

1 Habit formation viewed as structural
2 change in the behavioral network.

3

4 AUTHORS:

5 Kota Yamada^{1,2}, Koji Toda¹

6

7 AFFILIATIONS:

8 1 Department of Psychology, Keio University, Tokyo, JAPAN

9 2 Japan Society for Promotion of Science, Tokyo, JAPAN

10

11 CORRESPONDENCE:

12 Kota Yamada

13 Department of Psychology, Keio University

14 Mita 2-15-45, Minato-ku, Tokyo, JAPAN

15 Email: haroldthebarrel.yk@gmail.com

16

17 Abstract

18 Habit formation is a process in which an action becomes involuntary. While goal-directed
19 behavior is driven by its consequences, habits are elicited by a situation rather than its
20 consequences. Existing theories have proposed that actions are controlled by
21 corresponding two distinct systems. Although canonical theories based on such distinctions
22 are starting to be challenged, a few theoretical frameworks that implement goal-directed
23 behavior and habits within a single system. Here, we propose a novel theoretical framework
24 by hypothesizing that behavior is a network composed of several responses. With this
25 framework, we have shown that the transition of goal-directed actions to habits is caused by
26 a change in a single network structure. Furthermore, we confirmed that the proposed
27 network model behaves in a manner consistent with the existing experimental results
28 reported in animal behavioral studies. Our results revealed that habit could be formed under
29 the control of a single system rather than two distinct systems. By capturing the behavior as
30 a single network change, this framework provides a new perspective on studying the
31 structure of the behavior for experimental and theoretical research.

32

33

34 Author summary

35 To obtain the desired consequences, organisms need to respond based on the knowledge
36 of the consequences obtained by the response and the change in the environment caused
37 by it. Such a process is called goal-directed behavior, which is flexible, but requires high
38 computational cost. Once the same response is repeatedly performed under the same
39 environment, the response becomes automatic, and transforms into a habit. In the canonical
40 views, such a change from goal-directed response to habit was explained by the associative
41 structures between the corresponding systems, goal-directed, and habit systems. However,
42 the dichotomy in the mechanisms of behavior between goal-directed responses and habits
43 has recently been challenged. Here, we show that, instead of assuming two explicitly
44 distinguished mechanisms as in the canonical views, behavior is regarded as a network
45 consisting of multiple responses, and that changes in the structure of the network cause two
46 behavioral features, goal-directed behavior and habit. The transition from goal-directed
47 behavior to habit has been operationally defined by sensitivity to the reward obtained by the
48 response. We replicate such an experimental paradigm in the simulation and show that the
49 behavioral network model can reproduce the empirical results on habit formation obtained
50 from animal experiments. Our results demonstrate that habit formation can be explained in
51 terms of changes in the network structure of behavior without assuming explicitly distinct
52 systems and thus, provide a new theoretical framework to study the psychological, biological,
53 and computational mechanisms of the behavior.

54

55 Introduction

56 To behave flexibly in a given environment, organisms need to choose their actions based
57 on the consequences of the actions. This type of behavior is called goal-directed behavior.
58 As we keep repeating the same action under a certain situation, the action is elicited by the
59 situation rather than its consequences. This type of behavior is called a habit. Goal-directed
60 behavior requires high computational resources because organisms must process the
61 information about their external environment and how their actions affect it. In contrast, habit
62 shows a more stereotyped and less flexible behavior, requiring less computation. In this
63 sense, habit formation can be viewed as the optimization process of energy consumption by
64 the organism.

65 Existing theories about habit formation are based on evidence from experimental or
66 theoretical research in psychology and neuroscience. In the canonical view, responses are
67 controlled by two different systems: goal-directed and habit systems. Such theories
68 proposed that goal-directed and habit systems control responses by assigning different
69 weights, and the difference in the weights determines whether the response is goal-directed
70 or habit^{1, 2}. In this assumption, habit formation can be viewed as losing control by the
71 consequence of the response or reward sensitivity. However, some models explain habits
72 in a multistage Markov decision task and challenge the canonical dichotomy of goal-directed
73 and habits systems^{3, 4}. In addition, some researchers reviewed existing studies on habit
74 formation and cast doubt on the canonical framework of habit formation by showing the
75 possibilities that habits are also controlled by their consequences^{5, 6}.

76 In contrast to the canonical view, Dezfouli and Balleine⁷ proposed a new perspective
77 that habit formation can be viewed as shaping or acquiring response sequences. In their
78 model, an agent chooses their goal in a goal-directed manner and generates a response
79 sequence to reach there. Although habits are viewed as a lack of reward sensitivity in the
80 canonical view, their new model considers stereotyped behaviors as acquired response
81 sequences. To what extent could this model change the way of viewing accumulating
82 evidence of habit formation? Garr and D'Amato⁸ shows that rats acquired stereotyped
83 response sequences did not lose reward sensitivity. In a series of studies reported by
84 Dezfouli and Balleine^{7, 9, 10} dealt with only a few experiments on the reward sensitivity in free
85 operant situations¹¹⁻¹⁵. Another approach employs the planning process^{3, 4}. Pezzulo et al.³

86 stressed the importance of planning in goal-directed behaviors and built a single mixed-
87 controller model consisting of goal-directed behaviors and habits. Keramati et al.⁴ proposed
88 that the canonical goal-directed and habits systems can be viewed as edges of the spectrum
89 by building an integrated model of goal-directed planning and habits. Although application
90 of their models was limited to the multistage choice task, the model could serve as a basis
91 for a novel model with common assumptions and additional applicability in experiments on
92 reward sensitivity in free situations¹¹⁻¹⁵.

93 Here, instead of assuming two explicitly distinguished mechanisms as in the
94 canonical views, we consider behavior as a network consisting of multiple responses and
95 show that changes in the structure of the network cause two behavioral features, goal-
96 directed behavior and habit. By doing so, we could explain the lack of reward sensitivity in
97 habit formation, which is a characteristic of the canonical view on habits.

98 Behavioral network

99 There are two methodological approaches for studying animal behavior. One stream is an
100 in-laboratory psychological approach that studies the behavior of animals, including humans,
101 under experimentally controlled environments. Here, investigators measure only
102 experimentally defined responses of subjects (lever press, key peck, nose poke, freezing,
103 salivation, licking, eye blink, etc.) or put them into rigidly controlled situations where they can
104 only engage in the responses to the well-defined stimulus. Another stream is an ethological
105 approach that studies animal behavior under more natural and ecologically valid
106 environments¹⁷. In this case, behavior that the organism is engaged in the real world could
107 be observed, but the stimulus is difficult to control in terms of the strength, frequency, timing,
108 etc. Although these two approaches seem to conflict with each other, both are
109 complementary for understanding behavior and its biological substrates. Recent advances
110 in machine learning have allowed us to objectively measure the detailed structure of
111 behaviors¹⁷⁻¹⁹. Animals are engaged in more than lever press, key peck, or nose poke, they
112 approach and orient to the stimulus, and walk or sniff around and explore in the given
113 environment. Although the importance of observation and measurement of the behavior
114 during learning was attempted in classic behavioral studies²⁰⁻²³, current behavioral
115 quantification methods are expected to reveal the relationship between behavior and its
116 underlying mechanism in a way that integrates the different disciplines of psychology,

117 neuroscience, and ethology^{24, 25}. However, conventional views on behavior in psychology and
118 neuroscience are based on empirical results obtained from the approaches before the
119 appearance of such a new quantification technique of the behavior. Here, we present a new
120 theoretical novel framework that focuses on how behavior is organized and how its structure
121 brings specific characteristics to behavior.

122 Existing studies measured only specific experimenter-defined responses of animals
123 including humans, and ignored various responses that the animals actually engaged in.
124 However, there is considerable evidence that animals engage in various responses which
125 affect the learned responses. For example, animals engage in a specific response
126 immediately after the reward presentation²⁶⁻²⁸, engage in responses irrelevant to an
127 experiment²¹, show a specific response sequence between reward presentations²², or show
128 a specific response that counteracts learned responses²⁹. Theoretically, some characteristics
129 of operant responses are explained by assuming the existence of other responses³⁰⁻³⁴. These
130 experimental facts and theoretical assumptions indicate that animal responses do not exist
131 in isolation but are associated with other responses. We assume such relationships between
132 responses as a network in which responses and transitions between them are considered
133 nodes and edges, respectively.

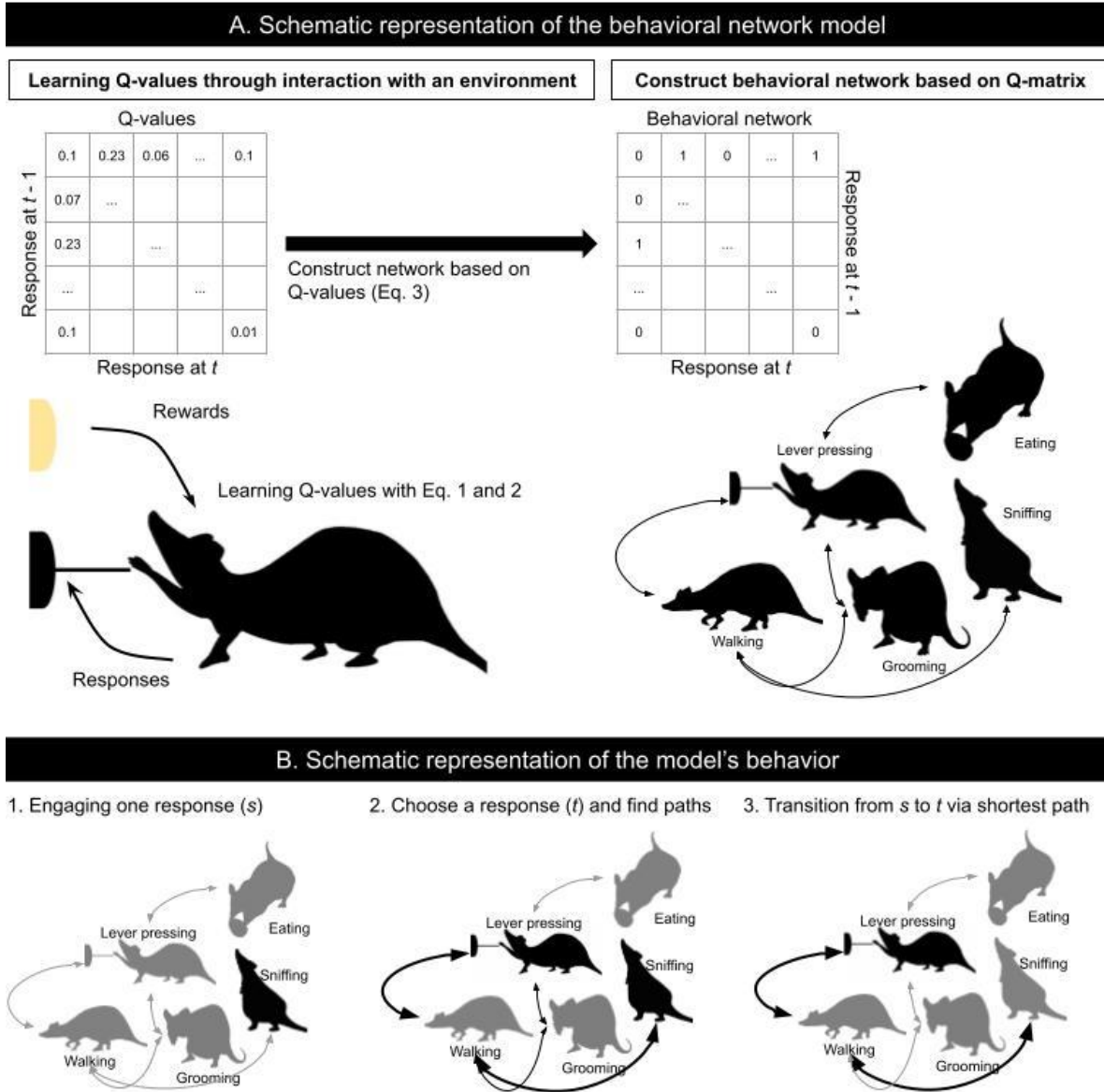
134 Network science emerged in the mid to late 1990s and has spread to a wide range
135 of fields. One of the important aspects of network science is handling the structure of the
136 network. For example, in a network in which individual nodes are randomly connected, the
137 distance between each node is large and the information transmission is slow. However, if
138 there is a node called a hub in the network, which has acquired a large number of edges
139 from other nodes, information can be rapidly transmitted through that node. This is like an
140 influencer sending out information on a social networking service, which attracts the
141 attention of a larger number of users and rapidly spreads the information. In this way, the
142 structure of the network is closely related to the behavior of the entire system. We introduce
143 this perspective of network structure to behavioral science. In this view, each response is
144 assumed as a node, and behavior could be captured as a network of interconnected nodes.
145 By doing so, we try to explain existing behavioral phenomena from a new perspective of the
146 overall structure of behavior. Introducing the concept of network science to experimental
147 analysis of behavior and the theory of habit formation has not been focused on so far.

148 Here, we provide a computational formulation of the behavioral network and explain
149 habit formation from the viewpoint of changes in the network structure. In simulation 1, we
150 generated an arbitrary network and examined what kind of structure forms habit and showed
151 that habit formation occurs when edges are concentrated on a specific response. In
152 Simulation 2, we examined whether the factors reported to promote or inhibit habit formation
153 from existing behavioral studies have similar effects on the proposed model. There are three
154 important factors on habit formation: 1) the amount of training^{11,12}, 2) the schedule of rewards¹³,
155 and 3) the presence or absence of choice^{14,15}. The effects of these factors on the proposed
156 model were consistent with the existing experimental results. These results imply that habit
157 formation can be explained not by the control of the two systems, but by a single system
158 constituting the change in the structure of the behavioral network. Furthermore, the results
159 demonstrate that all responses are goal-directed, rather than the conventional dichotomy of
160 goal-directed and habitual behaviors.

161

162 Results

163 We considered the behavior of an agent as a network consisting of different categories of
164 responses (e.g., lever pressing, grooming, stretching, etc.). Each response was assumed to
165 be a node, and the transition between responses was assumed to be an edge (Figure 1A).
166 The purpose of our agent was the same as the normal reinforcement learning setting of
167 reward maximization. To achieve it, the agent's behavior was modeled by choices based on
168 the values of rewards and the shortest path from the currently engaging response to the
169 chosen response. Although this modeling differed from the ordinary setting, it accounted for
170 the behavior of organisms in the natural environment. Our model reflected three facts
171 (Figure 1B). (1) Most organisms, including humans, engage in various responses in their
172 lives. For example, a rat in a free-operant experiment presses a lever in one moment and
173 grooms its hair or explores the experimental apparatus the next moment. (2) The responses
174 are associated with different types of rewards. Lever pressing is associated with food
175 presentation. Hair grooming is associated with removing disconformity. Exploring within the
176 apparatus is associated with escaping from the apparatus. (3) When an animal shifts from
177 the currently engaging response to another response, it may choose to reach the response
178 via relatively fewer responses. For example, if a rat engages in sniffing (Figure 1B left) and
179 then chooses to press a lever (Figure 1B center), two paths or response sequences are
180 available: walking to the front of the lever and pressing the lever or walking to the front of
181 the lever followed by grooming and then pressing the lever (Figure 1B center). Grooming
182 requires additional time and is redundant for pressing the lever. Thus, the rat may choose
183 the shortest path, i.e., walk to the front of the lever and press it (Figure 1B right). In a large
184 behavioral space, random search increases the time required to reach the desired response
185 and does not warrant reaching the desired response. In summary, the agent chooses one
186 available response associated with different rewards and reaches the chosen response by
187 following the shortest path from the currently engaging response. The agent loops through
188 this process in the behavioral network, which is composed of responses.



189

190 Figure 1. Scheme of the behavioral network

191 A. The schematic representation of the behavioral network model represents how agents
 192 learn the Q-values by interacting with the environment and generate a behavioral network
 193 based on these values. The behavioral network consists of multiple responses. B. The
 194 schematic representation of the model's behavior shows how the agents transit in the
 195 network. The left panel shows the initial state in which agents engage in a response. The
 196 center panel shows that agents choose a goal and search for the shortest path. The right
 197 panel shows that agents transit from the initial response to the goal via the shortest path.

