

Systematic Perturbation of an Artificial Neural Network:

A Step Towards Quantifying Causal Contributions in The Brain

Kayson Fakhar^{1*} and Claus C. Hilgetag^{1,2}

1. Institute of Computational Neuroscience, University Medical Center Eppendorf, Hamburg University, Germany.

2. Department of Health Sciences, Boston University, Boston, MA, USA.

*k.fakhar@uke.de

Abstract

Lesion inference analysis is a fundamental approach for characterizing the causal contributions of neural elements to brain function. Historically, it has helped to localize specialized functions in the brain after brain damage, and it has gained new prominence through the arrival of modern optogenetic perturbation techniques that allow probing the functional contributions of neural circuit elements at unprecedented levels of detail.

While inferences drawn from brain lesions are conceptually powerful, they face methodological difficulties due to the brain's complexity. Particularly, they are challenged to disentangle the functional contributions of individual neural elements because many elements may contribute to a particular function, and these elements may be interacting anatomically as well as functionally. Therefore, studies of real-world data, as in clinical lesion studies, are not suitable for establishing the reliability of lesion approaches due to an unknown, potentially complex ground truth. Instead, ground truth studies of well-characterized artificial systems are required.

Here, we systematically and exhaustively lesioned a small Artificial Neural Network (ANN) playing a classic arcade game. We determined the functional contributions of all nodes and links, contrasting results from single-element perturbations and perturbing multiple elements simultaneously. Moreover, we computed pairwise causal functional interactions between the network elements, and looked deeper into the system's inner workings, proposing a mechanistic explanation for the effects of lesions.

We found that not every perturbation necessarily reveals causation, as lesioning elements, one at a time, produced biased results. By contrast, multi-site lesion analysis captured crucial details that were missed by single-site lesions. We conclude that even small and seemingly simple ANNs show surprising complexity that needs to be understood for deriving a causal picture of the system. In the context of rapidly evolving multivariate brain-mapping approaches and inference methods, we advocate using *in-silico* experiments and ground-truth models to verify fundamental assumptions, technical limitations, and the scope of possible interpretations of these methods.

Author summary

The motto “*No causation without manipulation*” is canonical to scientific endeavors. In particular, neuroscience seeks to find which brain elements are causally involved in cognition and behavior of interest by perturbing them. However, due to complex interactions among those elements, this goal has remained challenging.

In this paper, we used an Artificial Neural Network as a ground-truth model to compare the inferential capacities of lesioning the system one element at a time against sampling from the set of all possible combinations of lesions.

We argue for employing more exhaustive perturbation regimes since, as we show, lesioning one element at a time provides misleading results. We further advocate using simulated experiments and ground-truth models to verify the assumptions and limitations of brain-mapping methods.

1 Introduction

2 One of the most challenging goals of neuroscience is to identify neural elements – brain regions,
3 populations, neuronal circuits, and large-scale networks – that pivot cognition and behavior[1].
4 During the past two decades, brain mapping flourished with the help of neuroimaging techniques
5 that associate elements and functions. Arguably though, the first method of mapping brain
6 function, i.e., by studying lesions, yet has an authoritative role in establishing causation since it
7 indicates the *necessity* of the element for a given function[2,3]. With this inferential capacity, though,
8 comes practical and methodological difficulties that might deliver deceiving results [4,5]. Crucially,
9 since the ground-truth causal processes in the brain are unknown, the limitations of how functional
10 contributions are mapped to interacting neural elements are not fully resolved, and thus
11 conventional lesion-based methods are left with unverified assumptions and unexplored
12 alternatives[5].

13 From a practical point of view, the scale of available human lesion datasets is nowhere on a par
14 with those used in and produced by correlative approaches. This is in particular problematic since,
15 as it is shown, even by focusing on single local lesions, mass-univariate lesion analysis provides
16 systematically biased maps while multivariate approaches require a considerable amount of data to
17 remedy the problem [2,6]. Additionally, with invasive approaches and in animal models, the sheer
18 number of elements in the brain makes it practically impossible to lesion all of them exhaustively
19 in all but very small nervous systems[7,8].

20 Practical issues aside, cognitive functions emerge from interactions of distributed neural elements
21 that make it challenging to isolate the functional contributions of individual units[5,9] while the
22 established approach assumes to disassemble such coalitions by removing individual elements and
23 assigning the resulting behavioral change as the elements’ contribution[3,10]. Historical cases of
24 lesion inference after brain damage, in patients such as Phineas Gage and Henry Molaison (‘HM’)
25 [11], as well as modern cutting-edge experimental tools employing opto- and chemogenetics that
26 temporarily perturb the brain with astonishing spatiotemporal precision [12,13], mostly follow the
27 same “Single-element Perturbation Analysis (SPA)” framework. It is important to note that a SPA
28 study might have a multivariate approach by incorporating many variables, e.g., lesion volume, but

29 one neural element is perturbed -or fed into a statistical model- at a time, whether the element is
30 single neurons, a local circuit, or a brain region[5,14]. Put differently, neural elements produce
31 behavior as spatially distributed, *interacting* coalitions[15–17] while the established methods mainly
32 map the observed effects on local processes. Consequently, the SPA framework might overlook
33 the subsequent effects that local lesions might have on the system as a whole[18]. Paradoxical
34 lesion effects and, in particular, the “Sprague Effect” are intriguing phenomena to illustrate
35 potential issues with this approach[19,20]. The Sprague effect describes a scenario in which
36 disruptions in behavior caused by a first lesion revert to normal after a second lesion[20,21]. In
37 other words, lesioning region i disrupts the behavior, providing apparently compelling evidence
38 for its “necessity” for the behavior, while a subsequent lesion to another region j restores the
39 behavior showing the redundancy or degeneracy of the contribution of i .

40 Different hypotheses have attempted to explain this unexpected result based on the inhibitory
41 relationship between competing regions[22–24] or neuronal plasticity and the increased excitatory-
42 to-inhibitory synaptic balance of the circuit[25]. Essentially, the Sprague effect points towards a
43 more complex causal relationship in the brain rather than a single neural element-to-single function
44 relationship, indicating how misleading it can be to assign functions to neural elements relying on
45 individual lesions[18].

46 To further emphasize on this point, Jonas and Kording performed an exhaustive SPA of every
47 transistor in a microprocessor to see if it reveals a meaningful causal picture of a system that we
48 have confound-free access to, virtually, every computational unit of it[26]. They found a subset of
49 transistors that perturbed, would disrupt the function of the microprocessor; however, they
50 declared the results “*grossly misleading*” since “*The transistors are not specific to any one behavior [...] but*
51 *rather implement simple functions*” [26]. Their results suggest that even by perturbing every relevant
52 unit of a system, one at a time, we are still far from a coherent causal understanding of what is
53 doing what and indeed prone to miss-attribute individual elements to a behavior that is emerged
54 from complex interactions of many units.

55 In this work, we use an alternative approach known as “Multi-perturbation Shapley value Analysis
56 (MSA)” that, in contrast to SPA, derives causal contributions of elements from permuting all
57 combinations of multi-element lesions [27,28]. MSA is based on Shapley value, a game-theoretical
58 metric that is used for fair distribution of costs, gains, or resources among players of a cooperative
59 game[29]. In the context of neuroscience, players are arbitrarily defined neural elements that are
60 ranked according to their contributions to an arbitrary quantified behavior or cognitive
61 function[30,31].

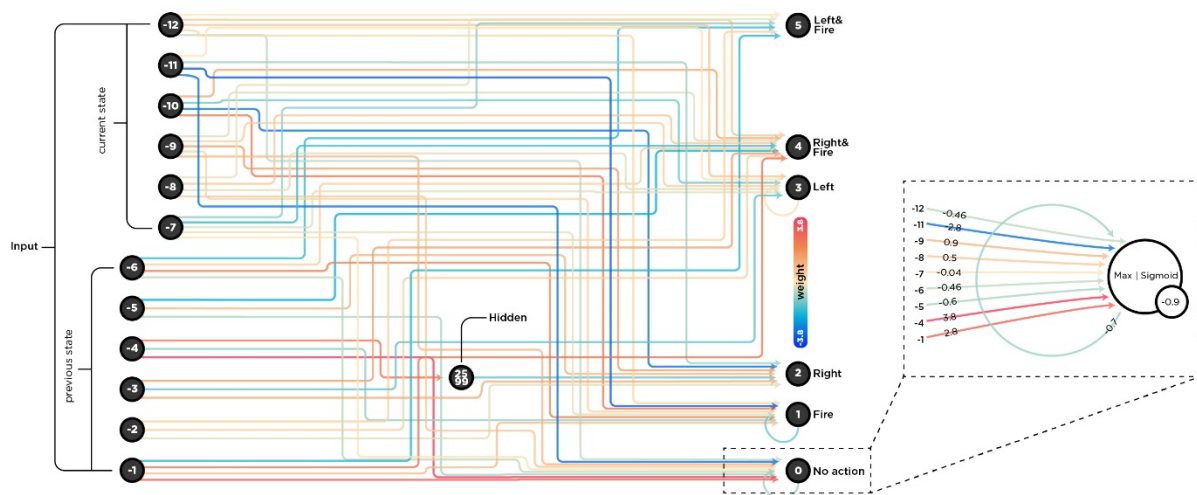
62 Inspired by the provocative findings of Jonas and Kording and further investigating the inferential
63 capabilities of SPA and MSA frameworks, we decided to use a ground-truth model and
64 systematically perturb its components. Therefore, we perturbed all neurons and connections of a
65 compact ANN, both one element at a time, that is, through SPA, or many combinations of
66 elements, that is, by MSA. We used an ANN instead of a microprocessor to capture the whole
67 spectrum of the behavioral performance instead of a binary state of disturbed versus functional
68 performance. Moreover, to train the network we specifically used an evolutionary algorithm
69 focused on the network’s topology to avoid handcrafting and potentially biasing its organization

70 and to see if an *in-silico* evolutionary process produces topologies with the functional motif of
71 inhibition between rivalrous elements.

72 Briefly, we found that not every perturbation necessarily revealed causation. Although data from
73 both lesioning regimes showed similarities, SPA missed a few of the key contributing elements and
74 miss-attributed their causal ranks. Therefore, it provided biased contributions for individual
75 elements, while the MSA captured these nuances more accurately. To further quantify the complex
76 interaction of elements within the system, we used an extension of MSA, here called Pairwise
77 Causal Interaction Analysis (PCIA)[27,28], and found a handful of pairs in which lesioning one
78 unit while the other is perturbed restored the disrupted behavior. Finally, we delved deeper into
79 the inner mechanisms of the network to identify why MSA ranked the units in the given way and
80 what these units do that SPA was insensitive to. We discuss the findings, the limitations of the
81 current approach and outline potential future questions to pursue.

82 Results

83 Our *in-silico* experimental setup was the ATARI arcade game Space Invaders, in which the agent,
84 located at the bottom of the environment, needs to defend itself from aliens descending from the
85 upper part of the screen using laser canons. The main objectives are to stay alive by avoiding alien
86 laser shots and scoring as many points as possible by eliminating aliens. On average, a human
87 subject obtains a score of 1652, and an algorithm that randomly selects actions can reach a score
88 of 148[32]. Other classic algorithms, such as an earlier implementation of a Deep Q-learning
89 Network (DQN), State–Action–Reward–State–Action (SARSA), and a refined DQN, reach 581,
90 271, and 1976, respectively[32,33].



