



18 as well as cleaner and client fish densities from the locations of cap-  
19 ture. Using Bayesian statistics to fit the model parameters to per-  
20 formance data revealed that cleaner fish most likely estimate future  
21 consequences of an action, while it appears unlikely that the removal  
22 of the ephemeral reward acts as psychological punishment (negative  
23 reinforcement). Incorporating future consequences also yields perfor-  
24 mances that can be considered the result of locally optimal decision-  
25 rules, in contrast to the negative reinforcement mechanism. We argue  
26 that the combination of computational models with data is a powerful  
27 tool to infer the mechanistic underpinnings of cognitive performance.

## 28 **Lay summary**

29 Cleaner fish eat ectoparasites off other fishes, so-called clients. It regularly  
30 happens that two clients seek a cleaner’s service simulatenously. Cleaners  
31 benefit from prioritising clients unwilling to wait, so they can feed on those  
32 willing to wait. To make the right choice, cleaners must somehow “look”  
33 into the future to anticipate consequences of current choices. By combining  
34 a learning model with data, we show that cleaners estimate the long-term  
35 value of their actions rather than using simpler heuristics. Estimating long-  
36 term value is a mechanism involved in human foresight.

37 **Keywords**— learning, behavior, cleaners, bayesian statisitics, behavioral  
38 mechanisms

## 39 Introduction

40 Often alternative cognitive mechanisms yield similar behavior and/or cog-  
41 nitive performances. This poses a problem for disentangling the mechanistic  
42 underpinnings of behavior. This is particularly clear in research aimed at  
43 discovering between species variation in *higher* cognitive abilities; or in other  
44 words, research on whether non-human animals show cognitive abilities be-  
45 lieved to be uniquely human. For instance, when researchers try to find  
46 *mental time travel-like* behavior, they usually come up with experiments to  
47 show the behavior displayed requires inferences made through past events  
48 (Dally, Emery, and Clayton 2006). However, they often face the challenge  
49 of alternative scenarios where simpler explanations, like classic associative  
50 learning, can bring about the observed behavioral outcome (Suddendorf and  
51 Corballis 2007). Similarly, attempts to demonstrate the presence of *theory*  
52 *of mind* in non-human animals face objections justified by alternative mech-  
53 anisms underpinning similar behavioral results (C. M. Heyes 1998). Such  
54 controversies are usually settled by using the principle of parsimony and its  
55 cognitive version, Lloyd Morgan’s cannon, which states that the simpler ex-  
56 planation (mechanism) should be accepted. Ideally, alternative hypotheses  
57 should be evaluated in light of their explanatory power.

58 Learning is a key overarching cognitive mechanism that allows individuals  
59 to associate rewards with environmental stimuli and thus behave adaptively  
60 (Staddon 2016; Shettleworth 2009). Associative learning, in particular, exists

61 in all major vertebrate taxa, and in many invertebrates as well (C. Heyes 2012;  
62 Macphail 1982; Staddon 2016; Behrens et al. 2008). Associative learning is  
63 not homogeneous throughout its taxonomic distribution, rather there are  
64 differences across and within species (Shettleworth 2009; Sih and Giudice  
65 2012). So presumably, the mechanistic underpinnings of learning have been  
66 modified by natural selection (Marler and Peters 1989).

67 One way to formalize the alternative mechanistic underpinnings of associa-  
68 tive learning is to develop quantitative models of learning processes. This ap-  
69 proach, which started within experimental psychology (Staddon 2016; Bran-  
70 don, Vogel, and Wagner 2002), has been very fruitful in disentangling the  
71 mechanistic structure of cognitive systems. More recently, the development  
72 of reinforcement learning theory (Sutton and Barto 2018), has allowed to  
73 evaluate the empirical support for alternative mechanistic hypotheses by  
74 providing quantitative predictions which are amenable to statistical tests  
75 (Farashahi et al. 2020, 2017). Interestingly, these learning models not only  
76 have received support from behavioral data, but also are consistent with the  
77 current view on reward processing in the brain (Schultz 2015).