198 We assumed that how nodes in a network and attachment of an edge between two
199 nodes depended on the history of past rewards experienced by the agent. We employed Q-
200 learning³⁵ to represent the history of rewards obtained when transitioning from one response
201 to another. In ordinary Q-learning, an agent learns the action-value in a state. However,
202 since our model dealt with transitions between responses, we treated the response of the
203 agent as a state. Thus, Q-learning in our model was represented by the following equation,
204 assigning the response a time point prior to the state:

$$205 \quad Q(a_{t-1}, a_t) \leftarrow Q(a_{t-1}, a_t) + \alpha \cdot \delta \quad (1)$$

206 In this equation, α denotes the learning rate, we set $\alpha = 0.1$ for all simulations; and δ is the
207 reward prediction error (or temporal difference error). The reward prediction error was
208 calculated as follows:

$$209 \quad \delta = R(a_t) + \gamma \cdot \max_{a_{t+1}} Q(a_t, a_{t+1}) - Q(a_{t-1}, a_t) \quad (2)$$

210 In this equation, γ denotes the discount rate of future rewards and we set $\gamma = 0.5$ for all
211 simulations. R_t denotes the reward obtained by a transition, and the reward functions are
212 different between simulations, which have been explained in detail in the Materials and
213 Methods section.

214 The probability that an edge is attached between any two nodes depends on the Q-
215 value and is calculated using the softmax function. The probability was calculated using the
216 following equation:

$$217 \quad p_{i,j} = \frac{e^{-\beta_n Q(i,j)}}{\sum_{j=1}^N e^{-\beta_n Q(i,j)}} \quad (3)$$

218 In this equation, N denotes the number of nodes in the network and all the responses that
219 the agent can engage in. β_n denotes the inverse temperature and we set $\beta_n = 50$ in all
220 simulations. We also sampled two edges according to Equation 3, such that every node had
221 at least two edges. We used “networkx,” a Python library for network analysis, to generate
222 the network.

223 The algorithm for the agent to choose a response contains two steps: 1) choice of
224 the response based on the value of the reward, and 2) searching the shortest path from the
225 current engaging response to the chosen node. In the choice of the response based on the
226 value of the reward, the probability of choosing a response is calculated by proportional
227 allocation of the reward value. The shortest path search includes selecting the shortest path

228 between the current response to the chosen response and the agent engaging in the
229 responses containing the path in sequence.

230 The probability of response i was calculated according to the following equation:

$$231 \quad p_i = \frac{r_i}{\sum_{j=1}^N r_j} \quad (4)$$

232 In this equation, r_i denotes the value of the reward obtained from response i . In our
233 simulation, the value of the reward obtained from the operant response was 1.0, and the
234 other response was 0.001.

235 The shortest path search is used to find the shortest path between any two nodes in
236 the network. We employed Dijkstra's algorithm³⁶ in all our simulations. If there were multiple
237 shortest paths between any two nodes, we randomly choose one of them. We implemented
238 the path search by using NetworkX³⁷.

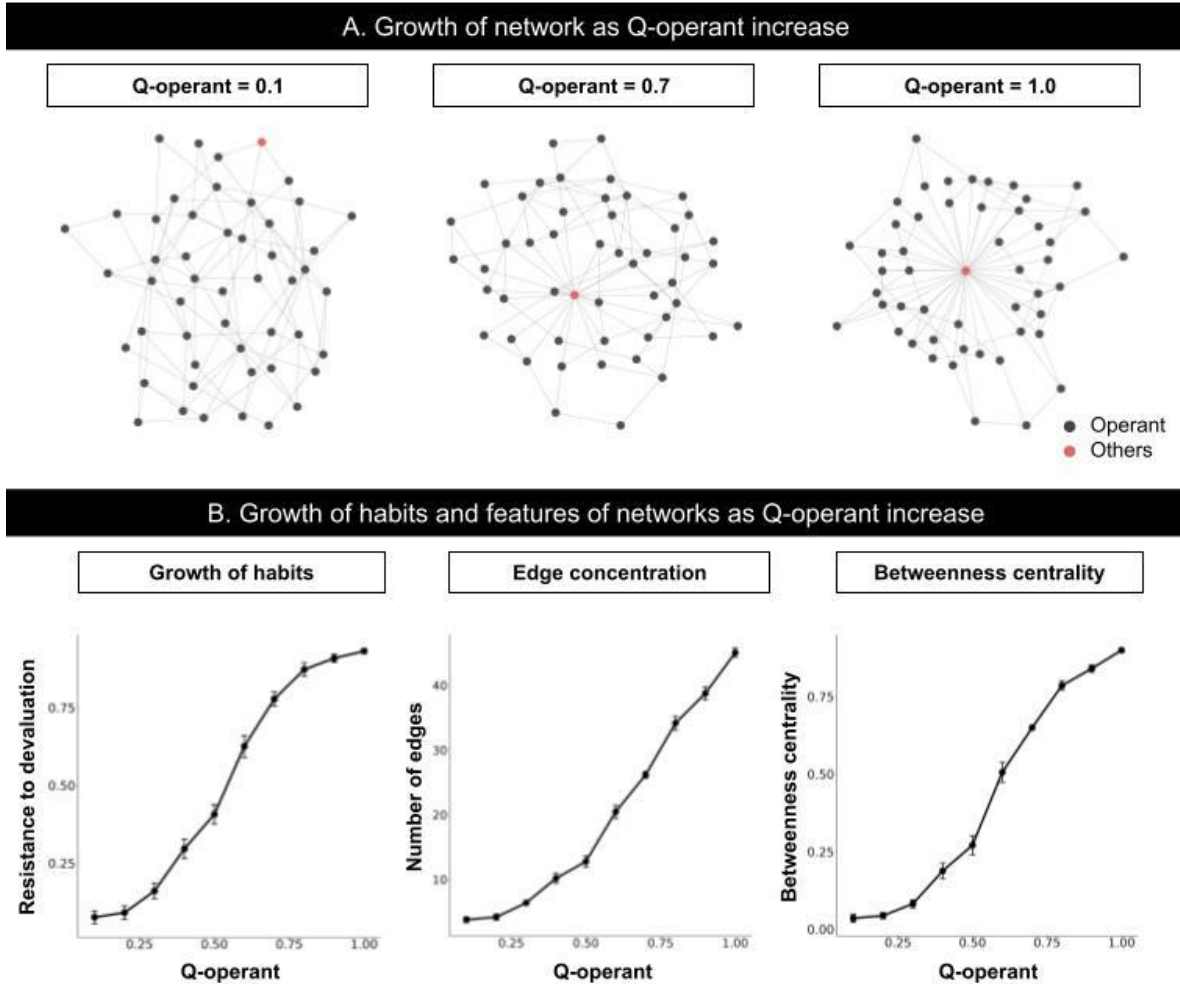
239 Simulation 1: Network structure and habit formation

240 In the Simulation 1, we searched for the structure of the network where habits formation
241 occurs. First, we generate a network based on the Q-matrix. We used an arbitrary Q-matrix
242 to operate the degree of the edge concentration on the operant response. The Q-matrix is
243 defined as the direct product of the Q-vector. The Q-vector contains scalars ranging from 0.
244 - 1. and each element corresponds to each response. More specifically, the first element
245 corresponds to the operant response and others correspond to the other responses. In
246 simulation 1, we fixed the value for the other responses to 0.001 and varied the value for the
247 operant response, Q-operant, from 0.0 - 1.0. To examine the degree of habit, we used the
248 reward devaluation procedure used in free-operant experimental situations. The earliest
249 demonstrations of habit formation¹¹⁻¹³ used the reward devaluation procedure. In this
250 procedure, the investigators train the animals to press the lever with a reward. After the
251 animal learned lever pressings to obtain the reward, the value of reward was reduced by
252 poisoning it with lithium chloride. In this procedure, animals learnt the reward value outside
253 the experiment. Subsequently, investigators examined if the animal pressed the lever
254 without reward deliveries, or an extinction test. Thus, the reward value for the animal was
255 not updated in the test. When the animal pressed the lever, the reward was poisonous, and
256 the responses were considered to be a habit. When the lever-presses decreased after
257 devaluation, the responses were considered to be goal-directed behavior. To reproduce the

258 procedure in the simulation setting, we set up the baseline and devaluation phases where
259 the value of reward obtained by the operant response is 1 and 0, respectively. As animals
260 had experienced reward devaluation outside the experiments in the experimental setting,
261 our agents did not update the reward value within the simulation but changed it from 1.0 to
262 0.0 before starting when moving from baseline to test phases. In both baseline and test
263 phases, the first response that the agent engaged was randomly determined. Then, the
264 agent chooses a response based on the reward value and searches for the shortest path to
265 the response from the current engaging response. They engage in responses contained in
266 the path and the agent reaches the chosen response. After the agent reaches the response,
267 it repeats this process again. After several loops, we calculated the proportion of the operant
268 response to the total number of responses to assess whether the operant response is habit
269 or not.

270 Simulation result

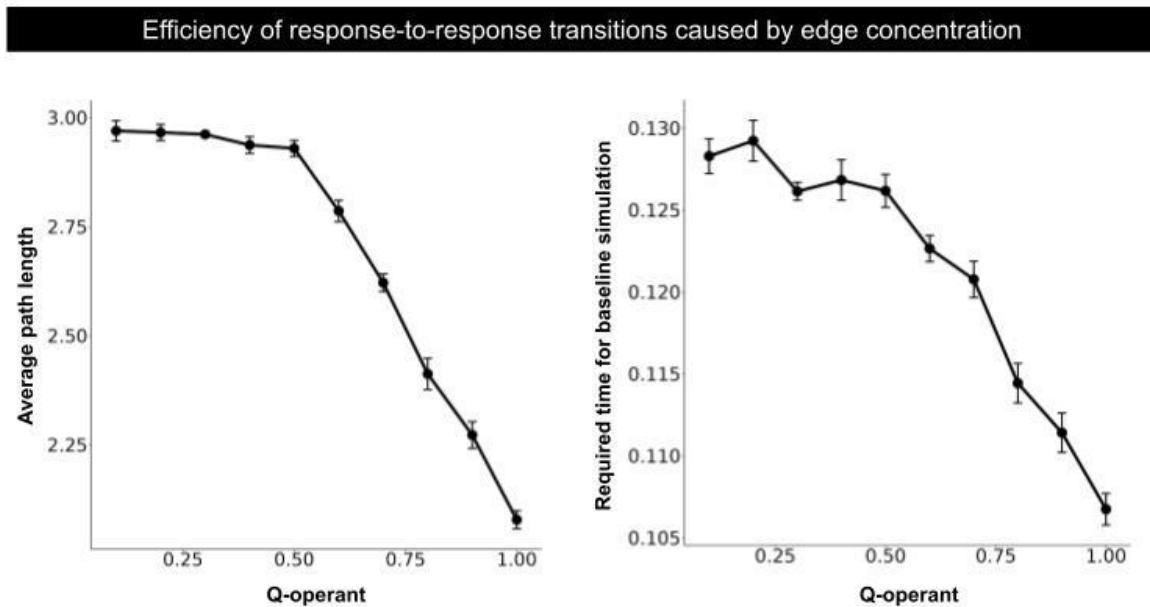
271 Figure 2A shows examples of generated networks under the Q-operant. Other responses
272 (black nodes) connected to the operant response (red node) as the Q-operant increased.
273 Figure 2B shows the resistance to devaluation (left panel), number of edges that the operant
274 response acquired (center panel), and betweenness centrality (right panel). The resistance
275 to devaluation was larger when the operant response did not decrease with reward
276 devaluation and higher Q-operant the resistance to devaluation were larger (Figure 2B left).
277 The number of edges that the operant response acquired increased as the Q-operant
278 increased (Figure 2A and Figure 2B center), implying that edges from other responses were
279 concentrated to the operant response. The betweenness centrality, i.e., the probability that
280 the operant response is included in the shortest path between two nodes in the network,
281 increased as the Q-operant increased (Figure 2B right). With edge concentration in the
282 operant response, distances between two nodes in the network decreased (Figure 3 left).
283 Furthermore, transitions made by agents in the simulations became efficient, and time
284 required for simulations shortened (Figure 3 right). These results were replicated in a wide
285 range of Q-operant and Q-others, in different numbers of nodes (Supplementary figure 2),
286 and with a different path search algorithm (Supplementary figure 3).



287

288 Figure 2. Results of simulation 1

289 (A) Change in network with an increased Q-operant. Each point denotes a response, with
290 black and red indicating other responses and the operant response, respectively. (B)
291 Change in resistance to devaluation and features of the network with an increased Q-
292 operant. The left panel shows resistance to devaluation, which indicates the decrease in the
293 operant response caused by reward devaluation and implies that the operant response
294 becomes a habit at higher values. The center panel shows the change in the number of
295 edges that the operant response acquired. The right panel shows the betweenness centrality,
296 i.e., the probability that the operant response is included in the shortest path connecting two
297 nodes in the generated network.



298

299 Figure 3. Reduced computation costs with habit formation

300 The left panel shows the average path length, i.e., the average of the shortest path between
301 two nodes in the network. When the path length is shorter, the transition from one response
302 to another becomes faster. The right panel shows the required time to simulate the baseline
303 phase. The required time is the real time, i.e., the duration from the start to the end of the
304 simulation. Since the number of loops is the same for all simulations, the decrease in
305 required time implies efficiency in shortest path search and transitions between responses.

306 Interim Discussion

307 In simulation 1, we examined the structure of the network and habit formation under arbitrary
308 Q-matrix and showed that habit formation occurred when edges from other responses were
309 concentrated in the operant response. By manipulating Q_{operant} systematically, the operant
310 response acquired most edges in the network (Figure 2A and Figure 2B center) and it
311 caused that increase in the resistance to devaluation (Figure 2B left). These results suggest
312 that habit formation can be viewed as the structural change in the behavioral network. In
313 particular, habits are considered as concentration of edges from other responses to the
314 operant response. This is because when agents move one response to another, the operant
315 response is included in the path between the two nodes (Figure 2B right). These results

316 were replicated in different settings of algorithms or parameters (Supplementary figure 1, 2,
317 and 3), suggesting these results were not limited to the specific setting.

318 Habits are efficient in the computational cost and transition^{7, 38}. In our model, these
319 features of habits were also found. Animal responses are constrained by some factors, such
320 as space and the animal's body. For example, an animal cannot eat food if the food is not
321 in front of it and if it cannot walk when it is sleeping. These examples imply that not all
322 responses are connected to each other and that the number of edges in the network is
323 limited. When the number of edges was constrained, the structure of the network promoted
324 that agent to engage in the desired response. When edges from other responses were
325 concentrated in the operant response, the average distance between two nodes was
326 shortened³⁸, and transitions made by agents became efficient (Figure 3). These results also
327 imply that agents can find the path between two nodes faster. Thus, habit formation, i.e.,
328 edge concentration to the response, reduces the computational cost and hastens the
329 transition under constraints.

