



# Heritability, selection, and the response to selection in the presence of phenotypic measurement error: Effects, cures, and the role of repeated measurements

Erica Ponzi,<sup>1,2</sup> Lukas F. Keller,<sup>1,3</sup> Timothée Bonnet,<sup>1,4</sup> and Stefanie Muff<sup>1,2,5</sup>

<sup>1</sup>Department of Evolutionary Biology and Environmental Studies, University of Zürich, Winterthurerstrasse 190, 8057 Zürich, Switzerland <sup>2</sup>Department of Biostatistics, Epidemiology, Biostatistics and Prevention Institute, University of Zürich, Hirschengraben 84, 8001 Zürich, Switzerland <sup>3</sup>Zoological Museum, University of Zürich, Karl-Schmid-Strasse 4, 8006 Zürich, Switzerland

<sup>4</sup>Division of Ecology and Evolution, Research School of Biology, The Australian National University, Acton, Canberra, ACT 2601, Australia

<sup>5</sup>E-mail: stefanie.muff@uzh.ch

Received January 12, 2018 Accepted July 12, 2018

Quantitative genetic analyses require extensive measurements of phenotypic traits, a task that is often not trivial, especially in wild populations. On top of instrumental measurement error, some traits may undergo transient (i.e., nonpersistent) fluctuations that are biologically irrelevant for selection processes. These two sources of variability, which we denote here as measurement error in a broad sense, are possible causes for bias in the estimation of quantitative genetic parameters. We illustrate how in a continuous trait transient effects with a classical measurement error structure may bias estimates of heritability, selection gradients, and the predicted response to selection. We propose strategies to obtain unbiased estimates with the help of repeated measurements taken at an appropriate temporal scale. However, the fact that in quantitative genetic analyses repeated measurements are also used to isolate permanent environmental instead of transient effects requires that the information content of repeated measurements is carefully assessed. To this end, we propose to distinguish "short-term" from "long-term" repeats, where the former capture transient variability and the latter help isolate permanent effects. We show how the inclusion of the corresponding variance components in quantitative genetic models yields unbiased estimates of all quantities of interest, and we illustrate the application of the method to data from a Swiss snow vole population.

**KEY WORDS**: Animal model, Breeder's equation, error variance, permanent environmental effects, quantitative genetics, Robertson–Price identity.

Quantitative genetic methods have become increasingly popular for the study of natural populations in the last decades, and they now provide powerful tools to investigate the inheritance of characters, and to understand and predict evolutionary change of phenotypic traits (Falconer and Mackay 1996; Lynch and Walsh 1998; Charmantier et al. 2014). At its core, quantitative genetics is a statistical approach that decomposes the observed phenotype *P* into the sum of additive genetic effects *A* and a residual component *R*, so that P = A + R. For simplicity, nonadditive genetic effects, such as dominance and epistatic effects, are ignored throughout this article, thus the residual component can be thought of as the sum of all environmental effects. This basic model can be extended in various ways (Falconer and Mackay 1996; Lynch and Walsh 1998), with one of the most common

being P = A + PE + R, where *PE* captures *dependent* effects, the so-called *permanent environmental effects*, while *R* captures the residual, *independent* variance that remains unexplained. Permanent environmental effects are stable differences among individuals above and beyond the permanent differences due to additive genetic effects. In repeated measurements of an individual, these effects create within-individual covariation. To prevent inflated estimates of additive genetic variance, these effects must therefore be modeled and estimated (Lynch and Walsh 1998; Kruuk 2004; Wilson et al. 2010).

This quantitative genetic decomposition of phenotypes is not possible at the individual level in nonclonal organisms, but under the crucial assumption of independence of genetic, permanent environmental, and residual effects, the phenotypic variance at the population level can be decomposed into the respective variance components as  $\sigma_P^2 = \sigma_A^2 + \sigma_{PE}^2 + \sigma_R^2$ . These variance components can then be used to understand and predict evolutionary change of phenotypic traits. For example, the additive genetic variance ( $\sigma_A^2$ ) can be used to predict the response to selection using the Breeder's equation. It predicts the response to selection  $R_{\text{BE}}$  of a trait *z* (bold face notation denotes vectors) from the product of the heritability ( $h^2$ ) of the trait and the strength of selection (*S*) as

$$R_{\rm BE} = h^2 \cdot S \tag{1}$$

(Lush 1937; Falconer and Mackay 1996), where  $h^2$  is the proportion of additive genetic to total phenotypic variance

$$h^2 = \frac{\sigma_A^2}{\sigma_P^2},\tag{2}$$

and *S* is the selection differential, defined as the mean phenotypic difference between selected individuals and the population mean or, equivalently, the phenotypic covariance  $\sigma_p(z, w)$  between the trait (*z*) and relative fitness (*w*). Besides the Breeder's equation, evolution can be predicted using the secondary theorem of selection, according to which evolutionary change is equal to the additive genetic covariance of a trait with relative fitness, that is,

$$R_{\rm STS} = \sigma_a(\boldsymbol{z}, \boldsymbol{w}) \tag{3}$$

(Robertson 1966; Price 1970). Morrissey et al. (2010, 2012) discuss the differences between the Breeder's equation and the secondary theorem of selection in detail. A major difference is that in contrast to  $R_{\rm BE}$ ,  $R_{\rm STS}$  only estimates the population's evolutionary trajectory, but does not measure the role of selection in shaping this evolutionary change.

One measure of the role of selection is the selection gradient, which quantifies the strength of natural selection on a trait. For a normally distributed trait (z), it is given as the slope  $\beta_z$  of the linear regression of relative fitness on a phenotypic trait (Lande and Arnold 1983), that is,

$$\beta_z = \frac{\sigma_p(z, \boldsymbol{w})}{\sigma_p^2(z)},\tag{4}$$

where  $\sigma_p^2(z)$  denotes the phenotypic variance of the trait, for which we only write  $\sigma_p^2$  when there is no ambiguity about what trait the phenotypic variance refers to.

The reliable estimation of the parameters of interest  $(h^2,$  $\sigma_p(z, w), \sigma_a(z, w), \text{ and } \beta_z)$ , and the successful prediction of evolution as  $R_{BE}$  or  $R_{STS}$ , require large amounts of data, often collected across multiple generations and with known relationships among individuals in the dataset. For many phenotypic traits of interest, data collection is not trivial, and multiple sources of error, such as phenotypic measurement error, pedigree errors (wrong relationships among individuals), or nonrandomly missing data may affect the parameter estimates. Several studies have discussed and addressed pedigree errors (e.g., Keller et al. 2001; Griffith et al. 2002; Senneke et al. 2004; Charmantier and Reale 2005; Hadfield 2008) and problems arising from missing data (e.g., Steinsland et al. 2014; Wolak and Reid 2017). In contrast, although known for a long time (e.g., Price and Boag 1987), the effects of phenotypic measurement error on estimates of (co-)variance components have received less attention (but see, e.g., Hoffmann 2000; Dohm 2002; Macgregor et al. 2006; van der Sluis et al. 2010; Ge et al. 2017). In particular, general solutions to obtaining unbiased estimates of (co-)variance parameters in the presence of phenotypic measurement error are lacking.

