

# The formal foundation of an evolutionary theory of reinforcement

*Borgstede, M. \*, Eggert, F. \*\**

\*Corresponding author:

University of Bamberg

Markusplatz 3

D-96047 Bamberg

[matthias.borgstede@uni-bamberg.de](mailto:matthias.borgstede@uni-bamberg.de)

\*\*Technische Universität Braunschweig

Spielmannstrasse 19

D-38106 Braunschweig

[f.eggert@tu-braunschweig.de](mailto:f.eggert@tu-braunschweig.de)

# 16 The formal foundation of an evolutionary theory of 17 reinforcement

## 18 Abstract

19 Reinforcement learning is often described by analogy to natural selection. However, there is no  
20 coherent theory relating reinforcement learning to evolution within a single formal model of  
21 selection. This paper provides the formal foundation of such a unified theory. The model is based on  
22 the most general description of natural selection as given by the Price equation. We extend the Price  
23 equation to cover reinforcement learning as the result of a behavioral selection process within  
24 individuals and relate it to the principle of natural selection via the concept of statistical fitness  
25 predictors by means of a multilevel model of behavioral selection.

26 The main result is the *covariance-based law of effect*, which describes reinforcement learning on a  
27 molar level by means of the covariance between behavioral allocation and a statistical fitness  
28 predictor. We further demonstrate how this abstract principle can be applied to derive theoretical  
29 explanations of various empirical findings, like conditioned reinforcement, blocking, matching and  
30 response deprivation.

31 Our model is the first to apply the abstract principle of selection to derive a unified description of  
32 reinforcement learning and natural selection within a single model. It provides a general analytical  
33 tool for behavioral psychology in a similar way that the theory of natural selection does for  
34 evolutionary biology. We thus lay the formal foundation of a general theory of reinforcement as the  
35 result of behavioral selection on multiple levels.

36 **Keywords:** *selection by consequences, behavioral selection, natural selection, reinforcement learning,*  
37 *Price equation, multilevel model of behavioral selection*

## 38 1 Introduction

39 It is a long held belief that reinforcement learning can be characterized by analogy to evolution by  
40 natural selection (e.g., Broadbent, 1961; D. T. Campbell, 1956; Gilbert, 1970; Herrnstein, 1964;  
41 Pringle, 1951; Skinner, 1966; Thorndike, 1900). For example, Staddon and Simmelhag (1971), state  
42 that the ‘Law of Effect [...] can best be understood by analogy with evolution by means of natural  
43 selection’ (p. 40). Skinner (1981) even claims that natural selection and reinforcement learning are  
44 two instances of the same underlying causal principle: *selection by consequences*. The selectionist  
45 account of reinforcement has also found its way into textbooks on behavioral psychology as the  
46 *Darwinian metaphor* (Baum, 2005; Staddon, 2016) and is now a popular theme in theoretical work on  
47 behavior analysis (e.g., Baum, 2017, 2018; Becker, 2019; Borgstede, 2020; Donahoe, 2011; Donahoe,  
48 Burgos, & Palmer, 1993; Hull, Langman, & Glenn, 2001; Richerson, 2019; Simon & Hesse, 2019).

49 The appeal of the Darwinian metaphor conceivably stems from its generality: given selection by  
50 consequences is a fundamental principle of behavior, it might constitute the foundation of a unified  
51 theory of behavior. Such a theory would provide a general analytical framework for behavioral  
52 psychology in a similar way that the theory of natural selection does for evolutionary biology. A  
53 theoretical description of fundamental behavioral principles that goes ‘beyond the collection of  
54 uniform relationships’ (Skinner, 1950, p.215) could offer a theory-driven explanation of the basic  
55 laws of learning. This would help to understand why the regularities in empirical findings of  
56 behavioral psychology – for instance, the matching law – are to be expected and might even  
57 generate new testable hypotheses.

58 However, the theoretical status of ‘selection by consequences’ has been subject to criticism (e.g., the  
59 open peer commentaries to Skinner, 1984; also Burgos, 2019; Pennypacker, 1992; Tonneau &  
60 Sokolowski, 2000 for more recent accounts). Apart from an ongoing debate about the ontological  
61 status of the Darwinian metaphor, there are three recurring questions concerning the adequacy of  
62 the analogy between natural selection and reinforcement learning: the first one addresses the

63 hereditary particles of behavioral selection<sup>1</sup> (i.e., a gene analogue), the second one tackles the  
64 problem of a fitness equivalent on the behavioral level (i.e., some kind of reinforcer value), and the  
65 third one is about the units of selection (i.e., the question of what is selected). While the question of  
66 a gene analogue for behavioral selection is surely interesting (cf. Donahoe et al., 1993), it is not  
67 essential for selectionist theory, because it is possible to describe evolutionary change on a  
68 phenotypic level without loss of generality (Frank, 1997, 1998; Grafen, 2014). The second and third  
69 question, however, are crucial for the Darwinian metaphor to make sense – if there is no fitness  
70 equivalent on a behavioral level, there is no criterion for selection, and if the units of selection are  
71 unclear, we do not know what is selected.

72 The question of a behavioral fitness equivalent has been clarified from a maximization perspective  
73 (Borgstede, 2020): if there is a behavioral maximand (‘reinforcer value’) that reinforcement selects  
74 for, and if maximization of this value leads to maximization of evolutionary fitness, reinforcer value  
75 must be proportional to marginal fitness (i.e., fitness change per unit change in behavioral  
76 allocation). However, Borgstede (2020) does not link the maximization principle to the dynamics of  
77 change. Hence, whilst providing a valid mathematical definition of reinforcer value, the implications  
78 for reinforcement learning as a process of selection by consequences remain open.

79 The issue of the units of selection in the context of reinforcement learning has been addressed  
80 explicitly by McDowell (2013), who models reinforcement learning by means of an evolutionary  
81 algorithm that is applied to a population of ‘behaviors’. In this view, behaviors relate to learning in  
82 the same way that individuals relate to evolution. However, in the context of learning, behavior is  
83 often treated as the *target* of selection, as well (Hull et al., 2001). Hence the analogy between  
84 learning and evolution is at least vague in this respect: do we conceptualize learning as selection *of*  
85 behaviors or as selection *for* behaviors (cf. Sober, 1984)?

---

<sup>1</sup> In this paper, we use the term ‘behavioral selection’ exclusively for behavioral adaptations by means of reinforcement learning. We do *not* refer to behavior being the target of evolution by natural selection.

86 In this paper, we aim to resolve the conceptual ambiguities of the analogy between learning and  
87 evolution by formally integrating reinforcement learning with natural selection in a unifying model  
88 that captures both levels of selection simultaneously. We build our model around the most abstract  
89 description of selection by means of the Price equation (Price, 1970, 1972). The main result is a molar  
90 account of reinforcement learning in terms of selection that applies to different levels on different  
91 time scales. In this view, reinforcement learning can be described as a Darwinian process where the  
92 units of selection are individuals showing behavioral variability and the target of selection is the  
93 relative allocation of behavior over time within a specified context. This process is universal in that it  
94 does not depend on the specific (molecular) mechanisms involved in learning but constitutes a  
95 general invariance principle which we call the *covariance based law of effect*. The dynamics of  
96 reinforcement are thus re-conceptualized in a molar way, shifting the focus from contiguity between  
97 single behavioral instances (Thorndike, 2010/1911) to the correlation between behavior and  
98 reinforcement (Baum, 1973). When applied to different experimental paradigms, the covariance  
99 based law of effect explains why *conditioned reinforcers* work (Skinner, 1969), why conditioning is  
100 sometimes *blocked* by previous reinforcement (Kamin, 1969), why *response deprivation* can establish  
101 reinforcers (Timberlake & Allison, 1974), and why individuals tend to *match* relative behavioral  
102 allocation to relative reinforcement in concurrent variable interval schedules (Herrnstein, 1961),  
103 Our model integrates these empirically well-established regularities in a mathematically rigorous way  
104 by means of a single theoretical principle that is derived from the theory of evolution by natural  
105 selection. We thus lay the foundation for a general theory of behavior that is of 'greater generality  
106 than any assemblage of facts' (Skinner, 1950, p.216).