330 Simulation 2: Devaluation and its effect on behavior under free-operant 331 situation

332 We examined if our model could reproduce the effects of factors that promote or disrupt
333 habit formation in free-operant situations¹¹⁻¹⁵. In simulation 2, we let an agent learn Q-values
334 under arbitrary experimental environments and examine whether habit formation occurs.
335 Under free-operant situations, there are three factors that lead to an operant response to
336 habit. The first is the amount of training, where one response is rewarded repeatedly under
337 one situation, and the response becomes habit^{11, 12}. The second factor is the rule, called
338 schedule of reinforcement, which determines the criteria for presentation of a reward for a
339 response³. Habit formation does not occur when reward presentation is determined by the
340 number of responses by the animal. In this environment, the presence/absence of a reward
341 is determined with a certain probability each time the animal presses a lever, e.g., in the
342 bandit task or slot machine use. Habit formation occurs when rewards are determined
343 according to the time elapsed since the previous reward. In this environment, the availability
344 of a reward is determined potentially at arbitrary time steps with a certain probability, and
345 the reward is presented at the first response after reward presentation becomes possible,

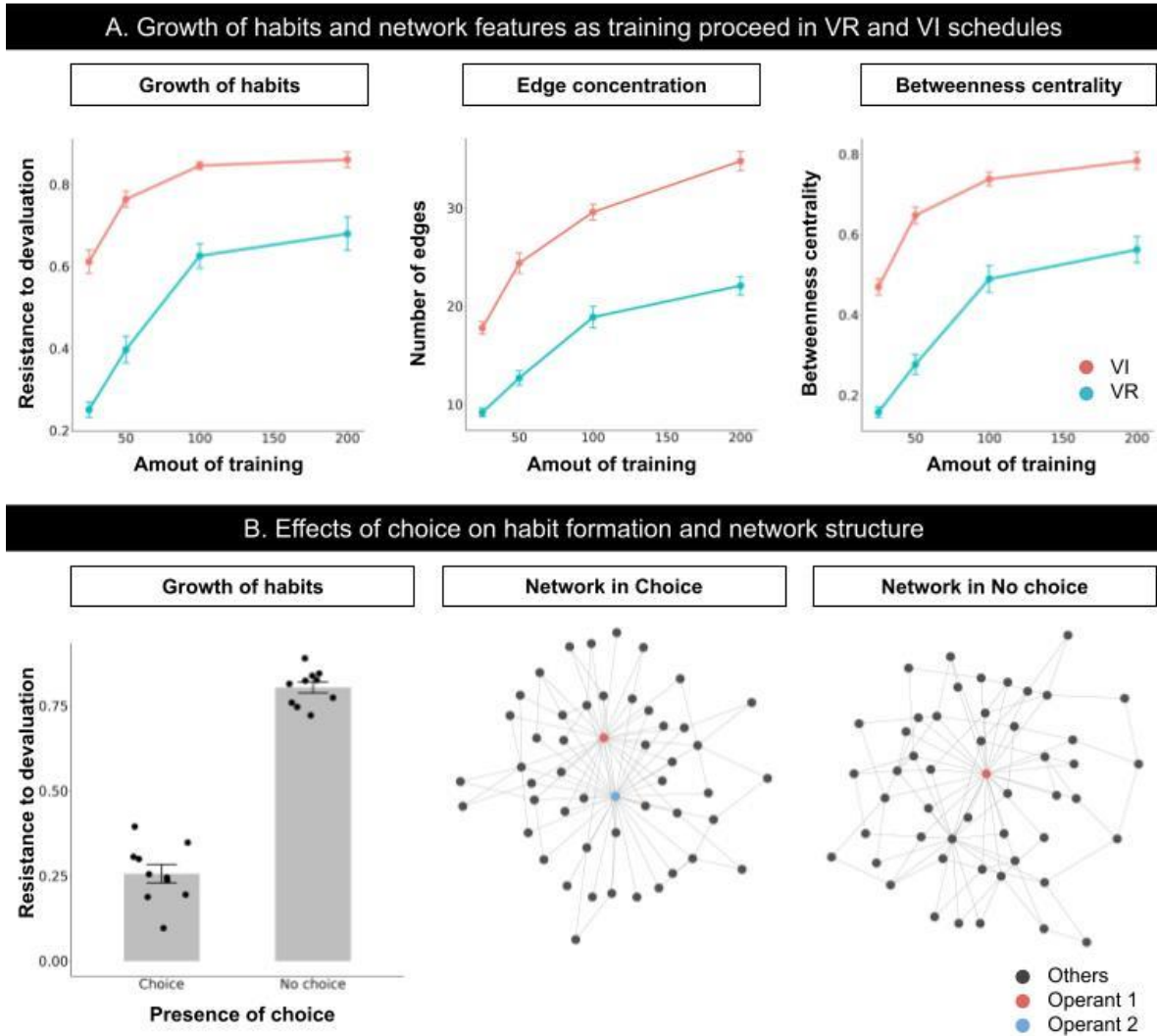
346 such as checking a mailbox. The former response-based rule is called the variable ratio (VR)
347 schedule, and the latter time-based rule is called the variable interval (VI) schedule. The
348 third factor is the presence of alternatives. If two alternatives are available under a situation
349 and different rewards are obtained from them (e.g., left lever → food, right lever → water),
350 the operant response does not become a habit^{14, 15}. Here, we reproduce the above
351 experimental settings and examine whether our model becomes a habit under these
352 environments.

353 The only difference between simulations 1 and 2 is whether the agent learns the Q-
354 values. Here, the agent experienced the training phase preceding the baseline phase, where
355 the agent learned Q-values through interaction with a given environment and constructed a
356 network based on them (more detail in Materials and Methods). After the training phase, the
357 agent experienced the baseline and devaluation phases in the same way as in Simulation
358 1.

359 Simulation result

360 Figure 4A shows the growth of resistance to devaluation (left), number of edges (center),
361 and betweenness centrality (right) with increased amounts of training in VI (time-based rule;
362 red line) and VR (response-based rule; blue line) schedules. All measures were larger in the
363 VI schedule than in the VR schedule. Figure 4B shows the resistance to devaluation (left)
364 and examples of networks learned in the choice (center) and no-choice situations (right).
365 The resistance to devaluation was larger in the no-choice situation than in the choice
366 situation (Figure 4B left). Two operant responses acquired almost the same number of
367 edges in the choice situation (Figure 4B center), while only one operant response acquired
368 the greatest number of edges in the network in the no-choice situation (Figure 4B right).
369 Figure 5 shows the Q-value for self-transition of the operant response. The Q-value
370 increased with an increased amount of training and was larger in the VR schedule than in
371 the VI schedule. These results were replicated in different experimental settings.
372 Supplementary Figure 2B shows the replicated results in different numbers of nodes (25, 50,
373 75, and 100). In simulation 2, agents received rewards every time they engaged in other
374 responses. In other words, we assigned fixed ratio (FR) 1 for other responses.
375 Supplementary Figure 4 shows the results when a different schedule was assigned to other
376 responses instead of FR 1. The results were almost the same. We examined if the results

377 remained similar when a different learning algorithm, SARSA, was employed and
378 Supplementary Figure 5 shows that similar results were obtained.



379

380 Figure 4. Results of simulations in VI and VR schedules and presence and absence of
381 choice

382 (A) Results of simulations manipulating the amount of training in the VI and VR schedules.

383 In all panels, the red and blue lines denote the VI and VR schedules, respectively. The left,

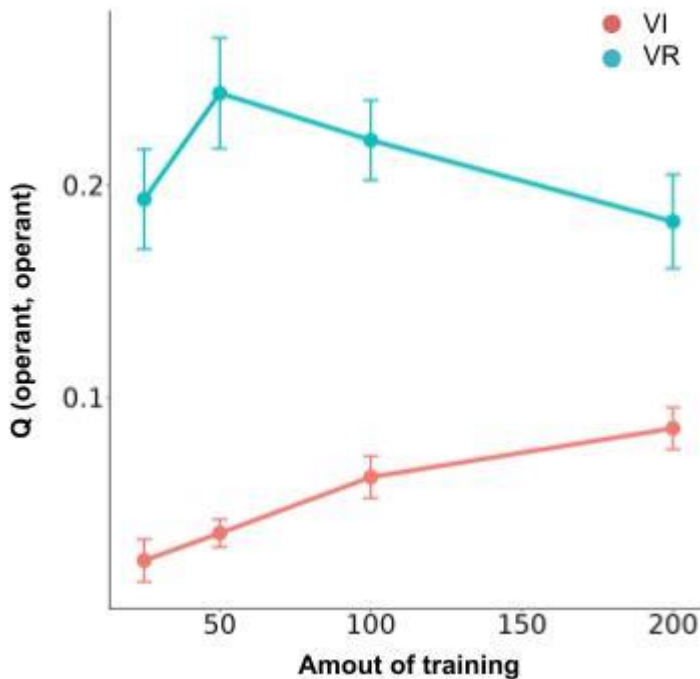
384 center, and right panels show the resistance to devaluation, number of edges, and

385 betweenness centrality, respectively. (B) Results of simulations in the choice and no-choice

386 simulations. The left panel shows the resistance to devaluation. The center and right panels

387 show the learned network in the choice and no-choice situations, respectively. In the network,

388 the red and blue nodes denote the operant response, and black nodes denote other
389 responses.



390

391 Figure 5. Q-value for self-transition of the operant response

392 Q-value of self-transition of the operant response. The red and blue lines denote the VI and
393 VR schedules, respectively.

394 Interim Discussion

395 In simulation 2, we examined whether our model shows similar behavior to real animals in
396 environments that affect habit formation, and our model reproduced the similar results
397 reported from the empirical studies. The resistance to devaluation increased with an
398 increased amount of training and was larger in the VI schedule than in the VR schedule
399 (Figure 4A left). As we have seen in simulation 1, the operant response acquired most of
400 the edges in the network under VI schedule, but not under VR schedule (Figure 4A center),
401 and it turned out that the betweenness centrality grew up under VI schedule (Figure 4A right).
402 These results imply that the VI schedule and a large amount of training promote habit
403 formation. The resistance to devaluation was lower in the choice situation than in the no-

404 choice situation (Figure 4B left), suggesting that the presence of explicit alternatives
405 disturbed habit formation.

406 The amount of training affects the structure of the network (Figure 4A), and as the
407 amount of training increases, the cohesion of edges in the operant response increases. The
408 smaller the amount of training, the smaller the Q-values of the transition from other
409 responses to the operant response. Consequently, the probability that an edge is attached
410 to the operant response is smaller. As shown in simulation 1, habit formation occurs when
411 the operant response acquires most of the edges in the network. Thus, the amount of
412 training affects habit formation.

413 The resistance to devaluation was larger in the VI schedule than in the VR schedule,
414 suggesting that habit formation was promoted in the VI schedule. The VR schedule is a
415 response-based rule of reward presentation. Therefore, all operant responses, independent
416 of the agent's engagement immediately before, were rewarded with constant probability. In
417 contrast, the VI schedule is a time-based rule and it causes that an operant response, longer
418 elapsed time from last operant response, is selectively rewarded. In other words, an operant
419 response emitted after a few periods was selectively rewarded and implied a transition from
420 the other responses to the operant response in our model. In summary, transitioning from
421 other responses to the operant response was selectively rewarded in the VI schedule and
422 resulted in edge concentration in the operant response and habit formation.

423 One might suspect that, contrary to the experimental facts that the response rate is
424 larger in VR schedule than VI schedule, if operant responses acquire more edge in the VI
425 schedule, then the response rate would be higher in the VI schedule as well. However,
426 Figure 5 shows the Q value of the self-transition of the operant response is larger in VR
427 schedule than VI schedule. It implies that once an agent starts to engage in an operant
428 response, it will repeat the same response over and over again. In fact, it has been
429 experimentally shown that the difference in response rate between VI and VR schedules is
430 caused by such a mechanism⁴⁰⁻⁴².

431 Although the operant response acquired most of the edges on the network under the
432 choice environment, the operant response did not become a habit. There are two reasons
433 for this. First, the agent chooses its response based on the value of the reward obtained
434 from the response. In the test phase, the value of the reward obtained from the operant
435 response was reduced, and that of the alternative response remained the same value as

436 the baseline. Thus, the agent chose the alternative response more in the test phase than in
437 the baseline phase. Second, if only the operant response acquired most edges, any shortest
438 path may contain the operant response. However, the alternative response acquired most
439 of the edges, so that any shortest path contained the alternative response. Thus, the operant
440 response no longer has a greater chance of being engaged, and habit formation does not
441 occur.

442 In the no-choice situation, the operant response acquired the most edges in the
443 network, but several other responses also acquired multiple edges (Figure 4B right),
444 resembling the scale-free network, which should be assessed by the distribution of degree.
445 However, habit formation occurred in the network. Therefore, although scale-free networks
446 were not compared with random or hub-and-spoke networks, habit formation might be
447 present in the scale-free-like network.

448 Simulation 3: Correlation-based account vs contiguity-based account of 449 habit formation

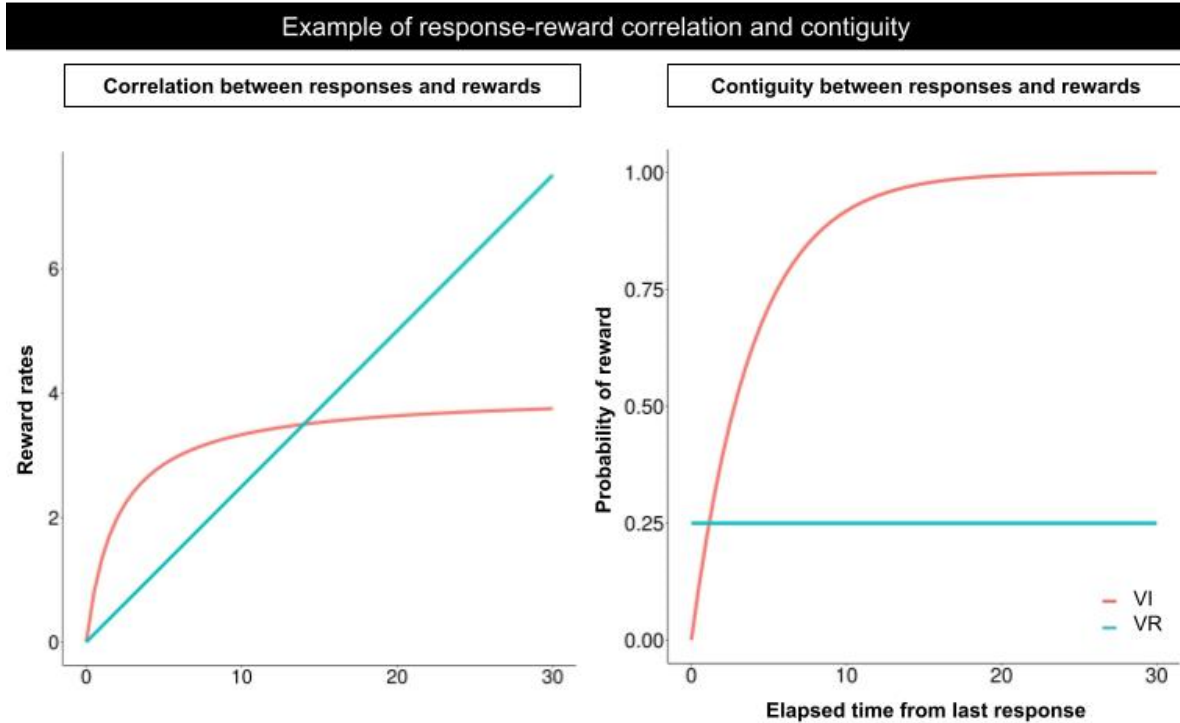
450 Here, we propose an experiment to directly test the response-reinforcer correlation, which
451 has been considered as a factor leading to habit formation in the past, and our model's
452 explanation: selective reinforcement of transitions from other behaviors to the operant
453 response and the resulting structural changes in the network. This is a new experiment
454 predicted by our model, which has not yet been examined in real animals, and will encourage
455 future theoretical tests.

456 From canonical view, response-reward correlation, the operant responses remain
457 goal-directed when animals experience a correlation between the operant responses and
458 rewards but become habits when they do not experience the correlation^{1, 43}. Under VR
459 schedule, the more they engage the operant response, the more rewards they can obtain.
460 It leads that they experience positive correlation between the operant response and rewards,
461 and the operant response remains goal-directed. In contrast, under VI schedule, since
462 rewards availability is governed by time, such correlation is collapsed, and they do not
463 experience it. It results that the operant response becomes habit.

464 In recent years, results have been reported that contradict the response-reward
465 correlation⁴⁴⁻⁴⁶. For example, De Russo, et al.⁴⁵ trained mice under VI and FI schedules. FI

466 and VI have a common molar relationship between response rate and rewards: in both
467 schedules, animals cannot obtain more than the determined number of rewards within a
468 certain duration, no matter how much they engage in the operant response. Under such a
469 condition, the response-reward correlation view predicts that both schedules guide the same
470 level of habit formation. However, the operant response of mice trained under FI schedule
471 remains goal-directed but under VI schedules, the operant response becomes habit.
472 DeRusso, et al.⁴⁵ conclude that the contiguity, which is defined by average temporal distance
473 between responses and successive rewards, disrupts habit formation. In the FI schedule,
474 animals tend to emit more response as they approach the time when rewards are presented.
475 In contrast, animals do not know when the reward becomes available, they emit responses
476 uniformly during inter-reward intervals in VI schedule. Thus, under the FI schedule, animals
477 emit many responses just before rewards and the contiguity of responses and rewards
478 becomes higher but, under the VI schedule, operant responses are distributed uniformly,
479 and the contiguity becomes lower.

480 A similar discussion has been made for VI-VR response rate difference and there
481 are two kinds of accounts. One explains the difference by the difference in interresponse
482 time that is likely to be rewarded^{47, 48}. In VI schedule, probability of reward availability
483 increases as the elapsed time from last response increases and it results that longer IRTs
484 are more likely to be rewarded than shorter ones. In contrast, such characteristics are not
485 found in the VR schedule or shorter IRTs are more likely to be reinforced. (Figure 6 right).
486 Thus, response rate is lower in VI schedule than VR schedule. Especially, the copyist model⁴⁷
487 explains the difference by average of inter-response times between successive rewards and
488 this is similar to contiguity-based account of habit formation^{45, 46}. Second account is based on
489 the molar relationship between response rate and reward rate^{49, 50}. The more animals emit
490 responses under VR schedule, the more rewards they can obtain (blue line in Figure 6 left).
491 In contrast, under VI schedule, animals cannot obtain more rewards than experimentally
492 defined, no matter how they emit responses under the schedule (red line in Figure 6 left).
493 This account underlies the response-reward correlation account of habit formation^{1, 43}.



