



OPEN

Sources of predictive information in dynamical neural networks

Madhavun Candadai^{1,2} & Eduardo J. Izquierdo^{1,2}✉

Behavior involves the ongoing interaction between an organism and its environment. One of the prevailing theories of adaptive behavior is that organisms are constantly making predictions about their future environmental stimuli. However, how they acquire that predictive information is still poorly understood. Two complementary mechanisms have been proposed: predictions are generated from an agent's internal model of the world or predictions are extracted directly from the environmental stimulus. In this work, we demonstrate that predictive information, measured using bivariate mutual information, cannot distinguish between these two kinds of systems. Furthermore, we show that predictive information cannot distinguish between organisms that are adapted to their environments and random dynamical systems exposed to the same environment. To understand the role of predictive information in adaptive behavior, we need to be able to identify where it is generated. To do this, we decompose information transfer across the different components of the organism-environment system and track the flow of information in the system over time. To validate the proposed framework, we examined it on a set of computational models of idealized agent-environment systems. Analysis of the systems revealed three key insights. First, predictive information, when sourced from the environment, can be reflected in any agent irrespective of its ability to perform a task. Second, predictive information, when sourced from the nervous system, requires special dynamics acquired during the process of adapting to the environment. Third, the magnitude of predictive information in a system can be different for the same task if the environmental structure changes.

Predictive coding is emerging as a strong candidate for its ability to provide a general framework for understanding the neural basis of behavior¹⁻⁴. The idea is that organisms encode information about future environmental stimuli in their neural activity based on their knowledge of the environment. Intuitively, an organism that is able to predict the consequences of its action on its future sensory experiences is more likely to be adapted to its environment. There are two prominent research fronts that study the role of predictive coding in behavior: the hierarchical predictive processing framework^{5,6} and the efficient coding principle^{7,8}. These two fronts are complementary because they address different aspects of how a nervous system acquires predictive information. The hierarchical predictive processing framework focuses on how predictions are generated in the organism's brain. The efficient coding principle focuses on how the nervous system extracts predictive information from environmental stimuli. Both theories have been supported by experimental evidence, primarily in the visual and auditory systems⁹⁻¹².

In living organisms, predictive information is likely acquired from a dynamically changing contribution of the environment and the agent's own internal dynamics². Consequently, although different systems may be equally predictive about their future stimuli, the operation of their nervous systems may be entirely different. Understanding the role of predictive information in behavior requires that the source of information is identified. In this paper, we address the following questions. How do we identify the source of predictive information and study its dynamics during a behavior? Does tracking the source of predictive information better explain an agent's ability to perform a task? What are the factors that influence the source and magnitude of predictive information encoded in a neural network?

In the first part of this paper, we demonstrate that predictive information will generate indistinguishable results for systems that are at the two extremes of potential agent-environment interaction: a system whose only source of predictive information is the nervous system and a system whose only source of predictive information is the environmental stimuli. Understanding the operation of neural circuits in agent-environment

¹Cognitive Science program, Indiana University, Bloomington, IN, USA. ²The Luddy School of Informatics, Computing, and Engineering, Indiana University, Bloomington, IN, USA. ✉email: edizquie@iu.edu

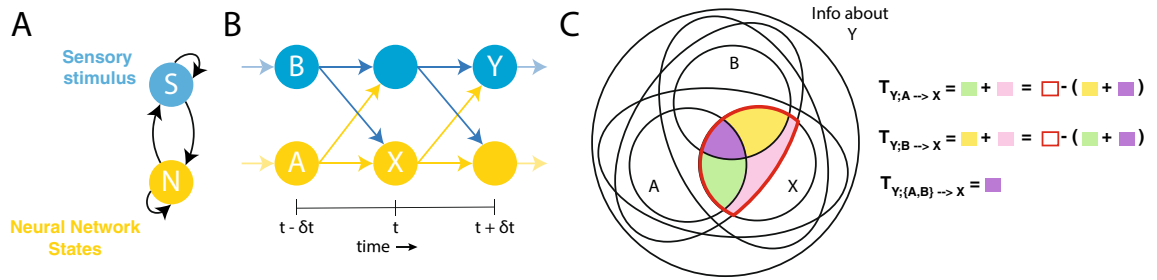


Figure 1. Predictive information source estimation based on idealized agent-environment interaction. (A) Sensory stimuli (S) and neural activity (N) are two coupled dynamical systems. (B) Agent-environment interaction unrolled over time. X represents current neural activity, $N(t)$, Y, future environmental 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.

systems involves understanding the information flow dynamics across its components^{13,14}. The principal contribution of this paper is the application of an information-theoretic framework, specifically Partial Information Decomposition¹⁵, to quantify the contributions from the nervous system and the contributions from the environmental stimuli to the total predictive information in an agent. First, we decompose the total predictive information in the neural system into information that was uniquely transferred from each source. In order to do this, we employ multivariate extensions to information theory¹⁵. Second, we unroll information over time to backtrack the origin of the source of predictive information and how they change over time. To validate the proposed theoretical framework, we examine it on a set of computational models of agent-environment systems, where the agent is driven by a dynamical recurrent neural network^{16,17}. The systems have been deliberately designed so that the source of predictive information is tractable and manipulable. We demonstrate how the proposed framework correctly reveals different sources of predictive information in systems with otherwise similar amounts of predictive information. Ultimately, we demonstrate how revealing the flow of information across the agent-environment system can help us to better understand the mechanisms underlying predictive coding.

Predictive information is studied in living organisms because it is considered a signature of their adaptive capacities^{5,8,9}. In the second part of this paper we study the relationship between a system’s ability to perform a task and its predictive information. In order to do this, we turn to a computational model of an agent that is required to process the received stimulus from the environment and make a decision based on it. Specifically, we study predictive information in the context of a relational categorization task^{18,19}. We generate model systems that are adapted to their environment and yet remain tractable to analysis by optimizing dynamical recurrent neural networks using an evolutionary algorithm to perform the task^{20,21}. We then proceed to analyze the resulting systems using predictive information and we compare the results against that of random systems that cannot solve the task. Counterintuitively, we observe that predictive information in trained neural networks is similar to predictive information in random neural networks. This suggests that predictive information alone is not sufficient to distinguish between living organisms that are adapted to their environments and non-adaptive systems. The rest of the paper focuses on an analysis of optimized and random systems using the framework proposed. Altogether, we demonstrate that decomposing predictive information across the components of an agent-environment system, and unrolling it over time reveals its true nature.