78 From an evolutionary perspective, mechanisms are likely selected because  
79 of how they allow individuals to respond to environmental variation. For  
80 example, biological market theory predicts that the exchange rate of goods  
81 and/or services traded between cooperative partners adjusts to the law of sup-  
82 ply and demand, when individuals have some degree of partner choice (Noë

83 and Hammerstein 1995). Supply and demand conditions, which typically  
84 depend on the abundance of the species involved, certainly vary in time and  
85 space. Therefore, natural selection should favor the ability to flexibly adjust  
86 decisions and behavioral output to current market conditions. Indeed, such  
87 adjustments have been documented (Axén, Leimar, and Hoffman 1996). One  
88 example of strategic adjustment in a biological market is the marine cleaning  
89 mutualism involving the cleaner fish *Labroides dimidiatus* and ‘client’ fishes.  
90 Client fishes seek cleaner fish services at their territory (so-called “cleaning  
91 station”) and offer themselves as food patches to get their ectoparasites re-  
92 moved, which provides cleaners with food and clients with improved health  
93 (Waldie et al. 2011; Ros et al. 2011; Triki et al. 2016; Demairé et al. 2020).  
94 Given the capacity of some client fish to swim larger distances and access mul-  
95 tiple cleaning stations while others access the only cleaning station in their  
96 territory, it is crucial to categorize clients as either “visitors” or “residents”,  
97 respectively. During cleaning interactions, a cleaner fish often faces a choice  
98 between a visitor and a resident client seeking its cleaning services simultane-  
99 ously. Visitors have the option to switch to another cleaner fish if being made  
100 to wait, while residents must wait for inspection. Indeed, visitors have been  
101 observed to use their partner choice option in that way (Bshary and Schäffer  
102 2002), which may explain why cleaners give visitors service priority in a field  
103 study in the Red Sea (Bshary 2001). Furthermore, in a lab based paradigm,  
104 design to mimic the resident-visitor choice (ephemeral reward task), cleaners  
105 learned to prefer the cue associated with the epheral food source (visitor)

106 and hence accessed both food sources, obtaining double the amount of food  
107 (Bshary and Grutter 2002). However, further exploration revealed that not  
108 all cleaner fish manage to develop a preference for the ephemeral option in  
109 the lab (Triki et al. 2018). Over the last decade, over a hundred wild-caught  
110 cleaner fish have been tested in the exact same paradigm of the ephemeral  
111 reward task (Salwiczek et al. 2012; Wismer et al. 2014; Triki et al. 2018,  
112 2019, 2020). These fish often come from different reef locations. Investigation  
113 of the local eco-sociological conditions revealed that cleaner and client fish  
114 population densities have a substantial impact on cleaner fish performance  
115 in the task. Cleaner fish from reef sites with relatively low densities were  
116 more likely to fail at solving the task (Triki et al. 2019, 2018; Wismer et  
117 al. 2014). This intra-specific variation is unlikely due to local genetic adap-  
118 tation, because cleaner fish are open water spawners and the environmental  
119 conditions can vary within the lifespan of a fish.

120 Mechanistic models explicitly designed to mimic the ephemeral reward task  
121 have shown that the simplest form of associative learning (operant condi-  
122 tioning) cannot account for a solution to the ephemeral reward task (Prat,  
123 Bshary, and Lotem 2022; Quiñones et al. 2019). Operant conditioning is  
124 a form of associative learning where individuals use short term reward to  
125 associate and choose actions. Such models allow varying the cognitive tool  
126 kit and evaluating which minimal kit is necessary to solve the task at hand  
127 (e.g. Dubois et al. (2021)). To be able to give visitors priority over resi-  
128 dents, cleaners need to be able to assess a client’s value separately for the

129 three possible scenarios (alone, paired with a fish with the same strategic  
 130 option, paired with a fish with the alternative strategic option) (Quiñones  
 131 et al. 2019). The ability to distinguish and value one stimulus differently  
 132 alone from compound versions of it has been termed configurational learning,  
 133 chunking, or segmentation (see references in Prat, Bshary, and Lotem (2022)).  
 134 In addition to configurational learning, cleaners also need to account for the  
 135 future consequences of current decisions. In the model by Quiñones et al.  
 136 (2019), this could be achieved in two non-mutually exclusive ways: through  
 137 low temporal discounting of future effects, also termed ‘chaining’ (Enquist,  
 138 Lind, and Ghirlanda 2016); and/or through perceiving a visitor client leav-  
 139 ing as psychological punishment (i.e. as a negative reinforcer). Chaining is  
 140 when individuals include in their valuation of an action the reward effects  
 141 that this will have in the future. This is done by combining in a single valua-  
 142 tion the reward obtained in the current time with all the reward that comes  
 143 after, discounting for how far in the future reward is accrued. ‘Chaining’  
 144 the reward of these different time steps allows individuals to take actions  
 145 that increase the long-term reward at the sacrifice of short term considera-  
 146 tions (Enquist, Lind, and Ghirlanda 2016). Even though, ‘chaining’ can be  
 147 readily implemented computationally in learning models (Enquist, Lind, and  
 148 Ghirlanda 2016; Sutton and Barto 2018), cognitively it seems to be a com-  
 149 plex adaptation (Suddendorf and Corballis 2007). On the other hand, using  
 150 client behavior as a negative reinforcer is, in principle, easier to implement.  
 151 Thus, the standard logic of Lloyd Morgan’s cannon demands that operant