494

495 Figure 6. Response rate and reward rate correlation (left) and reward probability as function
496 of elapsed time last response (right) in VR and VI schedules. In VR schedule (black line),
497 reward rate is proportional to response rate, in contrast, reward rates reach a plateau as
498 response rate increases in VI schedule (red line). Reward probability is constant
499 independent from elapsed time from last response in VR schedule, in contrast, it increases
500 exponentially as the time increases.

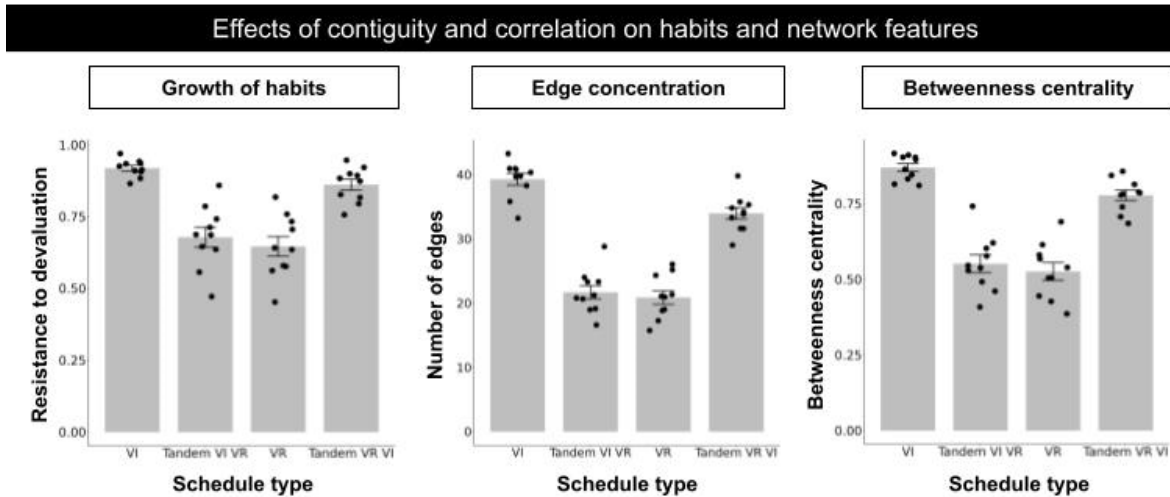
501 Our model is positioned similarly to the contiguity-based account in these
502 discussions. As we show in simulation 2, the VI-VR response rate difference can be
503 explained by which transitions are likely to be rewarded: In VI schedule, the transitions from
504 other responses to the operant response are more likely to be rewarded but not in VR
505 schedule (Figure 5). Viewing the cause of long IRTs as engagement in other responses^{33, 51},
506 differential reinforcement of long IRTs can be interpreted as differential reinforcement of the
507 transition from other response to the operant response. Considering these discussions, our
508 model suggests that the same discussions for VI-VR response rate difference can be applied
509 to habit formation.

510 Here, we mimic an experiment which is conducted to reveal that the VI-VR response
511 rate difference is caused by IRTs immediately followed by rewards⁵². In the experiment,
512 pigeons are trained under tandem VI VR and tandem VR VI schedules. The former schedule,
513 tandem VI VR, shares a molar relationship between response rate and reward rate with VI
514 schedule. However, VI schedule is immediately followed by short VR schedule and longer
515 IRTs are less likely to be rewarded than simple VI schedule. The later one is tandem VR VI,
516 it's molar relationship between response rate and reward rate is similar to the simple VR
517 schedule. However, since VR schedule is followed by VI schedule, longer IRTs are more
518 likely to be rewarded. In this schedule, pigeons showed higher response rate in tandem VI
519 VR schedule and lower in tandem VR VI schedule⁵². These findings contradict the account
520 based on response rate and reward rate correlation but well explained by differential
521 reinforcement of IRTs⁵⁷. Will habit formation occur under these schedules? From the view of
522 response-reward correlation, tandem VI VR schedule leads habit but not in tandem VR VI
523 schedule because there is lower response-reward correlation under the former schedule but
524 higher than the later one. In contrast, our model makes the opposite prediction that habit
525 formation will be guided under tandem VR VI schedule but not under tandem VI VR schedule.
526 This is because, in the former schedule, transitions from other responses to the operant
527 response are more likely to be rewarded, and the operant response acquired more edges.
528 In the later schedule, transitions from other response to the operant response and the self-
529 transition of the operant response are rewarded in the same probability so the operant
530 response acquired not so many edges.

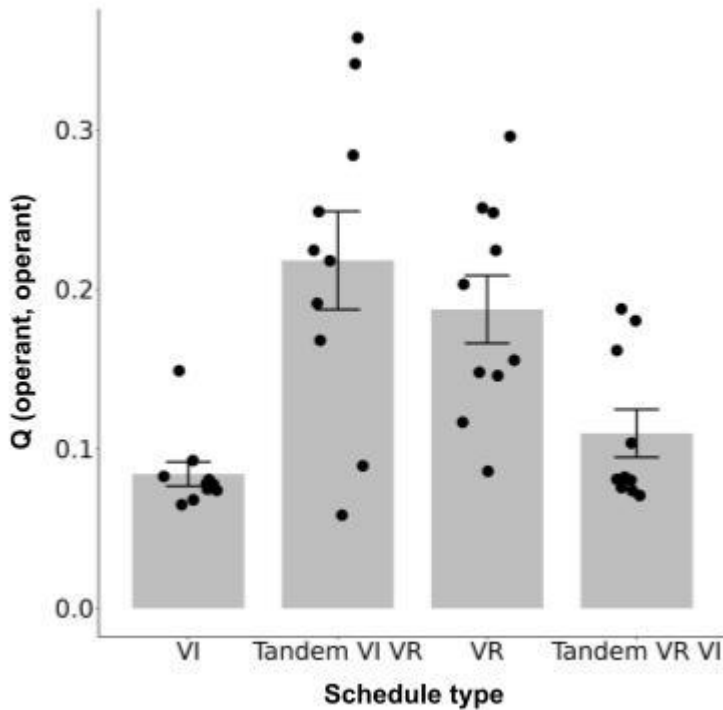
531 Simulation result

532 Figure 7 shows the resistance to devaluation, number of edges, and betweenness
533 centrality simulated under VI, tandem VI VR, VR, and tandem VR VI schedules. They were
534 higher under VI and tandem VR VI schedules than VR and tandem VI VR schedules.
535 Although the response-reward correlation account suggests that habit formation is disrupted
536 under tandem VR VI schedule and is promoted tandem VI VR schedule, the results were
537 the opposite, habit formation was promoted under tandem VR VI and but not under tandem
538 VI VR. The center of Figure 7 shows the number of edges that the operant response
539 acquired to the overall number of edges in the network and the operant response acquired
540 more edges under VI and tandem VR VI schedules. Figure 7 (right panel) shows the

541 betweenness centrality of the operant response. The betweenness centrality was larger in
542 the VI and tandem VR VI schedules than in the VR and tandem VI VR schedules. Figure 8
543 shows the Q-value of the operant response. It was larger in the VR and tandem VI VR
544 schedules than in the VI and tandem VR VI schedules.



545
546 Figure 7. The simulation results in tandem VI VR and tandem VR VI schedules.
547 The left panel shows the resistance to devaluation in VI, Tandem VI VR, VR, and Tandem
548 VR VI schedules. The center panel shows the number of edges that the operant response
549 acquired in each schedule. The right panel shows the centrality of the operant response.



550

551 Figure 8. Q-value for self-transition of the operant response.

552 Interim Discussion

553 In simulation 3, we mimicked the schedules employed by Peele et al.⁵² to reveal what
554 characteristics of schedules, response-reward correlation, or response reward contiguity,
555 promote habit formation. Traditional accounts suggest that lack of the response-reward
556 correlation promotes habit formation^{1, 43}. In contrast, other researchers suggest that the
557 response-reward contiguity is crucial for habit formation but not the correlation^{45, 46}. These two
558 accounts make different predictions in the schedules we employed here. Tandem VR VI
559 schedule has a common molar relationship between response rate and reward rate with
560 simple VR schedule (blue line in Figure 6 left) but it also has a time-dependent property,
561 which is found in VI schedule (red line in Figure 6 right), that the probability of obtaining
562 rewards increases as time elapses. In summary, the Tandem VR VI schedule has higher
563 response-reward correlation but lower response-reward contiguity, and the response-
564 reward correlation account predicts that habit formation is disrupted in the schedule. In
565 contrast, the tandem VI VR schedule lacks both a molar relationship between response rate

566 and reward rates and time-dependency (red line in Figure 6 left and blue line in Figure 6
567 right). In such schedule, animals cannot obtain more than the determined number of rewards
568 within a certain duration, no matter how much they engage in the operant response but the
569 transition from other responses to the operant response is less likely rewarded. In summary,
570 the tandem VR VI schedule had a higher response–reward correlation but a lower response-
571 reward contiguity, and the response-reward correlation account predicted that habit
572 formation was disrupted in the schedule. However, in contrast to the traditional view, our
573 model predicts that habit formation is more likely promoted in tandem VR VI schedule
574 (Figure 7 left). Because of time-dependency tandem VR VI schedule have, transition from
575 others response to the operant response is more likely to be rewarded in the schedule and
576 acquired more edges that simple VR schedule and tandem VI VR schedule (Figure 7 center).
577 Thus, as we showed in simulation 1 and 2, the probability that the operant response is
578 included in the shortest paths increased and habit formation occurred (Figure 7 right).

579 Our model supports the account that the contiguity between responses and rewards
580 promotes habit formation^{45, 46}. In tandem VI VR and simple VR schedule, the self-transition of
581 the operant response is more likely rewarded than transition from other responses to the
582 operant response. This is because, the operant response occurred as a bout, a burst of
583 responses is followed by long pauses, and this implies that animals emit more responses
584 just before reward presentation. In tandem VR VI and simple VI schedule, the self-transition
585 of the operant response is less likely rewarded because of the time-dependent property
586 between response and reward (red line in Figure 6 right). This result implies that animals
587 emit less response just before the reward presentation. Thus, response reward contiguity is
588 higher in the tandem VI VR and simple VR schedule than tandem VR VI and simple VI
589 schedule.

590

591 Discussion

592 In this research, we explain habit formation as changes in network structure by assuming
593 the behavior of organisms viewed as a network of responses. In simulation 1, we generated
594 arbitrary networks and examined the underlying structure of goal-directed behavior and
595 habits. We revealed that habit formation occurs when a particular response acquires most
596 of the edges from other responses. In Simulation 2, we simulated the environments that
597 were reported to promote or inhibit habit formation from existing studies and examined
598 whether the proposed model showed habit formation. These results were consistent with
599 experimental results reported by many laboratories, suggesting that our results demonstrate
600 habit formation as a structural change in the behavioral network. In simulation 3, we
601 analyzed the behavior of the proposed model in an experimental situation where the
602 canonical theory^{1, 43} and the proposed model make different predictions. The results suggest
603 that our model supports the view of reward-responses contiguity promoting habit formation⁴⁵.
604 ⁴⁶ but not the canonical view of reward-response correlation.

605 Relationship to other theoretical models of habit formation

606 Although there are many models of habit formation, most of them are viewed as goal-
607 directed behavior and habits as interactions between two distinctive associative structures.
608 Here, we succeeded in providing a novel explanation by taking a more molar view of
609 behavior. Specifically, the proposed model substantially differs from existing models in three
610 ways. First, the proposed model does not consider behavior as a single element, but as a
611 network of interconnected responses. Conventional views focus only on responses under
612 highly constrained experimental situations, such as lever pressings or button pushings, and
613 ignore the molar structure of behavior that the real organisms may have. Responses of
614 organisms, including humans, are not independent of each other, but they are
615 probabilistically conditioned by the preceding and succeeding responses. In the proposed
616 model, the structures of such responses are represented as a network, and habit formation
617 is explained as a change in the structure. Second, our model seems to have no state variable,
618 unlike previous models^{3, 4, 7, 53}. We treated the immediately prior response of the agent as a
619 state; thus, so there is no lack of state variables. This treatment of past responses as states
620 has often been employed in modeling animal behavior^{32, 33, 51, 54}. However, our model differs from

621 past models of habits. Many models of habits were built in consideration of the multistage
622 Markov decision task^{2,4, 7}. In the multistage Markov decision task, experimentally explicit
623 states, each choice point, exist. In contrast, we studied habits in free-operant situations in
624 which animals could engage in responses freely and repeatedly, and experimentally explicit
625 states were lacking. Previous models were applied to the free-operant situation in two
626 different ways. One way was to not assume the state (1), and the other way was to introduce
627 a hypothetical state^{7, 53}. We treated the immediately prior response as a state, similar to the
628 later one. Although our model seems to have no state variable, our approach was similar to
629 the previous one^{7, 53}. Third, some models of habits assumed two distinct systems
630 corresponding to goal-directed behavior and habits^{1, 2, 53}. Particularly, only the model that could
631 explain habits in free-operant situations assumed them explicitly (1). Although all responses
632 were assumed to be under goal-directed control, choices were based on reward values and
633 shortest path search, and results reported in free-operant situations were reproduced¹¹⁻¹⁵.
634 Recently, in the context of the multistage Markov decision task, several models showed no
635 distinct systems between goal-directed behavior and habits^{3, 4}. Our model also showed no
636 explicit distinction but that the idea could be applied to habits in free-operant situations.

637 Although the proposed model deals with experiments on habit formation in rodents'
638 operant situations¹¹⁻¹⁵, most of the experiments discussed here are also dealt in Perez and
639 Dickinson¹. Both models reproduce results that are consistent with the experimental results.
640 Perez and Dickison¹ provide an explanation based on reward-response correlations. In their
641 model, the lower the correlation between response and reward, the more habit formation is
642 promoted. On the other hand, the proposed model provides an explanation based on
643 contiguity between response and reward^{45, 46}. Contiguity is defined by the temporal distance
644 between the reward and the emitted response to obtain it. The lower the contiguity, the
645 longer the temporal distance between the response and the reward, the more habit
646 formation is promoted. Although the proposed model does not explicitly incorporate
647 contiguity as a variable in the model, it allows for a similar representation by dividing the
648 agent's behavior into the operant responses and other responses, and separating transitions
649 to the operant responses into self-transitions and transitions from other responses. For
650 example, in a schedule with low reward-response contiguity, such as the VI schedule,
651 transitions from other behaviors to the operant are more likely to be reinforced, while in a
652 VR schedule with high contiguity, transitions from other behaviors are less likely to be

653 reinforced. As a result, the operant response obtains more edges and promotes habit
654 formation in schedules with low contiguity. As an experiment in which these two factors can
655 be more clearly separated, we employed the procedure of Peele et al.⁵². Under this
656 procedure, correlation-based and contiguity-based explanations provide opposite
657 predictions. The proposed model reproduced the same results as predicted by the
658 contiguity-based explanation. Whether habit formation occurs under this experimental
659 procedure has yet to be examined, but it does provide useful insights for updating the theory
660 of habit formation.

661 The proposed model may seem similar to the model of Dezfouli and Balleine^{7, 9, 10}. In
662 fact, their model and our proposed model have two common assumptions. First, instead of
663 treating the agent's behavior as a single response, the two models explicitly assume other
664 responses. They explain habit formation in terms of the acquisition of those sequences or
665 the structure of the network. The second point is that the agent generates sequences or
666 searches for the shortest path based on the value of the reward. However, the models have
667 two differences. First, the targeting experimental situations differed. Their model was built
668 with the multistage Markov decision task, while our model was built to explain habit formation
669 in free-operant situations. The existing comprehensive theory in free-operant situations
670 assumed parallel control by two systems (1). A kind of response-chaining/action-chunking
671 models have limited applicability in free-operant situations. Second, the view of behavior
672 differed. Our model tried to overcome the limitation. In free-operant situations, animals could
673 engage in responses freely without explicit states defined experimentally. In the case of free-
674 operant situations, direct application of the idea of response-chaining or action-chunking
675 was difficult because no points corresponded to the start and end of trials. Instead of the
676 chunk or chain, we considered behavior as a network and the agent's behavior as a
677 transition within the network. In other words, by viewing behavior as a loop without a clear
678 start or end, we successfully modeled the behavior of free-operant situations.