Identifying the source of predictive information

Predictive information is the information encoded in neural activity about its future stimulus. Formally, it is defined as mutual information between current neural activity (N_t) and the stimulus at a future time ($S_{t'}$)^{9,22-25}. Predictive information, $I(S_{t'}, N_t)$, is given by:

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

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

The presence of predictive information in a neural network suggests there is a source where this information was generated. In an idealized agent-environment system (Fig. 1A), the source of information can be either the neural activity in the previous time step, the environmental stimulus in the previous time step, or both (Fig. 1B). Measuring predictive information as defined in Eq. (1) requires that we examine two variables: current neural

activity (N_t , henceforth X) and future stimulus ($S_{t+\delta t}$, henceforth Y). Identifying the source of this predictive information requires that we examine two additional variables: past neural network activity ($N_{t-\delta t}$, henceforth A) and past stimulus ($S_{t-\delta t}$, henceforth B). Such an analysis requires that we adopt multivariate extensions to information theory. We focus specifically on Partial Information Decomposition (PID)¹⁵, a method for decomposing multivariate mutual information into combinations of unique, redundant and synergistic contributions, as well as measures of information gain, loss and transfer^{15,26–35}. PID allows us to decompose the total information about Y in the combined set of variables A , B and X into the constituent unique, redundant and synergistic information atoms. In this work we utilize these partial information atoms to identify the source of predictive information by measuring the following information transfer terms: (a) information uniquely transferred from past environmental stimulus, $T_{Y:A \rightarrow X}$; (b) information uniquely transferred from past neural network activity, $T_{Y:B \rightarrow X}$; and (c) information redundantly transferred from past environment stimulus and past neural network activity, $T_{Y:\{A,B\} \rightarrow X}$.

In words, a systematic development of how PID atoms can be utilized to measure these terms is as follows:

- Current neural activity, X , has information about future environmental state, $Y - I(Y; X)$
- Past neural activity, A , has information about future stimulus, $Y - I(Y; A)$
- Past environmental state, B , has information about future stimulus, $Y - I(Y; B)$
- Sources A and B have synergistic and redundant information about Y
- The variable of interest, current neural activity X , has information about Y that is redundant with information from the combined sources A and $B - \Pi_R(Y; \{[A, B]\}\{X\})$ (red outlined region in Fig. 1C)
- Of *that* information, some of it uniquely came from A , (green in Fig. 1C); some of it uniquely came from B , (yellow in Fig. 1C); some of it is redundant between them, (purple in Fig. 1C); and some synergistic (pink in Fig. 1C).
- Information transferred from a given source is then measured as information that was uniquely provided by the source plus the information that was synergistically transferred by the combined sources³⁵.
- As shown by the equations in Fig. 1C, the sum of the unique and synergistic information is equivalent to the difference between the total redundant information and the unique contribution from the other source.

Mathematically, the three transfer terms namely, transfer from A , transfer from B , and redundantly from both sources can be expressed as follows:

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

where $\Pi_R(Y; \{A_1\}\{A_2\}..\{A_k\})$ is the redundant information that random variables A_1 through A_k have about the random variable Y and $[A, B]$ refers to a random variable that is a concatenation of A and B . In words, information about Y transferred uniquely from source A to X is estimated as the total redundant information from the combined sources $[A, B]$ minus the information that is redundant with the other source B . This decomposition of the total information into different contributions is typically represented using a PI-decomposition diagram (Fig. 1C). Several approaches have been proposed to measure redundant information, Π_R ^{26,36,37}. Here, we use I_{min} because this is the only approach that can guarantee non-negative information decomposition in a system with four random variables, as is the case here.

During the course of behavior, the flow of information in a system changes over time^{38,39}. In order to understand the source of predictive information for any agent-environment system, it is not enough to decompose information from multiple sources; we must also track its flow of information over time. Although information theoretic measures are typically averaged over time, the measures described above can be unrolled over time^{13,39}. This is done by measuring information transfer at each time-point using data collected across several trials thereby allowing us to study the dynamics of predictive information sources.

Disparate systems with similar predictive information

Neural systems can be predictive in fundamentally different ways: they can generate predictive information internally or they can extract it from environmental stimulus. We use computational models of two extreme conditions where the ground-truth predictive information source is known to be the environment in one condition and the neural network in the other, to demonstrate that (a) predictive information cannot distinguish between these different kinds of systems and (b) it is only through decomposing the information across sources and unrolling over time that we can distinguish the two systems based on their operation. The two conditions we consider are agent-environment interactions at two extremes of the range of possible interactions: a central pattern generator (CPG) and a passive perceiver (PP). In the CPG condition, the neural network influences the environment by producing spontaneous oscillatory activity but receives no input from the environment (Fig. 2A). In the PP condition, the neural network is influenced by input from the environment, but it does not affect the environment (Fig. 2B). We evolved 100 different dynamical recurrent neural network CPGs, and in each case, we fed the sum of the neurons' outputs to the environment (Fig. S1A,B). For the PPs, we generated 100 random neural networks and fed them an oscillatory input. In order to provide the same distribution of activity as the CPG condition, we provided the random neural networks with the same oscillatory environmental signal that CPGs generated (Fig. S1C). The environmental signal and neural data were recorded from each instance for 500 trials where, in each trial the environment started with a different initial condition. Although, the environmental signal and the neural activity exhibit oscillatory activity in both conditions, the key difference in the operation

