

# Synergistic information supports modality integration and flexible learning in neural networks solving multiple tasks

ALEXANDRA M. PROCA,<sup>1</sup> FERNANDO E. ROSAS,<sup>2,3,4,5</sup> ANDREA I. LUPPI,<sup>6,7,8</sup> DANIEL BOR,<sup>9,10</sup>  
MATTHEW CROSBY,<sup>11,\*</sup> AND PEDRO A.M. MEDIANO<sup>9,11,\*</sup>

<sup>1</sup>*Department of Computer Science, ETH Zürich*

<sup>2</sup>*Department of Informatics, University of Sussex*

<sup>3</sup>*Centre for Psychedelic Research, Department of Brain Sciences, Imperial College London*

<sup>4</sup>*Centre for Complexity Science, Imperial College London*

<sup>5</sup>*Centre for Eudaimonia and Human Flourishing, University of Oxford*

<sup>6</sup>*Department of Clinical Neurosciences and Division of Anaesthesia, University of Cambridge*

<sup>7</sup>*Leverhulme Centre for the Future of Intelligence, University of Cambridge*

<sup>8</sup>*The Alan Turing Institute*

<sup>9</sup>*Department of Psychology, University of Cambridge*

<sup>10</sup>*Department of Psychology, Queen Mary University of London*

<sup>11</sup>*Department of Computing, Imperial College London*

## ABSTRACT

Striking progress has recently been made in understanding human cognition by analyzing how its neuronal underpinnings are engaged in different modes of information processing. Specifically, neural information can be decomposed into *synergistic*, *redundant*, and *unique* features, with synergistic components being particularly aligned with complex cognition. However, two fundamental questions remain unanswered: (a) precisely how and why a cognitive system can become highly synergistic; and (b) how these informational states map onto artificial neural networks in various learning modes. To address these questions, here we employ an information-decomposition framework to investigate the information processing strategies adopted by simple artificial neural networks performing a variety of cognitive tasks in both supervised and reinforcement learning settings. Our results show that synergy increases as neural networks learn multiple diverse tasks. Furthermore, performance in tasks requiring integration of multiple information sources critically relies on synergistic neurons. Finally, randomly turning off neurons during training through dropout increases network redundancy, corresponding to an increase in robustness. Overall, our results suggest that while redundant information is required for robustness to perturbations in the learning process, synergistic information is used to combine information from multiple modalities — and more generally for flexible and efficient learning. These findings open the door to new ways of investigating how and why learning systems employ specific information-processing strategies, and support the principle that the capacity for general-purpose learning critically relies in the system’s information dynamics.

Corresponding author: Alexandra M. Proca  
[alexandra.proca.20@ucl.ac.uk](mailto:alexandra.proca.20@ucl.ac.uk)

\* M.C. and P.M. are joint senior authors.

## 1. INTRODUCTION

### 1.1. *Leveraging information theory and deep learning for neuroscience*

A central goal in cognitive neuroscience is to understand how the brain processes information to learn and behave intelligently; and a central goal in machine learning research is to recreate these processes on a computer. Historically, this close partnership between cognitive neuroscience and machine learning has been a fruitful symbiosis (Marblestone et al. 2016). Although artificial neural networks are not perfect models of biological neurons (Izhikevich 2007), they are a valuable tool to investigate how groups of neurons collectively represent and manipulate information (Kriegeskorte 2015). Overall, the aim of this interdisciplinary research effort is not so much to clarify the implementation details of a particular instantiation of successful distributed information-processing (e.g., the human brain), but to extract fundamental principles to allow a better design of a wide range of novel cognitive systems (Eliasmith and Anderson 2004).

Information theory provides an ideal conceptual framework for the study of distributed information processing, motivated by the goal of understanding how groups of neurons store, transfer, and modify information (Lizier 2012). One particularly relevant tool for the analysis of such processes is the recent framework of *Partial Information Decomposition*, or PID (Williams and Beer 2010), that distinguishes the information held by a set of sources about a target variable into qualitatively different components: unique (present in exactly one source), redundant (provided by multiple sources separately), and synergistic information (only available when considering multiple sources jointly). Recent work has revealed a strong relationship between human high-level cognition and synergistic information processing taking place in the so-called ‘central executive network,’ which involves the lateral prefrontal and parietal cortices, while redundant information has been found to dominate in cortical areas responsible for perception and low-level processing (Luppi et al. 2022). Synergy and related quantities, such as integrated information (Mediano et al. 2021), have also been found to dominate in complex information processing taking place within cellular automata (Rosas et al. 2018; Mediano et al. 2022), and its disruption has been associated with loss of consciousness (Luppi et al. 2020) and ageing (Gatica et al. 2021). However, despite these promising findings, the precise nature of the underlying mechanisms resulting in the emergence of synergistic information and its utility for computation remains unknown.

Considering that a crucial feature of human cognition is the ability to learn flexibly and generalize across many different — potentially novel — settings, here we hypothesize that synergistic information may be important for such general-purpose learning. Machine learning provides well-suited avenues for testing this hypothesis by investigating the computations associated with, and the possible utility of, synergy. Although current AI models have yet to reach the level of generality of humans, they have recently shown good multitask performance. For example, “Agent 57” can outperform humans across all 57 Atari games (Badia et al. 2020), and Gato, a multi-modal, multi-task, generalist policy, uses a single network to answer language questions, caption images, play Atari and 3D exploration games, and control a real robot arm (Reed et al. 2022). Although little is known about the information-theoretic properties of such networks, here we propose that a closer investigation will provide a better understanding for the role of synergy in general-purpose learning systems.

To develop these ideas, in this paper we employ simple artificial neural networks with several different architectures in both supervised and reinforcement learning settings as a testbed to investigate general information-processing principles related to learning. The main contributions of this work are (i) to propose functional roles played by information decomposition components in learning scenarios, and (ii) to establish a computational basis for the existing evidence of synergy’s importance for complex cognition, with a specific relation to general-purpose learning by supporting multi-modal integration.

### 1.2. Background: Partial Information Decomposition

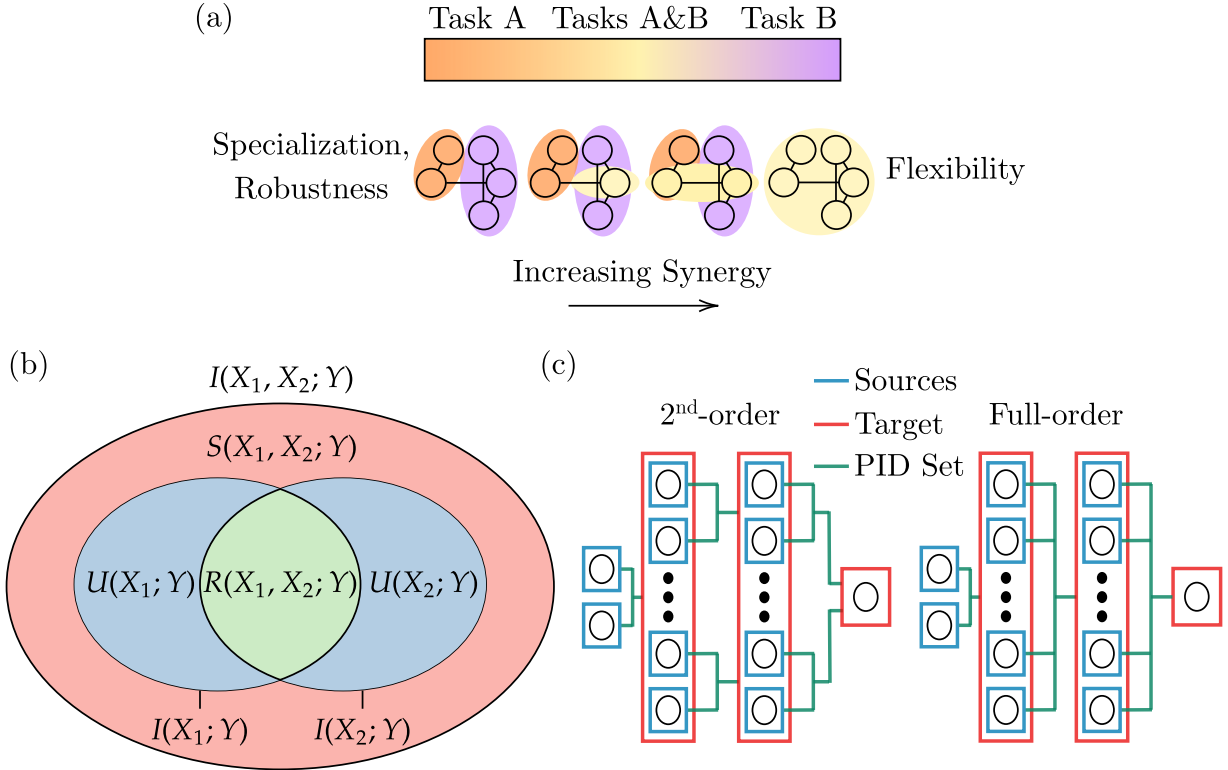
One of the central measures of information theory is Shannon’s mutual information  $I(X; Y)$ , which quantifies the amount of information a random variable  $X$  provides about another variable  $Y$  by measuring the extent to which knowing  $X$  reduces the uncertainty about the outcome of  $Y$ .

Extending beyond the bivariate case, for a set of random variables (*sources*)  $\mathbf{X} = (X_1, \dots, X_n)$  and another random variable (*target*)  $Y$ , the mutual information  $I(\mathbf{X}; Y)$  can be separated into distinct terms that describe the partial information contributed by subsets of sources about the target (Fig. 1b). As described earlier, these PID terms (or *atoms*) can either be unique ( $U$ ), redundant ( $R$ ), or synergistic ( $S$ ) and can be computed as described in Sec. 5.2.1. For the case with two sources, mutual information can be decomposed as

$$\begin{aligned} I(X_1; Y) &= R(X_1, X_2; Y) + U(X_1; Y) , \\ I(X_2; Y) &= R(X_1, X_2; Y) + U(X_2; Y) , \\ I(X_1, X_2; Y) &= R(X_1, X_2; Y) + U(X_1; Y) + U(X_2; Y) + S(X_1, X_2; Y) . \end{aligned}$$

Consider now a neural network learning two distinct tasks. The network can represent these tasks in different ways. It can assign a particular set of neurons to one task, and a separate set to the other (orange and purple neurons in Fig. 1a). Alternatively, it could use the same overlapping set of neurons to encode information about both tasks, distinguishing tasks by the collective behavior and interactions of such neurons (yellow neurons in Fig. 1a). The first method specializes its neurons by designating them for particular tasks, whereas the second method reuses its neurons for multiple tasks. In particular, the first method uses unique (and potentially redundant) information to solve the two tasks, and the second method uses synergistic information. Intuitively, the first approach conveys specialization and robustness, while the second approach provides potentially greater flexibility and reusability, as also suggested in prior work (Yang et al. 2019a).

Mutual information can be decomposed differently at different scales and is dependent on the selection of sources and targets. For example, each task-specific group of neurons can also vary in its decomposition, such that the group solving Task B could be more synergistic than the group used for Task A — even though the neuronal populations across tasks do not overlap. These types of scenarios can be studied by considering different sets of sources in a network over which to compute PID (Fig. 1c): one could use all neurons in a layer as a single set of sources and the joint state of the next layer as the target (full-order), or select all the combinations of pairs of neurons in a layer as sets of sources and average the resulting information atoms (2<sup>nd</sup>-order), or anything in