679 Dezfouli and Balline⁷ applied their model to the free operant situation and reproduced
680 the effect of amount of training on habit formation. However, they did not treat how other
681 factors, schedule types and presence of alternatives, affect habit formation. The proposed
682 model, which shares common assumptions with their model, can reproduce the results
683 reported in empirical study¹¹⁻¹⁵, suggesting that the idea of response-chaining or action-
684 chunking could be applied in free-operant situations. Moreover, the model clarifies the

685 difference between the canonical correlation-based account and common points with the
686 contiguity-based account. We also found common features with the recently proposed
687 models^{3,4}. In those models, goal-directed planning was employed, and the behavior of human
688 and rodents' multistage decision-making tasks, such as multistage Markov decision tasks
689 and tree-shaped maze, were explained. Pezzulo et al.³ built a mixed-controller model
690 consisting of goal-directed and habit behaviors in a single system. Keramati et al.⁴ proposed
691 that these two systems were not separated but placed in one spectrum. Our model also
692 considered these two systems to be not separated but coexisting in a single system and
693 placed in one spectrum, with only a difference in the structure of the network. However,
694 similar to many other models, their models targeted multistage decision-making tasks but
695 not free-operant situations. Our model shared common features, i.e., planning and
696 singularity of the system, with their models^{3,4} and successfully applied those features in free-
697 operant situations. From the canonical view, two distinct systems control a response in the
698 flat manner^{1,2}. This view has been challenged recently, and new models have been proposed
699 in the context of the multistage decision-making tasks. Although their applications are limited
700 to free-operant situations, our model adopted those ideas, i.e., response-chaining/action-
701 chunking, planning, and mono-systematicity, and explained habit formation in free-operant
702 situations, suggesting a link between the different experimental procedures and providing a
703 comprehensive understanding of habit formation.

704 Neural substrates of behavioral network

705 The corticostriatal network is involved in habit formation, and generates response patterns⁵⁶.
706 ⁵⁶. Especially, dorsolateral striatum (DLS) is known to be important in transition from goal-
707 directed behavior to habits⁵⁷. DLS activity changes as proceedings of training and responses
708 become habits^{58, 59}, and lesion of DLS turns habits into goal-directed behavior after extended
709 training⁵⁷. DLS also carries forming response sequences⁶⁰ and motor routines⁶¹. In addition to
710 its importance in the learned behavior, DLS also encodes innate response sequences⁶².
711 These facts imply that habit formation and the formation of response sequences have
712 common neural substrates.

713 A recent study reported that DLS encodes not only information about response
714 sequences but also more divergent information about behavior, which are topographically
715 categorized responses and transitions between them¹⁸. They recorded the DLS activities of

716 mice with fiber-photometry under an open-field situation and reported neural activities that
717 correlated with the behavior. The activities differed depending on the preceding and
718 succeeding responses, and DLS encoded a transition between the responses. Moreover,
719 the behavior of the mice with DLS lesions showed random transitions of the responses
720 compared to the sham-lesion group. These results imply that the information encoded in
721 DLS is the transition of the structure of behavior. Thus, the function of the DLS might be well
722 understood by considering the habit and goal-directed behavior from the viewpoint of the
723 behavioral network.

724 Corticostriatal circuits, the associative network, which consists of the prefrontal
725 cortex, dorsomedial, or ventral striatum, plays a role in goal-directed behavior⁶³. The
726 dorsomedial striatum (DMS) is known to be involved in the acquisition of goal-directed
727 behavior, maintaining sensitivity to outcomes, and expressing goal-directed behavior^{64, 65}. The
728 DMS receives excitatory inputs from the prefrontal cortex, whereas the DLS receives inputs
729 from the sensorimotor and premotor cortices⁶⁵. In the canonical dichotomous view of habit
730 formation, goal-directed behavior is replaced by habit after extensive training. After habit
731 formation, the contribution of DLS becomes more important than that of DMS^{57, 65}. However,
732 even after extensive training, many brain areas such as the prefrontal cortex, anterior
733 cingulate cortex, and ventral and dorsal striatum are modulated by anticipated rewards⁶⁶⁻⁷⁰. In
734 our model, any response emitted by an agent is considered goal-directed. Regardless of the
735 training stage, our agents choose their responses based on the value of the rewards.
736 Therefore, the fact that regions involving goal-directed behavior are modulated by
737 anticipated rewards even after extensive training, our assumptions do not contradict each
738 other. Combined with the fact that DLS is more responsible for sequential responding than
739 DMS⁷¹, the transition from DMS to DLS during habit formation might reflect the corresponding
740 behavioral sequence induced by changes in the behavioral network.

741 Neuronal circuits involving ventral striatum and hippocampus play key roles in spatial
742 navigation and are considered to be related to the planning^{72, 73}. Both spatial navigation and
743 planning are related to habits, and they share common neurobiological substrates^{3, 74-78}.
744 Although roles of hippocampus and planning in habits and goal-directed behavior in free-
745 operant situations remains unknown, our model sheds light on the role of planning and
746 related brain regions in habits in the free-operant situations.

747

748 Relationship to other behavioral phenomena

749 Animals engage in specific responses, such as orienting, approaching, and consummatory
750 behavior, just after the presentation of the reward. Specific action sequences are observed
751 during experiments, and learning is sometimes disrupted by innate responses. These
752 experimental and observational facts lead us to assume that behavior is a network
753 constructed from responses.

754 In our model, the structure of a network depends only on past experiences under a
755 given situation. In other words, our model does not consider the connections between
756 specific responses that real organisms may have. Thus, we could not reproduce this
757 phenomenon. However, our model can be further extended and modified to include this
758 phenomenon.

759 Schedule-induced behavior, observed under intermittent schedules of reinforcement,
760 is a behavioral phenomenon in which animals show aggression or water intake just after the
761 reward presentation²⁶⁻²⁸. This phenomenon can be attributed to the innate connections
762 between reward consumption and schedule-induced behavior. Because of these
763 connections, animals tend to engage in aggression or water intake immediately after reward
764 presentation. Similarly, terminal behavior, which occurs as approaches reward
765 presentations, can be explained by assuming an innate connection, which may explain the
766 fact that animals show a specific sequence of responses during the experiment.

767 To deal with such phenomena, we assume that it is possible to express the innate
768 susceptibility of edges as a prior distribution and impose constraints on the probabilities of
769 edges attached by learning. Furthermore, we can systematically treat phenomena such as
770 misbehavior and biological constraints on learning by examining differences in prior
771 distributions among species and environments. Thus, we can extend our model to a
772 comprehensive framework of behavior that incorporates the innate behavior of organisms
773 under natural settings.

774 Goal-directed behavior and habits are related to spatial navigation^{3, 74-78}. Pezzulo et
775 al.³ target an experiment with tree type maze and the task is similar in the abstract structures
776 to the multistage Markov decision task. Our model employed a planning process as the
777 model proposed in Pezzullo et al.³. However, planning is made in the real space in their model,
778 but planning is made in behavioral space in ours. Thus, the application of our model for

779 spatial navigation is limited. However, the idea of learning response sequence can be
780 applied to spatial navigation, such as learning a series of responses of turning to left and
781 then turning to the right. As we discussed in the above, the limitation is also related to the
782 experimental situations, multistage Markov decision tasks and free-operant situations. We
783 expect a more comprehensive view or model that targets both experimental situations in the
784 future.

785 Limitations and future directions

786 Our model has three major limitations. First, as we discussed in the previous section, our
787 model does not consider innate constraints that real organisms have, and we believe that
788 we can solve the problem by expressing the innate constraints as a prior distribution. Second,
789 our model could not treat the self-transition of each response. Third, it can only deal with
790 experiments on habit formation under free-operant situations.

791 Our model cannot treat the self-transition of responses because we employed the
792 shortest path search algorithm to generate response sequences. Any self-transition makes
793 paths between any two responses longer, and paths containing self-transitions must not be
794 the shortest paths. However, animals show a particular response pattern, which is called
795 bout-and-pause and characterized by phasic bursts of one response and pauses following
796 them. Such response patterns imply that the responses have self-transitions. To solve this
797 problem, it is necessary to employ a different algorithm to generate response sequences
798 that allow self-transition.

799 All our simulations deal with experiments in free-operant situations, and not with
800 recent experiments with the two-stage Markov decision task. This is not a specific problem
801 for our model; other existing models treat either of them. Although many experiments have
802 been conducted in both experimental tasks, the differences and identities of the procedures
803 and results among them have not been systematically examined. To obtain a more unified
804 understanding of habit formation, we need to conduct a systematic analysis of the
805 procedures and results employed and obtained from existing studies. Therefore, the
806 validation of our model is limited to habit formation in free-operant situations.

807 Recent advances in machine learning allow us to measure animal behavior more
808 objectively and precisely than ever before. However, behavior estimation technologies are
809 not well established at present, preventing us from validating some assumptions in our

810 model. In this field, no consensus has been reached on what timescale should be employed
811 to classify behavior and how finely behavior should be classified. For example, we assumed
812 that the behavioral network consisted of 50 nodes but did not know how many nodes
813 constitute the behavioral network of real animals. However, as shown in Supplementary
814 Figure 2, habit formation occurred in networks of a slightly smaller size, suggesting that our
815 explanation for habits could be applicable to the real behavioral network even if the size is
816 smaller than we assumed. In the future, such technologies and by utilizing these techniques,
817 it is possible to understand behavior on a macroscale rather than capture the behavior in
818 highly constrained experimental settings. Our model provided a novel perspective on how
819 behavior could be viewed on macroscale behavioral phenomena and raised questions that
820 could be answered by such techniques, which would further help us understand the function
821 of the brain in behavioral changes.

822 Conclusion

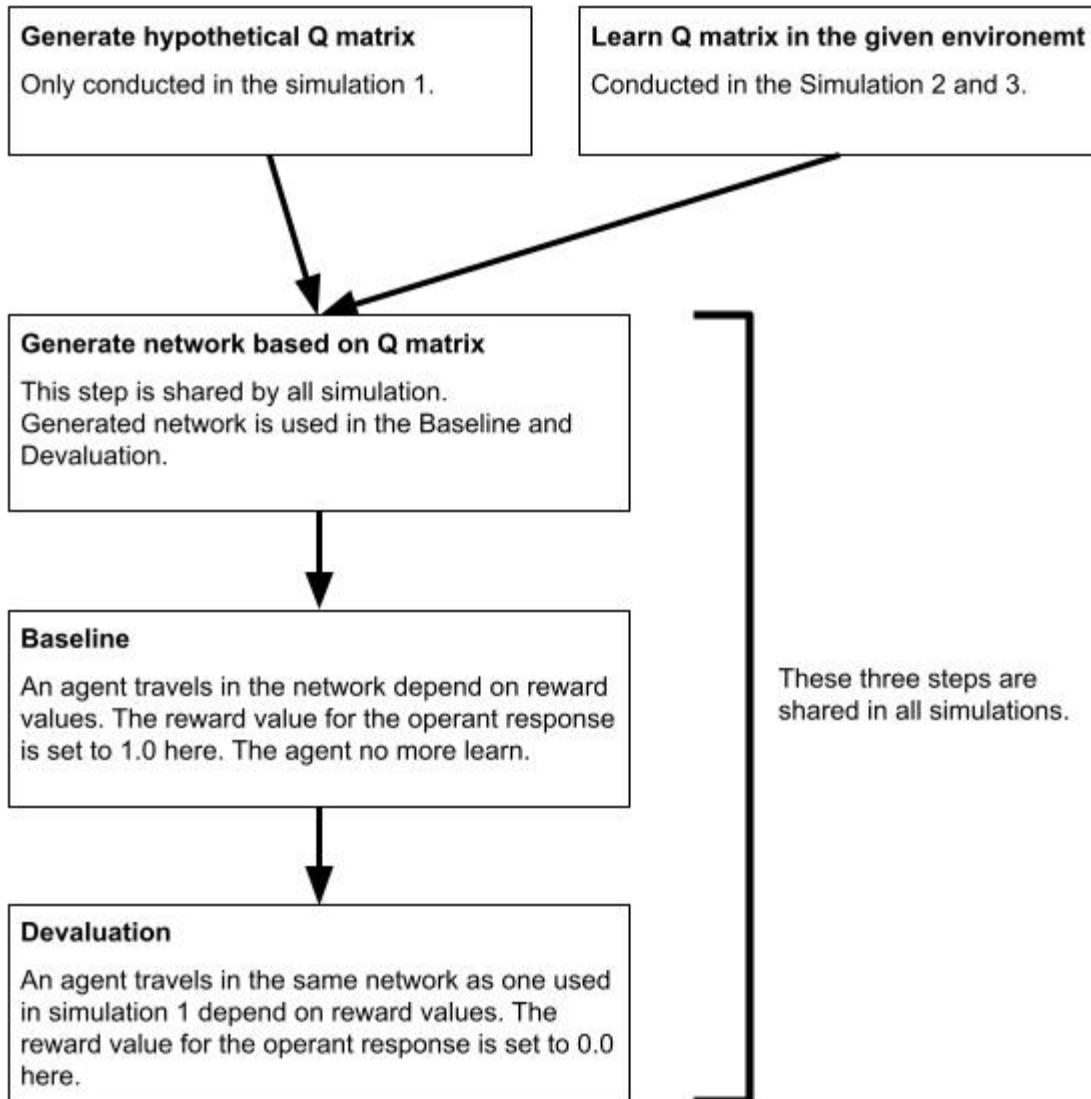
823 In this paper, we provide a novel perspective on habit formation by assuming behavior as a
824 network. In existing models, goal-directed behavior and habits are controlled by two distinct
825 systems. On the other hand, our model shows that although all responses are goal-directed,
826 both goal-directed and habits result from the structure of the network. It proposes that habit
827 formation is not caused by a change in the control of the two systems, but rather by a
828 continuous change in a single system. Furthermore, the most important feature of our model,
829 which differentiates it from other models, is that behavior is a network constructed from
830 responses. With this view, we have succeeded in providing a novel explanation for habit
831 formation. This implies that the possible algorithms can be changed depending on how one
832 views the behavior of organisms. Our study also suggests that changing the method of
833 capturing behavior could be a fundamental step in understanding the biological structure of
834 the behavior.

835

836 Materials and Methods

837 Overview

838 We conducted three simulations in this article, and they contain four steps (Figure 9). In the
839 first step, our agents are given a hypothetical Q matrix in the simulation 1 or learn the Q
840 matrix in the given environments in the simulation 2 and 3. In the second step, the agents
841 generate a network based on the Q matrix. The way to generate the network is the same in
842 all simulations. In the third step, baseline, the agents travel in the network and engage
843 responses. Here, the agents choose their responses based on reward values and the reward
844 value obtained by the operant response is set to 1.0. The agents no longer update the Q
845 matrix nor reconstruct the network. In the final step, devaluation, the agents behave in the
846 same way as the baseline. However, the reward value of operant response is reduced to
847 0.0. The only difference between baseline and devaluation is the reward value of operant
848 response. We explain the procedures conducted in the four steps in detail after sections.
849 Our simulation codes are available at: <https://github.com/7cm-diameter/hbnet>.



850

851 Figure 9. Overview of simulations.

852 **Generate hypothetical Q matrix**

853 Here, the agents are given a hypothetical Q-matrix instead of learning it through interactions

854 between an environment. First, we determine the number of nodes contained in a network.

855 We assign a scalar value for each node and it is represented by a vector. The first element

856 of the vector denotes the operant response and other elements denote other responses.

857 The values of the other responses are fixed to 0.01. The value of the operant response is

858 ranged from 0.01 to 1.0. The Q matrix is then defined as the direct product of the Q-vector.

859 Learn Q matrix in the given Environment

860 In simulation 2 and 3, agents learn the Q-matrix in an experimental environment. In the
861 simulation 2, we conducted simulations with variable interval (VI), variable ratio (VR),
862 concurrent VI VI, and VI with non-contingent rewards. Moreover, we changed the number
863 of rewards in the learning phase to examine the effect of training on habit formation. In
864 simulation 3, we conducted simulations with tandem VI VR and tandem VR VI.

