

Sources of predictive information in dynamical neural networks

Madhavun Candadai^{a,b} and Eduardo J. Izquierdo^{a,b,1}

^aProgram in Cognitive Science, Indiana University, Bloomington, U.S.A.; ^bSchool of Informatics, Computing and Engineering, Indiana University, Bloomington, U.S.A.

This manuscript was compiled on December 23, 2019

1 **Behavior involves the ongoing interaction between an organism and**
2 **its environment. One of the prevailing theories of adaptive behav-**
3 **ior is that organisms are constantly making predictions about their**
4 **future environmental stimuli. However, how they acquire that pre-**
5 **dictive information is still poorly understood. Two complementary**
6 **mechanisms have been proposed: predictions are generated from**
7 **an agent's internal model of the world or predictions are extracted di-**
8 **rectly from the environmental stimulus. In this work, we demonstrate**
9 **that predictive information, measured using mutual information, can-**
10 **not distinguish between these two kinds of systems. Furthermore,**
11 **we show that predictive information cannot distinguish between or-**
12 **ganisms that are adapted to their environments and random dynam-**
13 **ical systems exposed to the same environment. To understand the**
14 **role of predictive information in adaptive behavior, we need to be**
15 **able to identify where it is generated. To do this, we decompose in-**
16 **formation transfer across the different components of the organism-**
17 **environment system and track the flow of information in the system**
18 **over time. To validate the proposed framework, we examined it on**
19 **a set of computational models of idealized agent-environment sys-**
20 **tems. Analysis of the systems revealed three key insights. First,**
21 **predictive information, when sourced from the environment, can be**
22 **reflected in any agent irrespective of its ability to perform a task. Sec-**
23 **ond, predictive information, when sourced from the nervous system,**
24 **requires special dynamics acquired during the process of adapting**
25 **to the environment. Third, the magnitude of predictive information**
26 **in a system can be different for the same task if the environmental**
27 **structure changes.**

neural coding | predictive information | information theory

1 **P**redictive coding is emerging as a strong candidate for its
2 ability to provide a general framework for understanding
3 the neural basis of behavior (1–4). The idea is that organ-
4 isms encode information about future environmental stimuli
5 in their neural activity based on their knowledge of the envi-
6 ronment. Intuitively, an organism that is able to predict the
7 consequences of its action on its future sensory experiences
8 is more likely to be adapted to its environment. There are
9 two prominent research fronts that study the role of predic-
10 tive coding in behavior: the hierarchical predictive processing
11 framework (5, 6) and the efficient coding principle (7, 8). These
12 two fronts are complementary because they address different
13 aspects of how a nervous system acquires predictive informa-
14 tion. The hierarchical predictive processing framework focuses
15 on how predictions are generated in the organism's brain. The
16 efficient coding principle focuses on how the nervous system
17 extracts predictive information from environmental stimuli.
18 Both theories have been supported by experimental evidence,
19 primarily in the visual and auditory systems (9–12).

20 In living organisms, predictive information is likely acquired
21 from a dynamically changing contribution of the environment

22 and the agent's own internal dynamics (2). Consequently,
23 although different systems may be equally predictive about
24 their future stimuli, the operation of their nervous systems
25 may be entirely different. Therefore, understanding the role
26 of predictive information in behavior requires that the source
27 of information is identified. In this paper, we address the
28 following questions. How do we identify the source of predictive
29 information and study its dynamics during a behavior? Does
30 tracking the source of predictive information better explain an
31 agent's ability to perform a task? What are the factors that
32 influence the source and magnitude of predictive information
33 encoded in a neural network?

34 In the first part of this paper, we demonstrate that predic-
35 tive information will generate indistinguishable results for
36 systems that are at the two extremes of potential agent-
37 environment interaction: a system whose only source of pre-
38 dictive information is the nervous system and a system whose
39 only source of predictive information is the environmental
40 stimuli. In order to better understand how the nervous system
41 generates predictive information, we propose that it is essen-
42 tial to decompose information transfer across the different
43 components of the system and to track the flow of information
44 in the agent-environment system over time. The principal con-
45 tribution of this paper is an information-theoretic framework
46 to quantify the contributions from the nervous system and
47 the contributions from the environmental stimuli to the total
48 predictive information in an agent. First, we decompose the
49 total predictive information in the neural system into infor-
50 mation that was uniquely transferred from each source. In
51 order to do this, we employ multivariate extensions to infor-

Significance Statement

An organism's ability to predict the consequences of its actions on future stimuli is considered a strong indicator of its environmental adaptation. However, in highly structured natural environments, to what extent does an agent have to develop specialized mechanisms to generate predictions? To study this, we present an information theoretic framework to infer the source of predictive information in an organism: extrinsically from the environment or intrinsically from the agent. We find that predictive information extracted from the environment can be reflected in any agent and is therefore not a good indicator of behavioral performance. Studying the flow of predictive information over time across the organism-environment system enables us to better understand its role in behavior.

M.C. and E.J. designed research, performed research, contributed new analytic tools, analyzed data, and wrote the paper.

The authors declare no conflict of interest.

¹To whom correspondence should be addressed. E-mail: edizquie@indiana.edu

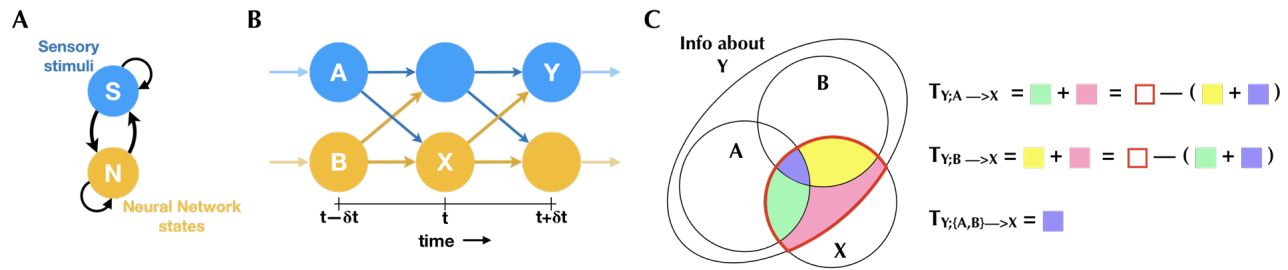


Fig. 1. Predictive information source estimation based on idealized agent-environment interaction. [A] Sensory stimuli (S) and neural activity (N) are two coupled dynamical systems. [B] Agent-environment interaction unrolled over time. X represents current neural activity, $N(t)$, Y , future environmental state, $S(t + \delta t)$, and A and B represent the sources, namely past neural activity $N(t - \delta t)$ and past environmental state, $s(t - \delta t)$ respectively. [C] Partial information diagram for calculating the sources of predictive information in an agent-environment system. The total information that X has about Y is a combination of information that is available uniquely from A alone (green), uniquely from B alone (yellow), synergistically from their combination $[A, B]$ (pink), and redundantly from both of them (purple). PID allows us to measure information transfer using these components. Alternatively, they can also be measured by estimating the total redundant information from both sources combined (red) and removing the information from the other source.

52 mation theory (13). Second, we unroll information over time
 53 to backtrack the origin of the source of predictive information
 54 and how they change over time. To validate the proposed
 55 theoretical framework, we examine it on a set of computa-
 56 tional models of agent-environment systems, where the agent
 57 is driven by a dynamical recurrent neural network (14, 15).
 58 The systems have been deliberately designed so that the source
 59 of predictive information is tractable and manipulable. We
 60 demonstrate how the proposed framework correctly reveals
 61 different sources of predictive information in systems with
 62 otherwise similar amounts of predictive information. Ulti-
 63 mately, we demonstrate how revealing the flow of information
 64 across the agent-environment system can help us to better
 65 understand the mechanisms underlying predictive coding.