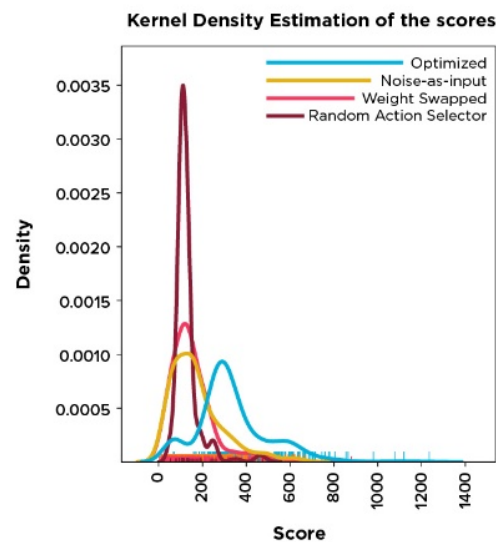
92 **Fig.1: The complete wiring diagram of the evolved ANN.** At each time point, the network received a compressed
93 version of the game-state as a vector of 12 features, six features per frame. It then chooses an action from six available
94 actions (output nodes). Due to its importance, which was revealed later in the analysis, we plotted node 0 separately
95 with more information on the right part of the figure. The aggregation function for this node is max, the activation
96 function is a sigmoid function, and the bias is -0.9. Note that these functions are different for each node (see section
97 *Evolutionary optimization*).

98 Instead of training deep networks using backpropagation in a predefined architecture, we evolved
99 a compact network using a Neural Architecture Search (NAS) algorithm called Neuro Evolution
100 of Augmenting Topologies (NEAT)[34]. Briefly, NEAT uses evolutionary principles such as cross-
101 over of genes (network topologies), speciation (preserving novelty), and incremental
102 complexification to find the “fittest” topology. This means the network’s architecture and
103 connectivity are not handcrafted, nor does the algorithm solely optimize connection weights.
104 Instead, the fittest network is evolved with respect to the environmental constraints, in this case,
105 to have the highest score by adjusting its topology according to a set of given limitations, for
106 instance, low probability of adding connections versus higher probability of removing them, see
107 section *Evolutionary optimization*.

108 In addition to these sets of hyperparameters, to further enforce a compact architecture, we
109 compressed the game frames using a deep auto-encoder and fed our network with two feature
110 vectors (12 features in total, neurons labeled with negative numbers in Fig.1) at each time point.
111 We fed two frames instead of one due to the non-Markovian structure of the game in which only
112 knowing the current position of laser beams does not provide enough information about the
113 beams’ directions.

114

115



116

Fig.2: Distribution of performances. Optimized network is the evolved network, which reached a good-enough performance. Noise-as-input is the same network that receives random values drawn from a uniform distribution [0, 1] as input instead of receiving game-states. Weight swapped network receives the game-states while the connection weights are shuffled. Finally, Random action selector is an algorithm that selects a random action, at each timepoint, regardless of the game-states.

117

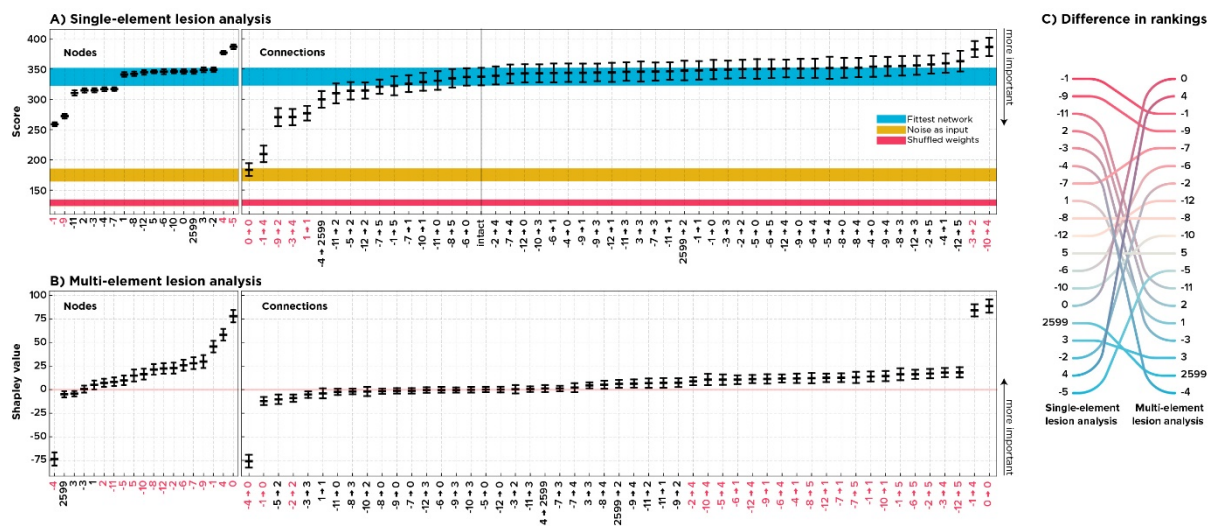
118 On average, our evolved network obtained a score of 337 that is significantly higher than a random
119 agent with a score of 148 (Mann-Whitney U statistics; MWUs = 39542, p-value <0.001, Fig.2). In
120 addition to the random agent and to ensure that the score is not higher merely because of innate
121 topological privileges, we compared the performance with the performance of two control
122 networks. In one, we kept the network as is and made it blind by feeding noise instead of features,

123 and in the other, while the network was receiving game-states, we shuffled the connection weights.
 124 Both control networks obtained substantially lower scores, i.e., from 175 (MWUs = 50919, p-value
 125 <0.001) to 129 (MWUs = 31157, p-value <0.001), respectively. Altogether, these results show that
 126 our compact network did learn the task to some degree and could reach a good enough score
 127 (Fig.2) that formed the basis of the subsequent perturbation analyses.

128

129 Perturbing all elements, one at a time

130 After evolving the network, we intervened to see if perturbing elements could reveal their causal
 131 importance for the behavior. We first silenced neurons one at a time and ran the simulation with
 132 the lesioned network. Conventionally, we searched for neurons, which, when lesioned, resulted in
 133 a considerably deteriorated performance, indicating their “necessity” for the behavior. As (Fig.3A)
 134 shows, lesioning either of two input neurons -1 and -9 had such a disruptive impact, while
 135 individually perturbing most other neurons had a negligible effect on the performance.
 136 Interestingly, lesions of two neurons, 4 and -5, improved the performance, suggesting their
 137 hindering role during normal functioning.



138

139 **Fig.3: Single-element Perturbation Analysis versus multi-perturbation Shapley value Analysis of the ANN.**
 140 This figure shows the result and the rank difference derived from a SPA (A; 512 samples per element) versus an MSA
 141 (B; 1,000 samples per element). On the left side, the nodes, and in the middle, the connections are sorted according
 142 to their inferred average contributions. For SPA, the lowest value means the most influential while the other way
 143 around applies to Shapley values, with the highest value means the most critical. Error bars are %95 Confidence
 144 Interval (CI; bootstrapped 10,000 times). The blue, yellow, and red strips show the %95 CI of the labeled control
 145 networks. Red labels on the x-axis show significant elements (alpha inflation is corrected using Bonferroni correction,
 146 see *Statistical inference in Materials and methods*). On the right-hand side, the node rankings are compared.

147

148 To account for the unique consequences of white matter lesions, also known as disconnection
 149 syndromes[11,35], we performed the same lesioning scheme on all connections. We wanted to see
 150 if severing individual connections among neurons instead of silencing a whole neuron with all its
 151 connections can further localize functional contributions in our ANN. For example, are neurons
 152 -1 and -9 essential elements for the behavior of the ANN, or are there connections of these

153 neurons such that the neurons only appear to be critical in the sense that lesioning them perturbed
154 those connections as well? Based on the single-node removal experiment results, we expected to
155 see either no specific connections to be causally crucial, showing that neurons are the actual units
156 of causation or a major disruption in behavior following lesions to the outgoing connections from
157 neurons -1 and -9.

158 Surprisingly, a loop from neuron 0 to itself (self-loop) appeared to be the most critical element
159 (Fig.3B). This observation indicates that, although SPA of all elements resulted in some degree of
160 coherence by first capturing neurons -1 and -9 as major players and then tracking their importance
161 to connections ($-1 \rightarrow 4$) and ($-9 \rightarrow 2$), another key aspect is downplayed. If neurons were the
162 essential elements, no single connection lesion would have had such devastating effects, or the
163 critical connections would be associated with the critical neurons. However, lesioning single
164 connections did impact the performance considerably. The critical connection is not a connection
165 from or to the most important neurons but a self-loop of a neuron that itself had a near-zero causal
166 contribution.

167 To summarize our point, results from the SPA of each neuron indicated that neuron 0 has little
168 impact on the performance while SPA of the self-loop ($0 \rightarrow 0$) disrupted the behavior the most.
169 Note that throughout the lesioning experiments, the network was fixed, and its architecture
170 determined its behavior. Therefore, we suspected a more complex interaction among neuron 0's
171 connections such that lesioning ($0 \rightarrow 0$), while those key connections were intact, disrupted the
172 behavior, and lesioning ($0 \rightarrow 0$) alongside them had no adverse effect. We suspect those
173 connections to be among other connections of neuron 0 since removing the node virtually
174 perturbed all its connections, which ended in no disruption in the behavior. Put simply, lesioning
175 connection ($0 \rightarrow 0$) alone caused the most damage while lesioning neuron 0 with all its 11
176 connections – including ($0 \rightarrow 0$) – did not show any behavioral impairment. In the next section,
177 we describe the MSA algorithm and elaborate on its results.

178

179 **Multi-perturbation Shapley value Analysis of all elements**

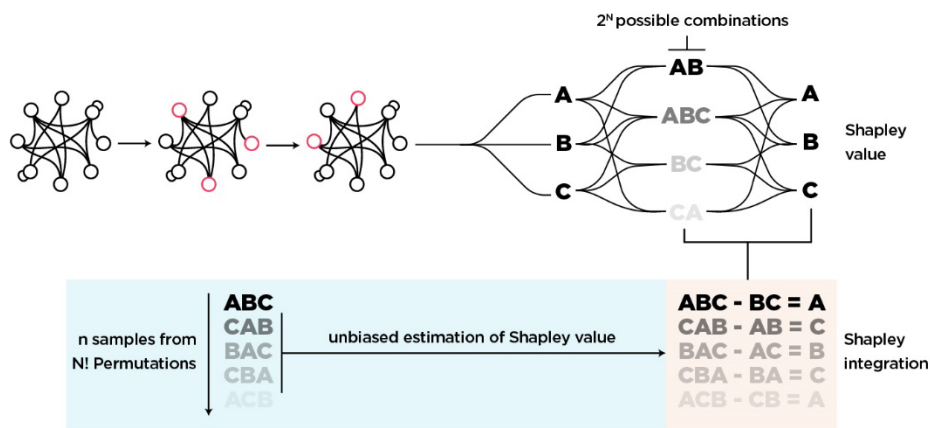
180 We next adopted a multi-element lesioning approach to perturb all neurons and all connections.
181 We used a rigorous, game-theoretical metric based on the Shapley value (γ) called MSA[27]. To
182 elaborate, the Shapley value accounts for the “*worth*” of an element in the grand coalition of all
183 elements, forming the entire system, in terms of the element's contribution to the system's
184 outcome, given its unique added contribution to all possible combinations of coalitions[27–29].
185 So far, Shapley values are the only values that mathematically proven to satisfy the following
186 axioms[29]:

- 187 1. **Symmetry:** If two elements are functionally interchangeable, then their contributions will
188 not differ by their labels.
- 189 2. **Null player property:** If an element does not contribute to the given function, its Shapley
190 value is zero.
- 191 3. **Additivity:** Summing the contributions of all elements results in the performance of the
192 grand coalition.