865 In all these simulations, the agent chooses a response and the environment provides
866 a reward based on the response. The agent chooses a response according to the softmax
867 function; $p_i = \frac{e^{\beta_c Q_i}}{\sum_{j=1}^N e^{\beta_c Q_j}}$, where β_c denotes the inverse temperature, and N denotes the
868 number of responses in the given environment. We set $\beta_c = 3.0$ in all simulations. Then, the
869 agent updates the Q-matrix according to the response and the reward. In all simulations, we
870 employ fixed ratio (FR) 1 for other responses, where the agent can obtain rewards every
871 time it engages in the responses and the reward values are 0.001 for all other responses.
872 These flows are the same in all simulations. The only difference between the simulations is
873 the schedule in which the environment gives rewards to the agents. Algorithm 1 describes
874 general flow of all simulations. In the following sections, we explain the differences in the
875 schedules for each simulation.

Algorithm 1 Learn Q-matrix in the given environment

INPUT: N , $schedule$, $amount_training$

OUTPUT: Q

Q = generate an $N \times N$ matrix with all elements set to zero.

a_{t-1} = choose a response randomly from $[0, 1, 2, \dots, N]$

$response_durations$ = sample N samples from exponential distribution with $\lambda = 1 / 2.5$.

$r_{operant} = 0$

while $r_{operant} < amount_training$:

a_t = choose a response with softmax function and Q-values

τ = engage the chosen response for $response_durations[a]$ seconds.

r_t = schedule receive a_t and τ and return a reward

 update $Q(a_t, a_{t+1})$ according to Eq. 1

 if $a_t == 0$ and $r_t == 1$:

$r_{operant} += 1$

$a_{t-1} = a_t$

return Q

876

877 VR VI comparison and amount of training

878 The VR schedule presents rewards depending on the response of the agent. At each
879 response, the reward is presented at a given probability, which is the same as in the
880 simulations. This means that reward presentation follows the Bernoulli process, and the
881 number of responses required to obtain rewards follows the geometric distribution. We
882 generate pseudo-random numbers following the distribution in order for the numbers to
883 converge to the distribution in all simulations. More specifically, we divided the interval
884 ranging from 0 to 1 into equal divisions according to the number of rewards, and the
885 percentile points of the distribution were calculated for each point. Algorithm 2 shows how
886 to generate the required number of responses that follow the geometric distribution. We
887 employ VR 15 in simulation 2.

Algorithm 2 Variable ratio schedule

INPUT: p , $rewards$, $amount_training$

OUTPUT: $schedule$

q = Divide the range 0 - 1 into N equal parts

$reward_count \leftarrow 0$

$required_response = q[reward_count]$ th quantile of geometry distribution with parameter p .

def $schedule(a_t, r)$:

 if $a_t == 0$:

$required_response -= 1$

 else:

 return $rewards[a_t]$

 if $required_response \leq 0$:

$reward_count += 1$

$required_response = q[reward_count]$ th quantile of geometry distribution with parameter p .

 return $rewards[a_t]$

 return $schedule$

888

889 The VI schedule presents rewards depending on the time lapse. However, the agent
890 must emit responses to obtain rewards. Reward availability is determined at each time step
891 according to a probability, and once the reward becomes available, it remains available until
892 the agent takes the response. Reward availability follows the Poisson process, and the
893 intervals between each reward follow an exponential distribution. Pseudo-random numbers
894 are generated following the distribution in the same manner as the VR schedule. Algorithm
895 3 shows how to generate inter-reward intervals that follow an exponential distribution.
896 Moreover, we examined the effect of the amount of training on habit formation by
897 manipulating the number of rewards in both schedules. We calculated the average of inter-
898 reward intervals in the VR schedule and used them as the parameter of VI schedule.

Algorithm 3 Variable interval schedule

INPUT: λ , *rewards*, *amout_training*

OUTPUT: *schedule*

q = Divide the range 0 - 1 into N equal parts

reward_count = 0

required_time = $q[\textit{reward_count}]$ th quantile of exponential distribution with parameter λ .

def *schedule*(a_t , τ):

required_time -= τ

 if $a_t == 0$ and *required_time* <= 0:

reward_count += 1

required_time = $q[\textit{reward_count}]$ th quantile of exponential distribution with parameter λ .

 return *rewards*[a_t]

 elif $a_t \neq 0$:

 return *rewards*[a_t]

return *schedule*

899

900 Comparison between choice and single schedule

901 To examine the degree of habit formation when an explicit alternative is given, we used an
902 environment that mimics the experiment conducted by Kosaki and Dickinson¹⁵, where the
903 effect of the presence or absence of the alternative on habit formation. Here, the agent can
904 engage in two operant responses, and different rewards are assigned to each response. For
905 example, two levers were inserted into the apparatus and pressing the left lever produced
906 food, and the right levers produced a sucrose solution. In addition, as a control condition,
907 we used an environment in which the agent can engage only one operant response, but the
908 reward unavailable from the operant response is presented independent of the agent
909 responses.

910 We mimicked these experiments. In the choice condition, two of the responses were
911 treated as operant responses, and assigned two VI schedules with the same value and the
912 reward values obtained from both were set to 1.0. In the no-choice condition, the operant
913 response was assumed to be one, but the reward was presented independently of the
914 response in order to control the reward amount. We assigned a variable time schedule to
915 the rewards that are presented independent from the agent responses. We employ
916 concurrent VI 60 VI 60 in the choice condition, and concurrent VI 60 VT 60 in the no choice
917 condition.

Algorithm 4 Concurrent VI VI schedule

INPUT: λ , *rewards*, *amout_training*

OUTPUT: *schedule*

q_1 = Divide the range 0 - 1 into N equal parts

q_2 = Divide the range 0 - 1 into N equal parts

reward_count_1 = 0

reward_count_2 = 0

required_time_1 = q_1 [*reward_count_1*] th quantile of exponential distribution with parameter λ .

required_time_2 = q_2 [*reward_count_2*] th quantile of exponential distribution with parameter λ .

def *schedule*(a_t , τ):

required_time_1 -= τ

required_time_2 -= τ

 if $a_t == 0$ and *required_time_1* <= 0:

reward_count_1 += 1

required_time_1 = q_1 [*reward_count_1*] th quantile of exponential distribution with parameter λ .

 return *rewards*[a_t]

 if $a_t == 1$ and *required_time_2* <= 0:

reward_count_2 += 1

required_time_2 = q_2 [*reward_count_2*] th quantile of exponential distribution with parameter λ .

 return *rewards*[a_t]

 elif $a_t != 0$:

 return *rewards*[a_t]

return *schedule*

Algorithm 5 Concurrent VI VT schedule

INPUT: λ , *rewards*, *amout_training*

OUTPUT: *schedule*

q_1 = Divide the range 0 - 1 into N equal parts

q_2 = Divide the range 0 - 1 into N equal parts

reward_count_1 = 0

reward_count_2 = 0

required_time_1 = $q_1[\textit{reward_count_1}]$ th quantile of exponential distribution with parameter λ .

required_time_2 = $q_2[\textit{reward_count_2}]$ th quantile of exponential distribution with parameter λ .

def *schedule*(a_t , τ):

required_time_1 -= τ

required_time_2 -= τ

 if $a_t == 0$ and *required_time_1* ≤ 0 :

reward_count_1 += 1

required_time_1 = $q_1[\textit{reward_count_1}]$ th quantile of exponential distribution with parameter λ .

reward = *rewards*[a_t]

 elif $a_t \neq 0$:

reward = *rewards*[a_t]

 if *required_time_2* ≤ 0 :

reward_count_2 += 1

required_time_2 = $q_2[\textit{reward_count_2}]$ th quantile of exponential distribution with parameter λ .

reward += *rewards*[1]

 return *schedule*

919

920 Tandem VI VR and tandem VR VI

921 The tandem schedule is a schedule that presents multiple schedules in temporal succession.

922 For example, tandem FR 5 VI 30 means that VI 30 will start after the agent has responded

923 5 times, and the reward will be presented at the end of VI 30. In addition, since tandem does

924 not provide any explicit cues about the components it consists of, the agent cannot know

925 which schedule it is under. In tandem VI VR, the agent is first placed under a VI schedule,

926 and after it is finished, it is moved to a VR schedule. In tandem VR VI the order of

927 components is reversed, starting with the VI schedule, and followed by the VR schedule.

928 We employ tandem VI 15 VR 3 and VR 10 VI 5.

Algorithm 6 Tandem VR VI schedule

INPUT: $p, \lambda, rewards, amount_training$

OUTPUT: $schedule$

q = Divide the range 0 - 1 into N equal parts

$reward_count = 0$

$required_response = q[reward_count]$ th quantile of geometry distribution with parameter p .

$required_time = q[reward_count]$ th quantile of geometry distribution with parameter λ .

def $schedule(a_t, \tau)$:

 if $a_t == 0$:

$required_response -= 1$

 else:

 return $rewards[a_t]$

 if $required_response \leq 0$:

$required_time -= \tau$

 if $required_time \leq 0$. and $a_t == 0$:

$reward_count += 1$

$required_response = q[reward_count]$ th quantile of geometry distribution with parameter p .

$required_time = q[reward_count]$ th quantile of geometry distribution with parameter λ .

 return $rewards[a_t]$

return $schedule$

929

Algorithm 7 Tandem VI VR schedule

INPUT: $p, \lambda, rewards, amount_training$

OUTPUT: $schedule$

q = Divide the range 0 - 1 into N equal parts

$reward_count = 0$

$required_response = q[reward_count]$ th quantile of geometry distribution with parameter p .

$required_time = q[reward_count]$ th quantile of geometry distribution with parameter λ .

def $schedule(a_t, \tau)$:

$required_time -= \tau$

 if $a_t == 0$ and $required_time \leq 0$:

$required_response -= 1$

 elif $a_t \neq 0$:

 return $rewards[a_t]$

 if $required_response \leq 0$:

$reward_count += 1$

$required_response = q[reward_count]$ th quantile of geometry distribution with parameter p .

$required_time = q[reward_count]$ th quantile of exponential distribution with parameter λ .

 return $rewards[a_t]$

return $schedule$

930

931 Baseline and Devaluation

932 The reward devaluation is a procedure to examine whether an operant response is goal-
933 directed or habit under free-operant situations. First, an animal learns that he or she can
934 obtain a reward, food, or sucrose solution by pressing the lever. Learning lever pressings,
935 the animal was placed in an experimental environment and trained to the operant response.
936 After the training, reward devaluation was done by poisoning it with lithium chloride and a
937 brief period was added where the animal can access the reward freely. Then, the animal
938 was put into the experimental environment again and examined whether the number of
939 operant responses decreased. If the number of responses does not change, it implies that
940 the response is no longer controlled by its consequence, and the response becomes a habit.
941 In contrast, if the number of responses decreases, the response is controlled by its
942 consequences, such as goal-directed behavior. In our simulation, to reproduce the
943 procedure, we reduced the value of the reward obtained from the operant response after the
944 baseline phase.

945 Baseline

946 In the baseline phase, an agent travels on a network by choosing a response following Eq.
947 4 and searching for the shortest path between a currently engaging response and the goal.
948 The simulation contains three steps: 1) choice of response based on reward values, 2)
949 searching for the shortest path between the current response and the goal, and 3) engaging
950 responses successively contained in the path. We calculated the proportion of an operant
951 response to the total number of responses after some loops of the above 3 three steps.
952 Algorithm 1 shows the pseudocode of the simulation in the baseline phase.

953 Devaluation

954 In the devaluation phase, the agent behaved in the same way as in the baseline phase. The
955 difference between the devaluation and baseline phases is only the value of the reward
956 obtained from the operant response. In the baseline phase, we set the value to 1, and in the
957 test phase, we set it to 0. At the baseline phase, we calculated the proportion of the number
958 of operant responses to the total number of responses. Algorithm 8 describes the procedure
959 of the baseline and devaluation phase.

Algorithm 8 Procedure of baseline and devaluation phases

INPUT: *Network, N, rewards, loop*
OUTPUT: *propotion_operant*

s = choose an initial response from 0 - *N* randomly
operant = 0
total = 0
for *_* in 1:*loop*:
 t = choose a response according to Eq. 4
 shortest_path = find a shortest path from *s* to *t* on the *Network*
 total += number of response contained in *shortest_path*
 if 0 in *shortest_path*:
 operant += 1
 s = *t*
return *operant* / *total*

961 Data availability

962 All relevant data are within the paper (Figure 2 - 8 and all Supplementary Figures) and the
963 data and figures were generated using author's scripts (See Code availability).

964 Code availability

965 All Python scripts written for the simulations and analysis are available at
966 <https://github.com/7cm-diameter/hbnet>.

967 Acknowledgement

968 This research was supported by JSPS KAKENHI 20J21568 (KY), 18KK0070(KT),
969 19H05316 (KT), 19K03385 (KT), 19H01769 (KT), 22H01105 (KT), Keio Academic
970 Development Fund (KT), Keio Gijuku Fukuzawa Memorial Fund for the Advancement of
971 Education and Research (KT).

972 Reference

- 973 1. Perez, O. D., & Dickinson, A. (2020). A theory of actions and habits: The interaction
974 of rate correlation and contiguity systems in free-operant behavior. *Psychological*
975 *Review*, 127(6), 945.
- 976 2. Daw, N. D., Gershman, S. J., Seymour, B., Dayan, P., & Dolan, R. J. (2011). Model-
977 based influences on humans' choices and striatal prediction errors. *Neuron*, 69(6),
978 1204-1215.
- 979 3. Pezzulo, G., Rigoli, F., & Chersi, F. (2013). The mixed instrumental controller: using
980 value of information to combine habitual choice and mental simulation. *Frontiers in*
981 *psychology*, 4, 92.
- 982 4. Keramati, M., Smittenaar, P., Dolan, R. J., & Dayan, P. (2016). Adaptive integration
983 of habits into depth-limited planning defines a habitual-goal-directed spectrum.
984 *Proceedings of the National Academy of Sciences*, 113(45), 12868-12873.
- 985 5. De Houwer, J. (2019). On how definitions of habits can complicate habit research.
986 *Frontiers in Psychology*, 10, 2642.
- 987 6. Kruglanski, A. W., & Szumowska, E. (2020). Habitual behavior is goal-driven.
988 *Perspectives on Psychological Science*, 15(5), 1256-1271.
- 989 7. Dezfouli, A., & Balleine, B. W. (2012). Habits, action sequences and reinforcement
990 learning. *European Journal of Neuroscience*, 35(7), 1036-1051.
- 991 8. Garr, E., & Delamater, A. R. (2019). Exploring the relationship between actions,
992 habits, and automaticity in an action sequence task. *Learning & Memory*, 26(4), 128-
993 132.
- 994 9. Dezfouli, A., & Balleine, B. W. (2013). Actions, action sequences and habits:
995 evidence that goal-directed and habitual action control are hierarchically organized.
996 *PLoS computational biology*, 9(12), e1003364.
- 997 10. Dezfouli, A., Lingawi, N. W., & Balleine, B. W. (2014). Habits as action sequences:
998 hierarchical action control and changes in outcome value. *Philosophical*
999 *Transactions of the Royal Society B: Biological Sciences*, 369(1655), 20130482.
- 1000 11. Adams, C. D. (1982). Variations in the sensitivity of instrumental responding to
1001 reinforcer devaluation. *The Quarterly Journal of Experimental Psychology Section B*,
1002 34(2b), 77-98.

- 1003 12. Dickinson, A., Balleine, B., Watt, A., Gonzalez, F., & Boakes, R. A. (1995).
1004 Motivational control after extended instrumental training. *Animal Learning and*
1005 *Behavior*, 23, 197–206.
- 1006 13. Dickinson, A., Nicholas, D. J., & Adams, C. D. (1983). The effect of the instrumental
1007 training contingency on susceptibility to reinforcer devaluation. *The Quarterly Journal*
1008 *of Experimental Psychology*, 35(1), 35-51.
- 1009 14. Colwill, R. M., & Rescorla, R. A. (1985). Instrumental responding remains sensitive
1010 to reinforcer devaluation after extensive training. *Journal of Experimental*
1011 *Psychology: Animal Behavior Processes*, 11(4), 520.
- 1012 15. Kosaki, Y., & Dickinson, A. (2010). Choice and contingency in the development of
1013 behavioral autonomy during instrumental conditioning. *Journal of Experimental*
1014 *Psychology: Animal Behavior Processes*, 36(3), 334.
- 1015 16. Tinbergen, N. (1951). *The study of instinct*. Pygmalion Press, an imprint of Plunkett
1016 Lake Press.
- 1017 17. Wiltschko, A. B., Johnson, M. J., Iurilli, G., Peterson, R. E., Katon, J. M., Pashkovski,
1018 S. L., ... & Datta, S. R. (2015). Mapping sub-second structure in mouse behavior.
1019 *Neuron*, 88(6), 1121-1135.
- 1020 18. Markowitz, J. E., Gillis, W. F., Beron, C. C., Neufeld, S. Q., Robertson, K., Bhagat,
1021 N. D., ... & Datta, S. R. (2018). The striatum organizes 3D behavior via moment-to-
1022 moment action selection. *Cell*, 174(1), 44-58.
- 1023 19. Mathis, A., Mamidanna, P., Cury, K. M., Abe, T., Murthy, V. N., Mathis, M. W., &
1024 Bethge, M. (2018). DeepLabCut: markerless pose estimation of user-defined body
1025 parts with deep learning. *Nature Neuroscience*, 21(9), 1281-1289.
- 1026 20. Guthrie, E. R., & Horton, G. P. (1946). Cats in a puzzle box.
- 1027 21. Skinner, B. F. (1948). 'Superstition' in the pigeon. *Journal of experimental psychology*,
1028 38(2), 168.
- 1029 22. Staddon, J. E., & Simmelhag, V. L. (1971). The "superstition" experiment: A
1030 reexamination of its implications for the principles of adaptive behavior.
- 1031 23. Jenkins, H. M., & Moore, B. R. (1973). THE FORM OF THE AUTO-SHAPED
1032 RESPONSE WITH FOOD OR WATER REINFORCERS 1. *Journal of the*
1033 *experimental analysis of behavior*, 20(2), 163-181.

