

Assembling hierarchies of action using sequencing and abstraction: studies and models of zero-shot learning

Gwydion Williams (✉ gwydion.williams.15@ucl.ac.uk)

University College London

Stefano Palminteri

École Normale Supérieure <https://orcid.org/0000-0001-5768-6646>

Patrick Haggard

University College London <https://orcid.org/0000-0001-7798-793X>

Article

Keywords:

Posted Date: February 3rd, 2022

DOI: <https://doi.org/10.21203/rs.3.rs-1312094/v1>

License:  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

1

2 **Assembling hierarchies of action using**
3 **sequencing and abstraction:**
4 **studies and models of zero-shot learning**

5

6

7 Gwydion Williams¹, Stefano Palminteri^{2,3,4}, Patrick Haggard¹

8

9 ¹ Institute of Cognitive Neuroscience, University College London, Alexandra House,
10 17 Queen Square, London WC1N 3AZ, UK

11 ² Laboratoire de Neurosciences Cognitives et Computationnelles, INSERM, Paris,
12 France

13 ³ Département d'Etudes Cognitives, ENS, PSL Research University, Paris, France

14 ⁴ Institute of Cognitive Neuroscience, HSE, Moscow, Federation of Russia

15

16 Corresponding author: Gwydion Williams; gwydion.williams.15@ucl.ac.uk

17 Abstract

18 Although the hierarchical structure of human action is widely acknowledged, we do
19 not fully understand how hierarchies of action are assembled. The standard view is
20 that low-level actions are *sequenced* to establish higher-level routines of behaviour.
21 Here we develop an alternative approach to building hierarchies, based on two
22 insights. First, we consider relations between sequence elements. Second, we
23 identify *abstract* features common to several such relations, and show how these
24 abstract features allow for flexible action sequence learning. We combine
25 *sequencing* and *abstraction* within a single model of hierarchical structure and test
26 this model in two distinct versions of a novel experimental paradigm. We
27 demonstrate that humans can learn entirely novel sequences of actions without
28 practice, by generalising learned sequence structures from one context to another.
29 Computational modelling showed that this ‘zero-shot learning’ of novel behaviours
30 was successfully captured by a hierarchical organisation of the kind we propose.

31 Introduction

32 As Lashley (1951) seminally observed, all actions we produce are component parts
33 of some sequence, and sequences of action are best understood through a
34 hierarchical lens. Learning to sequence actions, and then to flexibly arrange these
35 sequences into higher-level routines, is essential for many everyday tasks
36 (Rosenbaum, Kenny, & Derr, 1983; Yokoi & Diedrichsen, 2019). For example,
37 consider the hierarchy of behaviours required to brew coffee. To begin, we must
38 fetch a mug, some ground coffee, and a kettle of boiled water. To boil a kettle, we
39 must lift and move it to a tap, fill it with water, return it to its power source, and turn it
40 on. Each of these elements can in turn be further decomposed: to lift a kettle, we
41 must locate its handle in space, reach for the handle, and grasp it. This cascade of
42 ever lower-level representations of action establishes an entire hierarchy of
43 behaviours that, taken as a whole, will satisfy the original goal of making coffee (see
44 Figure 1). A hierarchical organisation is most often justified by its computational
45 efficiency; limited cognitive capacity can be dedicated to high-level features of action,
46 while low-level details are delegated to modular circuits at lower levels. In this study
47 we provide evidence and argue for a second benefit of hierarchy: it does not only
48 minimise cost, but it also maximises benefit by speeding rule discovery.

49 Classically, sequential action has been described as being a process of
50 building up chunks of behaviour by sequencing elementary or primitive actions
51 (Lashley, 1951). Under the classical view, each chunk of action would activate its
52 motor components in order and chunks themselves could be sequenced to form
53 progressively higher-level routines of action. This hierarchical sequencing of lower-
54 level actions to produce higher-level representations of order facilitates faster and
55 more accurate execution of primitive actions (Rosenbaum et al., 1983), provides a
56 more computationally efficient scheme to store and recall sequences of behaviour
57 (Ramkumar et al., 2016), and allows for entirely new sequences to be learned by
58 combining existing chunks in novel orders (Sakai, Kitaguchi, & Hikosaka, 2003).
59 Hierarchical sequencing has been observed in the study of sequential motor control
60 in humans (Cooper & Shallice, 2000; Fuster, 2008; Humphreys & Forde, 1998;
61 Miller, Galanter, & Pribram, 2017; Yokoi & Diedrichsen, 2019), and it is the
62 organisational principle used to arrange actions in hierarchical reinforcement
63 learning (see *temporal abstraction*: Botvinick, Niv, & Barto, 2009; Botvinick,
64 Weinstein, Solway, & Barto, 2015; Solway et al., 2014; Sutton, Precup, & Singh,
65 1999). According to this sequencing of low-level parts, high-level representations of
66 sequence therefore lack fine temporal detail of the low-level movements involved,
67 and they need only store information of the order in which constituent actions must
68 be initiated.

69 There is evidence for a second and more abstract mode of representation at
70 high levels; single neuron (Shima, Isoda, Mushiake, & Tanji, 2007) and population
71 (neuroimaging) data (Kornysheva et al., 2019) converge on the notion that, at high
72 levels, the relations between sequence elements are represented independently of
73 the elements themselves (see *abstraction* in Figure 1). These findings suggest that

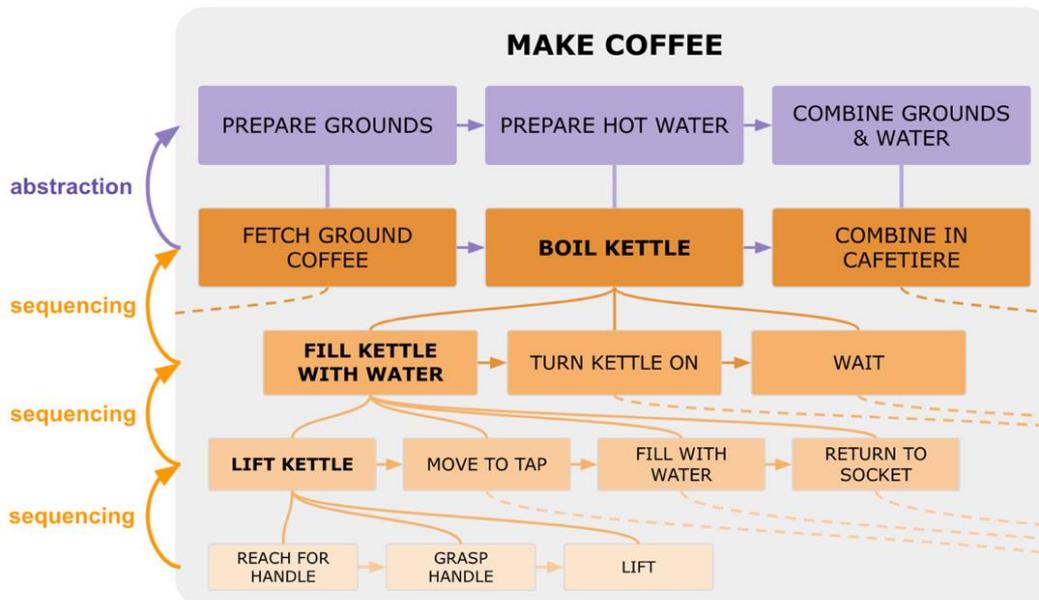


Figure 1 – Hierarchy of actions required to make coffee. Higher-level representations of action can come from two distinct operations: (1) sequencing low-level actions (e.g., reach for and grasp the handle of a kettle) can provide higher-level representations (e.g., lift kettle); and (2) abstracting over the individual actions in a sequence can provide abstract and relational representations of the relations between sequence elements independent of their content. This second method of abstraction can allow for the same relational representation (in purple) to produce distinctly different low-level sequences that adhere to the same relational structure (e.g., fetch ground coffee could be replaced with grind coffee beans to satisfy prepare grounds).

74 the human brain not only abstracts away the fine details of a motor command, but in
 75 some cases the actions themselves are lost in favour of a representation of the
 76 *relational structure* of the sequence (e.g., whether a given action should be repeated
 77 or not) independently of the specific action elements involved. Humans may use both
 78 *sequencing* and *abstraction* to form high-level representations of action, but it is
 79 unclear whether these two distinct operations are integrated under a single
 80 hierarchical framework to control behaviour, and the benefits doing so have not been
 81 explored.

82 We propose that *sequencing* and *abstraction* (see Figure 1) as two methods
 83 of building up higher-level routines of behaviour from lower-level actions are
 84 combined by the human brain under a single hierarchical organisation of behaviour.
 85 We further propose that this combination establishes an efficient, generalisable, and
 86 adaptive structure of human action. How might we detect this high-level organisation
 87 from low-level behavioural data? Movement patterns are typically silent about the
 88 generative processes that produce them. Further, observable movements represent
 89 the direct output of low-level modules, and recovering underlying higher-level
 90 structure is difficult because it is filtered by lower-level modules. Here we propose a
 91 new approach to extracting abstract hierarchical representations from behavioural
 92 data based on immediate generalisation of learned sequence structure to produce
 93 entirely novel action sequences that meet completely new challenges. We refer to
 94 this process as *zero-shot learning of novel behaviours*. We reasoned that, if people
 95 indeed form relational representations during learning complex action sequences,
 96 this should allow immediate generalisation to new action sequences that share the
 97 same relational properties but involve distinct low-level actions. For example,
 98 consider the abstract representation of the steps required to brew coffee in Figure 1.
 99 If one holds this abstract and relational representation of the steps required to brew

100 coffee, then when faced with a new coffee maker (say, a filter coffee machine), one
101 may be able to learn quickly how to brew coffee with the new apparatus by using this
102 abstract high-level representation of the steps required to produce an entirely novel
103 sequence of low-level actions. We propose this abstraction and generalisation of
104 structure to produce novel behaviours as a novel behavioural marker of latent
105 hierarchical structure.