193 As with the SPA framework, this approach aims to find elements that, when lesioned, most
 194 strongly impair the behavior. In this case, these elements have the highest Shapley value that is
 195 derived from permuting all combinations of multi-site lesions (many elements are lesioned at each
 196 time) such that the target element is once included in the lesioned coalition and once excluded
 197 from it. In other words, for each permutation, a set of elements are lesioned, the performance is
 198 quantified, the target element is then lesioned alongside the other elements in the coalition, and
 199 the performance is quantified again. The difference between these two conditions, both negative
 200 and positive, is what lesioning an element contributes to that specific group of lesioned elements
 201 (Fig.4). Note that the subsets have arbitrarily different sizes, which means the analysis is reduced
 202 to SPA if the coalition contains only one element, i.e., the target element, and is expanded to the
 203 whole network if the coalition contains all elements. Therefore, while focusing on the importance
 204 of one element, MSA incorporates the multivariate influences of lesioning other elements.
 205 Averaging over these contributions will then be the Shapley value of the target element, indicating
 206 its marginal causal contribution to the system's performance.

207 However, having all possible combinations of subsets explored can be computationally prohibitive
 208 in large sets. Therefore, we used an unbiased estimator of the Shapley value that samples coalitions
 209 from the space of 2^N possible combinations, where N is the number of all elements (see [27] for
 210 detailed information).

211



212

213 **Fig.4: Visual depiction of MSA algorithm.** Since there are 2^N possible combinations of coalitions, an analytical
 214 solution for the Shapley value is computationally prohibitive. Therefore, we sampled 1,000 random permutations from
 215 all $N!$ possible orderings and used those to dictate which coalitions to perturb. One sample of Shapley value for any
 216 element is then its contribution to one permutation, simply by calculating the score difference of the coalition with
 217 the element (e.g., {A, B, C}) and the score of the same coalition without the targeted element (i.e., {A, B} to isolate
 218 C). Note that permutations are order-invariant, which means the performance of coalitions {A, B, C} = {C, B, A}.

219 As mentioned, the Shapley value is additive and thus has an intuitive interpretation in which the
 220 highest possible Shapley value is the grand coalition's worth. For example, for our network, the
 221 Shapley value of the overall coalition is 337. This means an element with a Shapley value of 80
 222 accounts for a fraction of 23% of the network's performance. A negative Shapley value follows
 223 the same line of interpretation, that is, an element with a Shapley value of -80 on average prevents
 224 the network from an additional 23% increase in performance.

225 As depicted in (Fig.3B), MSA shows many noncritical nodes and connections, just as the single-
226 site lesion analysis did. Importantly, according to the MSA, neuron 0 is the most influential,
227 followed by many less critical nodes. Interestingly, neuron -4 is assigned a negative Shapley value,
228 indicating its proportionally large and inhibiting contribution to the system. This contradicts the
229 result obtained from SPA that pointed to -5 and 4 to have such an influence (Fig.3C).

230 As with SPA, we dissected nodes to their connections but this time using MSA to test if we can
231 further track the critical neurons' causal influence down to their connections. Again, we expected
232 to see either lesioning of no single connection to have drastic effects, indicating a distributed
233 regime of processing in which no lower-level unit is as critical, or to find that there are critical
234 connections, and they correspond to the influential nodes since lesioning a node here is the same
235 as lesioning all its connections.

236 MSA tracked the importance of -4 to a single connection from -4 to 0, and the same
237 correspondence applies to the elements with the highest Shapley value. The causal contribution of
238 neuron 0, for example, can be attributed to its connection ($0 \rightarrow 0$) since besides ($0 \rightarrow 0$) and (-4
239 $\rightarrow 0$), other connections of this neuron have negligible contributions (Fig.3B). As a sanity check,
240 we performed the same procedure on the blinded network. Here we expected no element to
241 contribute to the network's overall performance since, on average, the network had the same
242 baseline score. As shown in (Supplementary Figures 1), this is indeed the case.

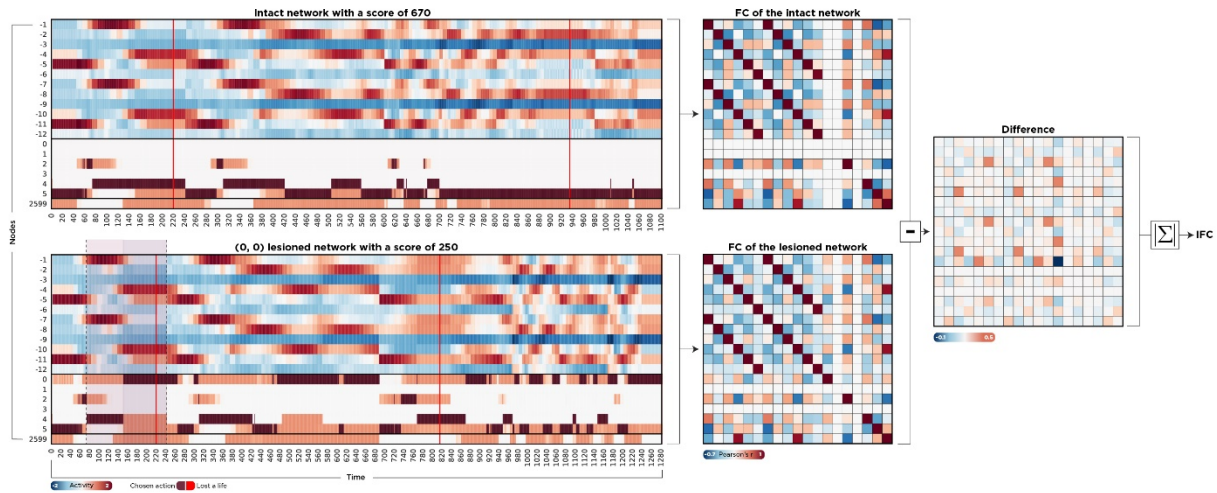
243 The most crucial difference between SPA and MSA was how they ranked connections ($0 \rightarrow 0$)
244 and ($-4 \rightarrow 0$). Remember, even data from the SPA showed ($0 \rightarrow 0$) as the most critical connection.
245 The missing piece was another link to neuron 0 that we suspected to have a Sprague effect-
246 inducing interaction with the self-loop ($0 \rightarrow 0$) and the reason was that by perturbing all 11
247 connections, including ($0 \rightarrow 0$), we had no adverse effect. MSA attributed a negative Shapley value
248 to the connection ($-4 \rightarrow 0$), while SPA assigned minor importance to this connection. This
249 discrepancy aligns with the Sprague effect's essence since at least two elements are required to be
250 lesioned for such a phenomenon to emerge.

251 Altogether, MSA and SPA found key elements to be a small and localized set. MSA dissociated
252 these and assigned the negative contribution to neuron -4 while SPA missed it. While SPA
253 excluded neuron 0, MSA ranked it as the most critical neuron and further dissected this importance
254 to the self-loop. It then showed that the incoming connection from -4 is the possible answer to
255 why lesioning neuron 0 has a near-zero impact.

256 **Impact of lesioning on functional connectivity**

257 In addition to their direct impact on the behavior of a system, lesions may also disrupt functional
258 connectivity (FC), and different features of the impact on FC are associated with behavioral
259 performance. Thus, FC forms a bridge, or 'intermediate phenotype' from structure to function
260 and behavior [35–38]. It was shown that lesions of critical brain regions in terms of FC, such as
261 hubs, have a greater impact on the dynamics of the whole brain[37]. To explore this aspect in our
262 *in-silico* model, we first calculated the FC of the intact network using Pearson's correlation. We
263 then employed a SPA framework for all units, that is, nodes and connections. To quantify the
264 impact of lesioning individual elements on global FC, we calculated the element-wise differences
265 between intact and lesioned FC matrices. The absolute sum of the resulted difference matrix was

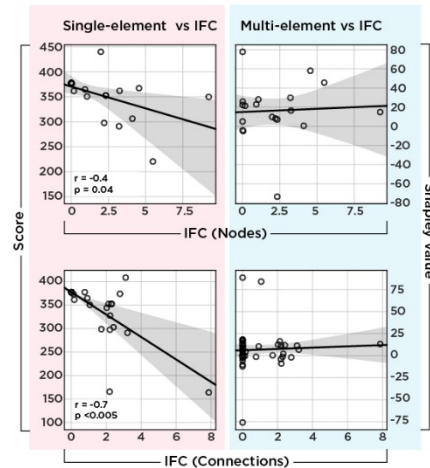
266 considered as the Impact of lesioning on Functional Connectivity (IFC; Fig.5). A larger IFC results
267 from a greater difference between FC of the intact network and FC of the lesioned network and
268 intuitively indicates the importance of elements, this time by their contribution to overall
269 functional connectivity instead of performance.



270

271 **Fig.5: Calculating the impact of lesions on functional connectivity.** We recorded the activity of all neurons to
272 compute the functional connectivity of the network. We exhaustively perturbed all units one by one and compared
273 the element-wise differences between intact and lesioned FC matrices. The absolute sum of this difference matrix
274 (IFC) quantifies how much a lesion caused the network dynamics to deviate from its uninterrupted state. On the left-
275 hand side, the activity of two scenarios is depicted. In the upper timeline, the network is intact, and the score is 670,
276 while in the lower timeline, the feedback loop (0 → 0) is lesioned, leading to a drastic decrease in performance. Red
277 vertical lines showed when the agent was shot and lost a life. Brown cells indicate the chosen action, and the dashed
278 window is the same time window that we zoomed in further in the section *Understanding the Paradoxical lesion*.

279 Interestingly, IFC is negatively correlated with both nodal and connection perturbation scenarios,
280 corroborating previous findings (Fig.6). However, IFC is not associated with Shapley values of
281 these elements. This means that, although SPA has internal coherency by identifying units that,
282 perturbed one by one, have the largest effect on both functional connectivity and the agent's
283 performance, these units are not the same as those captured by an MSA framework. In other
284 words, the bridge is formed. However, as shown in Fig.3, the actual players remained obscure. We
285 show why the rankings differ and propose a possible underlying mechanism that accounts for this
286 discrepancy in the next two sections.



287

288 **Fig.6: Correlation between IFC and single-site lesioning scheme.** The upper left scatterplot shows the
289 relationship between the impact of the SPA of nodes on functional connectivity and the agent's performance. The
290 lower left scatterplots show the same relationship but for each connection. Both show a negative correlation, which
291 means the larger the impact on functional connectivity, the lower the performance. However, this relationship is
292 absent from the right-hand side that compares the Shapley value of each element with their IFC. As with the left-hand
293 side, the x-axis shows the IFC of nodes (upper plot) and connections (lower plots), while here, the y-axis represents
294 Shapley value instead of raw performance.

295

296 Quantifying complex interactions between causal building blocks

297 In previous sections, we presented two causal rankings of elements from the same ground-truth
298 neural network model, one using a SPA framework and the other using MSA (Fig3.C). We found
299 that the changes in the inner dynamics of the system perturbed using SPA support this approach's
300 ranking, which mistakenly adds more certainty to the accuracy of the approach in finding critical
301 units. Here we show why these rankings differ by measuring the complex interactions of units.
302 Although MSA is a multivariate approach that accounts for a large variety of combinations of
303 units, it eventually describes the system in terms of how much, averaged over all combinations
304 with other units, *single units* contribute to the output. In other words, it isolates the average *individual*
305 contributions and not the nature of their interactions. Using an extension of MSA, here called
306 PCIA, we formalized and then quantified these interactions since the causal influence of one
307 element is intertwined with the state of others.