## 107 2 Natural selection and the Price equation

108 The Price equation provides a mathematical description of evolutionary processes on the most  
109 general level by partitioning the change in mean character value from one generation to the next  
110 generation into a covariance term and an expectation term (Price, 1970)<sup>2</sup>:

$$\bar{w}\Delta\bar{z} = \text{Cov}(w_i, z_i) + E(w_i\Delta z_i) \quad (1.)$$

111 Here,  $z_i$  refers to an arbitrary character (usually an allele frequency) and  $\bar{z}$  designates the arithmetic  
112 mean of  $z$  over all individuals  $i$ . There is nothing special about gene frequencies here.  $z$  may be any  
113 quantitative character (e.g., body size or parental investment). Also,  $i$  does not necessarily refer to  
114 individuals, but can designate the members of an arbitrary set (e.g., groups). In order to define  $\Delta\bar{z}$   
115 one needs to relate the first set (the 'parent population') to a second set (the 'offspring population'):  
116  $\Delta\bar{z}$  is defined as  $\sum q'_i z'_i - \sum q_i z_i$ , where  $q_i$  refers to the frequency of the value  $z_i$  in the parent  
117 population and  $q'_i$  refers to the frequency of value  $z'_i$  in the offspring population. The index  $i$  does  
118 not refer to the individuals in the offspring population but to the parent population. This means that  
119  $q'_i$  is the number of elements in the offspring population that originate from parents of type  $i$  and  $z'_i$   
120 is their corresponding character value.

121 Mathematically, the Price equation builds on the existence of a *right-total relation* (i.e., a *mapping*)  
122 between two sets. Using this relation it is possible to define the fitness  $w_i$  as the contribution of a  
123 type  $i$  parent to the offspring population, resulting in  $q'_i = q_i w_i / \bar{w}$ , where  $\bar{w}$  is the mean fitness of  
124 the parent population.  $z'_i$  is also defined with respect to the parent population, which means that it  
125 refers to the average character value  $z$  of descendants from parent type  $i$ . The change in  $z$  from  
126 parent to offspring is defined accordingly:  $\Delta z_i = z'_i - z_i$ . The Price equation is valid given these  
127 definitions (mathematical proofs can be found in Frank, 1998, Gardner, 2020 and elsewhere). The  
128 change in mean character value can always be partitioned into one covariance term and one

---

<sup>2</sup> Price used a different notation in his original paper. However, the notation adapted in this paper has become more common in the literature. Compare Luque (2017) for a review of the many different versions of the Price equation.

129 expectation term. In biological models, the covariance term captures the change in character value  
130 due to natural selection, whereas the expectation term refers to changes from parent to offspring  
131 due to imperfect transmission or environmental factors.

132 It is possible to extend the formalism to capture different genetic architectures, class structured  
133 populations, stochasticity and inclusive fitness (e.g., Frank, 1998; Grafen, 2000; Taylor, 1990).  
134 Moreover, the second term can be further partitioned by inserting the Price equation recursively into  
135 the expectation term. Taking individuals to be nested within groups and adding a new index  $g$  for  
136 these groups, one may partition the change in mean character value within groups,  $z_g$ , into a within  
137 group covariance term and a within group expectation term:

$$w_g \Delta z_g = \text{Cov}(w_{gi}, z_{gi}) + E_i(w_{gi} \Delta z_{gi}) \quad (2.)$$

138 Here,  $z_{gi}$  and  $w_{gi}$  stand for the individual character value and fitness of individual  $i$  in group  $g$ .  $z_g$   
139 and  $w_g$  refer to mean character value and mean fitness in group  $g$ . Taking the expectation of mean  
140 change  $w_g \Delta z_g$  over groups, one can extend the Price equation to capture selection within groups and  
141 selection between groups at the same time:

$$\bar{w} \Delta \bar{z} = \text{Cov}(w_g, z_g) + E_g \left( \text{Cov}(w_{gi}, z_{gi}) + E_i(w_{gi} \Delta z_{gi}) \right) \quad (3.)$$

142 This multilevel Price equation is useful to model fitness trade-offs between the individual and the  
143 group, thereby explaining, how a trait that is harmful to the individual can spread in a population by  
144 positively affecting average fitness on a group level (Price, 1972).

### 145 3 The Price equation and behavioral selection

146 Price himself noted that his formal account of selection was not restricted to natural selection acting  
147 on gene frequencies but might well be applied to other areas such as operant learning (Price, 1995,  
148 written ca. 1971). The Price equation has been applied to such diverse fields as probability theory,  
149 particle physics and information theory (Frank, 2017, 2018, 2020). Especially the application of a  
150 selectionist framework to the field of information theory seems to imply that there might be an

151 intricate relationship between selection and learning. Nevertheless, whilst there have been various  
152 attempts to apply the Price equation to cultural evolution (see Nettle, 2020 for an overview), its  
153 potential for behavioral psychology and reinforcement learning in particular remained largely  
154 unnoticed until a recent publication by Baum (2017).

155 Baum (2017) identifies the objects in the ‘parent population’ with operant behavior in a fixed time  
156 interval, and the objects in the ‘offspring population’ with operant behavior in a later time interval of  
157 equal length. These time intervals correspond to different trials in a behavioral experiment, where  
158 the individual is repeatedly confronted with the same reinforcement contingencies. In order to  
159 construct the necessary ‘parent-offspring’ relation between the two sets he argues that behaviors  
160 *recur* in the sense that a behavior emitted in interval one may be emitted in interval two, as well.  
161 Following this rationale, recurrence in behavioral selection means ‘to occur again’. However, this  
162 understanding of recurrence departs from the conceptual foundation of the Price equation: ‘to occur  
163 again’ is by no means sufficient to establish a right-total relation between two sets of behavioral  
164 episodes. Baum further elaborates his position by linking reinforced behavior to Phylogenetically  
165 Important Events (PIE<sup>3</sup>). Following Baum, a PIE like, e.g., the availability of food induces PIE-related  
166 behavior like feeding and foraging. Thus, behaviors in interval one recur because they co-vary with a  
167 PIE, which in turn induces the same behaviors in interval two (Baum, 2012, 2017).

168 Even though induction possibly plays an important role for the allocation of behavior, it does not  
169 account for recurrence in the sense of a ‘parent-offspring’ relation as required by the Price equation.  
170 The indices in the Price equation always refer to the parent population – this means that  $z_i$   
171 designates the mean character value of objects *descending from parent type i*. Therefore, in order to  
172 apply the mathematical apparatus of the Price equation it is necessary that every object in the  
173 offspring population can be *individually linked* to an object in the parent population. Applying Baum’s  
174 claim that behaviors recur due to induction, we need to identify which behavioral instance in interval

---

<sup>3</sup> A PIE is an event that directly affects an individual’s evolutionary fitness, e.g., the availability of food or the presence of a physical threat (Baum, 2012).



175 one induces which behavioral instance(s) in interval two. But this is not possible. Therefore, the  
176 presented account of behavioral recurrence by means of induction remains unclear.

177 Even though Baum's model is formally consistent for some special cases (namely, when there are no  
178 sources of behavioral change apart from selection), it is impossible to retrieve the original meaning  
179 of 'fitness' as the contribution from one generation to the next generation because the necessary set  
180 mapping (a right-total relation) does not exist. Moreover, formally treating behavioral selection as  
181 analogous to natural selection does not provide a functional relation between both levels of  
182 selection, thereby missing the opportunity to integrate learning and evolution within the same  
183 model.

#### 184 4 A multilevel model of behavioral selection

185 As shown in the previous section, behaviors do not recur in the sense that one could establish a  
186 'parent-offspring' relation between behaviors occurring in one time interval and behaviors occurring  
187 in a future time interval. Therefore, it is not possible to derive a coherent definition of fitness from  
188 the recurrence of behaviors themselves. We solve this issue by not defining fitness on the level of  
189 single behaviors, but on the level of the whole organism (i.e., in the standard way the Price equation  
190 is applied in biological models of natural selection). Thus, we do not conceptualize behavioral  
191 selection as a process *similar* to natural selection, but as *a part* of natural selection itself. In this view,  
192 behavioral selection is not a 'Darwinian metaphor' but a theoretically derived scientific fact:  
193 behavioral selection is *literally* a part of evolution by natural selection as described by the Price  
194 equation.

195 If the objects used in the Price equation are whole organisms instead of behaviors, several challenges  
196 arise. First, natural selection and behavioral selection are usually taken to act on different time  
197 scales. This poses a formal problem since the Price equation only deals with the change from one  
198 generation to the next. Second, natural selection and behavioral selection are mediated by different  
199 mechanisms of inheritance (or recurrence). Evolving characters may be influenced by both,

200 genetically transmitted components, as well as learned components (for which the mechanisms of  
201 transmission are yet to be understood). Third, since fitness is defined on the level of the individual  
202 (i.e., fitness refers to the contribution of an individual to the future population), it needs to be  
203 clarified how evolutionary fitness relates to reinforcement in order to make sense of assigning values  
204 of evolutionary fitness to *behavior*. We will provide solutions to each of the above problems by  
205 integrating natural selection and behavioral selection using a multilevel extension of the Price  
206 equation and linking evolutionary fitness to reinforcement via the concept of statistical fitness  
207 predictors. The general idea is that reinforcement learning can only be effective to the degree that  
208 the average outcome of reinforcement (i.e., the learned behavior) contributes to expected  
209 evolutionary fitness (Borgstede & Simon, under review).

