

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24

Estimating (non)linear selection on reaction norms:

A general framework for labile traits

Jordan S. Martin^{*1,2}, Yimen Araya-Ajoy³, Niels J. Dingemanse⁴,
Alastair J. Wilson⁵, & David Westneat⁶

*Corresponding author: jordanscott.martin@eawag.ch

*¹Evolutionary Ecology of Aquatic Ecosystems Laboratory, Fish Ecology and Evolution,
Eawag Swiss Federal Institute of Aquatic Science & Technology, Switzerland*

*²Human Ecology Group, Institute of Evolutionary Medicine,
University of Zurich, Switzerland*

*³Center for Biodiversity Dynamics, Department of Biology,
Norwegian University of Science and Technology, Norway*

*⁴Behavioral Ecology Unit, Department of Biology,
Ludwig Maximilian University of Munich, Germany*

*⁵Evolution Group, Centre for Biosciences,
University of Exeter, United Kingdom*

*⁶Department of Biology,
University of Kentucky, United States of America*

Abstract

25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49

1. Individual reaction norms describe how labile phenotypes vary as a function of organisms' expected trait values (intercepts) and plasticity across environments (slopes), as well as their degree of stochastic phenotypic variability or predictability (residuals). These reaction norms can be estimated empirically using multilevel, mixed-effects models and play a key role in ecological research on a variety of behavioral, physiological, and morphological traits. Many evolutionary models have also emphasized the importance of understanding reaction norms as a target of selection in heterogeneous and dynamic environments.
2. However, it remains difficult to empirically estimate nonlinear selection on reaction norms, inhibiting robust tests of adaptive theory and accurate predictions of phenotypic evolution. To address this challenge, we propose generalized multilevel models for estimating stabilizing, disruptive, and correlational selection on the reaction norms of labile traits, which can be applied to any repeatedly measured phenotype using a flexible Bayesian framework.
3. Our modelling approach avoids inferential bias by simultaneously accounting for uncertainty in reaction norm parameters and their potentially nonlinear fitness effects. We formally introduce these nonlinear selection models and provide detailed discussion on their interpretation and potential extensions. We then validate their application in a Bayesian framework using simulation-based calibration and power analyses.
4. We find that our models facilitate unbiased Bayesian inference across a broad range of effect sizes and desirable power for hypothesis tests with large sample sizes. Coding tutorials are further provided to aid empiricists in applying these models to any phenotype of interest using the Stan statistical programming language in R. The proposed modeling framework should, therefore, readily enhance tests of adaptive theory for a variety of labile traits in the wild.

50 **Keywords:**

51 phenotypic evolution, complex trait, multivariate, adaptation, personality, flexibility

Introduction

52

53 A population will evolve by natural selection whenever heritable variation occurs in fitness-
54 relevant phenotypes (Darwin 1859). Measuring the fitness consequences of individual differences
55 in highly labile behavioral, physiological, and morphological traits is, therefore, fundamental for
56 explaining their adaptive evolution. Across a variety of phenotypes and taxa, repeatable individual
57 differences have been observed in organisms' average trait values (Bell, Hankison, & Laskowski
58 2009; Fanson & Biro 2015; Cauchoux et al. 2018) and in their plasticity across environments
59 (Dingemanse et al. 2010; Stamps 2016; Arnold, Nicotra, & Kruuk 2019), with some individuals
60 consistently being more or less responsive to environmental change than others. In addition, it is
61 increasingly appreciated that individuals may repeatably differ in their degree of stochastic
62 phenotypic variability within a given environment (see **Box 1** below for a conceptual overview;
63 Biro & Adriaenssens 2013; Westneat, Schofield, & Wright 2013; Mitchell, Beckmann, & Biro
64 2021), a phenomenon which has often been ignored in ecological research (Hansen, Carter &
65 Pélabon 2006). These individual-specific patterns reflect distinct but potentially integrated
66 parameters (intercepts, slopes, and within-individual residuals) of the reaction norms (RNs, i.e.
67 state-dependent functions relating phenotype to environment, **Table 1**) evolving in a population
68 (**Figure 1**). RN models provide a highly generalizable, quantitative framework for investigating the
69 evolution and development of labile traits, with broad applications ranging from social behaviors
70 (Dingemanse & Araya-Ajoy 2015; McNamara & Leimar 2020; Martin, Jaeggi, & Koski, 2023) and
71 learning processes (Wright, Haaland, Dingemanse, & Westneat 2022) to thermal performance
72 curves (Svensson, Gomez-Llano, & Waller 2020) and extended phenotypes (Munar-Delgado,
73 Araya-Ajoy, & Edelaar, 2023), such as gall size in insect-host plant interactions (Weis & Gorman
74 1990). Interest in the evolutionary ecology of RNs has grown steadily across a diverse range of
75 fields in recent decades (e.g. Brommer, Kontiainen, & Pietiäinen 2012; Strickland et al. 2021;
76 Newediuk, Prokopenko, & Wal 2022), generating methodological innovations for estimating RNs
77 subject to measurement error (e.g. Nussey, Wilson, & Brommer 2007; Dingemanse &
78 Dochtermann 2013; Gomulkiewicz et al. 2018; O'Dea, Noble, & Nakagawa 2021; Martin & Jaeggi
79 2022), as well as theoretical models for explaining the selection pressures shaping and
80 maintaining individual variation in RNs within populations (e.g. Wolf & Weissing 2010; Dall &
81 Griffith 2014; Sih et al. 2015; Wright et al. 2019). Attention to RNs has also increased in related
82 fields of inquiry such as personality psychology (Denissen & Penke 2008; Nettle & Penke 2010)
83 and evolutionary anthropology (Jaeggi et al. 2016).

84 RN models are not only useful statistical tools for describing phenotypic variation. Classic
85 theoretical models often assumed that selection acted independently on phenotypes expressed
86 in discrete states of the environment (so-called *character states*), where the evolution of RN
87 parameters and thus phenotypic plasticity across environments was interpreted as a byproduct
88 of state-specific selection (Via & Lande 1985; Gomulkiewicz & Kirkpatrick 1992). Many biologists
89 disagreed with this perspective on empirical and theoretical grounds, resulting in historical
90 debates about whether RN parameters should be conceptualized as direct or indirect targets of
91 natural selection (Gavrilets & Sheiner 1993; Scheiner 1993a; Via et al. 1995; Nicoglou 2015).
92 Fortunately, this disagreement is now largely resolved (Futuyma 2021), with evolutionary
93 quantitative genetic theory demonstrating the mathematical equivalence and thus conceptually
94 complementarity of models emphasizing selection on expressed character states or the RNs
95 producing them (de Jong 1995). As such, many contemporary evolutionary frameworks
96 emphasize RNs parameters (intercepts, slopes, and residuals) and their underlying mechanisms
97 as putative targets of selection, leading to differential patterns of adaptation and extinction in
98 changing environments (Schlichting & Piglucci 1998; Ghalambor, McKay, Carroll, & Reznick
99 2007; Fox et al. 2019). For instance, evolutionary ecologists have long investigated the unique
100 role of both cue-induced and stochastic phenotypic plasticity in the colonization of novel habitats
101 (Caño et al. 2008; Volis, Ormanbekova, & Yermekbayev 2015; Hendry 2016; Wang & Althoff
102 2019). In addition, evolutionary geneticists have shown how plasticity in social environments can
103 magnify heritable variation in mean trait values, accelerating or inhibiting phenotypic evolution in
104 comparison to unresponsive phenotypes (Moore et al. 1997; Bijma & Wade 2008; McGlothlin et
105 al. 2010; Kazancıoğlu, Klug, & Alonzo 2012). Game theorists and behavioral ecologists have
106 further emphasized the importance of understanding selection on RNs due to the prevalence of
107 fluctuating density- and frequency-dependent selection in social environments (Araya-Ajoy,
108 Westneat, & Wright 2020; McNamara & Leimar 2020; Martin, Jaeggi, & Koski 2023), as well as
109 the role of dynamic environments more generally in selecting for learning mechanisms and
110 emotional states rather than specific behaviors per se (Skinner, 1966; Henrich & McElreath 2003;
111 McNamara & Houston 2009; Fawcett, Hamblin, & Giraldeau 2013; Nakahashi & Ohtsuki 2015;
112 Wright et al. 2022). Distinct genetic control of phenotypic stability and change has also been
113 experimentally demonstrated for diverse phenomena from cold tolerance (Ørsted, Rohde,
114 Hoffmann, Sørensen, & Kristensen 2018) to body size (Scheiner & Lyman, 1991) and various
115 forms of developmental polyphenism (Suzuki & Nijhout 2006; Projecto-Garcia, Biddle, Ragsdale
116 2017), suggesting that differential selection on heritable variation in RN intercepts, slopes, and
117 residuals, as well as differential patterns of genetic integration between RN parameters (Wagner,

118 Booth, & Bagheri-Chaichian, 1997; Tonsor, Elnaccash, & Scheiner, 2013), can in turn have
119 distinct consequences for the evolutionary response to selection (de Jong 1995; Martin et al.
120 2024). Accordingly, divergence has been observed in the RNs of many naturally occurring
121 populations, such as differential plasticity in the growth rates of phytoplankton (*Thalassiosira*
122 *pseudonana*; Schaum, Buckling, Smirnoff, & Yvon-Durocher 2022), ponderosa pine (*Pinus*
123 *ponderosa*; de la Mata et al. 2022) and single-leaf pinyon (*Pinus monophylla*; Vasey, Weisberg,
124 & Urza 2022) populations in response to temperature fluctuations and microhabitat heterogeneity.
125 Despite this strong theoretical emphasis and empirical basis, robust statistical methods have not
126 yet been developed for detecting complex patterns of selection on the RNs of labile traits.

127 Many of the phenotypes commonly studied by evolutionary ecologists are highly labile (i.e.
128 exhibit high degrees of reversible plasticity; Scheiner, 1993b) in response to the local
129 environment. This means that repeatable individual differences in the RN underlying these traits
130 tend to account for only a modest proportion of the total variation observed in measurements
131 across space and time (Bell, Hankison, & Laskowski 2009; Fanson & Biro 2015; Cauchoix et al.
132 2018). This is expected, given that labile traits are often adapted to facilitate flexible responses
133 toward fitness-relevant variation in the environment (Scheiner 1993b), such as by up-regulating
134 circulating testosterone in response to social challenges (Wingfield et al. 1990; Eisenegger,
135 Haushofer, & Fehr 2011), temporarily inducing a fear state in response to odor cues of predation
136 (Mathuru et al. 2012), or regulating alloparental care in response to the quality of the local
137 environment (Guindre-Parker & Rubenstein, 2018; Martin et al. 2020). Conversely, labile traits
138 may also be prone to maladaptive plasticity in response to novel or extreme environmental
139 stressors (e.g. Ghalambor et al. 2015). As such, single measures of labile phenotypes tend to
140 reflect within- rather than among-individual variation, potentially biasing empirical estimates of
141 trait (co)variances and selection gradients estimated across heterogeneous environments
142 (Brommer 2013; Dingemanse & Dochtermann 2013; Niemelä & Dingemanse 2018; Royauté et
143 al. 2018), leading to inaccurate inferences about adaptive evolution (Dingemanse, Araya-Ajoy, &
144 Westneat 2021; Martin & Jaeggi 2022). Classical approaches such as the Lande and Arnold
145 (1983) regression framework do not partition repeatable and non-repeatable differences across
146 phenotypic measurements and, as a consequence, may lead to downwardly biased estimates of
147 selection gradients for labile traits in field research (Dingemanse et al. 2021). Classical methods
148 can also be biased by unmeasured, within-individual environmental effects on fitness and
149 phenotype that generate spurious signals of selection (Scheiner et al. 2002; Stinchcombe et al.
150 2002). Using these methods to estimate selection on labile traits with single measures, averages

151 of raw data, or point estimates in multi-stage analyses can, therefore, increase the risk of biased
152 evolutionary inference (Hadfield et al. 2010), particularly when attempting to understand the
153 adaptation of RNs underlying observed phenotypes across environments.

154 Fortunately, generalized linear mixed-effects models (GLMMs) provide a flexible toolkit for
155 estimating RNs from empirical data, as well as for modelling the effects of RNs on fitness and
156 other biological outcomes of interest. Current variance-partitioning methods rely on the use of
157 multi-response/multivariate GLMMs with covarying random effects to model selection, which
158 effectively account for uncertainty in individuals' RNs and their estimated effects (Hadfield et al.
159 2010). This approach has been repeatedly introduced to selection studies of RNs in variety of
160 contexts, demonstrating its broad applicability (e.g. Brommer, Kontiainen, & Pietiäinen 2012;
161 Houslay & Wilson 2017; Arnold, Nicotra, & Kruuk 2019), and can be further extended to provide
162 a veritable treasure chest of biological insights (Blows 2007). For example, such models can be
163 used to identify trajectories of phenotypic conservation and divergence among closely related
164 populations (Royauté, Hedrick, & Dochtermann 2020), discover latent behavioral characters
165 among multiple traits (Araya-Ajoy & Dingemans 2014; Martin et al. 2019), or calculate genetic
166 responses to directional selection (Stinchcombe, Simonsen, & Blows 2014). Therefore, multi-
167 response GLMMs with covarying random effects can be used to accomplish many empirical goals
168 with relative ease, while also avoiding statistical bias due to uncertainty in RNs.