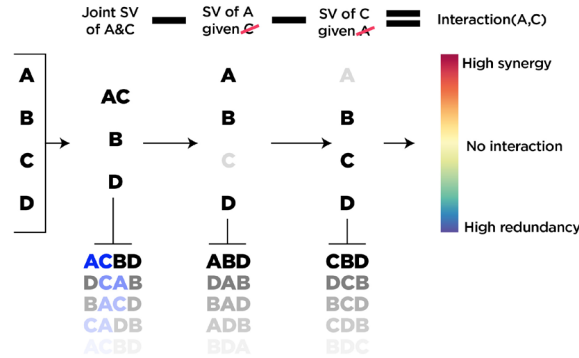
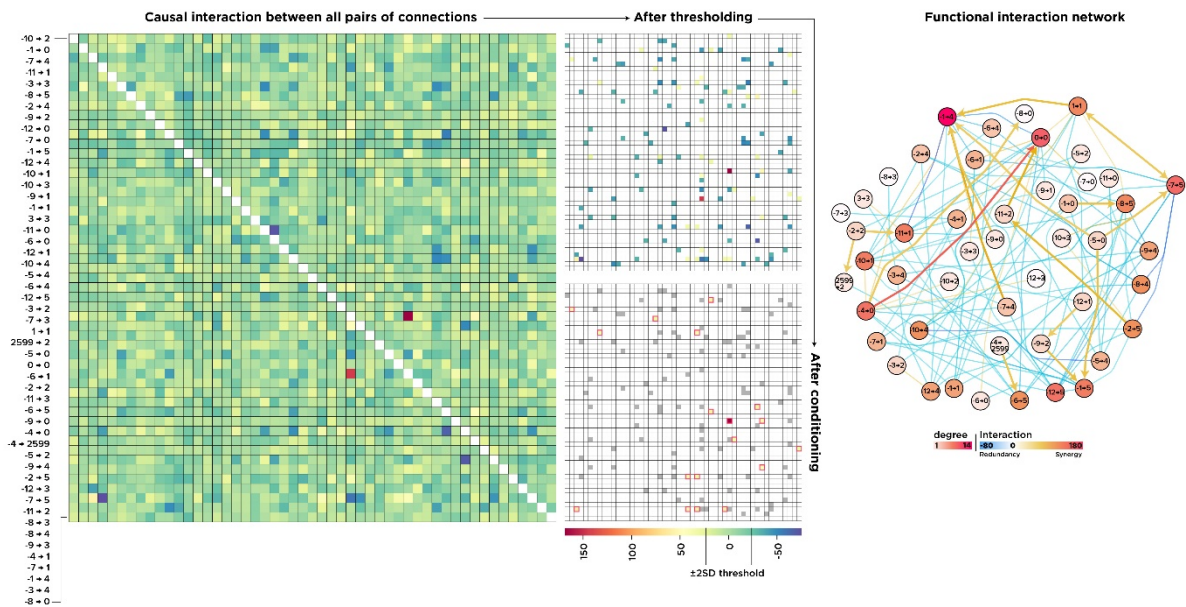


Fig.7: Visual depiction of the PCIA algorithm. At its core, PCIA comprises multiple MSAs. We first start with calculating the joint contribution of two elements, followed by the contribution of each, given the other is perturbed. The interaction term is then calculated by subtracting these values from each other, indicating how much the joint contribution of a pair of elements is bigger or smaller than the sum of their individual contributions. Like MSA, permutations are order-invariant.

308

309 At its core, PCIA is a chain of multiple MSAs in different conditions. To elaborate, quantifying
 310 the complex pairwise interaction of two elements i and j requires first to calculate the Shapley value
 311 of them both as a single compound element $\gamma_{(i,j)}$, followed by the Shapley value of each one given
 312 the other is perturbed $\gamma_{(i,j)}$ and $\gamma_{(i,j)}$ respectively. As Fig.7 shows, subtracting all three provides
 313 an interaction term that, if positive, indicates “*synergy*” between the pair and, if negative, shows
 314 “*redundancy*” or functional overlap. In other words, PCIA quantifies how much the causal
 315 contribution of a pair of units is bigger or smaller than the sum of their individual contributions.



316

Fig.8: Pairwise interactions among all connections. An interaction matrix resulted from the PCIA procedure in which warmer colors show greater synergy and cooler colors indicate functional overlap (left). We then excluded ± 2 SD and applied the “Sprague effect” condition to the thresholded matrix (middle). On the right-hand side, we plotted the interaction network in which the nodes represent connections in the actual network, and the edges are interactions among them. Arrows show paradoxical-lesion effects ($i \rightarrow j$).

321

322 Since PCIA involves the calculation of multiple MSAs, it is computationally even more expensive.
323 Therefore, we focused on the connections, and to calculate all pairs of them, we sampled 100
324 permutations per element instead of 1000, as in the case of MSA.

325 The results are shown in Fig.8, and as quickly stands out, there is a strong synergy between two
326 elements ($0 \rightarrow 0$) and ($-4 \rightarrow 0$), followed by a handful of strongly redundant and many minuscule
327 interactions in both directions. Therefore, the results from this analysis provide more evidence for
328 ($0 \rightarrow 0$) and ($-4 \rightarrow 0$) to have a unique form of interaction, which we next investigate with respect
329 to whether it is a paradoxical lesion-effect.

330 To do so, we formalized the Sprague effect as the difference between the average importance of
331 element i given the state of the element j . Specifically, the Sprague effect is defined as a scenario
332 in which element i has a negative Shapley value when element j is perturbed $\gamma_{(i,j)} < 0$, thus
333 hindering the performance and has a positive contribution when j is intact $I_{i,j} + \gamma_{(i,j)} > 0$. Put
334 simply, on average, element i disrupts the performance if element j is intact and improves if j is
335 lesioned[28,31].

336 To reduce the number of false-positive findings, we looked for this condition among a smaller set
337 of pairs with an interaction term above and below two standard deviations of the mean. The results
338 are shown in Fig.8, with connections indicating the interactions and arrows depicting a Sprague
339 effect between two elements (the stem of the arrow indicates the element i that has a negative
340 contribution when the pointed element j is lesioned.) As depicted, we found many paradoxical
341 lesion effects predominantly among synergistic interactions, with the interaction between ($0 \rightarrow 0$)
342 and ($-4 \rightarrow 0$) being the most prominent one. This network is a higher order “functional/interaction
343 network” in which its nodes represent connections in the “structural/actual network”.

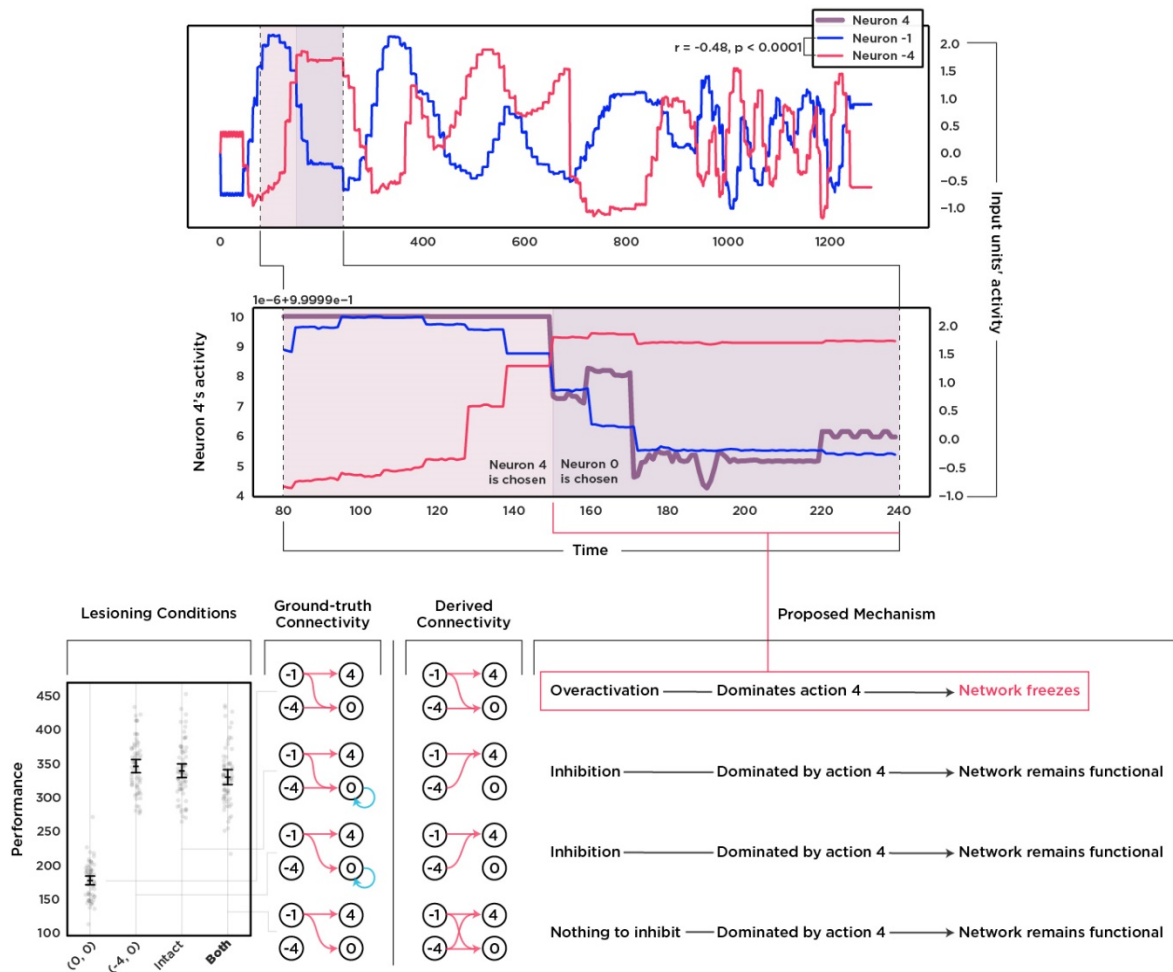
344 To summarize this part, we first quantified how much two elements’ causal importance is larger
345 or smaller than the sum of the individual elements. We then used this metric to classify the
346 modulatory effect of each element on the others, with a focus on paradoxical modulations, and
347 found a handful of elements in which lesioning one, while the other is perturbed, restored the
348 performance. The connections ($0 \rightarrow 0$) and ($-4 \rightarrow 0$) had the highest synergy, meaning that,
349 together as a whole, they functionally contribute much more than their summed individual
350 contributions. This unique synergy is also a paradoxical-lesion effect in that lesioning ($0 \rightarrow 0$) alone
351 disrupts the performance while lesioning it alongside ($-4 \rightarrow 0$) restores it. Note that the metric
352 captures what a SPA framework is insensitive to, specifically, complex pairwise causal interactions.
353 In other words, PCIA is built upon MSA that, as seen, extends SPA to lesioning combinations of
354 elements, and here, it is systematically bundled to quantify complex multivariate relationships that
355 elements might have. These interactions and insensitivity of SPA to them are what, we believe,
356 eventually leads to misattributing key elements in their ranking of causally critical units. By focusing
357 on the two connections ($0 \rightarrow 0$) and ($-4 \rightarrow 0$) in the next section, we show paradoxical lesion
358 effects might not be that unlikely and, quite contrary, they might be a direct result of perturbing a
359 simple and ubiquitous motif of connectivity, which explains why we found many such paradoxical
360 effects in this analysis.

361 Understanding the Paradoxical Lesion

362 The Sprague effect was first discovered in cats and later in humans, with its underlying mechanisms
363 still partly elusive[19,20]. One current theory suggests the phenomenon is caused by a reduction
364 of inhibition from a functionally competing region, and the deficit reverses when both are
365 lesioned[22]. To see if this is the case in our network, we focused on the two most prominent units
366 ($0 \rightarrow 0$) and ($-4 \rightarrow 0$). Note that the SPA also ranked ($0 \rightarrow 0$) among the most critical connections.
367 However, ($-4 \rightarrow 0$) was only captured by MSA and was the only unit with a large negative Shapley
368 value. The top plot in Fig.9 shows the activity of two input units, -1 and -4, over the trial in which
369 ($0 \rightarrow 0$) is lesioned (also see Fig.5). A Pearson's correlation analysis shows they are negatively
370 correlated.