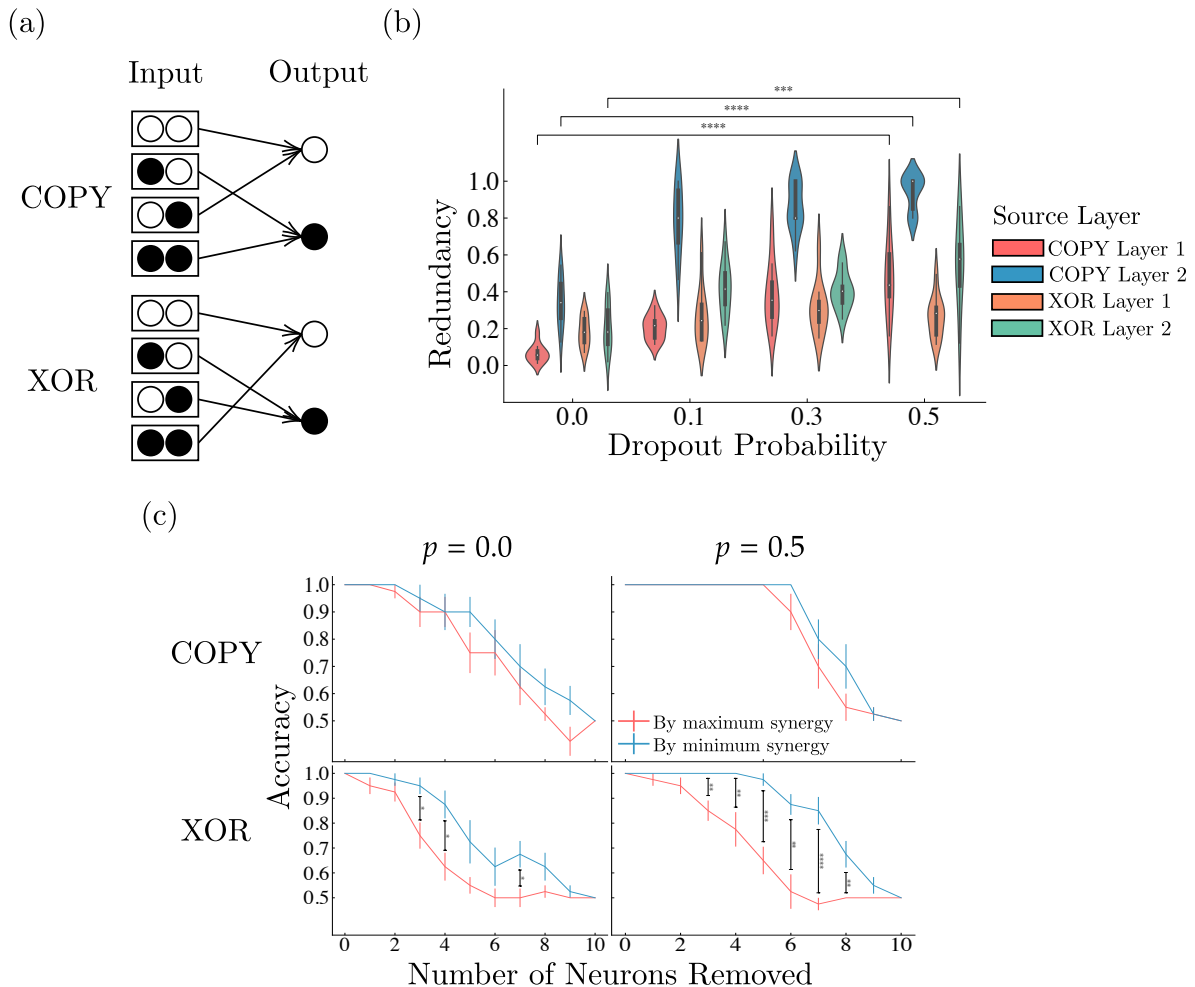


**Figure 1. Information decomposition in neural networks.** (a) A network solving two different tasks could potentially represent information about the two tasks in several ways: it can use distinct populations of neurons for each task (orange, purple) or have some combination of overlapping neuron populations of neurons (yellow) used for both tasks. (b) The decomposition of mutual information between two sources ( $X_1, X_2$ ) and a target ( $Y$ ). (c) Example of the set of sources and target considered at different orders in a neural network. The neurons in one layer comprise the set of sources for the target in the subsequent layer. In this setting, the entire layer of neurons is considered collectively as the target. Full-order PID refers to treatment of the entire set of neurons as the set of sources, whereas 2<sup>nd</sup>-order PID treats pairs of neurons as the set of sources of which all possible combinations are computed and averaged.

between. Studying different scales can contextualize information-processing behavior in terms of the interactions occurring between different sets and subsets of system parts. With the conceptual framework of PID and the numerical estimators presented in Sec. 5.2.1, one can properly investigate the information decomposition of neural networks, to which we now turn our attention.

## 2. RESULTS

We now present a series of experiments exploring the properties of synergy and redundancy in small neural networks with two hidden layers of ten neurons each. First, we look at information flow in networks solving logic gates, where the types of information needed to solve the task are known. We then study a more complex setting in which logic gates are embedded and extended in a 3D simulated environment and solved using reinforcement learning, either individually or in a set of tasks. Finally, we explore networks learning multiple tasks further using the NeuroGym suite of tasks inspired by cognitive neuroscience experiments (Molano-Mazón et al. 2022). Overall, our experiments converge



**Figure 2. Effects of lesions and dropout on network information profiles.** (a) White nodes represent 0 and black nodes represent 1, specifying the data used for each logic gate. (b) Dropout increases redundancy across the hidden layers of the network for both COPY and XOR tasks. By forcing the network to decrease its reliance on individual (sets of) neurons by randomly turning them off during training, dropout encourages redundancy to overrepresent important information. (c) Lesioning by permanently removing neurons during evaluation shows the causal role of synergy-rich neurons for task-performance, especially for tasks requiring integration of information (XOR). By increasing redundancy, dropout results in decreased reliance on individual neurons, allowing the network to be more robust to their removal. For the XOR gate, which requires the integration of information to solve, even after applying dropout performance quickly degrades when synergistic neurons are affected; this loss in performance is substantially attenuated if non-synergistic neurons are removed.

in supporting a relation of synergy with multi-modal integration and the learning of multiple tasks, and redundancy with robustness.

### 2.1. Functional roles of information atoms in simple learning problems

We first study redundancy and synergy in small feedforward networks learning a copy (COPY) or exclusive-or (XOR) logic gate involving two inputs (Fig. 2a), which are well-defined tasks with known informational requirements and few confounding factors. In particular, the processing done

by the COPY gate involves no synergy, as it requires no integration of information as it is solved by simply copying the first dimension of the input. Conversely, an XOR gate is solved by integrating information from both dimensions of its input as it reflects the parity of the set (i.e., if the inputs are similar or different). In fact, the XOR gate is known to be maximally synergistic (Rosas et al. 2019), as there is no reduction of uncertainty about the output unless all input sources are considered jointly. We perform analyses at both the 2<sup>nd</sup>-order and full-order scales to compare how information profiles vary depending on the number of sources considered; a discussion on our choice of order can be found in Sec. 5.2.3. We refer to 2<sup>nd</sup>-order measures in text and figures unless otherwise specified.

### 2.1.1. Dropout removes irrelevant input information and increases hidden layer redundancy

To study how learning pressures can influence the information profile of a network, for each logic gate we apply four levels of dropout ( $p = 0.0, 0.1, 0.3, 0.5$ ), a popular regularization method in deep learning that randomly omits different neurons (with some given probability) in each forward pass during training. This differs from “lesioning” in that it is applied during training, and a neuron is not permanently rendered inactive. This forces the network to be robust and be able to adapt to random perturbations applied to its neurons, rather than strongly dependent on particular sets of neurons. Furthermore, by applying dropout, we disrupt synergistic and unique information because they rely on individual neurons that may be randomly “turned off,” while redundant information persists in other neurons that remain “on” — making it an interesting paradigm for investigating the resulting learned information profiles.

By applying increasing amounts of dropout, we find that redundancy and synergy from the input significantly decrease for networks solving COPY gates, but are preserved for XOR gates (Supp. Fig. 8). Whereas the XOR gate requires both dimensions of the input to successfully complete the task, the COPY gate relies on only the first dimension. Correspondingly, the observed reduction of redundant and synergistic information from the COPY input is due to the loss of the second dimension, which is ignored. Thus, dropout encourages the removal of task-irrelevant information from the input. This may be due to the increased risk of information interference dropout yields — if an important neuron which modulates the input from an unimportant neuron is removed, the network could be influenced in a disadvantageous way. Extending this idea to successive inputs to each hidden layer of a deep network could additionally explain why dropout acts as an effective regularizer against overfitting, as noise or task-irrelevant features are minimized.

In addition to pruning irrelevant input information, dropout also causes networks to overrepresent important information. As shown in Fig. 2b, applying dropout significantly increases redundancy in the hidden layers for both logic gate tasks. As the risk of a neuron’s removal increases, the network must compensate by ensuring that a robust, redundant representation of important information remains, decreasing the reliance on individual (sets of) neurons (i.e., unique information). Although less efficient within the space of neuronal coding, redundant information acts as a robust representation. Thus, with limited information resources and the heightened risk of information loss, task-irrelevant information is more likely to be removed in favor of task-relevant information, which is instead overrepresented. Our finding provides an explicit measure of redundancy, complementing recent work (Lange et al. 2022) that has also suggested dropout to increase redundancy based on an increase in clustering driven by the similarity of neurons.



### 2.1.2. *Performance relies on synergistic neurons*

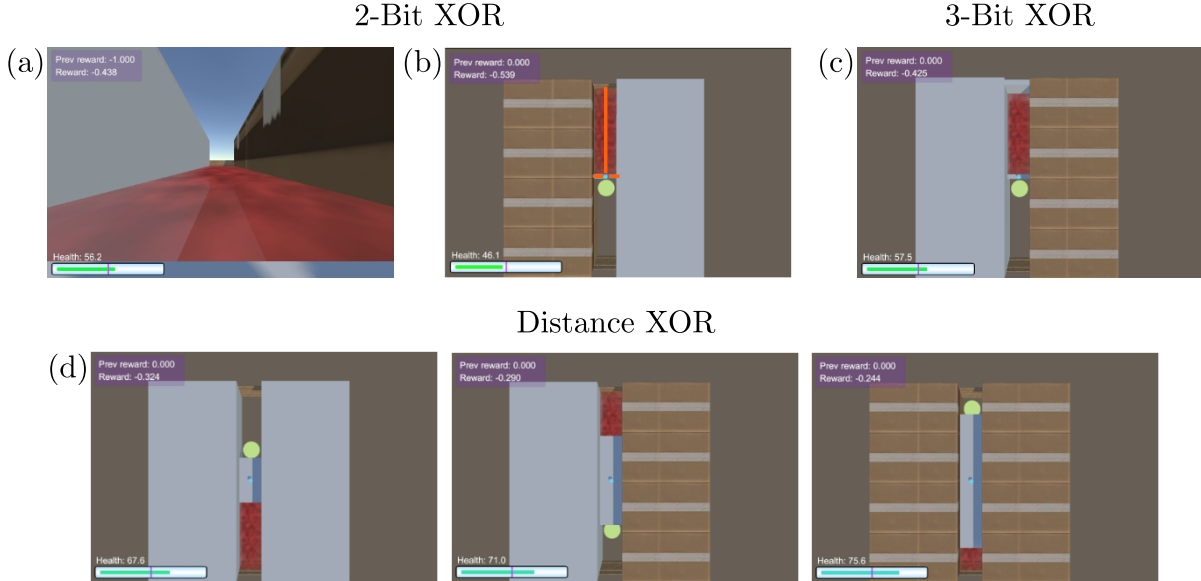
Using the trained logic gate networks, we perform lesioning experiments to evaluate whether a neuron’s synergy is predictive of its importance to the network. In lesioning experiments, each neuron’s pairwise synergy (average synergy with every other neuron in the same layer) is computed. We then permanently, iteratively remove the most (or least) synergistic neurons from each layer by setting their outgoing activations to 0 and evaluate subsequent performance. This differs from dropout in that lesioning is performed after training during evaluation, rather than during training, such that the network is unable to modify its parameters. By using networks with and without dropout applied, we can observe how dropout, and its resulting increased redundancy, influences the reliance on synergistic and non-synergistic neurons for performance.

Lesion experiments reveal that synergistic neurons are more critical for performance than non-synergistic neurons (Fig. 2c), especially for tasks requiring the integration of information (e.g., XOR). Synergistic neurons have less robust representations — the removal of one synergistic neuron can change the information carried by all of the sources acting synergistically with it, while the same is not true for redundant or unique information. Thus, synergy-rich neurons are more sensitive and their removal decreases performance more than synergy-poor neurons.

These results further show the effects of dropout on increasing robustness of the network via increased redundancy: with higher dropout, more minimally-synergistic neurons can be removed without disrupting performance as their information is overrepresented through an increase in redundancy. However, the XOR logic gate networks still remain highly sensitive to a disruption in synergistic neurons even after dropout is applied, exemplifying the importance of synergy for tasks requiring integration, as well as the vulnerability of such representations. Because the XOR logic gate requires the integration of input information to solve the task and because synergistic information between sets of sources is lost with the removal of a subset of sources, the XOR networks remain more sensitive to the removal of synergistic neurons. Conversely, because the COPY logic gate networks do not need to integrate information in order to solve the task, dropout instead reduces the number of synergistic neurons and the reliance on them such that their removal does not disrupt performance until over half of their neurons have been removed.