106 In the present study, we demonstrate that humans do indeed exhibit zero-shot
107 learning of novel behaviours. We report two experiments on goal-directed action
108 which use very different visual presentations and framing, but an identical underlying
109 structure. To earn reward in the tasks, participants needed to navigate to a sub-goal
110 before moving on to an end goal. For example, in one of the two tasks, subjects
111 would need to find a key (sub-goal) in one room before opening a chest (goal) in
112 another room. There were two possible sub-goal locations and two possible goal
113 locations, but only one of each was active on a given trial (i.e., the key would be
114 placed in one of two rooms on each trial, as would the chest). The locations of the
115 sub-goal and goal were associated, such that one could predict the location of the
116 goal from the location of the sub-goal. Importantly, participants were told only where
117 the sub-goal was located on each trial, and the central challenge was to learn to
118 predict the location of the goal from that information, so as to achieve reward in as
119 few steps as possible. Given that there were multiple different sub-goal locations,
120 any single association between sub-goal and goal could require distinct sequences
121 of action on different trials. Further, these sub-goal-to-goal associations could
122 change without warning. We observed that participants learned new associations
123 from only a single trial following a switch. Crucially, they also immediately
124 generalised knowledge of the new association to produce entirely novel sequences
125 of low-level actions on subsequent trials. To formally verify that this zero-shot
126 learning of novel behaviours is indeed evidence of a hierarchical system that
127 included both sequencing and abstraction, we used computational modelling to
128 explore the necessary cognitive components of this learning process. We found that
129 we could only replicate zero-shot learning with a system that (1) organised behaviour
130 hierarchically by sequencing lower-level parts to provide higher-level representations
131 of ordered elements, (2) made use of relational high-level representations of action
132 by *abstracting* over lower-level sequences, (3) abstracted learning about these
133 relational representations over multiple states, and (4) directed exploration at
134 appropriate hierarchical levels. In sum, we present novel findings and related
135 computational evidence showing that hierarchical reinforcement learning is useful
136 not only for an efficient storage of wide repertoires of behaviour, but also for
137 impressively fast adaptation to changes in the environment when it is combined with
138 highly abstract representations of action.

139 **Results**

140 **Behavioural protocol and task structure.**

141 Our behavioural paradigm sought evidence for specific hierarchical representations
 142 that specify the relations between actions within a sequence. Participants were to
 143 navigate around the state map seen in Figure 2A in search of a sub-goal location
 144 (SG on the map). Visiting the subgoal would then allow them to receive reward at a
 145 separate goal location (G on the map). We used this state map to build two tasks
 146 which appeared to be very different (see supplementary materials) but were in fact
 147 structurally identical. In a spatial version of the task, participants navigated a set of
 148 rooms in search of a key (SG) that would open a chest (G). In a procedural version
 149 of the task, participants solved a puzzle by moving a rod to a specific cube-face
 150 (SG), which would then unlock reward at another cube-face (G). At subsequent
 151 debriefing, none of the participants reported recognising any similarities between the
 152 two tasks despite their identical structure.

153 The state map underlying both tasks was designed to require a specific
 154 hierarchy of actions to navigate efficiently around it (see Figure 3 for the full

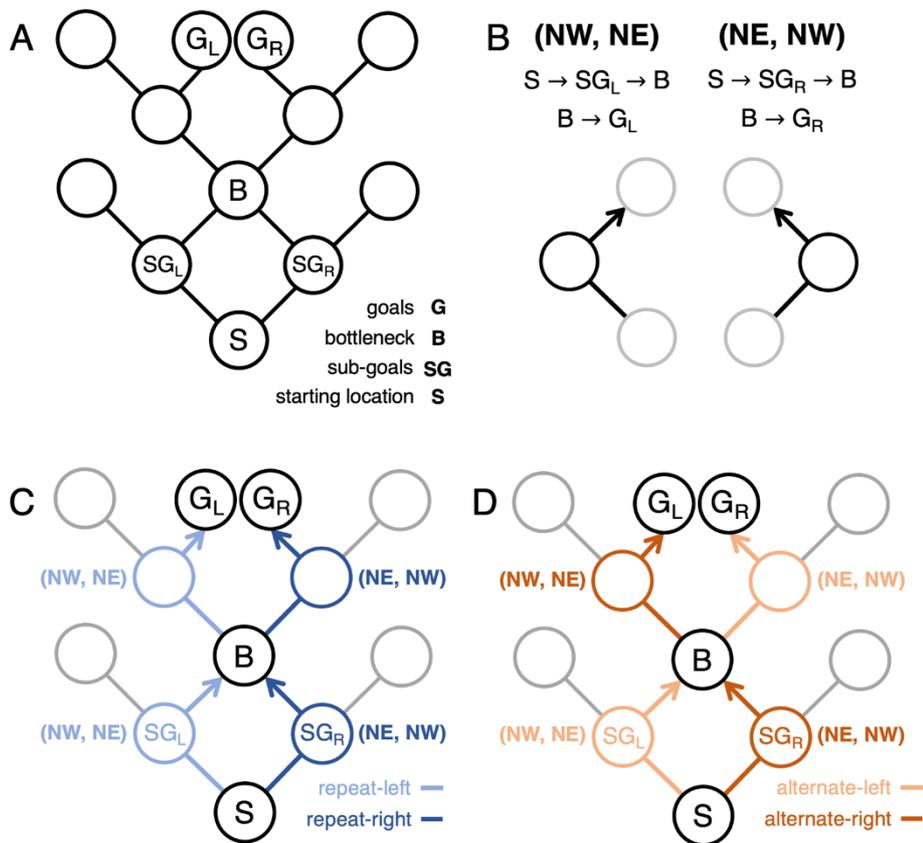


Figure 2 – (A) map of state space followed by both spatial and procedural tasks; (B) illustration of useful chunk of actions for navigating to/from the bottleneck state; (C) illustration of the two distinct sequences of action required if the association between SG and G is *repeat*; (D) illustration of the two distinct sequences of action required if the association between SG and G is *alternate*.

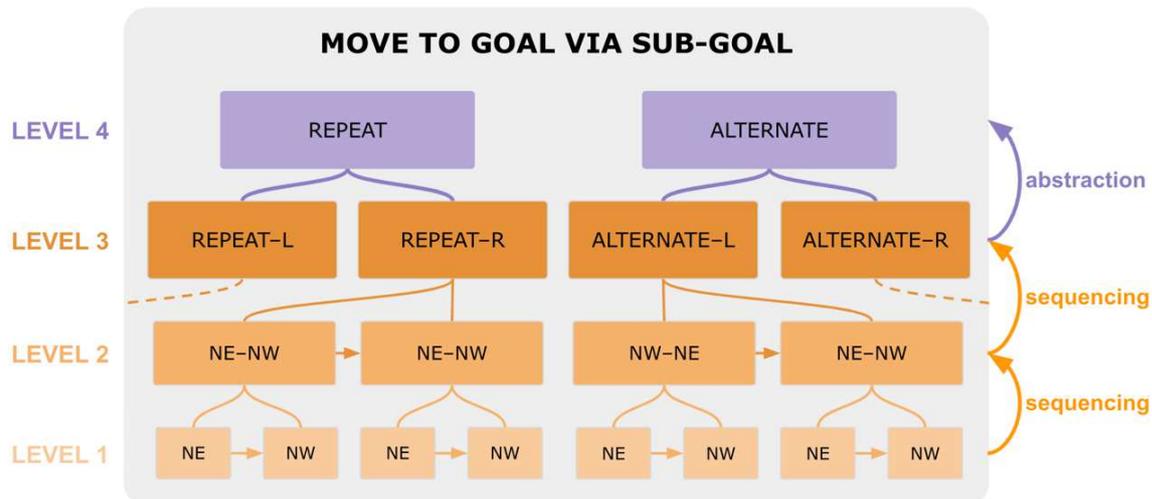


Figure 3 – Schematic of the hierarchy of actions targeted by the task design. Level 1 comprises the four primitive actions available in the task. Level 2 contains length-2 sequences that are useful for navigating to/from the central bottleneck (see Figure 1). Level 3 contains the full sequences of action required for an optimal solution of the four possible trial types (i.e., for all combinations of sub-goal–goal associations and sub-goal locations); and finally at level 4 we find abstractions over the two sub-goal–goal associations towards a relational representation of the actions involved.

155 hierarchy). The bottleneck in the centre of the map (see Figure 2A) needed to be
 156 traversed on all trials, and it needed to be traversed to move from the bottom half of
 157 the space to the top half, making it a useful target for behaviour. From the start
 158 position (S), either a sequence of (NW, NE) or a sequence of (NE, NW) would move
 159 participants from the starting location to the bottleneck (see Figure 2B). Given the
 160 symmetry between the bottom and top halves of the map, these same sequences
 161 were sufficient to then move from the bottleneck to each of the two possible goal
 162 locations. The four primitive actions (NE, NW, SE, SW) therefore occupy the lowest
 163 level (*level 1* in Figure 3) of our target behavioural hierarchy, and the chunks of 2
 164 sequential actions that are used for travelling to and from the bottleneck are one
 165 hierarchical level above the primitive actions (*level 2* in Figure 3; see Figure 2B for a
 166 demonstration).

167 On a given trial, only one of the two sub-goals (SGL or SGR) and one of the
 168 two goals (GL or GR) was active. For example, in the spatial task the participant
 169 would discover a key in only one of the two sub-goal rooms and a chest in only one
 170 of the two goal rooms. Participants were told at the start of a trial which of the two
 171 sub-goal states they should visit, and this therefore guided which of the two *level 2*
 172 sequences they should execute (see Figure 2B). Importantly, participants were not
 173 told which of the two goal locations was active, but the location of the goal could be
 174 predicted from the location of the sub-goal. Participants were told that they could
 175 predict where the goal would be from where the sub-goal was, but they were not told
 176 how to make this prediction. There were two possible associations between sub-goal
 177 and goal: (1) the goal could be on the same side as the sub-goal, or (2) the goal and
 178 sub-goals could be on different sides. We refer to the first of these two associations
 179 as *repeat*, and second as *alternate*. If a participant selected the correct *level 2*
 180 sequence such that they travelled to the bottleneck via the active sub-goal, then
 181 upon reaching the bottleneck they would need to decide between repeating the *level 2*
 182 sequence that got them there or alternating to execute the other of the two *level 2*
 183 sequences. The correct decision here would depend on the current association
 184 between sub-goal and goal: if the association was *repeat*, then the correct decision

185 is to repeat whatever *level 2* sequence used to reach the bottleneck, and if the
186 association is *alternate*, then one should alternate. This repetition of or alternation
187 between *level 2* sequences establishes four higher-level representations of the
188 sequences of actions required to solve the task (*level 3* in Figure 3): there are two
189 repetition sequences (one each for travelling via the left and right sub-goal, see
190 Figure 2C), and two alternation sequences (again, one each for travelling via the left
191 and right sub-goal, see Figure 2D). Finally, there is potential for an abstraction over
192 the *level 2* sequences being repeated in *level 3* such that our participants would
193 represent *repetition* and *alternation* independently of the *level 2* action sequences
194 being repeated or alternated (*level 4* in Figure 3). Crucially, participants were never
195 explicitly told whether they should repeat or alternate, but they could derive this
196 information by correctly learning and representing the relation between the sub-goal
197 and the goal, i.e., by representing the hierarchical and relational structure of the task.