210 The formal aspects of our model are very general in that they apply to any kind of behavior and that  
211 they are independent of specific learning algorithms. In fact, we do not attempt to model molecular  
212 mechanisms at all. Instead, we adopt a molar view that treats behavior as being extended over time  
213 and over contexts (Rachlin, 1978). Hence, when we speak of behavior, we mean the relative  
214 allocation of competing behavioral options over time within a certain context, which in itself is  
215 defined by a certain structure of contingencies. In the course of reinforcement learning, behavior  
216 becomes controlled by context-specific discriminative stimuli. This molar view implies that we  
217 analyze behavior in terms of allocated time per interval while in a certain context. Therefore, we do  
218 not need to worry about the measurement units of specific behaviors, which may very well vary  
219 between behaviors. We call these context-specific intervals 'behavioral episodes' to stress the fact  
220 that behavior is not measured at a point in time but extended over time within a specified context. A  
221 behavioral episode is thus defined with respect to the structure of contingencies that are effective in  
222 a specified environmental context. From the perspective of behavioral psychology, one may think of  
223 behavioral episodes as the trials in a reinforcement experiment, exposing the individual to the same  
224 contingencies again and again. However, the concept can also be applied to settings outside the  
225 laboratory. Here, we may identify behavioral episodes with recurring contexts the animal

226 encounters, for example, different food patches, whose contingency structures are signaled by the  
227 presence of certain discriminative stimuli.

228 To keep the mathematical derivation as simple as possible, we restrict our analysis to learning within  
229 only one type of behavioral episode. In other words, we treat learning separately in different  
230 contextual settings. Thus, when we calculate the mean behavioral allocation over several episodes,  
231 these need not be adjacent in time, but instead are recurring instances of the same contingency  
232 structure. Similarly, when we speak of statistical fitness predictors, we refer to predictors that are  
233 valid within the current class of behavioral episodes (i.e., they are context-dependent). This is in line  
234 with reinforcement learning as well as extinction being context-specific<sup>4</sup>.

#### 235 4.1 The covariance based law of effect

236 The first step to integrate natural selection and behavioral selection is to partition the expectation  
237 term of the Price equation using recursive expansion (cf. section 2). However, here the two levels of  
238 selection are not individuals nested within groups but *behavioral episodes nested within individuals*.  
239 Like in the simple Price equation, we designate individuals by the index  $i$ . Behavioral episodes are  
240 designated by the index  $j$ . We use the letter  $b$  to refer to the behavioral allocation within a  
241 behavioral episode. Behavioral allocation is measured as relative time spent at a certain behavior  
242 within a given behavioral episode<sup>5</sup>. For reasons of simplicity, we only deal with one type of behavior  
243 here. Hence,  $b_{ij}$  refers to the behavioral allocation of individual  $i$  in episode  $j$ .

244 To capture the rather small time scale of learning, we formally treat surviving individuals as if they  
245 were their own offspring. Even though this treatment of survival may be counter-intuitive, it is not  
246 uncommon in biological models, since it provides a mathematically simple way to capture the  
247 survival part of evolutionary fitness (cf. Taylor, 1990). Moreover, since we focus only on the learned

---

<sup>4</sup> It would be possible to include several contexts into the model by introducing a class structured version of the Price equation (cf. Taylor, 1990). However, this would inflate the mathematical formalism and distract the reader from the general import of the model.

<sup>5</sup> Actually, the model allows for different measures of behavior, as well, e.g., relative response rate, running speed, spatial position, or even neural activity. However, for reasons of consistency, we restrict ourselves to cases where behavior is measured as the relative allocation over time.

248 components of behavior, we restrict the analysis to the survival part of fitness, thereby excluding the  
249 actual offspring of the individuals<sup>6</sup>. This means that, strictly speaking,  $b$  does not refer to behavioral  
250 allocation per se, but to the part of behavioral allocation that is ‘transmitted’ within individuals from  
251 one set of behavioral episodes to a second set of behavioral episodes . The mechanisms of  
252 transmission certainly involve some kind of neural processing and might be referred to as ‘memory’.  
253 We do not imply the notion of a cognitive process of ‘storing’ and ‘retrieving’ here, but use the term  
254 ‘memory’ for any mechanism that integrates past experiences with present behavior. However, since  
255 the model does not depend on these supposed mechanisms, we prefer to speak of ‘behavior’  $b$ .  
256 Formally, this definition of  $b$  is an analogue to the biological concept of ‘breeding value’ or ‘additive  
257 genetic value’, which is defined as the component of the evolving character that is genetically  
258 transmitted. Using breeding value instead of the actual character comes without loss of generality  
259 when using the Price equation to describe natural selection (Frank, 1998). Effectively, treating  
260 surviving individuals as their own offspring provides a way to bridge the time scales between natural  
261 selection and reinforcement learning because no matter how short the interval between the two sets  
262 of behavioral episodes, if fitness is defined by surviving individuals, the Price equation still describes  
263 (a part of) natural selection. At the same time, restricting the analysis to the survival component of  
264 fitness circumvents the problem of different mechanisms of transmission.

265 We can thus model several time steps on an evolutionary scale within the life span of a single  
266 individual. There are no formal constraints on the choice of the time scale, hence we can choose the  
267 level of analysis to fit our experimental requirements. Conceptually, this corresponds to a multiscale  
268 view of behavior analysis as advocated by Baum (2018). Note, however, that no matter how small we  
269 choose our evolutionary time scale to be, behavior is still understood as extended over time with  
270 several behavioral episodes occurring in each (evolutionary) time step.

---

<sup>6</sup> Including transmission by reproduction would require a separation of hereditary mechanisms on the different levels of the Price equation. Although this might help to understand how the levels of selection interact and influence one another, it is not necessary in order to derive the general structure of behavioral selection.

271 Within this conceptual framework, we can now describe average behavioral change simultaneously  
272 for individuals and for the whole population by means of a multilevel Price equation:

$$\bar{w}\Delta\bar{b} = \text{Cov}(w_i, b_i) + E_i\left(\text{Cov}(w_{ij}, b_{ij}) + E_j(w_{ij}\Delta b_{ij})\right) \quad (4.)$$

273 Using behavioral allocation  $b$  as an evolving character and focusing on a time scale small enough to  
274 capture behavioral adaptations on an individual level, the two terms on the right hand side can be  
275 interpreted as follows: the covariance term over individuals  $i$  designates the survival part of natural  
276 selection; the expectation term over individuals  $i$  designates behavioral changes within the  
277 individuals. The latter term is further partitioned into one selection part on the level of behavioral  
278 episodes  $j$  and one expectation part referring to changes within behavioral episodes. Individual  
279 fitness  $w_i$  is defined in the usual way as the contribution of individuals to the future population  
280 (because the only 'offspring' in this model are the individuals themselves, this is essentially survival).  
281 Since we are mainly interested in the selection part, we can simplify our notation by treating the  
282 within individual expectation term as a residual term  $\delta$ , defining  $E_j(w_{ij}\Delta b_{ij}) = \delta$ . Due to the  
283 recursive extension of the Price equation, change in behavioral allocation within each individual can  
284 now be expressed as:

$$w_i\Delta b_i = \text{Cov}(w_{ij}, b_{ij}) + \delta \quad (5.)$$

285 Although fitness can be explicitly defined on an individual level, it is difficult to make sense of the  
286 'fitness' ascribed to behavioral episodes  $w_{ij}$ . Since it is always the whole organism that dies or  
287 survives (or, more formally, has a certain survival probability), it is not reasonable to attribute  
288 evolutionary fitness to a behavioral episode (note that we are talking about fitness in a literal way  
289 here – thus, it would be inconsistent to invoke a metaphorical interpretation). Moreover, it is unlikely  
290 that individuals adapt their behavior with regard to their *actual* fitness, since this would require  
291 information about their future survival. Therefore, it is reasonable to assume that individuals adapt  
292 their behavior with regard to *fitness proxies*  $p$ . These fitness proxies are essentially statistical  
293 predictors of evolutionary fitness. In some cases, they may coincide with the aforementioned PIEs,

294 but this is not necessarily the case. Formally, we predict fitness  $w$  by a context-dependent linear  
295 regression of the form  $w = \beta_0 + \beta_{wp}p + \varepsilon$ . The regression is calculated over all individuals that are  
296 exposed to the contextual factors that constitute the class of behavioral episodes. This means that  
297 we calculate separate regressions for each context in accordance with the definition of context by its  
298 contingency structure.