## 2.2. *Compositional tasks in 3D RL environment*

As a second step in our investigation, we extend the idea of solving logic gates to the context of reinforcement learning agents in Animal-AI (Crosby et al. 2020), a 3D environment with simulated physics used for assessing agents on cognitive common-sense physical reasoning tasks (Shanahan et al. 2020) (experimental details can be found in the Methods, Sec. 5.5). These experiments are motivated by the interest in studying synergy in the context of task-transfer and how agents may allocate their parameter space when learning a novel task after specializing on an initial task — and, more broadly, on how the structure of new tasks could influence information decomposition. In effect, contrasting with the previous tasks in which training happens directly on the desired input and output, agents in these reinforcement learning scenarios learn through trial-and-error interactions with their environment while trying to maximize a reward function. This significantly increases the difficulty of the task as reward is delayed, the environment contains a much larger state space, and



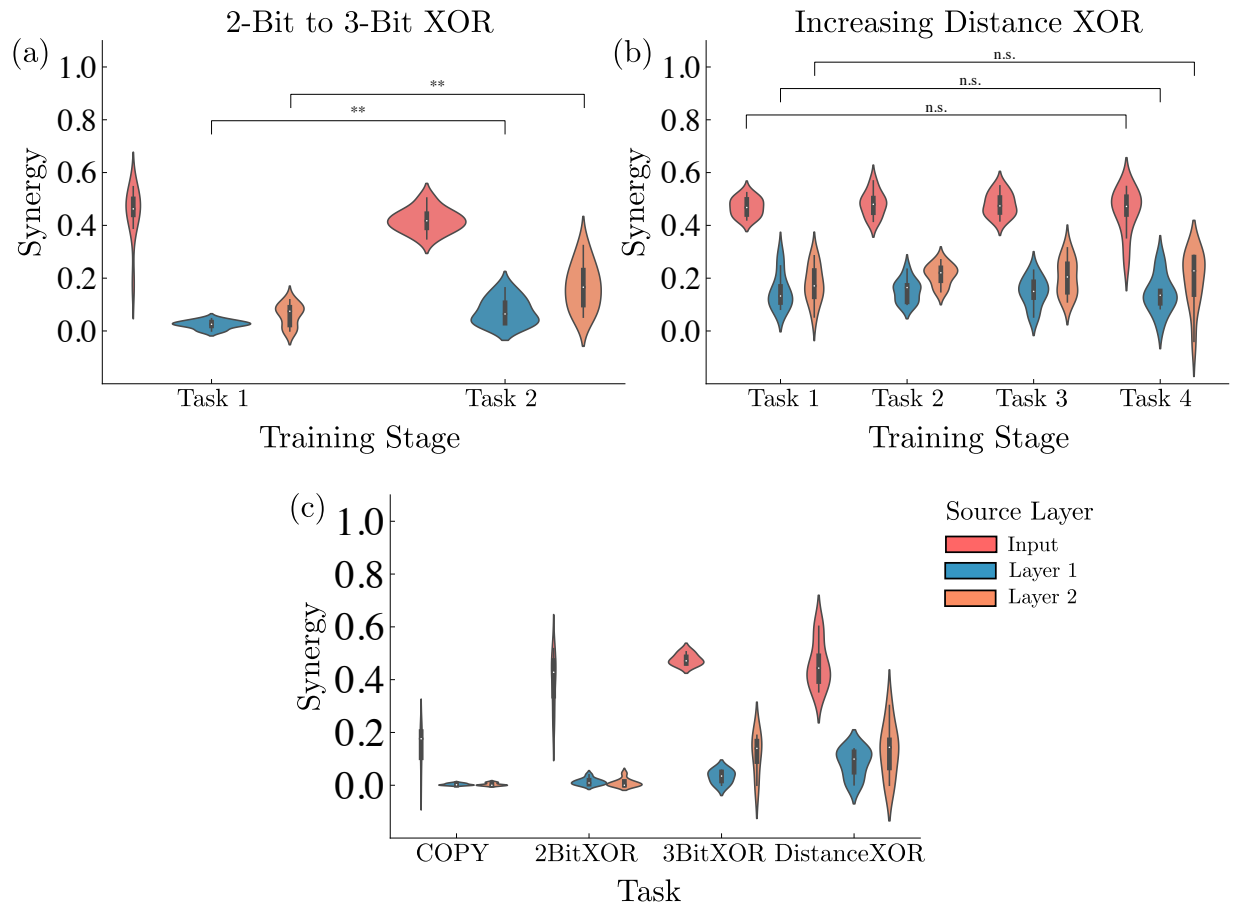
**Figure 3. Animal-AI tasks.** (a) An image visualization of the Animal-AI environment from the position of the agent for the ‘01’-input 2-Bit XOR configuration, where the ‘0’ gate output corresponds to a reward situated in the pit behind the agent. The agent receives raycast observations (rather than pixel inputs) of the environment from its position, which occlude the content of the pits. (b) An aerial view of the ‘10’-input 2-Bit XOR configuration, where the left cardboard barrier represents a ‘1’, the right wall barrier represents a ‘0,’ and the backward-relative-to-agent position of the green reward represents a ‘0’ output of the gate. The orange lines represent the orientation of the raycasts projected from the agent. (c) The configuration of the 3-Bit XOR task for ‘001’-input, where the additional barrier in front of the agent represents a third logic gate input. (d) Example configurations (each part of a separate task corresponding to platform-length) for the Distance XOR set.

the observation space includes additional (task-irrelevant) inputs about the environment. With these added challenges, the agent must use information about the arena walls to determine which platform direction (forward or backwards) to move to and act accordingly in order to retrieve the positive reward and solve the task.

### 2.2.1. Information profiles reflect specific task demands

We train models to perform either an individual task or a set of tasks in sequence, usually denoted in AI research as a *curriculum*. Each task consists of several environment configurations (each corresponding to a configuration of logic gate inputs) that are interleaved across episodes. In the curriculum experiments, models are trained on each task until reaching a reward threshold or a maximum number of training steps, after which they are trained on a new task for a set number of steps. Each considered task follows a similar design in which the agent aims to solve a problem based on the object-type of the barriers surrounding it (Fig. 3). Specifically, the agent receives as input three raycasts (indicating the type of the objects to its front, left, and right), and as output the agent can move either forward or backward into a pit to obtain a reward. The object-type of the barriers encodes three input bits (wall being 0, cardboard being 1), and the position of the reward encodes the correct output (forward being 0, backward being 1). To successfully solve the task, the agent must





**Figure 4. Relation of compositional tasks and synergy in Animal-AI.** (a) Synergy significantly increases from the end of training on the 2-Bit XOR task to the end of training on the 3-Bit XOR task, driven by the addition of a third source of information to integrate. (b) Synergy does not significantly increase in the Distance XOR curriculum, likely due to the fact that each subsequent task requires the same integration of sources and only varies in the number of times an action must be performed to reach the reward. This could be done easily by ignoring global position and object distance or by extending the number of environment states associated with an action, without additional integration of information. (c) The information decomposition of a network is influenced by the task it’s trained on and is qualitatively different across tasks. Distance XOR refers to Distance 10 XOR.

decide based on the barrier types which direction to move, in order to retrieve the positive reward. The direction of the reward corresponding to the configuration of the barrier types is determined by the logic gate task being performed, which includes the same gates as in Sec. 2.1 (2-Bit COPY and 2-Bit XOR; Fig. 3b), plus a 3-Bit extension of the XOR gate (where the correct output is the parity of all inputs; Fig. 3c), and a “Distance XOR” task where the length of the platform is increased, introducing a longer delay between action and reward (Fig. 3d).

Our results show that networks increase their synergy as they learn new tasks, and that the integration of an additional source of information specifically drives this behavior (see Fig. 4a). Synergy significantly increases from the 2-Bit XOR task to the 3-Bit XOR task, despite agents not

being able to learn the second task to perfect accuracy, whereas synergy remains constant across all Distance XOR tasks (Fig. 4b), even when learned accurately. Although the 3-Bit XOR alone does yield more synergy than the 2-Bit XOR task alone, even when the 3-Bit XOR task is not accurately learned in the curriculum, the learning process of the new task drives an increase in synergy.

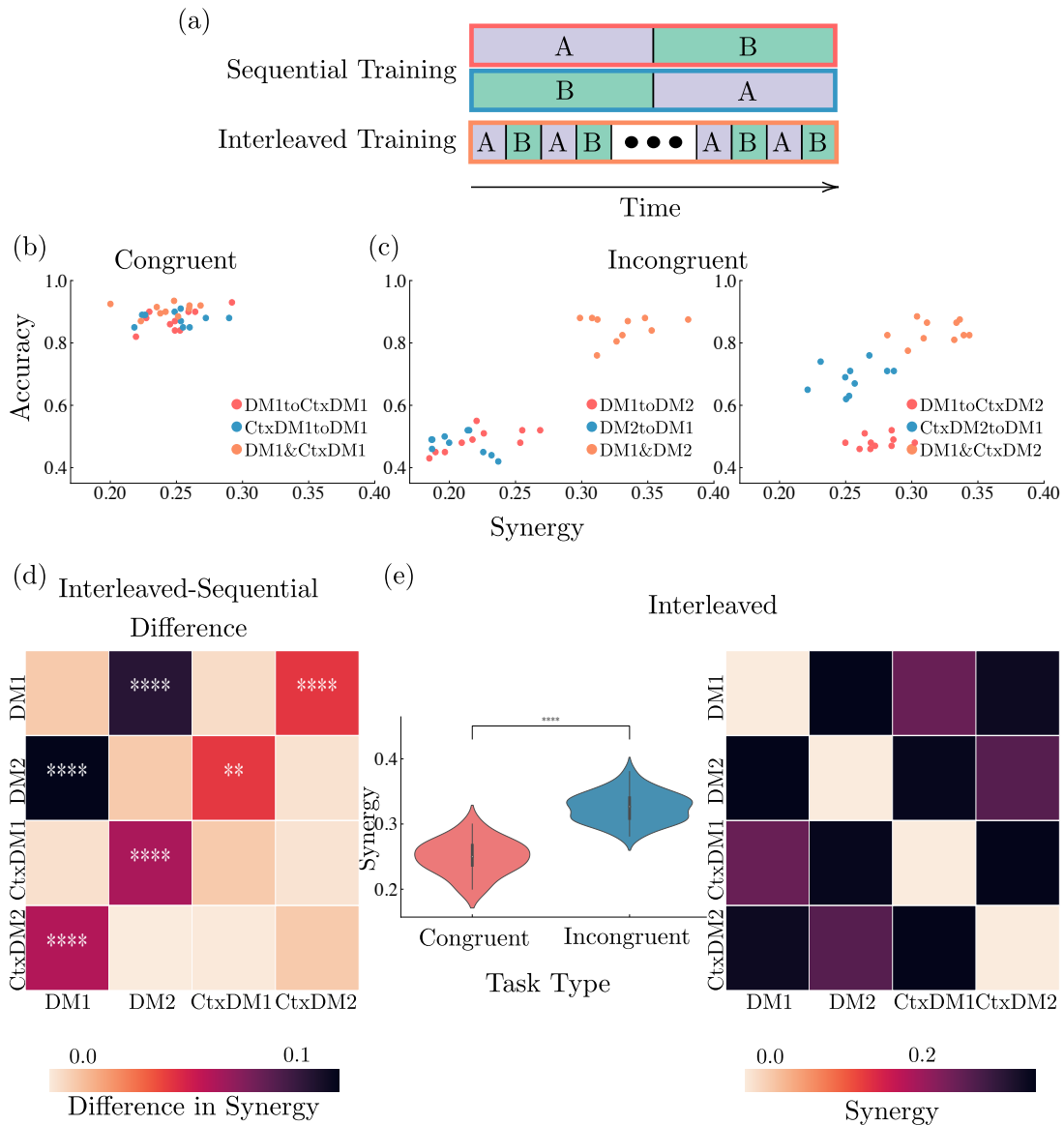
Conversely, the Distance XOR curriculum does not drive an increase in synergy, although agents successfully learn all tasks. The main factor distinguishing the 2-Bit to 3-Bit XOR curriculum from the Distance XOR curriculum is that the former requires integration of an additional source (3 bits instead of 2). Instead, the latter requires learning new tasks, but does not require the incorporation of any new sources of information or any modification in its processing — only the association of more states with a particular learned action using the same mapping from existing sources. Overall, these differences in synergy between curricula highlight a difference in the complexity of the set of tasks being learned, as sets of tasks that are more simple (e.g., using few sources of information in a specialized manner) do not drive an increase in synergy, while more complex tasks dependent on the integration of several sources do. This suggests that synergy is specifically related to the learning of multiple complex tasks in which multiple information sources have to be integrated to yield novel behavior.

### 2.3. *Effect of learning multiple diverse cognitive tasks on synergy*

Our final set of experiments seeks to analyze networks learning multiple cognitively-inspired tasks, and investigate how learning a set of tasks requiring the capacity to integrate different modalities compares to a set of tasks relying on a single modality. When trained on a non-stationary sequence of tasks (for example, a sequential curriculum), neural networks often suffer from interference as they update parameters to learn a new task that potentially disrupt setups that were important for solving previous tasks, resulting in a phenomenon known as ‘catastrophic forgetting’ (Kirkpatrick et al. 2017). Analyzing networks that can remember (rather than forget) their entire training curriculum allows us to study how information profiles are influenced by learning and solving multiple (versus single) tasks. With this motivation, we use two different training protocols for each pair of tasks: sequential training, which does not prevent forgetting the first task trained on, and interleaved training, which forces the network to retain the capacity to solve both tasks.

#### 2.3.1. *Studying RNNs learning decision making tasks*

We explore information decomposition in the hidden layer of recurrent neural networks (RNN) — a simple form of memory — on tasks requiring integration over time. This model is used to solve tasks taken from the NeuroGym environment (Molano-Mazón et al. 2022) — specifically, a subset of their collection of decision-making tasks, referred to as the ‘DM family’ (DM1, DM2, CtxDM1, CtxDM2) in Yang et al. (2019b). This set of tasks is based on various decision-making tasks used in neuroscience and psychophysics (Gold and Shadlen 2007; Mante et al. 2013; Raposo et al. 2014), involving the presentation of two simultaneous stimuli (e.g., numeric values of input dimensions), after which models must indicate which stimulus is stronger (higher in value). Depending on the specific task (i.e., in contextual decision-making tasks, CtxDM1, CtxDM2), stimulus information presented in a separate modality (set of input dimensions) may also need to be ignored. We refer to the reader to Molano-Mazón et al. (2022); Yang et al. (2019b) for a detailed description of all aforementioned tasks. For each task, a set of consecutive trials is given as input to the RNN. Each



**Figure 5. Synergy increases with the solving of incongruent tasks.** (a) In the sequential protocol, networks are trained on each task in a single block sequentially (red and blue in panels **b-c**) and forgetting can occur. In the interleaved protocol, networks are trained on both tasks simultaneously (orange in panels **b-c**) and thus are both solved. (b) Congruent tasks yield similar accuracy and synergy for both training protocols. (c) Incongruent tasks yield distinct accuracy-synergy clusters based on training protocol. Sequential training yields lower synergy corresponding to worse performance, while interleaved training performs better with higher synergy. (d) Training with interleaving yields significantly higher synergy than training sequentially for incongruent tasks, but not for congruent tasks. (e) Networks accommodating two incongruent tasks (via interleaving) yield significantly higher synergy than those accommodating two congruent tasks.

trial consists of an initial fixation phase, followed by stimuli, and ending with a decision-making period. To successfully solve a task, the model must integrate stimulus information over time to make its final decision.