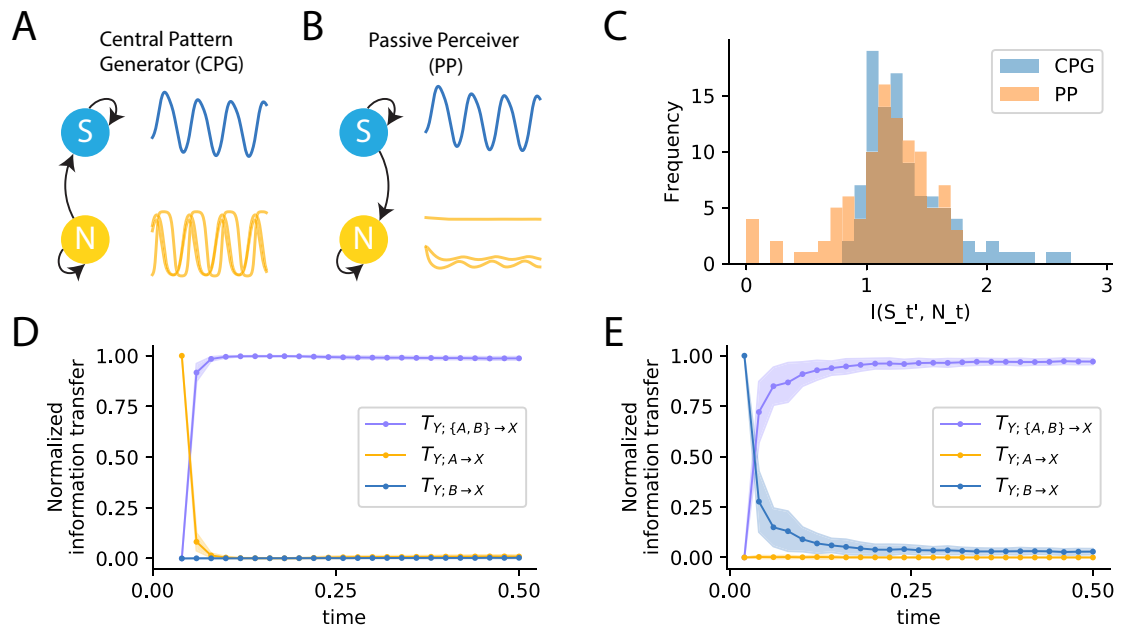


Figure 2. Predictive information in systems on the extremes of the range of possible agent-environment interactions (A) Schematic and traces of a Central Pattern Generator (CPG) that influences the environment through intrinsically generated oscillations. (B) Schematic and traces of a Passive Perceiver (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 Eq. (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 synchronized. (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).

of these systems is that in the CPGs the neural network drives its own activity and in the PP, the environment drives the neural network. Therefore, the neural network is the source of predictive information in the CPGs and the environment is the source of predictive information in the PPs.

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

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

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

In the PP condition, since the neural networks are driven by the environment (Fig. 2B), the only source of information about the future environmental signal is the stimulus from the environment itself. Accordingly, the dynamics of information transfer for PP systems reveals correctly that the environment is the source of predictive information (Fig. 2E). As opposed to the CPG systems, at the start of the interaction between the neural network and the environment, it is the environment that transfers unique information to the neural network.

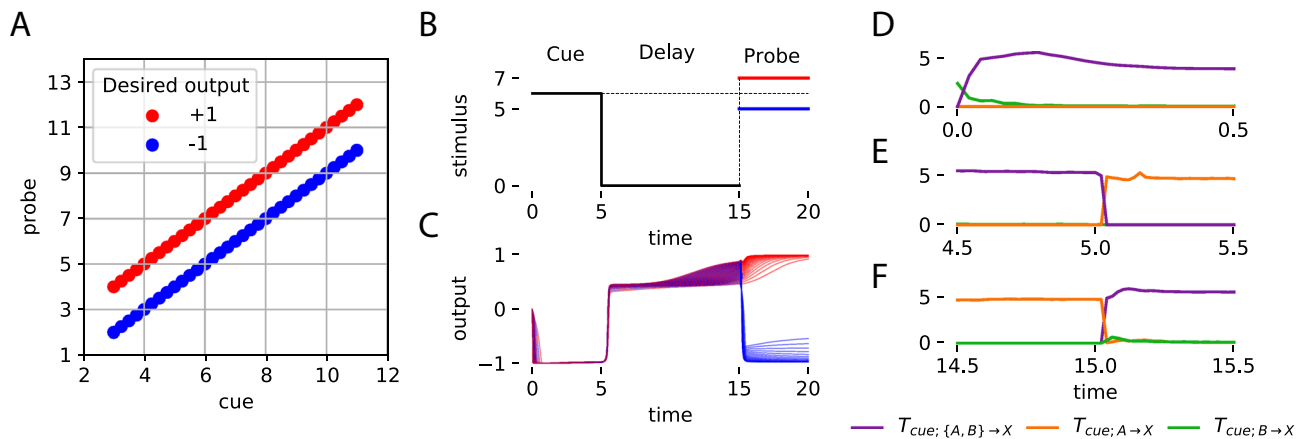


Figure 3. Predictive information source dynamics with structured stimuli. **(A)** Distribution of cue and corresponding probes in the relational categorization task. For each cue, the probe can be one of two values: greater, $cue + 1$, or lesser, $cue - 1$, with the expected outputs of $+1$ (red) and -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.

Subsequently, and similarly to the CPG condition, as the state of the neural network begins to encode the information from the environmental stimulus, the predictive information is redundantly 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 information alone cannot distinguish between two extremely different kinds of neural systems, both of which encode predictive information about the future of the environment. This is because when the entire time course of the data is considered, 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. However, in reality, both the environment and the neural network would have independent dynamics that would be continuously evolving. Thus, there will be transfer of information to or from the neural network (or both) at every time-step. We show a minimal version of such a setting where changes in the environment after synchronization results in new information transfer to the neural network, in the next section. 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^{2,40,41}. 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 objects based on the relative value of their attributes^{18,19}. 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 stimuli across three stages: cue, delay, and probe. In the cue stage, the neural network is provided with a stimulus of specific magnitude for a duration of time. This is followed by a delay stage, where no stimulus is provided. Finally, in the probe stage, the neural network is provided with a second stimulus. The magnitudes for the cue and probe stage stimuli are picked from a predesignated distribution (Fig. 3A). It is this distribution that defines the structure in the environment. For this study, we design it such that the stimulus in the probe stage can have a magnitude that is one of two values: smaller ($cue - 1$) or larger ($cue + 1$) than the stimulus provided during the cue stage (Fig. 3B). The goal of the neural network in this task to perform a relational categorization of “greater than” or “lesser than”

by producing an output of +1 or −1 respectively, during the probe phase. This task has been widely studied in a variety of contexts including in humans⁴², pigeons⁴³, rats⁴⁴, insects⁴⁵, as well as using computational models^{46,47}.