198 The tasks were organised into three blocks of at least 30 trials each. In the
199 first block, the sub-goal-to-goal association was fixed. In the two blocks that followed,
200 the association between sub-goal and goal would switch on one of the first 10 trials,
201 and participants would then complete 30 trials under the new association (see Figure
202 4A). We refer to the trials where these switches in association occur as *switch trials*.
203 Participants were informed in the instructions that the associations between sub-goal
204 and goal could occasionally change. A switch trial could occur on a trial where the
205 sub-goal was present in either of the right or left locations, and so participants first
206 experienced the new association along only one of two possible paths through the
207 environment. For example, one participant might have experienced a switch from
208 *repeat* to *alternate* on a trial where the sub-goal was on the right, and they could
209 then learn how to act under *alternate* when the sub-goal is on the right. When the
210 sub-goal is next on the left, although the sequence of actions required will adhere to
211 the same *alternate* structure learned via the right, it requires a completely novel
212 sequence of low-level actions (compare the two *alternate* paths in Figure 2D). We
213 refer to the first trial along this inexperienced path following a switch in sub-goal-to-
214 goal associations the *novel post-switch trial*, and we refer to the novel sequences of
215 actions required on these trials as the *novel paths*. Given that the sub-goal is
216 randomly allocated to the right or left trial by trial, the novel post-switch trial might not
217 necessarily follow immediately after the switch: in our dataset the maximum number
218 of trials between a switch trial and its associated novel post-switch trial was four).

219 Immediate Acquisition of Novel Sequences

220 On both spatial and procedural tasks, all participants learned within the first
221 nine trials how to travel to the correct goal via the active sub-goal in an optimal four
222 moves (for the spatial task, median number of trials taken to make the optimal four
223 moves to goal was 3.5, inter-quartile range = 6.25; for the procedural task, median =
224 4.5, inter-quartile range = 2.75). Learning was slightly slower on the procedural task
225 (see the shallower rate of learning in Figure 4B), though behaviour did nevertheless
226 converge on the optimum of four moves to goal. The slower rate of learning on the
227 procedural task (learning rates found by fitting exponential models of learning to
228 each participant's data for each domain were significantly slower for procedural than

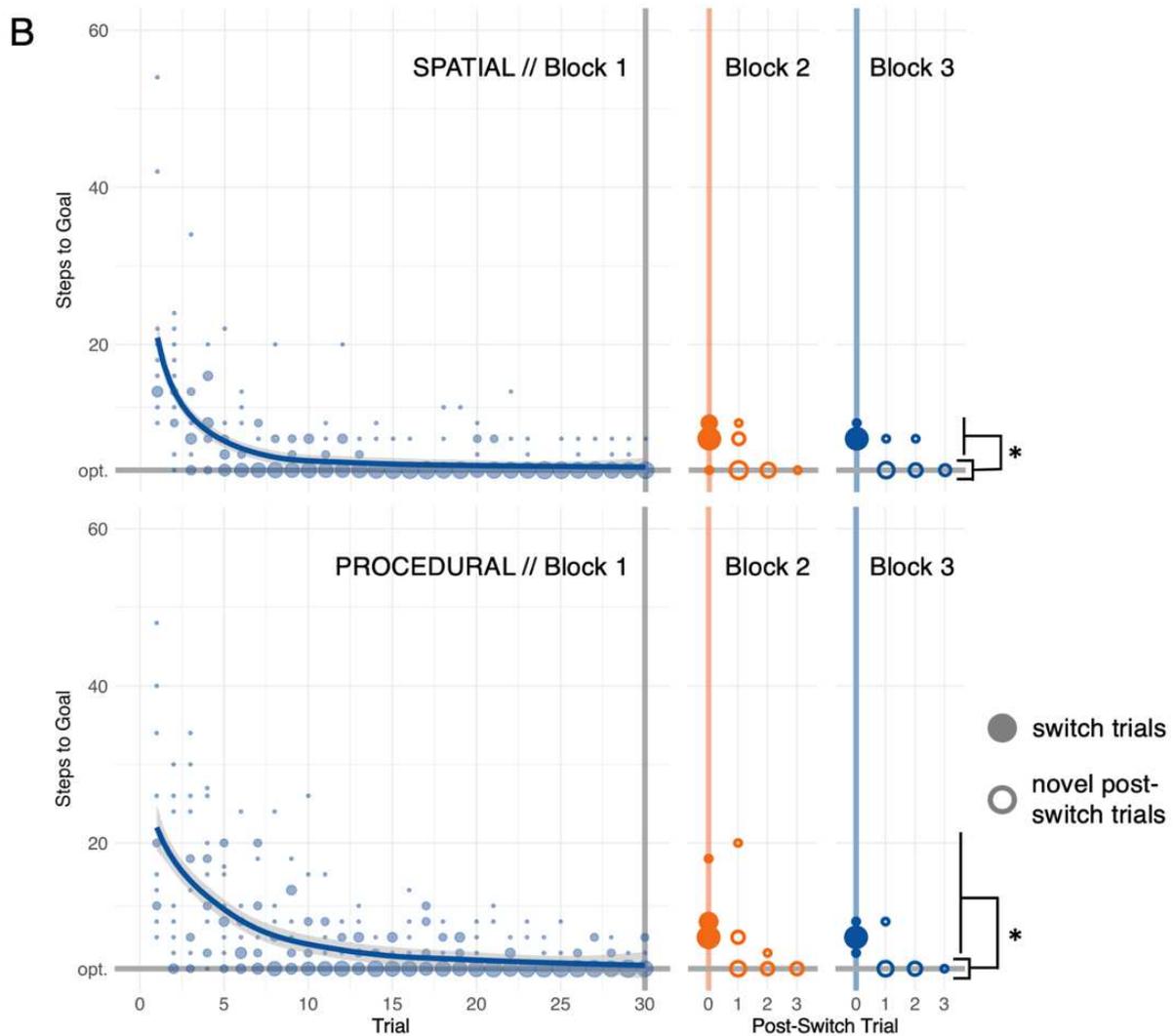


Figure 4 – (A) An example of the procedure followed by each of the tasks (note that the order of SG-G associations was counterbalanced over participants); (B) Observed behaviour of 12 subjects on each of the spatial and procedural tasks. The first column plots behaviour of all 12 subjects in the first block of each task to demonstrate an initial phase of learning and an eventual convergence onto the optimal solution to both tasks. The following two columns present recovery after a switch in SG-G association. The vertical orange/blue bars are the switch trials (these correspond to the underlined switch trials in A), and the hollowed-out points that follow plot behaviour on the novel post-switch trial for all twelve participants (these correspond to the underlined novel post-switch trials in A). Across the board, for any number of trials in between the switch and novel post-switch trials, participants were more likely than not to exhibit optimal behaviour, and this was true in both spatial and procedural tasks.

229 for spatial, $t(11) = 2.61, p = .024, d = 0.75$) may be due to the unfamiliar setting.
 230 Once participants found the optimal solution, they generally continued to perform
 231 optimally (see the stable optimal behaviour in block 1 of Figure 4B), with only minor
 232 and infrequent deviations, presumably reflecting lapses in attention.

233 Our central interest here was in how quickly our participants could recover
 234 from a switch in the associations between sub-goal and goal. Specifically, we wanted

235 to ask whether participants would behave optimally on novel post-switch trials
236 despite having no experience of travelling along the corresponding novel path. That
237 is, we were searching for zero-shot learning. This would require a high-level
238 relational representation of alternation and repetition (as in *level 4* of Figure 3) that
239 participants could use to adaptively generate completely novel sequences of
240 behaviour that followed these relational structures. For example, participants could
241 learn that alternating via the left sub-goal following the switch would mean that they
242 should also alternate via the right sub-goal, and upon first visiting the right sub-goal
243 they would know immediately how to solve the task. We found that most of our
244 participants selected the optimal path on novel post-switch trials; of a total of 48
245 novel post-switch trials, behaviour on 37 of these trials was optimal ($\chi^2(1) =$
246 $14.08, p < .001$). Further, the proportions of post-switch trials that were optimal for
247 each subject deviated significantly from a conservative chance level of 0.5 ($t(11) =$
248 $3.22, p = .008, d = 0.93$). The number of intervening trials in between the switch and
249 novel post-switch trials (see Figure 4) had no significant effect (the number of
250 intervening trials did not predict a significant portion of variance in steps to goal on
251 novel post-switch trials, $F(1, 46) = 0.94, p = .336$). Finally, there was no evidence to
252 support an association between optimality on novel post-switch trials and task
253 domain (results from chi-squared test for association: $\chi^2(1) = 0.12, p = .731$),
254 indicating that the ability to learn immediately how to act was general and not tied to
255 any individual task. That is, our participants spontaneously generalised learned sub-
256 goal–goal associations to produce entirely novel and optimal sequences of
257 behaviour, an observation which we refer to as *zero-shot learning*. Crucially, the
258 level 1 actions on switch and novel post-switch trials are entirely different, which
259 requires an abstraction over the sequences produced on switch trials to later
260 produce a novel sequence of behaviour that follows the same relational structure.

261 Note that 0.5 is a very conservative chance level for the likelihood of
262 mistakenly performing zero-shot learning of the novel path following a switch. In
263 reality, if our participants understood nothing of the high-level relations between sub-
264 goal and goal, then there would be no reason to think that any change in association
265 between the left sub-goal and its corresponding goal location should result in a
266 change in association between the right sub-goal and its corresponding goal
267 location. As a result, if a switch to *alternate* trial fell on a trial where the sub-goal was
268 on the left, when next encountering a trial where the sub-goal was on the right, the
269 rational choice would be to follow the association that was active before the switch
270 (*repeat*) and not the new association learned via the left sub-goal (*alternate*), making
271 the chance level for *alternating* via the right sub-goal 0. In fact, we found a mean
272 proportion of 0.77 (SD = 0.29) of novel post-switch trials being optimal over our
273 subjects, providing strong evidence for the ability to perform zero-shot learning of
274 novel sequences of action. We planned to derive estimates of the true chance level
275 for zero-shot learning from our simulations, which we discuss in greater detail below.

276 Computational Models

277 To illustrate how zero-shot learning of novel behaviours arises from hierarchically
278 organised behaviour that makes use of abstract relational representations of action
279 and to search for any other necessary cognitive components of the process, we built
280 a systematically organised set of four different RL models that aimed to capture our
281 participants' behavioural data (see Table 1 for a summary of differences between the
282 four models). The first and simplest model (Model 1 or *flat-history*) is the only non-

283 hierarchical model included, meaning it only has access to the four primitive actions
 284 (see Figure 3). It makes use of memory to solve the task (which is required given
 285 that the task is non-Markovian), where the remaining four models use hierarchically
 286 organised action to solve the task. We used standard Q-learning over temporal
 287 difference prediction errors (Sutton, 1988; Watkins & Dayan, 1992), and modelled
 288 participants' choices using a softmax function. This first model provided a non-
 289 hierarchical baseline against which we could compare the performance of our more
 290 complex hierarchical RL models.