To better understand the relationship between synergy and the curriculum of tasks being learned, we compare networks trained on pairs of tasks we define to be congruent or incongruent. In this context, congruence refers to the similarity between both tasks and whether learning one task may aid in learning of the other task (such that transfer of performance is possible). In our experimental design, congruent tasks are defined as the pairing of decision-making tasks using the same modality (DM1&CtxDM1; DM2&CtxDM2) and incongruent tasks are those using different modalities (DM1&DM2; DM1&CtxDM2; DM2&CtxDM1; CtxDM1&CtxDM2). Learning congruent tasks requires attending to and integrating information from a single input modality and ignoring stimuli in the other modality, while learning incongruent tasks requires switching between attending to one and ignoring the other of two modalities, depending on the task. We further compare networks learning a curriculum of tasks sequentially and networks learning both tasks simultaneously through interleaving (Fig. 5a). We do not use any continual learning methods to prevent catastrophic forgetting in either case. Therefore, sequential learning presents a condition where forgetting a previously experienced task could potentially occur while interleaved learning forces the network to solve both tasks.

### 2.3.2. Solving diverse sets of tasks increases synergy

Our results show that networks trained on congruent tasks yield similar levels of synergy and accuracy in both sequential and interleaved training regimes (Fig. 5b). The similarity of congruent tasks allows for easy transfer, as their learned parameters can be reused across tasks without strong interference, achieving comparable performance in both sequential and interleaved protocols. Because both tasks are similar and only require attending to and integrating from a single modality, networks trained sequentially can reuse their representations for both tasks and achieve adequate performance, without accommodating additional information. We suggest that because two congruent tasks require information to be integrated from only a single modality, levels of synergy are lower for all training settings than for two incongruent tasks.

In contrast, networks trained on incongruent tasks fall in different clusters depending on the training regime (Fig. 5c). In effect, interleaved training yields clusters with higher synergy and accuracy, whereas sequential training yields clusters with lower synergy and accuracy. This occurs because sequential training of incongruent tasks leads to some forgetting of the first task due to interference, leading to lower accuracy, while interleaved training forces the network to solve both tasks resulting in higher accuracy. Furthermore, networks trained with interleaving have an increased capacity to integrate information from both modalities in order to adequately perform both incongruent tasks, resulting in a higher level of synergy compared to networks trained sequentially.

Overall, in networks trained with interleaving, the amount of synergy is significantly higher for incongruent tasks compared to congruent tasks, and is consistently higher for individual pairs of incongruent tasks than pairs of congruent tasks (see Fig. 5e). This suggests that the capacity to integrate and use information from several modalities (as in incongruent tasks) results in a higher proportion of synergistic information compared to that yielded by a single modality (as in congruent tasks). In other words, synergy is related to a system’s ability to combine different sources of information flexibly for distinct tasks. We speculate that this insight may help to explain recent empirical findings showing a higher synergy associated with the brain’s associative cortices, which

integrate information from multiple sensory systems, in contrast to unimodal brain areas such as sensory or motor cortices (Luppi et al. 2022).

Furthermore, by contrasting networks with interleaved vs. sequential training we can attribute this synergy increase specifically to the capability of simultaneously solving multiple incongruent tasks (Fig. 5d). Sequential training of incongruent tasks results in the forgetting of the first task and a corresponding drop in performance as well as a drop in synergy (compared to interleaved training), which is related to the fact that the network is only using information from one modality at a time. This doesn't hold for congruent tasks, however, for which both interleaved and sequential training result in similar levels of both performance and synergy. This results in synergy being specifically linked with scenarios where the network needs flexibility to adapt between tasks requiring integration of different modalities — and, in these cases, synergy results to be highly correlated with performance.

### 3. DISCUSSION

This paper presents a series of experiments using neural networks in a variety of tasks and learning settings, and examines their internal representations using the Partial Information Decomposition (PID) framework (Williams and Beer 2010). Based on our results, we draw several interpretations of the functional roles played by different forms of information, and suggest their relation to learning in artificial and biological neural networks.

#### 3.1. *Functional roles of information atoms*

We start by delineating the functional roles of redundancy, unique information, and synergy in learning contexts.

*Redundancy*—Although redundant information makes less efficient use of neuronal capacity, it grants robustness to the network. In effect, over-representing important information (i.e., encoding it in multiple units) is an effective way to ensure that it will be propagated through the network, even in the presence of possible interferences. Applying perturbations to the network during training — such as dropout — incentivizes such reliability, which naturally leads to an increase in redundancy. Furthermore, after training with dropout, these networks resist performance drops caused by lesions to a much higher degree than networks trained without dropout. In the human brain, redundancy dominates lower-level cortical regions, particularly in sensory and motor areas (Luppi et al. 2022) — which could be due to a similar need to overrepresent important sensory information in order to extract critical features of incoming data and resist perturbations (e.g., noise).

These findings have interesting parallels with previous analyses of neural networks based on the information bottleneck principle (Tishby et al. 1999), which uses mutual information to bound optimal networks trading off information loss due to compression and information preserved about a desired output (Tishby and Zaslavsky 2015; Shwartz-Ziv and Tishby 2017; Saxe et al. 2018). In agreement with that line of work, the results of our logic gate experiments show that dropout acts as an information bottleneck through its increase of redundant information and pruning of task-irrelevant information — the mutual information between the network and the input is decreased (via removal of task-irrelevant features), while the mutual information about the desired output is

increased (via higher redundancy). However, our results also reveal phenomena beyond the usual formulation of the bottleneck: our work shows that dropout not only changes the total mutual information, but it also changes the structure of that information by altering its composition. This finding emphasizes the fact that it is not only the content of a task, but also how it is learned, that affects the information processing strategy adopted by a cognitive system.

*Unique information*—Unique representations provide specialized encoding of information, whereby encoding can take place in a single neuron rather than requiring a set of neurons to operate as a distributed representation (as is the case for synergy). Such specialized representations are particularly efficient when a network consistently performs the same task or several tasks with the same substructure, or when it is not required to integrate multiple sources of information.

We speculate that the utility of unique information could be harnessed by functionally specialized circuits yielded by evolution and early development, especially in sensory cortices. The one-to-one mapping of simple receptive fields in the primary visual cortex is a clear example of such specialization in the brain. High-level cognition, in contrast, requires the integration of information from several cortical regions being in principle more diverse and less stable than processes associated with low-level feature extraction. Nevertheless, regardless of what cognitive process is occurring in higher cortical regions based on visual input, the visual cortex remains a specialized and stable region for the subtask of visual feature extraction, rather than other sensory functions. For these types of tasks, the flexibility that may be provided by synergistic representations and the associated learning processes are not necessarily beneficial. The empirical evaluation of these conjectures is an interesting avenue for future investigations.

*Synergy*—Perhaps one of the greatest functional advantages of synergy is that it can encode more information than other information atoms for a given population size (Rosas et al. 2020). In effect, in contrast with redundancy and unique information, synergy relies on combinations of neurons, which make its informational capacity to be exponential with the system size. Synergistic information, however, is also more vulnerable to noise (Mediano et al. 2022), because a disruption in a single source could disrupt information synergistically held together with other sources. This vulnerability to noise may partially explain why networks exhibit higher levels of synergy at comparatively lower orders rather than higher orders (Tax et al. 2017) as an attempt to minimize the effect of losing synergistic information with several neurons.

Our experiments find that the removal of maximally synergistic neurons yields a larger drop in performance compared to minimally synergistic neurons. Building in the above discussion, this finding can be explained along three — not mutually exclusive — lines of reasoning: (1) synergistic neurons may be encoding more information collectively than other forms of information are individually, (2) this information is necessary for integrating several sources of information, and (3) this total information is more vulnerable to perturbations because it can be altered by a disruption in a



single source. Future work may seek to disentangle the relevance of each of these potential causes in driving this effect.

### 3.2. *Synergy facilitates flexible general-purpose learning*

In addition to the functional roles highlighted above, an important overarching result is the association between synergy and the learning of multiple tasks. Results across all our experiments demonstrate that synergy supports integration in networks solving multiple tasks. Specifically, our logic gate experiments show that performance relies on synergy-rich neurons for tasks requiring the integration of multiple sources of information, and that these neurons are more sensitive to perturbations. Results in the Animal-AI experiments show that the incorporation of additional sources of information when learning multiple tasks drives an increase in synergy, as opposed to learning multiple tasks that do not substantially change the processing of information sources. Thus, the complexity (loosely understood here as relying on integrating various information sources in flexible and diverse ways) of a task set relates to the degree of synergistic information processing it elicits. Finally, results in NeuroGym decision-making tasks show that synergy increases with a network’s capacity to simultaneously integrate information from several modalities for different tasks. Taken together, these findings support a link between synergistic information processing and the ability to perform multiple complex tasks, which require the flexibility to integrate and process various sources of information in different ways. In addition to this, in the following we hypothesize several other functional advantages synergy could provide in learning systems performing several tasks.

The first hypothesis is that, in addition to providing additional capacity for modality integration, synergy could be a response to the learning pressure of having to encode more task-relevant information overall in the neuronal information space — which, in the case of our experiments, is severely constrained by size. Thus, in order to successfully solve two or more tasks requiring the integration of diverse sources, representations may need to be encoded in increasingly efficient manners — a feature which is provided more readily by relying on synergy than unique or redundant information. Performing increasingly complex tasks could have a similar effect as the number of task-parameters and variables increases, and a higher number of information sources and proportion of information must be integrated and encoded in a network. In this way, a complex task could just be seen as a collection of many smaller and simpler tasks.

A second hypothesized utility of synergy is its ability to represent information in a structurally different way. Whereas unique information can provide information along a single dimension for each feature encoded by a source, synergy could provide higher-dimensional representations across several neurons signifying distance between different features and neuronal encodings. This could aid in representing structure and similarity across and within tasks, providing increased flexibility for generalization. Alternatively, synergistic information could be used for integrating other (possibly specialized) representations, potentially occupying some low-dimensional subspace.

### 3.3. *Implications for cognitive and computational neuroscience*

Our results and hypothesized functions for synergy in the context of general-purpose learning are complementary to those observed in the brain. Within cognitive science, the division of complex tasks into simpler sub-tasks has been proposed for many decades as an important mechanism for

solving almost any cognitive goal (Newell 1990). In addition, this parsing of tasks into sub-tasks has been associated with the prefrontal parietal network in brain-scanning experiments (Bor et al. 2003; Bor and Owen 2006; Duncan 2010, 2013). This brain network has been shown to preferentially support synergistic forms of information processing. For instance, Luppi et al. (2022) found that synergistic information is higher in regions of the brain that are especially responsible for complex human cognition (particularly association cortices in frontal and parietal regions, including the default mode and executive control networks), and that redundant information is higher in areas responsible for perception and low-level cognition (particularly primary motor and sensory cortices). The high levels of synergy in higher cortical regions may be explained by their need and capacity to integrate information from other brain regions, encode and learn a vast information space throughout life, and reuse or relate this information flexibly in order to generalize to new settings (Duncan et al. 2000; Miller and Cohen 2001; Stokes et al. 2013).

In addition, just as synergy could contribute to creating higher level representations and structure in the context of artificial neural networks, synergistic brain networks tightly overlap with regions that integrate information (Luppi et al. 2020). Whereas low-level brain regions may benefit from having redundant and specialized representations that are more robust and, correspondingly, comparatively less adaptable after development, higher cortical areas continue to have the responsibility of learning throughout life, reflected in their information decomposition.

Related to this paper, deep learning has been previously used to study representations elicited by networks performing multiple tasks. Prior work has shown that distinct functional units emerge in networks trained on multiple tasks, becoming specialized to specific sub-task features (Yang et al. 2019b), that networks learning multiple tasks naturally produce abstract representations (Johnston and Fusi 2021), and that networks performing context-dependent tasks use task representations lying on a low-dimensional and orthogonal manifold (Flesch et al. 2022). It remains an open question how and where information decomposition fits into these and other prior findings, which future work should investigate.

In addition to neurons that respond selectively to particular stimuli (pure selectivity), there also exist, especially in higher-order cortical regions such as the prefrontal cortex and hippocampus, neurons that respond to diverse sets of stimuli and tasks, rather than performing a single specialized function. These sets of neurons are said to have (nonlinear) mixed selectivity (NMS), exhibiting complex responses to different task parameters. Recent work has shown that NMS neurons support flexible behavior and complex cognition (Rigotti et al. 2013). We suggest that synergistic information processing may be closely related to NMS neurons in the brain. In particular, we predict that neurons exhibiting NMS are synergistic and neurons exhibiting pure selectivity have predominantly unique and redundant representations. Future work should study the relation of information decomposition to neuronal selectivity, as it could provide new approaches for understanding the information processing of various neuronal populations.

#### 3.4. *Limitations and future work*

PID is a relatively recent theoretical framework, and as such its practical applications are often faced with certain limitations. Perhaps the most important of these is that the number of PID atoms

grows super-exponentially with the number of sources (Williams and Beer 2010) — a problem that we bypass here by averaging across small subsets of neurons. Furthermore, it seems increasingly clear that there is no universal redundancy function, and that different formulations capture different aspects of multivariate information. In this paper we address these issues by validating our analyses by using two different redundancy functions ( $I_{\text{MMI}}$  (Barrett 2015), and  $I_{\text{min}}$  (Williams and Beer 2010)). Future work should investigate ways of scaling PID to larger systems and clarify the relationship between different redundancy functions.