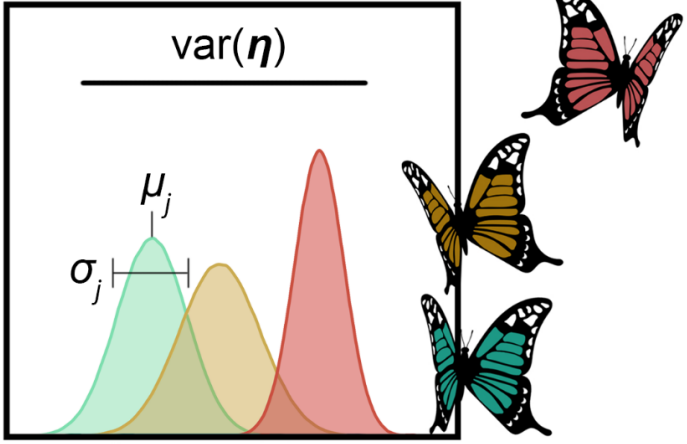
169 Despite their benefits, these commonly used GLMMs cannot detect nonlinear selection on
170 RNs (i.e. disruptive, stabilizing, and correlational selection) because the random effect covariance
171 is defined as an average measure of linear dependency among fitness and phenotype. By failing
172 to describe the curvature of the adaptive landscape, and thus the ecological phenomena
173 generating fitness saddles, ridges, domes, and cliffs (Lande & Arnold, 1983; Blows & Brooks,
174 2003; Blows 2007; Vercken et al., 2012), random effect models can provide an incomplete and
175 potentially misleading perspective on the biological processes driving and constraining
176 multivariate evolution. In non-randomized experiments or field settings, ignoring nonlinear
177 selection can further generate biased estimates of directional selection gradients, in addition to
178 biased predictions of the evolutionary response to selection on the expectations and
179 (co)variances of RN parameters (Arnold et al., 2001; Morrissey et al., 2012; Pick et al., 2022).
180 Therefore, despite their clear utility, current covarying random effects models can also limit robust
181 tests of adaptive theory, which often predicts that stabilizing, disruptive, and/or correlational
182 selection will shape RN evolution (e.g. Wagner et al., 1997; Gavrillets & Hastings, 1994). This

183 inhibits accurate predictions of phenotypic evolution more generally (Bulmer 1971; Lande &
184 Arnold 1983; Arnold, Pfrender, & Jones, 2001; Villemereuil et al., 2020).

185 Here we address this challenge by introducing multi-response/multivariate GLMMs for
186 unbiased estimation of nonlinear selection on RNs, building on well-established approaches to
187 estimating linear selection (e.g. Brommer, Kontiainen, & Pietiäinen 2012; Houslay & Wilson 2017;
188 Arnold, Nicotra, & Kruuk 2019; Araya-Ajoy, Dingemanse, Westneat, & Wright 2023). The
189 proposed GLMMs are applicable to any labile and repeatedly measured phenotype. We begin by
190 reviewing so-called double hierarchical GLMMs for estimating RNs from longitudinal, repeated
191 measures data (Westneat, Schofield, & Wright, 2013; O’Dea et al. 2021) and formally introduce
192 multi-response/multivariate models estimating linear and nonlinear selection on RNs, applicable
193 to both Gaussian and non-Gaussian measurements. We then consider their implementation in a
194 Bayesian framework, using a simulation-based calibration procedure to validate that the proposed
195 models are unbiased for statistical inference. We also explore statistical power for Bayesian
196 hypothesis tests across a range of sampling designs and selection effect sizes. Guided tutorials
197 are further provided (see **data availability**) to aid researchers in implementing and interpreting
198 these models for their own data using the Stan statistical programming language (Carpenter et
199 al. 2017).

200 **Figure 1. Empirical estimation of reaction norms.** Repeatable among-individual differences $\text{var}(\eta)$ (*top*
201 *left*) in the expected value μ and dispersion σ of observed phenotype \mathbf{z} can be predicted with a RN model
202 (*top right*) using link functions \mathbf{g} and three (or more) distinct parameters: RN intercept parameters μ_0
203 describing each individual’s average phenotype across a mean-centered environment or in the absence of
204 the environment (i.e. when the environmental state $x = 0$); RN slope parameters β_x describing each
205 individual’s systematic change in phenotype across environmental states \mathbf{x} ; and RN residual parameters
206 σ_0 reflecting each individual’s degree of stochastic variability (or, conversely, their predictability/precision)
207 in phenotype within a given environment. See **Eq. 1** for index rather than matrix notation. These parameters
208 will be unknown in empirical research and must be estimated using raw measurements (teal circles) across
209 environmental states (*bottom left*). An example is shown for a simple linear RN with a log-link on the
210 dispersion of a normal distribution, so that an individual’s residual parameter, expressed as a variance on
211 the squared log scale $\text{sqrt}(\exp(\sigma_0 + \sigma_{0j}))$, is proportional to (\propto) the spread of observed residuals on the
212 original data scale, as shown here by a 95% credible interval. Failure to account for uncertainty around
213 point estimates of individual j ’s RN parameters (*bottom right*) leads to anti-conservative inference and
214 increased risk of false positives (Hadfield et al. 2010).

Repeatable among-individual differences

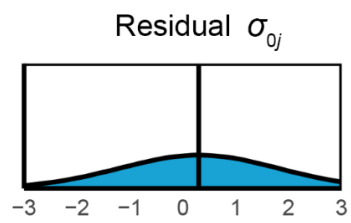
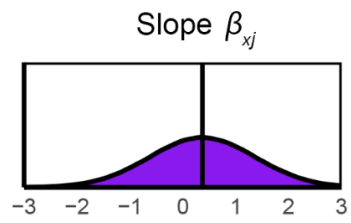
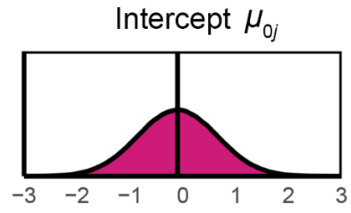
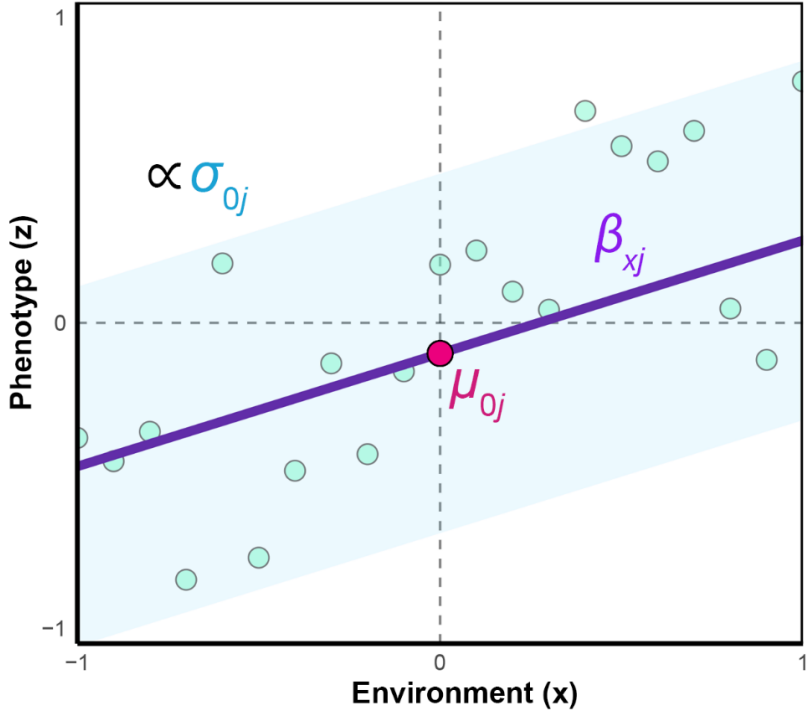


Reaction norm (RN) model

$$z \sim f(\mu, \sigma)$$
$$g(\mu) = \mu_0 + \beta_x \cdot x$$
$$g(\sigma) = \sigma_0$$

RN uncertainty for individual j

Linear RN point estimates for individual j



216 **Table 1.** Notation and terminology.

Term	Symbol	Description
Individual reaction norm (RN)	$f(\mu, \sigma)$	A probabilistic function f with parameters predicting the expectation μ and dispersion σ of an individual's phenotype in response to a measurable aspect of the environment.
RN intercept	μ_0, μ_{0j}	The expected phenotype in the average environment or in the absence of an environmental factor. Individual RN intercept μ_{0j} is expressed as a deviation from population RN intercept μ_0 .
RN slope	β_x, β_{xj}	The expected change in phenotype in response to a measured environment x . Individual RN slope β_{xj} is expressed as a deviation from the population average slope β_x .
RN residual	σ_0, σ_{0j}	The magnitude of stochastic variability in phenotype within a given environment, i.e. the inverse of predictability (O'Dea et al., 2021) and precision (Hansen et al., 2006). Individual RN residual parameter σ_{0j} is expressed as a deviation from population average residual parameter σ_0 , which together determine the magnitude of variation in observed residuals.
RN trait value/ character state	η_{jt}	The repeatable trait value predicted by individual j 's reaction norm being expressed within the environmental state at time t . This context-specific trait value is also referred to as a character state in quantitative genetics.
Repeatable among-individual differences	$\text{var}(\eta)$	The total amount of among-individual variation in the phenotype available to natural selection over the sampling period, which reflects consistent individual differences in RN expression across environments (i.e. the variance of character states).
Link functions	$g_\mu, g_\sigma, g_\theta$	Transformations that facilitate modelling of non-Gaussian phenotypes and fitness measures on a linear scale.
Fitness	$W, f(\theta, \delta)$	A measure of an individual's observed survival, reproduction, and/or performance W , as predicted by the expectation θ and dispersion δ parameters of distribution f . These quantifiable 'fitness components' are used to approximate the repeatable, differential rate of zygote production across individuals.
Directional selection	\mathbf{b}, β	Selection gradients β quantify the magnitude of direct selection on the population means of reaction norm parameters. Regression coefficients \mathbf{b} approximate these effects on the transformed scale of a GLMM.
Quadratic selection	\mathbf{q}, γ	Selection gradients γ quantify the magnitude of direct selection on the (co)variances of reaction norm parameters. Regression coefficients \mathbf{q} approximate these effects on the transformed scale of a GLMM.
Fluctuating selection	$\Delta\beta, \Delta\gamma$	Environmental change that shifts the magnitude of selection on RNs $\Delta\beta, \Delta\gamma$ across space and/or time

217 **Modelling nonlinear selection on labile traits**

218 The models we propose in this section are straightforward extensions of the multi-
219 response/multivariate random effects GLMMs discussed above. Our trait-based approach shifts
220 estimation of fitness effects from random effect covariances to flexibly parameterized linear and
221 nonlinear selection coefficients. This approach builds on a long tradition of measurement error
222 models in biostatistics (Loken & Gelman 2017; Ponzi et al. 2018; Martin & Jaeggi 2022), also
223 known as structural equation (Bollen & Noble 2011; Araya-Ajoy & Dingemanse 2014; Martin et
224 al. 2019) or errors-in-variables models (Dingemanse et al. 2021), which allow for latent trait values
225 such as RN intercept, slope, and residual parameters to simultaneously affect multiple response
226 models. The basic structure of these models has been previously introduced in the broader
227 context of phenotypic selection analysis by Ponzi et al. (2018), Dingemanse et al. (2021), and
228 Araya-Ajoy et al. (2023), who considered Gaussian models of selection on repeatable trait values.
229 Here, we generalize and extend these models to allow for estimating (non)linear selection on RN
230 intercepts, slopes, and residuals (and any other distributional parameters of interest), as well as
231 to estimate directional and quadratic selection gradients on RN parameters with non-Gaussian
232 phenotype and fitness measures.

233 **Reaction norm model**

234 The first step in any selection analysis is to define the trait of interest. For repeatedly
235 expressed traits that exhibit plasticity, the ‘traits’ of interest may be latent properties of a RN,
236 which researchers can estimate as functional parameters. As shown in **Figure 1**, individual
237 variation in a linear RN can be decomposed into underlying repeatable differences in individuals’
238 RN intercept μ_0 , slope β_x , and residual parameters σ_0 . Note that we use β_x here to reference
239 any slope defined over a non-social environmental state (see Martin & Jaeggi, 2022 for a
240 treatment of social effects). **Table 1** provides a glossary of formal notation and terminology used
241 throughout the paper. GLMMs effectively describe the RNs of non-Gaussian phenotypes using
242 additive linear functions on a transformed latent scale (Bolker et al., 2009; Villemereuil et al.,
243 2016). Extensive prior work has been done on appropriate study design and GLMM
244 implementation for RN research in evolutionary ecology (e.g. see Nussey, Wilson, & Brommer
245 2007; Martin, Nussey, Wilson, & Réale, 2010; Dingemanse & Dochtermann 2013; O’Dea et al.
246 2021 among others). Therefore, we avoid reviewing this material in detail here, instead focusing
247 on the introduction of a general form and notation for RN models of any labile trait.