299 We now use these regression coefficients to obtain statistical estimates of the expected fitness for  
300 each behavioral episode. Formally, we substitute the fitness of each behavioral episode  $w_{ij}$  with the  
301 corresponding predicted value from the regression model:

$$w_i \Delta b_i = \text{Cov}(\beta_0 + \beta_{wp} p_{ij}, b_{ij}) + \delta \quad (6.)$$

302 Since the regression coefficients are calculated on a between individuals level, they can be treated as  
303 constants within individuals. Therefore we can simplify to get<sup>7</sup>:

$$w_i \Delta b_i = \beta_{wp} \text{Cov}(p_{ij}, b_{ij}) + \delta \quad (7.)$$

304 This means that the change in behavior for an individual  $i$  equals the covariance between behavioral  
305 allocation  $b$  and a linear fitness predictor  $p$ , weighted by the statistical regression effect of the fitness  
306 proxy  $p$  on evolutionary fitness, with  $\delta$  being a residual term capturing all changes in behavior that  
307 are not caused by selection. We call this the *covariance based law of effect*, since it provides the  
308 conditions under which behavior is changed by reinforcement learning: *the change in mean*  
309 *behavioral allocation due to selection is proportional to the covariance between the behavior and a*  
310 *reinforcer*. Furthermore, the coefficient  $\beta_{wp}$  acts as a weighting factor to scale the covariance term.  
311 This means that behavioral change is also proportional to the degree to which a reinforcer is  
312 predictive for evolutionary fitness. Therefore,  $\beta_{wp}$  is called the *reinforcing power* of a  $p$  (cf.  
313 Borgstede, 2020).

---

<sup>7</sup> See Appendix 1 for the complete derivation.

## 314 5 Application of the model

315 The covariance based law of effect is a theoretical description of reinforcement learning on the most  
316 abstract level. It provides a quantitative account of behavioral selection that formally links the level  
317 of individual learning to the level of natural selection. We will now demonstrate how this abstract  
318 principle can be applied to various experimental paradigms. Using the method of path analysis, we  
319 will partition the reinforcing effects of behavioral selection into different components<sup>8</sup>. The first  
320 application deals with the implications of our model for steady state behavior. The second  
321 application is concerned with the factors that constitute a reinforcer.

### 322 5.1 Path analysis

323 The covariance based law of effect describes reinforcement by means of behavioral selection. The  
324 covariance term stresses the understanding of learning as a selection process as described by the  
325 Price equation. However, it is useful to translate the law of effect into a form that is more easily  
326 tractable in practical applications. Therefore, we introduce an alternative formulation of the same  
327 law by means of a statistical path model.

328 Because, by definition,  $\beta_{pb} = \frac{\text{Cov}(p_{ij}, b_{ij})}{\text{Var}(b_{ij})}$ , the covariance based law of effect can be equivalently  
329 stated as:

$$w_i \Delta b_i = \beta_{wp} \beta_{pb} \text{Var}(b_{ij}) + \delta \quad (8.)$$

330 Here,  $\beta_{pb}$  and  $\beta_{wp}$  are the partial linear effects from a path model where the effect of behavior  $b$  on  
331 evolutionary fitness  $w$  is fully mediated by a fitness proxy  $p$  (see Figure 1). The parameter  $\beta_{wp}$  refers  
332 to the reinforcing power of the fitness proxy. It corresponds to the expected change in evolutionary  
333 fitness per unit change in reinforcement. The parameter  $\beta_{pb}$  is the slope of the feedback function of

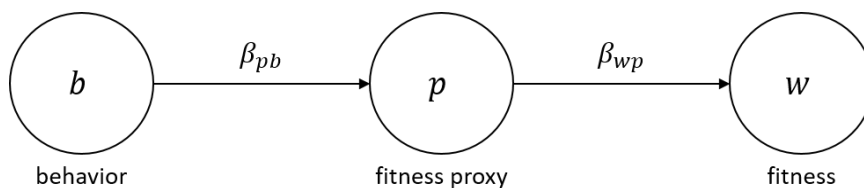
---

<sup>8</sup> Path analysis is routinely applied to partition different fitness effects on evolutionary fitness in biological models of selection (cf. Okasha and Otsuka, 2020; Scheiner and Gurevitch, 2001). The difference here is that, apart from the fitness effects of the fitness proxies, we use within individual regressions to calculate the contextual effects that are effective in the current class of behavioral episodes.

334 the schedule of reinforcement. It thus captures the expected gain in reinforcement per unit change  
335 in behavior at the current point of behavioral allocation.

336 Using the technique of statistical path analysis, we can calculate the total effect  $\beta_{wb}$  of behavior  $b$  on  
337 fitness  $w$  by multiplication of the two partial regression slopes, which gives  $\beta_{wb} = \beta_{wp}\beta_{pb}$ . This total  
338 effect corresponds to the marginal value of a change in behavioral allocation in terms of evolutionary  
339 fitness. According to Borgstede (2020), this may be regarded as marginal *reinforcer value*  $r(b)$  of  
340 behavior  $b$ . We may thus state the covariance based law of effect in terms of reinforcer value<sup>9</sup>:  
341 *behavioral selection equals the product of reinforcer value and behavioral variance.*

342 If we calculate the path model using standardized variables, behavioral variance will be one, leaving  
343 us with the standardized total effect of behavior on fitness. We can thus investigate different  
344 components of reinforcement by partitioning the reinforcer value of a behavior into partial  
345 regression effects. Similarly, we can start with a given partitioning and retrieve the corresponding  
346 expression of  $r(b)$  by summing up the products of the partial regression coefficients for each path  
347 (Shiple, 2016). This method can be applied to analyze the structure of different contexts from the  
348 perspective of behavioral selection.



349  
350 *Figure 1: Path diagram depicting the statistical relation between behavior  $b$ , fitness proxy  $p$  and evolutionary fitness  $w$ .  $\beta_{pb}$*   
351 *and  $\beta_{wp}$  designate the slopes of linear regression models that predict  $p$  from behavior  $b$ , and fitness  $w$  from  $p$  controlling*  
352 *for  $b$ , respectively. Statistically, the total effect of  $b$  on  $w$  equals the product of  $\beta_{pb}$  and  $\beta_{wp}$ .*

---

<sup>9</sup> Note that in Borgstede (2020),  $r(b)$  refers to the *partial* effect of behavior on fitness, whereas here, we refer to the *total* effect, which is the sum of the partial effects over all paths from  $b$  to  $w$ .



## 353 5.2 Steady state behavior

354 Our first application deals with several mutually exclusive behaviors that compete for time within  
355 behavioral episodes. Formally, behavioral allocation of behaviors  $b_1, b_2, \dots, b_n$  is subject to a fixed time  
356 budget constraint such that the relative time spent at each behavior sums up to one. From this it  
357 follows that any increase in mean behavioral allocation towards one behavior will result in an equal  
358 amount of decrease in the sum over all other behaviors.

359 In the following, we will focus on the dynamics of only one behavior  $b$ . Like before, the amount of  
360 behavioral selection is determined by the corresponding reinforcer value. In contrast to the simple  
361 mediation model of the previous section, we now have to incorporate the time constraint when  
362 calculating the total effect of behavior  $b$  on fitness  $w$ . Figure 2 depicts the simplest case with only  
363 one additional behavior (we call this second behavior  $b'$ )<sup>10</sup>. In the path model, we account for the  
364 time constraint by adding a path from behavior  $b$  to  $b'$  (dashed line in Figure 2). The corresponding  
365 partial regression coefficient  $\beta_{b'b}$  expresses the expected change in behavior  $b'$  per unit change in  $b$ .  
366 Due to the above budget constraint, it holds that  $\beta_{b'b} = -1$ . Consequently, we can calculate the  
367 marginal reinforcer value of behavior  $b$  as:

$$r(b) = \beta_{wp}\beta_{pb} - \beta_{wpr}\beta_{p'b'} \quad (9.)$$

368 This means that we have to subtract the partial fitness effect of behavior  $b'$  from the partial fitness  
369 effect of behavior  $b$  to predict behavioral selection on  $b$ .

370 We can use this result to derive a general equilibrium condition for behavioral selection. Behavioral  
371 equilibrium (or 'steady state behavior') is characterized by a constant behavioral allocation (i.e., the  
372 absence of change due to reinforcement). Formally, this means that behavior  $b$  is at a steady state, if

---

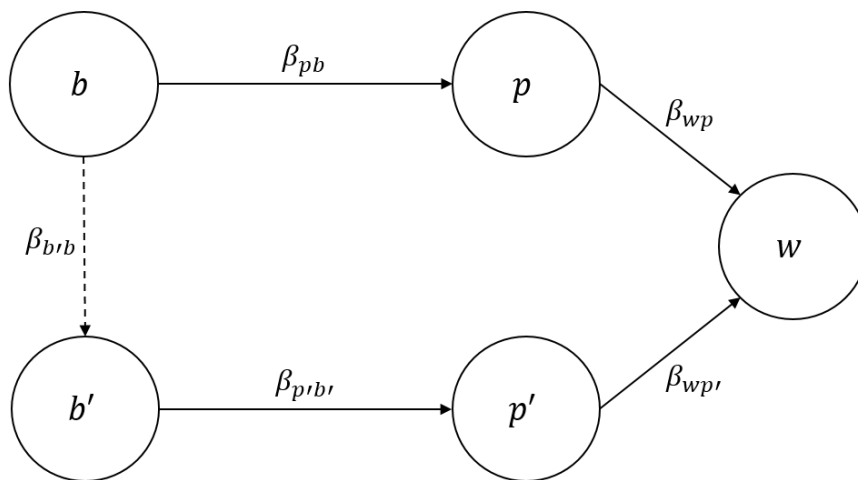
<sup>10</sup> We restrict ourselves to this minimal example to avoid unnecessary complexity of the path model. This does not restrict the generality of the analysis, however, because we only consider the point of behavioral equilibrium here. This means that, even if behaviors  $b$  and  $b'$  do not take up the whole time interval at the beginning of the experiment, in the absence of selection on other behaviors, they will eventually compete for the whole available time.