In this section, we show results from analysis of neural networks performing the relational categorization task. We demonstrate that decomposing information across the sources and unrolling over time reveals that the environment is structured by appropriately attributing the observed predictive information to either the environment or the dynamics of the neural network. Furthermore, we demonstrate that encoding predictive information alone is not indicative of task performance and that the magnitude and source of predictive information can change during the course of a behavior depending on environmental structure and neural network dynamics.

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

In order to better understand how a neural network performed this task, we can characterize the flow of information across the agent–environment system. To this end, we decomposed 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 environment 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 becomes 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.

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

Environmental regularities induces predictive information in any neural network. Since optimized neural networks encode information about the probe merely by encoding the cue, does any neural network that encodes the cue also encode information about the probe, and therefore have similar predictive information? In order to study this, we created 100 random neural networks and presented them with the same task. Although these neural networks were not able to perform the relational categorization task (Fig. S2B), they encoded similar amounts of total predictive information as the trained neural networks (Fig. 4A). Specifically, they encode the same amount of information about the probe during the cue stage (Fig. 4B). Furthermore, decomposing that information revealed that the information originated from the environmental stimulus and that the neural network dynamics had no role in its encoding of predictive information in both random and optimized neural networks (Fig. 4C). Thus, predictive information alone is not sufficient to distinguish neural networks optimized to perform specific tasks from random neural networks that are merely reflecting the information provided by the environment.

Information decomposition distinguishes between random and optimized neural networks. Unlike CPG and PP that were distinguished based on having different information sources, random and optimized neural networks in the relational categorization task have the same information sources. Even under this condition, decomposing the total information across sources and unrolling over time helps distinguish them by revealing differences in the magnitude of information transferred from each source over time. Specifically, predictive information sourced by the neural network during the delay stage is markedly different between random and optimized neural networks. As discussed in the previous section, optimized neural networks preserve information about the cue (and hence predictive information about the probe) during the delay stage. In contrast, random neural networks tend to lose that information. As a consequence, the amount of unique information provided by the neural network at the end of the delay period is higher for the trained neural networks than for the random neural networks (Fig. 5A). This difference disappears when information is measured across time, and can only be observed by unrolling it over time.

Statistics of the environment influences magnitude of predictive information. Encoding the cue results in encoding information about the probe in this task because of the relationship between them. How does changing this relationship impact predictive information in the neural networks? In order to study this, without changing the nature of the relational categorization task we merely changed the structure in the environ-