248 Consider a GLMM for repeated measure t of individual j , who expressed labile phenotype
 249 z_{jt} in environmental state x_{jt} . The distribution of measurements can be predicted using a
 250 probability function $f(\mu, \sigma)$ with mean, location, or rate parameter μ and dispersion, shape, or
 251 scale parameter σ (e.g. as with normal, gamma, and beta distributions). Link functions g_μ and g_σ
 252 are used for modelling the vectors $\boldsymbol{\mu}$ and $\boldsymbol{\sigma}$ across observations so that the RN parameters $\boldsymbol{\mu}_0$,
 253 $\boldsymbol{\beta}_x$, and $\boldsymbol{\sigma}_0$ can be expressed as additive linear effects on a transformed scale, irrespective of the
 254 assumed distribution of the raw data. For instance, $g_\mu = \text{identity}(\mu)$ and $g_\sigma = \log(\sigma^2)$ are
 255 sensible choices for a Gaussian measure. The generalized form of the model is given by

$$256 \quad z_{jt} \sim f(\mu_{jt}, \sigma_j) \quad (1)$$

$$257 \quad g_\mu(\mu_{jt}) = \mu_0 + \mu_{0j} + (\boldsymbol{\beta}_x + \boldsymbol{\beta}_{xj})x_t$$

$$258 \quad g_\sigma(\sigma_j) = \sigma_0 + \sigma_{0j}$$

$$259 \quad [\boldsymbol{\mu}_0^\top, \boldsymbol{\beta}_x^\top, \boldsymbol{\sigma}_0^\top]^\top \sim \text{MVN}(\mathbf{0}, \mathbf{P}): \mathbf{P} = \begin{bmatrix} \text{var}(\boldsymbol{\mu}_0) & \dots & \dots \\ \text{cov}(\boldsymbol{\beta}_x, \boldsymbol{\mu}_0) & \text{var}(\boldsymbol{\beta}_x) & \vdots \\ \text{cov}(\boldsymbol{\sigma}_0, \boldsymbol{\mu}_0) & \text{cov}(\boldsymbol{\sigma}_0, \boldsymbol{\beta}_x) & \text{var}(\boldsymbol{\sigma}_0) \end{bmatrix}$$

260 where \top indicates the transpose operator. Here μ_0 , β_x , σ_0 are the average values for the RN
 261 intercept, slope, and residual parameters in the population, expressed on the scale of the link
 262 functions. Repeatable individual differences in RN parameters are in turn estimated as deviations
 263 from these averages using random effects μ_{0j} , β_{xj} , and σ_{0j} . For simplicity, the model assumes
 264 environmental exposures \mathbf{x} are randomized across individuals, but it may be necessary in non-
 265 experimental contexts to center covariates within individuals for appropriate scaling of RN slopes
 266 (Schaeffer, 2004; van de Pol & Wright 2009; Araya-Ajoy, Mathot, & Dingemanse, 2015; Westneat
 267 et al., 2020; Fay, Martin, & Plard 2022). The magnitude of among-individual (co)variance in these
 268 RN parameters is described by the \mathbf{P} matrix. See **Box 1** for further discussion of the RN residual
 269 parameter.

Box 1. Interpreting among-individual differences in RN residuals.

271 The functional role of the RN residual parameters σ_0 can be ambiguous because these individual
272 effects are modelled on the dispersion σ of the phenotypic distribution, rather than the expectation
273 μ (Eq. 1). Phenotypic variance due to dispersion is generally interpreted as noise or measurement
274 error around individuals' repeatable mean trait values (Brommer 2013), which are determined by
275 the expression of RN intercepts μ_0 and slopes β_x across measured environments. However, the
276 residuals of labile traits may also contain repeatable and fitness-relevant variation in how
277 organisms intrinsically regulate their phenotype (Westneat, Wright, & Dingemanse 2015), such
278 as in their assessment and response toward developmental noise within a given environment
279 (Gavrilets & Hastings, 1994; Hansen et al., 2006; Mitchell et al. 2021). Such repeatable *among-*
280 individual differences in *within*-individual variation, described by σ_0 , may arise from a variety of
281 mechanisms regulating patterns of stochastic expression in behavior or other labile traits
282 (Prentice, Houslay, Martin & Wilson, 2020). For instance, stochasticity can be generated through
283 the repeatable activities of the organism, such as by random sampling of the environment, which
284 can be shaped via reinforcement and punishment to facilitate adaptive learning in novel or
285 uncertain environments (Niv et al. 2002; Barrett 2011; Wright et al., 2022). Consequently, intrinsic
286 variability may evolve in conjunction with learning mechanisms to track unpredictable shifts in
287 fitness optima during development (Borenstein, Feldman, & Aoki 2008). Predation may also select
288 for greater variability in movement, so as to reduce predators' capacity to predict prey escape
289 trajectories (Hugie, 2003; Moore et al. 2017), while reduced variability may instead be adaptive
290 for reputation formation and trust in repeated social interactions (McNamara & Leimar, 2010).
291 Stochasticity may also result from exogenous factors, such that individual differences in σ_0 reflect
292 how organisms regulate in response to the environment. For example, when environmental states
293 fluctuate rapidly in an unpredictable and uncontrollable manner, negative selection may act on
294 the RN residual parameter to promote phenotypic canalization, decreasing susceptibility of the
295 phenotype to developmental perturbation (Flatt 2005; Siegal & Leu 2014; Westneat et al., 2015).

296 In empirical research, it will often be challenging to distinguish variance in residuals due to
297 intrinsically stochastic variability or unmeasured processes of cue-induced plasticity and
298 individual-by-environment interaction (Westneat et al. 2015; Prentice et al., 2020). Estimates of
299 $\text{var}(\sigma_0)$ in the field may, for example, reflect repeatable functional interactions between
300 unmodelled RN slopes and stochastic environmental exposures. Therefore, caution is warranted
301 when inferring the mechanistic underpinnings of σ_0 outside of well-controlled experiments. Poorly
302 specified statistical models, in which predicted residual processes do not accurately describe

303 observed phenotypic variance, will also inhibit accurate biological inference of RNs (Mitchell,
 304 Dujon, Beckmann, & Biro, 2020; Ramakers, Visser, & Gienapp, 2020). Nonetheless, to the degree
 305 that individual differences in residuals are repeatable across time and not due to unbalanced
 306 sampling or pseudo-repeatability (Dingemanse & Dochtermann 2013), selection can still shape
 307 RN residuals, irrespective of whether within-individual deviations arise from mechanisms of
 308 intrinsically stochastic or cue-induced trait expression. Therefore, we suggest that researchers in
 309 both observational and experimental systems focus their attention on functionally interpreting and
 310 operationally defining RN residual parameters with respect to theoretically motivated RN slopes,
 311 defined over measured dimensions of environmental change (Figure 1).

312

313 **Box 2. Repeatable among-individual differences due to RNs.**

314 Selection on the RNs of labile traits can only occur if individuals differ in their intercepts, slopes,
 315 and residual parameters across time. The covariance matrix \mathbf{P} in Eq. 1 describes these
 316 repeatable among-individual differences and, therefore, ultimately determines the total amount of
 317 trait (co)variation available to natural selection on phenotype \mathbf{z} over the sampling period of
 318 interest, given that RN parameters μ_0 , β_x , and σ_0 predict how organisms will repeatedly express
 319 their phenotype within and across environments. We denote the total magnitude of repeatable
 320 among-individual differences in \mathbf{z} due to RNs as $\text{var}(\boldsymbol{\eta})$, which in the general case sets an upper
 321 limit on the heritability of a phenotype due to direct genetic effects (see Bijma, 2011 for social
 322 traits) and thus provides a useful phenotypic proxy of the evolvability of a trait (Martin et al., 2023).
 323 The trait values $\boldsymbol{\eta}$ represent the repeatable character states that organisms are expected to
 324 express within and across sampled environments, as predicted by their individual RNs (Fig. 1 top
 325 left). Conversely, any variance in observed trait values \mathbf{z} due to non-repeatable effects $\text{var}(\boldsymbol{\xi}) =$
 326 $\text{var}(\mathbf{z}) - \text{var}(\boldsymbol{\eta})$ introduces noise into the estimation of selection gradients defined across
 327 sampled environments. Failure to distinguish non-repeatable $\text{var}(\boldsymbol{\xi})$ and repeatable $\text{var}(\boldsymbol{\eta})$
 328 variance in measured phenotypes $\text{var}(\mathbf{z}) = \text{var}(\boldsymbol{\eta}) + \text{var}(\boldsymbol{\xi})$ can thus lead to biased estimates of
 329 directional $\boldsymbol{\beta}^*$ and quadratic $\boldsymbol{\gamma}^*$ selection gradients (Figure 2). For evolutionary ecologists,
 330 correlations between fitness and phenotype that are repeatable over time and potentially heritable
 331 across generations will generally be of primary interest, motivating partitioning of $\text{var}(\boldsymbol{\eta})$ from
 332 $\text{var}(\mathbf{z})$ with a GLMM (Martin & Jaeggi 2022).

333 O'Dea et al. (2022) and de Villemereuil et al. (2016), among others, provide exact analytic
 334 solutions and numeric methods for calculating $\text{var}(\boldsymbol{\eta})$ with many commonly used GLMMs. For the

335 general case, $\text{var}(\boldsymbol{\eta})$ can always be approximated on the original data scale, irrespective of model
 336 complexity, by using simulation to compare the variance of model predicted phenotypic
 337 distributions in the presence $\text{var}(\mathbf{z}_{\text{pred}})_{\boldsymbol{\eta}}$ and absence $\text{var}(\mathbf{z}_{\text{pred}})_{-\boldsymbol{\eta}}$ of repeatable individual
 338 effects $\boldsymbol{\mu}_0$, $\boldsymbol{\beta}_x$, and $\boldsymbol{\sigma}_0$, using a large number of random samples.

$$339 \quad \text{var}(\boldsymbol{\eta}) \approx \text{var}(\mathbf{z}_{\text{pred}})_{\boldsymbol{\eta}} - \text{var}(\mathbf{z}_{\text{pred}})_{-\boldsymbol{\eta}} \quad (2)$$

340 Model predictions can also be used to approximate the total repeatability of among-individual
 341 differences in the phenotype on the original data scale for any GLMM

$$342 \quad R_{\boldsymbol{\eta}} \approx \frac{\text{var}(\boldsymbol{\eta})}{\text{var}(\mathbf{z}_{\text{pred}})_{\boldsymbol{\eta}}} \quad (3)$$

343 The bias of estimated selection gradients will increase as the $R_{\boldsymbol{\eta}}$ of a phenotype decreases and
 344 $\text{var}(\boldsymbol{\xi})$ in turn increases (Spearman, 1904; Searle, 1961). Therefore, failure to remove non-
 345 repeatable causes of variation from observed phenotypic measures is a particularly serious issue
 346 when estimating selection on labile traits across heterogeneous environments (Figure 2;
 347 Dingemanse et al. 2021).

348 (Non)linear selection model

349 To model selection on the individual-specific RN parameters μ_{0j} , β_{xj} , and σ_{0j} , the RN
 350 GLMM in [Eq. 1](#) can be expanded to include an additional response model predicting measure t
 351 of fitness component or proxy W . Linear \mathbf{b} and quadratic \mathbf{q} selection coefficients, as well as other
 352 more complex forms of nonlinear selection, can then be estimated directly for the RN parameters.

$$353 \quad z_{jt} \sim f(\mu_{jt}, \sigma_j) \quad (4)$$

$$354 \quad g_\mu(\mu_{jt}) = \mu_0 + \mu_{0j} + (\beta_x + \beta_{xj})x_t$$

$$355 \quad g_\sigma(\sigma_j) = \sigma_0 + \sigma_{0j}$$

$$356 \quad [\boldsymbol{\mu}_0^\top, \boldsymbol{\beta}_x^\top, \boldsymbol{\sigma}_0^\top]^\top \sim \text{MVN}(\mathbf{0}, \mathbf{P})$$

$$358 \quad W_{jt} \sim f(\theta_{jt}, \delta)$$

$$359 \quad g_\theta(\theta_{jt}) = W_0 + W_{0j} + b_1\mu_{0j} + b_2\beta_{xj} + b_3\sigma_{0j}$$

$$360 \quad + q_1\mu_{0j}^2 + q_2\beta_{xj}^2 + q_3\sigma_{0j}^2 + q_4\mu_{0j}\beta_{xj} + q_5\mu_{0j}\sigma_{0j} + q_6\beta_{xj}\sigma_{0j}$$

$$361 \quad \mathbf{W}_0 \sim \text{N}(0, \text{sd}(\mathbf{W}_0))$$

362 Fitness W for individual j at measurement t is described by a GLMM with expectation parameter
 363 θ and dispersion parameter δ . The full model thus estimates the RN parameters and their
 364 accompanying selection coefficients in the fitness model simultaneously using a multivariate
 365 analysis. [Figure 2](#) visualizes this model structure and explains how it avoids bias by partitioning
 366 repeatable sources of (non)linear association between phenotype and fitness. Parameter W_0 is
 367 the average fitness on the transformed scale given by link function g_θ . When repeated fitness
 368 measures t are available, an individual random effect W_{0j} should be estimated to capture
 369 repeatable among-individual differences in fitness that are not due to the modelled phenotypes
 370 (i.e. unexplained selection). If only a single fitness measure is available, $\text{sd}(\mathbf{W}_0)$ cannot be
 371 identified separately from fitness residual dispersion δ , so these effects should instead be
 372 excluded from the analysis.

373 The polynomial regression in [Eq. 4](#) can be used to infer short-term population trajectories
 374 on the adaptive landscape, under the assumption that a quadratic function effectively
 375 approximates the local shape of the individual selection surface on the latent transformed scale
 376 ([Lande & Arnold 1983](#); [Phillips & Arnold 1989](#)). However, the values of the \mathbf{b} and \mathbf{q} regression

377 coefficients should only be interpreted as measures of directional and quadratic selection
 378 gradients when fitness is a mean-scaled Gaussian response, after appropriately scaling the
 379 coefficients (see [Stinchcombe et al. 2008](#); [Dingemanse et al. 2021](#) for details). Analytic
 380 expressions can also be used for direct interpretation of coefficients in a log-normal fitness model
 381 ([Bollen, Morrissey, & Kruuk 2019](#)). However, in the general case, it will be necessary to further
 382 process regression coefficients from the fitness model before making quantitative inferences
 383 about directional and quadratic selection on the scale of the original data, which is generally of
 384 greater biological interest.