373 and only if behavioral selection on  $b$  equals zero (i.e.,  $w_i \Delta b_i = 0$ ). The equilibrium condition thus  
374 becomes:

$$\beta_{wp} \beta_{pb} = \beta_{wp'} \beta_{p'b} \quad (10.)$$

375 Hence, behavior is stable if and only if the marginal fitness effects of competing behaviors are equal.  
376 This result is equivalent to the one derived in Borgstede (2020). However, we did not assume a  
377 tendency to maximize reinforcer value, but derived the same condition for steady state behavior  
378 from the covariance based law of effect.

379 It has been shown that in concurrent variable interval schedules of reinforcement, the condition of  
380 equal marginal reinforcer value coincides with the *matching law* (Baum, 1981). However, this is not  
381 the case for all schedules of reinforcement, since the resulting behavioral allocation depends on the  
382 shape of the corresponding feedback functions. Hence, the matching law is best understood as a  
383 special case of the above equilibrium condition and therefore follows from the covariance based law  
384 of effect.



385

386 *Figure 2: Path diagram for behavioral selection on a behavior  $b$  that is constrained by a temporal budget. In addition to the*  
387 *direct path from behavior  $b$  to fitness  $w$ , the time constraint induces a negative correlation with behavior  $b'$ . At the point of*  
388 *behavioral equilibrium, this results in a second path from  $b$  to  $w$  that is mediated by  $b'$ .*

389 5.3 The nature of reinforcement

390 Let us now consider a case where behavioral allocation of behavior  $b$  has reached an equilibrium  
391 state (we may call this the 'baseline condition'). Effectively, this means that there is no behavioral  
392 change due to selection. Given behavioral variance does not equal zero, this is equivalent to a  
393 marginal reinforcer value of zero, yielding  $r(b) = 0$ . In other words, there is no benefit in terms of  
394 evolutionary fitness when behavioral allocation is changed.

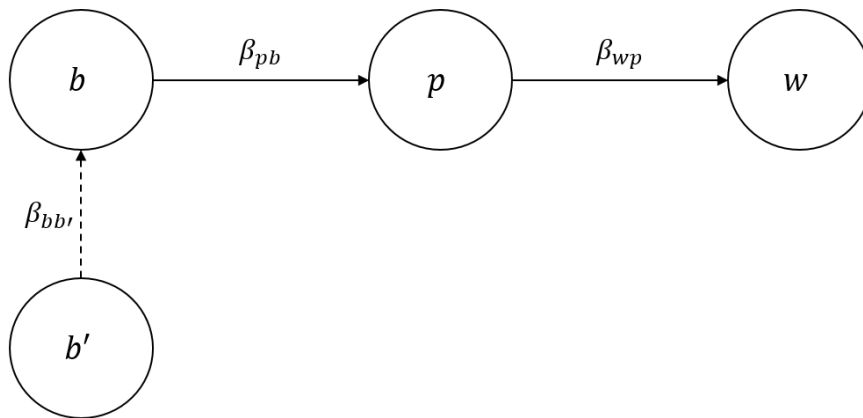
395 We now imagine a second behavior  $b'$  that we try to reinforce by pairing it with  $b$ , thereby  
396 establishing a positive covariance between  $b$  and  $b'$  (this is indicated by the dashed line in Figure 3).  
397 We can easily see that this will not affect  $b'$  because, under baseline conditions, the marginal  
398 reinforcer value of  $b'$  is:

$$r(b') = \beta_{bb'}r(b) = 0 \quad (11.)$$

399 We can now constraint the behavior  $b$ , thereby disturbing the equilibrium state we observed under  
400 baseline conditions. Following the principle of diminishing returns (i.e., assuming that the marginal  
401 effect of a change in behavioral allocation becomes smaller for higher values of  $b$ ), this will result in a  
402 positive reinforcer value for behavior  $b$ . Therefore, if we make  $b$  contingent on a second behavior  $b'$ ,  
403 the reinforcer value of  $b'$  becomes positive, as well, which in turn results in behavioral selection. This  
404 explains why *response deprivation* can establish constrained behaviors as reinforcers (Timberlake  
405 & Allison, 1974).

406 The covariance based law of effect provides us with a general definition of a reinforcer as a context-  
407 dependent fitness predictor. This implies that behaviors may become reinforcers, given they predict  
408 an expected gain in fitness, which will be the case for most constrained behaviors. Under the  
409 assumption that high-probability behaviors are often constrained by the environment (i.e.,  
410 individuals would engage even more in high-probability behaviors, if they could), this also explains  
411 Premack's principle, which states that usually, behaviors that occur at a high probability under

412 baseline conditions, function as reinforcers for behaviors that occur with a lower probability  
 413 (Premack & Premack, 1963).



414

415 *Figure 3: Path diagram for the reinforcing effects of a behavior  $b$  for a different behavior  $b'$ . Assuming diminishing returns of*  
 416 *fitness proxy  $p$  for increasing values of  $b$ , constraining behavior  $b$  below baseline will raise the value of  $\beta_{pb}$ . This, in turn, will*  
 417 *establish the constrained behavior as a reinforcer for  $b'$ .*

418 We can extend the method of path analysis to even more complex scenarios, such as the  
 419 establishment of conditioned reinforcers. In this paradigm, a formerly neutral stimulus (e.g., a  
 420 flashing light) is repeatedly paired with the availability of food (or any other reinforcer). The pairing  
 421 of a stimulus with a reinforcer corresponds to establishing an empirical covariance between the  
 422 stimulus  $s$  and a fitness proxy  $p$ . Given there is variation in  $s$ , we can thus calculate a linear  
 423 regression of  $p$  on  $s$ . Since the stimulus is made contingent on behavior  $b$ , this results in a mediation  
 424 model with  $s$  being a mediating variable between  $b$  and  $p$ . The total fitness effect of behavior  $b$  (i.e.,  
 425 its reinforcer value) thus becomes:

$$r_0(b) = \beta_{wp}\beta_{ps}\beta_{sb} \quad (12.)$$

426 When a second discriminative stimulus  $s'$  is added to the experiment, we establish a new path in the  
 427 predictive model (this is indicated by the dashed lines in Figure 4). Given  $s$  has already been  
 428 established as a discriminative stimulus, the statistical effect of behavior  $b$  on the new stimulus  $s'$  is  
 429 completely mediated by  $s$ . Therefore, the total fitness effect of the behavior can be partitioned as  
 430 follows:

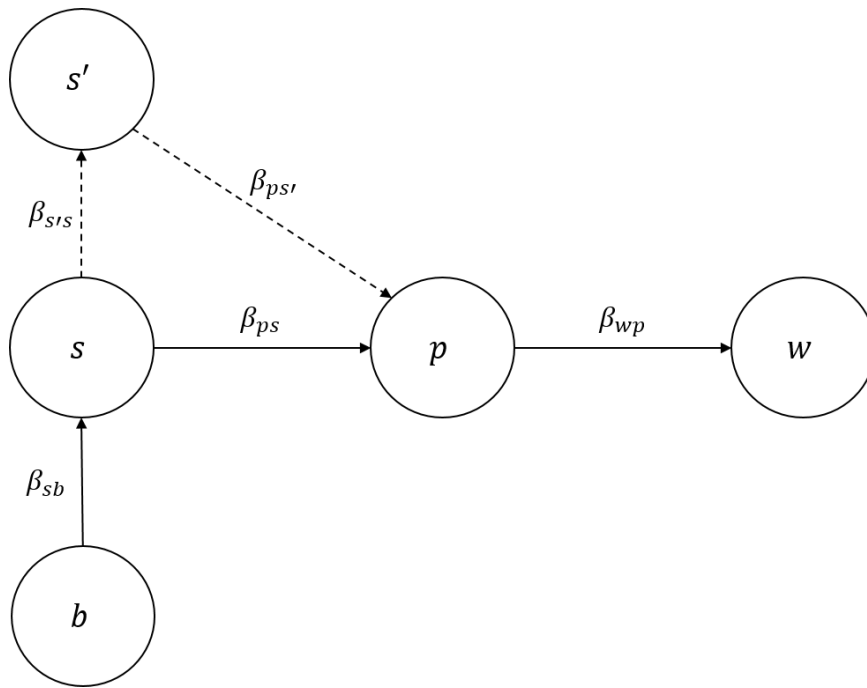
$$r_1(b) = \beta_{wp}\beta_{ps}\beta_{sb} + \beta_{wp}\beta_{ps'}\beta_{s's}\beta_{sb} \quad (13.)$$

