. & Tenenbaum, J. B. Computational rationality: A converging
1049 paradigm for intelligence in brains, minds, and machines. *Science* **349**, 273–278 (2015).
- 1050 28. Glimcher, P. W. Efficiently irrational: deciphering the riddle of human choice. *Trends*
1051 *Cogn. Sci.* **26**, 669–687 (2022).
- 1052 29. Padamsey, Z. & Rochefort, N. L. Paying the brain's energy bill. *Curr. Opin. Neurobiol.* **78**,
1053 102668 (2023).
- 1054 30. Harris, J. J., Jolivet, R. & Attwell, D. Synaptic energy use and supply. *Neuron* **75**, 762–777
1055 (2012).
- 1056 31. Isler, K. & van Schaik, C. P. The Expensive Brain: A framework for explaining evolutionary
1057 changes in brain size. *J. Hum. Evol.* **57**, 392–400 (2009).
- 1058 32. Heldstab, S. A., Isler, K., Graber, S. M., Schuppli, C. & van Schaik, C. P. The economics of
1059 brain size evolution in vertebrates. *Curr. Biol.* **32**, R697–R708 (2022).
- 1060 33. Chalk, M., Marre, O. & Tkačik, G. Toward a unified theory of efficient, predictive, and
1061 sparse coding. *Proc. Natl. Acad. Sci.* **115**, 186–191 (2018).
- 1062 34. Lewicki, M. S. Efficient coding of natural sounds. *Nat. Neurosci.* **5**, 356–363 (2002).
- 1063 35. Olshausen, B. A. & Field, D. J. Natural image statistics and efficient coding. *Netw.*
1064 *Comput. Neural Syst.* **7**, 333–339 (1996).
- 1065 36. Barlow, H. Redundancy reduction revisited. *Network* **12**, 241–253 (2001).

- 1066 37. Louie, K. & Glimcher, P. W. Efficient coding and the neural representation of value. *Ann.*
1067 *N. Y. Acad. Sci.* **1251**, 13–32 (2012).
- 1068 38. Vishwanath, A. A. *et al.* Mitochondrial Ca²⁺ efflux controls neuronal metabolism and
1069 long-term memory across species. Preprint at <https://doi.org/10.1101/2024.02.01.578153>
1070 (2024).
- 1071 39. Padamsey, Z., Katsanevaki, D., Dupuy, N. & Rochefort, N. L. Neocortex saves energy by
1072 reducing coding precision during food scarcity. *Neuron* **110**, 280-296.e10 (2022).
- 1073 40. Attwell, D. & Laughlin, S. B. An Energy Budget for Signaling in the Grey Matter of the
1074 Brain. *J. Cereb. Blood Flow Metab.* **21**, 1133–1145 (2001).
- 1075 41. Levy, W. B. & Baxter, R. A. Energy efficient neural codes. *Neural Comput.* **8**, 531–543
1076 (1996).
- 1077 42. Chen, L., Li, X., Tjia, M. & Thapliyal, S. Homeostatic plasticity and excitation-inhibition
1078 balance: The good, the bad, and the ugly. *Curr. Opin. Neurobiol.* **75**, 102553 (2022).
- 1079 43. Navlakha, S., He, X., Faloutsos, C. & Bar-Joseph, Z. Topological properties of robust
1080 biological and computational networks. *J. R. Soc. Interface* **11**, 20140283 (2014).
- 1081 44. Kitano, H. Towards a theory of biological robustness. *Mol. Syst. Biol.* **3**, 137 (2007).
- 1082 45. Srinivasan, S. & Stevens, C. F. Robustness and fault tolerance make brains harder to
1083 study. *BMC Biol.* **9**, 46 (2011).
- 1084 46. Chen, G., Kang, B., Lindsey, J., Druckmann, S. & Li, N. Modularity and robustness of
1085 frontal cortical networks. *Cell* **184**, 3717-3730.e24 (2021).
- 1086 47. Ikeda, M. *et al.* Circuit Degeneracy Facilitates Robustness and Flexibility of Navigation
1087 Behavior in *C.elegans*. Preprint at <https://doi.org/10.1101/385468> (2018).
- 1088 48. Hoefler, T., Alistarh, D., Ben-Nun, T., Dryden, N. & Peste, A. Sparsity in Deep Learning:
1089 Pruning and growth for efficient inference and training in neural networks. *J. Mach. Learn.*
1090 *Res.* **22**, 1–124 (2021).