291 The three remaining models were all hierarchical. All three follow the options
 292 framework (Sutton et al., 1999), which supplements the primitive actions available in
 293 standard, flat RL with temporally-abstract *options*, corresponding to superordinate
 294 chunks of behaviour. Our three HRL models hold the options required to furnish
 295 particular subsets of the behavioural hierarchy outlined in Figure 3. Models 2
 296 (*simple-hierarchical*) and 3 (*structured-hierarchical*) hold the first three levels of the
 297 hierarchy, but only model 4 (*abstract hierarchical*) holds the abstract and relational
 298 representations of repetition and alternation. Model 4 also abstracts learning over
 299 trials where the sub-goal is on the right and trials where the sub-goal is on the left.

300 We hypothesised that a preference to explore at high rather than low levels is
 301 central to the ability to quickly learn and use high-level relational rules, and to test
 302 this we implemented a specific modification of the softmax function in models 3 and
 303 4. Whereas standard softmax would include all actions/options no matter their
 304 hierarchical level, our structured-softmax chooses between only the highest-level
 305 options available given the current state of the agent. Further, where standard
 306 softmax uses a temperature parameter that is insensitive to any features of
 307 hierarchically organised action, our structured-softmax modifies the value of its
 308 temperature parameter as a function of the hierarchical level of the options under
 309 consideration. In practice, models 3 and 4 therefore choose only between the
 310 highest-level actions available in a given state and choice in general is biased
 311 towards exploration at higher hierarchical levels, and exploitation at lower levels.

Table 1 – Key differences between our four models.

	1. Flat w/History	2. Simple Hierarchical	3. Structured Hierarchical	4. Abstract Hierarchical
Levels available (see Figure 3)	1 (primitive actions only)	1-3 (hierarchy with no abstraction)	1-3 (hierarchy with no abstraction)	1-4 (hierarchy with abstraction)
History representation	Learns $Q(a, s_t, SG)$	Implicit in option execution	Implicit in option execution	Implicit in option execution
Hierarchical operations performed	Sequencing	Sequencing	Sequencing	Sequencing; Abstraction
Policy	Softmax	Softmax	Structured-Softmax	Structured-Softmax
Selects between...	All available actions	All available actions	Highest-level available actions	Highest-level available actions
Performs state abstraction?	No	No	No	Yes (between SG-R & SG-L trials)

312 In the first step of our computational analyses, we used simulations to see to
313 what extent our models could generate zero-shot learning of novel paths, as
314 observed in the behaviour of our participants. In the second step, we fit estimable
315 versions of these models to behaviour to move beyond the few trials where learning
316 of novel paths could take place and to investigate the global process of learning to
317 solve the entire task.

318 The Necessary Components of Zero-Shot Learning

319 To estimate how frequently each of our four models could reproduce zero-shot
320 learning by behaving optimally on novel post-switch trials, we simulated the
321 behaviour from each model for a range of parameter values. We manipulated
322 learning rate (alpha) and temperature (beta) to establish a grid of parameter values
323 (each of these two parameters could occupy any of the following values: 0.2, 0.4,
324 0.6, 0.8, 1.0), and for each combination of learning rate and temperature within this
325 grid we simulated behaviour on the task 20 times. From these simulations, we
326 computed the proportion of novel post-switch trials where behaviour was optimal. We
327 use these simulations not only to investigate the success of our models, but also to
328 estimate the true chance level of producing zero-shot learning under our various
329 models and the hypotheses they test.

330 We found that only model 4 (abstract-hierarchical) exhibited proportions of
331 zero-shot learning close to those observed empirically. As expected, our non-
332 hierarchical baseline produced almost no zero-shot learning, and this model provides
333 a good estimate of the true chance level of behaving optimally on novel post-switch
334 trials if we make no assumptions about the structure of behaviour. Models 2 (simple-
335 hierarchical) and 3 (structured hierarchical) lead to modest incremental
336 improvements as the organisation of behaviour becomes more sophisticated.
337 However, these two models only produce zero-shot learning by chance; they must
338 explore the options available to them on novel post-switch trials, and should they
339 happen to explore by selecting the newly-optimal high-level routine of action, we
340 would then see zero-shot learning. Model 4 offers a qualitative change in this
341 process, as it is able to learn from one context how to behave in another. That is,
342 model 4 can learn the abstract relations between sub-goal and goal from experience
343 with only one of the two sub-goal locations, and it can apply these learnings to guide
344 behaviour when it next encounters the other sub-goal. Unsurprisingly, therefore, the
345 success of model 4 in capturing zero-shot learning grows monotonically with learning
346 rate (see Figure 5B). A higher learning rate allows the model to learn new abstract
347 relations between sub-goal and goal from only a single trial. In summary, models 1,
348 2, and 3 fail to capture zero-shot learning of novel behaviours, but model 4
349 succeeds. Neither hierarchical organisation nor a preference for high-level
350 exploration alone were sufficient to capture zero-shot learning, but when combined
351 with abstract relational representations of action and an ability to abstract learning
352 over distinct states, all four components allowed model 4 to exhibit a near-human
353 ability to generalise learned structure to produce entirely novel sequences of action.

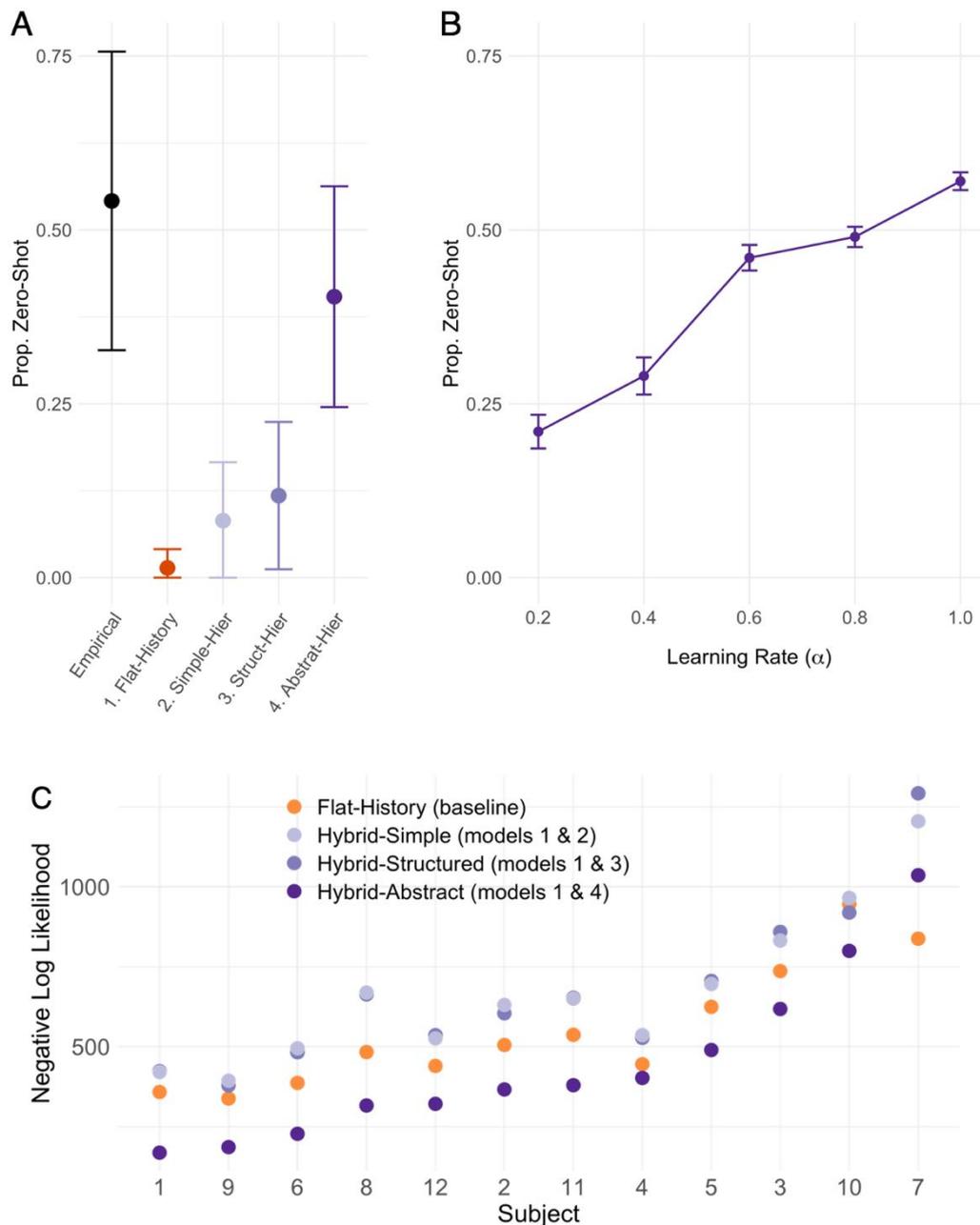


Figure 5 – (A) Mean (\pm SD) proportion of replications that exhibited zero-shot learning for a range of learning rates and temperatures for all four models, with the empirical means plotted for comparison. We see incremental improvements as we increase the complexity of our hierarchical models, but only model 4 is capable of reaching near-human performance. (B) Plot of how the ability of model 4 (our most successful model from (A) to capture zero-shot learning varies with learning rate – we find a monotonic increase in success with learning rate. (C) Fits of our flat-history (baseline) and three hybrid models to all participants. For 11 of 12 participants, the hybrid-abstract model clearly fits best, with the one remaining participants being fit best by our baseline flat model.