66 Predictive information is studied in living organisms be-
 67 cause it is considered a signature of their adaptive capaci-
 68 ties (5, 8, 9). In the second part of this paper we study the
 69 relationship between a system’s ability to perform a task and
 70 its predictive information. In order to do this, we turn to
 71 a computational model of an agent that is required to pro-
 72 cess the received stimulus from the environment and make a
 73 decision based on it. Specifically, we study predictive infor-
 74 mation in the context of a relational categorization task (16, 17).
 75 We generate model systems that are adapted to their envi-
 76 ronment and yet remain tractable to analysis by optimizing
 77 dynamical recurrent neural networks using an evolutionary
 78 algorithm to perform the task (18, 19). We then proceed to
 79 analyze the resulting systems using predictive information and
 80 we compare the results against that of random systems that
 81 cannot solve the task. Counterintuitively, we observe that
 82 predictive information in trained neural networks is similar
 83 to predictive information in random neural networks. This
 84 suggests that predictive information alone is not sufficient to
 85 distinguish between living organisms that are adapted to their
 86 environments and non-adaptive systems. The rest of the paper
 87 focuses on an analysis of optimized and random systems using
 88 the framework proposed. Altogether, we demonstrate that
 89 decomposing predictive information across the components
 90 of an agent-environment system, and unrolling it over time
 91 reveals its true nature.

Identifying the source of predictive information

92 Predictive information is the information encoded in neural
 93 activity about its future stimulus. Formally, it is defined as
 94 mutual information between current neural activity (N_t) and
 95 the stimulus at a future time ($S_{t'}$) (9, 20–23), according to:
 96

$$I(S_{t'}, N_t) = \sum_{s_{t'}, n_t} P_N(n_t) P(s_{t'} | n_t) \log_2 \frac{P(s_{t'} | n_t)}{P_S(s_{t'})} \quad [1]$$

97 where $t' = t + \delta t$ with $\delta t > 0$, P_S is the distribution of environ-
 98 mental stimuli, P_N is the distribution of neural activity across
 99 the entire experiment, $P(s_{t'} | n_t)$ is the conditional probability
 100 that the stimulus is s at a future time t' given that we have
 101 observed a neural activity of n at time t . When this measure
 102 is estimated using the stimulus and neural activity across all
 103 data points separated in time by some δt , it is a measure of
 104 reduction in uncertainty in future stimulus given the current
 105 neural activity.
 106

107 The presence of predictive information in a neural network
 108 suggests there is a source where this information was gener-
 109 ated. In an idealized agent-environment system (Fig. 1A),
 110 the source of information can be either the neural activity
 111 in the previous time step, the environmental stimulus in the
 112 previous time step, or both (Fig. 1B). Measuring predictive
 113 information as defined in equation 1 requires that we exam-
 114 ine two variables: current neural activity (N_t , henceforth X)
 115 and future stimulus ($S_{t+\delta t}$, henceforth Y). Identifying the
 116 source of this predictive information requires that we exam-
 117 ine two additional variables: past neural network activity
 118 ($N_{t-\delta t}$, henceforth A) and past stimulus ($S_{t-\delta t}$, henceforth B).
 119 Such an analysis requires that we adopt multivariate exten-
 120 sions to information theory. We focus specifically on Partial
 121 Information Decomposition (PID) (13), a method for decom-
 122 posing multivariate mutual information into combinations of
 123 unique, redundant and synergistic contributions, as well as
 124 measures of information gain, loss and transfer (13, 24–32).
 125 In order to identify the source of predictive information, we
 126 can decompose the total information that the current neural
 127 activity has about the future stimulus into three components:
 128 (a) information uniquely transferred from past environmental
 129 stimulus, $T_{Y;A \rightarrow X}$; (b) information uniquely transferred from
 130 past neural network activity, $T_{Y;B \rightarrow X}$; and (c) information

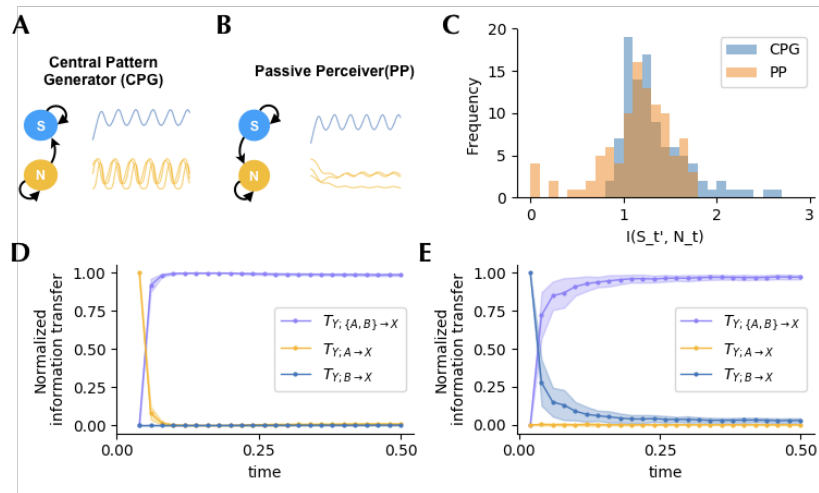


Fig. 2. Predictive information in systems on the extremes of the range of possible agent-environment interactions [A] Schematic and traces of a Central Pattern Generator (CPG) that influences the environment through intrinsically generated oscillations. [B] Schematic and traces of a Passive Perceiver (PP) that is driven by oscillatory inputs from the environment (in this case, by the environmental signals recorded from the CPGs) [C] Estimating total predictive information as shown in equation 1 shows that CPG and PP models encode similar amounts of predictive information about environmental state in the next time-step. [D] Decomposing that total information into information that came from the environment and the neural network consistently showed that information about the next time-step in the CPG originated in the neural network (yellow) before becoming redundant (purple) as the environment and the neural network synchronize. [E] Conversely, with PPs, the environment was consistently shown to be the source of information (blue) before they environment and neural network synchronize and become redundant (purple).

131 redundantly transferred from past environment stimulus and
 132 past neural network activity, $T_{Y;\{A,B\} \rightarrow X}$, according to:

$$\begin{aligned}
 T_{Y;A \rightarrow X} &= \Pi_R(Y; \{[A, B], X\}) - \Pi_R(Y; \{B, X\}) \\
 T_{Y;B \rightarrow X} &= \Pi_R(Y; \{[A, B], X\}) - \Pi_R(Y; \{A, X\}) \quad [2] \\
 T_{Y;\{A,B\} \rightarrow X} &= \Pi_R(Y; \{A, B, X\})
 \end{aligned}$$

134 where $\Pi_R(Y; \{A_1, A_2, \dots, A_k\})$ is the redundant information that
 135 random variables A_1 through A_k have about the random
 136 variable Y and $[A, B]$ refers to a random variable that is a
 137 concatenation of A and B . In words, information about Y
 138 transferred uniquely from source A to X is estimated as the
 139 total redundant information from the combined sources $[A, B]$
 140 minus the information that is redundant with the other source
 141 B . This decomposition of the total information into different
 142 contributions is typically represented using a PI-decomposition
 143 diagram (Fig. 1C). Several approaches have been proposed to
 144 measure redundant information, Π_R (24, 33, 34). Here, we
 145 use I_{min} because this is the only approach that can guarantee
 146 non-negative information decomposition in a system with four
 147 random variables, as is the case here.

148 During the course of behavior, the flow of information in a
 149 system changes over time (35, 36). In order to understand the
 150 source of predictive information for any agent-environment sys-
 151 tem, it is not enough to decompose information from multiple
 152 sources; we must also track its flow of information over time.
 153 Although information theoretic measures are typically averaged
 154 over time, the measures described above can be unrolled
 155 over time (36, 37). This is done by measuring information
 156 transfer at each time-point using data collected across several
 157 trials thereby allowing us to study the dynamics of predictive
 158 information sources.

159 Disparate systems with similar predictive information

160 Neural systems can be predictive in fundamentally different
 161 ways: they can generate predictive information internally or
 162 they can extract it from environmental stimulus. We use
 163 computational models of two extreme conditions where the
 164 ground-truth predictive information source is known to be
 165 the environment in one condition and the neural network
 166 in the other, to demonstrate that (a) predictive information
 167 cannot distinguish between these different kinds of systems

168 and (b) it is only through decomposing the information across
 169 sources and unrolling over time that we can distinguish the
 170 two systems based on their operation. The two conditions we
 171 consider are agent-environment interactions at two extremes of
 172 the range of possible interactions: a central pattern generator
 173 (CPG) and a passive perceiver (PP). In the CPG condition,
 174 the neural network influences the environment by producing
 175 spontaneous oscillatory activity but receives no input from
 176 the environment (Fig. 2A). In the PP condition, the neural
 177 network is influenced by input from the environment, but it
 178 does not affect the environment (Fig. 2B). We evolved 100
 179 different dynamical recurrent neural network CPGs, and in
 180 each case, we fed the sum of the neurons' outputs to the
 181 environment (Fig. S1A,B). For the PPs, we generated 100
 182 random neural networks and fed them an oscillatory input.
 183 In order to provide the same distribution of activity as the
 184 CPG condition, we provided the random neural networks
 185 with the same oscillatory environmental signal that CPGs
 186 generated (Fig. S1C). The environmental signal and neural
 187 data were recorded from each instance for 500 trials where,
 188 in each trial the environment started with a different initial
 189 condition. Although, the environmental signal and the neural
 190 activity exhibit oscillatory activity in both conditions, the key
 191 difference in the operation of these systems is that in the CPGs
 192 the neural network drives its own activity and in the PP, the
 193 environment drives the neural network. Therefore, the neural
 194 network is the source of predictive information in the CPGs
 195 and the environment is the source of predictive information in
 196 the PPs.