431  $\beta_{ps'}$  corresponds to the partial effect of  $s'$  on  $p$  when controlled for the effect of  $s$ . As long as  $s'$   
432 always occurs together with  $s$ , there is no additional variance in  $p$  explained when  $s'$  is added to the  
433 regression. Consequently,  $\beta_{ps'}$  will be zero. If we now stop presenting the first stimulus, the expected  
434 fitness gain of behavior  $b$  reduces to the second path:

$$r_2(b) = \beta_{wp}\beta_{ps'}\beta_{s's}\beta_{sb} \quad (14.)$$

435 Since  $\beta_{ps'} = 0$ , the reinforcer value of  $b$  becomes zero. Therefore, we do not expect the second  
436 stimulus to affect behavior in the absence of the first one.

437 This phenomenon has been demonstrated repeatedly in classical and operant conditioning  
438 experiments and is known as the *blocking effect* (Kamin, 1969). Note that applying the covariance  
439 based law of effect to this special case does not only explain the blocking effect, but also makes a  
440 quantitative prediction about the amount of reinforcer value when  $s$  and  $s'$  are not perfectly  
441 correlated (i.e., if  $\beta_{ps'} \neq 0$ ). By partitioning the predictive effects of a behavior in a given context,  
442 we can thus explain how conditioned reinforcers acquire control over an individual's behavior. The  
443 general rule is that new stimuli affect the current reinforcer value (and thus the amount of  
444 reinforcement) to the degree that they provide additional information about the availability of  
445 fitness proxies (cf. Rescorla & Wagner, 1972). Thus, we can understand the well-known Rescorla-  
446 Wagner-Model as a molecular model that follows the general principle of reinforcement as described  
447 by the covariance based law of effect.



448

449 *Figure 4: Path diagram to illustrate the principle of conditioned reinforcement and the blocking effect. If a stimulus  $s$  is*  
450 *paired with a fitness proxy  $p$ , it will become predictive of fitness  $w$  and can then act as a reinforcer itself. However, if we pair*  
451 *an additional stimulus  $s'$  with the fitness proxy  $p$ , there will be no gain in predictive power, as long as  $s'$  is presented*  
452 *alongside with the previously established conditioned reinforcer  $s$  (i.e., the second stimulus is 'blocked' by the first one).*

## 453 6 Discussion

454 This paper deals with the question how reinforcement learning can be formally described as a  
455 Darwinian process. We present an evolutionary model of behavioral selection using the formalism of  
456 the Price equation, thereby clarifying some common conceptual ambiguities. The model treats  
457 behavioral selection (by means of reinforcement learning) as a part of natural selection, rather than a  
458 process that is merely analogous to natural selection. Hence, the 'Darwinian metaphor' advocated by  
459 many behavioral scientists is replaced by the view that behavioral selection literally *is* a Darwinian  
460 process (i.e., a part of evolution by natural selection). Therefore, in contrast to existing approaches to  
461 behavioral selection (e.g., Baum, 2017; Donahoe et al., 1993; McDowell, 2004), the presented  
462 multilevel model of behavioral selection is more than a re-statement of standard learning theories  
463 using a selectionist vocabulary – it is a true integration of behavioral psychology and evolutionary  
464 biology.

465 The model is used to derive the *covariance based law of effect*, which describes reinforcement  
466 learning as a selection process that is proportional to the covariance between behavioral allocation  
467 and a fitness proxy, weighted by the corresponding statistical fitness effect. We further showed how  
468 this abstract principle of behavioral selection can be applied to various experimental paradigms,  
469 thereby integrating such diverse empirical regularities as conditioned reinforcement, response  
470 deprivation, the blocking effect and the matching law. We thus lay the formal foundation for a  
471 general theory of reinforcement that is grounded in the theory of evolution by natural selection.

472 The covariance based law of effect also provides a formal definition of a ‘reinforcer’: a reinforcer is  
473 anything that is statistically predictive of evolutionary fitness. Due to the statistical nature of  
474 reinforcers, it does not matter whether we conceive them as external events (e.g., the availability of  
475 food), activities of the individual (e.g., eating), or perceptions (e.g., taste). Hence, we might just as  
476 well attribute the reinforcing power of a  $p$  to the feeding behavior that it induces, or to the resulting  
477 taste perception. A reinforcer is not a ‘thing’ that somehow changes the individual, it is essentially a  
478 (context-dependent) statistical fitness predictor.

479 Since the multilevel model of behavioral selection translates behavioral change into statistical terms  
480 like covariances, variances and regression coefficients, behavioral selection is inherently linked to the  
481 concept of *prediction*. In fact, it has been shown that natural selection as described by the Price  
482 equation maximizes Fisher Information<sup>11</sup> (Frank, 2009). This means that evolutionary change can be  
483 understood as a mechanism that maximizes predictive power with regard to the environment.

484 Adopting this view, we interpret behavioral selection as a mechanism to optimize the predictions of  
485 an organism about the environment. This is in line with the theory of *predictive coding* (Helmholtz,  
486 1909) and the *Bayesian brain* hypothesis (Clark, 2013). Although the formal integration of the  
487 ‘predictive brain’ hypothesis and the principle of behavioral selection has yet to be accomplished, the  
488 link between the Price equation and information theory provides a promising approach towards a

---

<sup>11</sup> Fisher information is a statistical measure of how much information observations provide about an unknown parameter of a probability distribution.

489 general theory of learning that explains the structure of learning processes by means of a universal  
490 selection principle.

491 Of course, the high level of generality of the presented model is limited by the underlying  
492 assumptions. First, we did not model the hereditary mechanisms that mediate behavioral selection.  
493 Treating surviving individuals as their own offspring and leaving aside the contribution to the  
494 population by reproduction, it was formally possible to leave the question of hereditary mechanisms  
495 open. However, this means that the 'offspring population' only consists of surviving individuals,  
496 thereby ignoring all behavioral changes that stem from (genetical) transmission and reproduction.  
497 This means that newborn individuals are not counted to calculate evolutionary fitness. Therefore,  
498 strictly speaking, the model only refers to the survival component of evolutionary fitness and thus  
499 does not capture the whole effect of natural selection. However, since reinforcement occurs only  
500 within individuals, omitting the reproduction part of evolutionary fitness does not pose a major  
501 problem. Nevertheless, it would be interesting to investigate the effects of different hereditary  
502 mechanisms on different levels of selection. Apart from the genetic part, in social species, cultural  
503 mechanisms (like imitation, model learning, or verbal instruction) partly mediate the transmission of  
504 behavior from parent (in the usual meaning) to offspring. Moreover, individual learning invokes at  
505 least one more mechanism of transmission (some kind of 'memory'). Disentangling these effects  
506 requires a careful mathematical treatment and will be targeted in subsequent work.

507 Finally, it should be mentioned that the Price equation in its original form implies a homogeneous  
508 population, thereby assuming inter-individual variation negligible. This assumption has been made  
509 here to keep the mathematical notation as simple as possible. Inter-individual variation can be added  
510 using a class structured version of the Price equation. This introduces additional weighting factors for  
511 the classes, where each class refers to a 'type' of individual, defined by a certain combination of  
512 characteristics (Taylor, 1990). In population biology, these weighting factors are the *reproductive*  
513 *values* of the different types of individuals and depend on the demography and long-term dynamics  
514 of the population (cf. Caswell, 2001). If the aim is to make specific predictions about adaptive



515 behavior, at least some kind of population structure has to be modelled. However, as long as one is  
516 mainly concerned with the fundamental principles of selection, class structure can be ignored.  
517 Therefore, here we stick to the most general form of the Price equation, resulting in a most general  
518 account of behavioral selection. Incorporating inter-individual variation will be the objective of  
519 subsequent work.

520 Despite these limitations, this paper provides a consistent quantitative account of reinforcement  
521 learning on a molar level. It integrates behavioral selection with natural selection and provides new  
522 insights into the quantitative relation between reinforcement, behavioral allocation and evolutionary  
523 fitness. This is an important step towards a general account of learning and of behavior in general,  
524 based on the theory of evolution by natural selection.