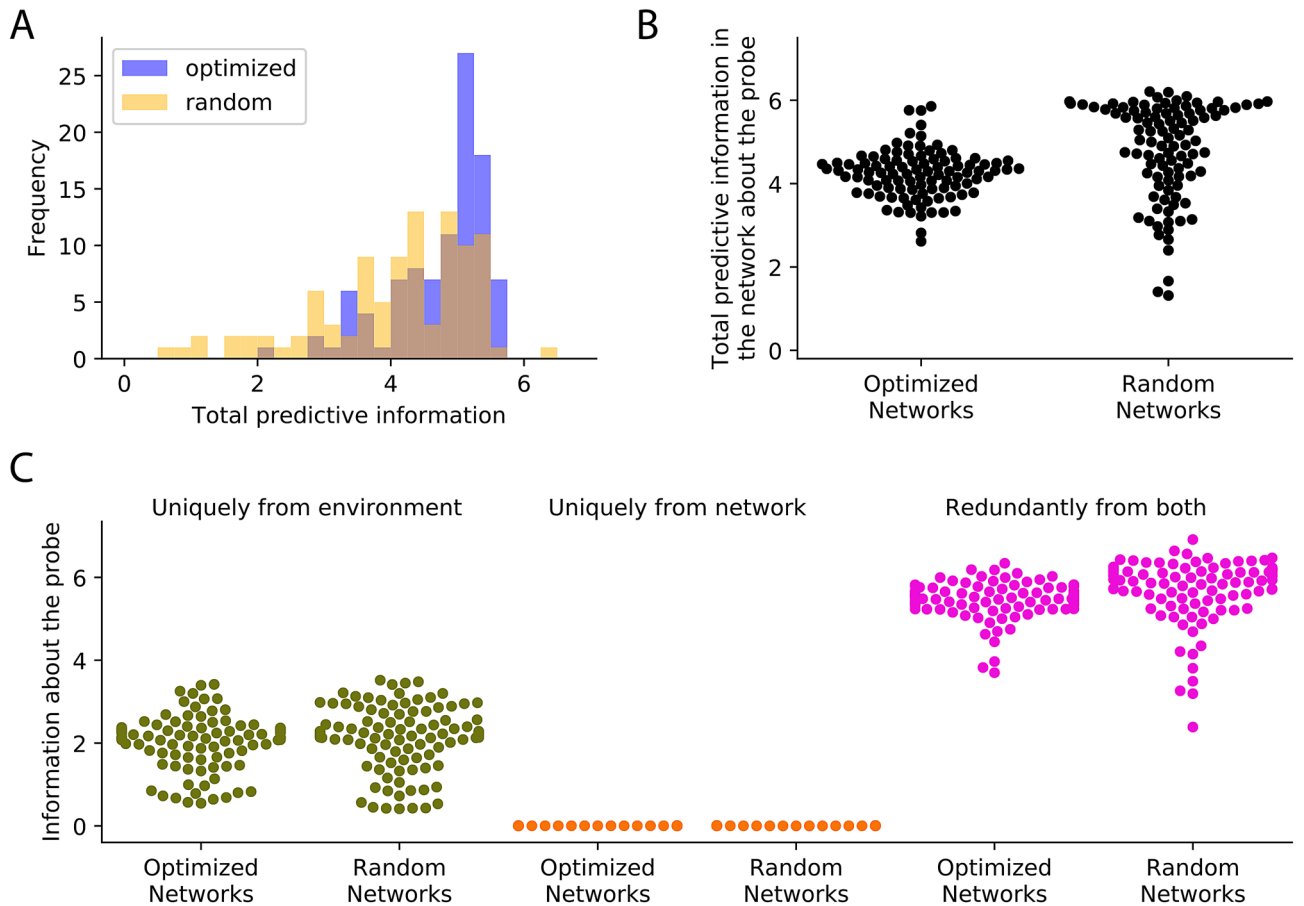


Figure 4. Comparison of predictive information sources in optimized and random neural networks. (A) Total predictive information estimated by averaging over the entire course of the task is similar in random and optimized neural networks. (B) Total predictive information about the probe averaged across 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).

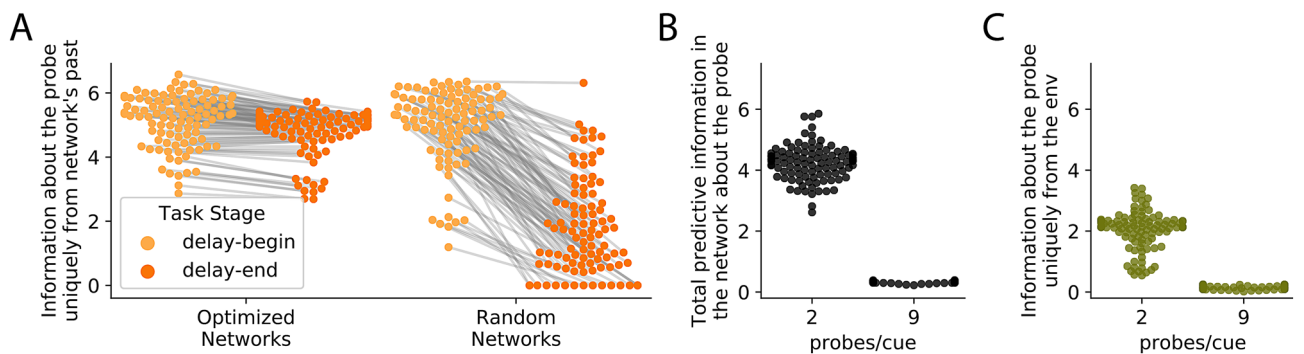


Figure 5. Influence of neural network and environmental properties on predictive Information (A) Both random and optimized neural networks have similar levels of information about the probe at the beginning of the delay stage (light orange), but unlike optimized neural networks, random neural networks 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.

ment. This was achieved by modifying the task such that the probe could be one of 9 possible values for a given cue, rather than one of two possible values (Fig. S4B). Reduction in uncertainty about the probe's value given the cue is now much less compared to the original environmental structure (Fig. S4D,E). This will be reflected in the information that the cue can provide about the probe. However, this came at no cost to performance because the neural networks were still encoding the cue just as well. The same ensemble of optimized neural networks were able to perform this task successfully without any more training (Fig. S4E). Information dynamics was then measured using data recorded under this 9-probe condition. Measuring the total information in the neural network during the cue stage about the probe revealed that there was significantly less information in the neural network in 9 probes per cue condition (Fig. 5B). The reduction in total predictive information can be wholly attributed to the reduction in information about the probe (Fig. 5C). Thus, differences in environmental structure can result in significantly different amounts of predictive information encoded in neural networks without any behavioral differences.

Discussion

The study of predictive coding and its relevance to behavior has been studied from multiple perspectives in the literature: with regards to the source of information, predictive information can be generated by the neural network^{5,6} and predictive information can be provided by the environment^{7,22}. In this work, using computational models where the ground-truth about the source of information was known, we demonstrate that predictive information can originate from either the environment or the neural network or both, and that the source of information can dynamically change during the course of a behavior. In order to do this, we first presented a theoretical framework based on multivariate information theory that allows us to infer the source of predictive information and its dynamics. This involved decomposing the total information that neural networks encode about a future stimulus into information transferred uniquely from the neural network, uniquely from the environment and redundantly from both sources. We validated this framework using the CPG and PP models where information is known to originate from the neural network and the environment respectively. Second, using the more structured relational categorization task, we demonstrated that (a) amount of predictive information encoded in a neural network is not indicative of its performance; (b) the source of information about a future stimulus can change during the course of the task; and (c) the source of information about a future stimulus can change within the same task depending on the regularities of the environment. Thus, predictive information might be necessary but is not sufficient to explain the neural basis of a behavior. Decomposing information across sources and studying its dynamics over time takes us one step further in understanding the role of predictive information in a behavior.

Although there are other methods of measuring the amount of information transferred from a source to a target, information transfer measured according to our approach captures crucial aspects relevant to predictive coding that others do not. Two widely-used alternative approaches are input-corrected Active Information Storage⁴⁸ and Transfer Entropy⁴⁹. Input-corrected Active Information Storage measures the information from the past of one source that is actively in use by the target for the current time-step. In order to apply this framework to identifying the source of predictive information, one can measure mutual information about the target variable, Y , in the current neural activity, X , conditioned on one source, say A , namely $I(Y; X|A)$. This assumes that any information in X that does not come from A must have come from B . Alternatively, one can find $I(A; Y|B)$, which measures the information in A about Y conditioned on the other source B . This measure, however, does not capture the information in current neural activity, X which is the information quantity we are interested in. In other words, this latter conditional mutual information measure assumes perfect transfer of all information from A to X . The difference between this approach and ours is that essentially these approaches consider the three variable decomposition whereas we consider the four variable decomposition. While this approach can be expanded to include other variables, it requires that we know all the other factors, our approach explicitly by only decomposing the information in X , about Y , that came from the sources A, B . Another relevant measure is Transfer Entropy (TE) which allows us to measure the information transferred from the past environmental state, B , to the current neural network state, X . Similarly, from past neural network state, A , to the current neural network state, X . However, conventional TE measures information in each source about the target, and not information in the target about the variable of interest, Y . There exists multivariate extensions to TE⁵⁰ that allow us to measure polyadic interactions between multiple sources. The equivalence between dyadic TE and three variable PID has been shown in Williams and Beer, 2011³⁵. For the four variable case as shown in this work, multivariate extensions to TE potentially can be shown to be equivalent, an interesting research direction that could unite the different methodologies. Further, TE is also often estimated at several delays beyond just the previous time step^{51,52}. Similarly, in the future, our proposed approach could be evaluated at different delays for conditions when there is a delay in the interaction between environment and neural network. In this work, we have already utilized different delays by evaluating information in the current neural activity about an environmental state further in the future rather than the next time step.