197 As a first step in the analysis of these two systems, we used
 198 the recorded data to measure predictive information in the
 199 neural network about the environmental signal in the next
 200 time-step ($\delta t = 0.02s$). To calculate predictive information,
 201 data distributions were constructed using all tuples of neural
 202 activity at time t and environmental signal at time $t + \delta t$,
 203 averaged across time and trials. The analysis revealed that the
 204 neural networks, in these two otherwise diametrically opposed
 205 systems, encoded similar levels of information about stimulus
 206 in the next time step (Fig. 2C). From this first experiment, we
 207 conclude that predictive information is not sufficient to distin-
 208 guish systems that generate their own predictive information
 209 from systems that encode the information available from the
 210 environmental stimuli.

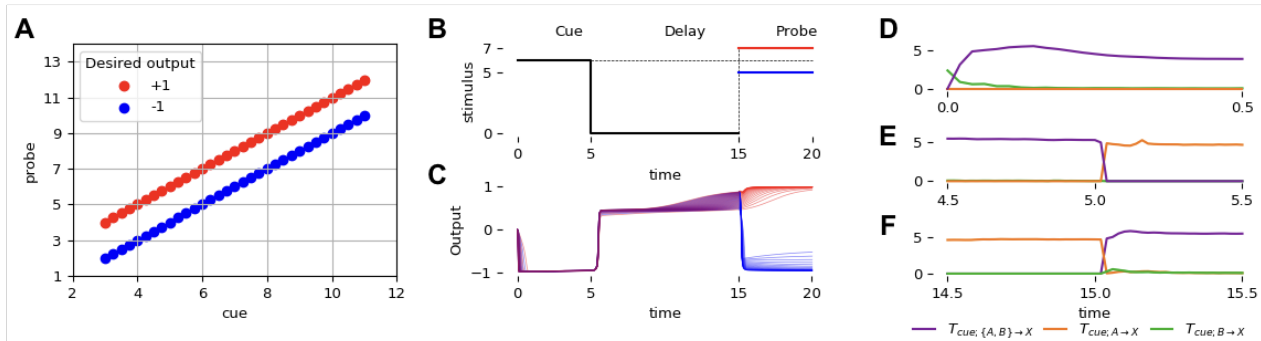


Fig. 3. Predictive information source dynamics with structured stimuli. [A] Distribution of cue and corresponding probes in the relational categorization task. For each cue, the probe can be one of two values: greater, $cue + 1$, or lesser, $cue - 1$, with the expected outputs of +1 (red) and -1 (blue) respectively. [B] One trial of the relational categorization task. The cue stimulus is presented till $t=5$, followed by a delay period with no stimulus ($t=5$ to $t=15$) and then a probe that is greater (red) or lesser (blue) than the cue is provided. [C] Behavior of the best out of 100 dynamical neural networks optimized to perform this task showing perfect categorization of the relational value from 35 trials where the probe was greater (red) and 35 where the probe was lesser (blue). [D] Dynamics of information about the cue during the cue stage show information uniquely provided by the environment (green) initially, but becoming redundantly available in the neural network and environment (purple) as it encoded the cue. [E] Towards the end of the cue stage, information is entirely redundant (purple). When the stimulus stops being provided at $t=5$, the neural network is the unique source of information about the cue (orange). [F] Dynamics of information about the cue just before the probe arrives showing that the neural network continues to retain information about the cue (orange). At $t=15$, when the probe is provided, information quickly becomes redundant (purple) denoting that the probe has information about the cue.

211 To understand what makes these two neural systems differ-
 212 ent, it is necessary to identify the source of their predictive
 213 information. As a next step in our analysis, we decomposed
 214 the information in the neural system about the future stimuli
 215 across the different possible sources and we unrolled the analy-
 216 sis over time. At each time-point, we measured information in
 217 the neural network about the environmental signal in the next
 218 time-step that was uniquely transferred from the environment,
 219 uniquely transferred from the neural network and redundantly
 220 from both.

221 In the CPG condition, since the neural networks are not
 222 influenced by the environment (Fig. 2A), the only source of
 223 information about the future environmental signal is from the
 224 neural network itself. Accordingly, the dynamics of information
 225 transfer for CPG systems reveals correctly that the neural
 226 network is the source of predictive information (Fig. 2D). At
 227 the start of the interaction between agent and environment,
 228 the neural network uniquely transfers information about the
 229 future environmental state to the environment. Following
 230 that, the environment quickly becomes synchronized with the
 231 neural activity. This means that the state of the environment
 232 becomes informative of its own future state. This results in
 233 the environment and the neural network becoming redundant
 234 sources of predictive information. Crucially, however, the
 235 environment never provides any unique information to the
 236 neural network about its future stimulus.

237 In the PP condition, since the neural networks are driven
 238 by the environment (Fig. 2B), the only source of information
 239 about the future environmental signal is the stimulus from the
 240 environment itself. Accordingly, the dynamics of information
 241 transfer for PP systems reveals correctly that the environment
 242 is the source of predictive information (Fig. 2E). As opposed
 243 to the CPG systems, at the start of the interaction between
 244 the neural network and the environment, it is the environment
 245 that transfers unique information to the neural network. Sub-
 246 sequently, and similarly to the CPG condition, as the state
 247 of the neural network begins to encode the information from
 248 the environmental stimulus, the predictive information is re-

dundantly transferred by both the neural network and the
 environmental stimulus. Consistent with our expectation, the
 neural network never provides any unique information to itself
 about the future of the stimulus.

In summary, in this section we show that predictive in-
 formation alone cannot distinguish between two extremely
 different kinds of neural systems, both of which encode pre-
 dictive information about the future of the environment. This
 is because when the entire time course of the data is consid-
 ered, the environment and neural network are synchronized
 for a majority of the time. Information uniquely transferred
 from any source is only detectable within a short time window
 before they synchronize. In this section, we have shown that
 decomposing information across sources and unrolling over
 time allows us to study information source dynamics at every
 perturbation to the agent-environment interaction and hence
 reveals the source of predictive information.

Predictive information with structured stimuli

The natural environment is not uniformly random but is in
 fact highly structured with spatial and temporal regularities
 (2, 38, 39). This structure is reflected in the stimulus
 that agents receive from the environment. Accordingly, this
 is emulated in most preparations in neuroscience, where a
 neural system is presented with artificial stimuli with some
 underlying structure designed by the experimenter. We posit
 that the structure in the environment will strongly influence
 the amount of predictive information encoded by the neural
 network and its sources. In order to study this, we examined
 the flow of information in a neural network model trained to
 solve a relational categorization task.