- 1091 49. Pessiglione, M. & Daunizeau, J. Bridging across functional models: The OFC as a value-
1092 making neural network. *Behav. Neurosci.* **135**, 277–290 (2021).
- 1093 50. Hunt, L. T. *et al.* Triple dissociation of attention and decision computations across
1094 prefrontal cortex. *Nat. Neurosci.* **21**, 1471–1481 (2018).
- 1095 51. McGinty, V. B. & Lupkin, S. M. Behavioral read-out from population value signals in
1096 primate orbitofrontal cortex. *Nat. Neurosci.* **26**, 2203–2212 (2023).
- 1097 52. Fine, J. M. *et al.* Abstract Value Encoding in Neural Populations But Not Single Neurons.
1098 *J. Neurosci.* **43**, 4650–4663 (2023).
- 1099 53. Hunt, L. T., Dolan, R. J. & Behrens, T. E. J. Hierarchical competitions subserving multi-
1100 attribute choice. *Nat. Neurosci.* **17**, 1613–1622 (2014).
- 1101 54. Stone, C., Mattingley, J. B. & Rangelov, D. On second thoughts: changes of mind in
1102 decision-making. *Trends Cogn. Sci.* **26**, 419–431 (2022).
- 1103 55. O’Doherty, J. P., Rutishauser, U. & Iigaya, K. The hierarchical construction of value. *Curr.*
1104 *Opin. Behav. Sci.* **41**, 71 (2021).
- 1105 56. Padoa-Schioppa, C. & Conen, K. E. Orbitofrontal Cortex: A Neural Circuit for Economic
1106 Decisions. *Neuron* **96**, 736–754 (2017).
- 1107 57. Padoa-Schioppa, C. & Assad, J. A. Neurons in the orbitofrontal cortex encode economic
1108 value. *Nature* **441**, 223–226 (2006).
- 1109 58. Yamins, D. L. K. *et al.* Performance-optimized hierarchical models predict neural
1110 responses in higher visual cortex. *Proc. Natl. Acad. Sci. U. S. A.* **111**, 8619–8624 (2014).
- 1111 59. Caucheteux, C. & King, J. R. Brains and algorithms partially converge in natural language
1112 processing. *Commun. Biol.* **2022 51 5**, 1–10 (2022).
- 1113 60. Ballesta, S. & Padoa-Schioppa, C. Economic Decisions through Circuit Inhibition. *Curr.*
1114 *Biol.* **29**, 3814–3824.e5 (2019).
- 1115 61. Jamadar, S. D., Behler, A., Deery, H. & Breakspear, M. The metabolic costs of cognition.
1116 *Trends Cogn. Sci.* **0**, (2025).

- 1117 62. Yang, D.-P., Zhou, H.-J. & Zhou, C. Co-emergence of multi-scale cortical activities of
1118 irregular firing, oscillations and avalanches achieves cost-efficient information capacity.
1119 *PLoS Comput. Biol.* **13**, e1005384 (2017).
- 1120 63. Hirokawa, J., Vaughan, A., Masset, P., Ott, T. & Kepecs, A. Frontal cortex neuron types
1121 categorically encode single decision variables. *Nature* **576**, 446–451 (2019).
- 1122 64. Nadal, J.-P. & Parga, N. Nonlinear neurons in the low-noise limit: a factorial code
1123 maximizes information transfer. *Netw. Comput. Neural Syst.* **5**, 565–581 (1994).
- 1124 65. Moneta, N., Grossman, S. & Schuck, N. W. Representational spaces in orbitofrontal and
1125 ventromedial prefrontal cortex: task states, values, and beyond. *Trends Neurosci.* **47**, 1055–
1126 1069 (2024).
- 1127 66. Landron, T. *et al.* Dissociation of Value and Confidence Signals in the Orbitofrontal
1128 Cortex during Decision-Making: An Intracerebral Electrophysiology Study in Humans. *J.*
1129 *Neurosci.* **45**, (2025).
- 1130 67. Hunt, L. T. *et al.* Mechanisms underlying cortical activity during value-guided choice.
1131 *Nat. Neurosci.* **15**, 470–476 (2012).
- 1132 68. Maierov, V. & Pinkus, A. Lower bounds for approximation by MLP neural networks.
1133 *Neurocomputing* **25**, 81–91 (1999).
- 1134 69. Guliyev, N. J. & Ismailov, V. E. Approximation capability of two hidden layer feedforward
1135 neural networks with fixed weights. *Neurocomputing* **316**, 262–269 (2018).
- 1136 70. Caucheteux, C., Gramfort, A. & King, J.-R. Evidence of a predictive coding hierarchy in
1137 the human brain listening to speech. *Nat. Hum. Behav.* **7**, 430–441 (2023).
- 1138 71. Daunizeau, J., Adam, V. & Rigoux, L. VBA: A Probabilistic Treatment of Nonlinear Models
1139 for Neurobiological and Behavioural Data. *PLOS Comput. Biol.* **10**, e1003441 (2014).
- 1140 72. Daunizeau, J. The variational Laplace approach to approximate Bayesian inference.
1141 Preprint at <https://doi.org/10.48550/arXiv.1703.02089> (2018).

- 1142 73. Blei, D. M., Kucukelbir, A. & McAuliffe, J. D. Variational Inference: A Review for
1143 Statisticians. *J. Am. Stat. Assoc.* **112**, 859–877 (2017).
- 1144 74. Kriegeskorte, N. Relating population-code representations between man, monkey, and
1145 computational models. *Front. Neurosci.* **3**, 363–373 (2009).