In the simplest case, and the case considered here, phenotypic measurement error is assumed to be independent and additive, that is, instead of the actual phenotype z, an error-prone version

$$z^{\star} = z + e, \qquad e \sim \mathsf{N}(\mathbf{0}, \sigma_{e_m}^2 \mathbf{I})$$
 (5)

is measured, where *e* denotes an error term with independent correlation structure I and error variance  $\sigma_{e_m}^2$  (see Lynch and Walsh 1998, p. 121). As a consequence, the *observed* phenotypic variance of the measured values is  $\sigma_p^2(z^*) = \sigma_p^2(z) + \sigma_{e_m}^2$ , and thus larger than the *actual* phenotypic variance. The error variance  $\sigma_{e_m}^2$ thus must be disentangled from  $\sigma_p^2(z)$  to obtain unbiased estimates of quantitative genetic parameters. However, most existing methods for continuous trait analyses that acknowledged measurement error have modeled it as part of the residual component, and thus implicitly as part of the total phenotypic value (e.g., Dohm 2002; Macgregor et al. 2006; van der Sluis et al. 2010). This means that in the decomposition of a phenotype, P = A + PE + R, measurement error is absorbed in *R*, thus  $\sigma_{e_m}^2$  is absorbed by  $\sigma_R^2$ . This practice effectively *downwardly* biases measures that are proportions of the phenotypic variance, in particular  $h^2$  and  $\beta_z$ . To see why, let us denote the biased measures as  $h_{\star}^2$  and  $\beta_z^{\star}$ . The biased version of heritability is then given as

$$h_{\star}^{2} = \frac{\sigma_{A}^{2}}{\sigma_{P}^{2} + \sigma_{e_{m}}^{2}} \leq \frac{\sigma_{A}^{2}}{\sigma_{P}^{2}}$$
(6)

because under the assumption taken here that measurement error is independent of the actual trait value, measurement error is also independent of additive genetic differences and therefore leaves the estimate of the additive genetic variance  $\sigma_A^2$  unaffected. This was already pointed out, for example, by Lynch and Walsh (1998, p. 121) or Ge et al. (2017). Equation (6) directly illustrates that  $h_*^2$ is attenuated by a factor  $\lambda = \sigma_P^2 / (\sigma_P^2 + \sigma_{e_m}^2)$ , denoted as reliability ratio (e.g., Carroll et al. 2006). Using the same argument, one can show that  $\beta_z^* = \lambda \beta_z$ , but also  $R_{BE}^* = \lambda R_{BE}$ , as will become clear later.

To obtain unbiased estimates of  $h^2$ ,  $\beta_z$ , or any other quantity that depends on unbiased estimates of  $\sigma_P^2$ , it is thus necessary to disentangle  $\sigma_{e_m}^2$  from the actual phenotypic variance  $\sigma_P^2$ , and particularly from its residual component  $\sigma_R^2$ . Importantly, however, purely mechanistic measurement imprecision is often not the only source of variation that may be considered irrelevant for the mechanisms of inheritance and selection in the system under study. Here, we therefore follow Ge et al. (2017) and use the term "transient effects" for the sum of measurement errors *plus* any biological short-term changes of the phenotype itself that are not considered relevant for the selection process, briefly denoted as "irrelevant fluctuations" of the actual trait.

As an example, if the trait is the mass of an adult animal, repeated measurements within the same day are expected to differ even in the absence of instrumental error, simply because animals eat, drink, and defecate (for an example of the magnitude of these effects see Keller and Van Noordwijk 1993). Such short-term fluctuations might not be of interest for the study of evolutionary dynamics, if the fluctuations do not contribute to the selection process in a given population. Under the assumption that these fluctuations are additive and independent among each other and of the actual trait value, they are mathematically indistinguishable from pure measurement error. In the remainder of the article, we therefore do not introduce a separate notation to discriminate between (mechanistic) measurement error and biological shortterm fluctuations, but treat them as a single component (e) with a total "error" variance  $\sigma_{e_m}^2$ . Consequently, we may sometimes refer to "measurement error" when in fact we mean transient effects as the sum of measurement error and transient fluctuations.

The aim of this article is to develop general methods to obtain unbiased estimates of heritability, selection, and response to selection in the presence of measurement error and irrelevant fluctuations of a trait, building on the work by Ge et al. (2017). We start by discussing the meaning and information content of repeated phenotypic measurements on the same individual. The type of phenotypic trait we have in mind is a relatively plastic trait, such as milk production or an animal's mass, which are expected to undergo changes across an individual's life span that are relevant for selection. We point out that repeated measures taken over different time intervals can help separate transient effects from more stable (permanent) environmental and genetic effects, and that based on such a variance decomposition it is possible to formulate models that yield unbiased estimates of heritability, selection, and the response to selection. We illustrate these approaches with empirical quantitative trait analyses of body mass measurements taken in a population of snow voles in the Swiss alps (Bonnet et al. 2017).

# Theory short- and long-term repeated measurements

Table 1 gives an overview of how the different parameters considered here are (or are not) affected by the presence of measurement error. To retrieve unbiased estimates of all quantities given in Table 1, we must be able to appropriately model and estimate the measurement error variance  $\sigma_{e_m}^2$ , which can be achieved with repeated measurements. These repeated measurements must be taken in close temporal vicinity, that is, on a time scale where the focal trait is not actually undergoing any phenotypic changes that are considered relevant for selection. We introduce the notion of a measurement session for such short-term time intervals. In other words, a measurement session can be defined as a sufficiently short period of time during which the investigator is willing to assume that the residual component is constant. On the other hand, measurements are often repeated across much longer periods of time, such as months, seasons, or years, during which phenotypic change is not expected to be solely due to transient effects, and the resulting trait variation is often relevant for selection. Thus, long-term repeats, taken across different measurement sessions, help separating permanent environmental effects from residual components (e.g., Wilson et al. 2010).

The distinction between short- and long-term repeats, and thus the definition of a measurement session, may not always be obvious or unique for a given trait. In the introduction, we employed the example of an animal's mass that transiently fluctuates within a day. Depending on the context, such fluctuations might not be of interest, and the "actual" phenotypic value would correspond to the average daily mass. A reasonable measurement session could then be a single day, and within-day repeats can thus be used to estimate  $\sigma_{e_m}^2$ . If however *any* fluctuations in body mass are of interest, irrespective of how persistent they are, much shorter measurement sessions, such as seconds or minutes, would be appropriate to ensure that only the purely mechanistic measurement error variance is represented by  $\sigma_{e_m}^2$ .



**Figure 1.** Schematic representation of three study designs, where one individual is measured (A) multiple times across multiple measurement sessions, (B) multiple times in one single measurement measurement session, or (C) one single time across multiple measurement sessions. Only case (A) allows to disentangle the measurement error variance  $\sigma_{e_m}^2$  and the permanent environmental effects  $\sigma_{PE}^2$  from  $\sigma_R^2$ , whereas case (B) allows to separate only the measurement error variance and case (C) only allows to disentangle permanent environmental from residual effects.

## REPEATED MEASUREMENTS IN THE ANIMAL MODEL

In the following, we show how measurement error can be incorporated in the key tool of quantitative genetics, the *animal model*, a special type of (generalized) linear mixed model, which is commonly used to decompose the phenotypic variance of a trait into genetic and nongenetic components (Henderson 1976; Lynch and Walsh 1998; Kruuk 2004).