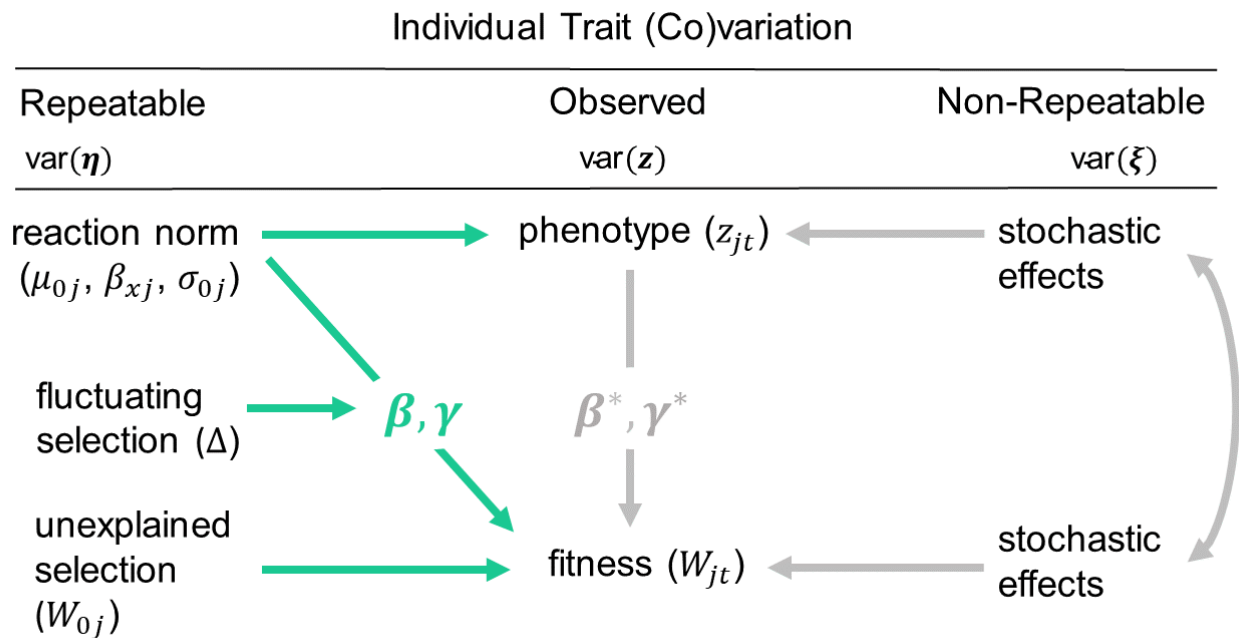
385 Following [Lande and Arnold \(1983\)](#) and [Morrissey and Sakrejda \(2013\)](#), directional β and
 386 quadratic γ selection gradients can be numerically calculated for any GLMM by taking the first ∂
 387 and second ∂^2 partial derivatives of the estimated fitness function with respect to the expected
 388 population-level RN parameters $\bar{\mu}_0$, $\bar{\beta}_x$, and $\bar{\sigma}_0$.

$$389 \quad \beta_{\mu_0} = \frac{\partial E(\bar{W}, \bar{\mu}_0)}{\partial \bar{\mu}_0} \bar{W}^{-1} \dots \quad \gamma_{\mu_0} = \frac{\partial^2 E(\bar{W}, \bar{\mu}_0)}{\partial \bar{\mu}_0 \delta \bar{\mu}_0} \bar{W}^{-1} \dots \quad \gamma_{\beta_x \sigma_0} = \frac{\partial^2 E(\bar{W}, \bar{\beta}_x)}{\partial \bar{\beta}_x \delta \bar{\sigma}_0} \bar{W}^{-1} \quad (5.1)$$

390 where \bar{W} is the expected population fitness on the original data scale, as predicted by the fitness
 391 function defined with \mathbf{b} and \mathbf{q} coefficients on the link scale in [Eq. 4](#). The directional gradients β_{μ_0} ,
 392 β_{β_x} , and β_{σ_0} indicate the direction and magnitude of selection on the expected values of
 393 population RN parameters, with respect to the original untransformed scale of the data. Quadratic
 394 selection gradients γ_{μ_0} , γ_{β_x} and γ_{σ_0} in turn indicate convex or concave curvature in the selection
 395 surface shaping the variance of RN parameters ([Stinchcombe et al. 2008](#)); and $\gamma_{\mu_0 \beta_x}$, $\gamma_{\mu_0 \sigma_0}$, and
 396 $\gamma_{\beta_x \sigma_0}$ indicate further curvature due to the presence of correlational selection between RN
 397 parameters ([Blows & Brooks 2003](#)). These gradients can be expressed in standardized units for
 398 effect size comparison between traits and parameters using the appropriate variances and
 399 standard deviations ([Lande & Arnold 1983](#))

$$400 \quad \beta_{\mu_0}^{\text{sd}} = \beta_{\mu_0} \text{sd}(\mu_0) \dots \quad \gamma_{\mu_0}^{\text{sd}} = \gamma_{\mu_0} \text{var}(\mu_0) \dots \quad \gamma_{\beta_x \sigma_0}^{\text{sd}} = \gamma_{\beta_x \sigma_0} \text{sd}(\beta_x) \text{sd}(\sigma_0) \quad (5.2)$$

401 Standardized gradients are particularly useful for GLMMs because the magnitude of variances
 402 may differ appreciably between the distinct transformed link scales used for estimating RNs and
 403 selection, which makes it challenging to meaningfully distinguish between small and large effect
 404 sizes across models.



406

406 **Figure 2. Removing non-repeatable effects from selection gradients.** The diagram shows causal
 407 pathways (directional arrows) by which repeatable (green) and non-repeatable (grey) effects can
 408 influence selection gradients of fitness (W) on phenotype (z). Non-repeatable, stochastic effects
 409 influence both fitness and phenotype (directional arrows) and may be correlated (double-headed
 410 arrow), introducing statistical noise into the selection analysis. This leads to biased directional β^* and
 411 quadratic gradients γ^* when observed variance in the phenotype $\text{var}(\mathbf{z})$ is used to estimate selection
 412 across environments. However, if the (non)linear relationships between phenotype and fitness are
 413 modelled independently of stochastic effects on the phenotype $\text{var}(\boldsymbol{\xi})$, using RN parameters $\mu_0, \beta_x,$
 414 and σ_0 (Eq. 1-4), unbiased selection gradients β and γ can be estimated (Eq. 5) directly for
 415 repeatable among-individual differences in the phenotype $\text{var}(\boldsymbol{\eta})$ (see Box 2). Spatiotemporal
 416 fluctuations Δ in these selection gradients can also be described by additional coefficients (Eq. 6), and
 417 any repeatable among-individual differences in fitness unexplained by RN parameters can be
 418 estimated with random effects W_0 when repeated fitness measures are available (Eq. 4).

419

Model extensions

420 Simplified models are presented above ([Eq. 1](#), [Eq. 4](#)) to aid interpretation, but it will often
421 be necessary to specify more complex models for explaining empirically observed variation in
422 fitness and phenotype. Various model extensions can be straightforwardly accomplished using
423 the basic toolkit of GLMMs and related regression frameworks, along with appropriate study
424 design and sufficient repeated sampling for reliable estimation. Below we briefly consider three
425 key areas for model extension and provide references for further consideration. Implementation
426 for social traits and interactions is discussed by [Martin and Jaeggi \(2022\)](#).

427 Fluctuating selection

428 Fluctuating selection on RNs may occur due to variation in the density of mates and
429 competitors, resource availability and seasonality, bodily condition and age, the availability of local
430 niches, or any other state that modulate the fitness costs and benefits of labile traits ([Houston &
431 McNamara, 1999](#); [Sih et al., 2015](#)). Fluctuating selection is also expected to be a key mechanism
432 for explaining patterns of macroevolutionary stasis ([Estes & Arnold, 2007](#)), the adaptive evolution
433 of phenotypic plasticity ([de Jong, 1995](#); [King & Hadfield, 2019](#); [Martin et al., 2024](#)), and the
434 evolutionary maintenance of individual and genetic variation within populations more generally
435 (e.g. [Sasaki & Ellner, 1997](#); [Dingemans & Wolf, 2010](#); [Wolf & Weissing, 2010](#); [Wright et al.,
436 2019](#); [Abdul-Rahman, Tranchina, & Gresham, 2021](#); [Martin et al., 2023](#)). As previously noted,
437 quantitative genetic theory has demonstrated the mathematical equivalence of models for
438 selection on character states and RNs. A key finding from this theoretical work is that fluctuating
439 selection on character states expressed within environments generates directional and quadratic
440 selection on RN parameters across environments ([Gavrilets & Scheiner, 1993](#); [de Jong, 1995](#)).
441 Therefore, estimating non-zero directional and quadratic selection on a RN parameter via [Eq. 4-
442 5](#) implies that selection on the phenotype is fluctuating with respect to the environment over which
443 the RN parameter is defined ([Martin et al., 2024](#)). For example, the degree to which density-
444 dependent selection on character states fluctuates across the environments encountered by a
445 population will be proportional the average directional and quadratic selection on the RN slope β_x
446 defining how the phenotype changes with respect to density x . In general, this means that the RN
447 selection model can be used to infer the presence of fluctuating phenotypic selection with much
448 fewer parameters than an equivalent character state model ([de Jong, 1995](#)).

449 These considerations suggest that RN selection analyses will often not require estimating
 450 additional parameters beyond the main linear \mathbf{b} and nonlinear \mathbf{q} effects on RN parameters to
 451 accurately describe patterns of fluctuating selection on the expressed phenotype. However, in the
 452 presence of environmental change, the magnitude and pattern of fluctuating character state
 453 selection experienced by a population may also vary across space and time, which is expected
 454 to result in fluctuating selection on RN parameters (**Figure 2**). In some systems, the putative
 455 environmental causes of fluctuating selection will be directly measured, while in others, it may be
 456 informative to estimate spatiotemporal heterogeneity in RN selection even if the underlying
 457 causes are not directly measured (Reynolds, de Los Campos, Egan, & Ott 2016). For example,
 458 long-term field studies can be used to investigate the adaptive maintenance of RN variation by
 459 yearly fluctuations in selection, even if the mechanisms underpinning these fluctuations remain
 460 unclear (de Villemereuil et al., 2020; Mouchet et al., 2021). To incorporate such effects, the basic
 461 fitness model (**Eq. 4**) can be extended by including fixed or random interaction effects on the
 462 selection coefficients, which will estimate continuous or discrete fluctuations $\Delta\beta$ and $\Delta\gamma$ (**Figure**
 463 **2**) across space and time. For example,

$$464 \quad g_{\theta}(\theta_{jt}) = W_0 + W_{0j} + (b_1 + b_{1x}x_t + u_{tb_1})\mu_{0j} + \dots (q_1 + q_{1x}x_t + u_{tq_1})\mu_{0j}^2 + \dots \quad (6)$$

465 where b_{1x} and q_{1x} describe how the (non)linear selection coefficients change as a function of x_t ,
 466 and u_{tb_1} and u_{tq_1} describe changes due to a random factor at time t .

467 **Adjusted and nonlinear effects**

468 As with any regression analysis, additional fixed and random effects may need to be
 469 adjusted for to facilitate appropriate biological inference. Predation may, for instance, cause
 470 differential mortality as a function of repeatable differences in behavior across sex and age
 471 classes, but this selection will not generate an evolutionary response on behavioral variation
 472 within sexes or age classes. This motivates estimating repeatable individual variation adjusted for
 473 the effects of sex and age, among other commonly studied factors such as size and morphology
 474 (Bolnick et al., 2003). Unadjusted environmental effects on fitness and phenotype can also bias
 475 estimates of selection and among-individual variation in both field and laboratory settings
 476 (Scheiner et al. 2002; Stinchcombe et al., 2022; Kinsler et al., 2023; Munar-Delgado et al., 2023).
 477 It is, therefore, often useful to include additional environmental covariates (e.g. average
 478 temperature and rainfall, date within season, resource availability), including potential interaction
 479 effects, and random factors (e.g. nesting site, spatial position, batch, observer identity) to adjust

480 fitness variation during the selection analysis. As discussed in **Box 2**, model predictions can
 481 always be used to quantify and better understand how adjusting for these effects changes the
 482 repeatable variation available to selection in any multivariate GLMM.

483 Relationships between fitness, phenotype, and the local environment may also be best
 484 described by additional terms beyond quadratic regression coefficients. For example, RN slopes
 485 of thermoregulatory and life history traits such as growth rate are often highly nonlinear in
 486 response to temperature (Oomen & Hutchings, 2022), violating the assumption of **Eq. 4** that
 487 individuals' phenotypic deviations from the linear RN slope β_x are multivariate normally
 488 distributed. Polynomials (Henderson, 1982; Yamahira, Kawajiri, Takeshi, & Irie, 2007) or
 489 generalized additive effects such as splines or Gaussian processes (Schluter & Nychka, 1994;
 490 Sigourney, Munch, & Letcher, 2012; Pederson, Miller, Simpson, & Ross, 2019; Catalina, Bürkner,
 491 & Vehtari, 2020) can be used to account for nonlinearity in the population RN and ensure the
 492 statistical model more accurately predicts observable phenotypic and fitness variation. In the
 493 general case, the basic model (**Eq. 4**) can be expanded to include any generalized additive
 494 function $s()$ describing how expected phenotypic μ_{jt} or fitness values θ_{jt} change in response to
 495 the environment

$$496 \quad g_{\mu}(\mu_{jt}) = \mu_0 + \mu_{0j} + s(x_t) + \beta_{xj}x_t \quad (7)$$

$$497 \quad g_{\theta}(\theta_{jt}) = W_0 + W_{0j} + s(x_t) + b_1\mu_{0j} + b_2\beta_{xj} + b_3\sigma_{0j} \dots$$

498 Extensive tutorials for incorporating such nonlinear effects into Bayesian regression
 499 models in Stan are freely available online (see <https://mc-stan.org/documentation/case-studies>
 500 for worked examples of fitting splines and Gaussian processes). Code from Stan models
 501 constructed using familiar R syntax in the brms package (Bürkner, 2019) also provides a helpful
 502 reference point for getting started. By allowing for arbitrarily complex average RN shapes across
 503 subjects, individual deviations from the average slope for phenotype as well as for fitness are
 504 much more likely to exhibit multivariate normality. This general approach allows researchers to
 505 accurately describe trait change across complex and dynamic environments, while still using
 506 standard theory from quantitative genetics to quantify selection gradients and predict short-term
 507 evolutionary responses.