- 1034 24. Datta, S. R., Anderson, D. J., Branson, K., Perona, P., & Leifer, A. (2019).
1035 Computational neuroethology: a call to action. *Neuron*, 104(1), 11-24.
- 1036 25. Leon, A., Hernandez, V., Lopez, J., Guzman, I., Quintero, V., Toledo, P., ... &
1037 Escamilla, E. (2021). Beyond single discrete responses: An integrative and
1038 multidimensional analysis of behavioral dynamics assisted by Machine Learning.
1039 *bioRxiv*.
- 1040 26. Falk, J. L. (1966). Schedule-induced polydipsia as a function of fixed interval length
1041 1. *Journal of the Experimental Analysis of Behavior*, 9(1), 37-39.
- 1042 27. Gentry, W. D. (1968). FIXED-RATIO SCHEDULE-INDUCED AGGRESSION 1.
1043 *Journal of the Experimental Analysis of Behavior*, 11(6), 813-817.
- 1044 28. Levitsky, D., & Collier, G. (1968). Schedule-induced wheel running. *Physiology &*
1045 *Behavior*, 3(4), 571-573.
- 1046 29. Breland, K., & Breland, M. (1961). The misbehavior of organisms. *American*
1047 *Psychologist*, 16(11), 681.
- 1048 30. Guthrie, E. R. (1930). Conditioning as a principle of learning. *Psychological review*,
1049 37(5), 412.
- 1050 31. Herrnstein, R. J. (1970). On the law of effect 1. *Journal of the Experimental Analysis*
1051 *of Behavior*, 13(2), 243-266.
- 1052 32. Killeen, P. R., & Fetterman, J. G. (1988). A behavioral theory of timing. *Psychological*
1053 *Review*, 95(2), 274.
- 1054 33. Baum, W. M. (2012). Rethinking reinforcement: Allocation, induction, and
1055 contingency. *Journal of the experimental analysis of behavior*, 97(1), 101-124.
- 1056 34. Yamada, K., & Kanemura, A. (2020). Simulating bout-and-pause patterns with
1057 reinforcement learning. *PLoS One*, 15(11), e0242201.
- 1058 35. Watkins, C. J., & Dayan, P. (1992). Q-learning. *Machine learning*, 8(3-4), 279-292.
- 1059 36. Dijkstra, E. W. (1959). *Communication with an automatic computer* (Doctoral
1060 dissertation, Excelsior).
- 1061 37. Hagberg, A., Swart, P., & S Chult, D. (2008). *Exploring network structure, dynamics,*
1062 *and function using NetworkX* (No. LA-UR-08-05495; LA-UR-08-5495). Los Alamos
1063 National Lab.(LANL), Los Alamos, NM (United States).

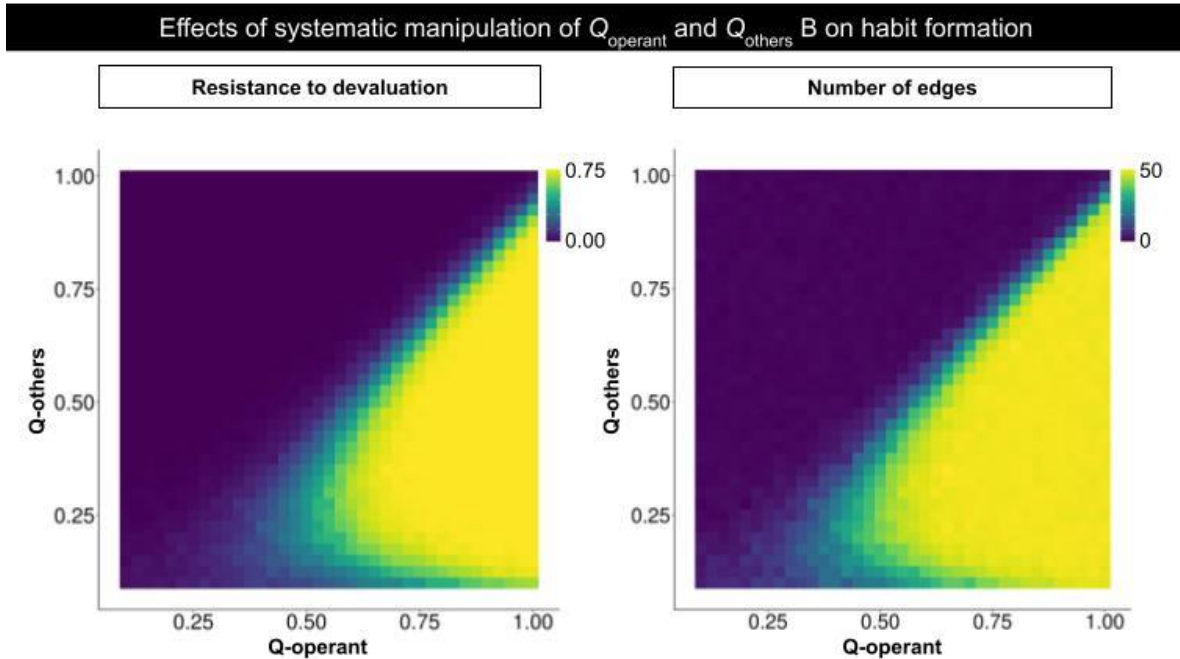
- 1064 38. Keramati, M., Dezfouli, A., & Piray, P. (2011). Speed/accuracy trade-off between the
1065 habitual and the goal-directed processes. *PLoS computational biology*, 7(5),
1066 e1002055.
- 1067 39. Albert, R., Jeong, H., & Barabási, A. L. (1999). Diameter of the world-wide web.
1068 *nature*, 401(6749), 130-131.
- 1069 40. Shull, R. L., Gaynor, S. T., & Grimes, J. A. (2001). Response rate viewed as
1070 engagement bouts: Effects of relative reinforcement and schedule type. *Journal of*
1071 *the experimental analysis of behavior*, 75(3), 247-274.
- 1072 41. Tanno, T. (2016). Response-bout analysis of interresponse times in variable-ratio
1073 and variable-interval schedules. *Behavioural processes*, 132, 12-21.
- 1074 42. Matsui, H., Yamada, K., Sakagami, T., & Tanno, T. (2018). Modeling bout-pause
1075 response patterns in variable-ratio and variable-interval schedules using hierarchical
1076 Bayesian methodology. *Behavioural processes*, 157, 346-353.
- 1077 43. Dickinson, A. (1985). Actions and habits: the development of behavioural autonomy.
1078 *Philosophical Transactions of the Royal Society of London. B, Biological Sciences*,
1079 308(1135), 67-78.
- 1080 44. Corbit, L. H., Chieng, B. C., & Balleine, B. W. (2014). Effects of repeated cocaine
1081 exposure on habit learning and reversal by N-acetylcysteine.
1082 *Neuropsychopharmacology*, 39(8), 1893-1901.
- 1083 45. DeRusso, A., Fan, D., Gupta, J., Shelest, O., Costa, R. M., & Yin, H. H. (2010).
1084 Instrumental uncertainty as a determinant of behavior under interval schedules of
1085 reinforcement. *Frontiers in integrative neuroscience*, 4, 17.
- 1086 46. Garr, E., Bushra, B., Tu, N., & Delamater, A. R. (2020). Goal-directed control on
1087 interval schedules does not depend on the action-outcome correlation. *Journal of*
1088 *Experimental Psychology: Animal Learning and Cognition*, 46(1), 47.
- 1089 47. Wearden, J. H., & Clark, R. B. (1988). Interresponse-time reinforcement and
1090 behavior under aperiodic reinforcement schedules: A case study using computer
1091 modeling. *Journal of Experimental Psychology: Animal Behavior Processes*, 14(2),
1092 200.
- 1093 48. Tanno, T., & Silberberg, A. (2012). The copyist model of response emission.
1094 *Psychonomic Bulletin & Review*, 19(5), 759-778.

- 1095 49. Baum, W. M. (1973). The correlation-based law of effect 1. *Journal of the*
1096 *experimental analysis of behavior*, 20(1), 137-153.
- 1097 50. Baum, W. M. (1981). Optimization and the matching law as accounts of instrumental
1098 behavior. *Journal of the experimental analysis of behavior*, 36(3), 387-403.
- 1099 51. Shull, R. L. (2011). Bouts, changeovers, and units of operant behavior. *European*
1100 *Journal of Behavior Analysis*, 12(1), 49-72.
- 1101 52. Peele, D. B., Casey, J., & Silberberg, A. (1984). Primacy of interresponse-time
1102 reinforcement in accounting for rate differences under variable-ratio and variable-
1103 interval schedules. *Journal of experimental psychology: Animal behavior processes*,
1104 10(2), 149.
- 1105 53. Daw, N. D., Niv, Y., & Dayan, P. (2005). Uncertainty-based competition between
1106 prefrontal and dorsolateral striatal systems for behavioral control. *Nature*
1107 *neuroscience*, 8(12), 1704-1711.
- 1108 54. Sanabria, F., Daniels, C. W., Gupta, T., & Santos, C. (2019). A computational
1109 formulation of the behavior systems account of the temporal organization of
1110 motivated behavior. *Behavioural processes*, 169, 103952.
- 1111 55. Graybiel, A. M. (1998). The basal ganglia and chunking of action repertoires.
1112 *Neurobiology of learning and memory*, 70(1-2), 119-136.
- 1113 56. Graybiel, A. M. (2008). Habits, rituals, and the evaluative brain. *Annu. Rev. Neurosci.*,
1114 31, 359-387.
- 1115 57. Yin, H. H., Knowlton, B. J., & Balleine, B. B. (2004). Lesions of dorsolateral striatum
1116 preserve outcome expectancy but disrupt habit formation in instrumental learning.
1117 *European Journal of Neuroscience*, 19, 181-189.
- 1118 58. O'Hare, J. K., Ade, K. K., Sukharnikova, T., Van Hooser, S. D., Palmeri, M. L., Yin,
1119 H. H., & Calakos, N. (2016). Pathway-specific striatal substrates for habitual behavior.
1120 *Neuron*, 89(3), 472-479.
- 1121 59. Tang, C., Pawlak, A. P., Prokopenko, V., & West, M. O. (2007). Changes in activity
1122 of the striatum during formation of a motor habit. *European Journal of Neuroscience*,
1123 25(4), 1212-1227.
- 1124 60. Yin, H. H. (2010). The sensorimotor striatum is necessary for serial order learning.
1125 *Journal of Neuroscience*, 30(44), 14719-14723.

- 1126 61. Jurado-Parras, M. T., Safaie, M., Sarno, S., Louis, J., Karoutchi, C., Berret, B., &
1127 Robbe, D. (2020). The dorsal striatum energizes motor routines. *Current Biology*,
1128 30(22), 4362-4372.
- 1129 62. Aldridge, J. W., & Berridge, K. C. (1998). Coding of serial order by neostriatal
1130 neurons: a “natural action” approach to movement sequence. *Journal of*
1131 *Neuroscience*, 18(7), 2777-2787.
- 1132 63. Balleine, B. W., & O’doherly, J. P. (2010). Human and rodent homologies in action
1133 control: corticostriatal determinants of goal-directed and habitual action.
1134 *Neuropsychopharmacology*, 35(1), 48-69.
- 1135 64. Ostlund, S. B., & Balleine, B. W. (2005). Lesions of medial prefrontal cortex disrupt
1136 the acquisition but not the expression of goal-directed learning. *Journal of*
1137 *Neuroscience*, 25(34), 7763-7770.
- 1138 65. Yin, H. H., Ostlund, S. B., Knowlton, B. J., & Balleine, B. W. (2005). The role of the
1139 dorsomedial striatum in instrumental conditioning. *European Journal of*
1140 *Neuroscience*, 22(2), 513-523.
- 1141 66. Niki, H., & Watanabe, M. (1979). Prefrontal and cingulate unit activity during timing
1142 behavior in the monkey. *Brain research*, 171(2), 213-224.
- 1143 67. Schultz, W., Apicella, P., Scarnati, E., & Ljungberg, T. (1992). Neuronal activity in
1144 monkey ventral striatum related to the expectation of reward. *Journal of*
1145 *neuroscience*, 12(12), 4595-4610.
- 1146 68. Shidara, M., & Richmond, B. J. (2002). Anterior cingulate: single neuronal signals
1147 related to degree of reward expectancy. *Science*, 296(5573), 1709-1711.
- 1148 69. Watanabe, M. (1996). Reward expectancy in primate prefrontal neurons. *Nature*,
1149 382(6592), 629-632.
- 1150 70. Toda, K., Sugase-Miyamoto, Y., Mizuhiki, T., Inaba, K., Richmond, B. J., & Shidara,
1151 M. (2012). Differential encoding of factors influencing predicted reward value in
1152 monkey rostral anterior cingulate cortex. *PloS one*, 7(1), e30190.
- 1153 71. Turner, K. M., Svegborn, A., Langguth, M., McKenzie, C., & Robbins, T. (2021).
1154 Opposing roles of the dorsolateral and dorsomedial striatum in the acquisition of
1155 skilled action sequencing. *bioRxiv*.
- 1156 72. Chersi, F., & Burgess, N. (2015). The cognitive architecture of spatial navigation:
1157 hippocampal and striatal contributions. *Neuron*, 88(1), 64-77.