Let us assume that phenotypic measurements of a trait are blurred by measurement error following model (5), and that measurements have been taken both across and within multiple measurement sessions, as indicated in Fig. 1A. Denoting by  $z_{ijk}^{*}$  the k th measurement of individual i in session j, it is possible to fit a model that decomposes the trait value as

$$z_{ijk}^{\star} = \mu + \boldsymbol{x}_{ijk}^{\top} \boldsymbol{\beta} + a_i + id_i + R_{ij} + e_{ijk},$$
(7)

where  $\mu$  is the population intercept,  $\beta$  is a vector of fixed effects, and  $x_{ijk}$  is the vector of covariates for measurement k in session j of animal i. The remaining components are the random effects, namely the breeding value  $a_i$  with dependency structure  $(a_1, \ldots, a_n)^T \sim N(\mathbf{0}, \sigma_A^2 \mathbf{A})$ , an independent, animal-specific permanent environmental effect  $id_i \sim N(0, \sigma_{PE}^2)$ , an independent Gaussian residual term  $R_{ij} \sim N(0, \sigma_R^2)$ , and an independent error term  $e_{ijk} \sim N(0, \sigma_{e_m}^2)$  that absorbs any transient effects captured by the within-session repeats. The dependency structure of the breeding values  $a_i$  is encoded by the additive genetic relatedness matrix A (Lynch and Walsh 1998), which is traditionally derived from a pedigree, but can alternatively be calculated from genomic data (Meuwissen et al. 2001; Hill 2014). The model can be further expanded to include more fixed or random effects, such as maternal, nest, or time effects, but we omit such terms here without loss of generality. Importantly, model (7) does not require that all individuals have repeated measurements in each session to obtain an unbiased estimate of the variance components in the presence of measurement error. In fact, even if there are, on average, fewer than two repeated measurements per individual within sessions, it may be possible to separate the error variance from the residual variance, as long as the total number of within-session repeats over all individuals is reasonably large. We will in the following refer to model (7) as the "error-aware" model.

If, however, a trait has not been measured across different time scales (i.e., either only within or only across measurement sessions), not all variance components are estimable. In the first case, when repeats are only taken within a single measurement session for each individual, as depicted in Fig. 1B, an error term can be included in the model, but a permanent environmental effect cannot. The model must then be reduced to

$$z_{ik}^{\star} = \mu + \boldsymbol{x}_{ik}^{\top} \boldsymbol{\beta} + a_i + R_i + e_{ik}, \qquad (8)$$

thus it is possible to estimate the error variance  $\sigma_{e_m}^2$  and to obtain unbiased estimates of  $\sigma_A^2$  and  $h^2$ , whereas the residual variance  $\sigma_R^2$ then also contains the permanent environmental variance. In the second case, when repeated measurements are only available from across different measurement sessions, as illustrated in Fig. 1C, the error variance cannot be estimated. Instead, an animal-specific permanent environmental effect can be added to the model, which is then given as

$$z_{ij}^{\star} = \mu + \boldsymbol{x}_{ij}^{\top} \boldsymbol{\beta} + a_i + i d_i + R_{ij}$$
(9)

for the measurement in session j for individual i. Interestingly, this last model mirrors the types of repeats that motivated quantitative geneticists to isolate  $\sigma_{PE}^2$ , which may otherwise be confounded not only with  $\sigma_R^2$ , but also with  $\sigma_A^2$ . This occurs because the repeated measurements across sessions induce an increased within-animal correlation (i.e., a similarity) that may be absorbed by  $\sigma_A^2$  if not modeled appropriately (Kruuk and Hadfield 2007; Wilson et al. 2010).

### MEASUREMENT ERROR AND SELECTION

Selection occurs when a trait is correlated with fitness, such that variations in the trait values lead to predictable variations among the same individuals in fitness. The leading approach for measuring the strength of directional selection is the one developed by Lande and Arnold (1983), who proposed to estimate the selection gradient  $\beta_z$  as the slope of the regression of relative fitness  $\boldsymbol{w}$  on the phenotypic trait  $\boldsymbol{z}$ :

$$\boldsymbol{w} = \alpha + \beta_z \cdot \boldsymbol{z} + \boldsymbol{\epsilon}, \tag{10}$$

with intercept  $\alpha$  and residual error vector  $\epsilon$ . This model can be further extended to account for covariates, such as sex or age. If the phenotype z is measured with error (which may again encompass any irrelevant fluctuations), such that the observed value is  $z^* = z + e$  with error variance  $\sigma_{e_m}^2$  as in (5), the regression of w against  $z^*$  leads to an attenuated version of  $\beta_z$  (Fuller 1987; Mitchell-Olds and Shaw 1987; Carroll et al. 2006). Using that  $\hat{\beta}_z = \frac{\sigma_p(z,w)}{\sigma_p^2(z)}, \sigma_p^2(z^*) = \sigma_p^2(z) + \sigma_{e_m}^2$ , and the assumption that the error in  $z^*$  is independent of w, simple calculations show that the error-prone estimate of selection is

$$\hat{\boldsymbol{\beta}}_{z}^{\star} = rac{\sigma_{p}(\boldsymbol{z}^{\star}, \boldsymbol{w})}{\sigma_{p}^{2}(\boldsymbol{z}^{\star})} = rac{\sigma_{p}(\boldsymbol{z}, \boldsymbol{w})}{\sigma_{p}^{2}(\boldsymbol{z}) + \sigma_{e_{m}}^{2}} \leq \hat{\boldsymbol{\beta}}_{z}$$

Hence, the quantity that is estimated is  $\beta_z^{\star} = \lambda \beta_z$  with  $\lambda =$  $\sigma_p^2(z)/(\sigma_p^2(z) + \sigma_{e_m}^2)$ , thus  $\beta_z$  suffers from exactly the same bias as the estimate of heritability (see again Table 1). To obtain an unbiased estimate of selection it may thus often be necessary to account for the error by a suitable error model. Such erroraware model must rely on the same type of short-term repeated measurements as those used in (7) or (8), but with the additional complication that z is now a covariate in a regression model, and no longer the response. To estimate an unbiased version of  $\beta_{z}$ , we therefore rely on the interpretation as an error in variables problem for classical measurement error (Fuller 1987; Carroll et al. 2006). To this end, we propose to formulate a Bayesian hierarchical model, because this formulation, together with the possibility to include prior knowledge, provides a flexible way to model measurement error (Stephens and Dellaportas 1992; Richardson and Gilks 1993). To obtain an error-aware model that accounts for error in selection gradients, we need a three-level hierarchical model: the first level is the regression model for selection, and the second level is given by the error model of the observed covariate  $z^*$  given its true value z. Third, a so-called

*exposure model* for the unobserved (true) trait value is required to inform the model about the distribution of z, and it seems natural to employ the animal model (9) for this purpose. Again using the notation for an individual i measured in different sessions j and with repeats k within sessions, the formulation of the three-level hierarchical model is given as

$$w_{ij} = \alpha + \beta_z z_{ij} + \boldsymbol{x}_{ij}^{\top} \boldsymbol{\beta} + \boldsymbol{\epsilon}_{ij}, \, \boldsymbol{\epsilon}_{ij} \sim N\left(0, \, \sigma_{\boldsymbol{\epsilon}}^2\right)$$
  