Relational categorization is the ability to discriminate ob-
 jects based on the relative value of their attributes (16, 17).
 This task allows us to specify the inherent structure in the
 environment by changing the distribution of objects whose
 attributes are compared thus making it especially suited for
 studying the influence of environmental structure on predictive
 information. It involves providing the neural network with

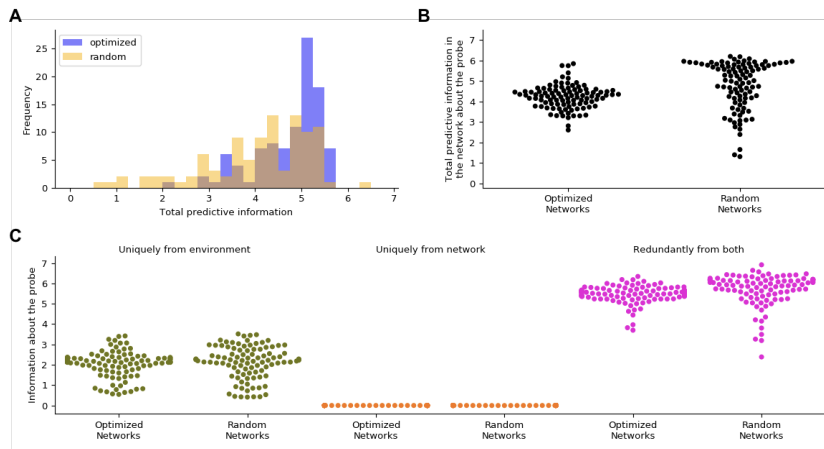


Fig. 4. Comparison of predictive information sources in optimized and random neural networks. [A] Total predictive information estimated by averaging over the entire course of the task is similar in random and optimized neural networks. [B] Total predictive information about the probe averaged across the cue stage of the task, is the same in random and optimized neural networks. [C] Decomposition of that total predictive information showing that information about the probe in both random and optimized neural networks was from the environment (green), eventually becoming redundant as they both encoded the cue stimulus (pink). The neural network had no role to play in its encoding of predictive information about the probe during the cue stage (orange).

286 stimuli across three stages: cue, delay, and probe. In the
 287 cue stage, the neural network is provided with a stimulus of
 288 specific magnitude for a duration of time. This is followed
 289 by a delay stage, where no stimulus is provided. Finally, in
 290 the probe stage, the neural network is provided with a second
 291 stimulus. The magnitudes for the cue and probe stage stimuli
 292 are picked from a predesignated distribution (Fig. 3A). It is
 293 this distribution that defines the structure in the environment.
 294 For this study, we design it such that the stimulus in the
 295 probe stage can have a magnitude that is one of two values:
 296 smaller ($cue - 1$) or larger ($cue + 1$) than the stimulus
 297 provided during the cue stage (Fig. 3B). The goal of the neural
 298 network in this task to perform a relational categorization
 299 of “greater than” or “lesser than” by producing an output
 300 of $+1$ or -1 respectively, during the probe phase. This task
 301 has been widely studied in a variety of contexts including in
 302 humans (40), pigeons (41), rats (42), insects (43), as well as
 303 using computational models (44, 45).

304 In this section, we show results from analysis of neural
 305 networks performing the relational categorization task. We
 306 demonstrate that decomposing information across the sources
 307 and unrolling over time reveals that the environment is struc-
 308 tured by appropriately attributing the observed predictive
 309 information to either the environment or the dynamics of
 310 the neural network. Furthermore, we demonstrate that en-
 311 coding predictive information alone is not indicative of task
 312 performance and that the magnitude and source of predic-
 313 tive information can change during the course of a behavior
 314 depending on environmental structure and neural network
 315 dynamics.

316 **Characterizing information source dynamics in the best opti-
 317 mized neural network.** Dynamical recurrent neural networks
 318 were optimized using an evolutionary algorithm to perform
 319 the relational categorization task. A total of 100 independent
 320 evolutionary runs yielded an ensemble of 100 different neural
 321 networks that could successfully perform the task (Fig. S2A).
 322 The best neural network from this ensemble achieved a perfor-
 323 mance of 93.12%. Although this neural network correctly
 324 classified all probes, the performance score was not perfect
 325 due to slight deviations in the output (Fig. 3C).

326 In order to better understand how a neural network per-
 327 formed this task, we can characterize the flow of information
 328 across the agent-environment system. To this end, we decom-
 329 posed the total information that the best neural network from

the ensemble had about the cue into information uniquely
 transferred from the environment, uniquely transferred from
 the neural network, and redundantly from both, during the
 course of the task. During the cue stage, the environment
 was initially the unique source of information about the cue
 (Fig. 3D). As the neural network encoded the stimulus, the
 source became redundant. During the delay stage, the envi-
 ronment ceases to be a source of information. As the neural
 network had already encoded information about the cue, it
 becomes the unique source (Fig. 3E). Crucially, the neural
 network preserves this information throughout the delay stage.
 Finally, during the probe stage, the environment once again be-
 comes a source, and therefore the source is redundant (Fig. 3F).
 Note that when the environment provides the probe stimulus
 it became the source of information about the cue. Since the
 neural network already contained information about the cue,
 the neural network and the environment both redundantly act
 as the source.

348 As explained previously, predictive information in this task
 349 arises from the relationship between cue and probe stimuli.
 350 Encoding information about the cue automatically results in
 351 encoding information about the probe (and vice versa). This
 352 is because knowing the cue significantly reduces uncertainty
 353 about the probe; the probe can only be one of two values
 354 given a cue. Predictive information that the neural network
 355 has about the probe and its sources is qualitatively similar to
 356 the information it has encoded about the cue (Fig. S3A). The
 357 neural network encodes information about the probe stimulus
 358 upon receiving the cue, and retains that predictive informa-
 359 tion during the delay stage. This is merely a consequence of
 360 encoding and retaining the cue. The entire ensemble of neural
 361 networks optimized to perform this task consistently exhibit
 362 this phenomenon of encoding information about the probe
 363 transferred uniquely from the cue stimulus (Fig. S3B) and is
 364 even robust to noise in the neural network (Fig. S5).

365 **Environmental regularities induces predictive information in
 366 any neural network.** Since optimized neural networks encode
 367 information about the probe merely by encoding the cue, does
 368 any neural network that encodes the cue also encode informa-
 369 tion about the probe, and therefore have similar predictive
 370 information? In order to study this, we created 100 random
 371 neural networks and presented them with the same task. Al-
 372 though these neural networks were not able to perform the
 373 relational categorization task (Fig. S2B), they encoded similar

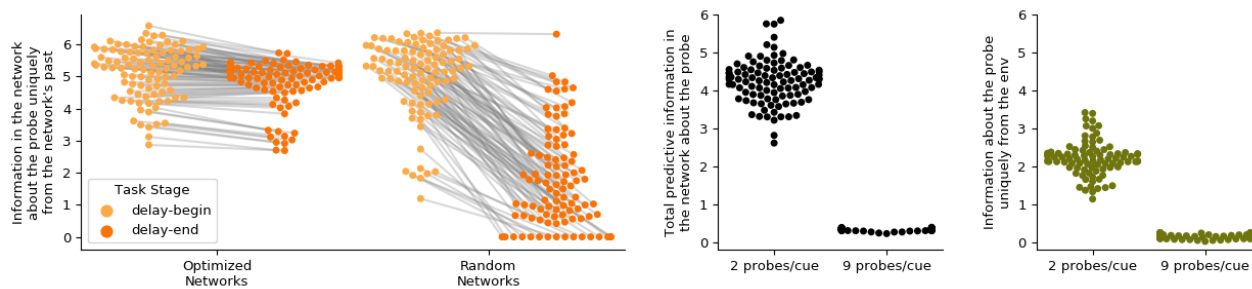


Fig. 5. Influence of neural network and environmental properties on predictive Information [A] Both random and optimized neural networks have similar levels of information about the probe at the beginning of the delay stage (light orange), but unlike optimized neural networks, random neural networks lose that information by the end of the delay stage (dark orange). [B] Total predictive information in the optimized neural networks about the probe during the cue stage showed a significant drop upon changing environmental statistics from 2 probes/cue to 9 probes/cue. [C] Drop in total information show in B can be attributed to the drop in information uniquely from the environment about the probe in the 9 probes/cue setting.

374 amounts of total predictive information as the trained neural
 375 networks (Fig. 4A). Specifically, they encode the same amount
 376 of information about the probe during the cue stage (Fig. 4B).
 377 Furthermore, decomposing that information revealed that the
 378 information originated from the environmental stimulus and
 379 that the neural network dynamics had no role in its encoding
 380 of predictive information in both random and optimized neural
 381 networks (Fig. 4C). Thus, predictive information alone is not
 382 sufficient to distinguish neural networks optimized to perform
 383 specific tasks from random neural networks that are merely
 384 reflecting the information provided by the environment.