## 525 7 Author contributions

526 XX conceived the original idea, developed the theoretical formalism, performed the mathematical  
527 derivations and wrote the original draft. XY verified the analytical methods, provided additional  
528 ideas and critical feedback and helped shape the research, analysis and interpretation. Both authors  
529 contributed to the final version of the manuscript.

## 530 8 Competing interests

531 The authors declare no competing interests.

## 532 9 Appendix

### 533 9.1 Derivation of the covariance based law of effect

534 We start with the elementary Price equation for mean individual behavioral allocation  $b_i$ :

$$535 \quad \bar{w}\Delta\bar{b} = \text{Cov}(w_i, b_i) + E_i(w_i\Delta b_i)$$

536 Applying the logic of the multilevel Price equation, the expectation term can be separated into a  
537 within individual covariance term and a within individual expectation term:

538 
$$w_i \Delta b_i = \text{Cov}(w_{ij}, b_{ij}) + E_j(w_{ij} \Delta b_{ij})$$

539 By substitution, we arrive at the multilevel Price equation for behavioral selection:

540 
$$\bar{w} \Delta \bar{b} = \text{Cov}(w_i, b_i) + E_i \left( \text{Cov}(w_{ij}, b_{ij}) + E_j(w_{ij} \Delta b_{ij}) \right)$$

541 We now assume a statistical fitness predictor  $p$  of the form  $w = \beta_0 + \beta_{wp} p + \varepsilon$  and substitute the

542  $w_{ij}$  with the predicted values from this regression:

543 
$$\bar{w} \Delta \bar{b} = \text{Cov}(w_i, b_i) + E_i \left( \text{Cov}(\beta_0 + \beta_{wp} p_{ij}, b_{ij}) + E_j((\beta_0 + \beta_{wp} p_{ij}) \Delta b_{ij}) \right)$$

544 This can be rearranged to:

545 
$$\bar{w} \Delta \bar{b} = \text{Cov}(w_i, b_i) + E_i \left( \text{Cov}(\beta_0, b_{ij}) + \text{Cov}(\beta_{wp} p_{ij}, b_{ij}) + E_j((\beta_0 + \beta_{wp} p_{ij}) \Delta b_{ij}) \right)$$

546 Since  $\beta_0$  is a constant,  $\text{Cov}(\beta_0, b_{ij})$  equals 0. This results in:

547 
$$\bar{w} \Delta \bar{b} = \text{Cov}(w_i, b_i) + E_i \left( \text{Cov}(\beta_{wp} p_{ij}, b_{ij}) + E_j((\beta_0 + \beta_{wp} p_{ij}) \Delta b_{ij}) \right)$$

548 Rearrangement yields:

549 
$$\bar{w} \Delta \bar{b} = \text{Cov}(w_i, b_i) + E_i \left( \beta_{wp} \text{Cov}(p_{ij}, b_{ij}) + E_j((\beta_0 + \beta_{wp} p_{ij}) \Delta b_{ij}) \right)$$

550 Because the expectation term in the Price equation is taken over the fitness weighted changes in

551 behavioral allocation  $w_i \Delta b_i$ , which has been separated into  $\text{Cov}(w_{ij}, b_{ij}) + E_j(w_{ij} \Delta b_{ij})$  above, all

552 rearrangements of the multilevel Price equation within this expectation term equally apply to  $w_i \Delta b_i$ .

553 Therefore, we can write the change in behavioral allocation within individuals as:

554 
$$w_i \Delta b_i = \beta_{wp} \text{Cov}(p_{ij}, b_{ij}) + E_j((\beta_0 + \beta_{wp} p_{ij}) \Delta b_{ij})$$

555 Defining  $\delta = E_j((\beta_0 + \beta_{wp} p_{ij}) \Delta b_{ij})$  we arrive at the covariance based law of effect:

556 
$$w_i \Delta b_i = \beta_{wp} \text{Cov}(p_{ij}, b_{ij}) + \delta$$

557 10 References

- 558 Baum, W. M. (1973). The correlation-based law of effect. *Journal of the Experimental Analysis of*  
559 *Behavior*, 20, 137–153.
- 560 Baum, W. M. (1981). Optimization and the matching law as accounts of instrumental behavior.  
561 *Journal of the Experimental Analysis of Behavior*, 36, 387–403.
- 562 Baum, W. M. (2005). *Understanding behaviorism: Behavior, culture, and evolution* (2nd edition).  
563 Malden, MA: Blackwell Publishing.
- 564 Baum, W. M. (2012). Rethinking reinforcement: Allocation, induction, and contingency. *Journal of the*  
565 *Experimental Analysis of Behavior*, 97(1), 101–124. <https://doi.org/10.1901/jeab.2012.97-101>
- 566 Baum, W. M. (2017). Selection by consequences, behavioral evolution, and the price equation.  
567 *Journal of the Experimental Analysis of Behavior*, 107(3), 321–342.  
568 <https://doi.org/10.1002/jeab.256>
- 569 Baum, W. M. (2018). Multiscale behavior analysis and molar behaviorism: An overview. *Journal of the*  
570 *Experimental Analysis of Behavior*, 110(3), 302–322. <https://doi.org/10.1002/jeab.476>
- 571 Becker, A. M. (2019). The flight of the locus of selection: Some intricate relationships between  
572 evolutionary elements. *Behavioural Processes*, 161, 31–44.  
573 <https://doi.org/10.1016/j.beproc.2018.01.002>
- 574 Borgstede, M. (2020). An evolutionary model of reinforcer value. *Behavioural Processes*, 104109.
- 575 Borgstede, M., & Simon, C. (under review). A conceptual synthesis of learning and evolution.
- 576 Broadbent, D. E. (1961). *Behaviour*. London: Methuen.
- 577 Burgos, J. E. (2019). Selection by reinforcement: A critical reappraisal. *Behavioural Processes*, 161,  
578 149–160. <https://doi.org/10.1016/j.beproc.2018.01.019>
- 579 Campbell, D. T. (1956). Adaptive behavior from random response. *Behavioral Science*, 1(2), 105–110.  
580 <https://doi.org/10.1002/bs.3830010204>
- 581 Caswell, H. (2001). *Matrix Population Models. Construction, analysis, and interpretation. 2nd ed.*  
582 Sunderland: Sinauer Associates.
- 583 Clark, A. (2013). Whatever next? Predictive brains, situated agents, and the future of cognitive  
584 science. *The Behavioral and Brain Sciences*, 36(3), 181–204.  
585 <https://doi.org/10.1017/S0140525X12000477>
- 586 Donahoe, J. W. (2011). Selectionism. In K. A. Lattal & P. A. Chase (Eds.), *Behavior theory and*  
587 *philosophy* (Vol. 33, pp. 103–128). New York, London: Springer. [https://doi.org/10.1007/978-1-](https://doi.org/10.1007/978-1-4757-4590-0_6)  
588 [4757-4590-0\\_6](https://doi.org/10.1007/978-1-4757-4590-0_6)
- 589 Donahoe, J. W., Burgos, J. E., & Palmer, D. C. (1993). A selectionist approach to reinforcement.  
590 *Journal of the Experimental Analysis of Behavior*, 60(1), 17–40.  
591 <https://doi.org/10.1901/jeab.1993.60-17>
- 592 Frank, S. A. (1997). The Price equation, Fisher's fundamental theorem, kin selection, and causal  
593 analysis. *Evolution; International Journal of Organic Evolution*, 51(6), 1712–1729.  
594 <https://doi.org/10.1111/j.1558-5646.1997.tb05096.x>
- 595 Frank, S. A. (1998). *Foundations of social evolution. Monographs in behavior and ecology*. Princeton,  
596 NJ: Princeton Univ. Press.
- 597 Frank, S. A. (2009). Natural selection maximizes Fisher information. *Journal of Evolutionary Biology*,  
598 22(2), 231–244. <https://doi.org/10.1111/j.1420-9101.2008.01647.x>