354 Model Fits to Behaviour

355 Our models were designed to capture one key behaviour of interest, namely zero-
 356 shot learning at the novel post-switch trial. However, zero-shot learning corresponds
 357 to a single sequence of actions within a much larger sequence of navigational or
 358 problem-solving actions (i.e., the entire task). We therefore additionally fitted these
 359 models to behaviour in the task, to investigate their generality, in addition to their
 360 local fit. In fact, our hierarchical models learned to use built-in options that were

361 designed to meet the demands of the task, while in reality the agent first needs to
362 learn from experience with the task what these useful options might be. Practically,
363 this means that our hierarchical models are unable to capture the initial period of
364 learning how to solve the task. We reasoned that this reflects the intuition that an
365 agent in a novel environment must first explore the outcomes of their low-level
366 actions and learn the structure of their environment, and only thereafter can they
367 build a hierarchical structure able to exploit the relational and structural features of
368 the task. We therefore decided to hybridise our models by using model 1 (*flat-*
369 *history*) again as a non-hierarchical baseline and to additionally combine this non-
370 hierarchical element with each of the hierarchical models in turn, producing three
371 distinct hybrid flat + hierarchical models. To specify the hybridisation process itself,
372 we included an arbitration process that apportioned control of behaviour between flat
373 and hierarchical systems, via an additional parameter, ω . When $\omega > 0.5$,
374 the flat system predominates, while for $\omega < 0.5$ the hierarchical system
375 predominates. The value of ω decays exponentially over time reflecting a shift,
376 with experience, from a flat system of behavioural control to a hierarchical
377 organisation of action. The agent must begin the task with a flat organisation of
378 behaviour (as it does not yet know the structure of the task) but with time discovers a
379 useful hierarchy of actions. The hybridising approach was intended to capture this
380 transition from flat to hierarchical behaviour while still allowing us to compare the
381 performance of each of our hierarchical models against experimental data.

382 We fit our hybrid models to behaviour using standard maximum likelihood
383 estimation. All four models were fit with only two free parameters – learning rate, and
384 temperature. The hybrid models set ω (governing arbitration between flat and
385 hierarchical systems) and its decay parameter to be fixed at values of 0.9 and 0.95
386 respectively. Fixing these parameters was necessary, as in order to fit our hybrid
387 model to behaviour, we had to permit occasional errors in behaviour to be attributed
388 to the flat system included in the model whatever the value of ω . In practice,
389 this means that we would occasionally allow the flat system to take control despite
390 the value for ω being below the threshold that would allow this to take place as
391 per the model specification. This slight deviation from the specification was
392 necessary because once the hierarchical system takes control (i.e., once ω
393 decays to a value below 0.5), the hierarchical models that use our modified
394 structured-softmax policy (models 3 and 4) can no longer account for actions that do
395 not conform to one of the highest-level representations of action available to these
396 models. This would result in infinitely poor fits. The errors observed empirically at this
397 late stage of the task were presumably due to lapses in attention and they are not of
398 central interest here, and so we allow for this slight deviation from the model
399 specification to avoid this issue. This was the case for all hybrid models, and so it
400 does not impair comparison between them.

401 We found that our hybrid version of model 4 provided the best fit to most of
402 our participants (see Figure 5C for fits, and Figure 6 for average predicted behaviour
403 given best-fitting parameters). Of our twelve participants: eleven were fit best by
404 hybrid model 4, and one was fit best by our flat baseline (*flat-history*). Not only did
405 hybrid model 4 provide the best fits to behaviour, but it did so with parameters which
406 we demonstrated to consistently exhibit zero-shot learning. We showed that learning

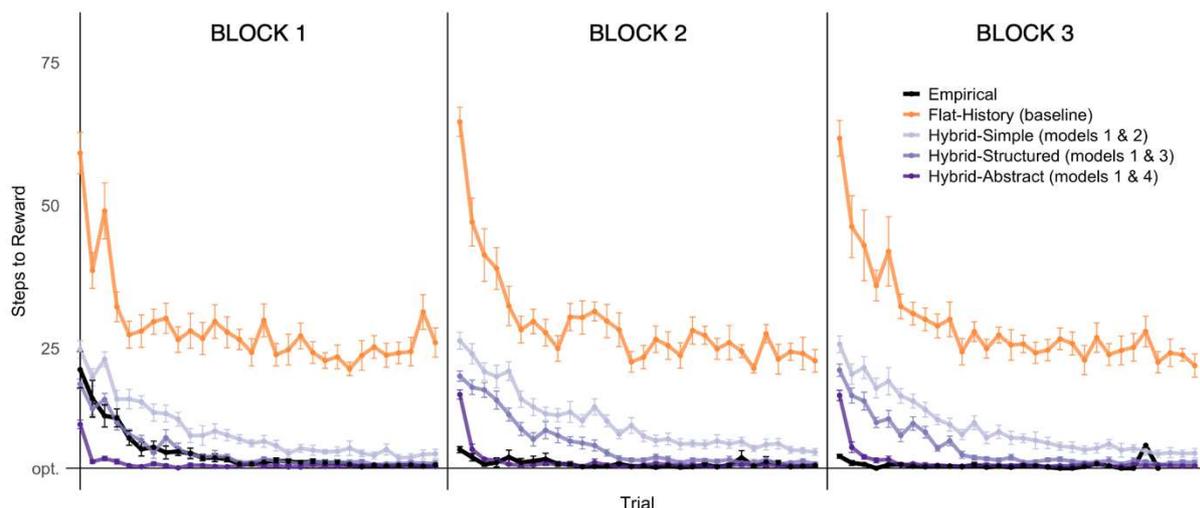


Figure 6 – Average simulated behaviour for each model with best fitting parameters to all 12 participants. For each model, we simulated the behaviour under the best fitting parameters to each participant 20 times, and we averaged over all replications and over all participants. We compare this with the empirically observed average behaviour (in black), which is taken over all 12 participants in both contexts (spatial and procedural).

407 rate was the primary factor determining the success of model 4 in capturing zero-
 408 shot learning. Further, zero-shot learning was best captured by high learning rates
 409 (see Figure 5B). The best fitting learning rates of the hybrid model to our 12
 410 participants were close to 1 (mean = 0.97, SD = 0.06), and these learning rates did
 411 not deviate significantly from the learning rate which we found to produce zero-shot
 412 learning most consistently in model 4 (no significant deviation from the optimal
 413 learning rate of 1, $t(11) = -1.74, p = .110$). Thus, the hybrid model achieved
 414 generality while still capturing our key behaviour of interest. Our hybrid model
 415 therefore fits well to behaviour, and it does so by fitting parameters that we have
 416 demonstrated to facilitate immediate generalisation of learned patterns of behaviour
 417 to generate completely novel sequences of action.

418 Hybridised versions of models 2 and 3 performed poorly, and were
 419 outperformed by the flat model for all but one participant. This reflects their inability
 420 to capture zero-shot learning. Given that these models cannot reliably capture zero-
 421 shot learning, not only are individual instances of zero-shot learning unlikely, but all
 422 following trials that perform the same sequence of actions are unlikely because
 423 these models receive no opportunity to unlearn the previous association between
 424 sub-goal and goal. For example, if the original association is repeated, model 3 will
 425 learn to solve the task by learning the two paths that implement repetition via the two
 426 sub-goals. However, consider a participant that experiences a switch to alternate via
 427 the right sub-goal, and then learns immediately what to do via the left sub-goal (i.e.,
 428 a participant that performs zero-shot learning of the left alternate path). Model 3
 429 would, in this case, be offered no opportunity to learn that repeating via the right is
 430 no longer rewarding, and given that model 3 learns only by experience with its
 431 environment (and it cannot learn by generalising abstract knowledge), it would
 432 expect repetition to be more likely than it was in reality because it expects repetition
 433 to still lead to reward. This was an unexpected finding: the hierarchical organisation
 434 used by models 2 and 3 was detrimental to their fits to behaviour, and this was owed
 435 to the inflexibility of these hierarchies and the omission of the *abstraction* step we
 436 outlined in Figure 1.

437 Discussion

438 Humans readily learn and produce action sequences based on high-level relational
439 features that cannot easily be accounted for by simple chaining or flat reinforcement
440 learning models. Here, we presented a novel and purely behavioural marker of the
441 otherwise latent hierarchical structure of behaviour; we found that human
442 participants were able to apply learned structural knowledge to generate completely
443 novel sequences of behaviour that met the demands of an evolving environment.
444 This ability to learn to produce novel sequences of behaviour without practice was
445 only captured by a (1) hierarchical reinforcement learning model that contained high-
446 level and (2) relational representations of action, similar to those observed in primate
447 prefrontal cortex (Shima et al., 2007), as well as an (3) ability to abstract learning
448 over multiple states and a (4) preference to explore at high levels of representation.
449 Simpler models lacking hierarchical structure could not capture this aspect of
450 performance, nor could hierarchical models that lacked relational representations of
451 action; all four components listed were necessary. Further, we found that immediate
452 generalisation of the structure of behaviour from one context to another also
453 depended on fast learning rates, and the best fits to behaviour were found by this
454 same model of abstract hierarchy (paired with a flat system to describe initial phases
455 of learning) with near-perfect learning. Learning how to behave in complex and
456 dynamic environments involves progressively building the hierarchies of behaviour
457 necessary to navigate through them, and we suggest that human agents do this not
458 only by *sequencing* lower-level actions towards higher-level representations of order,
459 but also by *abstracting* over actions in order to achieve a flexible, efficient, and
460 adaptive organisation of action.

461 Hierarchical Organisation, Relational Abstraction

462 To reproduce the immediate acquisition of novel behaviours, we found that a
463 hierarchical organisation of action was necessary, and we found that this should
464 include representations of the relations between sequence elements and not only
465 simpler chunks of primitive actions. These two components combine insights from
466 the study of motor control in human and non-human primates to provide a more
467 complete view of hierarchical control. Studies investigating the sequencing of action
468 suggest that the brain holds representations for action at several distinct levels of
469 detail (Botvinick, 2007; Koechlin, Ody, & Kouneiher, 2003; Lashley, 1951; Yokoi &
470 Diedrichsen, 2019). For example, representations of individual actions, of chunks of
471 actions, and of sequences of chunks have been found in the motor and premotor
472 areas of the human brain (Yokoi & Diedrichsen, 2019). Separately, *abstraction* has
473 been observed in the shape of relational representations of action that hold
474 information about the relations between the elements of a sequence (such as their
475 position or whether they will be repeated) independently of the actions that make up
476 the sequence found in primate prefrontal cortex (Shima et al., 2007) and in human
477 parahippocampal and cerebellar areas (Kornysheva et al., 2019).

478 Here, we found relational representations (e.g., repeat vs. alternate) over
479 sequences (e.g., repeat-left vs. repeat-right) composed of chunks (e.g., (NE, NW) vs.
480 (NW, NE)) of primitive actions (e.g., NE, NW, SE, SW). This organisation involves

481 sequencing of lower-level chunks to establish higher-level representations of order,
482 and abstraction to represent the relations between the lower-level sequences.
483 Evidence has been presented for both of these operations in isolation: (Yokoi &
484 Diedrichsen, 2019) recorded representations of individual movements, chunks of
485 movements, and sequences of chunks (sequencing); and Shima and colleagues
486 (2007) identified individual neurons in primate prefrontal cortex that responded to
487 any sequence containing an alternation between individual, primitive actions
488 (abstraction). Our results imply the use of relational representation of repetition of or
489 alternation between *chunks* of action, which requires sequencing to form the chunk
490 and abstraction to form the relational representation. Our research therefore ties
491 together sequencing and abstraction to demonstrate that both are used in tandem to
492 generate progressively higher-level representations of action and to produce
493 adaptive and flexible hierarchies of behaviour.