385 **Information decomposition distinguishes between random**
 386 **and optimized neural networks.** Unlike CPG and PP that were
 387 distinguished based on having different information sources,
 388 random and optimized neural networks in the relational cate-
 389 gorization task have the same information sources. Even un-
 390 der this condition, decomposing the total information across
 391 sources and unrolling over time helps distinguish them by
 392 revealing differences in the magnitude of information trans-
 393 ferred from each source over time. Specifically, predictive
 394 information sourced by the neural network during the delay
 395 stage is markedly different between random and optimized
 396 neural networks. As discussed in the previous section, opti-
 397 mized neural networks preserve information about the cue
 398 (and hence predictive information about the probe) during
 399 the delay stage. In contrast, random neural networks tend
 400 to lose that information. As a consequence, the amount of
 401 unique information provided by the neural network at the end
 402 of the delay period is higher for the trained neural networks
 403 than for the random neural networks (Fig. 5A). This differ-
 404 ence disappears when information is measured across time, and
 405 can only be observed by unrolling it over time.

406 **Statistics of the environment influences magnitude of predic-**
 407 **tive information.** Encoding the cue results in encoding infor-
 408 mation about the probe in this task because of the relationship
 409 between them. How does changing this relationship impact
 410 predictive information in the neural networks? In order to
 411 study this, without changing the nature of the relational cate-
 412 gorization task we merely changed the structure in the envi-
 413 ronment. This was achieved by modifying the task such that
 414 the probe could be one of 9 possible values for a given cue,
 415 rather than one of two possible values (Fig. S4B). Reduction
 416 in uncertainty about the probe's value given the cue is now

417 much less compared to the original environmental structure
 418 (Fig. S4D,E). This will be reflected in the information that
 419 the cue can provide about the probe. However, this came at
 420 no cost to performance because the neural networks were still
 421 encoding the cue just as well. The same ensemble of optimized
 422 neural networks were able to perform this task successfully
 423 without any more training (Fig. S4E). Information dynamics
 424 was then measured using data recorded under this 9-probe
 425 condition. Measuring the total information in the neural net-
 426 work during the cue stage about the probe revealed that there
 427 was significantly less information in the neural network in 9
 428 probes per cue condition (Fig. 5B). The reduction in total
 429 predictive information can be wholly attributed to the reduction
 430 in information about the probe (Fig. 5C). Thus, differences
 431 in environmental structure can result in significantly different
 432 amounts of predictive information encoded in neural networks
 433 without any behavioral differences.

434 Discussion

435 The study of predictive coding and its relevance to behavior has
 436 been studied from multiple perspectives in the literature with
 437 regards to the source of information: predictive information
 438 can be generated by the neural network (5, 6) and predictive
 439 information can be provided by the environment (7, 20). In
 440 this work, using computational models where the ground-truth
 441 about the source of information was known, we demonstrate
 442 that predictive information can originate from either the envi-
 443 ronment or the neural network or both, and that the source
 444 of information can dynamically change during the course of a
 445 behavior. In order to do this, we first presented a theoretical
 446 framework based on multivariate information theory that al-
 447 lows us to infer the source of predictive information and its
 448 dynamics. This involved decomposing the total information
 449 that neural networks encode about a future stimulus into infor-
 450 mation transferred uniquely from the neural network, uniquely
 451 from the environment and redundantly from both sources. We
 452 validated this framework using the CPG and PP models where
 453 information is known to originate from the neural network
 454 and the environment respectively. Second, using the more
 455 structured relational categorization task, we demonstrated
 456 that (a) amount of predictive information encoded in a neural
 457 network is not indicative of its performance; (b) the source
 458 of information about a future stimulus can change during the
 459 course of the task; and (c) the source of information about a

460 future stimulus can change within the same task depending
 461 on the regularities of the environment. Thus, predictive infor-
 462 mation might be necessary but is not sufficient to explain the
 463 neural basis of a behavior. Decomposing information across
 464 sources and studying its dynamics over time takes us one step
 465 further in understanding the role of predictive information in
 466 a behavior.

467 The framework presented here for inferring the source of
 468 predictive information takes us beyond general correlations
 469 that information theoretic measures are known to capture by
 470 capturing the effects of perturbation on the neural system.
 471 Identifying the sources of predictive information requires that
 472 the system under study be perturbed. The presentation, re-
 473 moval or sudden change of a stimulus is a perturbation. This
 474 causes the system to break the redundant encoding observed
 475 in a steady-state. It is during such a perturbation that we
 476 can use partial information decomposition to determine the
 477 source of information in a coupled system. Once the neural
 478 network and the environment settle into the next steady-state
 479 after the transient due to the perturbation, information once
 480 again becomes redundant between them. Thus, through the
 481 combination of information decomposition, time-unrolling and
 482 perturbation we are able to infer the ground-truth causal
 483 influences in the models we have analyzed.

484 The framework presented here can be applied to experi-
 485 mental data across multiple scales. In fact, it can be applied
 486 to any time-series data spanning multiple trials corresponding
 487 to several perturbations from the steady state. However, in
 488 this work, we focus on open-loop systems. Specifically, we
 489 focus on agent-environment systems where the agent influ-
 490 ences its environment or where the agent is influenced by the
 491 environment. Such an open-loop setup is typical in experi-
 492 ments in neuroscience, where the subject receives a stimulus,
 493 but does not have the ability to influence the future stimulus
 494 through their state or actions. In natural behavior, the agent
 495 and environment are in closed-loop interaction. The analysis
 496 of closed-loop systems introduces an added complexity. The
 497 regularities of the environment can be generated by the regu-
 498 larities of the neural network's dynamics and vice-versa. As
 499 a result, the distribution of environmental stimuli and the
 500 distribution of the neural activity are dependent on each other,
 501 unlike the open-loop setup where one of them is independent
 502 of the other. As it is, the framework requires that one of
 503 the distributions be fixed across time in order to make fair
 504 comparisons of information at different time-points. Future
 505 work in this direction will involve extending the framework
 506 and designing the experimental setting that would allow us to
 507 infer the source of predictive information in a freely moving
 508 animal.

509 Materials and Methods

510 In the agent-environment models used throughout this paper, the
 511 agents were modeled using dynamical recurrent neural networks.
 512 The parameters of the neural network were optimized using an
 513 evolutionary algorithm such that it was able to perform the required
 514 task. In this section, we specify implementation details about the
 515 neural network model, the tasks, and the optimization algorithm.

516 **Neural network model.** A Continuous-Time Recurrent Neural Net-
 517 work (CTRNN) was used as the model neural network (14, 15).
 518 The neural network consisted of three layers: the input layer which
 519 was connected by a set of feed-forward weights to the interneuron
 520 layer; the interneuron layer was a CTRNN which fed into the output

521 layer; the output layer produced the output of the neural network
 522 which was given by a weighted combination of the interneurons'
 523 output. The dynamics of each interneuron was governed based on
 524 state equations given by

$$\tau_i \frac{dy_i}{dt} = -y_i + \sum_{j=1}^N w_{ij} o_j + w_i^{in} I \quad [3]$$

$$o_j = \sigma(y_j + \theta_j) \quad [4]$$

527 where y_i refers to the internal state of neuron i ; τ_i , the time-constant;
 528 w_{ij} , the strength of connection from neuron j to neuron i ; o_j , the
 529 output of the neuron; I , the input and w_i^{in} , the weight from the
 530 input to the neuron. Based on the state of the neuron its output
 531 is given by equation 4, where $\sigma()$ refers to the sigmoid activation
 532 function given by $\sigma(x) = 1/(1 + e^{-x})$, and θ_j refers to the bias
 533 of neuron j . The output of the network at any time t , $O(t)$, is
 534 estimated as a weighted sum of the outputs of each neuron (weights
 535 given by w_i^o), passed through a sigmoid function and scaled to be
 536 in the range $[-1, 1]$.

$$O(t) = 2 * \sigma \left(\sum_{i=1}^N w_i^o o_i(t) \right) - 1 \quad [5]$$

538 All neural networks described in this paper were made up of
 539 $N = 3$ neurons. The tunable parameters of such a model include
 540 the weights between the neurons (w_{ij}), the input weights (w_i^{in}), the
 541 output weights (w_i^o), time-constants (τ_i) and biases (θ_{ij}) of each
 542 neuron. The model was simulated using Euler integration with a
 543 step-size of 0.02.