Selection model (11a)

$$z_{ijk}^{\star} = z_{ij} + e_{ijk} , e_{ijk} \sim \mathsf{N}\left(0, \sigma_{e_m}^2\right)$$
  
Error model (11b)

$$z_{ij} = \mu + \boldsymbol{x}_{ij}^{\top} \boldsymbol{\gamma} + a_i + id_i + R_{ij}, R_{ij} \sim N\left(0, \sigma_R^2\right)$$
  
Exposure model (11c)

where  $w_{ij}$  is the measurement of relative fitness for individual *i*, usually taken only once per individual and having the same value for all measurement sessions j,  $\beta$  is a vector of fixed effects,  $x_{ii}$ is the vector of covariates for animal i in measurement session j,  $\beta_z$  is the selection gradient, and  $\alpha$  and  $\epsilon_{ij}$  are, respectively, the intercept and the independent residual term from the linear regression model. The classical independent measurement error term is given by  $e_{iik}$ , and the vector of fixed effects in the exposure model is now denoted as  $\gamma$  to discriminate it from  $\beta$  in the selection model. This formulation as a hierarchical model gives an unbiased estimate of the selection gradient  $\beta_z$ , because the lower levels of the model properly account for the error in z by explicitly modeling it. It might be helpful to see that the second and third levels are just a hierarchical representation of model (7). Model (11a-c) can be fitted in a Bayesian setup, see for instance Muff et al.(2015) for a description of the implementation using integrated nested Laplace approximations (INLA) (Rue et al. 2009) via its R interface R-INLA.

Note that the selection model (11a) is formulated here for directional selection. Although the explicit discussion of alternative selection mechanisms, such as stabilizing or disruptive selection, is beyond the scope of the present article, we note that error modeling for these cases is straightforward: The only change is that the linear selection model is replaced by the appropriate alternative, for example by including quadratic or any other kind of nonlinear terms (*e.g.* Fisher 1930; Lande and Arnold 1983). Moreover, (11a) can be replaced by any other regression model, for example by one that accounts for nonnormality of fitness (see e.g., Morrissey and Sakrejda 2013; Morrissey and Goudie 2016). Similarly, it is conceptually straightforward to replace the Gaussian error and

Parameter	Notation	Effect of ME	Biased parameter
Additive genetic variance	$\sigma_A^2$	Unbiased	_
Permanent environmental variance	$\sigma_{PE}^2$	Unbiased	-
Residual variance	$\sigma_R^2$	Biased	$\sigma_R^2 + \sigma_e^2$
Heritability	$h^2$	Biased	$\lambda h^2$
Selection gradient	$\beta_z$	Biased	$\lambda \beta_z$
Phenotypic covariance	$\sigma_p(\boldsymbol{z}, \boldsymbol{w}) = S$	Unbiased	-
Response to selection (STS)	$\sigma_a(\boldsymbol{z}, \boldsymbol{w}) = R_{STS}$	Unbiased	-
Response to selection (BE)	$R_{BE}$	Biased	$\lambda R_{BE}$
Evolvability	Ι	Unbiased	_

Table 1. Overview of the effects of measurement error and transient fluctuations (abbreviated as ME) on important quantitative genetic parameters.

The table indicates for each parameter whether it is biased or unbiased. For biased parameters the quantities that are estimated when ignoring transient effects in the quantitative genetic models are given.  $\lambda$  is the reliability ratio, defined as  $\lambda = \frac{\sigma_p^2}{\sigma_p^2 + \sigma_{em}^2}$ . For notation see the main text.

exposure models, if there is reason to believe that the normal assumptions for the error term  $e_{ijk}$  or the residual term  $R_{ij}$  are unrealistic, for example if z is a count or a binary variable. In fact, equation (10) to estimate selection does not actually assume a specific distribution for z, however, the interpretation of  $\beta_z$  as a directional selection gradient to predict evolutionary change may be lost for non-Gaussian traits (Lande and Arnold 1983). Finally and importantly, although multivariate selection is not covered in the present article, it is possible to extend the hierarchical model (11a–c) to the multivariate case.

# MEASUREMENT ERROR AND THE RESPONSE TO SELECTION

#### The Breeder's equation

Evolutionary response to a selection process on a phenotypic trait can be predicted either by the Breeder's equation (1) or by the Robertson-Price identity (3), and these two approaches are equivalent only when the respective trait value (in the univariate model) is the sole causal factor affecting fitness (Morrissey et al. 2010, 2012). Even if the Breeder's equation is formulated for multiple traits, the implicit assumption still is that all correlated traits causally related to fitness are included in the model. Given that fitness is a complex trait that usually depends on many unmeasured variables (Møller and Jennions 2002; Peek et al. 2003), it is not surprising that the Breeder's equation is often not successful in predicting evolutionary change in natural systems (Hadfield 2008; Morrissey et al. 2010), in contrast to (artificial) animal breeding situations, where, thanks to the control over the process, all the traits affecting fitness are known and included in the models (Lush 1937; Falconer and Mackay 1996; Roff 2007).

To understand how transient effects affect the estimate of  $R_{\rm BE} = h^2 \cdot S$ , we must understand how the components  $h^2$  and S are affected. We have seen that  $h_*^2 = \lambda h^2$ . On the other hand, the selection differential  $S^* = \sigma_p(z^*, w)$  is an unbiased estimate of  $\sigma_p(z, w)$ , because under the assumption of independence of the error vector e and fitness w,

$$\sigma_p(\boldsymbol{z}^{\star}, \boldsymbol{w}) = \sigma_p(\boldsymbol{z} + \boldsymbol{e}, \boldsymbol{w}) = \sigma_p(\boldsymbol{z}, \boldsymbol{w}) + \underbrace{\sigma_p(\boldsymbol{e}, \boldsymbol{w})}_{=0} = \sigma_p(\boldsymbol{z}, \boldsymbol{w}).$$
(12)

Consequently, the bias in  $h_{\star}^2$  directly propagates to the estimated response to selection, that is,  $R_{\rm BE}^{\star} = \lambda R_{\rm BE}$  (Table 1).

#### The Robertson–Price identity

Response to selection can also be predicted using the secondary theorem of selection. Specifically, the additive genetic covariance of the relative fitness  $\boldsymbol{w}$  and the phenotypic trait  $\boldsymbol{z}$ ,  $\sigma_a(\boldsymbol{w}, \boldsymbol{z})$ , can be estimated from a bivariate animal model. If interest centers around the evolutionary response of a single trait, the model for the response vector including the (error-prone) trait values  $\boldsymbol{z}^*$  and relative fitness values  $\boldsymbol{w}$  is bivariate with

$$\begin{bmatrix} z^* \\ w \end{bmatrix} = \mu + X\beta + \mathbf{D}a + \mathbf{Z}r, \qquad (13)$$

where  $\mu$  is the intercept vector,  $\beta$  the vector of fixed effects, X the corresponding design matrix, **D** is the design matrix for the breeding values a, and **Z** is a design matrix for additional random terms r. These may include environmental and/or error terms, depending on the structure of the data, that may correspond to the univariate cases of equations (7)–(9), or again to other random terms such as maternal or nest effects. The actual component

of interest is the vector of breeding values, which is assumed multivariate normally distributed with