494 Hierarchy and relational structure also form a bridge between the study of
495 sequential motor control (Lashley, 1951) and hierarchical reinforcement learning
496 (Botvinick et al., 2009). To the best of our knowledge, hierarchical reinforcement
497 learning has considered only temporal abstraction (as in the options framework,
498 Sutton et al., 1999) as a method for building higher-level representations of action
499 from lower-level parts. This involves building increasingly high-level representations
500 of action by sequencing together lower-level actions. However, we have
501 demonstrated that relationally abstract representations of action similar to those
502 identified in the brain (Kornysheva et al., 2019; Shima et al., 2007) can and should
503 be included in hierarchical reinforcement learning models to accurately capture
504 human behaviour. Relational abstraction led to a powerful and impressively fast
505 ability to generalise behaviour between contexts in our participants, and it may
506 therefore be of computational benefit for HRL. In particular, relational abstraction
507 appeared essential for the key behavioural target of this paper: zero-shot learning, or
508 the ability to produce entirely novel sequences of action by generalising learned
509 relational representations to new contexts. In this way, a hierarchical organisation led
510 not only to an efficient storage of action that minimised computational cost, but also
511 to a beneficial ability to learn quickly how to adapt that maximised reward earned.

512 State Abstraction

513 Abstract representation of action was useful for our HRL models only because
514 of a third component we identified as necessary for immediate acquisition of novel
515 behaviours: state abstraction (Abel, 2019; Andre & Russell, 2002; Botvinick et al.,
516 2009; Radulescu, Niv, & Ballard, 2019). We allowed our most complex HRL model
517 (model 4: abstract hierarchical) to generalise whatever it learned from one context to
518 other relevant contexts. In our task, this meant being able to generalise learning
519 between trials where the sub-goal was on the right and trials where the sub-goal was
520 on the left. Abstraction over behaviour and the generalisation of learning over states
521 are tightly linked. Abstract representations of behaviour are useful because we often
522 want to execute sequences of action that are structurally similar but differ in the low-
523 level details. However, sequences will differ in low-level details only when they are
524 performed in different contexts. Thus, abstraction over behaviour is only useful if we
525 can apply whatever we learn about abstract behaviour to other contexts where it is
526 relevant and useful. For example, I do not need to learn how to brew a coffee anew
527 each time I visit a new kitchen – I can reapply my learnings from one kitchen to
528 another, i.e., I can abstract over states. Further, if the layout of a new kitchen is

529 different to any I have encountered before, I can still make coffee so long as I
530 represent the order of the high-level steps involved divorced from the low-level
531 actions that would implement those steps (i.e., I hold an abstract representation of
532 the sequence) such that I can adapt the precise low-level actions to match the new
533 layout. To summarise, we argue that abstraction over behaviour and abstraction of
534 learning over states together offer a powerful, adaptive, and efficient framework for
535 learning how to behave. In this study, we show how these two crucial cognitive
536 elements coexist in complex goal-directed action sequences.

537 Preference for High-Level Exploration

538 The fourth and final component we identified as necessary for zero-shot
539 learning of novel behaviours was a preference to explore at high levels. In our
540 models, this constraint was a directive, rather than a preference – the two models
541 with our novel structured softmax function were required to select only between the
542 highest-level options available to them in a given state. Exploration at high-levels will
543 generally be more valuable and efficient than low-level exploration. This is most
544 apparent at the extremes: there is little to no value in exploring new methods of
545 reaching out and grasping the handle of a kettle (low-level), but there may well be
546 value in exploring alternative coffee machines or sources of coffee beans. While
547 exploration-exploitation trade-offs are well-established in psychology (Mehlhorn et
548 al., 2015), their interaction with hierarchical representation has not been explicitly
549 considered. In our task, the changes in the environment that prompt exploration are
550 relevant for high-level representations of behaviour, so we cannot disentangle a
551 genuine preference to explore at higher-levels of abstraction from a preference for
552 level-appropriate exploration. Future research could change environments in
553 different ways, prompting a need to explore at distinct levels of abstraction, to clarify
554 this point. However, we see the cognitive efficiency of high-level exploration as a
555 prima facie advantage for a genuine preference for exploration at higher levels.
556 Whether this is correct or not, it seems the case that pruning the action space by
557 exploring at appropriate or high-levels would be beneficial for effectively and
558 efficiently resolving the exploration-exploitation trade-off.

559 Limitations & Future Directions

560 Although we identified these four components as necessary for reliably
561 producing zero-shot learning of novel sequences of action, they were not sufficient
562 by themselves to exactly match human behaviour. In fact, our participants showed
563 more frequent zero-shot learning than any of our models. We suggest this limitation
564 arises because our models lack a sophisticated mode of directed exploration that
565 would build on the level-appropriate exploration outlined above. Our participants
566 presumably immediately recognised that the rules of the task had changed upon
567 visiting a goal location they knew to previously hold reward, only to find no reward
568 upon reaching it. They could then rule out the association they previously believed to
569 be true and engage with directed exploration of alternatives, rather than exploring
570 either of the two high-level options they have available to them. This requires
571 incorporating a more sophisticated logic into how our agents explored alternative
572 actions in response to changes in the environment. As already discussed, future
573 research might investigate how more sophisticated exploration interacts with a
574 hierarchical organisation of behaviour to efficiently and adaptively guide choice.

575 Hierarchical models provided the best account of the key specific target
576 behaviour of zero-shot learning, from the set of models that we compared. However,
577 hierarchical models alone are insufficient to explain all behaviour for the simple fact
578 that in order to form a hierarchy of behaviour one must understand the structure of
579 the environment, and in order to understand the structure of the environment, one
580 must have some experience with it. To resolve this we needed to integrate our
581 hierarchical systems which captured the stable and optimal behaviour observed for
582 the majority of the task with a flat system that could capture the initial phase of
583 learning and any subsequent lapses. In effect, these hybrid models capture the
584 transition people must make from a flat system of behavioural control to a
585 hierarchical one, and our simple arbitration process represents the gestation of the
586 high-level options people come to use. This recalls the “option discover problem” in
587 hierarchical reinforcement learning (Botvinick et al., 2009; Stolle & Precup, 2002),
588 which remains largely unsolved. Although we captured a general transition from
589 memory-based flat control to hierarchical control, we have not explored mechanisms
590 to explain how hierarchies emerge from flat memory-based systems. Recent
591 developments in computational RL describe hierarchical memory systems that
592 divides the past into chunks for efficient recall of goal-relevant events (Lampinen,
593 Chan, Banino, & Hill, 2021). Hierarchical memory suggests a plausible intermediate
594 step between our rather simplistic flat system and our more sophisticated
595 hierarchical agent; memory could be chunked and explored in such a way that
596 associated chunks of behaviour can then be consolidated. Further research is
597 required to investigate how action hierarchies emerge from memory.

598 Conclusion

599 To conclude, we present a novel method of measuring the latent hierarchical
600 structure of action from behavioural data alone. Our findings support a novel view of
601 how hierarchies of action are formed in the human brain. Our key result was that
602 people are able to learn completely novel sequences of behaviour with no practice, a
603 process we refer to as zero-shot learning. We combined insights from sequential
604 motor control with hierarchical reinforcement learning to develop a model of goal-
605 directed hierarchical behaviour that could describe zero-shot learning and which
606 showed a number of interesting cognitive properties. First, we demonstrated that a
607 hierarchical organisation itself was necessary for zero-shot learning, as were
608 relational representations of action. This confirms our initial hypothesis that both
609 *sequencing* and *abstraction* were used to build hierarchies of behaviour in the
610 human brain. Second, we demonstrated that abstraction of learning between
611 different contexts goes hand in hand with relational representations of action to allow
612 an efficient, flexible, and adaptive organisation. Third, we showed that adding
613 hierarchical structure to action has important implications for how the exploration-
614 exploitation trade-off is negotiated. In sum, we provided direct behavioural evidence
615 for our proposed latent hierarchical structure of action sequences, and we identified
616 two unexpected additional components that were necessary to explain our
617 behavioural marker of this structure: (1) abstraction of learning between different
618 contexts and (2) level-appropriate exploration. Future research may shed further light
619 on the interactions between hierarchy and exploration, may describe more precisely
620 how we transition from flat memory-based behavioural control to hierarchical control,
621 and may expand further on the benefits of a hierarchical organisation of behaviour
622 that go beyond a mere minimisation of computational cost.

623 Methods

624 Subjects

625 Twelve subjects (mean age = 21.08 years, SD = 2.47; 5 males, 7 females) were
626 recruited to complete both tasks in one sitting. The only inclusion criteria were that
627 subjects were to be aged between 18-35, and of the twelve subjects seven were
628 female, and five were male. The probabilities of observing zero-shot learning by
629 chance under the null hypotheses of no hierarchical organisation/no relational
630 representations were identified by simulation. The chance probabilities were found to
631 be low, ranging from 0.01 (0.03) (Mean/SD) for our flat model to 0.12 (0.11) for our
632 non-abstract hierarchical models. We adopted a highly conservative estimate of 0.5
633 for zero-shot learning to occur by chance, given that zero-shot learning is ultimately
634 a binary choice between paths and so under a conservative atheoretical view this
635 choice becomes analogous to a coin flip. We performed a power calculation to
636 calculate the sample size required to detect a large effect (Cohen's $d=0.8$) of zero-
637 shot learning occurrence exceeding this chance estimate, with an alpha level of 0.05
638 and a beta (power) of 0.8. The large effect and relatively low power here are justified
639 by the functional nature of the test for zero-shot learning; we are testing for capacity,
640 which if present will be highly expressed, and if absent will not. This showed a
641 sample size of 12 participants. Subjects were all told that they would be paid an
642 amount that depended on their performance. In both tasks, performing well meant
643 moving from the starting location to the correct goal in as few moves as possible (the
644 optimum being four). All subjects consented to take part and the study was approved
645 by the relevant ethics committee.

646 Behavioural Task

647 Subjects performed both spatial and procedural tasks in one sitting. The order of the
648 tasks was counterbalanced across participants such that six of the twelve would
649 complete the spatial task first, and the other half would complete the procedural task
650 first. Each task began with an initial tutorial section which introduced subjects to the
651 rules of the tasks, how they could navigate around the environments, and the
652 incentive structure.