371 Unit -1 is one of the key input units to neuron 4, which itself is one of the most frequently chosen
372 actions by the intact network. Input unit -4, however, has a major influence on neuron 0 (Fig.1)
373 that is inhibited by the negative feedback loop, causing neuron 0 to be silent in the intact network.
374 Since neuron 0 is the action “no action,” the intact network always chooses an action, either 4
375 (right and fire) or 5 (left and fire). As depicted in Fig.5, lesioning the feedback loop disrupts the
376 inhibition that leads to hyperactivation of neuron 0. Interestingly, although neuron 0 is now
377 competing with neuron 4, it takes roughly 150 timesteps to be selected as the chosen action. The
378 middle part of (Fig9) shows how the decaying activity of unit -1 at around that timepoint causes
379 neuron 4's activity to follow and eventually lose to neuron 0 in the lesioned network. Naturally,
380 the behavioral consequence of excessively choosing “no action” is gaining a substantially lower
381 score. By lesioning the input from -4 to 0 with or without the feedback loop, the node never
382 reaches the critical threshold to dominate other actions, and thus, in both conditions, the
383 performance remains uninterrupted (Fig.9).

384 Altogether by looking deeper into the inner dynamics of these units that MSA distinguished, we
385 see a simple motif of connectivity among only four units is enough to produce a paradoxical lesion
386 effect. The key nodes are neurons -4 and 0; the key connections are ($0 \rightarrow 0$) and ($-4 \rightarrow 0$). The
387 input from -4 to unit 0 has a large negative Shapley value because in coalitions without ($0 \rightarrow 0$), it
388 over-activates neuron 0 and causes the network to freeze. The feedback loop ($0 \rightarrow 0$) has a high
389 positive Shapley value because it prevents this over-activation, and removing it causes the network
390 to freeze. Interestingly, the input from -1 to 4 has the next highest Shapley value because, without
391 it, unit 4 is dominated by other units, especially an over-activated “no action” node.



392

393 **Fig.9: Focusing on the critical elements discovered by MSA.** The upper timeline shows the negative correlation
 394 between the activity of two input units -1 and -4. This anticorrelation leads to competition between downstream units
 395 4 and 0. In an intact network, unit 4 is dominant due to the inhibitory feedback loop of unit 0. The middle plot shows
 396 how unit 4 loses to 0 after the inhibitory loop is lesioned since it is tightly following the input from -1 while neuron 0
 397 is driven by the input from -4. The bottom-left part shows the implications of this rivalry on the performance and
 398 how it produces the paradoxical lesion effect. Lesioning the feedback loop disrupts the performance while lesioning
 399 it alongside the input from -4 restores the deficit since neuron 0 stays dominated. The bottom-middle part shows the
 400 discrepancy between the actual flow of information and the inferred flow by an mTE analysis. Notice the absence of
 401 connection between -4 to 0 in the intact network due to the self-inhibition of the target neuron.

402 A crucial side effect of the functional contribution of silenced nodes is that it becomes very difficult
 403 to infer their causal relationship relying on time-series analyses. Here we used a Multivariate
 404 Transfer Entropy (mTE) analysis on the four key players in three lesioning conditions and the
 405 intact network to see how well they infer information flow in the circuit. As Fig.9 shows, in
 406 conditions that neuron 0 is inhibited, mTE missed the information flow even though the node
 407 receives input from both -1 and -4.

408 To conclude this section, we showed that a paradoxical lesion effect could emerge from a simple
 409 inhibitory motif. In our case, the inhibition is a negative feedback loop, and the competition is
 410 between two output neurons, 4 and 0. We then used mTE analysis to infer the causal relationships
 411 that resulted in a critical relationship between -4 and 0 to be overlooked. This shows the necessity
 412 of employing systematic lesioning alongside methods relying, for example, on the analysis of time-
 413 series dynamics. Altogether, we show that, even in a simple agent, finding which elements are

414 causally relevant for behavior and how, is extremely difficult to answer with confidence. In the
415 next section, we discuss our results, limitations, and future improvements.

416

417 **Discussion**

418 In this work, we defined causation not as events prior to effects nor as entities that raise their
419 probability of occurrence but as *contributors to the effect*. Having this definition of causation, we aimed
420 to understand an ANN in terms of its components' causal influence over its performance. We
421 initially lesioned both its neurons and connections one at a time. We then showed that even with
422 such an exhaustive analysis, which is yet to be reached *in-vivo*, the results are persistently biased.
423 We then formed a bridge from structure to function and eventually to behavior by measuring the
424 impact of single-element lesioning on global functional connectivity. The results supported the
425 ranking from the SPA and added more confidence to the biased conclusion about which units are
426 critical. In other words, our SPA confirms the results from Jonas and Kording's work[26], and we,
427 too, ended up with structured but biased results.

428 We then used MSA, a rigorous game-theoretical algorithm, and found the causal ranking to be
429 different. For example, neuron 0 had the highest causal contribution even though it has no major
430 role according to SPA. MSA then identified crucial connections and ranked $(0 \rightarrow 0)$ the most
431 causally important. It also found neuron -4 to hinder the system and tracked the disruptive element
432 to be the connection $(-4 \rightarrow 0)$. Next, using an extension of MSA, we first quantified the complex
433 pairwise interaction of all causal building blocks (connections) and, after formalizing the Sprague
434 effect, found lesioning connections $(0 \rightarrow 0)$ and $(-4 \rightarrow 0)$ to have such an effect. Lastly, we looked
435 into these two units and found the rivalrous interaction to be the potential mechanism.

436 Two points to bear in mind are 1. our network was fixed throughout the experiments, leaving no
437 space for plasticity, and 2. the network is a simple ANN with no excitatory-to-inhibitory synaptic
438 dynamics. It is indeed possible that these physiological mechanisms underlie paradoxical lesion
439 effects in the living brain[20]. However, we did not include them in our model; therefore, we
440 believe the paradoxical effects observed here result from none of these mechanisms. We found
441 functional inhibition between competing units sufficient to produce a Sprague effect, as also
442 investigated before ([22,23] and see [31] for a fixed artificial network).

443 Besides that, further research is needed to compare different mechanisms using biologically
444 plausible neural network models since understanding the phenomenon also relies on different
445 analytical approaches as we used PCIA while Sajid et al. [25], for example, used a dual-lesion
446 scheme. On the same line, since our results point towards a type of interaction that is possibly
447 rooted in the pattern of connectivity in a very rudimentary system compared to the human brain,
448 comparative studies can shine a light on how deep the motif is embedded in the evolution of
449 nervous systems and what, if there is any, are the adaptive values. Interestingly, in our model, the
450 motif is more costly and sub-optimal because instead of simply removing the input from -4 to 0,
451 the evolutionary process added a negative feedback loop to cancel the disruptive influence
452 producing the motif that leads to a paradoxical lesion effect.

453 Overall, our results, first and foremost, show the inferential limitations of SPA. We believe many
454 aspects of a system can indeed be investigated and understood by lesioning its elements one at a
455 time. However, it is important to know which aspects cannot. The example we used was the
456 Sprague effect, which was argued to be either noise or exceptional[39]. We speculate that if our
457 compact network with 19 neurons and 51 connections evolved at least one of such effects, then it
458 might not be a rare event (as also argued in [5]) but an indication of complex multivariate functional
459 motifs of computation as proposed in [24].

460 A substantial challenge in depicting a mechanistic blueprint of any system is to have a solid causal
461 understanding of it. The conventional approach perturbs its elements and pinpoints those resulting
462 in a disrupted behavior[10]. These elements were then called necessary causes of the observed
463 effect since they serve as critical substrates for an intact behavior[40]. However, there have been
464 arguments against the classification of neural components as such (see [41–43]). Supporting those
465 arguments, we propose that one step towards a solid causal understanding of the brain is to instead
466 *quantify the degree* to which its neural elements contribute to cognition and behavior.

467 To put it into perspective, for behaviors to emerge, many neural circuits coordinate, cooperate,
468 and form coalitions that boiled down to a single “necessary” entity, resulting in losing crucial
469 information of the brain’s inner workings[6,42]. This was the case in the contributions we derived
470 from SPA. Thus, we used MSA to capture the whole spectrum of causation instead. Shapley value
471 results from a mathematically sound analysis of all possible combinations in which units can form
472 coalitions and produce the behavior, either flawlessly or disrupted. In its essence, Shapley value is
473 the *fair* share of the elements in producing the function so that the most important elements
474 assigned the highest share followed by a continuum of importance to zero for independent
475 elements and negative values for hindering ones. Therefore, it provides a rigorous and intuitive
476 way that neural elements can be ranked according to their causal contributions to the under-
477 investigated behavior.

478 Although powerful and intuitive, it is important to emphasize what Shapley value is not (see [44]
479 for a more technical perspective). For example, Shapley value by default does not reveal
480 mechanisms neither it shows what computations were done by individual elements. It shows *how*
481 *much each element is functionally contributing to the underlying mechanistic processes*. As mentioned, we believe
482 this is *the first step* towards a more comprehensive mechanistic description of the brain, illuminating
483 which elements to focus on next. We, too, did so by focusing on the few key elements that
484 summarize why the intact and lesioned networks behave such and why MSA chooses these units
485 as causally relevant.

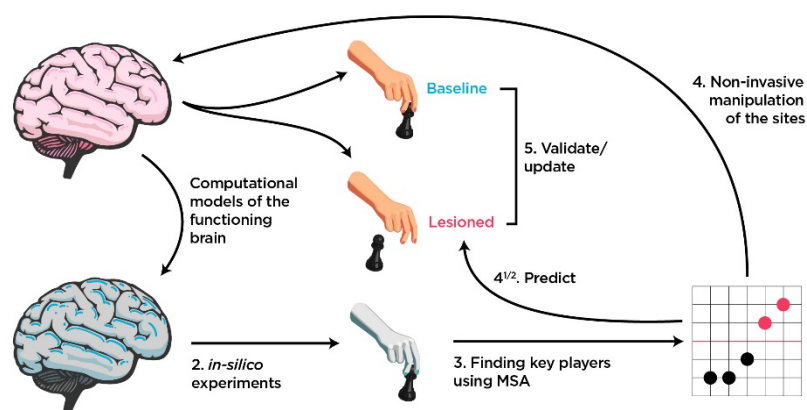
486 We can gain profound insight into the system by incorporating MSA in more complicated
487 analytical pipelines. For instance, the elements’ Shapley value will vary according to the behavior
488 of interest. In principle, one can produce a multidimensional map of each element, knowing how
489 much they are involved in various behaviors. This is relevant for neuroscience since it is shown
490 that brain regions are multifunctional and play different roles in different coalitions[45,46]. Having
491 negligible Shapley value in a task is not an indication of inutility but an indicator of independence
492 since the same element might have a considerably large share in the emergence of a different
493 function. This feature can be used to decompose and dissociate roles that neural elements play in
494 different tasks. For example, in this study, we used the system’s overall performance as our metric

495 that is the product of many behavioral primitives and found specific elements to be the most
496 critical.

497 Further analyses can decompose the behavior, as we did with the system itself from its nodes to
498 its connections, and calculate the causal share of the elements in each behavioral component, such
499 as, in our case, actions that construct the learned strategy. Therefore, we can expand our knowledge
500 of how elements dynamically form coalitions to solve sub-tasks of the given task, providing a
501 detailed description of the system's inner mechanism. In other words, given an elegant experiment
502 in which behavior and its components can be measured, Shapley value is a robust method to
503 unravel how neuronal units adaptively join communities and produce hierarchies in the brain. We
504 believe MSA is a powerful tool that can be used to understand the system far deeper than we
505 attempted to do here since, as described above, it has many favorable features and provides
506 intuitive results.