508 **Additional individual effects**

509 The RN model presented in the main text (**Eq. 1**) does not account the fact that phenotypic
 510 dispersion σ may also be plastic across environments, a phenomenon broadly referred to as
 511 ‘malleability’ (see [O’Dea, Noble, & Nakagawa 2021](#) for discussion). Malleability in residuals can
 512 be estimated by including population- and individual-level slopes in the linear predictor of the
 513 dispersion parameter ([Westneat et al., 2013](#)). For example,

$$514 \quad g_{\sigma}(\sigma_{jt}) = \sigma_0 + \sigma_{0j} + (\rho + \rho_j)x_t \quad (8)$$

515 if observation-level variation in environmental measure \mathbf{x} is expected to have effect ρ on average
 516 differences in phenotypic residuals. Malleability can then be treated as a further RN parameter
 517 that is also potentially under selection. Some statistical distributions such as the Poisson lack an
 518 explicit dispersion parameter, due to deterministic mean-variance relationships, and thus at first
 519 glance only provide scope for selection on the RN intercepts and slopes of expected values.
 520 However, in many empirical datasets, there is more variance observed in the phenotype than
 521 predicted by these distributions (overdispersion), which can be accounted for through the
 522 inclusion of further random effects capturing stochastic, observation-level deviations from model
 523 expectations (i.e. residuals; [Harrison, 2014](#)). The dispersion of these observation-level random
 524 effects can then be modelled as a function of individual-level intercepts and slopes, similar to a
 525 standard Gaussian model, providing scope for estimating selection on phenotypic variability using
 526 a broad range of RN GLMMs.

527 More generally, any theoretically relevant component of a statistical distribution may be
 528 modelled as a function of further individual-level effects and conceptualized as a RN parameter
 529 regulating the expression of phenotypes within and across environments. Hurdle models, for
 530 example, combine multiple distributions together to distinguish effects on the presence/absence
 531 of trait expression from effects on the subsequent magnitude or intensity of trait expression
 532 ([Mullahy 1986](#); [Heilbron 1994](#)). This is particularly useful for phenotypes such as allogrooming
 533 behavior in primates, which can vary repeatably among individuals both in its probability of
 534 occurring as well as its intensity and duration once expressed ([Silk et al., 2017](#)). These processes
 535 are interdependent but may nonetheless be subject to distinct selection pressures (e.g. whom
 536 should be groomed and how much), which can be investigated by estimating separate RN
 537 intercepts and/or slopes on both model components.

538

Statistical inference

539 Bayesian estimation

540 The proposed models cannot currently be estimated using popular GLMM software
541 packages, due to the need for latent RN parameters to be simultaneously estimated with random
542 and fixed effects across different response models. Fortunately, the Stan statistical programming
543 language (Carpenter et al. 2017), which relies on cutting-edge and computationally efficient
544 Markov Chain Monte Carlo (MCMC) sampling algorithms, provides the flexibility needed for
545 estimating these novel GLMMs within a Bayesian framework. Researchers unfamiliar with the
546 general motivations of Bayesian inference are encouraged to see McElreath (2020) and Gelman
547 et al. (2020) for helpful tips on developing an effective workflow for data analysis. The brms
548 package (Bürkner, 2018) is also a very helpful bridge for writing complex (non)linear Bayesian
549 GLMMs in Stan using familiar R formula syntax. We provide guided tutorials (see [data](#)
550 [availability](#)) for various implementations of the models presented here in Stan.

551 Prior distributions need to be specified for all the population-level parameters in a
552 Bayesian model. While flat or highly diffuse priors are often recommended in the literature (e.g.
553 Ellison 2004; Villemereuil et al. 2016; Houslay and Wilson 2017), weakly informative or
554 regularizing priors, which place relatively low probability on extreme effect sizes, facilitate more
555 robust inferences with limited sample sizes and should generally be preferred over flat priors
556 (Gelman & Tuerlinckx 2000; Lemoine 2019; McElreath 2020). This does not necessarily require
557 strong a priori assumptions; general-purpose priors can be used to increase the generalizability
558 and robustness of parameter estimates, even in a state of relative ignorance about the true effect
559 size. See Lemoine (2019) for more detailed discussion and recommendations.

560 Model validation

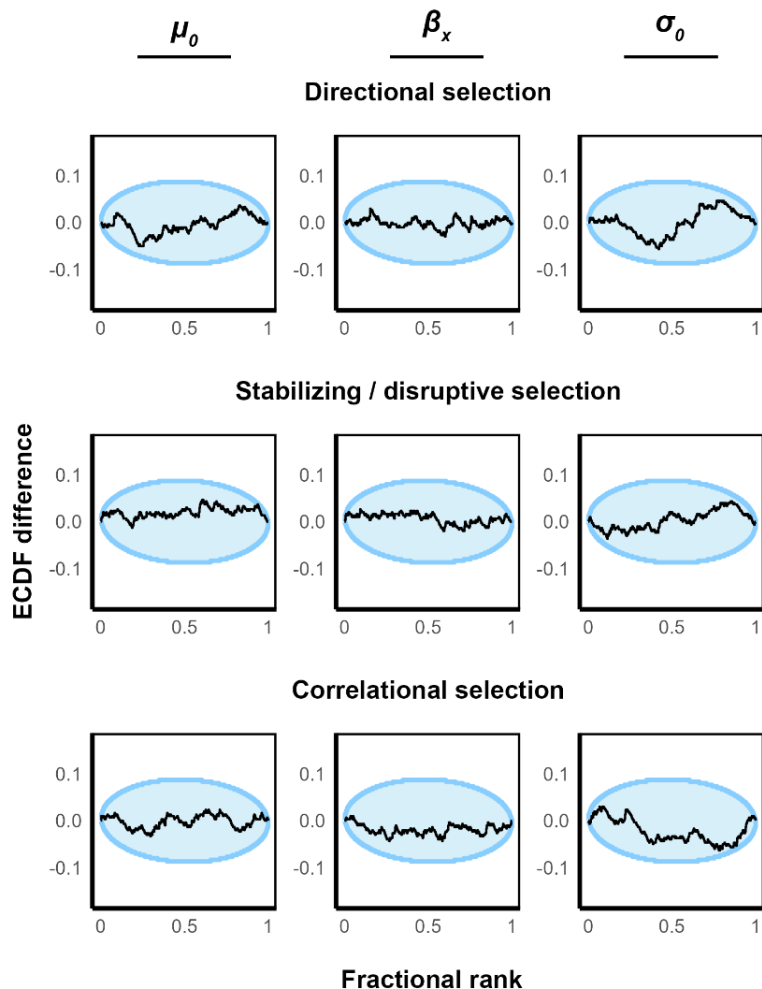
561 Previous work has validated the performance of our general approach in Stan for modest
562 effect sizes, showing robust estimates of directional selection on RN intercepts and slopes with
563 many repeated measures and sample sizes of $N = 100 - 300$ (Martin & Jaeggi, 2022). To provide
564 more general validation, we further conducted a simulation-based calibration (SBC; Talts et al.
565 2018; Säilynoja, Bürkner, & Vehtari, 2022) procedure to assess whether the proposed models are
566 unbiased estimators of nonlinear selection under a broader range of scenarios. SBC is a

567 procedure for validating the performance of any Bayesian algorithm across many possible
568 parameter values, as defined by the prior distributions of a generative model. This approach
569 removes the arbitrariness of setting a limited range of fixed parameter values for assessing
570 performance, which can lead to unexpected sources of bias being overlooked in uninvestigated
571 regions of parameter space (e.g. rare but possible combinations of phenotypic variances and
572 selection coefficients). Instead, random parameter values are repeatedly sampled across many
573 simulated datasets. Visual inspection of the correspondence between the generative distributions
574 used to simulate datasets and the subsequent posterior distributions inferred from these datasets
575 allows for detecting sources of bias such as overdispersion, overestimation, or inconsistent model
576 performance for extreme values. A GLMM validated through SBC is thus an unbiased Bayesian
577 estimator with respect to the range of effect sizes described by the prior generative model.

578 Particular attention was given to the estimation of directional and quadratic selection
579 coefficients during SBC, using 300 simulated datasets assuming conditions of very minimal
580 sampling effort (N = 100 subjects with 3 repeated phenotypic measurements and 2 repeated
581 fitness measures). Parameters were simulated such that
582 $\mu_0, \beta_x, \sigma_0, \mathbf{b}, \mathbf{q} \sim N(0,1)$, $\text{sd}([\mu_0, \beta_x, \sigma_0, \mathbf{W}_0])$, $\delta \sim \text{exponential}(2)$, and $\text{cor}([\mu_0, \beta_x, \sigma_0]) \sim \text{LKJ}(2)$. Note
583 that LKJ refers to the Lewandowski-Kurowicka-Joe distribution, which is useful for generating
584 positive-definite correlation matrices (Gelman et al., 2013). These priors led to a broad range of
585 very small to large selection effect sizes, as well as very small to large effects for the standard
586 deviations and correlations of RNs and the residual fitness standard deviation (δ). Phenotype and
587 fitness were assumed to be Gaussian for computational efficiency, with mean fitness fixed to 1.
588 Following the recommendations of Säilynoja et al. (2022), we computed and visualized the
589 difference in expected cumulative distribution functions between the generative and inferred
590 parameters to perform a quantitative graphical test of the model's performance. As shown in
591 **Figure 3**, our results demonstrated with probability ≥ 0.95 that the posterior distributions of
592 inferred selection coefficients were not systematically higher or lower than the prior distributions
593 used to generate expected selection coefficients. The proposed model thus provides unbiased
594 inference of nonlinear selection on RNs across a broad range of effect sizes, even under
595 conditions of minimal sampling effort.

596

597



598

599 **Figure 3. Simulation-based calibration of the nonlinear selection model.** Results are shown for
600 analyses of 300 simulated datasets ($N = 100$ subjects, 3 repeated phenotype measures and 2 repeated
601 fitness measures) generated from prior distributions defined over the parameters of a Gaussian nonlinear
602 selection model for RNs (Eq. 4). Plots show the difference between the expected cumulative density
603 functions (y-axis) for directional and quadratic selection gradients, based on their generative prior
604 distributions $N(0,1)$, and the estimated cumulative density functions based on inferred posterior
605 distributions. The x-axis indicates the ordered fractional ranks across posterior samples used for computing
606 these comparisons. Blue circles show 90% Bayesian credible intervals for regions of concordance between
607 the estimated and expected parameter distributions, and the black line reflects the observed difference
608 between the expected and inferred distribution (a perfectly horizontal line would thus indicate perfect
609 concordance with the simulated parameters in every dataset). Consistent deviations of the black line
610 beyond the blue region would provide evidence of systematic inferential bias during model estimation. Note
611 that due to stochasticity, fluctuations of the black line within the blue circle are expected at computationally
612 efficient sample sizes.