653 In the spatial task, subjects could move between rooms by clicking on the
654 door through which they would like to travel. Each room was identical but for a rune
655 which was present in the middle of the room. Each room held a different rune, and
656 these runes were static over all trials such that subjects could learn to place
657 themselves within the map by learning which room was associated with which rune.
658 The objective on each trial was to find a key and to then use that key to open a
659 chest. On each trial, participants were told in which of the two sub-goal rooms they
660 could find the key by providing them with the rune associated with the active sub-
661 goal room. Once they found the key, the queue for its location would disappear from
662 the screen, and the participants would then need to find the chest (i.e., the goal)
663 without any prompt. Once participants found the chest, the trial would end with the
664 delivery of reward. Participants would earn a set number of points for reaching the
665 goal, and they would earn a number of points based on how many doors they

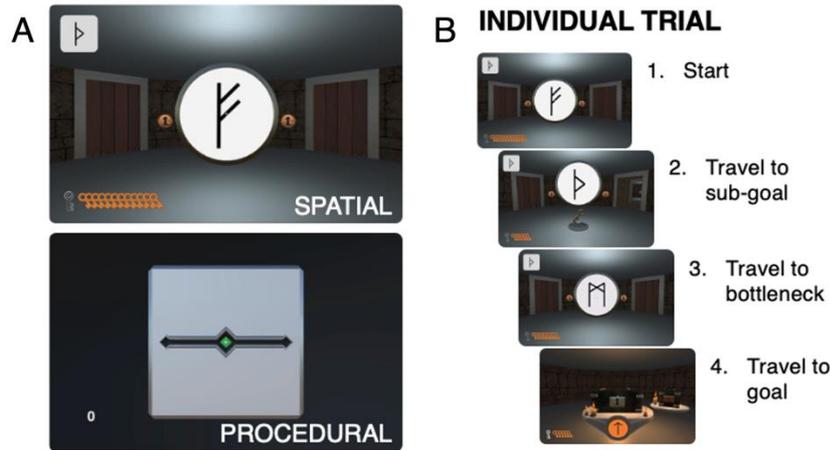


Figure S1 – (A) Initial starting states for spatial and procedural tasks. In the spatial task, participants navigated around a maze of rooms; in the procedural task, participants moved a rod around the faces of a cube. (B) Steps required to complete each individual trial.

666 opened and travelled through in the environment. If they opened all doors, they
 667 would earn no extra points, and they would earn 1 point per door left closed at the
 668 end of the trial. This incentivised participants to travel to the goal in as few moves as
 669 possible.

670 In the procedural task, subjects were to move a rod around the faces of a
 671 cube by clicking on the edge of the cube-face to which they would like to move. The
 672 objective was to move the rod to a cube face of a particular colour (sub-goal) before
 673 moving it to a golden cube face (goal). The target sub-goal colour was instructed at
 674 the start of the trial. Upon doing so, the trial would end, and points would be earned
 675 according to the number of moves taken to complete it. The objective was again to
 676 take as few moves as possible.

677 In both tasks, the location of the sub-goal was randomly allocated on each
 678 trial, though it was kept balanced such that there was always an even number of left
 679 and right sub-goal trials. From the location of the sub-goal, participants would need
 680 to learn to predict where the goal could be found. The sub-goal and goal could be
 681 associated in one of two ways: (1) under *repeat*, the sub-goal and goal would be on
 682 the same side; and under (2) *alternate*, they would be on different sides. Participants
 683 would start with one of these two associations being fixed for the 30 trials that make
 684 up the first block. Then, at some point during the first 10 trials of block 2 the
 685 association would switch and a fixed 30 trials under the new association would
 686 follow, and the same process would happen again on block 3. The order of repeat-
 687 alternate-repeat or alternate-repeat-alternate for blocks 1, 2, and 3 was
 688 counterbalanced over our twelve participants.

689 Computational Models

690 Our flat reinforcement learner followed the algorithm outlined in box 1, and our
 691 hierarchical reinforcement learning models all followed the algorithm outline in box 2.

Box 1 – Flat RL model specification

Initialise for all a :

$$Q(S, a) = 0$$

$$S_{origin} \leftarrow S_0$$

Repeat until $r = 1$:

$$A \leftarrow \text{softmax}(Q(S_t))$$

$$S_{t+1} \leftarrow T(S_t, O_{\text{primitive}})$$

$$r \leftarrow R(S_{t+1})$$

$$Q(S_t, S_{origin}, A) \leftarrow Q(S_t, S_{origin}, A) + \alpha \cdot (r + \gamma \cdot \max_a Q(S_{t+1}, S_{origin}, a) - Q(S_t, S_{origin}, A))$$

Box 2 – Hierarchical RL model specification

Initialise for all o and all corresponding S_{init} :

$$Q(S_{\text{init}}, o) = 0$$

Repeat until $r = 1$:

$$h \leftarrow 0$$

$$O(h) \leftarrow \text{structured-softmax}(Q(S_t))$$

$$S_{\text{init}} \leftarrow S_t$$

$$r \leftarrow 0$$

while S_t is not $S_{\text{term}, O(h)}$:

while $O(h)$ is not primitive:

$$h \leftarrow h + 1$$

$$O(h) \leftarrow \pi_{O_{h-1}}(S_t)$$

$$S_{t+1} \leftarrow T(S_t, O_{\text{primitive}})$$

$$r \leftarrow R(S_{t+1})$$

$$Q(S_{\text{init}}, O(0)) \leftarrow Q(S_{\text{init}}, O(0)) + \alpha \cdot (r + \gamma \cdot \max_o Q(S_{t+1}, o) - Q(S_{\text{init}}, O(0)))$$

692 Model Simulations

693 To simulate the behaviour of our four models, we established a grid of parameter
694 values for all learning rates in [0.2, 0.4, 0.6, 0.8, 1.0] and all temperatures in [0.2,
695 0.4, 0.6, 0.8, 1.0]. For each combination of learning rate and temperature, we
696 simulated the behaviour of our four models 20 times. Of these 20 simulated
697 datasets, we then investigated how often zero-shot learning of novel paths following
698 a switch in sub-goal–goal association occurred. The simulations included two blocks
699 of 100 trials, with the sub-goal alternating between right and left every other trial. The
700 switch in association would fall on the first trial of the second block, meaning that

701 these models had only a single trial to learn the new association before needing to
702 apply any learnings to guide behaviour on the novel post-switch trial.

703 Model Fitting Procedure

704 To fit our models to data, we used maximum likelihood estimation. To optimise, we
705 took the negative summed log likelihoods of each individual action given our model,
706 its parameters, and all “experience” up to that action. We minimised this value by
707 adjusting the relevant free parameters for each model using a limited memory BFGS
708 method of parameter estimation (Saputro & Widyaningsih, 2017).

709 Model Recovery

710 To ensure that our modelling and fitting procedure was sound and unbiased, we
711 simulated behaviour from our hybrid model given the best fitting parameters for each
712 subject. We then re-used the fitting procedure to fit our hybrid model to these now
713 simulated data to recover the parameters used in the simulation. We repeated this
714 process three. For most participants, we could recover the parameters used to
715 simulate the data with only minor deviations from ground truth (see Figure S2). The
716 only exception was participant 7: the fits for this participant were characterised by a
717 high temperature (beta). Higher temperature means that models explore their
718 environment more, leading to greater noise in the simulated datasets and therefore
719 more noise in the recovery process. With this exception, our simulations accurately
720 recovered the parameters used to simulate data from the hybrid model.

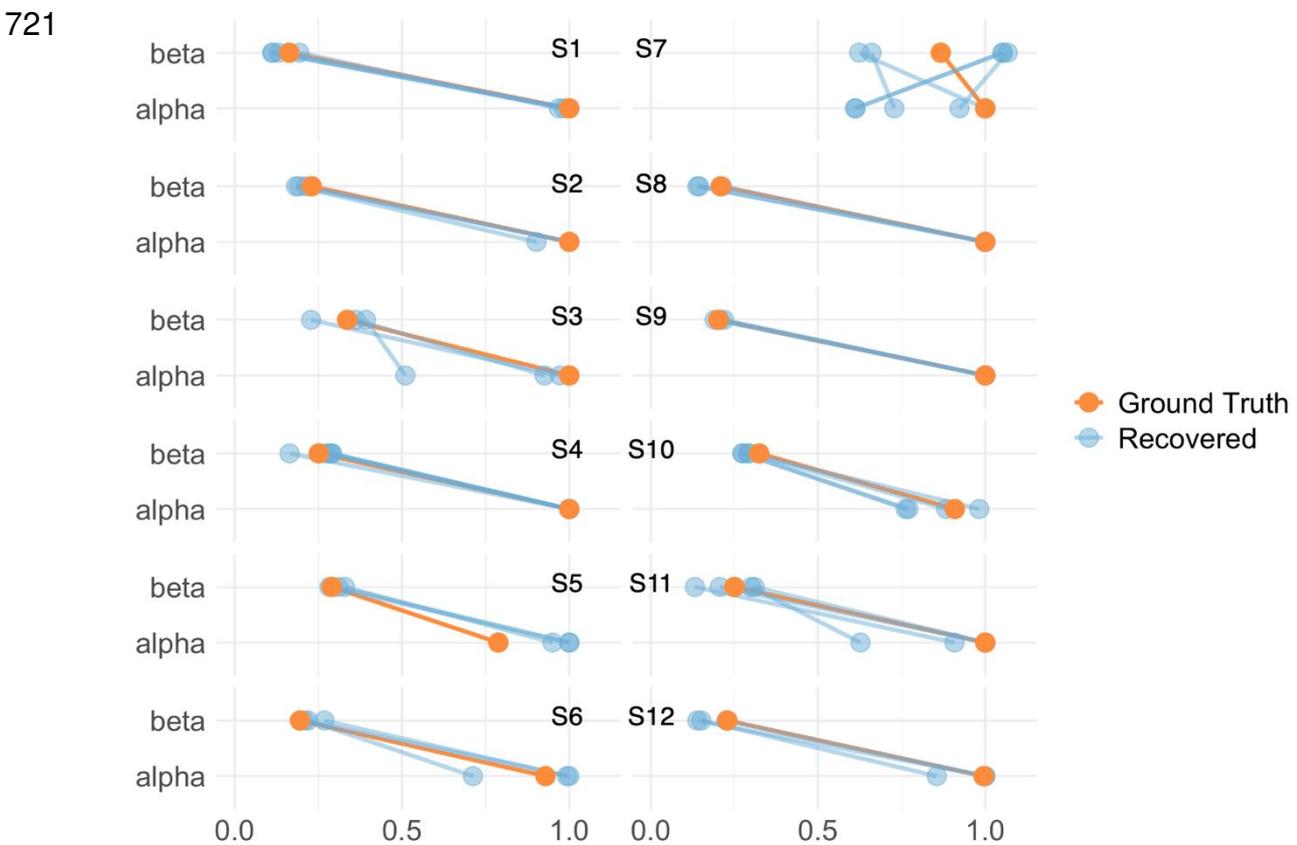


Figure S2 – Ground-truth alongside best-fitting parameters to data simulated from the ground-truth parameters. We simulated three datasets for each subject from our hybrid model with the best-fitting parameters to each subject’s empirical data, and then attempted to recover those ground-truth best-fitting parameters three times (corresponding to the three blue lines per participant).