544 **CPG task.** The neural network model described above is capable
 545 of intrinsically producing oscillations. To create Central Pattern
 546 Generators (CPGs), neural networks were optimized to produce
 547 oscillations from a range of initial conditions. The neural network
 548 was started at 100 different initial conditions by systematically
 549 setting the neuron outputs in the range $[0, 1]$. For each condition, the
 550 neural activity was recorded for 10 simulation seconds. The ability
 551 to generate oscillations was assessed by measuring the absolute
 552 difference in each neuron's as well as the neural network's output
 553 in consecutive time-steps across all time-points in a trial, and then
 554 across trials. The neural network's output was fed to an environment
 555 governed by

$$\tau \frac{ds}{dt} = -s + O \quad [6]$$

557 where s refers to the state of the environment, τ refers to its time-
 558 constant which was set to 0.5, and O refers to the output of the
 559 neural network given by equation 5.

560 **Relational categorization task.** We adapted the relational catego-
 561 rization task to provide neural networks with structured stimuli
 562 (16, 17, 44). This task involves first providing the neural network
 563 with a cue stimulus in the range $[3, 11]$ for 5 units of time. This is
 564 followed by a delay period when no stimulus is provided for 10 units
 565 of time. Finally, a probe stimulus that is of magnitude greater or
 566 less than the cue is provided for 5 units of time. The goal of the task
 567 is for the neural network to distinguish probes that were larger than
 568 the cue or smaller than the cue, by producing an output of $+1$ or -1
 569 respectively. In the first version of this task, the probe can take one
 570 of only two values, either $cue + 1$ or $cue - 1$. In the second version of
 571 the task, the probe can take any value in $[3, 11]$. While the goal of
 572 the task remains the same in both versions, the distribution of the
 573 probes given the cue, and therefore information that the cue gives
 574 about the probe is significantly different (Fig. S4). Performance of
 575 a neural network in this task was estimated by measuring absolute
 576 deviation of the network's output from the desired output of $+1$
 577 or -1 during the probe stage. Time-averaged deviation was also
 578 averaged across all trials of cue-probe values, to obtain a score in
 579 the range $[0, 1]$.

580 **Neural network optimization.** Neural network models described pre-
 581 viously were optimized to perform the relational categorization
 582 task using an evolutionary algorithm (46, 47). This optimization

583 methodology involves instantiating a population of 100 random
584 solutions that evolves over several generations to produce solutions
585 capable of performing the task. A generation is defined as the pro-
586 cess of creating a new population of solutions that has improved in
587 “fitness” (task performance) from the last. Each solution, referred
588 to as a genotype, is an N dimensional vector corresponding to the
589 parameters to be optimized. The parameters were encoded to be in
590 the range $[0, 1]$ and scaled to produce the neural network that the
591 genotype encoded. In each generation, the fitness of every genotype
592 is evaluated and a new population is created using a fitness-based
593 selection and mutation strategy as follows: The genotypes that
594 perform in the top 1% were retained as is for the next generation.
595 The rest of the individuals were created by selecting two genotypes
596 preferentially in proportion to their fitness and combining them. To
597 these offspring, Gaussian mutation noise with mean 0 and standard
598 deviation 0.01 was added before being added to the population
599 of genotypes for the next generation. After a fixed number of
600 generations, the best individual in the population was selected as
601 the representative solution from that optimization run. 100 such
602 runs were conducted to obtain an ensemble of 100 neural network
603 models that successfully performed each task. For the relational
604 categorization task, optimization was carried out for 500 genera-
605 tions. In the case of the CPG task, at the end of 50 generations
606 the optimization process was terminated and deemed successful if
607 the best agent in the population reached a fitness of 30 or greater.
608 This was repeated until 100 CPGs were produced. See supporting
609 information (Figs. S1 and S2) for training curves, behavior of best
610 optimized neural network, distribution of fitness of best models
611 from 100 runs, and sample neural traces.

612 **Random neural networks.** Matched random neural networks were
613 created for the relational categorization task by shuffling the pa-
614 rameters of the optimized neural networks. All parameter groups,
615 namely time-constants, input weights, recurrent weights, output
616 weights, and biases were randomly shuffled within themselves rather
617 than across groups. Thus, the ranges of parameters were preserved
618 in each group but their associations with neurons were randomly
619 shuffled.

620 **Measuring information transfer.** To identify the source of informa-
621 tion over time, information transfer measures were estimated inde-
622 pendently at each time point. For any given time step, data for
623 environmental stimulus at the previous time step, neural activity
624 of previous time step, current neural activity, and stimulus at a
625 future time step, was collected across multiple trials. Probability
626 densities were estimated from this data using a kernel density esti-
627 mation technique known as average shifted-histograms (48) with 7
628 shifted binnings of 100 bins along each dimension of the data space.
629 These probability density estimates were then used to measure the
630 redundant information terms in equation 2. Similar results were ob-
631 served with 5 and 11 shifts and with 50 and 200 bins per dimension
632 (Fig. S6). All information theoretic quantities were estimated from
633 raw data using the *infotheory* package (49).

- 634 1. A Clark, Whatever next? predictive brains, situated agents, and the future of cognitive science. *Behav. brain sciences* **36**, 181–204 (2013).
- 635 2. Y Huang, RP Rao, Predictive coding. *Wiley Interdiscip. Rev. Cogn. Sci.* **2**, 580–593 (2011).
- 636 3. RP Rao, DH Ballard, Predictive coding in the visual cortex: a functional interpretation of some
637 extra-classical receptive-field effects. *Nat. neuroscience* **2**, 79 (1999).
- 638 4. MV Srinivasan, SB Laughlin, A Dubs, Predictive coding: a fresh view of inhibition in the retina.
639 *Proc. Royal Soc. London. Ser. B. Biol. Sci.* **216**, 427–459 (1982).
- 640 5. K Friston, S Kiebel, Predictive coding under the free-energy principle. *Philos. Transactions*
641 *Royal Soc. B: Biol. Sci.* **364**, 1211–1221 (2009).
- 642 6. K Friston, The free-energy principle: a unified brain theory? *Nat. reviews neuroscience* **11**,
643 127 (2010).
- 644 7. W Bialek, I Nemenman, N Tishby, Predictability, complexity, and learning. *Neural computation*
645 **13**, 2409–2463 (2001).
- 646 8. S Still, Information-theoretic approach to interactive learning. *EPL (Europhysics Lett.)* **85**,
647 28005 (2009).
- 648 9. SE Palmer, O Marre, MJ Berry, W Bialek, Predictive information in a sensory population. *Proc.*
649 *Natl. Acad. Sci.* **112**, 6908–6913 (2015).
- 650 10. KS Chen, CC Chen, C Chan, Characterization of predictive behavior of a retina by mutual
651 information. *Front. computational neuroscience* **11**, 66 (2017).
- 652 11. ZC Chao, K Takaura, L Wang, N Fujii, S Dehaene, Large-scale cortical networks for hierar-
653 chical prediction and prediction error in the primate brain. *Neuron* **100**, 1252–1266 (2018).
- 654 12. AJ Sederberg, JN MacLean, SE Palmer, Learning to make external sensory stimulus pre-
655 dictions using internal correlations in populations of neurons. *Proc. Natl. Acad. Sci.* **115**,
656 1105–1110 (2018).
- 657