The framework presented here for inferring the source of predictive information takes us beyond general correlations that information theoretic measures are known to measure by capturing the effects of perturbation on the neural system. Identifying the sources of predictive information requires that the system under study be perturbed. The presentation, removal or sudden change of a stimulus is a perturbation. This causes the system to break the redundant encoding observed in a steady-state. It is during such a perturbation that we can use partial information decomposition to determine the source of information in a coupled system. Note that, like any other information theoretic measure, the use of PID alone does not automatically guarantee inference of causal effect^{31,34,53,54}. It is through the combination of information decomposition, time-unrolling and perturbation that we are able to infer the directed causal influences in the models we have analyzed. However, in the strictest

sense, inferring causal relationships require that the system is studied beyond perturbations involving changes in the stimuli, but through selective manipulation of one variable at a time⁵⁵. While, our information theoretic analyses and causality coincide in the models we have analyzed, it cannot be assumed in the general case.

The framework presented here can be applied to experimental data across multiple scales. In fact, it can be applied to any time-series data spanning multiple trials corresponding to several perturbations from the steady state. However, in this work, we focus on open-loop systems. Specifically, we focus on agent-environment systems where the agent influences its environment or where the agent is influenced by the environment. Such an open-loop setup is typical in experiments in neuroscience, where the subject receives a stimulus, but does not have the ability to influence the future stimulus through their state or actions. In natural behavior, the agent and environment are in closed-loop interaction. The analysis of closed-loop systems introduces an added complexity. The regularities of the environment can be generated by the regularities of the neural network's dynamics and vice-versa. As a result, the distribution of environmental stimuli and the distribution of the neural activity are dependent on each other, unlike the open-loop setup where one of them is independent of the other. As it is, the framework requires that one of the distributions be fixed across time in order to make fair comparisons of information at different time-points. Future work in this direction will involve extending the framework and designing the experimental setting that would allow us to infer the source of predictive information in a freely moving animal.

Methods

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

Neural network model. A Continuous-Time Recurrent Neural Network (CTRNN) was used as the model neural network^{16,17}. The neural network consisted of three layers: the input layer which was connected by a set of feed-forward weights to the interneuron layer; the interneuron layer was a CTRNN which fed into the output layer; the output layer produced the output of the neural network which was given by a weighted combination of the interneurons' output. The dynamics of each interneuron was governed based on state equations given by

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

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

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

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

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

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

$$\tau \frac{ds}{dt} = -s + O \quad (6)$$

where s refers to the state of the environment, τ refers to its time-constant which was set to 0.5, and O refers to the output of the neural network given by Eq. (5).

Relational categorization task. We adapted the relational categorization task to provide neural networks with structured stimuli^{18,19,46}. This task involves first providing the neural network with a cue stimulus in the range $[3, 11]$ for 5 units of time. This is followed by a delay period when no stimulus is provided for 10 units of

time. Finally, a probe stimulus that is of magnitude greater or less than the cue is provided for 5 units of time. The goal of the task is for the neural network to distinguish probes that were larger than the cue or smaller than the cue, by producing an output of +1 or -1 respectively. In the first version of this task, the probe can take one of only two values, either $cue + 1$ or $cue - 1$. In the second version of the task, the probe can take any value in [3, 11]. While the goal of the task remains the same in both versions, the distribution of the probes given the cue, and therefore information that the cue gives about the probe is significantly different (Fig. S4). Performance of a neural network in this task was estimated by measuring absolute deviation of the network's output from the desired output of +1 or -1 during the probe stage. Time-averaged deviation was also averaged across all trials of cue-probe values, to obtain a score in the range [0, 1].

Neural network optimization. Neural network models described previously were optimized to perform the relational categorization task using an evolutionary algorithm^{56,57}. This optimization methodology involves instantiating a population of 100 random solutions that evolves over several generations to produce solutions capable of performing the task. A generation is defined as the process of creating a new population of solutions that has improved in “fitness” (task performance) from the last. Each solution, referred to as a genotype, is an N dimensional vector corresponding to the parameters to be optimized. The parameters were encoded to be in the range [0, 1] and scaled to produce the neural network that the genotype encoded. In each generation, the fitness of every genotype is evaluated and a new population is created using a fitness-based selection and mutation strategy as follows: The genotypes that perform in the top 1% were retained as is for the next generation. The rest of the individuals were created by selecting two genotypes preferentially in proportion to their fitness and combining them. To these offspring, Gaussian mutation noise with mean 0 and standard deviation 0.01 was added before being added to the population of genotypes for the next generation. After a fixed number of generations, the best individual in the population was selected as the representative solution from that optimization run. 100 such runs were conducted to obtain an ensemble of 100 neural network models that successfully performed each task. For the relational categorization task, optimization was carried out for 500 generations. In the case of the CPG task, at the end of 50 generations the optimization process was terminated and deemed successful if the best agent in the population reached a fitness of 30 or greater. This was repeated until 100 CPGs were produced. See supporting information (Figs. S1 and S2) for training curves, behavior of best optimized neural network, distribution of fitness of best models from 100 runs, and sample neural traces.