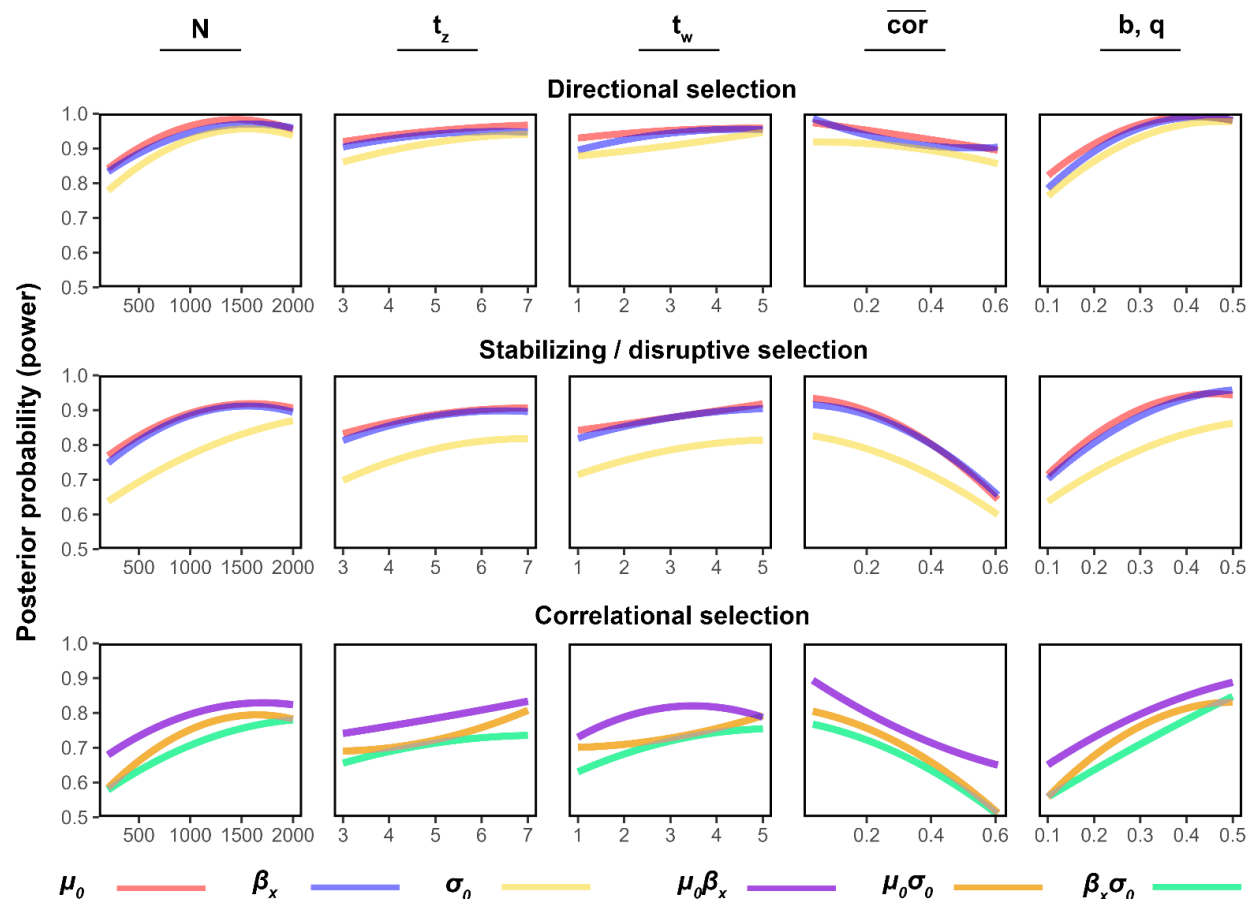
613 Power analysis

614 The SBC procedure demonstrated that our model facilitates unbiased Bayesian estimation
 615 across a broad range of parameter values (**Figure 3**). We also conducted an additional simulation
 616 study to provide concrete guidelines for empiricists designing studies to assess nonlinear
 617 selection on RNs, investigating how the power to detect the direction of selection gradients is
 618 influenced by the number of subjects and repeated measures per subject for phenotypes and
 619 fitness proxies. For simplicity and ease of effect size comparison, we modelled Gaussian
 620 phenotype and fitness measures. Fitness effects for the nonlinear selection model were simulated
 621 such that $\mathbf{b}, \mathbf{q} \sim U(0.1, 0.5)$, resulting in selection effects ranging from statistically weak to strong
 622 in strength, with a mean effect size of $|\mathbf{0.3}|$ across datasets. For simplicity, we assumed $W_0 = \delta =$
 623 1 and $\mu_0 = \beta_0 = 0$. Continuous environmental variation (\mathbf{x}) for quantifying reaction norm slopes
 624 was treated as a standardized variable drawn from $\mathbf{x} \sim N(0,1)$. Repeatable among-individual
 625 differences in RNs were fixed to $\text{sd}([\boldsymbol{\mu}_0, \boldsymbol{\beta}, \boldsymbol{\sigma}_0]) = 0.55$ with correlations drawn from
 626 $\text{cor}([\boldsymbol{\mu}_0, \boldsymbol{\beta}, \boldsymbol{\sigma}_0]) \sim LKJ(5)$, and the residual standard deviation of the phenotype was fixed to
 627 $\text{sqrt}(\exp(\sigma_0)) = 0.77$, so that repeatable and residual random effect variances were 0.3 and 0.6
 628 respectively. This resulted in each RN parameter exhibiting modest repeatability, $R = 0.2 =$
 629 $\frac{0.3}{3(0.3)+0.6}$ in the absence of phenotypic correlations. Unexplained selection was also fixed to
 630 $\text{sd}(\mathbf{W}_0) = 0.55$ for the fitness model.

631 Power to detect the appropriate direction of selection coefficients was explored with 1000
 632 datasets of varying size drawn from $N \sim U(200, 1000)$ subjects with $t_z \sim U(3, 7)$ repeated
 633 phenotype and $t_w \sim U(1, 5)$ repeated fitness measures per subject. Classical frequentist methods
 634 define power with respect to a binary decision rule based on the desired significance level of a
 635 null hypothesis test. In Bayesian analysis, ‘power’ is not precisely defined but may instead refer
 636 to the continuous level of support provided for a direct (rather than null) hypothesis test, such as
 637 the posterior probability of positive selection occurring on a trait. The power of a Bayesian analysis
 638 thus reflects how confident a model is likely to be in the existence and direction of a true selection
 639 effect, with $p = 0.5$ indicating no confidence (+ and – values are equally likely) and $p = 1.0$
 640 indicating complete confidence in the effect. We herein use ‘power’ in this sense to refer to the
 641 expected posterior probability supporting positive directional and quadratic selection effects on
 642 RN parameters.

643 Power for detecting selection across simulated scenarios is visualized in **Figure 4**, with
 644 second-order polynomial lines plotted across datasets to infer general patterns expected in

645 empirical research. As expected, we find that Bayesian power for inferring directional and
646 quadratic selection increases with a greater number of subjects (N) and repeated phenotype (t_z)
647 and fitness measures (t_w), as well as with greater selection effect sizes (\mathbf{b}, \mathbf{q}), while larger
648 absolute phenotypic correlations among RN parameters ($\overline{\text{cov}}$) reduce power, particularly for
649 detecting quadratic selection. Power to detect quadratic selection is lower than for directional
650 selection across small to moderate sample and effect sizes, with power for correlational selection
651 also being relatively lower than stabilizing/disruptive selection except under ideal conditions. This
652 implies that research particularly focused on detecting correlational selection of RNs will require
653 larger samples to attain confident inferences. Power is also consistently lower for detecting all
654 types of selection on RN residual parameters in comparison to RN intercepts and slopes,
655 indicating a need for greater sampling effort in selection studies on phenotypic variability. As with
656 any multivariate selection model, these results show that large sample sizes and sufficient
657 repeated measurements are crucial for robust hypothesis testing, particularly in the presence of
658 weak selection. As a rule of thumb, sample sizes of at least $N = 500-1000$ will be desirable to
659 appropriately reduce the risk of false negatives, particularly in the absence of many repeated
660 phenotype and/or fitness measures. The negative effect of RN parameter correlations on power
661 also shows that (non)linear selection will be much easier to detect when RN parameters vary
662 quasi-independently among individuals within a population.



663

664 **Figure 4. Bayesian power analysis of the nonlinear selection model.** Results are shown for
 665 directional hypothesis tests of selection effects across 1000 simulated datasets used to estimate
 666 the nonlinear selection model for RNs (Eq. 4) with Gaussian phenotype and fitness measures.
 667 Plots show the expected posterior probability ('power', y-axis) supporting selection effects as a
 668 function of variation in sampling conditions across simulated datasets (x-axis): the number of
 669 subjects/sample size (N), the number of phenotypic measures per subject (t_z), the number of
 670 fitness measures per subject (t_w), the mean absolute correlation among RN parameters ($\overline{\text{cor}}$), and
 671 the size of linear (b) and nonlinear (q) selection effects. General patterns were inferred using
 672 second-order polynomials across conditions, which are color-coded by RN parameter (red =
 673 intercepts, blue = slopes, yellow = residuals, purple = intercepts x slopes, orange = intercepts x
 674 residuals, and green = slopes x residuals).

Conclusion

675

676 Studying selection on highly labile traits is essential for explaining how and why organisms
677 adapt to environmental change. RN models are a crucial tool for characterizing such phenotypes,
678 but their application to selection analysis remains hampered by the limitations of current methods.
679 A major challenge is to avoid inferential bias caused by non-repeatable, stochastic effects and
680 other sources of measurement error in RNs and their fitness effects (Hadfield et al. 2010; **FIGURE**
681 **1-2**). A common solution is to use multi-response/multivariate random effect GLMMs to account
682 for uncertainty in selection on RNs. However, this approach restricts analyses to focus on linear
683 effects and directional selection. Ignoring quadratic selection caused by nonlinear effects
684 fundamentally inhibits researchers' capacity to study the adaptive landscape of labile traits
685 (Bulmer 1971; Arnold et al., 2001; Blows & Brooks, 2003).

686 To overcome this limitation, we proposed a novel Bayesian GLMM framework for studying
687 complex patterns of nonlinear selection on RNs, which we validated over a broad range of
688 possible parameter values using a simulation-based calibration approach (**Figure 3**). We also
689 found that these models exhibited desirable statistical power under reasonable sampling
690 conditions for many long-term field research projects (**Figure 4**). This modeling framework
691 synthesizes the well-established Lande and Arnold (1983) approach to error-free selection
692 analysis with measurement error or error-in-variables models (Ponzi et al. 2018; Dingemans et
693 al. 2021; Martin & Jaeggi 2022) and double hierarchical (Westneat et al. 2013; O'Dea et al. 2021),
694 multi-response GLMMs (Brommer et al. 2012; Houslay & Wilson 2017; Arnold et al. 2019). These
695 models can be applied to estimate directional and quadratic selection irrespective of the
696 distribution of the data and the potential nonlinearity of the RN or fitness function, allowing
697 researchers to construct more realistic models of the processes underlying their measurements.
698 This focuses attention on accurate description of observed data rather than the restrictive
699 assumptions of linear regression. With the analytic toolkit of quantitative genetics (Lande & Arnold
700 1983; Morrissey & Sakrejda 2013), estimates from these models can also be transformed to
701 quantify selection gradients, visualize multivariate selection, and predict ongoing adaptation. The
702 proposed modeling framework should, therefore, readily enhance tests of adaptive theory for
703 labile traits in the wild.

704

Data availability statement

705 R and Stan code with detailed tutorials for implementing the models presented in this paper are
706 available online through a GitHub public repository [HTTPS://GITHUB.COM/JORDAN-SCOTT-](https://github.com/Jordan-Scott-Martin/selection-on-rns)
707 [MARTIN/SELECTION-ON-RNS](https://github.com/Jordan-Scott-Martin/selection-on-rns).

708

Acknowledgements

709 JSM would like to thank Adrian Jaeggi, Adam Hunt, Camila Scaff, and Gabriel Šaffa for their
710 helpful feedback on previous versions of this manuscript, as well as the University of Zurich
711 Candoc/Forschungskredit PhD grant FK-20-034 and Statistical Quantification of Individual
712 Differences (SQuID) educational group postdoctoral fellowship for financial support.

References

- Abdul-Rahman, F., Tranchina, D., & Gresham, D. (2021). Fluctuating environments maintain genetic diversity through neutral fitness effects and balancing selection. *Molecular Biology and Evolution*, *38*(10), 4362–4375.
- Araya-Ajoy, Y. G., & Dingemanse, N. J. (2014). Characterizing behavioural ‘characters’: An evolutionary framework. *Proceedings of the Royal Society B*, *281*, 20132645.
- Araya-Ajoy, Y. G., Dingemanse, N. J., Westneat, D. F., & Wright, J. (2023). The evolutionary ecology of variation in labile traits: selection on its among- and within-individual components. *Evolution*, *77*, 2246–2256.
- Araya-Ajoy, Y. G., Mathot, K. J., & Dingemanse, N. J. (2015). An approach to estimate short-term, long-term and reaction norm repeatability. *Methods in Ecology and Evolution*, *6*(12), 1462–1473.
- Arnold, S. J., Pfrender, M. E., & Jones, A. G. (2001). The adaptive landscape as a conceptual bridge between micro- and macroevolution. *Genetica*, *112–113*, 9–32.
- Arnold, P. A., Nicotra, A. B., & Kruuk, L. E. (2019). Sparse evidence for selection on phenotypic plasticity in response to temperature. *Philosophical Transactions of the Royal Society B*, *374*, 20180185.
- Barrett, L. (2011). *Beyond the brain: How body and environment shape animal and human minds*. Princeton University Press.
- Bell, A. M., Hankison, S. J., & Laskowski, K. L. (2009). The repeatability of behaviour: A meta-analysis. *Animal Behaviour*, *77*, 771–783.
- Bijma, P. (2011). A general definition of the heritable variation that determines the potential of a population to respond to selection. *Genetics*, *189*(4), 1347–1359.
- Bijma, P., & Wade, M. J. (2008). The joint effects of kin, multilevel selection and indirect genetic effects on response to genetic selection. *Journal of Evolutionary Biology*, *21*(5), 1175–1188.
- Biro, P. A., & Adriaenssens, B. (2013). Predictability as a personality trait: Consistent differences in intraindividual behavioral variation. *The American Naturalist*, *182*, 621–629.

- Blows, M. W. (2007). A tale of two matrices: Multivariate approaches in evolutionary biology. *Journal of Evolutionary Biology*, *20*, 1–8.
- Blows, M. W., & Brooks, R. (2003). Measuring nonlinear selection. *The American Naturalist*, *162*, 815–820.
- Bolker, B. M., Brooks, M. E., Clark, C. J., Geange, S. W., Poulsen, J. R., Stevens, M. H. H., & White, J.-S. S. (2009). Generalized linear mixed models: a practical guide for ecology and evolution. *Trends in Ecology & Evolution*, *24*(3), 127–135.
- Bollen, K. A., & Noble, M. D. (2011). Structural equation models and the quantification of behavior. *Proceedings of the National Academy of Sciences*, *108*, 15639–15646.
- Bollen, T., Morrissey, M. B., & Kruuk, L. E. (2019). Estimation of genetic variance in fitness, and inference of adaptation, when fitness follows a log-normal distribution. *Journal of Heredity*, *110*, 383–395.
- Bolnick, D. I., Svanbäck, R., Fordyce, J. A., Yang, L. H., Davis, J. M., Hulsey, C. D., & Forister, M. L. (2003). The ecology of individuals: incidence and implications of individual specialization. *The American Naturalist*, *161*(1), 1–28.
- Oomen, R. A., & Hutchings, J. A. (2022). Genomic reaction norms inform predictions of plastic and adaptive responses to climate change. *The Journal of Animal Ecology*, *91*, 1073–1087.
- Borenstein, E., Feldman, M. W., & Aoki, K. (2008). Evolution of learning in fluctuating environments: When selection favors both social and exploratory individual learning. *Evolution*, *62*, 586–602.
- Brommer, J. E. (2013). On between-individual and residual (co) variances in the study of animal personality: Are you willing to take the 'individual gambit'? *Behavioral Ecology and Sociobiology*, *67*, 1027–1032.
- Brommer, J. E., Kontiainen, P., & Pietiäinen, H. (2012). Selection on plasticity of seasonal life-history traits using random regression mixed model analysis. *Ecology and Evolution*, *24*, 695–704.
- Bulmer, M. G. (1971). The effect of selection on genetic variability. *The American Naturalist*, *201*-211.