In addition to the PID-specific issues above, there is a more general difficulty that arises when computing information-theoretic quantities from data: to estimate these, one needs to know (or accurately approximate) the probability distribution of the observed data (Lizier 2014). This is particularly challenging in neural networks with non-linearities, such as rectified linear units (ReLUs). In the specific case of neural networks, this issue has caused extensive debate (Shwartz-Ziv and Tishby 2017; Saxe et al. 2018). Here we mitigate this problem by verifying that our results are consistent with two estimators, both discrete and continuous, and with different hyperparameter settings. Nonetheless, future work should elaborate on this direction by either using more sophisticated estimators (Lizier 2014), or using networks where distributions are easier to calculate analytically (e.g., deep linear networks (Saxe et al. 2014, 2019)).

Finally, it is worth mentioning that all the networks used here are small (on the order of tens of neurons) compared to state-of-the-art networks (typically on the order of many thousands, or more). Thus, future work should investigate to what extent the results presented here generalize to larger networks — as the ones often used nowadays in a range of applications. That being said, the consistency of our results across training regimes (supervised learning, and reinforcement learning with and without recurrent networks), as well as across experiment suites (logic gates, Animal-AI, NeuroGym) constitute encouraging preliminary evidence. After the limitations of scaling and approximation above have been overcome, future work should try to replicate these results with larger networks, more complex tasks, and different architectures and training hyperparameters.

#### 4. CONCLUSION

In this work we used information decomposition to analyze how artificial neural networks process information in a variety of experimental settings. By studying the learning of logic gates, we found that performance depends on synergistic neurons in tasks requiring the integration of information, and that randomly turning off neurons during training with dropout increases redundancy and robustness while minimizing task-irrelevant features. Using a 3D environment with simulated physics, we showed that synergy is driven by the integration of additional sources of information in a complex reinforcement learning setting. Finally, by studying decision-making tasks inspired by cognitive neuroscience, we found that synergy is specifically increased by the solving of multiple incongruent tasks and the capacity to integrate information from several modalities. Based on these findings, we suggest specific functional roles for PID atoms in the context of learning: redundancy for robustness, unique information for specialization, and synergy for modality integration, flexibility, and efficient encoding. These results lay down foundations to study how learning scenarios modulate information processing modalities, while providing insights into existing cognitive neuroscience results — where synergy is especially high in the most functionally flexible cortical regions (association cortices in

frontal and parietal regions), and redundancy has been found in the most functionally specialized and robust areas of the cortex, including sensory and motor regions.

- 1 F.E.R. is supported by the Ad Astra Chandaria foundation. A.I.L. is funded by the Gates Cambridge
- 2 Trust. D.B. is funded by the Wellcome Trust (grant no. 210920/Z/18/Z).

## REFERENCES

- Badia, A. P., Piot, B., Kapturowski, S., Sprechmann, P., Vitvitskyi, A., Guo, D., and Blundell, C. (2020). Agent57: Outperforming the atari human benchmark. *International Conference on Machine Learning*.
- Barrett, A. B. (2015). Exploration of synergistic and redundant information sharing in static and dynamical gaussian systems. *Physical Review E*, 91:052802.
- Bor, D., Duncan, J., Wiseman, R. J., and Owen, A. M. (2003). Encoding strategies dissociate prefrontal activity from working memory demand. *Neuron*, 37(2):361–367.
- Bor, D. and Owen, A. M. (2006). A Common Prefrontal–Parietal Network for Mnemonic and Mathematical Recoding Strategies within Working Memory. *Cerebral Cortex*, 17(4):778–786.
- Crosby, M., Beyret, B., Shanahan, M., Hernández-Orallo, J., Cheke, L., and Halina, M. (2020). The Animal-AI testbed and competition. In *Proceedings of the NeurIPS 2019 Competition and Demonstration Track*, volume 123 of *Proceedings of Machine Learning Research*, pages 164–176. PMLR.
- Duncan, J. (2010). The multiple-demand (md) system of the primate brain: mental programs for intelligent behaviour. *Trends in Cognitive Sciences*, 14:172–179.
- Duncan, J. (2013). The structure of cognition: Attentional episodes in mind and brain. *Neuron*, 80(1):35–50.
- Duncan, J., Seitz, R. J., Kolodny, J. A., Bor, D., Herzog, H., Ahmed, A., Newell, F. N., and Emslie, H. (2000). A neural basis for general intelligence. *Science*, 289 5478:457–60.
- Eliasmith, C. and Anderson, C. H. (2004). Neural engineering: Computation, representation, and dynamics in neurobiological systems. *IEEE Transactions on Neural Networks*, 15:528–529.
- Flesch, T., Juechems, K., Dumbalska, T., Saxe, A. M., and Summerfield, C. (2022). Orthogonal representations for robust context-dependent task performance in brains and neural networks. *Neuron*, 110:1258 – 1270.e11.
- Gatica, M., Cofré, R., Mediano, P. A., Rosas, F. E., Orio, P., Diez, I., Swinnen, S. P., and Cortes, J. M. (2021). High-order interdependencies in the aging brain. *Brain Connectivity*, 11(9):734–744.
- Gold, J. I. and Shadlen, M. N. (2007). The neural basis of decision making. *Annual Review of Neuroscience*, 30:535–74.
- Ince, R. A., Giordano, B. L., Kayser, C., Rousset, G. A., Gross, J., and Schyns, P. G. (2017). A statistical framework for neuroimaging data analysis based on mutual information estimated via a Gaussian copula. *Human Brain Mapping*, 38(3):1541–1573.
- Izhikevich, E. M. (2007). *Dynamical Systems in Neuroscience*. MIT Press.
- James, R. G., Ellison, C. J., and Crutchfield, J. P. (2018). dit: a Python package for discrete information theory. *The Journal of Open Source Software*, 3(25):738.
- Johnston, W. J. and Fusi, S. (2021). Abstract representations emerge naturally in neural networks trained to perform multiple tasks. *bioRxiv*.
- Kirkpatrick, J., Pascanu, R., Rabinowitz, N. C., Veness, J., Desjardins, G., Rusu, A. A., Milan, K., Quan, J., Ramalho, T., Grabska-Barwinska, A., Hassabis, D., Clopath, C., Kumaran, D., and Hadsell, R. (2017). Overcoming catastrophic forgetting in neural networks. *Proceedings of the National Academy of Sciences*, 114:3521 – 3526.
- Kriegeskorte, N. (2015). Deep neural networks: A new framework for modeling biological vision and brain information processing. *Annual Review of Vision Science*, 1:417–446.

- Lange, R. D., Rolnick, D., and Kording, K. (2022). Clustering units in neural networks: Upstream vs downstream information. *Transactions on Machine Learning Research*.
- Lizier, J. T. (2012). *The Local Information Dynamics of Distributed Computation in Complex Systems*. Springer Science & Business Media.
- Lizier, J. T. (2014). JIDT: An information-theoretic toolkit for studying the dynamics of complex systems. *Frontiers in Robotics and AI*, 1:11.
- Luppi, A., Mediano, P., Rosas, F., Holland, N., Fryer, T., O'Brien, J., Rowe, J., Menon, D., Bor, D., and Stamatakis, E. (2022). A synergistic core for human brain evolution and cognition. *Nature Neuroscience*.
- Luppi, A. I., Mediano, P. A. M., Rosas, F. E., Allanson, J., Pickard, J. D., Carhart-Harris, R. L., Williams, G. B., Craig, M. M., Finoia, P., Owen, A. M., Naci, L., Menon, D. K., Bor, D., and Stamatakis, E. A. (2020). A synergistic workspace for human consciousness revealed by integrated information decomposition. *bioRxiv*.
- Mante, V., Sussillo, D., Shenoy, K. V., and Newsome, W. T. (2013). Context-dependent computation by recurrent dynamics in prefrontal cortex. *Nature*, 503:78 – 84.
- Marblestone, A. H., Wayne, G., and Kording, K. P. (2016). Toward an integration of deep learning and neuroscience. *Frontiers in Computational Neuroscience*, 10:94.
- Mediano, P. A., Rosas, F. E., Farah, J. C., Shanahan, M., Bor, D., and Barrett, A. B. (2022). Integrated information as a common signature of dynamical and information-processing complexity. *Chaos: An Interdisciplinary Journal of Nonlinear Science*, 32(1):013115.
- Mediano, P. A., Rosas, F. E., Luppi, A. I., Carhart-Harris, R. L., Bor, D., Seth, A. K., and Barrett, A. B. (2021). Towards an extended taxonomy of information dynamics via Integrated Information Decomposition. *arXiv*.
- Miller, E. K. and Cohen, J. D. (2001). An integrative theory of prefrontal cortex function. *Annual Review of Neuroscience*, 24:167–202.
- Molano-Mazón, M., Barbosa, J., Pastor-Ciurana, J., Fradera, M., Zhang, R.-Y., Forest, J., del Pozo, J., Ji-An, L., Cueva, C., de la Rocha, J., Narain, D., and Yang, G. R. (2022). NeuroGym: An open resource for developing and sharing neuroscience tasks.
- Newell, A. (1990). *Unified Theories of Cognition*. Harvard University Press.
- Panzeri, S., Senatore, R., Montemurro, M. A., and Petersen, R. S. (2007). Correcting for the sampling bias problem in spike train information measures. *Journal of Neurophysiology*, 98 3:1064–72.
- Raffin, A., Hill, A., Gleave, A., Kanervisto, A., Ernestus, M., and Dormann, N. (2021). Stable-Baselines3: Reliable reinforcement learning implementations. *Journal of Machine Learning Research*, 22(268):1–8.
- Raposo, D., Kaufman, M. T., and Churchland, A. K. (2014). A category-free neural population supports evolving demands during decision-making. *Nature Neuroscience*, 17:1784 – 1792.
- Reed, S., Zolna, K., Parisotto, E., Colmenarejo, S. G., Novikov, A., Barth-Maron, G., Gimenez, M., Sulsky, Y., Kay, J., Springenberg, J. T., Eccles, T., Bruce, J., Razavi, A., Edwards, A., Heess, N., Chen, Y., Hadsell, R., Vinyals, O., Bordbar, M., and de Freitas, N. (2022). A generalist agent. *arXiv*.
- Rigotti, M., Barak, O., Warden, M. R., Wang, X.-J., Daw, N. D., Miller, E. K., and Fusi, S. (2013). The importance of mixed selectivity in complex cognitive tasks. *Nature*, 497:585–590.
- Rosas, F., Mediano, P. A., Ugarte, M., and Jensen, H. J. (2018). An information-theoretic approach to self-organisation: Emergence of complex interdependencies in coupled dynamical systems. *Entropy*, 20(10):793.
- Rosas, F. E., Mediano, P. A. M., Gastpar, M., and Jensen, H. J. (2019). Quantifying high-order interdependencies via multivariate extensions of the mutual information. *Physical Review E*, 100:032305.
- Rosas, F. E., Mediano, P. A. M., Rassouli, B., and Barrett, A. (2020). An operational information decomposition via synergistic disclosure. *Journal of Physics A: Mathematical and Theoretical*, 53.

- Saxe, A. M., Bansal, Y., Dapello, J., Advani, M., Kolchinsky, A., Tracey, B. D., and Cox, D. D. (2018). On the information bottleneck theory of deep learning.
- Saxe, A. M., McClelland, J. L., and Ganguli, S. (2014). Exact solutions to the nonlinear dynamics of learning in deep linear neural networks. *International Conference on Learning Representations*.
- Saxe, A. M., McClelland, J. L., and Ganguli, S. (2019). A mathematical theory of semantic development in deep neural networks. *Proceedings of the National Academy of Sciences*, 116:11537 – 11546.
- Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and Klimov, O. (2017). Proximal policy optimization algorithms. *CoRR*.
- Shanahan, M., Crosby, M., Beyret, B., and Cheke, L. (2020). Artificial intelligence and the common sense of animals. *Trends in Cognitive Sciences*, 24(11):862–872.
- Shwartz-Ziv, R. and Tishby, N. (2017). Opening the black box of deep neural networks via information. *CoRR*.
- Stokes, M., Kusunoki, M., Sigala, N., Nili, H., Gaffan, D., and Duncan, J. (2013). Dynamic coding for cognitive control in prefrontal cortex. *Neuron*, 78(2):364–375.
- Tax, T. M., Mediano, P. A., and Shanahan, M. (2017). The partial information decomposition of generative neural network models. *Entropy*, 19(9).
- Tishby, N., Pereira, F. C., and Bialek, W. (1999). The information bottleneck method. In *Proc. of the 37-th Annual Allerton Conference on Communication, Control and Computing*, pages 368–377.
- Tishby, N. and Zaslavsky, N. (2015). Deep learning and the information bottleneck principle. In *2015 IEEE Information Theory Workshop (ITW)*, pages 1–5.
- Varley, T. F. and Hoel, E. P. (2022). Emergence as the conversion of information: A unifying theory. *Philosophical Transactions. Series A, Mathematical, physical, and engineering sciences*, 380.
- Williams, P. L. and Beer, R. D. (2010). Nonnegative decomposition of multivariate information. *arXiv*.
- Yang, G. R., Cole, M. W., and Rajan, K. (2019a). How to study the neural mechanisms of multiple tasks. *Current Opinion in Behavioral Sciences*, 29:134–143.
- Yang, G. R., Joglekar, M. R., Song, H. F., Newsome, W. T., and Wang, X.-J. (2019b). Task representations in neural networks trained to perform many cognitive tasks. *Nature Neuroscience*, 22:297–306.