$$\boldsymbol{a} = \begin{bmatrix} \boldsymbol{a}(z^{\star}) \\ \boldsymbol{a}(w) \end{bmatrix} \sim \mathsf{N} \left( \boldsymbol{0}, \begin{bmatrix} \sigma_a^2(z^{\star})\mathbf{A} & \sigma_a(\boldsymbol{w}, \boldsymbol{z}^{\star})\mathbf{A} \\ \sigma_a(\boldsymbol{w}, \boldsymbol{z}^{\star})\mathbf{A} & \sigma_a^2(\boldsymbol{w})\mathbf{A} \end{bmatrix} \right), (14)$$

where  $a(z^*)$  and a(w) are the respective subvectors for the trait and fitness, and **A** is the relationship matrix derived from the pedigree. An estimate of the additive genetic covariance  $\sigma_a(w, z^*)$ is extracted from this covariance matrix. An interesting feature of the additive genetic covariance, and consequently estimates of the response to selection using the STS, is that it is unbiased by independent error in the phenotype. This can be seen by reiterating the exact same argument as in equation (12), but replacing the phenotypic with the genetic covariance.

We confirmed all these theoretical expectations with a simulation study, where we analyzed the effects of measurement error on the estimates of interest by adding error terms with different variances to the phenotypic traits. Details and results of the simulations are given in Appendix S2, whereas the code for their implementation is reported in Appendix S3.

# Empirical Example data and models for body mass of snow voles

The empirical data we use here stem from a snow vole population that has been monitored between 2006 and 2014 in the Swiss Alps (Bonnet et al. 2017). The genetic pedigree is available for 937 voles, together with measurements on morphological and life-history traits. Thanks to the isolated location, it was possible to monitor the whole population and to obtain high recapture probabilities ( $0.924 \pm 0.012$  for adults and  $0.814 \pm 0.030$ for juveniles). Details of the study are given in Bonnet et al. (2017).

Our analyses focused on the estimation of quantitative genetic parameters for the animals' body mass (in grams). The dataset contained 3266 mass observations from 917 different voles across nine years. Such measurements are expected to suffer from classical measurement error, as they were taken with a spring scale, which is prone to measurement error under field conditions. In addition, the actual mass of an animal may contain irrelevant within-day fluctuations (eating, defecating, digestive processes), but also unknown pregnancy conditions in females, which cannot reliably be determined in the field. Repeated measurements were available, both recorded within and across different seasons. In each season two to five "trapping sessions" were conducted, which each lasted four consecutive nights. Although this definition of measurement session was based purely on operational aspects driven by the data collection process, we used this time interval to estimate  $\sigma_{e_m}^2$ . It is arguably possible that four days might be undesirably long, and that variability in such an interval includes more than purely transient effects, but the data did not allow for a finer time resolution. However, to illustrate the importance of the measurement session length, we also repeated all analyses with measurement sessions defined as a calendar month, which is expected to identify a larger (and probably too high) proportion of variance as  $\sigma_{e_m}^2$ . The number of four-day measurement sessions per individual was on average 3.02 (min = 1, max = 24) with 1.15 (min = 1, max = 3) number of short-term repeats on average, whereas there were 2.37 (min = 1, max = 13) one-month measurement sessions on average, with 1.41 (min = 1, max = 6) short-term repeats per measurement session.

## Heritability

Bonnet et al. (2017) estimated heritability using an animal model with sex, age, Julian date (JD), squared Julian date, and the twoand three-way interactions among sex, age, and Julian date as fixed effects. The inbreeding coefficient was included to avoid bias in the estimation of additive genetic variances (de Boer and Hoeschele 1993). The breeding value  $(a_i)$ , the maternal identity  $(m_i)$ , and the permanent environmental effect explained by the individual identity  $(id_i)$  were included as individual-specific random effects.

If no distinction is made between short-term (within measurement session) and long-term (across measurement sessions) repeated measurements, the model that we denote as the *naive* model is given as

$$z_{ijk}^{\star} = \mu + \boldsymbol{x}_{ijk}^{\top} \boldsymbol{\beta} + a_i + m_i + id_i + R_{ijk}, \qquad (15)$$

where  $z_{ijk}^{\star}$  is the mass of animal *i* in measurement session *j* for repeat *k*. This model is prone to underestimate heritability, because it does not separate the variance  $\sigma_{e_m}^2$  from the residual variability, and  $\sigma_{e_m}^2$  is thus treated as part of the total phenotypic trait variability. To isolate the measurement error variance, the model expansion

$$z_{ijk}^{\star} = \mu + \boldsymbol{x}_{ijk}^{\top} \boldsymbol{\beta} + a_i + m_i + id_i + R_{ij} + e_{ijk},$$

with  $R_{ij} \sim N(0, \sigma_R^2)$  and  $e_{ijk} \sim N(0, \sigma_{e_m}^2)$  leads to what we denote here as the *error-aware* model. Under the assumption that the length of a measurement session was defined in an appropriate way, and that the error obeys model (5), this model yields an unbiased estimate of  $h^2$ , calculated as  $\frac{\sigma_A^2}{\sigma_A^2 + \sigma_M^2 + \sigma_{PE}^2 + \sigma_R^2}$  (in agreement with Bonnet et al. 2017), where  $\sigma_{e_m}^2$  is explicitly estimated and thus not included in the denominator. Both models were implemented in MCMCglmm Hadfield (2010) and are reported in Appendix S4. Inverse gamma priors IG(0.01, 0.01), parameterized with shape and rate parameters, were used for all variances in all models,

Model	$\hat{h}^2$	$\hat{\sigma}_A^2$	$\hat{\sigma}_{PE}^2$	$\hat{\sigma}_M^2$	$\hat{\sigma}_R^2$	$\hat{\sigma}_{e_m}^2$
Naive	0.14	3.40	6.09	1.16	12.40	_
	[0.07, 0.25]	[1.41, 6.15]	[4.33, 8.51]	[0.56, 2.84]	[11.78, 13.21]	
Error-aware						
(four-day measurement session)	0.23	3.97	5.62	1.48	6.58	6.07
	[0.09, 0.33]	[1.46, 6.06]	[3.68, 7.68]	[0.57, 2.73]	[5.76, 7.82]	[5.54, 7.05]
Error-aware						
(one-month measurement session)	0.24	3.82	4.78	1.58	5.77	7.91
	[0.10, 0.37]	[1.17, 5.84]	[3.16, 7.21]	[0.61, 2.86]	[4.78, 6.71]	[7.15, 8.38]

Table 2. Estimates of quantitative genetic parameters of body mass in snow voles using naive and error-aware models.

The posterior modes of variance components and heritability are given, together with their 95% credible intervals (in brackets).

whereas N(0, 10<sup>12</sup>) (i.e., default MCMCglmm) priors were given to the fixed effect parameters. Analyses were repeated with varying priors on  $\sigma_{e_m}^2$  for a sensitivity check, but results were very robust (results not shown).