152 conditioning as the simpler explanation is to be accepted by default. Ide-  
153 ally, however, the two mechanisms should be evaluated in light of how well  
154 they explain the available data. Note that different fields interested in cogni-  
155 tion and decision making use different words to refer to negative reinforcers  
156 (Quiñones et al. 2019; Sutton and Barto 2018). Here, for the sake of sim-  
157 plicity and clarity, we will use the word ‘penalty’ to refer to this mechanism  
158 which includes a negative reinforcer.

159 In here we used the field and experimental data to fit the parameters of a  
160 reinforcement learning model to infer the cognitive mechanism that cleaners  
161 use in their interaction with clients. Specifically, our approach of fitting the  
162 computational model to the empirical data aimed at: (i), determining which  
163 mechanism cleaner fish use to incorporate future consequences of current de-  
164 cisions by testing whether chaining, penalty, or a combination of both best  
165 explains their performance; (ii) determining whether the two mechanisms  
166 differed with respect to the ecological conditions that are likely to cause high  
167 versus low performance in the ephemeral reward task. Additionally, we as-  
168 sessed which mechanism yields optimal performance patterns. Relying on the  
169 logic of biological market theory, we predicted that appropriate performance  
170 is to show a high preference for visitors only under high local cleaner-to-client  
171 ratio.



## 172 **Methods**

### 173 **The model**

174 The model consists of a set of individual-based simulations where individuals,  
175 representing cleaner fish, face a series of choices between two options, which  
176 simulate the natural conditions of the cleaning market. Individuals experi-  
177 ence a series of discrete time points in which they face different ‘states’, de-  
178 fined by the number and category of client fish (visitor or residents) inviting  
179 for cleaning services. There are six possible states: zero clients, one resi-  
180 dent, one visitor, resident-resident, visitor-visitor, and resident-visitor. The  
181 probability of each state is largely determined by the relative abundance of  
182 cleaner fish, residents and visitors, but to some degree by cleaner fish choices  
183 when it faces the resident-visitor combination. This is because residents are  
184 willing to queue for cleaning service; while visitors leave the queue (with a  
185 certain probability) when made to wait. Individuals obtain a fixed reward  
186 from cleaning a client fish regardless of the category. Every time individuals  
187 face and make a choice they update the probability of making that same  
188 choice. The update is based on the difference between the expected value  
189 and the obtained reward - the prediction error ( $\delta_t$ ) - (Sutton and Barto 2018;  
190 Rescorla and Wagner 1972). Formally, the prediction error is given by

$$\delta_t = R_t - V_t(S_t) + \gamma V_t(S_{t+1}), \quad (1)$$

191 where  $R_t$  is the sum different reward sources at time  $t$ ;  $V_t(S_t)$  is the estimated  
 192 value at time  $t$  of the the state faced by the agent at time  $t$ ; similarly  $V_t(S_{t+1})$   
 193 is the estimated value of the state to come in the following time-step,  $\gamma$  is the  
 194 discount factor for future rewards. When the estimated value of the current  
 195 state ( $V_t(S_t)$ ) is equal to the sum of short-term ( $R_t$ ) and future discounted  
 196 reward ( $\gamma V_t(S_{t+1})$ ) learning stops for that state. If  $\gamma = 0$  the estimates  
 197 made by the agent only capture short-term reward. We assume short-term  
 198 reward to two components: positive reward determined by the amount of  
 199 food obtained from cleaning a client; and negative reward triggered when  
 200 by a client leaving the station without being cleaned. Formally, we let total  
 201 reward be given by  $R_t = P_t - \eta_t$ . Where  $\eta$  is a parameter of the model  
 202 that determines the the size of the negative reward triggered by unattended  
 203 clients leaving the station.