## 5. METHODS

### 5.1. Model architectures

For each task and setting, an ensemble of 10 models is trained, with each network initialized using a different random seed. Additionally, all models use either rectified linear unit (ReLU) or leaky ReLU (for computing continuous measures) activation functions between each linear layer to avoid the compression of mutual information associated with double-saturating nonlinearities, as described in Saxe et al. (2018). Network architectures varied slightly between experiments (see details in Sec. 5.4 to 5.6), but in general they consisted of either one or two layers with ten neurons each.

### 5.2. Quantifying information decomposition

We compute information decomposition over the sampled activations of a network during testing, with its weights frozen (i.e., after first training it on a task). For curriculum tasks, models are tested on all configurations within the curriculum they are trained on, ensuring a fair comparison between different training points within a curriculum. The activations sampled during the testing phase are then used to compute distributions over the activity of the network, used for quantifying redundancy and synergy. For our experiments performed in NeuroGym, we use a Gaussian copula for the information-theoretic estimation, which is a continuous estimator better-suited for the large and continuous observation space. We did not use a Gaussian copula for the other experiments (logic gate and Animal-AI experiments; see below) because their discrete observation space is incompatible with this method. Therefore, for these experiments we use a discrete estimator and discretize activations via binning.

Although the discretization and selection of sources-target pairing differ, the methods for information decomposition calculations remain the same across all settings. We refer to a source as being either an individual neuron, a dimension of input, or several dimensions of input grouped as a single random variable. Thus, a set of sources refers to either a set of neurons in a layer, a set of input dimensions, or a set of several dimensions of input that are considered as separate sources. In all cases, the target corresponds to the subsequent layer of neurons considered jointly.

#### 5.2.1. Discrete measures

For a given set of sources and target, their corresponding discretized activations are used to compute a probability distribution by counting the number of occurrences of each joint sources-target state and using the plug-in estimator (Panzeri et al. 2007). We use the `dit` library James et al. (2018) to create the distribution and compute the measures of interest.

Although PID proposes the distinction of unique, redundant, and synergistic information, it does not specify a method for computing these measures. Consequently, a number of different formulae have been proposed, although there is currently no general agreement on a particular measure. Thus, for completeness, we compute all measures using two different redundancy functions:  $I_{\min}$  (Williams and Beer 2010) and  $I_{\text{MMI}}$  (Barrett 2015). We find both measures to be consistent with each other across all experimental settings, with  $I_{\text{MMI}}$  yielding slightly higher synergy values. For the purposes of display, we only include  $I_{\text{MMI}}$  measures in the body of the text and refer the reader to Appendix B for all figures replicated using the  $I_{\min}$  redundancy function.

For a set of sources  $\mathbf{X} = \{X_1, X_2, \dots, X_M\}$  and a target  $Y$  with  $N$  possible values, redundancy and synergy for  $I_{\text{MMI}}$  are defined as:

$$R_{\text{MMI}}(\mathbf{X}; Y) = \min_i I(X_i; Y), \quad (1)$$

$$S_{\text{MMI}}(\mathbf{X}; Y) = I(\mathbf{X}; Y) - \max_{\mathbf{A} \in \mathbf{X}: |\mathbf{A}|=M-1} I(\mathbf{A}; Y). \quad (2)$$

Similarly, redundancy and synergy for  $I_{\text{min}}$  are defined as:

$$R_{\text{min}}(\mathbf{X}; Y) = \sum_{j=1}^N p(y_j) \min_i I(X_i; Y = y_j) \quad (3)$$

$$S_{\text{min}}(\mathbf{X}; Y) = I(\mathbf{X}; Y) - I_{\text{max}}(\mathbf{A} \in \mathbf{X} : |\mathbf{A}| = M - 1; Y), \quad (4)$$

where  $I_{\text{max}}$  is defined as:

$$I_{\text{max}}(\mathbf{X}; Y) = \sum_{j=1}^N p(y_j) \max_i I(X_i; Y = y_j). \quad (5)$$

### 5.2.2. Continuous measures

To compute PID in the NeuroGym experiments we use the Gaussian Copula Mutual Information (GCMI) estimator by [Ince et al. \(2017\)](#), which can deal with some of the nonlinearities introduced by the neurons’ activation function. In the Gaussian case,  $I_{\text{MMI}}$  and  $I_{\text{min}}$  are known to be very similar (in fact proven to be identical in some cases ([Barrett 2015](#))), so for simplicity we run all analyses with  $I_{\text{MMI}}$ . We compute the average 2<sup>nd</sup>-order synergy over a random sample of 45 pairs.

### 5.2.3. Full- vs 2<sup>nd</sup>-order decomposition

All calculations are performed over sets of sources of size  $K$ , where  $K$  is either the cardinality of the full set of sources (either the whole layer or the whole input space; referred to as full-order); or  $K = 2$  (2<sup>nd</sup>-order) (Fig. 1c). For 2<sup>nd</sup>-order measures, all calculations are performed between pairs of sources (i.e., over subsets of the source set — a hidden layer or input — of cardinality 2). Thus, a 2<sup>nd</sup>-order value is computed using only 2 elements of a set as sources, rather than the full set. Performing this operation over all possible combinations of pairs and computing the mean gives the average 2<sup>nd</sup>-order measure. When the system grows too large to efficiently compute all possible combinations, this value can additionally be approximated by uniformly sampling pairs of combinations.

We show that full-order and 2<sup>nd</sup>-order measures exhibit similar qualitative behavior in response to dropout and task (Supp. Figs. 6 & 7). However, both our results and prior work ([Tax et al. 2017](#)) suggest that redundancy and synergy are more prevalent at smaller orders, especially for small networks. By computing average 2<sup>nd</sup>-order synergy, we can partially capture how synergistically-biased ([Varley and Hoel 2022](#)) a set of sources is — with higher 2<sup>nd</sup>-order synergy, the PID lattice will have more synergistically-interacting atoms than redundantly-interacting atoms. Given these properties, we use 2<sup>nd</sup>-order measures in the remainder of our experiments.

All of the PID measures shown in the text are 2<sup>nd</sup>-order and are normalized by mutual information. Thus, we are specifically showing the proportion of mutual information occupied by each measure — an increase in (normalized) synergy is in favor of either redundancy or unique information, which must in turn be reduced.

### 5.3. *Statistical analyses*

We perform independent samples t-tests when comparing different models and paired samples t-test when comparing the same models at different points during training. Additionally, we perform a Benjamini-Hochberg False Discovery Rate correction to account for multiple comparisons made in our lesion experiments and when comparing interleaved and sequential protocols across pairs of tasks in NeuroGym.

### 5.4. *Logic gate experiments*

The data used for the COPY and XOR logic gate experiments are generated as a two-dimensional binary input with a binary output. The label of each COPY gate input corresponds to the copying of the first input and the label of the XOR gate input corresponds to the parity of both inputs.

Our models are small feedforward networks, with two layers consisting of ten neurons each. Dropout is only applied during initial training and not during testing or lesioning evaluation. Each model is trained to convergence. We subsequently test and compute various information decomposition measures.

Each activation sampled during testing is discretized using 3 bins in the range of [0,5]. We use 3 bins to ensure a sufficient number of samples in each source-target pair. The range of bins is chosen based on empirical observations of the network activations being heavily concentrated within this range.

For the lesioning experiments, we additionally compute the average pairwise synergy for each neuron. For a particular neuron, this is performed by computing all 2<sup>nd</sup>-order synergy values that include the neuron as one of the sources and calculating the mean.

### 5.5. *Animal-AI experiments*

The experiments conducted in the Animal-AI Environment (Crosby et al. 2020) are performed with proximal policy optimization (PPO) models (Schulman et al. 2017) using `Stable-Baselines3` (Raffin et al. 2021). The actor-critic networks of the models consist of two feedforward layers with ten neurons in each layer, identical to those used in our logic gate experiments.

During training, we evaluate and compute synergy for each model at each task threshold. For all but the last task trained on, the threshold is chosen as the point at which the model successfully reached the maximum reward or the maximum number of steps per task. The last task is trained on for the maximum number of steps per task.

We constrain the observation space to three object-oriented raycasts, each being a one-hot vector indicating the type of object hit by the raycast and its distance normalized by the size of the

arena, and an additional vector relaying information about the agent’s health, velocity, and global position. The raycasts are projected 90 degrees apart — directly in front of the agent, to its left, and to its right. Thus, the agent has full access to the information required to solve the task at the first time step of the episode, preventing the addition of bias that could be introduced by input bits being occluded and the need for integration of input information over time. The raycast observation space occludes the contents of both pits (positive/negative reward), which would otherwise be visible to an agent receiving the full pixel observation space. It also allows for the use of small feedforward networks, rather than larger convolutional layers necessitated by a full pixel input. The small observation space also facilitates our models’ (with small network parameter spaces) ability to solve the given tasks in a complex three-dimensional environment.

We use the same basic design shown in Fig. 3 for all tasks created. In the pit corresponding to the correct output lies an occluded object with a reward of 4 and in the other pit lies an occluded ‘death zone’ with a reward of -1, both of which terminate the episode upon being reached. The agent’s movement is constrained to the platform and one pit and therefore successful completion of the task is contingent on using the information relayed by the bit-representing barriers.

Using this design, we modify the placement of the positive and negative rewards according to the logic gate task being performed. Agents are placed on a short platform to restrict their possible state space and simplify the task, serving as a minimal baseline for solving logic gates in a RL setting. The 3-Bit XOR task explored the effect of integrating an additional source by placing a third input-source barrier in front of the agent (Fig. 3c). Finally, we create three tasks (Distance XOR-10, 20, and 30) by using the basic task design and elongating the platform to lengths of 10, 20, and 30 arena units, increasing the distance between the agent and the reward as a method of adding more difficulty to the logic gate task without the addition of sources (Fig. 3d). Our curriculum tasks consist of the combination of 2-Bit XOR to 3-Bit XOR; and 2-Bit XOR to Distance 10, 20, and 30 XOR.

Synergy is computed in the actor network of the PPO model. Because the observation space exceeds 20 dimensions, synergy cannot be efficiently computed over the entire input using our measures. However, due to the modularity of the raycast input, grouping dimensions based on object-related information is likely to yield more interpretable measurements. Thus, to compute synergy from the input sources to the first linear layer target, we treat each raycast as a single source with the vector dimension for normalized distance being discretized with 3 bins from  $[0,1]$ . Additionally, the global position is also treated as a single source and discretized using 5 bins from  $[0,40]$  (40 being the length of the arena). We compute the average 2<sup>nd</sup>-order synergy between all combinations of source pairs of raycasts (shown in the main text) and source pairs of each raycast and the global position (shown in Supp. Fig. 9), which yield similar values. The average 2<sup>nd</sup>-order synergy is then computed for the rest of the network.

### 5.6. *NeuroGym experiments*

Our RNN models consist of a single recurrent unit with a hidden layer size of ten neurons. During training, decision-making actions are weighted by a favor of 20 in the cross entropy loss compared to the action corresponding to fixation, due to their relative scarcity in the training process. The

observation space is modified to include a binary indicator signifying the task being performed. We test and compute synergy after training for 80,000 total steps (40,000 per task in the sequential protocol; 80,000 total in the interleaved protocol). Each task is trained using supervised learning and we modify tasks with variable episode timing to be of fixed length.

Unlike our other sets of experiments, the tasks in NeuroGym are stochastic and have a much larger task-space. To ensure sufficient activation sampling, we test models on a total of 100 trials. We use both the input and hidden layer at each time step as sources and the hidden layer at the following time step as the target.

## APPENDIX

## A. ADDITIONAL FIGURES

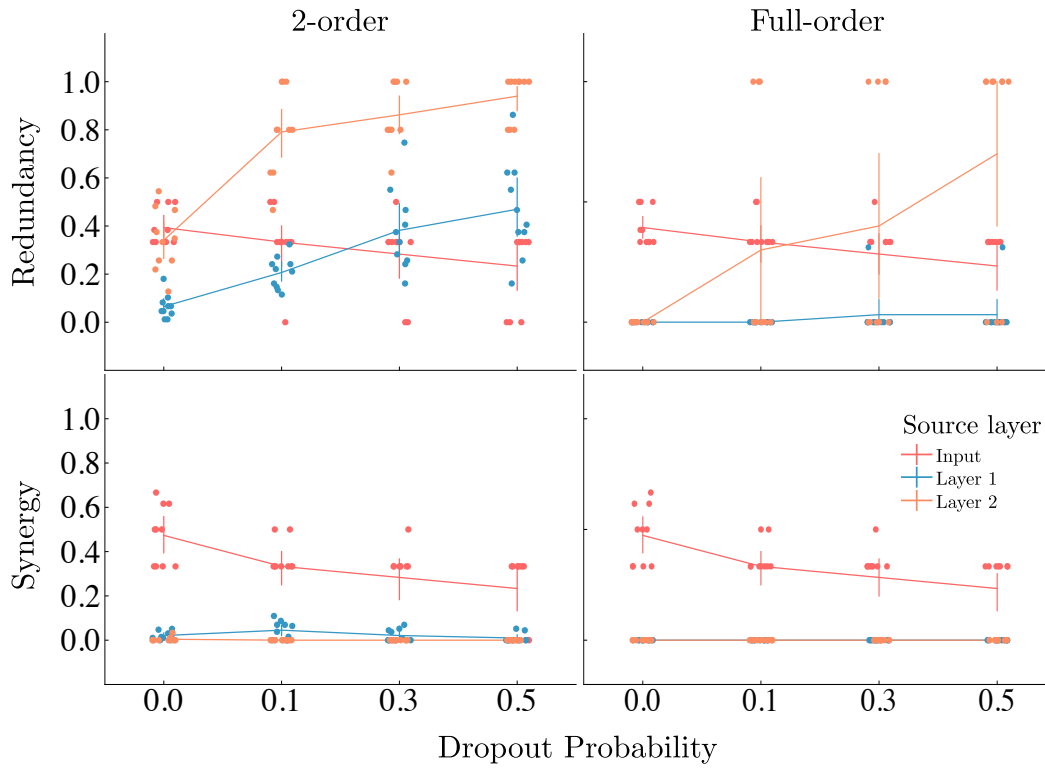


Figure 6. Full-order and 2<sup>nd</sup>-order measures exhibit similar behavior for the COPY task.