507 In this work, we used a version of Evolutionary Autonomous Agent models advocated by [47] to
508 be nifty tools for neuroscientists. Using NEAT, we allowed the network's topology to evolve with
509 respect to the environmental constraints instead of modeling the architecture ourselves and
510 optimize the weights or readout units. This way, we liberated ourselves from further assumptions
511 about the network's connectivity and structure. It is important to note that NEAT itself produces
512 simple networks that can do simple things. However, more advanced NAS algorithms such as
513 Hyper-NEAT[48] are gaining popularity in the AI community since they produce larger networks
514 that are not limited to the experimenter's design[49].

515 Interestingly, in some cases, genetic algorithms rival the conventional Gradient Descent-based
516 methods in non-trivial tasks[50]. This shows a potential role for such algorithms in neuroscience
517 since one can evolve arbitrary architectures to solve an ecologically valid task, e.g., foraging in a
518 patchy environment[51], and compare their topological features with brains evolved in such
519 environments. This extends the toolboxes available to computational neuroscientists,
520 neuroethologists, and behavioral ecologists to more realistic *in-silico* models and experiments.



521

522 **Fig.10: How MSA can be incorporated into the causal brain-mapping toolbox.** Since multiple *in-vivo* lesioning
523 is beyond the reach, we suggest connectivity-aware or neural network models of functioning brains to fill the gap. *In-*
524 *silico* experiments then can be performed to predict both key elements and their contributions to the behavior. These
525 predictions can then be tested *in-vivo* by the method of choice.

526 More cognitively and clinically oriented, *in-silico* multi-element lesioning experiments can be used
527 as a predictive tool to guide non-invasive brain stimulation experiments. For example, human brain
528 connectivity can be used as the backbones of ANNs trained to solve cognitive tasks[52–55]. These
529 connectivity-aware ANNs can then be investigated thoroughly using MSA to predict the critical
530 regions and the corresponding behavioral deficits. The predictions further can be used as testable
531 hypotheses about which regions to perturb *in-vivo*. In other words, connectivity-aware ANNs,
532 neural network models of cognitive processes[56], and large-scale models of functioning brains[57]
533 can add a unique value to the repertoire of ground-truth models to test brain-mapping tools and
534 their limitations (Fig.10).

535 The main limitation that is needed to be addressed is MSA’s computational complexity. Having an
536 analytical solution for Shapley values of large systems is an NP-complete problem[58]. Therefore,
537 heuristics[59], predictors[30], and estimators[28,60] are used and are under development to solve
538 this issue. Interestingly, Shapley value has found a unique spot in the field of explainable machine
539 learning[61] and is used to understand deeper and more complicated neural network
540 architectures[61], prune the unnecessary elements[62], and even correct biased networks[60].

541 Another limitation here is thresholding the “interaction matrix” (Fig.8). As mentioned, even
542 reliably estimating all elements’ pairwise interaction can quickly become impossible since the
543 number of elements is now squared, and three Shapley values are needed for each interaction.
544 Therefore, we reduced the number of samples from 1000 to 100, which means less certainty in the
545 estimated results. To partially account for this problem, we excluded two standard deviations
546 above and below the mean. A decision that directly influences the number of discovered
547 paradoxical-lesion effects. Therefore, a central interest is to address this issue using either better
548 thresholding criteria or estimation methods.

549 **Conclusion**

550 A common way of characterizing the causal contributions of elements in a system is to perturb
551 them and measure the effect. We showed that not every perturbation reveals causation since
552 lesioning elements, one at a time, produced coherent but biased results. We then used MSA and
553 captured the crucial details missed when we lesioned each site independently. We then found a
554 motif of functional inhibition among competing units to be the underlying mechanism of the
555 paradoxical lesion effects in our network. We believe this effect is the main contributor to the bias
556 in a single-site lesion analysis since, by definition, it emerges from a condition with at least two
557 lesions. This showed that even compact ANNs show surprising complexity that is needed to be
558 addressed to have a step towards a comprehensive causal picture of the system.

559 Lastly, in the context of rapidly evolving sophisticated uni-and-multivariate brain-mapping
560 methods, we advocate using *in-silico* experiments, and ground-truth models, especially neural
561 network models verify fundamental assumptions, technical limitations, and extent of
562 interpretations of these methods.

563 **Materials and methods**

564 In this section, we explain the methods and materials used in this research. The codes and
565 generated datasets are publicly available in the following repository:

566 <https://github.com/kuffmode/ANNLesionAnalysis>

567 Briefly, we first trained a deep autoencoder to compress the screen pixels to a handful of features
568 per frame. We then evolved a controller network to, based on these features, choose a proper
569 action. After having both networks, we started the lesioning experiments.

570

571 **Evolutionary optimization**

572 We used the NEAT-Python toolbox[63] to evolve a network from an initial stage of randomly
573 connected 12 input and six output nodes. During the evolutionary process, the algorithm was
574 optimizing many parameters, including the choice of activation functions, aggregation functions,
575 adding or removing hidden neurons, adjusting connection weights and node biases, and adding or
576 removing connections (see Table1 for a summary and the file AEconfig-SI.txt for the
577 complete list of hyperparameters). There were no restrictions on the connectivity pattern so that
578 a recurrent architecture could evolve from the initial feed-forward stage. We chose the probability
579 of removing connections to be slightly higher than adding (0.6 versus 0.5) to encourage sparsity.
580 We then ran the evolutionary processes 32 times to have 32 candidates. Each time the process
581 ended either after 128 trials or one member reached the fitness criterion of 1200 points. In each
582 trial, the generation comprised of 128 members that were instantiated from the same initial stage
583 and would play the ATARI game independently. After each step, the algorithm mutated the
584 genome according to the given probabilities and performed the cross-over among the top %30
585 networks to produce the next 128 members. At the end of the training phase, 32 candidate
586 networks reached either the generation limit or the fitness criterion. We then chose the one with
587 the highest score of 1300 points to move forward with the lesion experiments.

NEAT Hyperparameters	Value
Fitness Threshold	1200
Population Size	128
Activation Function's Mutation Rate	0.05
Aggregation Function's Mutation Rate	0.05
Probability of Linking Nodes	0.5
Probability of Removing Links	0.6
Probability of Adding Nodes	0.6
Probability of Removing Nodes	0.4
Number of Input Neurons	12
Initial Number of Hidden Neurons	0
Number of Output Neurons	6
Survival Threshold	0.3

588 **Table1: A summary of relevant NEAT hyperparameters.** NEAT produces a large variety of networks, all from a
589 set of constraints and probabilities. Since our goal was to produce a good-enough network, we did not tune these

590 parameters for maximum performance and either used the default values or adjusted them according to the
591 experimental objectives, e.g., sparse connectivity.

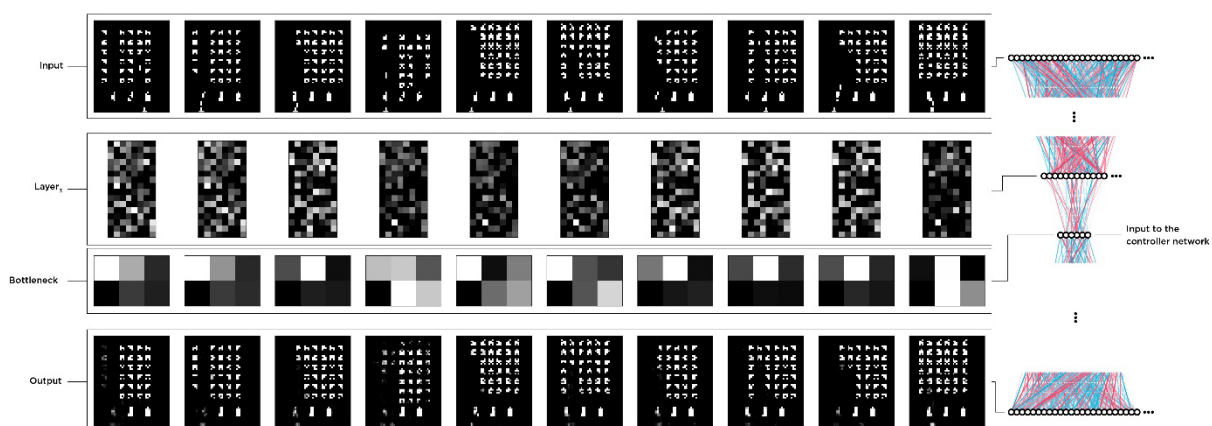
592

593 The preprocessing steps and the Autoencoder

594 We used OpenAI Gym[64] ATARI environment as our game environment. The game screen
595 generates an array with the size of (210, 160, 3); since the screen is 210×160 pixels, each contains
596 three color values, red, green, and blue. Throughout the whole work, the pixels passed through a
597 preprocessing pipeline first that would:

- 598 1. crop-out the unrelated parts of the screen such as scores and the ground,
- 599 2. convert colors to the monochrome gray-scale, therefore reducing the 3D space of red,
600 green, and blue values to one intensity value representing brightness of each pixel,
- 601 3. binarize the pixel values to either an “on pixel” or “off pixel”,
- 602 4. and finally, flatten the outcome into a vector with a size of 2679 pixels.

603 This vector represented the game with a series of zeros and ones that were then fed to the
604 Autoencoder (Fig.11). The Autoencoder was a Keras model[65] trained independently from the
605 controller network. We first recorded 43,200 frames from the game played by a random agent,
606 shuffled the frame orders, and used 28,800 frames ($\approx 65\%$ of the dataset) to train and the rest for
607 testing the Autoencoder. The architecture was designed with four encoding layers and four
608 decoding, and a bottleneck of six features.



609

610 **Fig.11: Visualization of the Autoencoder’s inputs, latent features, and decoded outputs.** The Autoencoder was
611 trained separately from the controller and received recorded frames from a random action selector agent. We then
612 used the encoder half to reduce the pixel space to six features per frame and fed the controller with two feature
613 vectors.

614 We used the ADAM optimizer, a binary cross-entropy loss function, 64 epochs, and a batch size
615 of 512. Since input frames are binarized, we used Rectified Linear Unit (ReLU) activation functions
616 for all layers except the last decoding layer, for which we used a Sigmoid function instead. After
617 the training session and accuracy of $\approx 98.8\%$, we kept the encoder network and fed the latent space
618 to the controller network throughout all experiments and the evolutionary process of the controller

619 (Fig.11). Together the Autoencoder and the controller network formed our agent. However, we
620 did not perturb the Autoencoder and focused solely on the controller during the experiments.

621

622 **Lesion Analysis**

623 We first pruned our network by pruning the already “disabled” connections. Briefly, connections
624 in the network are either enable, meaning they multiply the incoming value with the weights and
625 pass it to the receiver node, or disabled that pass zero. During the evolutionary process, these
626 disabled connections serve as “pseudogenes” *in-vivo* that can reactivate in later generations due to
627 mutation. Initially, the controller had 7 of them that, after pruning, we had 51 enabled connections
628 to target. We used the same attribute to lesion the connections by virtually disabling them from
629 passing values from source neurons to receivers. In other words, a lesion in our experiments means
630 a severed connection in which, technically, would disrupt the flow of information from the source
631 node to the receiver node. To lesion nodes, we then disabled the incoming/outgoing connections.
632 For example, to lesion a neuron that sends information to three other neurons, we set those three
633 connections to zero, which virtually silences the node.

634 Each lesion experiment started with silencing the targeted neuron or connection as described. All
635 experiments consisted of 512 trials in which the network played the game 16 times per trial. The
636 score of each trial was calculated by averaging these 16 scores, leading to a distribution of 512
637 scores per lesion experiment.