204 The prediction error (Eq. ??) is used to update the value of each one of the  
 205 states the agent faces, as well as the preference for the resident and visitor  
 206 options. The value update is simple the product of the prediction error  
 207 and the parameter for the speed of learning ( $\Delta V(S_t) = \alpha \delta_t$ ). The change  
 208 in preference between the resident and visitor is given by  $\Delta(\theta_v - \theta_r)_t =$   
 209  $\alpha \delta_t 2(1 - \pi_v)$ , where  $\theta_i$  represent the preference for one of two options and the  
 210 difference captures the total change relative to one another;  $\pi_v$  corresponds to  
 211 the current probability of choosing the visitor.  $p_i$  is determined by applying  
 212 the logistic function to the difference in preferences between the two mutually  
 213 exclusive options ( $\pi_v = \frac{1}{1+e^{-(\theta_v-\theta_r)}}$ ). This amounts to a preference update

214 that is carried in the direction that leads to more reward being obtained,  
215 given the new information. In the long run, the probability of choosing a  
216 visitor over a resident converges in the model. To which probability the  
217 model will converge depends on the relative abundance of cleaners, visitors  
218 and residents; as well as on the probability of visitors leaving the cleaning  
219 station when made to wait. Further details of the model implementation can  
220 be found in Quiñones et al. (2019).

221 The model shows that agents need to find a way to incorporate future conse-  
222 quences of current choices. In the model, this could be achieved with either of  
223 two parameters that could also work together. First,  $\gamma$  measures how much  
224 individuals include future rewards in their decision updates. If  $\gamma = 0$ , indi-  
225 viduals only use the immediate reward obtained from a cleaning interaction.  
226 As  $\gamma$  increases, individuals include more of the reward obtained from the sub-  
227 sequent choices. That amounts to estimating and using for decision making  
228 the future expected rewards of an action (chaining). Second,  $\eta$  measures how  
229 much individuals include in their reward the fleeing behavior of visitors as  
230 a negative component (penalty). Both of these parameters allow individuals  
231 to use in their estimates the future effects of their choices.

## 232 **Empirical data**

233 The empirical data were collected between 2010 and 2019 always during the  
234 austral winter months June to August from a total of five study reef sites

235 (Corner Beach-CB, Horseshoe-HS, Mermaid Cove-MC, Northern Horseshoe-  
 236 NHS, and The Crest-TC) at Lizard Island ( $14.6682^{\circ}S, 145.4604^{\circ}E$ ), Great  
 237 Barrier Reef, Australia. The data consist of three sets: fish censuses, field  
 238 observations of cleaner-client interactions to quantify the probability of visi-  
 239 tors leaving if made to wait, and the performance of wild-caught cleaner fish  
 240 in the ephemeral reward test. In total, we have twelve site/year data sets for  
 241 fish censuses and corresponding performance in the lab test. Thus, some sites  
 242 were sampled more than once. To estimate the population density of cleaner  
 243 fish and their clients at a given site in a given year, Triki et al. (2019) used a  
 244 series of ten transects of  $30m$  each. Observers swam along the transect lines  
 245 placed on the reef and first counted the visible large-bodied adult fish (species  
 246 with total length  $TL \geq 10cm$ ) including cleaner fish on a width of  $5m$ , and  
 247 then on the return individuals of small-bodied fish species ( $TL < 10cm$ ) on  
 248 a width of  $1m$  (see Triki et al. (2018) for further details on fish censuses  
 249 data collection). Total length estimates were done by the observer. We then  
 250 scaled the counts of cleaner fish, small-bodied, and large-bodied clients fish  
 251 densities per  $100m^2$ .

252 The field observation data consisted of video recordings/encodings of the  
 253 cleaner-client cleaning interactions. There were videos from eight cleaners  
 254 per site/year of a duration of  $30min$  each. Triki et al. extracted information  
 255 from every event wherein a visitor client was made to wait in favor of another  
 256 client (visitor or resident), and noted whether or not the visitor left or queued  
 257 for the cleaning service (Triki et al. 2019, 2020).