722 Acknowledgments

723 We thank Neil Burgess, Jesse Geerts, and Talfan Evans for helpful and insightful
724 conversations during the early stages of this research. The research was funded by
725 a UCL departmental studentship to GW. SP is supported by a research grant funded
726 by the Ministry of Science and Higher Education of the Russian Federation (grant ID:
727 075-15-2020-928) and the French National Agency of Research (ANR; FrontCog
728 ANR-17-EURE-0017).

729

730 References

- 731 Abel, D. (2019). A theory of state abstraction for reinforcement learning. *33rd AAAI*
732 *Conference on Artificial Intelligence, AAAI 2019, 31st Innovative Applications of*
733 *Artificial Intelligence Conference, IAAI 2019 and the 9th AAAI Symposium on*
734 *Educational Advances in Artificial Intelligence, EAAI 2019.*
735 <https://doi.org/10.1609/aaai.v33i01.33019876>
- 736 Andre, D., & Russell, S. J. (2002). State abstraction for programmable reinforcement
737 learning agents. *Proceedings of the National Conference on Artificial*
738 *Intelligence.*
- 739 Botvinick, M. M. (2007). Multilevel structure in behaviour and in the brain: a model of
740 Fuster's hierarchy. *Philosophical Transactions of the Royal Society B: Biological*
741 *Sciences*, 362(1485), 1615–1626. <https://doi.org/10.1098/rstb.2007.2056>
- 742 Botvinick, M. M., Niv, Y., & Barto, A. C. (2009). Hierarchically organized behavior
743 and its neural foundations: A reinforcement learning perspective. *Cognition*,
744 113(3), 262–280. <https://doi.org/10.1016/j.cognition.2008.08.011>
- 745 Botvinick, M. M., Weinstein, A., Solway, A., & Barto, A. (2015). Reinforcement
746 learning, efficient coding, and the statistics of natural tasks. *Current Opinion in*
747 *Behavioral Sciences*, 5, 71–77. <https://doi.org/10.1016/j.cobeha.2015.08.009>
- 748 Bullock, D. (2004). Adaptive neural models of queuing and timing in fluent action.
749 *Trends in Cognitive Sciences.* <https://doi.org/10.1016/j.tics.2004.07.003>
- 750 Bullock, D., & Rhodes, B. (2003). Competitive queuing for planning and serial
751 performance. *The Handbook of Brain Theory and Neural Networks.*
- 752 Cooper, R., & Shallice, T. (2000). CONTENTION SCHEDULING AND THE
753 CONTROL OF ROUTINE ACTIVITIES. *Cognitive Neuropsychology*, 17(4), 297–
754 338. <https://doi.org/10.1080/026432900380427>
- 755 Diuk, C., Tsai, K., Wallis, J., Botvinick, M. M., & Niv, Y. (2013). Hierarchical learning
756 induces two simultaneous, but separable, prediction errors in human basal
757 ganglia. *The Journal of Neuroscience : The Official Journal of the Society for*
758 *Neuroscience*, 33(13), 5797–5805. [https://doi.org/10.1523/JNEUROSCI.5445-](https://doi.org/10.1523/JNEUROSCI.5445-12.2013)
759 12.2013

- 760 Doll, B. B., Simon, D. A., & Daw, N. D. (2012). The ubiquity of model-based
761 reinforcement learning. *Current Opinion in Neurobiology*.
762 <https://doi.org/10.1016/j.conb.2012.08.003>
- 763 Fuster, J. M. (2008). The Prefrontal Cortex. In *The Prefrontal Cortex*.
764 <https://doi.org/10.1016/B978-0-12-373644-4.X0001-1>
- 765 Glimcher, P. W. (2011). Understanding dopamine and reinforcement learning: The
766 dopamine reward prediction error hypothesis. *Proceedings of the National
767 Academy of Sciences of the United States of America*.
768 <https://doi.org/10.1073/pnas.1014269108>
- 769 Humphreys, G. W., & Forde, E. M. E. (1998). Disordered action schema and action
770 disorganisation syndrome. *Cognitive Neuropsychology*.
- 771 Koechlin, E., Ody, C., & Kouneiher, F. (2003). The Architecture of Cognitive Control
772 in the Human Prefrontal Cortex. *Science*.
773 <https://doi.org/10.1126/science.1088545>
- 774 Kornysheva, K., Bush, D., Meyer, S. S., Sadnicka, A., Barnes, G., & Burgess, N.
775 (2019). Neural Competitive Queuing of Ordinal Structure Underlies Skilled
776 Sequential Action. *Neuron*, 101(6), 1166-1180.e3.
777 <https://doi.org/10.1016/j.neuron.2019.01.018>
- 778 Lampinen, A. K., Chan, S. C. Y., Banino, A., & Hill, F. (2021). *Towards mental time
779 travel: a hierarchical memory for reinforcement learning agents*. Retrieved from
780 <http://arxiv.org/abs/2105.14039>
- 781 Lashley, K. (1951). The problem of serial order in behavior. In *Cerebral mechanisms
782 in behavior*.
- 783 Mehlhorn, K., Newell, B. R., Todd, P. M., Lee, M. D., Morgan, K., Braithwaite, V. A.,
784 ... Gonzalez, C. (2015). Unpacking the exploration-exploitation tradeoff: A
785 synthesis of human and animal literatures. *Decision*.
786 <https://doi.org/10.1037/dec0000033>
- 787 Miller, G. A., Galanter, E., & Pribram, K. H. (2017). Plans and the structure of
788 behaviour. In *Systems Research for Behavioral Science: A Sourcebook*.
789 <https://doi.org/10.2307/411065>
- 790 Niv, Y. (2009). Reinforcement learning in the brain. *Journal of Mathematical
791 Psychology*. <https://doi.org/10.1016/j.jmp.2008.12.005>
- 792 Radulescu, A., Niv, Y., & Ballard, I. (2019). Holistic Reinforcement Learning: The
793 Role of Structure and Attention. *Trends in Cognitive Sciences*.
794 <https://doi.org/10.1016/j.tics.2019.01.010>
- 795 Ramkumar, P., Acuna, D. E., Berniker, M., Grafton, S. T., Turner, R. S., & Kording,
796 K. P. (2016). Chunking as the result of an efficiency computation trade-off.
797 *Nature Communications*. <https://doi.org/10.1038/ncomms12176>
- 798 Rhodes, B. J., Bullock, D., Verwey, W. B., Averbeck, B. B., & Page, M. P. A. (2004).
799 Learning and production of movement sequences: Behavioral,
800 neurophysiological, and modeling perspectives. *Human Movement Science*.
801 <https://doi.org/10.1016/j.humov.2004.10.008>
- 802 Ribas-Fernandes, J. J. F., Solway, A., Diuk, C., McGuire, J. T., Barto, A. G., Niv, Y.,
803 & Botvinick, M. M. (2011). A Neural Signature of Hierarchical Reinforcement

804 Learning. *Neuron*, 71(2), 370–379. <https://doi.org/10.1016/j.neuron.2011.05.042>

805 Rosenbaum, D. A., Kenny, S. B., & Derr, M. A. (1983). Hierarchical control of rapid
806 movement sequences. *Journal of Experimental Psychology: Human Perception*
807 *and Performance*. <https://doi.org/10.1037/0096-1523.9.1.86>

808 Sakai, K., Kitaguchi, K., & Hikosaka, O. (2003). Chunking during human visuomotor
809 sequence learning. *Experimental Brain Research*.
810 <https://doi.org/10.1007/s00221-003-1548-8>

811 Saputro, D. R. S., & Widyaningsih, P. (2017). Limited memory Broyden-Fletcher-
812 Goldfarb-Shanno (L-BFGS) method for the parameter estimation on
813 geographically weighted ordinal logistic regression model (GWOLR). *AIP*
814 *Conference Proceedings*. <https://doi.org/10.1063/1.4995124>

815 Schultz, W., Dayan, P., & Montague, P. R. (1997). A neural substrate of prediction
816 and reward. *Science*. <https://doi.org/10.1126/science.275.5306.1593>

817 Shima, K., Isoda, M., Mushiake, H., & Tanji, J. (2007). Categorization of behavioural
818 sequences in the prefrontal cortex. *Nature*, 445(7125), 315–318.
819 <https://doi.org/10.1038/nature05470>

820 Solway, A., Diuk, C., Córdova, N., Yee, D., Barto, A. G., Niv, Y., & Botvinick, M. M.
821 (2014). Optimal Behavioral Hierarchy. *PLoS Computational Biology*.
822 <https://doi.org/10.1371/journal.pcbi.1003779>

823 Stolle, M., & Precup, D. (2002). Learning options in reinforcement learning. *Lecture*
824 *Notes in Computer Science (Including Subseries Lecture Notes in Artificial*
825 *Intelligence and Lecture Notes in Bioinformatics)*. [https://doi.org/10.1007/3-540-](https://doi.org/10.1007/3-540-45622-8_16)
826 [45622-8_16](https://doi.org/10.1007/3-540-45622-8_16)

827 Sutton, R. S. (1988). Learning to Predict by the Methods of Temporal Differences.
828 *Machine Learning*. <https://doi.org/10.1023/A:1022633531479>

829 Sutton, R. S., & Barto, A. G. (1998). Reinforcement Learning: An Introduction. *IEEE*
830 *Transactions on Neural Networks*. <https://doi.org/10.1109/tnn.1998.712192>

831 Sutton, R. S., Precup, D., & Singh, S. (1999). Between MDPs and semi-MDPs: A
832 framework for temporal abstraction in reinforcement learning. *Artificial*
833 *Intelligence*, 112(1–2), 181–211. [https://doi.org/10.1016/S0004-3702\(99\)00052-](https://doi.org/10.1016/S0004-3702(99)00052-1)
834 [1](https://doi.org/10.1016/S0004-3702(99)00052-1)

835 Watkins, C. J. C. H., & Dayan, P. (1992). Q-learning. *Machine Learning*.
836 <https://doi.org/10.1007/bf00992698>

837 Yokoi, A., & Diedrichsen, J. (2019). Neural Organization of Hierarchical Motor
838 Sequence Representations in the Human Neocortex. *Neuron*, 1–13.
839 <https://doi.org/10.1016/j.neuron.2019.06.017>

840