638

639 **Multi-perturbation Shapley Value Analysis**

640 MSA is a rigorous method based on a Game-theoretical metric called Shapley value, here γ that
641 indicates how much an element is important for the grand coalition. To elaborate, assume the
642 marginal importance of an element i to a set of elements S , with $i \notin S$ is:

$$643 \quad \Delta_i(S) = v(S \cup \{i\}) - v(S)$$

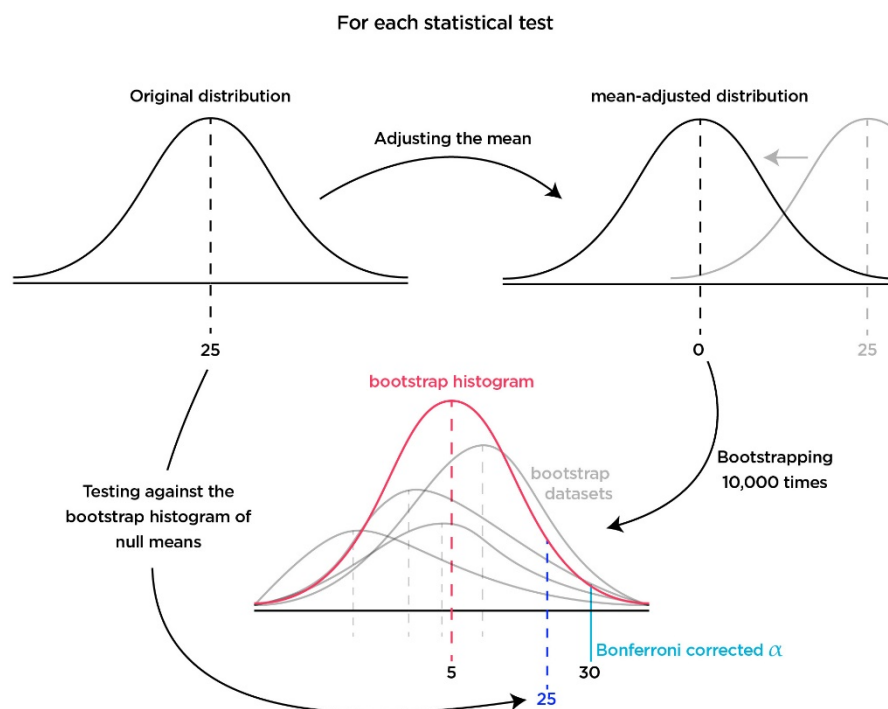
644 With v being the worth or importance of the element i , and S a coalition of elements. Then γ_i
645 while $i \in N$ is defined as:

$$646 \quad \gamma_i(N, v) = \frac{1}{n!} \sum_{R \in \mathcal{R}} \Delta_i(S_i(R))$$

647 With \mathcal{R} is the set of all $n!$ orderings of N , and $S_i(R)$ is the set of elements preceding i , in the
648 ordering R . We estimated γ of each neuron and connections by sampling 1000 orders from the
649 permutation space of $19!$ for neurons and $51!$ for connections. These 1000 permutations then
650 dictate which combinations of elements should be lesioned (Fig4). After selecting the target
651 elements, we used the same perturbation approach as the single-site lesion and disabled the
652 corresponding connections. The agent played the game 16 times, and the average score would be
653 the score of that random permutation, providing a γ distribution of 1000 data points for each
654 element. Altogether, we had around 70,000 unique combinations of lesions to estimate γ from.

655 Statistical Inference

656 Besides testing the performance of the intact network against the random agent, blind, and weight-
657 shuffled networks in which we used the non-parametric Mann-Whitney U test, we used bootstrap
658 hypothesis testing to find significant statistics throughout the study. We first generated a synthetic
659 null distribution for each statistical test by shifting the original distribution towards the H_0 's mean
660 value, either zero or an arbitrary number. For instance, to compare a distribution against a null
661 distribution centered around zero, such as Shapley values, we subtracted the average from each
662 data point, centered synthetic distributions around zero. In cases in which we tested distributions
663 against a second distribution that is not centered around zero, such as the performance of the
664 single-lesioned network versus the performance of the intact network, we shifted the synthetic
665 distributions toward the H_0 's mean, in this example, around 337 by adding the mean to each data
666 point.



667

668 **Fig.12: Visual diagram of the hypothesis testing process.** For each test, we first made a null distribution by
669 adjusting the mean. Then we resampled the synthetic distribution and kept track of the averages in the bootstrap
670 histogram. Lastly, we checked if the original mean falls below or above the Bonferroni corrected p-value.

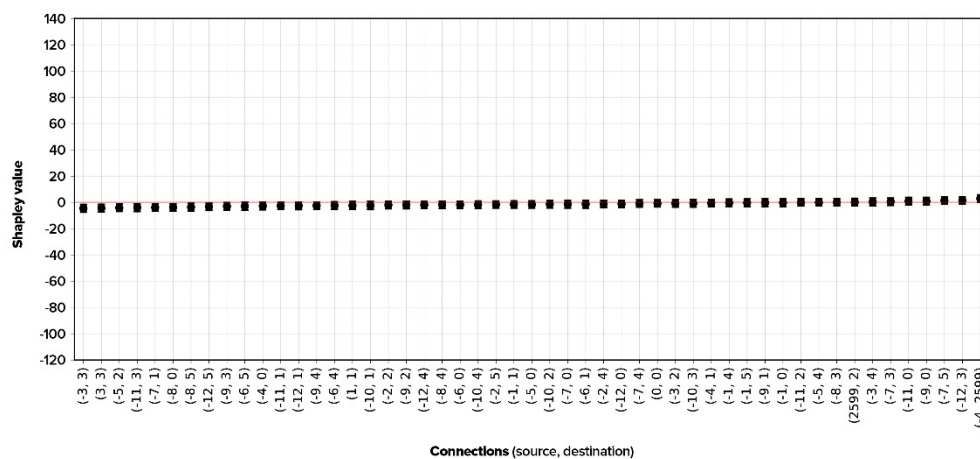
671 We then performed the bootstrapping and resampled the mean-adjusted distributions N times
672 with replacement, with N being the number of original samples, e.g., 512 for single-site lesions.
673 This generated a bootstrap dataset centered around the H_0 's mean (Fig.12). We then calculated the
674 bootstrap dataset's mean and repeated the process 10,000 times to generate the bootstrap
675 histogram of the means. In other words, the bootstrap histogram is a distribution of means if they
676 were from a null hypothesis. We then checked if the mean values of our distributions fall above
677 or below the p-value that is corrected for multiple comparisons using the Bonferroni correction
678 method ($0.05/\text{Number of tests}$).

679

680 Multivariate Transfer Entropy Analysis

681 We used the Information Dynamics Toolkit xl (IDTxl; [66]) to analyze mTE between a set of targets
682 (nodes 0 and 4) and sources (-1 and -4) in four conditions. First, the intact network, then the
683 feedback loop from 0 to itself is lesioned, then the input from -4 to 0, and lastly, both the feedback
684 loop and the input were lesioned. For each condition, we simulated 50 trials in which each trial
685 had 1200 samples. We enforced this number by discarding trials with fewer samples and cutting
686 the excessive samples from trials with more than 1200. Due to the quasi-binary dynamics of the
687 target nodes, we used the Kraskov estimator instead of Granger causality to infer multivariate
688 transfer entropy among the sources and targets. We further added information about the chosen
689 action to the time series of the target nodes. If the node is chosen at time point t , then the value
690 of the chosen node will be the value +1, and if not, just the raw data point (between 0 and 1) was
691 stored. The reason was to account for saturation of the target nodes since, at some points, the
692 actual values are very close to one another. Lastly, we injected a small amount of noise into the
693 estimator (noise level = $1e-7$). Both the minimum and maximum lag were set to 1 although we
694 explored maximum lags of two and three. Eventually, we discarded the resulted lags and only
695 reported the existence of TE between the pair of source and target since we found lags to be
696 irrelevant for this analysis. To account for multiple comparisons, we set the number of omnibus
697 permutations to 1000 and used the Bonferroni correction method to adjust the p-value ($0.05/8$),
698 which sets the adjusted value to around 0.005.

699 Supplementary Figures



700

701 **FigS1: Shapley Values of the blinded network.** As a sanity check, we performed the MSA on the optimized network
702 connections while feeding it noise instead of game-states. The procedure is explained in the section: *Multi-perturbation*
703 *Shapley value Analysis*. We found no connection with considerable causal importance since the network cannot perform
704 properly.

705

706 Conflict of Interest

707 The authors declare no conflict of interest.

708 **Acknowledgments**

709 Funding is gratefully acknowledged: KF: Deutsche Forschungsgemeinschaft, Germany (SFB
710 936/A1; TRR169/A2), SMARTSTART, the joint training program in computational neuroscience
711 of the Bernstein Network and the Volkswagen Foundation. CCH: Deutsche
712 Forschungsgemeinschaft, Germany (SFB 936/A1; TRR169/A2; SFB 1461/A4; SPP 2041/HI
713 1286/7-1, HI 1286/6-1), the Human Brain Project, EU (SGA2, SGA3). Authors thank Alexandros
714 Goulas, Fabrizio Damicelli, Fatemeh Hadaeghi, Joseph Lizier, Patricia Wollstadt, and Caroline
715 Malherbe for their valuable comments and insights.