- 599 Frank, S. A. (2017). Universal expressions of population change by the Price equation: Natural  
600 selection, information, and maximum entropy production. *Ecology and Evolution*, 7(10), 3381–  
601 3396. <https://doi.org/10.1002/ece3.2922>
- 602 Frank, S. A. (2018). The Price Equation Program: Simple Invariances Unify Population Dynamics,  
603 Thermodynamics, Probability, Information and Inference. *Entropy*, 20(12), 978.  
604 <https://doi.org/10.3390/e20120978>
- 605 Frank, S. A. (2020). Simple unity among the fundamental equations of science. *Philosophical*  
606 *Transactions of the Royal Society of London. Series B, Biological Sciences*, 375(1797), 20190351.  
607 <https://doi.org/10.1098/rstb.2019.0351>
- 608 Gardner, A. (2020). Price's equation made clear. *Philosophical Transactions of the Royal Society of*  
609 *London. Series B, Biological Sciences*, 375(1797), 20190361.  
610 <https://doi.org/10.1098/rstb.2019.0361>
- 611 Gilbert, R. M. (1970). Psychology and biology. *Canadian Psychologist/Psychologie Canadienne*, 11(3),  
612 221–238. <https://doi.org/10.1037/h0082574>
- 613 Grafen, A. (2000). Developments of the Price equation and natural selection under uncertainty.  
614 *Proceedings. Biological Sciences*, 267, 1223–1227.
- 615 Grafen, A. (2014). The formal darwinism project in outline. *Biology & Philosophy*, 29(2), 155–174.  
616 <https://doi.org/10.1007/s10539-013-9414-y>
- 617 Helmholtz, H. von (1909). *Treatise on physiological optics. Vol. III 3rd edition*. Hamburg: Voss.
- 618 Herrnstein, R. J. (1961). Relative and absolute strength of response as a function of frequency of  
619 reinforcement. *Journal of the Experimental Analysis of Behavior*, 4, 267–272.  
620 <https://doi.org/10.1901/jeab.1961.4-267>
- 621 Herrnstein, R. J. (1964). Will. *Proceedings of the American Philosophical Society*, 108(6), 455–458.
- 622 Hull, D. L., Langman, R. E., & Glenn, S. S. (2001). A general account of selection: Biology, immunology,  
623 and behavior. *The Behavioral and Brain Sciences*, 24(3), 511–528.  
624 <https://doi.org/10.1017/S0140525X01004162>
- 625 Kamin, L. J. (1969). Predictability, surprise, attention and conditioning. In B. A. Campbell & R. M.  
626 Church (Eds.), *Punishment and aversive behavior* (pp. 279–296). New York.
- 627 Luque, V. J. (2017). One equation to rule them all: a philosophical analysis of the Price equation.  
628 *Biology & Philosophy*, 32(1), 97–125. <https://doi.org/10.1007/s10539-016-9538-y>
- 629 McDowell, J. J. (2004). A computational model of selection by consequences. *Journal of the*  
630 *Experimental Analysis of Behavior*, 81(3), 297–317. <https://doi.org/10.1901/jeab.2004.81-297>
- 631 McDowell, J. J. (2013). A quantitative evolutionary theory of adaptive behavior dynamics.  
632 *Psychological Review*, 120(4), 731–750. <https://doi.org/10.1037/a0034244>
- 633 Nettle, D. (2020). Selection, adaptation, inheritance and design in human culture: The view from the  
634 Price equation. *Philosophical Transactions of the Royal Society of London. Series B, Biological*  
635 *Sciences*, 375(1797), 20190358. <https://doi.org/10.1098/rstb.2019.0358>
- 636 Okasha, S., & Otsuka, J. (2020). The Price equation and the causal analysis of evolutionary change.  
637 *Philosophical Transactions of the Royal Society of London. Series B, Biological Sciences*, 375(1797),  
638 20190365. <https://doi.org/10.1098/rstb.2019.0365>
- 639 Pennypacker, H. S. (1992). Is behavior analysis undergoing selection by consequences? *The American*  
640 *Psychologist*, 47(11), 1491–1498.
- 641 Premack, D., & Premack, A. J. (1963). Increased eating in rats deprived of running. *Journal of the*  
642 *Experimental Analysis of Behavior*, 6, 209–212. <https://doi.org/10.1901/jeab.1963.6-209>

- 643 Price, G. R. (1970). Selection and Covariance. *Nature*, 227(5257), 520–521.  
644 <https://doi.org/10.1038/227520a0>
- 645 Price, G. R. (1972). Extension of covariance selection mathematics. *Annals of Human Genetics*, 35(4),  
646 485–490. <https://doi.org/10.1111/j.1469-1809.1957.tb01874.x>
- 647 Price, G. R. (1995, written ca. 1971). The nature of selection. (Written circa 1971, published  
648 posthumously). *Journal of Theoretical Biology*, 175(3), 389–396.  
649 <https://doi.org/10.1006/jtbi.1995.0149>
- 650 Pringle, J. (1951). On the Parallel Between Learning and Evolution. *Behaviour*, 3(1), 174–214.  
651 <https://doi.org/10.1163/156853951X00269>
- 652 Rachlin, H. (1978). A molar theory of reinforcement schedules. *Journal of the Experimental Analysis*  
653 *of Behavior*, 30(3), 345–360. <https://doi.org/10.1901/jeab.1978.30-345>
- 654 Rescorla, R. A., & Wagner, A. R. (1972). A theory of Pavlovian conditioning: Variations in the  
655 effectiveness of reinforcement and nonreinforcement. In A. H. Black & W. F. Prokasy (Eds.),  
656 *Classical conditioning II: Current research and theory* (pp. 64–99). New York: Appleton-Century-  
657 Crofts.
- 658 Richerson, P. J. (2019). An integrated bayesian theory of phenotypic flexibility. *Behavioural Processes*,  
659 161, 54–64. <https://doi.org/10.1016/j.beproc.2018.02.002>
- 660 Scheiner, S. M., & Gurevitch, J. (2001). *Design and analysis of ecological experiments* (2nd ed.).  
661 Oxford, New York: Oxford University Press.
- 662 Shipley, B. (2016). *Cause and correlation in biology: A user's guide to path analysis, structural*  
663 *equations and causal inference with R* (2nd edition). Cambridge: Cambridge University Press.  
664 <https://doi.org/10.1017/CBO9781139979573>
- 665 Simon, C., & Hesse, D. O. (2019). Selection as a domain-general evolutionary process. *Behavioural*  
666 *Processes*, 161, 3–16. <https://doi.org/10.1016/j.beproc.2017.12.020>
- 667 Skinner, B. F. (1950). Are theories of learning necessary? *Psychological Review*, 57(4), 193–216.
- 668 Skinner, B. F. (1966). The phylogeny and ontogeny of behavior. Contingencies of reinforcement  
669 throw light on contingencies of survival in the evolution of behavior. *Science (New York, N.Y.)*,  
670 153(3741), 1205–1213. <https://doi.org/10.1126/science.153.3741.1205>
- 671 Skinner, B. F. (1969). *Contingencies of reinforcement: A theoretical analysis*. The century psychology  
672 series. Englewood Cliffs, NJ: Prentice-Hall.
- 673 Skinner, B. F. (1981). Selection by consequences. *Science (New York, N.Y.)*, 213(4507), 501–504.
- 674 Skinner, B. F. (1984). Selection by consequences. *The Behavioral and Brain Sciences*, 7(04), 477.  
675 <https://doi.org/10.1017/S0140525X0002673X>
- 676 Sober, E. (1984). *The nature of selection: Evolutionary theory in philosophical focus*. Cambridge,  
677 Mass.: Bradford Books/MIT Press. Retrieved from  
678 [http://search.ebscohost.com/login.aspx?direct=true&scope=site&db=nlebk&db=nlabk&AN=9097](http://search.ebscohost.com/login.aspx?direct=true&scope=site&db=nlebk&db=nlabk&AN=909783)  
679 83
- 680 Staddon, J. E. R. (2016). *Adaptive behavior and learning* (Second edition). Cambridge: Cambridge  
681 University Press. Retrieved from <http://dx.doi.org/10.1017/CBO9781139998369>  
682 <https://doi.org/10.1017/CBO9781139998369>
- 683 Staddon, J. E. R., & Simmelhag, V. L. (1971). The "superstition" experiment: A reexamination of its  
684 implications for the principles of adaptive behavior. *Psychological Review*, 78(1), 3–43.  
685 <https://doi.org/10.1037/h0030305>
- 686 Taylor, P. D. (1990). Allele-Frequency Change in a Class-Structured Population. *The American*  
687 *Naturalist*, 135(1), 95–106.

- 688 Thorndike, E. L. (1900). The associative processes in animals. *Biological Lectures from the Marine*  
689 *Biological Laboratory of Woods Holl, 1899*, 69–91.
- 690 Thorndike, E. L. (2010/1911). *Animal intelligence; experimental studies*. New Brunswick, NJ:  
691 Transaction Publishers. <https://doi.org/10.5962/bhl.title.55072>
- 692 Timberlake, W., & Allison, J. (1974). Response deprivation: An empirical approach to instrumental  
693 performance. *Psychological Review, 81*(2), 146–164. <https://doi.org/10.1037/h0036101>
- 694 Tonneau, F., & Sokolowski, M. B. C. (2000). Pitfalls of Behavioral Selectionism. In F. Tonneau & N. S.  
695 Thompson (Eds.), *Perspectives in Ethology: Vol. 13. Perspectives in Ethology: Evolution, Culture,*  
696 *and Behavior* (Vol. 13, pp. 155–180). Boston, MA: Springer. [https://doi.org/10.1007/978-1-4615-](https://doi.org/10.1007/978-1-4615-1221-9_6)  
697 [1221-9\\_6](https://doi.org/10.1007/978-1-4615-1221-9_6)