- Bürkner, P. (2018). Advanced Bayesian multilevel modeling with the R package brms. *The R Journal*, 10, 395–411.
- Caño, L., Escarré, J., Fleck, I., Blanco-Moreno, J. M., & Sans, F. X. (2008). Increased fitness and plasticity of an invasive species in its introduced range: a study using *Senecio pterophorus*. *The Journal of Ecology*, 96(3), 468–476.
- Carpenter, B., Gelman, A., Hoffman, M. D., Lee, D., Goodrich, B., Betancourt, M., & A. Riddell... (2017). Stan: A probabilistic programming language. *Journal of Statistical Software*, 74.
- Catalina, A., Bürkner, P. C., & Vehtari, A. (2020). Projection predictive inference for generalized linear and additive multilevel models. *arXiv*. <http://arxiv.org/abs/2010.06994>
- Cauchoix, M., Chow, P. K. Y., Horik, J. O. V., Atance, C. M., Barbeau, E. J.,...G. B.-J., & Cauchard, L. (2018). The repeatability of cognitive performance: A meta-analysis. *Philosophical Transactions of the Royal Society B*, 373, 20170281.
- Dall, S. R. X., & Griffith, S. C. (2014). An empiricist guide to animal personality variation in ecology and evolution. *Frontiers in Ecology and Evolution*, 14, 3.
- Darwin, C. (1859). *On the origin of species by means of natural selection*. London: John Murray.
- Denissen, J. J. A., & Penke, L. (2008). Motivational individual reaction norms underlying the five-factor model of personality: First steps toward a theory-based conceptual framework. *Journal of Research in Personality*, 69, 1285–1302.
- Dingemanse, N. J., & Araya-Ajoy, Y. G. (2015). Interacting personalities: behavioural ecology meets quantitative genetics. *Trends in Ecology & Evolution*, 30(2), 88–97.
- Dingemanse, N. J., Araya-Ajoy, Y. G., & Westneat, D. F. (2021). Most published selection gradients are underestimated: Why this is and how to fix it. *Evolution, Early View*.
- Dingemanse, N. J., & Dochtermann, N. A. (2013). Quantifying individual variation in behaviour: Mixed-effect modelling approaches. *Journal of Animal Ecology*, 82, 39–54.
- Dingemanse, N. J., Kazem, A. J., Réale, D., & Wright, J. (2010). Behavioural reaction norms: Animal personality meets individual plasticity. *Trends in Ecology and Evolution*, 25, 81–89.
- Dingemanse, N. J., & Wolf, M. (2010). Recent models for adaptive personality differences: a review. *Philosophical Transactions of the Royal Society B: Biological Sciences*, 365, 3947–3958.

- Eisenegger, C., Haushofer, J., & Fehr, E. (2011). The role of testosterone in social interaction. *Trends in Cognitive Sciences*, 15, 263–271.
- Ellison, A. M. (2004). Bayesian inference in ecology. *Ecology Letters*, 7, 509–520.
- Estes, S., & Arnold, S. J. (2007). Resolving the paradox of stasis: models with stabilizing selection explain evolutionary divergence on all timescales. *The American Naturalist*, 169, 227–244.
- Fanson, K. V., & Biro, P. A. (2015). Meta-analytic insights into factors influencing the repeatability of hormone levels in agricultural, ecological, and medical fields. *American Journal of Physiology-Regulatory, Integrative and Comparative Physiology*, 316, R101–R109.
- Fawcett, T. W., Hamblin, S., & Giraldeau, L.-A. (2013). Exposing the behavioral gambit: the evolution of learning and decision rules. *Behavioral Ecology*, 24, 2–11.
- Fay, R., Martin, J., & Plard, F. (2022). Distinguishing within from between individual effects: How to use the within-individual centering method for quadratic pattern. *Journal of Animal Ecology*, 91, 8–19.
- Flatt, T. (2005). The evolutionary genetics of canalization. *The Quarterly Review of Biology*, 80, 287–316.
- Fox, R. J., Donelson, J. M., Schunter, C., Ravasi, T., & Gaitán-Espitia, J. D. (2019). Beyond buying time: the role of plasticity in phenotypic adaptation to rapid environmental change. *Philosophical Transactions of the Royal Society B*, 374, 20180174.
- Futuyma, D. (2021). How does phenotypic plasticity fit into evolutionary theory? In D. W. Pfennig (Ed.), *Phenotypic plasticity & evolution: Causes, consequences, controversies* (pp. 349–366).
- Gavrilets, S., & Hastings, A. (1994). A quantitative-genetic model for selection on developmental noise. *Evolution*, 48, 1478–1486.
- Gavrilets, S., & Scheiner, S. M. (1993). The genetics of phenotypic plasticity. VI. Theoretical predictions for directional selection. *Journal of Evolutionary Biology*, 6, 49–68.
- Gelman, A., & Tuerlinckx, F. (2000). Type s error rates for classical and bayesian single and multiple comparison procedures. *Computational Statistics*, 15, 373–390.

- Gelman, A., Carlin, J. B., Stern, H. S., Dunson, D. B., Vehtari, A., & Rubin, D. B. (2013). *Bayesian Data Analysis* (3rd ed.). Chapman and Hall/CRC.
- Gelman, A., Vehtari, A., Simpson, D., Margossian, C. C., Carpenter, B., Yao, Y., & M. Modrák... (2020). Bayesian workflow. *arXiv Preprint*, *arXiv:2011.01808*.
- Ghalambor, C. K., McKay, J. K., Carroll, S. P., & Reznick, D. N. (2007). Adaptive versus non-adaptive phenotypic plasticity and the potential for contemporary adaptation in new environments. *Functional Ecology*, *21*(3), 394–407.
- Ghalambor, C. K., Hoke, K. L., Ruell, E. W., Fischer, E. K., Reznick, D. N., & Hughes, K. A. (2015). Non-adaptive plasticity potentiates rapid adaptive evolution of gene expression in nature. *Nature*, *525*(7569), 372–375.
- Gomulkiewicz, R., Kingsolver, J. G., Carter, P. A., & Heckman, N. (2018). Variation and evolution of function-valued traits. *Annual Review of Ecology, Evolution, and Systematics*, *49*, 139–164.
- Gomulkiewicz, R., & Kirkpatrick, M. (1992). Quantitative genetics and the evolution of reaction norms. *Evolution*, *46*, 390–411.
- Guindre-Parker, S., & Rubenstein, D. R. (2018). Multiple benefits of alloparental care in a fluctuating environment. *Royal Society Open Science*, *5*, 172406.
- Hadfield, J. D., Wilson, A. J., Garant, D., & Sheldon, B. C. (2010). The misuse of BLUP in ecology and evolution. *The American Naturalist*, *175*, 116–125.
- Hansen, T. F., Carter, A. J. R., & Pélabon, C. (2006). On adaptive accuracy and precision in natural populations. *The American Naturalist*, *168*(2), 168–181.
- Harrison, X. A. (2014). Using observation-level random effects to model overdispersion in count data in ecology and evolution. *PeerJ*, *2*, e616.
- Heilbron, D. C. (1994). Zero-altered and other regression models for count data with added zeros. *Biometrical Journal*, *36*, 531–547.
- Hendry, A. P. (2016). Key questions on the role of phenotypic plasticity in eco-evolutionary dynamics. *The Journal of Heredity*, *107*, 25–41.
- Henrich, J., & McElreath, R. (2003). The evolution of cultural evolution. *Evolutionary Anthropology*, *12*(3), 123–135.

- Houslay, T. M., & Wilson, A. J. (2017). Avoiding the misuse of BLUP in behavioural ecology. *Behavioral Ecology*, 28, 948–952.
- Houston, A. I., & McNamara, J. M. (1999) *Models of adaptive behaviour*. Cambridge, MA: Cambridge University Press.
- Hugie, D. M. (2003). The waiting game: a “battle of waits” between predator and prey. *Behavioral Ecology*, 14(6), 807–817.
- Jaeggi, A. V., Boose, K. J., White, F. J., & Gurven, M. (2016). Obstacles and catalysts of cooperation in humans, bonobos, and chimpanzees: Behavioural reaction norms can help explain variation in sex roles, inequality, war and peace. *Behaviour*, 153, 1015–1052.
- de Jong, G. (1995). Phenotypic plasticity as a product of selection in a variable environment. *The American Naturalist*, 145, 493–512.
- Kazancıoğlu, E., Klug, H., & Alonzo, S. H. (2012). The evolution of social interactions changes predictions about interacting phenotypes. *Evolution*, 66, 2056–2064.
- King, J. G., & Hadfield, J. D. (2019). The evolution of phenotypic plasticity when environments fluctuate in time and space. *Evolution Letters*, 3, 15–27.
- Kinsler, G., Schmidlin, K., Newell, D., Eder, R., Apodaca, S., Lam, G., Petrov, D., & Geiler-Samerotte, K. (2023). Extreme sensitivity of fitness to environmental conditions: Lessons from #1BigBatch. *Journal of Molecular Evolution*, 91, 293–310.
- Lande, R., & Arnold, S. J. (1983). The measurement of selection on correlated characters. *Evolution*, 37, 1210–1226.
- Lemoine, N. P. (2019). Moving beyond noninformative priors: Why and how to choose weakly informative priors in bayesian analyses. *Oikos*, 128.
- Loken, E., & Gelman, A. (2017). Measurement error and the replication crisis. *Science*, 355, 584–585.
- Martin, J. G., Nussey, D. H., Wilson, A. J., & Réale, D. (2011). Measuring individual differences in reaction norms in field and experimental studies: a power analysis of random regression models. *Methods in Ecology and Evolution*, 2(4), 362-374.
- Martin, J. S., & Jaeggi, A. V. (2022). Social animal models for quantifying plasticity, assortment, and selection on interacting phenotypes. *Journal of Evolutionary Biology*, 35, 520-538.

- Martin, J. S., Jaeggi, A. V., & Koski, S. E. (2023). Social evolution of individual differences: Future directions for a comparative science of personality in social behavior. *Neuroscience & BioBehavioral Reviews*, *144*, 104980.
- Martin, J. S., Massen, J. J., Šlipogor, V., Bugnyar, T., Jaeggi, A. V., & Koski, S. E. (2019). The EGA+ GNM framework: An integrative approach to modelling behavioural syndromes. *Methods in Ecology and Evolution*, *10*, 245–257.
- Martin, J. S., Ringen, E. J., Duda, P., & Jaeggi, A. V. (2020). Harsh environments promote alloparental care across human societies. *Proceedings of the Royal Society B*, *287*, 20200758.
- Martin, J. S., Westneat, D., Wilson, A. J., Dingemanse, N. J., & Araya-Ajoy, Y. (2024). Frequency-dependence favors social plasticity and facilitates socio-eco-evolutionary feedback in fluctuating environments. *EcoEvoRxiv*. Doi: [10.32942/X2KW5R](https://doi.org/10.32942/X2KW5R)
- de la Mata, R., Zas, R., Bustingorri, G., Sampedro, L., Rust, M., Hernandez-Serrano, A., & Sala, A. (2022). Drivers of population differentiation in phenotypic plasticity in a temperate conifer: A 27-year study. *Evolutionary Applications*, *15*, 1945–1962.
- Mathuru, A. S., Kibat, C., Cheong, W. F., Shui, G., Wenk, M. R., Friedrich, R. W., & Jesuthasan, S. (2012). Chondroitin fragments are odorants that trigger fear behavior in fish. *Current Biology*, *22*, 538–554.
- McElreath, R. (2020). *Statistical rethinking: A Bayesian course with examples in r and stan* (2nd ed.). CRC Press.
- McGlothlin, J. W., Moore, A. J., Wolf, J. B., & Brodie, E. D., 3rd. (2010). Interacting phenotypes and the evolutionary process. III. Social evolution. *Evolution*, *64*(9), 2558–2574.
- McNamara, J. M., & Houston, A. I. (2009). Integrating function and mechanism. *Trends in Ecology & Evolution*, *24*(12), 670–675.
- McNamara, J. M., & Leimar, O. (2010). Variation and the response to variation as a basis for successful cooperation. *Philosophical Transactions of the Royal Society of London. Series B, Biological Sciences*, *365*, 2627–2633.
- McNamara, J. M., & Leimar, O. (2020). *Game theory in biology*. Oxford University Press.