716 References

1. Adolphs R. The unsolved problems of neuroscience. *Trends Cogn Sci.* 2015;19: 173–175. doi:10.1016/j.tics.2015.01.007
2. Mah YH, Husain M, Rees G, Nachev P. Human brain lesion-deficit inference remapped. *Brain.* 2014;137: 2522–2531. doi:10.1093/brain/awu164
3. Adolphs R. Human Lesion Studies in the 21st Century. *Neuron.* 2016;90: 1151–1153. doi:10.1016/j.neuron.2016.05.014
4. Bartolomeo P. The quest for the “critical lesion site” in cognitive deficits: Problems and perspectives. *Cortex.* 2011;47: 1010–1012. doi:10.1016/j.cortex.2010.11.007
5. Xu T, Jha A, Nachev P. The dimensionalities of lesion-deficit mapping. *Neuropsychologia.* 2018;115: 134–141. doi:10.1016/j.neuropsychologia.2017.09.007
6. Nachev P. The first step in modern lesion-deficit analysis. *Brain.* 2015;138: e354. doi:10.1093/brain/awu275
7. Ramaswamy V. An Algorithmic Barrier to Neural Circuit Understanding. *bioRxiv.* 2019; 639724. doi:10.1101/639724
8. Towilson EK, Barabási A-L. Synthetic ablations in the *C. elegans* nervous system. *Netw Neurosci.* 2020;4: 200–216. doi:10.1162/netn_a_00115
9. Petersen SE, Sporns O. Brain Networks and Cognitive Architectures. *Neuron.* 2015;88: 207–219. doi:10.1016/j.neuron.2015.09.027
10. Vaidya AR, Pujara MS, Petrides M, Murray EA, Fellows LK. Lesion Studies in Contemporary Neuroscience. *Trends Cogn Sci.* 2019; 1–19. doi:10.1016/j.tics.2019.05.009
11. Thiebaut De Schotten M, Dell’Acqua F, Ratiu P, Leslie A, Howells H, Cabanis E, et al. From phineas gage and monsieur leborgne to H.M.: Revisiting disconnection syndromes. *Cereb Cortex.* 2015;25: 4812–4827. doi:10.1093/cercor/bhv173
12. Yizhar O, Fenno LE, Davidson TJ, Mogri M, Deisseroth K. Optogenetics in Neural Systems. *Neuron.* 2011. pp. 9–34. doi:10.1016/j.neuron.2011.06.004
13. Roth BL. DREADDs for Neuroscientists. *Neuron.* 2016. pp. 683–694. doi:10.1016/j.neuron.2016.01.040
14. Karnath H-O, Sperber C, Wiesen D, de Haan B. Lesion-Behavior Mapping in Cognitive Neuroscience: A Practical Guide to Univariate and Multivariate Approaches. *Molecular Imaging in the Clinical Neurosciences.* 2019. pp. 345–357. doi:10.1007/7657_2019_18
15. Saxena S, Cunningham JP. Towards the neural population doctrine. *Curr Opin Neurobiol.* 2019;55: 103–111. doi:10.1016/j.conb.2019.02.002
16. Vyas S, Golub MD, Sussillo D, Shenoy KV. Computation Through Neural Population Dynamics. *Annu Rev Neurosci.* 2020;43: 249–275. doi:10.1146/annurev-neuro-092619-094115
17. Mišić B, Sporns O. From regions to connections and networks: new bridges between brain and behavior. *Curr Opin Neurobiol.* 2016;40: 1–7. doi:10.1016/j.conb.2016.05.003
18. Young MP, Hilgetag CC, Scannell JW. On imputing function to structure from the behavioural effects of brain lesions. *Philos Trans R Soc Lond B Biol Sci.* 2000;355: 147–161. doi:10.1098/rstb.2000.0555
19. Sprague JM. Interaction of cortex and superior colliculus in mediation of visually guided behavior in the cat. *Science.* 1966;153: 1544–1547. doi:10.1126/science.153.3743.1544
20. Valero-Cabré A, Toba MN, Hilgetag CC, Rushmore RJ. Perturbation-driven paradoxical facilitation of visuo-spatial function: Revisiting the ‘Sprague effect.’ *Cortex.* 2019;2: 10–39. doi:10.1016/j.cortex.2019.01.031
21. Kapur N, Pascual-Leone A, Ramachandran V, Cole J, Sala SD, Manly T, et al. The paradoxical brain. 2011. doi:10.1017/CBO9780511978098
22. Hilgetag CC. Spatial neglect and paradoxical lesion effects in the cat- A model based on midbrain connectivity. *Neurocomputing.* 2000. doi:10.1016/S0925-2312(00)00246-0
23. Hilgetag CC, Kotter R, Young MP. Inter-hemispheric competition of sub-cortical structures is a crucial mechanism in paradoxical lesion effects and spatial neglect. *Prog Brain Res.* 1999. doi:10.1016/s0079-6123(08)63071-x
24. Toba MN, Godefroy O, Rushmore RJ, Zavaglia M, Maatoug R, Hilgetag CC, et al. Revisiting ‘brain modes’ in a new computational era: approaches for the characterization of brain-behavioural associations. *Brain.* 2019. doi:10.1093/brain/awz343
25. Sajid N, Parr T, Gajardo-Vidal A, Price CJ, Friston KJ. Paradoxical lesions, plasticity and active inference. *Brain Communications.* 2020; 1–10. doi:10.1093/braincomms/fcaa164
26. Jonas E, Kording KP. Could a Neuroscientist Understand a Microprocessor? *Diedrichsen J, editor. PLoS Comput Biol.* 2017;13: e1005268. doi:10.1371/journal.pcbi.1005268
27. Keinan A, Hilgetag CC, Meilijson I, Ruppin E. Causal localization of neural function: the Shapley value method. *Neurocomputing.* 2004;58–60: 215–222. doi:10.1016/j.neucom.2004.01.046
28. Keinan A, Sandbank B, Hilgetag CC, Meilijson I, Ruppin E. Axiomatic scalable neurocontroller analysis via the Shapley value. *Artif Life.* 2006;12: 333–352. doi:10.1162/artl.2006.12.3.333
29. Shapley LS. A value for n-person games. In: H. W. Kuhn & A. W. Tucker, editor. *Contributions to the theory of games.* Princeton, NJ: Princeton University Press; 1953. pp. 307–317.
30. Toba MN, Zavaglia M, Rastelli F, Valabrègue R, Pradat-Diehl P, Valero-Cabré A, et al. Game theoretical mapping of causal interactions underlying visuo-spatial attention in the human brain based on stroke lesions. *Hum Brain Mapp.* 2017;3471: 3454–3471. doi:10.1002/hbm.23601
31. Keinan A, Sandbank B, Hilgetag CC, Meilijson I, Ruppin E. Fair attribution of functional contribution in artificial and biological networks. *Neural Comput.* 2004;16: 1887–1915. doi:10.1162/0899766041336387
32. Mnih V, Kavukcuoglu K, Silver D, Rusu AA, Veness J, Bellemare MG, et al. Human-level control through deep reinforcement learning. *Nature.* 2015;518: 529–533. doi:10.1038/nature14236

33. Hausknecht M, Lehman J, Miikkulainen R, Stone P. A Neuroevolution Approach to General Atari Game Playing. 2014;6: 355–366.
34. Stanley KO, Miikkulainen R. Evolving Neural Networks through Augmenting Topologies. *Evol Comput.* 2002;10: 99–127. doi:10.1162/106365602320169811
35. Thiebaut de Schotten M, Foulon C, Nachev P. Brain disconnections link structural connectivity with function and behaviour. *Nat Commun.* 2020;11. doi:10.1038/s41467-020-18920-9
36. Aerts H, Fias W, Caeyenberghs K, Marinazzo D. Brain networks under attack: robustness properties and the impact of lesions. *Brain.* 2016;139: 3063–3083. doi:10.1093/brain/aww194
37. Alstott J, Breakspear M, Hagmann P, Cammoun L, Sporns O. Modeling the impact of lesions in the human brain. *PLoS Comput Biol.* 2009;5: e1000408. doi:10.1371/journal.pcbi.1000408
38. O'Reilly JX, Croxson PL, Jbabdi S, Sallet J, Noonan MP, Mars RB, et al. Causal effect of disconnection lesions on interhemispheric functional connectivity in rhesus monkeys. *Proc Natl Acad Sci U S A.* 2013;110: 13982–13987. doi:10.1073/pnas.1305062110
39. Sperber C, Karnath H-O. Inhibition between human brain areas or methodological artefact? *Brain.* 2020;143: e38–e38. doi:10.1093/brain/awaa092
40. Brennan A. Necessary and Sufficient Conditions. The Stanford Encyclopedia of Philosophy. 2017. Available: <https://plato.stanford.edu/archives/sum2017/entries/necessary-sufficient/>
41. Yoshihara M, Yoshihara M. “Necessary and sufficient” in biology is not necessarily necessary - confusions and erroneous conclusions resulting from misapplied logic in the field of biology, especially neuroscience. *J Neurogenet.* 2018;32: 1–12. doi:10.1080/01677063.2018.1468443
42. Gomez-Marin A. Causal circuit explanations of behavior: Are necessity and sufficiency necessary and sufficient? 2016; 1–22.
43. Wolff SBE, Ölveczky BP. The promise and perils of causal circuit manipulations. *Curr Opin Neurobiol.* 2018;49: 84–94. doi:10.1016/j.conb.2018.01.004
44. Zavaglia M, Forkert ND, Cheng B, Gerloff C, Thomalla G, Hilgetag CC. Technical considerations of a game-theoretical approach for lesion symptom mapping. *BMC Neurosci.* 2016;17: 1–15. doi:10.1186/s12868-016-0275-6
45. Faskowitz J, Esfahlani FZ, Jo Y, Sporns O, Betzel RF. Edge-centric functional network representations of human cerebral cortex reveal overlapping system-level architecture. *Nature Neuroscience.* 2020. doi:10.1038/s41593-020-00719-y
46. de Reus MA, Saenger VM, Kahn RS, van den Heuvel MP. An edge-centric perspective on the human connectome: Link communities in the brain. *Philos Trans R Soc Lond B Biol Sci.* 2014;369. doi:10.1098/rstb.2013.0527
47. Ruppín E. Evolutionary autonomous agents: A neuroscience perspective. *Nat Rev Neurosci.* 2002;3: 132–141. doi:10.1038/nrn729
48. Stanley KO, D'Ambrosio DB, Gauci J. A hypercube-based encoding for evolving large-scale neural networks. *Artif Life.* 2009. doi:10.1162/artl.2009.15.2.15202
49. Stanley KO, Clune J, Lehman J, Miikkulainen R. Designing neural networks through neuroevolution. *Nature Machine Intelligence.* 2018;1: 24–35. doi:10.1038/s42256-018-0006-z
50. Such FP, Madhavan V, Conti E, Lehman J, Stanley KO, Clune J. Deep neuroevolution: Genetic algorithms are a competitive alternative for training deep neural networks for reinforcement learning. *arXiv.* 2017. Available: <http://arxiv.org/abs/1712.06567>
51. Risi S, Hughes CE, Stanley KO. Evolving plastic neural networks with novelty search. *Adapt Behav.* 2010;18: 470–491. doi:10.1177/1059712310379923
52. Su LE, Richards BA, Lajoie G, Misić B. Learning function from structure in neuromorphic networks. 2020.
53. Damicelli F, Hilgetag CC, Goulas A. Brain Connectivity meets Reservoir Computing. 2021; 1–14.
54. Goulas A, Damicelli F, Hilgetag CC. Bio-instantiated recurrent neural networks. 2021.
55. van Albada SJ, Morales-Gregorio A, Dickscheid T, Goulas A, Bakker R, Bludau S, et al. Bringing Anatomical Information into Neuronal Network Models. *arXiv.* 2020. Available: <http://arxiv.org/abs/2007.00031>
56. Kietzmann TC, Spoerer CJ, Sörensen LKA, Cichy RM, Hauk O, Kriegeskorte N. Recurrence is required to capture the representational dynamics of the human visual system. *Proceedings of the National Academy of Sciences.* 2019; 201905544. doi:10.1073/pnas.1905544116
57. Eliasmith C, Stewart TC, Choo X, Bekolay T, DeWolf T, Tang Y, et al. A Large-Scale Model of the Functioning Brain. *Science.* 2012;338: 1202–1205. doi:10.1126/science.1225266
58. Fatima SS, Wooldridge M, Jennings NR. A linear approximation method for the Shapley value. *Artif Intell.* 2008;172: 1673–1699. doi:10.1016/j.artint.2008.05.003
59. van Campen T, Hamers H, Husslage B, Lindelauf R. A new approximation method for the Shapley value applied to the WTC 9/11 terrorist attack. *Social Network Analysis and Mining.* 2018;8: 1–12. doi:10.1007/s13278-017-0480-z
60. Ghorbani A, Zou J. Neuron Shapley: Discovering the responsible neurons. *arXiv.* 2020. Available: <http://arxiv.org/abs/2002.09815>
61. Lundberg S, Lee S-I. A unified approach to interpreting model predictions. *arXiv [cs.AI].* 2017. pp. 426–430. Available: <http://arxiv.org/abs/1705.07874>
62. Stier J, Gianini G, Granitzer M, Ziegler K. Analysing Neural Network Topologies: A Game Theoretic Approach. *Procedia Comput Sci.* 2018;126: 234–243. doi:10.1016/j.procs.2018.07.257
63. Silva AMMKCGMCF da. NEAT-Python. Available: <https://github.com/CodeReclaimers/neat-python>
64. Brockman G, Cheung V, Pettersson L, Schneider J, Schulman J, Tang J, et al. OpenAI Gym. 2016; 1–4. Available: <http://arxiv.org/abs/1606.01540>
65. Chollet F. keras. 2015.
66. Wollstadt P, Lizier J, Vicente R, Finn C, Martinez-Zarzuela M, Mediano P, et al. IDTxl: The Information Dynamics Toolkit xl: a Python package for the efficient analysis of multivariate information dynamics in networks. *J Open Source Softw.* 2019;4: 1081. doi:10.21105/joss.01081