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

Measuring information transfer. To identify the source of information over time, information transfer measures were estimated independently at each time point. For any given time step, data for environmental stimulus at the previous time step, neural activity of previous time step, current neural activity, and stimulus at a future time step, was collected across multiple trials. Probability densities were estimated from this data using a kernel density estimation technique known as average shifted-histograms⁵⁸ with 7 shifted binnings of 100 bins along each dimension of the data space. These probability density estimates were then used to measure the redundant information terms in Eq. (2). Similar results were observed with 5 and 11 shifts and with 50 and 200 bins per dimension (Fig. S6). All information theoretic quantities were estimated from raw data using the *infotheory* package⁵⁹.

Received: 8 March 2020; Accepted: 7 September 2020
Published online: 09 October 2020

References

1. Clark, A. Whatever next? predictive brains, situated agents, and the future of cognitive science. *Behav. Brain Sci.* **36**, 181–204 (2013).
2. Huang, Y. & Rao, R. P. Predictive coding. *Wiley Interdiscip. Rev. Cognit. Sci.* **2**, 580–593 (2011).
3. Rao, R. P. & Ballard, D. H. Predictive coding in the visual cortex: a functional interpretation of some extra-classical receptive-field effects. *Nat. Neurosci.* **2**, 79 (1999).
4. Srinivasan, M. V., Laughlin, S. B. & Dubs, A. Predictive coding: a fresh view of inhibition in the retina. *Proc. R. Soc. Lond. Ser. B. Biol. Sci.* **216**, 427–459 (1982).
5. Friston, K. & Kiebel, S. Predictive coding under the free-energy principle. *Philos. Trans. R. Soc. B Biol. Sci.* **364**, 1211–1221 (2009).
6. Friston, K. The free-energy principle: a unified brain theory? *Nat. Rev. Neurosci.* **11**, 127 (2010).
7. Bialek, W., Nemenman, I. & Tishby, N. Predictability, complexity, and learning. *Neural Comput.* **13**, 2409–2463 (2001).
8. Still, S. Information-theoretic approach to interactive learning. *Europhys. Lett.* **85**, 28005 (2009).
9. Palmer, S. E., Marre, O., Berry, M. J. & Bialek, W. Predictive information in a sensory population. *Proc. Nat. Acad. Sci.* **112**, 6908–6913 (2015).
10. Chen, K. S., Chen, C.-C. & Chan, C. Characterization of predictive behavior of a retina by mutual information. *Frontiers Comput. Neurosci.* **11**, 66 (2017).
11. Chao, Z. C., Takaura, K., Wang, L., Fujii, N. & Dehaene, S. Large-scale cortical networks for hierarchical prediction and prediction error in the primate brain. *Neuron* **100**, 1252–1266 (2018).