- 583 13. PL Williams, RD Beer, Nonnegative decomposition of multivariate information. *arXiv preprint*
584 *arXiv:1004.2515* (2010).
- 585 14. Ki Funahashi, Y Nakamura, Approximation of dynamical systems by continuous time recur-
586 rent neural networks. *Neural networks* **6**, 801–806 (1993).
- 587 15. RD Beer, On the dynamics of small continuous-time recurrent neural networks. *Adapt. Behav.*
588 **3**, 469–509 (1995).
- 589 16. D Gentner, KJ Kurtz, Relational categories. (2005).
- 590 17. AB Markman, CH Stiwell, Role-governed categories. *J. Exp. & Theor. Artif. Intell.* **13**, 329–
591 358 (2001).
- 592 18. RD Beer, JC Gallagher, Evolving dynamical neural networks for adaptive behavior. *Adapt.*
593 *behavior* **1**, 91–122 (1992).
- 594 19. D Floreano, P Dürri, C Mattiussi, Neuroevolution: from architectures to learning. *Evol. intelli-*
595 *gence* **1**, 47–62 (2008).
- 596 20. W Bialek, *Biophysics: searching for principles*. (Princeton University Press), (2012).
- 597 21. F Rieke, D Warland, RdR Van Steveninck, WS Bialek, et al., *Spikes: exploring the neural*
598 *code*. (MIT press Cambridge) Vol. 7, (1999).
- 599 22. TM Cover, JA Thomas, *Elements of information theory*. (John Wiley & Sons), (2012).
- 600 23. CE Shannon, A mathematical theory of communication. *Bell system technical journal* **27**,
601 379–423 (1948).
- 602 24. N Bertschinger, J Rau, E Olbrich, J Jost, N Ay, Quantifying unique information. *Entropy* **16**,
603 2161–2183 (2014).
- 604 25. SP Faber, NM Timme, JM Beggs, EL Newman, Computation is concentrated in rich clubs of
605 local cortical networks. *Netw. Neurosci.* **3**, 384–404 (2019).
- 606 26. RG James, N Barnett, JP Crutchfield, Information flows? a critique of transfer entropies.
607 *Phys. review letters* **116**, 238701 (2016).
- 608 27. RG James, CJ Ellison, JP Crutchfield, Anatomy of a bit: Information in a time series observa-
609 tion. *Chaos: An Interdiscip. J. Nonlinear Sci.* **21**, 037109 (2011).
- 610 28. R James, J Crutchfield, Multivariate dependence beyond shannon information. *Entropy* **19**,
611 531 (2017).
- 612 29. J Lizier, N Bertschinger, J Jost, M Wibral, Information decomposition of target effects from
613 multi-source interactions: perspectives on previous, current and future work (2018).
- 614 30. N Timme, W Alford, B Flecker, JM Beggs, Synergy, redundancy, and multivariate information
615 measures: an experimentalist’s perspective. *J. computational neuroscience* **36**, 119–140
616 (2014).
- 617 31. M Wibral, V Priesemann, JW Kay, JT Lizier, WA Phillips, Partial information decomposition as
618 a unified approach to the specification of neural goal functions. *Brain cognition* **112**, 25–38
619 (2017).
- 620 32. PL Williams, RD Beer, Generalized measures of information transfer. *arXiv preprint*
621 *arXiv:1102.1507* (2011).
- 622 33. RG James, J Emenheiser, JP Crutchfield, Unique information via dependency constraints. *J.*
623 *Phys. A: Math. Theor.* **52**, 014002 (2018).
- 624 34. V Griffith, C Koch, Quantifying synergistic mutual information in *Guided Self-Organization:*
625 *Inception*. (Springer), pp. 159–190 (2014).
- 626 35. EJ Izquierdo, PL Williams, RD Beer, Information flow through a model of the *c. elegans*
627 klinotaxis circuit. *PLoS one* **10**, e0140397 (2015).
- 628 36. PL Williams, RD Beer, Information dynamics of evolved agents in *International Conference*
629 *on Simulation of Adaptive Behavior*. (Springer), pp. 38–49 (2010).
- 630 37. PL Williams, Ph.D. thesis (PhD thesis, Indiana University) (2011).
- 631 38. H Barlow, The exploitation of regularities in the environment by the brain. *Behav. Brain Sci.*
632 **24**, 602–607 (2001).
- 633 39. D Graham, D Field, Statistical regularities of art images and natural scenes: Spectra, sparse-
634 ness and nonlinearities. *Spatial vision* **21**, 149–164 (2007).
- 635 40. KJ Kurtz, O Boukrina, Learning relational categories by comparison of paired examples in
636 *Proceedings of the Annual Meeting of the Cognitive Science Society*. Vol. 26, (2004).
- 637 41. S Willis, Relational learning in pigeons? *The Q. J. Exp. Psychol. Sect. B* **52**, 31–52 (1999).
- 638 42. EL Saldanha, ME Bitterman, Relational learning in the rat. *The Am. J. Psychol.* **64**, 37–53
639 (1951).
- 640 43. M Giurfa, S Zhang, A Jenett, R Menzel, MV Srinivasan, The concepts of ‘sameness’ and
641 ‘difference’ in an insect. *Nature* **410**, 930 (2001).
- 642 44. PL Williams, RD Beer, M Gasser, An embodied dynamical approach to relational categoriza-
643 tion in *Proceedings of the Annual Meeting of the Cognitive Science Society*. Vol. 30, (2008).
- 644 45. E Izquierdo-Torres, I Harvey, Learning to discriminate between multiple possible environ-
645 ments: an imprinting scenario in *Memory and Learning Mechanisms in Autonomous Robots*
646 *Workshop (ECAL 2005)*. (2005).
- 647 46. M Mitchell, *An introduction to genetic algorithms*. (MIT press), (1998).
- 648 47. DE Goldberg, JH Holland, Genetic algorithms and machine learning. (1988).
- 649 48. DW Scott, Averaged shifted histograms: effective nonparametric density estimators in several
650 dimensions. *The Annals Stat.*, 1024–1040 (1985).
- 651 49. M Candadai, EJ Izquierdo, infotheory: A c++/python package for multivariate information
652 theoretic analysis. *arXiv preprint arXiv:1907.02339* (2019).
- 653 654 655 656 657

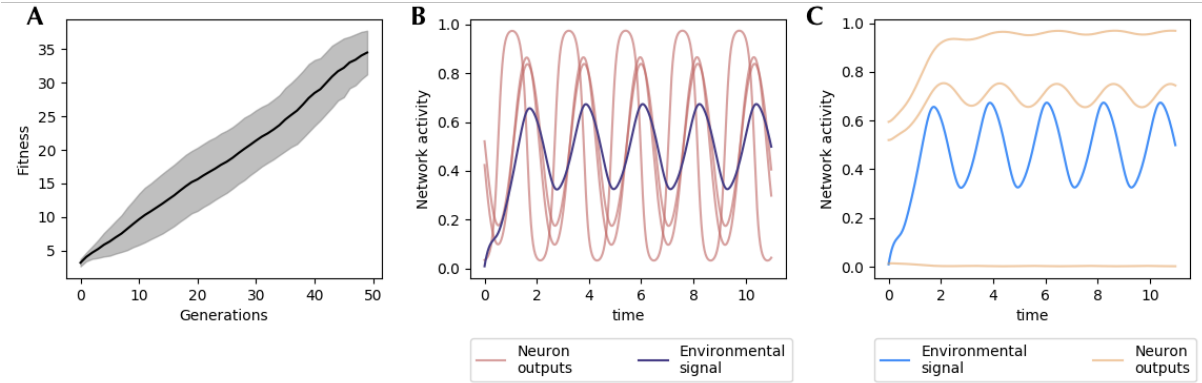


Fig. S1. Optimization and neural traces of CPG and PP. [A] Fitness over time for 100 valid runs of optimizing a CPG model. Only runs that achieved a fitness greater than 30 were deemed valid. [B] Neural traces from one trial of the best CPG demonstrating that all neurons (red) as well as the neural network output (blue) oscillate. [C] Neural traces (orange) when the output from the CPG shown in panel B was fed to a random neural network in the PP condition demonstrating input driven oscillation in the random neural network.

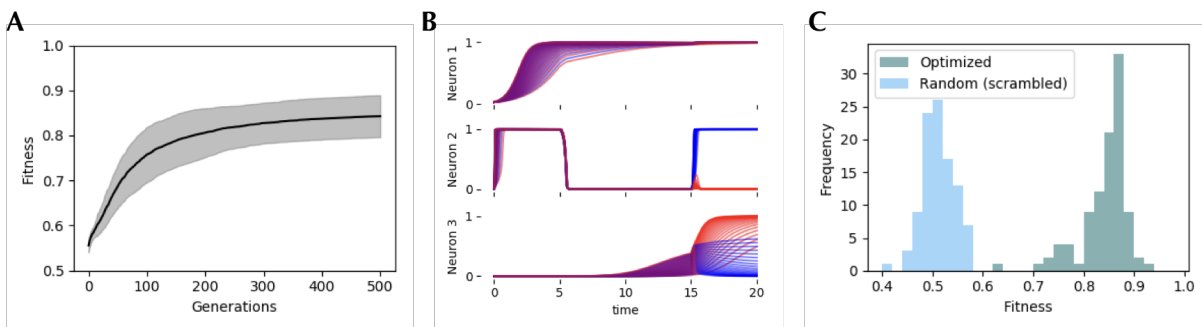


Fig. S2. Optimizing neural networks to perform relational categorization. [A] 100 independent runs all converged to near-perfect performance with deviation from a perfect score only due to small deviations from expected output and not mis-categorization. [B] Neural activity in the CTRNN of the best optimized agent over 35 trials where probe was larger than the cue (red) and 35 trials where the probe was lesser than the cue (blue). [C] Neural networks whose weights and time-constants were scrambled lost their ability to perform the task.