258 The cognitive performance data was from a total of 120 cleaners (10 individ-  
 259 uals per 12 site/year) tested in the ephemeral reward task (Triki et al. 2019,  
 260 2020). Authors housed all captured cleaners individually in glass aquaria  
 261 ( $62\text{cm} \times 27\text{cm} \times 37\text{cm}$ ) and provided them with PVC pipes ( $10\text{cm} \times 1\text{cm}$ ) as  
 262 shelters. The task consisted of exposing the cleaner fish to substitute models  
 263 of client fish in the form of two *Plexiglas* plates offering the same amount  
 264 of food (one item of mashed prawn). The two plates differed in colour and  
 265 pattern (horizontal green stripes or vertical pink stripes) but had equal size  
 266 ( $10\text{cm} \times 7\text{cm}$ ). Importantly, the two plates played different roles as either a  
 267 visitor (ephemeral food source) or resident (permanent food source). That is,  
 268 if a cleaner fish inspected the resident plate first, the experimenter withdrew  
 269 the visitor plate out of the aquarium as a consequence. Choosing first the  
 270 visitor plate, however, granted access to both plates. The equal size of the  
 271 plates forced cleaner fish to decide based solely on the association between  
 272 the behaviour and the collar/pattern cue (Wismer et al. 2019). Triki et al.  
 273 (2019, 2020) tested the fish for a maximum of 200 trials with 20 trials a day,  
 274 10 trials in the morning and 10 trials in the afternoon. They randomized  
 275 and counterbalanced the plates' spatial location (i.e. left or right) between  
 276 trials. Similarly, they counterbalanced the plates' decoration (colour and pat-  
 277 tern) and the plates' role (visitor or resident) between the tested fish. In the  
 278 original studies, once a fish reached a learning criterion, that is, performing  
 279 significantly above chance level in a binomial test ( $p - \text{value} \leq 0.05$ ), they  
 280 passed to a reversal version of the task where the roles of the visitor/resident

281 Plexiglas plates were swapped. The reversal phased stopped when the fish  
282 performed significantly above chance, or the fish completed the 200 trial to-  
283 gether with the initial; see Triki et al.(2019, 2020). Here, we are interested in  
284 explaining the total frequency of visitor choices using the model, rather than  
285 just the achievement of the criterion. Total frequency of visitor choices nat-  
286 urally comes out of the model, and allows us to use all the variation among  
287 cleaners, instead of reducing that to a binomial variable. Thus, we used  
288 instead a subset of these data to estimate the final cleaner fish preferences  
289 for the visitor plate, even if they do not reach the learning criterion within  
290 200 trials. To do so, we first extracted the trial-by-trial outcomes from the  
291 last two sessions (20 trials) of those who never reached the learning criterion  
292 for visitor plate ( $N = 45$  cleaner fish). For those who reached the learning  
293 criterion at some point during the test and passed to a reversal phase, we  
294 extracted the trial-by-trial outcomes from the last session (10 trials) before  
295 passing to reversal and the last 10 trials they were exposed to in the test  
296 ( $N = 75$  cleaner fish).

297 We chose a combination of initial and reversal to quantify preference for the  
298 visitor client. However, it could be argued that using only the initial phase  
299 gives a better estimation of the cleaner fish preference for the visitor. In  
300 the supplementary material (Fig. ??) we show how using initial and reversal  
301 (a), and only initial (b) maps to the previously used criteria. The initial and  
302 reversal match better the criteria chosen in previous analysis of the ephemeral  
303 reward task experimental set-up (Triki et al. 2019, 2020).

## 304 Statistical analysis

305 The aim of the analyses is to fit the key model parameters  $\gamma$  and  $\eta$ , to the  
306 empirical data from Triki et al. (2019, 2020) to test whether each or a  
307 combination of these effects is a better explanation for the pattern seen in  
308 the data. We used the ecological variables: cleaners, visitor clients, resident  
309 clients abundances and visitor clients leaving probability,  
310 as input to the models. As the response variable, we used the frequency  
311 with which cleaners chose the visitor option in the ephemeral reward task.  
312 Finally, we used the probability with which agents in the model simulations  
313 choose the visitor in resident-visitor options as the prediction for the response  
314 variable. We kept all other parameter values used for the model simulations  
315 constant, see Table ??.

316 To capture with the model the relationship between the ecological variables  
317 and cleaner fish preferences for visitors, we needed to scale the absolute popu-  
318 lation densities of cleaner fish from the empirical data to a measure of relative  
319 abundances that captures client visitation patterns. This is because, in the  
320 model, relative abundances of clients define not only the probability of resi-  
321 dents and visitors but also how often the cleaning station is empty (e.g. there  
322 are no clients to be cleaned). The frequency with which clients visit the sta-  
323 tion is another variable influencing the station occupancy, which in nature  
324 may vary among different client species depending on their ectoparasite loads.  
325 We do not have field estimates for species-specific parasite loads, especially