## Selection

Selection gradients were estimated from the regression of relative fitness (w) on body mass  $(z^*)$ . Relative fitness was defined as the relative lifetime reproductive success (rLRS), calculated as the number of offspring over the lifetime of an individual, divided by the population mean LRS. The naive estimate of the selection gradient was obtained from a linear mixed model (i.e., treating rLRS as continuous trait), where body mass, sex, and age were included as fixed effects, plus a cohort-specific random effect. The error-aware version of the selection gradient  $\beta_z$  was estimated using a three-layer hierarchical error model as in (11a-c) that also included a random effect for cohort in the regression model. Sex and age were also included as fixed effects in the exposure model, plus breeding values, permanent environmental and a residual term as random effects. The hierarchical model used to estimate the error-aware  $\beta_z$  was implemented in INLA and is described in Appendix S1, with R code given in Appendix S5. Again, IG(0.01, 0.01) priors were assigned to all variance components, whereas independent  $N(0, 10^2)$  priors were used for all slope parameters. Because rLRS is not actually a Gaussian trait, P values and CIs of the estimate for  $\beta_{z}$  from the linear regression model are, however, incorrect. Although recent considerations indicate that selection gradients could directly be extracted from an overdispersed Poisson model (Morrissey and Goudie 2016), we followed the original analysis of Bonnet et al. (2017) and extracted P values from an overdispersed Poisson regression model with absolute LRS as a count outcome, both for the (naive) model without error modeling and for the hierarchical error model, where the linear model (11c) was replaced by an overdispersed Poisson regression model (Appendices S1 and S5 include the model description and code for both models).

#### Response to selection

Response to selection on body mass was estimated with rLRS using the Breeder's equation (1) and the secondary theorem of selection (3), both for the naive and the error-aware versions of the model. The naive and error-aware versions of  $R_{\rm BE}$  were estimated by substituting either the naive  $h_{\star}^2$  or the error-aware estimates of  $h^2$  into the Breeder's equation, where the selection differential was calculated as the phenotypic covariance between mass and rLRS. On the other hand,  $R_{\rm STS}$  was estimated from the bivariate animal model, implemented in MCMCglmm using the same fixed and random effects as those in equation (15). Again IG(0.01, 0.01) priors were used for the variance components. No residual component was included for the fitness trait, as suggested by Morrissey et al. (2012), and its error variance was fixed at 0, because no error modeling is required. Appendix S6 contains the respective R code.

## SNOW VOLES RESULTS

#### Heritability

As expected from theory (Table 1), transient effects in the measurements of body mass biased some, but not all, quantitative genetic estimates in our snow vole example (Table 2). The estimates and confidence intervals of the additive genetic variance  $\sigma_A^2$ , as well as the permanent environmental variance  $\sigma_{PE}^2$  and the maternal variance (denoted as  $\sigma_M^2$ ) were only slightly corrected in the error-aware models. Residual variances, however, were much lower when measurement error was accounted for in the models. The error-aware model separated residual and transient (error) variance so that  $\hat{\sigma}_R^2 + \hat{\sigma}_{e_m}^2$  corresponded approximately to  $\hat{\sigma}_R^2$  from the naive model. The overestimation of the residual variance by nearly 40% when measurement error was ignored ( $\hat{h}^2 = 0.14$  in the naive model and  $\hat{h}^2 = 0.23$  in the error-aware model).

As expected, the estimated measurement error variance was larger when a measurement session is defined as a full month  $(\hat{\sigma}_{e_m}^2 = 7.91)$  than as a 4-day interval  $(\hat{\sigma}_{e_m}^2 = 6.07;$  Table 2),

Model	$\hat{\boldsymbol{\beta}}_z$	<i>P</i> -value
Naive	0.065	< 0.001
Error-aware (four-day	0.104	< 0.001
measurement session)		
Error-aware (one-month	0.104	< 0.001
measurement session)		

**Table 3.** Estimates of selection gradients  $(\hat{\beta}_z)$  for body mass in snow voles derived from naive (ML estimate) and error-aware models (posterior means).

For all models, Bayesian *P* values were derived from zero-inflated Poisson regressions.

because the trait then has more time and opportunity to change. As a consequence, heritability is even slightly higher ( $\hat{h}^2 = 0.24$ ) when the longer measurement session definition is used. This example is instructive because it underlines the importance of defining the time scale at which short-term repeats are expected to capture only transient, and not biologically relevant variability of the phenotypic trait. In the case of the mass of a snow vole, most biologists would probably agree that changes in body mass over a one-month measurement session may well be biologically meaningful and relevant (i.e., body fat accumulation, pregnancy in females, etc.), while it is less clear how much of the fluctuations within a 4-day measurement session are transient, and what part of it would be relevant for selection. Within-day repeats might be the most appropriate for the case of mass, because within-day variance is likely mostly transient, but because the data were not collected with the intention to quantify such effects, within-day repeats were not available in sufficient numbers in our example dataset.

# Selection

As expected, estimates of selection gradients ( $\hat{\beta}_z$ ) obtained with the naive models provided nearly 40% lower estimates of selection than the error-aware model (Table 3). The two measurement session lengths yielded similar results. With and without measurement error modeling, the *P* values of the zero-inflated Poisson models confirmed the presence of selection on body mass in snow voles (*P* < 0.001 in all models).

## Response to selection

In line with theory, estimates of the response to selection using the Breeder's equation were nearly 40% lower when transient effects were not incorporated in the quantitative genetic models using four-day measurement sessions ( $\hat{R}_{BE} = 0.10$  in the naive model and  $\hat{R}_{BE} = 0.16$  in the error-aware model; Table 4). As in the case of heritability, the one-month measurement session definition resulted in even slightly higher estimates of the response to selection

 $(\hat{R}_{BE} = 0.17)$ . In contrast, response to selection measured by the secondary theorem of selection  $\hat{R}_{STS}$  did not show evidence of bias, and the error-aware model with a four-day measurement session definition estimated the same value ( $\hat{R}_{STS} = -0.17$ ) as the naive model (Table 4). With a one-month measurement session, we obtained a slightly attenuated value ( $\hat{R}_{STS} = -0.14$ ), although the difference was small in comparison to the credible intervals (Table 4).

This example illustrates that the Breeder's equation is generally prone to underestimation of the selection response in real study systems when measurement error in the phenotype is present (Table 1). The results also confirm that estimates for response to selection may differ dramatically between the Breeder's equation and the secondary theorem of selection. As already noticed by Bonnet et al. (2017), the predicted evolutionary response derived from the Breeder's equation points in the opposite direction in the snow vole data than the estimate derived from the secondary theorem of selection (e.g., naive estimates  $\hat{R}_{\text{BE}} =$ 0.10 vs.  $\hat{R}_{\text{STS}} = -0.17$ , with nonoverlapping credible intervals; Table 4).

# Discussion

This study addressed the problem of measurement error and transient fluctuations in continuous phenotypic traits in quantitative genetic analyses. We have shown that measurement error and transient fluctuations can lead to substantial bias in estimates of several important quantitative genetic parameters, including heritability, selection gradients, and the response to selection (Table 1). We introduced modeling strategies to obtain unbiased estimates in these parameters in the presence of measurement error and transient fluctuations. These strategies rely on the distinction between variability from stable effects that are part of the biologically relevant phenotypic variability, and transient effects, which are the sum of mechanistic measurement error and biological fluctuations that are considered irrelevant for the selection process. We argue that ignoring the distinction between stable and transient effects may not only lead to an underestimation of the heritability due to inflated estimates of the residual variance,  $\sigma_R^2$ , but also to bias in the estimates of selection gradients and the response to selection. Measurements of the same individual repeated at appropriate time scales allow the variance from such transient effects to be partitioned, and thus prevent such bias.