- Mitchell, D. J., Beckmann, C., & Biro, P. A. (2021). Understanding the unexplained: The magnitude and correlates of individual differences in residual variance. *Ecology and Evolution*, *11*, 7201–7210.
- Mitchell, D. J., Dujon, A. M., Beckmann, C., & Biro, P. A. (2020). Temporal autocorrelation: a neglected factor in the study of behavioral repeatability and plasticity. *Behavioral Ecology*, *31*, 222–231.
- Moore, A. J., Brodie, E. D., 3rd, & Wolf, J. B. (1997). Interacting phenotypes and the evolutionary process: I. Direct and indirect genetic effects of social interactions. *Evolution*, *51*(5), 1352–1362.
- Moore, T. Y., Cooper, K. L., Biewener, A. A., & Vasudevan, R. (2017). Unpredictability of escape trajectory explains predator evasion ability and microhabitat preference of desert rodents. *Nature Communications*, *8*, 1–9.
- Morrissey, M. B., Parker, D. J., Korsten, P., Pemberton, J. M., Kruuk, L. E. B., & Wilson, A. J. (2012). The prediction of adaptive evolution: Empirical application of the secondary theorem of selection and comparison to the breeder's equation. *Evolution*, *66*, 2399–2410.
- Morrissey, M. B., & Sakrejda, K. (2013). Unification of regression-based methods for the analysis of natural selection. *Evolution*, *67*, 2094–2100.
- Mouchet, A., Cole, E. F., Matthysen, E., Nicolaus, M., Quinn, J. L., Roth, A. M., Tinbergen, J. M., van Oers, K., van Overveld, T., & Dingemanse, N. J. (2021). Heterogeneous selection on exploration behavior within and among West European populations of a passerine bird. *Proceedings of the National Academy of Sciences*, *118*.
- Mullahy, J. (1986). Specification and testing of some modified count data models. *Journal of Econometrics*, *33*, 341–365.
- Munar-Delgado, G., Araya-Ajoy, Y. G., & Edelaar, P. (2023). Estimation of additive genetic variance when there are gene–environment correlations: Pitfalls, solutions and unexplored questions. *Methods in Ecology and Evolution*, *14*, 1245–1258.
- Nakahashi, W., & Ohtsuki, H. (2015). When is emotional contagion adaptive? *Journal of Theoretical Biology*, *380*, 480–488.
- Nettle, D., & Penke, L. (2010). Personality: Bridging the literatures from human psychology and behavioural ecology. *Philosophical Transactions of the Royal Society B*, *365*, 4043–4050.

- Newediuk, L., Prokopenko, C. M., & Wal, E. V. (2022). Individual differences in habitat selection mediate landscape level predictions of a functional response. *Oecologia*, 1–12.
- Nicoglou, A. (2015). The evolution of phenotypic plasticity: genealogy of a debate in genetics. *Studies in History and Philosophy of Biological and Biomedical Sciences*, 50, 67–76.
- Niemelä, P. T., & Dingemanse, N. J. (2018). Meta-analysis reveals weak associations between intrinsic state and personality. *Proceedings of the Royal Society B*, 285, 20172823.
- Niv, Y., Joel, D., Meilijson, I., & Ruppin, E. (2002). Evolution of reinforcement learning in uncertain environments: A simple explanation for complex foraging behaviors. *Adaptive Behavior*, 10, 5–24.
- Nussey, D. H., Wilson, A. J., & Brommer, J. E. (2007). The evolutionary ecology of individual phenotypic plasticity in wild populations. *Journal of Evolutionary Biology*, 20, 831–844.
- O’Dea, R. E., Noble, D. W., & Nakagawa, S. (2021). Unifying individual differences in personality, predictability and plasticity: A practical guide. *Methods in Ecology and Evolution*, 13, 278-293.
- Ørsted, M., Rohde, P. D., Hoffmann, A. A., Sørensen, P., & Kristensen, T. N. (2018). Environmental variation partitioned into separate heritable components. *Evolution*, 72, 136–152.
- Phillips, P. C., & Arnold, S. J. (1989). Visualizing multivariate selection. *Evolution*, 43, 1209–1222.
- Pick, J. L., Lemon, H. E., Thomson, C. E., & Hadfield, J. D. (2022). Decomposing phenotypic skew and its effects on the predicted response to strong selection. *Nature Ecology & Evolution*, 6, 774–785.
- Pol, M. van de, & Wright, J. (2009). A simple method for distinguishing within- versus between-subject effects using mixed models. *Animal Behaviour*, 77, 753–758.
- Ponzi, E., Keller, L. F., Bonnett, T., & Muff, S. (2018). Heritability, selection, and the response to selection in the presence of phenotypic measurement error: Effects, cures, and the role of repeated measurements. *Evolution*, 72, 1992–2004.

- Prentice, P. M., Houslay, T. M., Martin, J. G. A., & Wilson, A. J. (2020). Genetic variance for behavioural “predictability” of stress response. *Journal of Evolutionary Biology*, *33*(5), 642–652.
- Projecto-Garcia, J., Biddle, J. F., & Ragsdale, E. J. (2017). Decoding the architecture and origins of mechanisms for developmental polyphenism. *Current Opinion in Genetics & Development*, *47*, 1-8.
- Ramakers, J. J. C., Visser, M. E., & Gienapp, P. (2020). Quantifying individual variation in reaction norms: Mind the residual. *Journal of Evolutionary Biology*, *33*, 352–366.
- Reynolds, R. J., de Los Campos, G., Egan, S. P., & Ott, J. R. (2016). Modelling heterogeneity among fitness functions using random regression. *Methods in Ecology and Evolution*, *7*, 70–79.
- Royauté, R., Berdal, M. A., Garrison, C. R., & Dochtermann, N. A. (2018). A meta-analysis of the pace-of-life syndrome hypothesis. *Behavioral Ecology and Sociobiology*, *72*, 1–10.
- Royauté, R., Hedrick, A., & Dochtermann, N. A. (2020). Behavioural syndromes shape evolutionary trajectories via conserved genetic architecture. *Proceedings of the Royal Society B*, *287*, 20200183.
- Säilynoja, T., Bürkner, P.-C., & Vehtari, A. (2022). Graphical test for discrete uniformity and its applications in goodness-of-fit evaluation and multiple sample comparison. *Statistics and Computing*, *32*, 32.
- Sasaki, A., & Ellner, S. (1997). Quantitative genetic variance maintained by fluctuating selection with overlapping generations: Variance components and covariances. *Evolution*, *51*, 682–696.
- Schaum, C.-E., Buckling, A., Smirnoff, N., & Yvon-Durocher, G. (2022). Evolution of thermal tolerance and phenotypic plasticity under rapid and slow temperature fluctuations. *Proceedings. Biological Sciences / The Royal Society*, *289*(1980), 20220834.
- Scheiner, S. M., & Lyman, R. F. (1991). The genetics of phenotypic plasticity. II. Response to selection. *Journal of Evolutionary Biology*, *4*, 23-50.
- Scheiner, S. M. (1993a). Plasticity as a selectable trait: Reply to via. *The American Naturalist*, *142*, 371–373.

- Scheiner, S. M. (1993b). Genetics and evolution of phenotypic plasticity. *Annual Review of Ecology and Systematics*, 24, 35–68.
- Scheiner, S. M., Donohue, K., Mazer, L. A. D. S. J., & Wolfe, L. M. (2002). Reducing environmental bias when measuring natural selection. *Evolution*, 56, 2156–2167.
- Schlichting, C. D., & Pigliucci, M. (1998). *Phenotypic evolution: a reaction norm perspective*. Sinauer Associates: Sunderland, MA.
- Schluter, D., & Nychka, D. (1994). Exploring fitness surfaces. *The American Naturalist*, 143, 597–616.
- Searle, S. R. (1961). Phenotypic, genetic and environmental correlations. *Biometrics*, 17, 474–480.
- Sigourney, D. B., Munch, S. B., & Letcher, B. H. (2012). Combining a Bayesian nonparametric method with a hierarchical framework to estimate individual and temporal variation in growth. *Ecological Modelling*, 247, 125–134.
- Siegal, M. L., & Leu, J. Y. (2014). On the nature and evolutionary impact of phenotypic robustness mechanisms. *Annual Review of Ecology, Evolution and Systematics*, 45, 495–517.
- Sih, A., Mathot, K. J., Moirón, M., Montiglio, P. O., Wolf, M., & Dingemanse, N. J. (2015). Animal personality and state–behaviour feedbacks: A review and guide for empiricists. *Trends in Ecology and Evolution*, 30, 50–60.
- Silk, J. B., Roberts, E. R., Barrett, B. J., Patterson, S. K., & Strum, S. C. (2017). Female–male relationships influence the form of female–female relationships in olive baboons, *Papio anubis*. *Animal Behaviour*, 131, 89–98.
- Skinner, B. F. (1966). The phylogeny and ontogeny of behavior. *Science*, 153, 1205–1213.
- Spearman, C. (1904). The proof and measurement of association between two things. *The American Journal of Psychology*, 15, 72–101.
- Stamps, J. A. (2016). Individual differences in behavioural plasticities. *Biological Reviews*, 91, 534–567.

- Stinchcombe, J. R., Agrawal, A. F., Hohenlohe, P. A., Arnold, S. J., & Blows, M. W. (2008). Estimating nonlinear selection gradients using quadratic regression coefficients: Double or nothing? *Evolution*, *68*.
- Stinchcombe, J. R., Rutter, M. T., Burdick, D. S., Tiffin, P., Rausher, M. D., & Mauricio, R. (2002). Testing for environmentally induced bias in phenotypic estimates of natural selection: theory and practice. *The American Naturalist*, *160*, 511–523.
- Stinchcombe, J. R., Simonsen, A. K., & Blows, M. W. (2014). Estimating uncertainty in multivariate responses to selection. *Evolution*, *68*.
- Strickland, K., Mitchell, D. J., Delmé, C., & Frère, C. H. (2021). Repeatability and heritability of social reaction norms in a wild agamid lizard. *Evolution*, *75*.
- Suzuki, Y., & Nijhout, H. F. (2006). Evolution of a polyphenism by genetic accommodation. *Science*, *311*, 650–652.
- Svensson, E. I., Gomez-Llano, M., & Waller, J. T. (2020). Selection on phenotypic plasticity favors thermal canalization. *Proceedings of the National Academy of Sciences USA*, *117*.
- Talts, S., Betancourt, M., Simpson, D., Vehtari, A., & Gelman, A. (2018). Validating Bayesian inference algorithms with simulation-based calibration. *arXiv*.
<http://arxiv.org/abs/1804.06788>
- Tonsor, S. J., Elnaccash, T. W., & Scheiner, S. M. (2013). Developmental instability is genetically correlated with phenotypic plasticity, constraining heritability, and fitness. *Evolution*, *67*, 2923–2935.
- Vasey, G. L., Weisberg, P. J., & Urza, A. K. (2022). Intraspecific trait variation in a dryland tree species corresponds to regional climate gradients. *Journal of Biogeography*, *49*, 2309–2320.
- Vercken, E., Wellenreuther, M., Svensson, E. I., & Mauroy, B. (2012). Don't fall off the adaptation cliff: when asymmetrical fitness selects for suboptimal traits. *PloS One*, *7*, e34889.
- Via, S. (1993). Adaptive phenotypic plasticity: target or by-product of selection in a variable environment? *The American Naturalist*, *142*, 352–365.

- Via, S., Gomulkiewicz, R., Jong, G. D., Scheiner, S. M., Schlichting, C. D., & Tienderen, P. H. V. (1995). Adaptive phenotypic plasticity: Consensus and controversy. *Trends in Ecology and Evolution*, *10*, 212–217.
- Via, S., & Lande, R. (1985). Genotype-environment interactions and the evolution of phenotypic plasticity. *Evolution*, *39*, 505–522.
- de Villemereuil, P., Charmantier, A., Arlt, D., Bize, P., Brekke, P., Brouwer, L.,..., & Chevin, L. M. (2020). Fluctuating optimum and temporally variable selection on breeding date in birds and mammals. *Proceedings of the National Academy of Sciences*, *117*, 31969–31978.
- Villemereuil, P. de, Schielzeth, H., Nakagawa, S., & Morrissey, M. (2016). General methods for evolutionary quantitative genetic inference from generalized mixed models. *Genetics*, *204*, 1281–1294.
- Volis, S., Ormanbekova, D., & Yermekbayev, K. (2015). Role of phenotypic plasticity and population differentiation in adaptation to novel environmental conditions. *Ecology and Evolution*, *5*(17), 3818–3829.
- Wagner, G. P., Booth, G., & Bagheri-Chaichian, H. (1997). A population genetic theory of canalization. *Evolution*, *51*, 329–347.
- Wang, S. P., & Althoff, D. M. (2019). Phenotypic plasticity facilitates initial colonization of a novel environment. *Evolution*, *73*(2), 303–316.
- Weis, A. E., & Gorman, W. L. (1990). Measuring selection on reaction norms: An exploration of the eurosta-solidago system. *Evolution*, *44*, 820–831.
- Westneat, D. F., Araya-Ajoy, Y. G., Allegue, H., Class, B., Dingemans, N., Dochtermann, N. A., ... & Schielzeth, H. (2020). Collision between biological process and statistical analysis revealed by mean centring. *Journal of Animal Ecology*, *89*(12), 2813–2824.
- Westneat, D. F., Schofield, M., & Wright, J. (2013). Parental behavior exhibits among-individual variance, plasticity, and heterogeneous residual variance. *Behavioral Ecology*, *24*, 598–604.
- Westneat, D. F., Wright, J., & Dingemans, N. J. (2015). The biology hidden inside residual within-individual phenotypic variation. *Biological Reviews*, *90*, 729–743.

- Wingfield, J. C., Hegner, R. E., Jr., A. M. D., & Ball, G. F. (1990). The 'challenge hypothesis': Theoretical implications for patterns of testosterone secretion, mating systems, and breeding strategies. *The American Naturalist*, *136*, 829–846.
- Wolf, M., & Weissing, F. J. (2010). An explanatory framework for adaptive personality differences. *Philosophical Transactions of the Royal Society B*, *365*, 3959–3968.
- Wright, J., Bolstad, G. H., Araya-Ajoy, Y. G., & Dingemanse, N. J. (2019). Life-history evolution under fluctuating density-dependent selection and the adaptive alignment of pace-of-life syndromes. *Biological Reviews*, *94*, 230–247.
- Wright, J., Haaland, T. R., Dingemanse, N. J., & Westneat, D. F. (2022). A reaction norm framework for the evolution of learning: how cumulative experience shapes phenotypic plasticity. *Biological Reviews*. <https://doi.org/10.1111/brv.12879>
- Yamahira, K., Kawajiri, M., Takeshi, K., & Irie, T. (2007). Inter- and intrapopulation variation in thermal reaction norms for growth rate: evolution of latitudinal compensation in ectotherms with a genetic constraint. *Evolution*, *61*, 1577–1589.