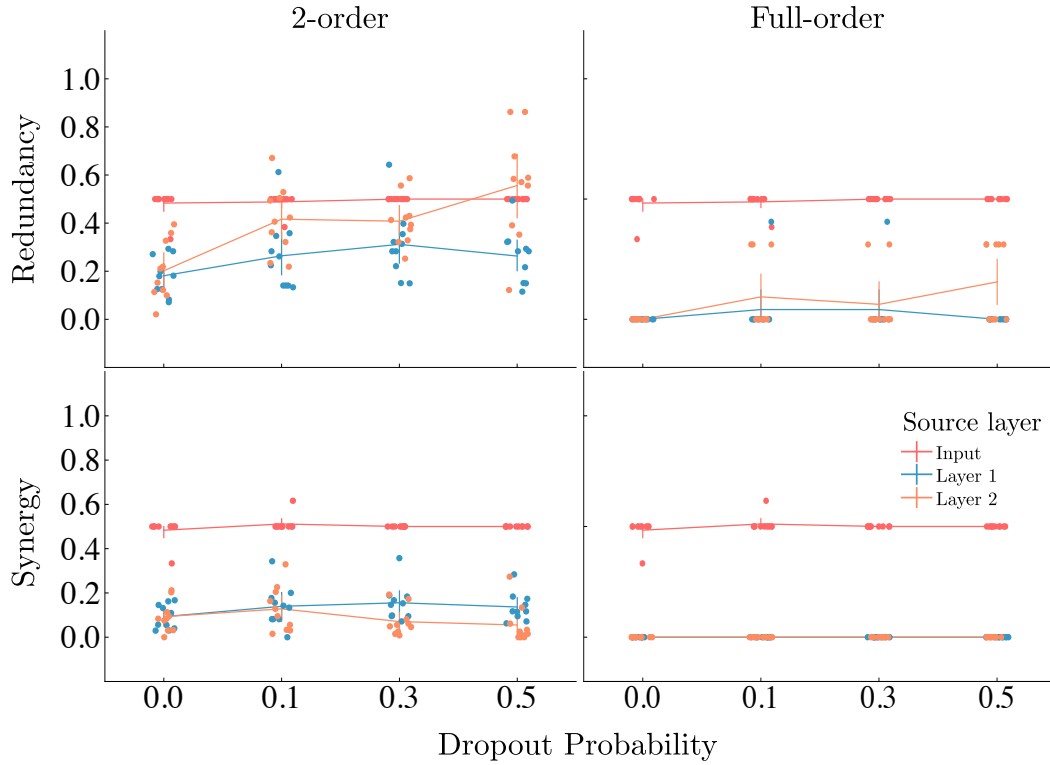


Figure 7. Full-order and 2<sup>nd</sup>-order measures exhibit similar behavior for the XOR task.

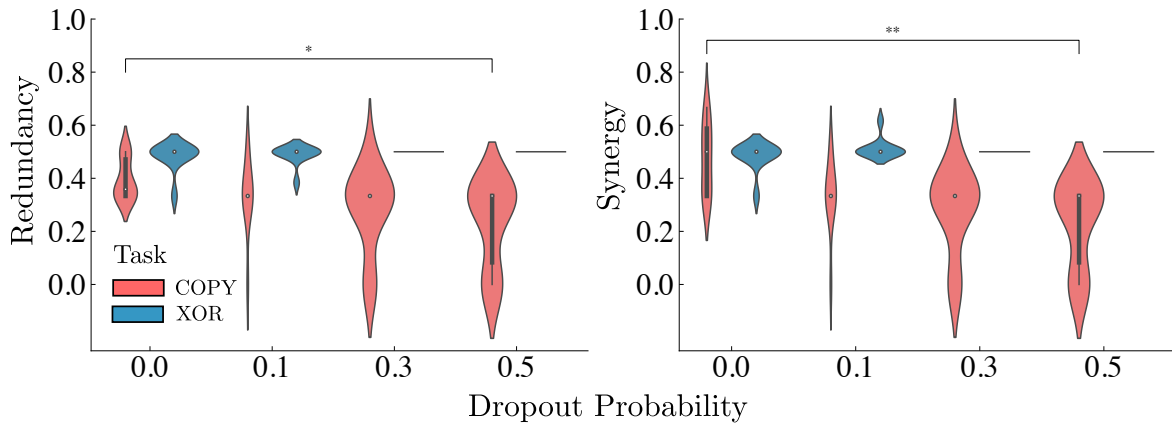


Figure 8. Dropout removes irrelevant redundant and synergistic information about the input in the COPY task, but not the XOR task.

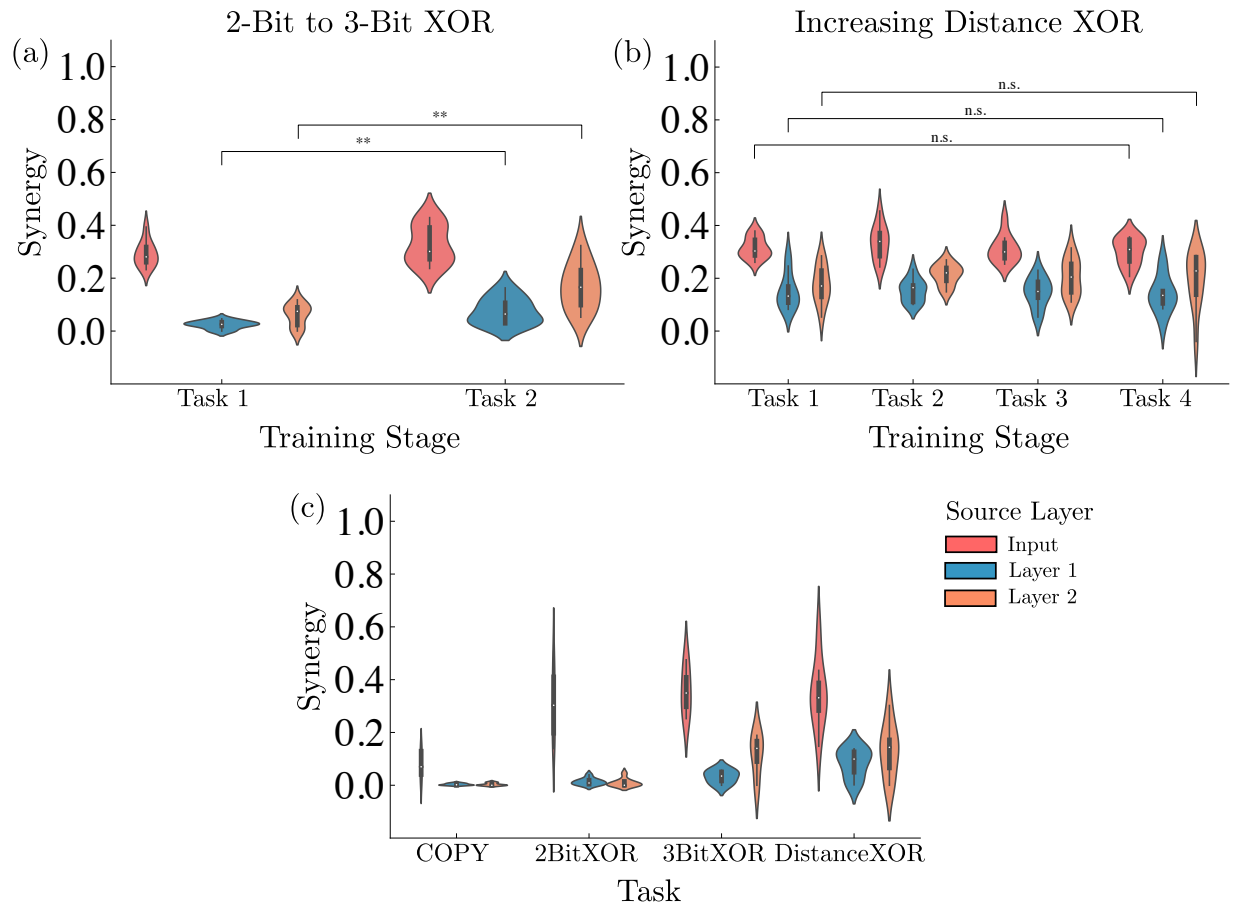


Figure 9. Relation of compositional tasks and synergy in Animal-AI using pairwise raycast-position sources.

B. FIGURES REPLICATED USING  $I_{\min}$  REDUNDANCY FUNCTION

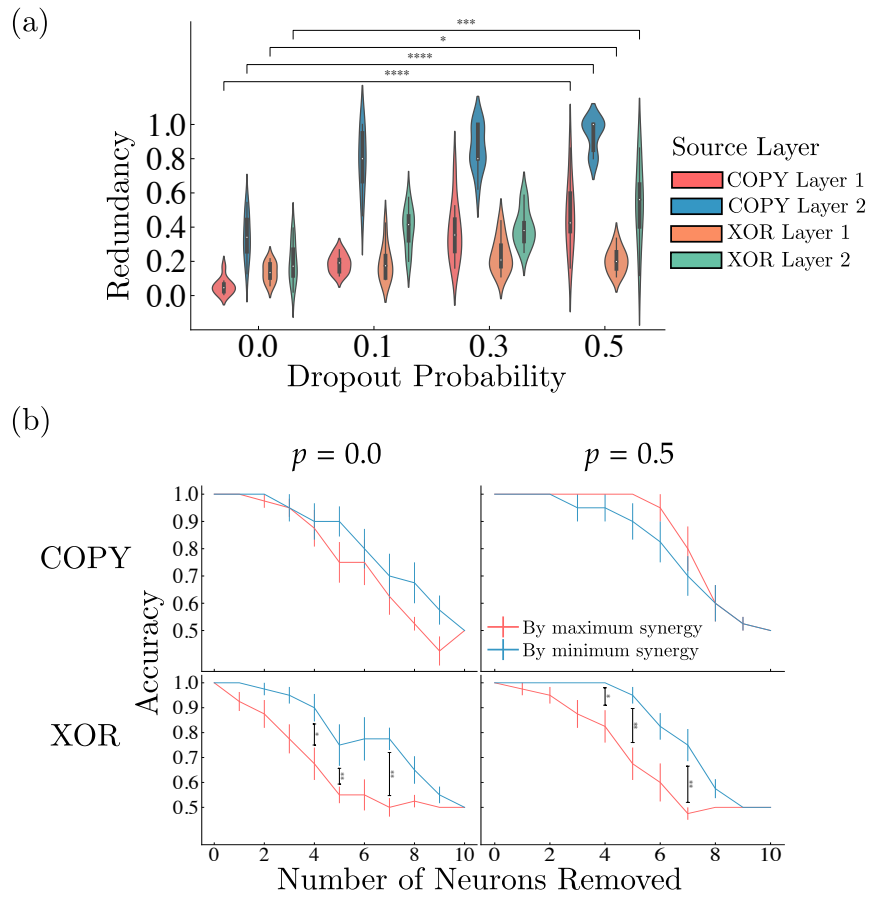
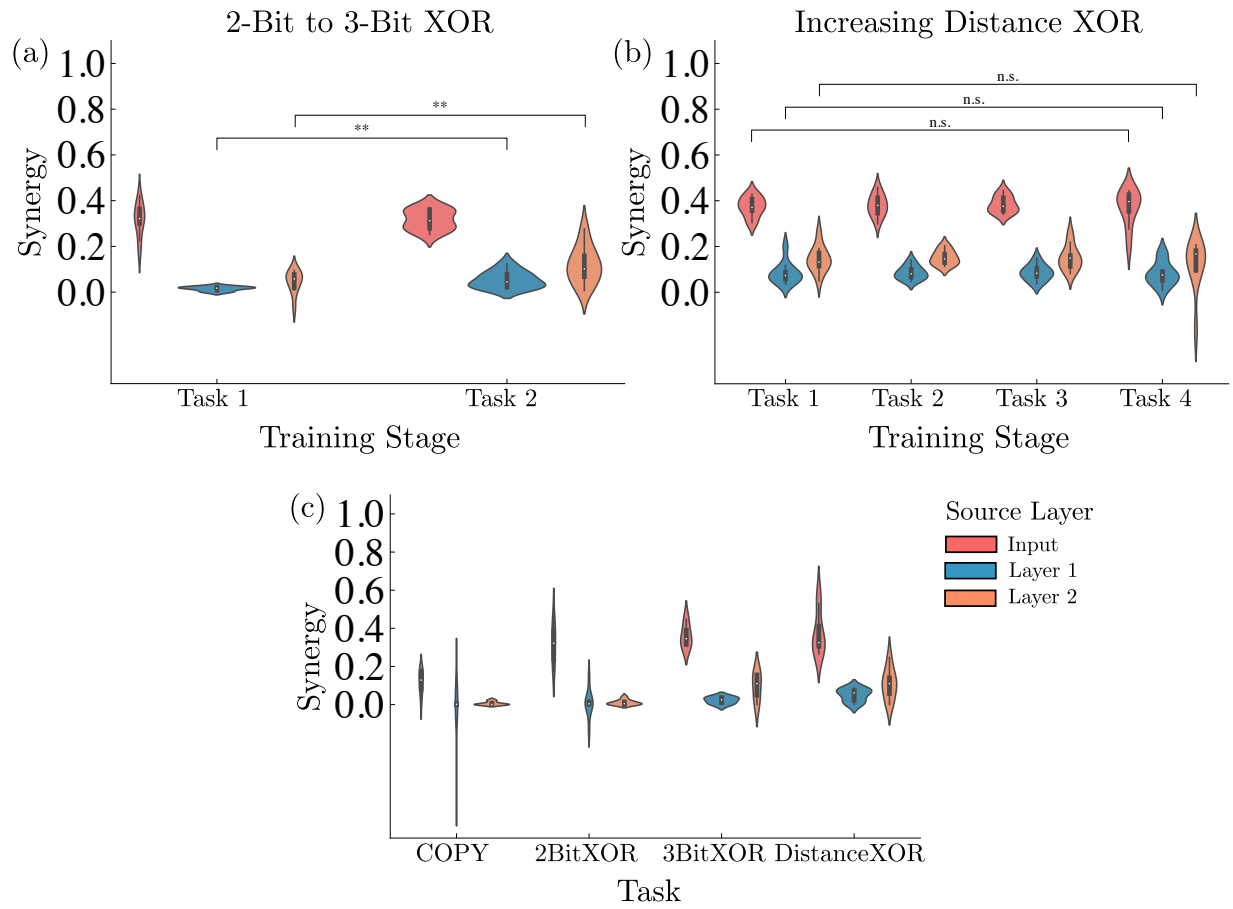