How can repeated measurements be used to prevent an *un*derestimation of heritability, selection, and response to selection, while permanent environment effects are required in quantitative genetic models of repeated measures to avoid an upward bias of  $\sigma_A^2$  and, hence, an *over*estimation of  $h^2$  (Wilson et al. 2010)? The fact that repeated measurements are used to prevent opposite

Model	$\hat{R}_{ m STS}$	95% CI	$\hat{R}_{ ext{BE}}$	95% CI
Naive	-0.17	[-0.54, 0.18]	0.10	[0.05, 0.17]
Error-aware (four-day measurement session)	-0.17	[-0.51, 0.19]	0.16	[0.06, 0.23]
Error-aware (one-month measurement session)	-0.14	[-0.53, 0.17]	0.17	[0.07, 0.26]

**Table 4.** Response to selection for body mass in snow voles (posterior modes and 95% credible intervals) estimated with the Breeder's equation ( $\hat{R}_{BE}$ ) and with the secondary theorem of selection ( $\hat{R}_{STS}$ ).

Results are shown for the naive and the error-aware models.

biases in heritability estimates makes it apparent that the information content in what is termed "repeated measurements" in both cases is very different. The crucial aspect is that it matters at which temporal distance the repeats were taken, and that the relevance of this distance depends on the kind of trait under study. Repeats taken on the same individual at different life stages ("long-term" repeats, e.g., across what we call measurement sessions here) can be used to separate the animal-specific permanent environmental effect from both genetic and residual variances. On the other hand, repeats taken in temporal vicinity ("short-term" repeats, e.g., within a measurement session) help disentangle any transient from the residual effects. Only by modeling both types of repeats, that is, across different relevant time scales, is it practically feasible to separate all variance components. To do so, the quantitative genetic model for the trait value, typically the animal model, needs extension to three levels of measurement hierarchy (eq. 7): the individual (i), the measurement session (j within i), and the repeat (k within j within i). As highlighted with the snow vole example, it may not always be trivial to determine, in a particular system, an appropriate distinction between short- and long-term repeats, and consequently how to define a measurement session. This decision must be driven by the definition of short-term variation as a variation that is not "seen" by the selection process (see, e.g., Price and Boag 1987, p. 279 for a similar analogy) in contrast to persistent effects that are potentially under selection. This distinction ultimately depends on the trait, on the system under study and on the research question that is asked, because some traits may fluctuate on extremely short time scales (minutes or days), whereas others remain constant across an entire adult's life.

The application to the snow vole data, where we varied the measurement session length from four days to one month, illustrated that longer measurement sessions automatically capture more variability, that is, the estimated error variance  $\hat{\sigma}_{e_m}^2$  increased. Consequently, unreasonably long measurement sessions may lead to overcorrected estimates of the parameters of interest. On the other hand, considering measurement sessions that are too short may lead to an insufficient number of within-session repeats, or they may fail to identify transient variability that is biologically irrelevant. This makes clear that a careful definition of measurement.

ment session length is important already at the design stage of a study.

If one is uncertain whether repeated measurements capture effects relevant to selection or not, would averaging over repeats result in better estimates of quantitative genetic measures? Averaging methods have been proposed specifically to reduce bias that emerges due to measurement error and transient effects (Carbonaro et al. 2009; Zheng et al. 2016). Although averaging will alleviate bias by reducing the error variance in the mean, it will not eliminate it completely. This can be seen from the fact that averaging over *K* within-session repeats for all animals and measurement sessions, the variance  $\sigma_{e_m}^2$  is reduced to  $\sigma_{\bar{e}_m}^2 = \sigma_{e_m}^2/K$ , assuming independence of the error term. Unless *K* is large,  $\sigma_{e_m}^2$  will not approach zero. Moreover, this practice only works if all animals have the same number of repeats within all measurement sessions, but it will not work in the unbalanced sampling design so common in studies of natural populations.

We approached the problem of measurement error and transient fluctuations by assuming a dichotomous distinction between short- and long-term repeats. This requires determining the temporal scale of measurement sessions, something that is likely highly trait and context specific. An alternative perspective of within-animal repeated measurements could take a continuous view, recalling that repeated measurements are usually correlated, even when taken across long time spans, and that the correlation increases the closer in time the measurements were taken. A more sophisticated model could thus take into account that the residual component in the model changes continuously, and introduce a time-dependent correlation structure instead of simply distinguishing between short- and long-term repeats. Such a model might be beneficial if repeats were not taken in clearly defined measurement sessions, although such a temporal correlation term introduces another level of model complexity, and thus entails other challenges.

It may sometimes not be possible to take multiple measurements on the same individual, or to repeat a measurement within a session. However, it may still be feasible to include an appropriate random effect in the absence of short-term repeats, provided that knowledge about the error variance is available, for example, from previous studies that used the same measurement devices, from a subset of the data, or from other "expert" knowledge. The Bayesian framework is ideal in this regard, because it is straightforward to include random effects with a very strong (or even fixed) prior on the respective variance component. Such Bayesian models provide error-aware estimates that are equivalent to those illustrated in Table 1, but with the additional advantage that posterior distributions naturally reflect all uncertainty that is present in the parameters, including the uncertainty that is incorporated in the prior distribution of the error variance.

Measurement error and transient fluctuations bias some, but not all quantitative genetic inferences. When  $\sigma_{e_m}^2 > 0$ , the naive estimates of  $h^2$ ,  $\beta_z$ , and  $R_{\rm BE}$  are attenuated by the same factor  $\lambda < 1$ , but other components, such as the selection differential *S* or  $R_{\text{STS}}$ , are not affected (Table 1). The robustness of the secondary theorem of selection to measurement error can certainly be seen as an advantage over the Breeder's equation. Nevertheless, the Robertson-Price identity does not model selection explicitly, and thus says little about the selective processes. The Robertson-Price equation can be used to check the consistency of predictions made from the Breeder's equation, but the Breeder's equation remains necessary to test hypothesis about the causal nature of selection (Morrissey et al. 2012; Bonnet et al. 2017). Another quantity that is unaffected by independent transient effects, which we however did not further elaborate on here, is evolvability, defined as the squared coefficient of variation  $I = \sigma_4^2 / \overline{z}^2$ , where  $\overline{z}$  denotes the mean phenotypic value (Houle 1992). Evolvability is often used as an alternative to heritability, and is interpreted as the opportunity for selection (Crow 1958). Not only  $\sigma_A^2$ , but also  $\overline{z}$  can be consistently estimated using  $z^*$ , namely because the expected values  $E[z^*] = E[z]$  due to the independence and zero mean of the error term. For completeness, we added evolvability to Table 1.

A critical assumption of our models was that the error components are independent of the phenotypic trait under study, but also independent of fitness or any covariates in the animal model or the selection model. Although the small changes in  $\hat{R}_{STS}$  that we observed in the snow vole application with one-month measurement sessions could be due to pure estimation stochasticity, an alternative interpretation is that the measurement error in the data are not independent of the animal's fitness. At least two processes could lead to a correlation between the measurement error in mass and fitness in snow voles. First, pregnant females will experience temporally increased body mass, and we expect the positive deviation from the true body mass to be correlated with fitness, because a pregnant animal is likely to have a higher expected number of offspring over its entire lifespan. And second, some of the snow voles were not fully grown when measured, and juveniles are more likely to survive if they keep growing, so that deviations from mean mass over the measurement session period would be nonrandomly associated with life-time fitness.