- 1158 73. Stoianov, I. P., Pennartz, C. M., Lansink, C. S., & Pezzulo, G. (2018). Model-based
1159 spatial navigation in the hippocampus-ventral striatum circuit: A computational
1160 analysis. *PLoS computational biology*, *14*(9), e1006316.
- 1161 74. Packard, M. G. (1999). Glutamate infused posttraining into the hippocampus or
1162 caudate-putamen differentially strengthens place and response learning.
1163 *Proceedings of the National Academy of Sciences*, *96*(22), 12881-12886.
- 1164 75. Packard, M. G., & McGaugh, J. L. (1996). Inactivation of hippocampus or caudate
1165 nucleus with lidocaine differentially affects expression of place and response
1166 learning. *Neurobiology of learning and memory*, *65*(1), 65-72.
- 1167 76. Dolan, R. J., & Dayan, P. (2013). Goals and habits in the brain. *Neuron*, *80*(2), 312-
1168 325.
- 1169 77. Cushman, F., & Morris, A. (2015). Habitual control of goal selection in humans.
1170 *Proceedings of the National Academy of Sciences*, *112*(45), 13817-13822.
- 1171 78. Corbit, L. H. (2018). Understanding the balance between goal-directed and habitual
1172 behavioral control. *Current opinion in behavioral sciences*, *20*, 161-168.
- 1173

1174 Supplementary materials

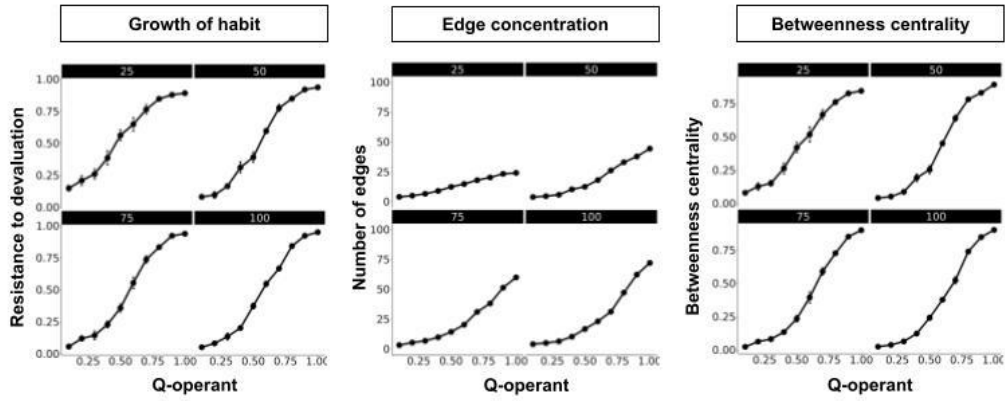


1176 Supplementary figure 1. Effects of systematic manipulation of Q_{operant} and Q_{others} on habit
1177 formation.

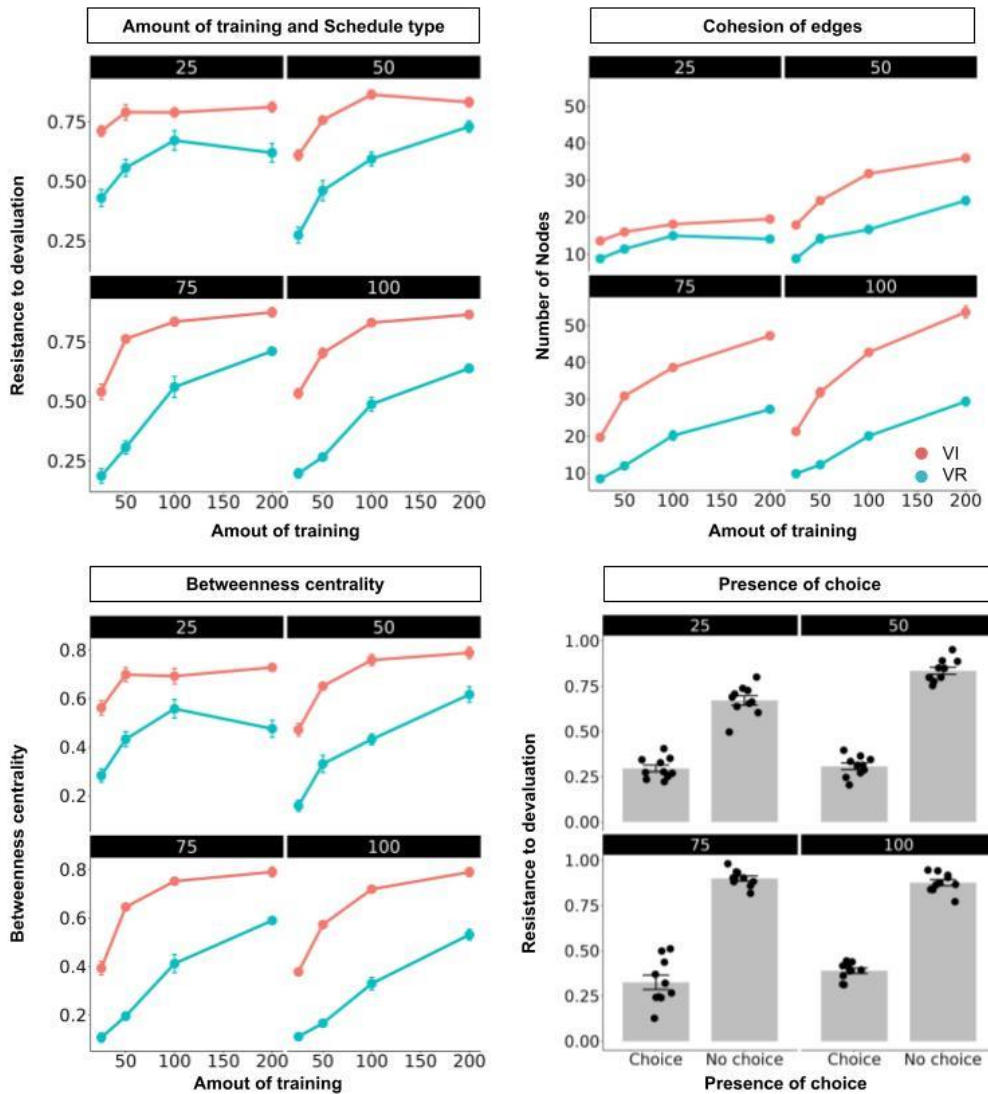
1178 The dependencies of the resistance to devaluation (left) and number of edges that the
1179 operant response acquired (right) on the Q_{operant} and Q_{others} . As the Q_{operant} increased, resistance
1180 to devaluation and number of edges increased, suggesting we confirmed the same result in
1181 the Simulation 1.

1182

A. Reproducibility of simulation results for different number of nodes in Simulation 1

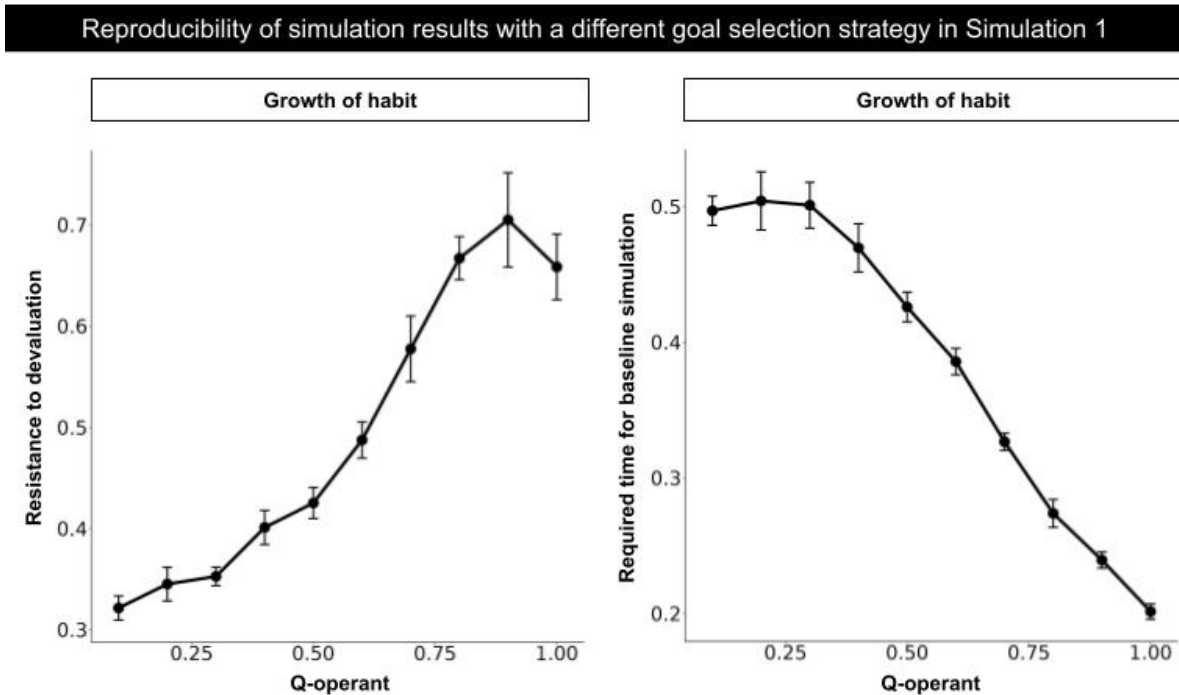


B. Reproducibility of simulation results for different number of nodes in Simulation 2



1184 Supplementary figure 2. Simulation results replicating Figure 2 (simulation 1), with the
1185 different numbers of nodes (25–100).

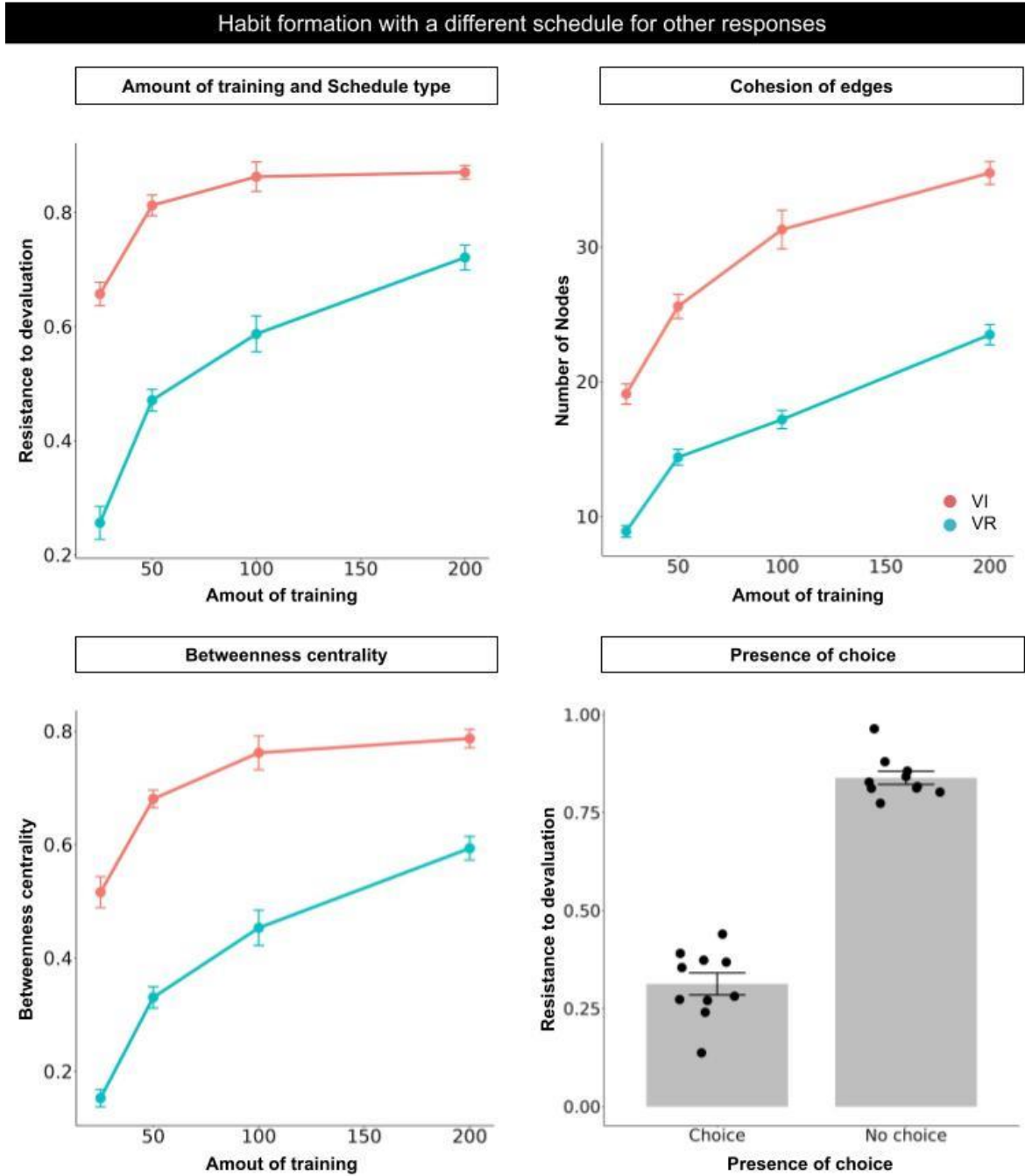
1186 We manipulate the number of nodes, 25, 50, 75, and 100, to confirm the results of our
1187 simulation are replicated in different numbers of nodes and all results are replicated.



1188

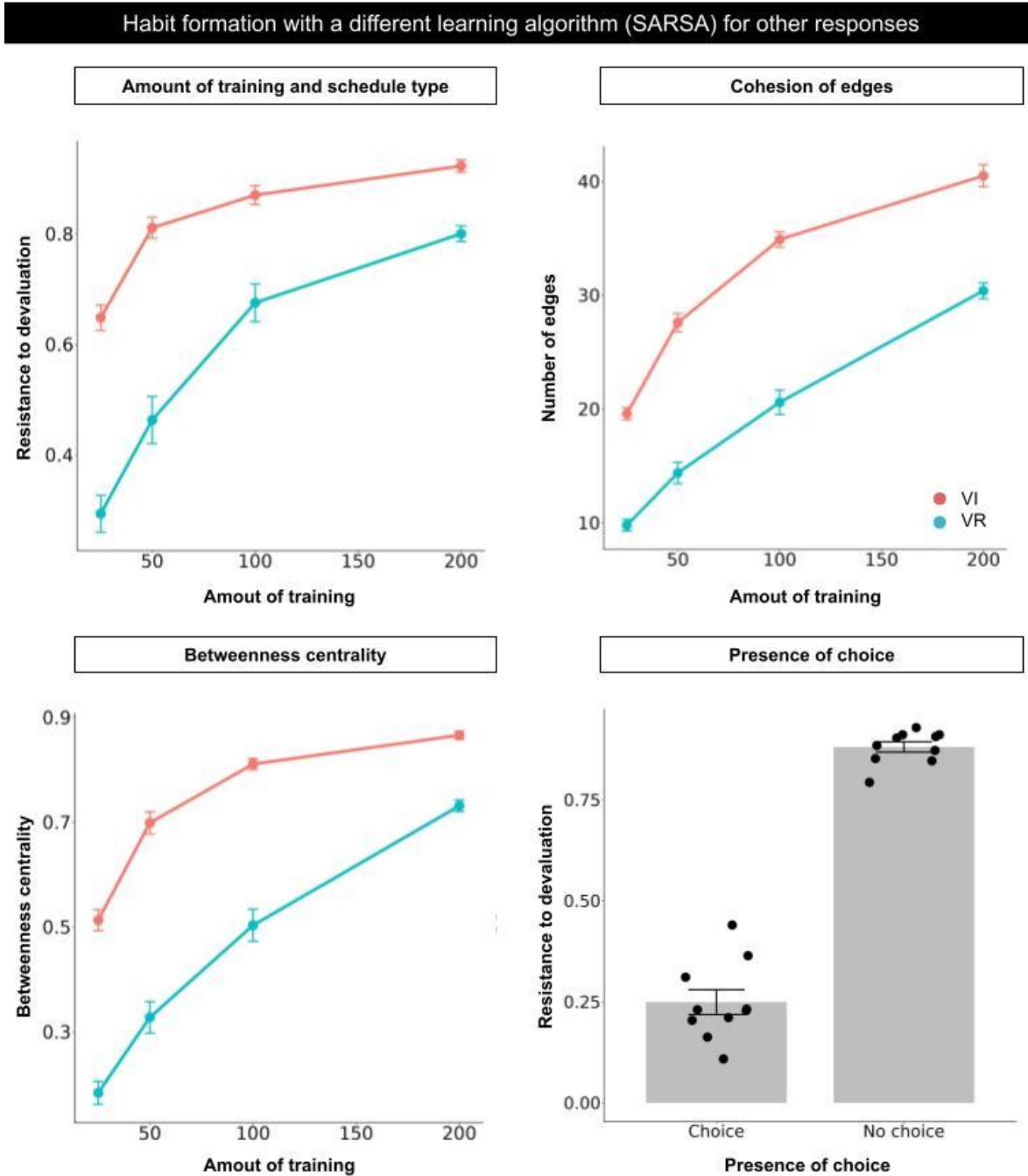
1189 Supplementary figure 3. Reproducibility of the results of Simulation 1 with a different
1190 response sequence generation algorithm.

1191 In the Simulation 1, response sequences were generated by a shortest path search,
1192 Dijkstra's algorithm. We employed another algorithm that is more weakly constrained and
1193 not the shortest path searching algorithm. In the new algorithm, an agent chooses a
1194 response randomly if a response chosen as a goal is not connected to the current engaging
1195 response. If the goal response is connected to the current engaging response, the agent
1196 chooses the response. In other words, the agent searches the goal response locally in the
1197 new algorithm. Resistance to devaluation, Edge concentration and betweenness centrality,
1198 all of features are replicated with the new algorithm, suggesting habit formation does not
1199 depend on the shortest path search as long as the response sequences are generated goal-
1200 directed.



1201

1202 Supplementary figure 4. Simulation results replicating Figure 3 (simulation 2), with a different
1203 schedule for other responses from the original simulation. We employed the VI 360 s
1204 schedule instead of FR 1 for other responses. We set their reward values as 1 / 50.



1205

1206 Supplementary figure 5. Simulation results replicating Figure 3 (simulation 3), with the
1207 different algorithm SARSA from the original algorithm Q-learning.