Figure 10. Effects of lesions and dropout on network information profiles.



**Figure 11. Relation of compositional tasks and synergy in Animal-AI using pairwise raycast sources.**

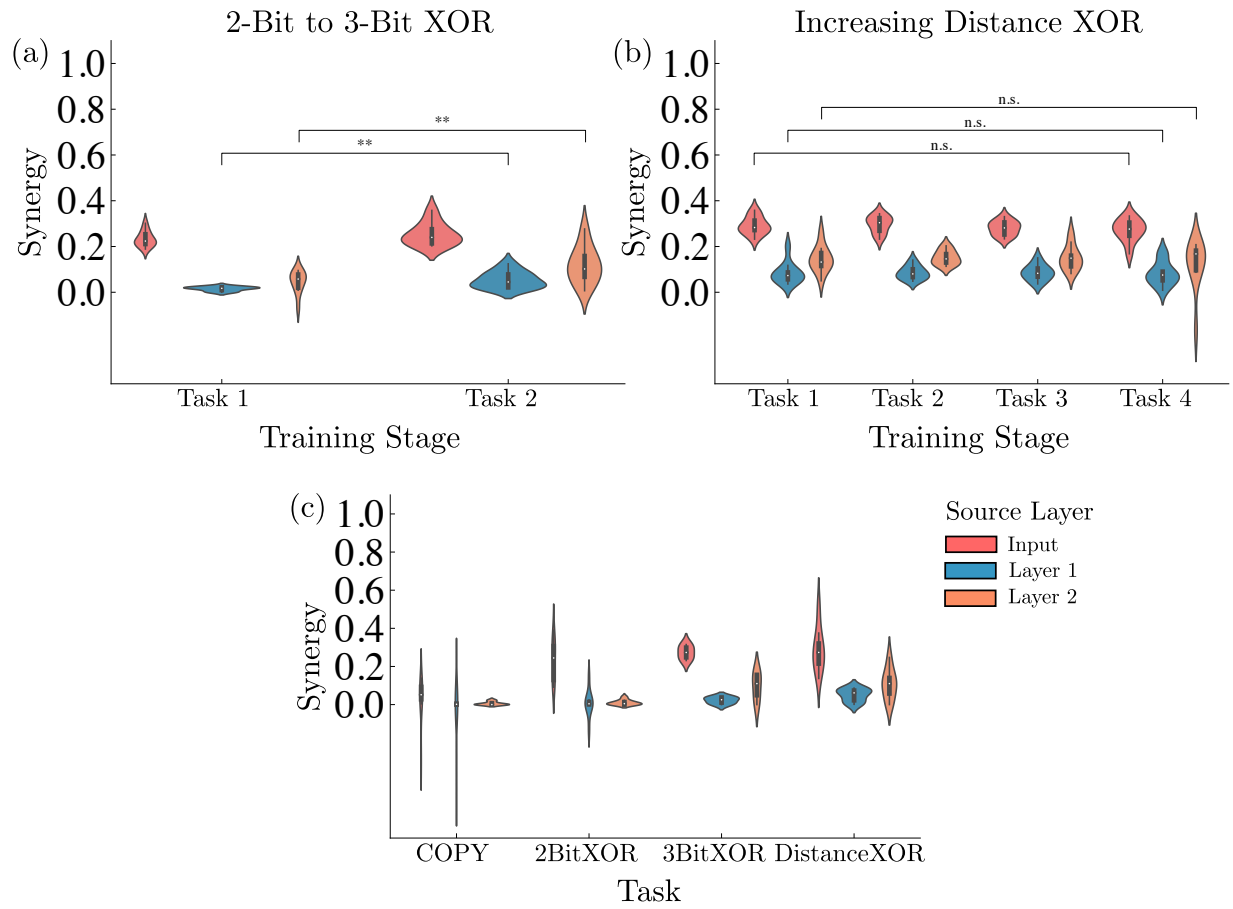


Figure 12. Relation of compositional tasks and synergy in Animal-AI using pairwise raycast-position sources.

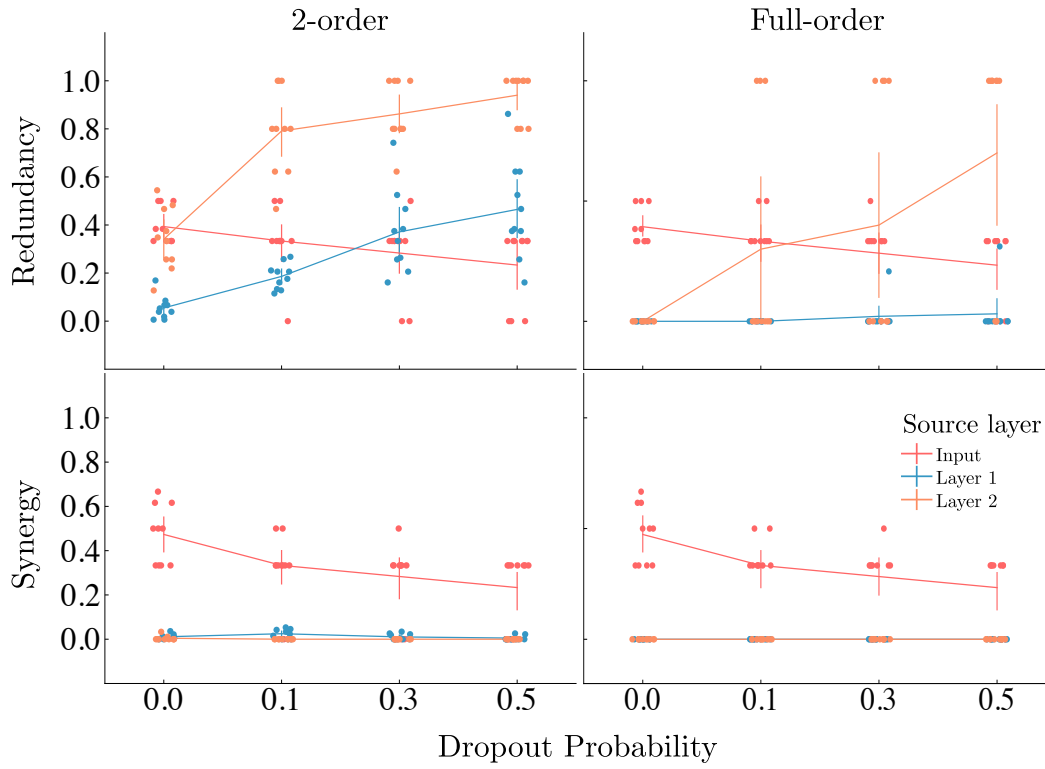


Figure 13. Full-order and 2<sup>nd</sup>-order measures exhibit similar behavior for the COPY task.

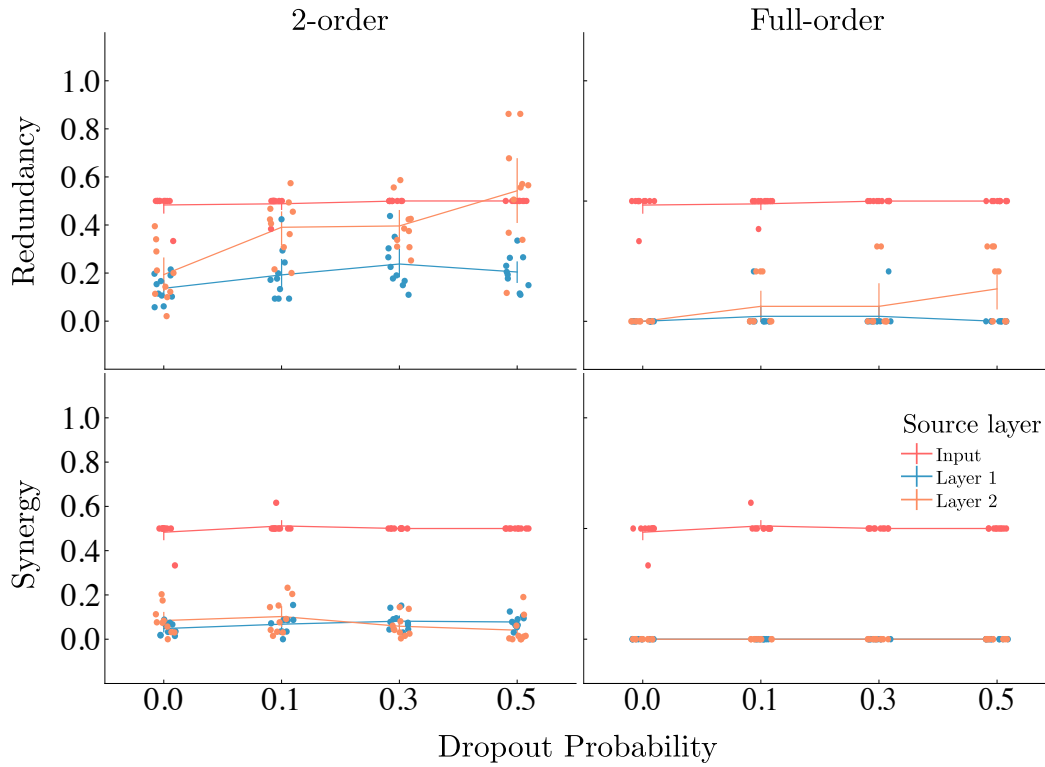


Figure 14. Full-order and 2<sup>nd</sup>-order measures exhibit similar behavior for the XOR task.

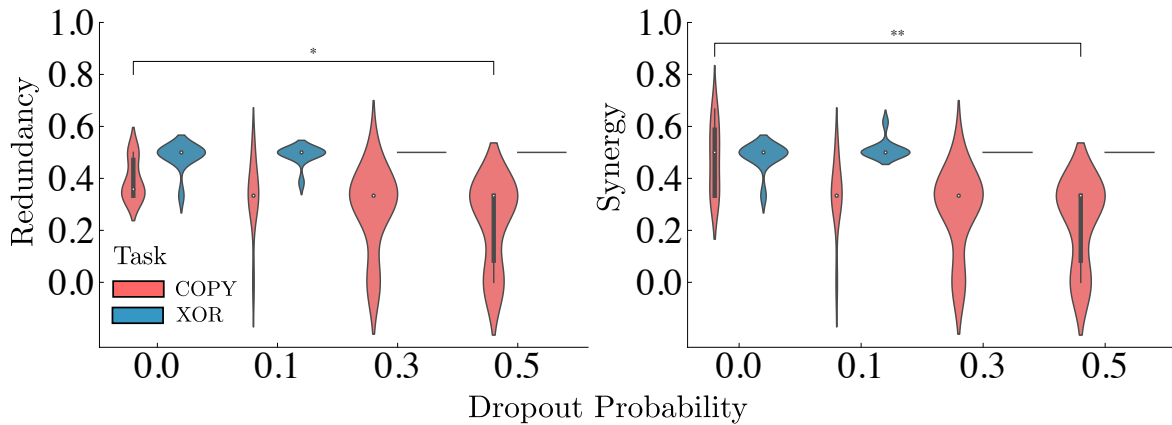


Figure 15. Dropout removes irrelevant redundant and synergistic information about the input in the COPY task, but not the XOR task.