12. Sederberg, A. J., MacLean, J. N. & Palmer, S. E. Learning to make external sensory stimulus predictions using internal correlations in populations of neurons. *Proc. Nat. Acad. Sci.* **115**, 1105–1110 (2018).
13. Williams, P. L. Information dynamics: Its theory and application to embodied cognitive systems. Ph.D. thesis, PhD thesis, Indiana University (2011).
14. Beer, R. D. & Williams, P. L. Information processing and dynamics in minimally cognitive agents. *Cogn. Sci.* **39**(1), 1–38 (2015).
15. Williams, P. L. & Beer, R. D. Nonnegative decomposition of multivariate information. arXiv preprint <https://arxiv.org/1004.2515> (2010).
16. Funahashi, K.-i. & Nakamura, Y. Approximation of dynamical systems by continuous time recurrent neural networks. *Neural Netw.* **6**, 801–806 (1993).
17. Beer, R. D. On the dynamics of small continuous-time recurrent neural networks. *Adapt. Behav.* **3**, 469–509 (1995).
18. Gentner, D. & Kurtz, K. J. Relational categories. In *APA decade of behavior series. Categorization inside and outside the laboratory: Essays in honor of Douglas L. Medin* (eds Ahn, W.-k., Goldstone, R. L., Love, B. C., Markman, A. B. & Wolff, P.) 151–175 (American Psychological Association, 2005). <https://doi.org/10.1037/11156-009>.
19. Markman, A. B. & Stilwell, C. H. Role-governed categories. *J. Exp. Theor. Artif. Intell.* **13**, 329–358 (2001).
20. Beer, R. D. & Gallagher, J. C. Evolving dynamical neural networks for adaptive behavior. *Adapt. Behav.* **1**, 91–122 (1992).
21. Floreano, D., Dürr, P. & Mattiussi, C. Neuroevolution: from architectures to learning. *Evol. Intel.* **1**, 47–62 (2008).
22. Bialek, W. *Biophysics: Searching for Principles* (Princeton University Press, Princeton, 2012).
23. Rieke, F., Warland, D., Van Steveninck, R. d. R., Bialek, W. S. *et al. Spikes: Exploring the Neural Code* Vol. 7 (MIT Press, Cambridge, 1999).
24. Cover, T. M. & Thomas, J. A. *Elements of Information Theory* (Wiley, Hoboken, 2012).
25. Shannon, C. E. A mathematical theory of communication. *Bell Syst. Tech. J.* **27**, 379–423 (1948).
26. Bertschinger, N., Rauh, J., Olbrich, E., Jost, J. & Ay, N. Quantifying unique information. *Entropy* **16**, 2161–2183 (2014).
27. Faber, S. P., Timme, N. M., Beggs, J. M. & Newman, E. L. Computation is concentrated in rich clubs of local cortical networks. *Netw. Neurosci.* **3**, 384–404 (2019).
28. James, R. G., Barnett, N. & Crutchfield, J. P. Information flows? A critique of transfer entropies. *Phys. Rev. Lett.* **116**, 238701 (2016).
29. James, R. G., Ellison, C. J. & Crutchfield, J. P. Anatomy of a bit: Information in a time series observation. *Chaos Interdiscip. J. Nonlinear Sci.* **21**, 037109 (2011).
30. James, R. & Crutchfield, J. Multivariate dependence beyond Shannon information. *Entropy* **19**, 531 (2017).
31. Lizier, J. T. & Prokopenko, M. Differentiating information transfer and causal effect. *Eur. Phys. J. B* **73**, 605–615 (2010).
32. Lizier, J. T., Bertschinger, N., Jost, J. & Wibral, M. Information decomposition of target effects from multi-source interactions: Perspectives on previous, current and future work. *Entropy* **20**, 307 (2018).
33. Timme, N., Alford, W., Flecker, B. & Beggs, J. M. Synergy, redundancy, and multivariate information measures: an experimentalist's perspective. *J. Comput. Neurosci.* **36**, 119–140 (2014).
34. Wibral, M., Priesemann, V., Kay, J. W., Lizier, J. T. & Phillips, W. A. Partial information decomposition as a unified approach to the specification of neural goal functions. *Brain Cognit.* **112**, 25–38 (2017).
35. Williams, P. L. & Beer, R. D. Generalized measures of information transfer. arXiv preprint <https://arxiv.org/1102.1507> (2011).
36. James, R. G., Emenheiser, J. & Crutchfield, J. P. Unique information via dependency constraints. *J. Phys. A Math. Theor.* **52**, 014002 (2018).
37. Griffith, V. & Koch, C. Quantifying synergistic mutual information. In *Guided Self-Organization: Inception*, 159–190 (Springer, Berlin, 2014).
38. Izquierdo, E. J., Williams, P. L. & Beer, R. D. Information flow through a model of the c. elegans klinotaxis circuit. *PLoS ONE* **10**, e0140397 (2015).
39. Williams, P. L. & Beer, R. D. Information dynamics of evolved agents. In: *International Conference on Simulation of Adaptive Behavior*, 38–49 (Springer, Berlin, 2010).
40. Barlow, H. The exploitation of regularities in the environment by the brain. *Behav. Brain Sci.* **24**, 602–607 (2001).
41. Graham, D. & Field, D. Statistical regularities of art images and natural scenes: Spectra, sparseness and nonlinearities. *Spat. Vis.* **21**, 149–164 (2007).
42. Kurtz, K. J. & Boukrina, O. Learning relational categories by comparison of paired examples. In: *Proceedings of the Annual Meeting of the Cognitive Science Society* **26**, (2004).
43. Wills, S. Relational learning in pigeons? *Quart. J. Exp. Psychol. Sect. B* **52**, 31–52 (1999).
44. Saldanha, E. L. & Bitterman, M. E. Relational learning in the rat. *Am. J. Psychol.* **64**, 37–53 (1951).
45. Giurfa, M., Zhang, S., Jenett, A., Menzel, R. & Srinivasan, M. V. The concepts of sameness and difference in an insect. *Nature* **410**, 930 (2001).
46. Williams, P. L., Beer, R. D. & Gasser, M. An embodied dynamical approach to relational categorization. In: *Proceedings of the Annual Meeting of the Cognitive Science Society* **30**, (2008).
47. Izquierdo-Torres, E. & Harvey, I. Learning to discriminate between multiple possible environments: an imprinting scenario. In: *Memory and Learning Mechanisms in Autonomous Robots Workshop (ECAL 2005)* (2005).
48. Obst, O., Boedecker, J., Schmidt, B. & Asada, M. On active information storage in input-driven systems. arXiv preprint <https://arxiv.org/1303.5526> (2013).
49. Vicente, R., Wibral, M., Lindner, M. & Pipa, G. Transfer entropy—a model-free measure of effective connectivity for the neurosciences. *J. Comput. Neurosci.* **30**, 45–67 (2011).
50. Novelli, L., Wollstadt, P., Mediano, P., Wibral, M. & Lizier, J. T. Large-scale directed network inference with multivariate transfer entropy and hierarchical statistical testing. *Netw. Neurosci.* **3**, 827–847 (2019).
51. Wibral, M. *et al.* Measuring information-transfer delays. *PLoS ONE* **8**, e55809 (2013).
52. Ito, S. *et al.* Extending transfer entropy improves identification of effective connectivity in a spiking cortical network model. *PLoS One* **6**, e27431 (2011).
53. Chicharro, D. & Ledberg, A. When two become one: the limits of causality analysis of brain dynamics. *PLoS One* **7**, e32466 (2012).
54. Mehler, D. M. A. & Kording, K. P. The lure of causal statements: Rampant mis-inference of causality in estimated connectivity. arXiv preprint <https://arxiv.org/1812.03363> (2018).
55. Pearl, J. An introduction to causal inference. *Int. J. Biostat.* **6**(2), 7. <https://doi.org/10.2202/1557-4679.1203> (2010).
56. Mitchell, M. *An Introduction to Genetic Algorithms* (MIT press, Cambridge, 1998).
57. Goldberg, D. E. & Holland, J. H. Genetic algorithms and machine learning. *Mach. Learn.* **3**, 95–99 (1988).
58. Scott, D. W. Averaged shifted histograms: effective nonparametric density estimators in several dimensions. *Ann. Stat.* **33**, 1024–1040 (1985).
59. Candadai, M. & Izquierdo, E. J. infotheory: A c++/python package for multivariate information theoretic analysis. arXiv preprint <https://arxiv.org/1907.02339> (2019).

Acknowledgements

The authors would like to thank John M. Beggs, Randall D. Beer and Ehren L. Newman for helpful discussions and feedback over the course of this work. The authors would like to acknowledge funding received from NSF as

part of grant NSF-1845322. M.C. was partially funded by NSF-NRT grant 1735095 “Interdisciplinary Training in Complex Networks and Systems”. This research was supported in part by Lilly Endowment, Inc., through its support for the Indiana University Pervasive Technology Institute.

Author contributions

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

Competing interests

The authors declare no competing interests.

Additional information

Supplementary information is available for this paper at <https://doi.org/10.1038/s41598-020-73380-x>.

Correspondence and requests for materials should be addressed to E.J.I.

Reprints and permissions information is available at www.nature.com/reprints.

Publisher’s note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article’s Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article’s Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2020