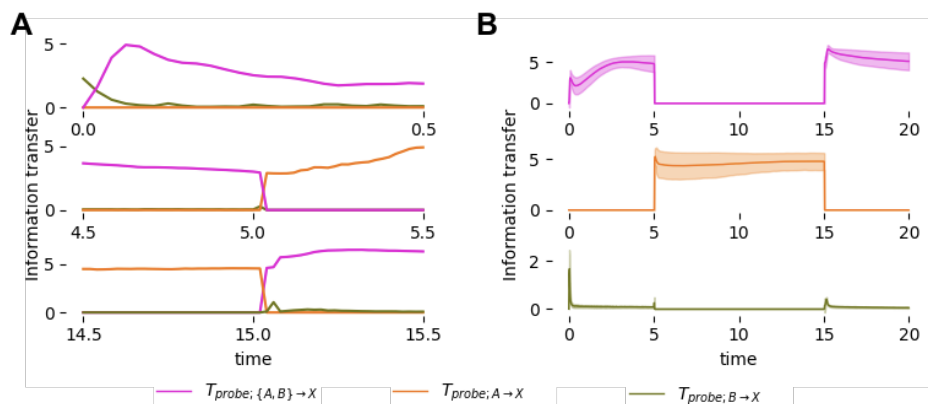


Fig. S3. Predictive information source dynamics is consistent and similar with information about the probe. [A] At the start of the cue stage (top), information about the probe arrives from the environment (green) as the cue is provided, and becomes redundant as the cue is encoded (pink). Towards the start of the delay stage (middle), the neural network becomes the source of information about the probe (orange) as it retains information about the cue, and since the environment ceases to provide that information. As the probe is provided (bottom), the environment once again becomes a source of information in addition to the neural network and they are both redundantly sources of information (pink) [B] Predictive information source dynamics are consistent across all 100 optimized neural networks during all three stages of the task. Their mean value is shown in bold and the shaded region represents one standard deviation around it.

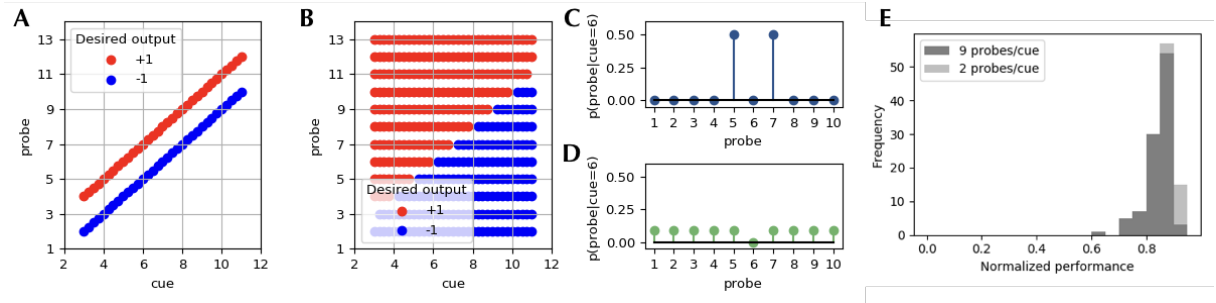


Fig. S4. Different environmental structures within the relational categorization task [A] Relational categorization task with highly structured stimuli; for each cue probe is one of two possible values. [B] Relational categorization task with minimal structure in stimuli; probe can be one of 9 values for a given cue. [C] Conditional probability of probes given a cue for environmental structure shown in panel A, demonstrating the significant reduction in uncertainty of the probe given the cue. [D] Conditional probability of probe values given a cue under the environmental structure in panel B shows that probe values still have a nearly uniform distribution, and hence very less reduction in uncertainty. [E] Neural networks optimized to perform under the distribution shown in panel A perform just as well under the distribution shown in panel B.

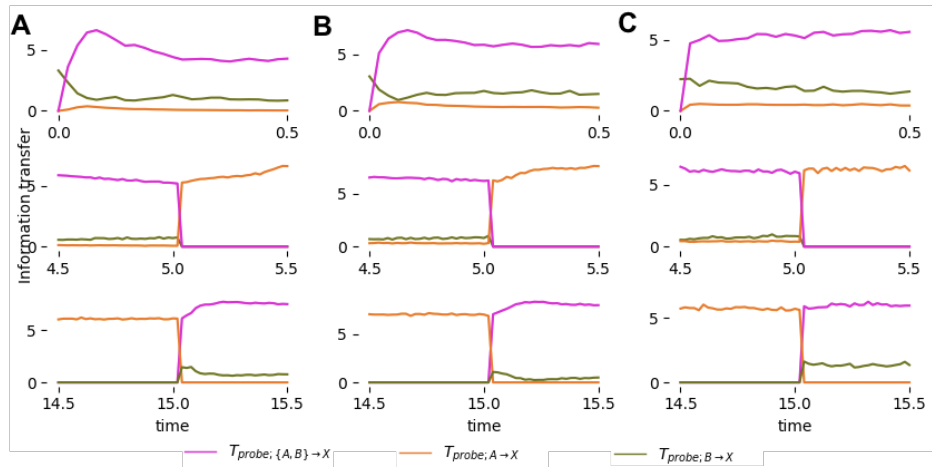


Fig. S5. Inferring the source of predictive information is robust to zero-mean Gaussian noise with standard deviation [A] 0.01, [B] 0.05 and [C] 0.1. Results are qualitatively similar to results from fig. S3A for cue (top row), delay (middle row) and probe (bottom row) stages of the task.

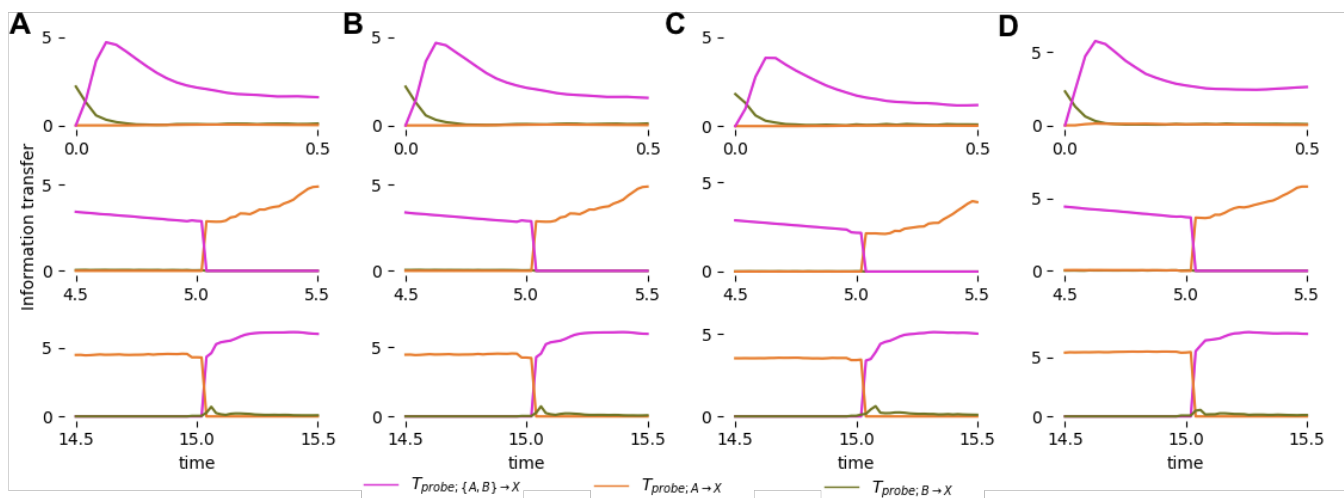


Fig. S6. Inferring the source of predictive information with different binning and shifted-histograms. Results are qualitatively similar to results from fig. S3A after changing [A] number of shifted bins to 3 [B] number of shifted bins to 11 [C] number of bins per dimension to 50 and [D] number of bins per dimension to 200, for cue (top row), delay (middle row) and probe (bottom row) stages of the task.