So far, we have focused on traits that can change relatively quickly throughout the life of an individual, such as body mass, or physiological and behavioral traits. Traits that remain constant after a certain age facilitate the isolation of measurement error, because the residual variance term is then indistinguishable from the error term, given that a permanent environmental (i.e., individualspecific) effect is included in the model. In such a situation it is sufficient to estimate  $\sigma_R^2$ , which then automatically corresponds to the measurement error variance, whereas  $\sigma_{PE}^2$  captures all the environmental variability. However, not many traits will fit that description. The majority of traits, even seemingly stable traits such as skeletal traits, are in fact variable over time (Price and Grant 1984; Smith et al. 1986).

We have shown that dealing appropriately with measurement error and transient fluctuations of phenotypic traits in quantitative genetic analyses requires the inclusion of additional variance components. Quantitative genetic analyses often differ in the variance components that are included to account for important dependencies in the data (Meffert et al. 2002; Kruuk and Hadfield 2007; Palucci et al. 2007; Hadfield et al. 2013). Besides the importance of separating the right variance components, it has been widely discussed which of the components are to be included in the denominator of heritability estimates, although the focus has been mainly on the proper handling of variances that are captured by the fixed effects (Wilson 2008; de Villemereuil et al. 2018). We hope that our treatment of measurement error in quantitative genetic analyses sparks new discussions of what should be included in the denominator when heritability is calculated.

The methods presented in this article have been developed and implemented for continuous phenotypic traits. Binary, categorical or count traits may also suffer from measurement error, which is then denoted as misclassification error (Copas 1988; Magder and Hughes 1997; Küchenhoff et al. 2006), or as miscounting error (e.g., Muff et al. 2018). Models for non-Gaussian traits are usually formulated in a generalized linear model framework (Nakagawa and Schielzeth 2010; de Villemereuil et al. 2016) and require the use of a link function (e.g., the logistic or log link). In these cases, it will often not be possible to obtain unbiased estimates of quantitative genetic parameters by adding an error term to the linear predictor as we have done here for continuous traits. Obtaining unbiased estimates of quantitative genetic parameters in the presence of misclassification and miscounting error will require extended modeling strategies, such as hierarchical models with an explicit level for the error process.

We hope that the concepts and methods provided here serve as a useful starting point when estimating quantitative genetics parameters in the presence of measurement error or transient, irrelevant fluctuations in phenotypic traits. The proposed approaches are relatively straightforward to implement, but further generalizations are possible and will hopefully follow in the future.

#### **AUTHOR CONTRIBUTIONS**

SM and LFK conceived the research idea. EP designed and conducted the simulations and analyses. TB collected and compiled the vole data. EP and SM wrote the manuscript, with inputs from LFK and TB. All authors gave final approval.

#### ACKNOWLEDGMENTS

This work would not have been possible without generous funding from the Faculty of Science and the Forschungskredit of the University of Zurich, grant no. FK-16-097. The snow vole monitoring was supported by the Claraz-Donation, and by the Swiss National Science Foundation project grants 31003*A*\_141110 and 31003*A*\_159462/1 to E. Postma. We thank E. Postma, G. Camenisch, and P. Wandeler for data collection and conceptual input. The thoughtful inputs of two anonymous reviewers and the associate editor, M. Morrissey, greatly improved the manuscript. The authors declare they have no competing interests.

#### **DATA ARCHIVING**

Data has been fully archived in https://doi.org/10.1371/journal.pbio. 1002592

#### LITERATURE CITED

- Bonnet, T., P. Wandeler, G. Camenisch, and E. Postma. 2017. Bigger is fitter? Quantitative genetic decomposition of selection reveals an adaptive evolution decline of body mass in a wild rodent population. PLOS Biol. 15:e1002592.
- Carbonaro, F., T. Andrew, D. A. Mackey, T. L. Young, T. D. Spector, and C. J. Hammond. 2009. Repeated measures of intraocular pressure result in higher heritability and greater power in genetic linkage studies. Invest. Ophthalmol. Vis. Sci. 50:5115–5119.
- Carroll, R. J., D. Ruppert, L. A. Stefanski, and C. M. Crainiceanu. 2006. Measurement error in nonlinear models, a modern perspective. Chapman and Hall, Boca Raton, FL.
- Charmantier, A., and D. Reale. 2005. How do misassigned paternities affect the estimation of heritability in the wild? Mol. Ecol. 14:2839–2850.
- Charmantier, A., D. Garant, and L. E. B. Kruuk. 2014. Quantitative genetics in the wild. Oxford Univ. Press, Oxford, U.K.
- Copas, J. B. 1988. Binary regression models for contaminated data (with discussion). J. R. Stat. Soc. B (Stat. Methodol.) 50:225–265.
- Crow, J. F. 1958. Some possibilities for measuring selection intensities in man. Hum. Biol. 30:1–13.
- de Boer, I. J. M., and I. Hoeschele. 1993. Genetic evaluation methods for populations with dominance and inbreeding. Theor. Appl. Genet. 86:245– 258.
- de Villemereuil, P., H. Schielzeth, S. Nakagawa, and M. B. Morrissey. 2016. General methods for evolutionary quantitative genetic inference from generalized mixed models. Genetics 204:1281–1294.
- de Villemereuil, P., M. B. Morrissey, S. Nakagawa, and H. Schielzeth. 2018. Fixed effect variance and the estimation of the heritability: issues and solutions. J. Evol. Biol. 31:621–632.
- Dohm, M. R. 2002. Repeatability estimates do not always set an upper limit to heritability. Funct. Ecol. 16:273–280.
- Falconer, D. S. and T. F. C. Mackay. 1996. Introduction to quantitative genetics. Pearson, Essex, England.
- Fisher, R. A. 1930. The genetical theory of natural selection. Oxford Univ. Press, Oxford, U.K.
- Fuller, W. A. 1987. Measurement error models. John Wiley & Sons, New York.

- Ge, T., A. J. Holmes, R. L. Buckner, J. W. Smoller, and M. Sabuncu. 2017. Heritability analysis with repeat measurements and its application to resting-state functional connectivity. PNAS 114:5521–5526.
- Griffith, S. C., I. P. F. Owens, and K. A. Thuman. 2002. Extrapair paternity in birds: a review of interspecific variation and adaptive function. Mol. Ecol. 11:2195–2212.
- Hadfield, J. D. 2008. Estimating evolutionary parameters when viability selection is operating. Proc. R. Soc. Lond. B Biol. Sci. 275:723–734.
- 2010. MCMC methods for multi-response greneralized linear mixed models: The MCMCglmm R package. J. Stat. Softw. 33(2):1–22.
- Hadfield, J. D., E. A. Heap, F. Bayer, E. A. Mittell, and N. M. Crouch. 2013. Disentangling genetic and prenatal sources of familial resemblance across ontogeny in a wild passerine. Evolution 67:2701–2713.
- Henderson, C. R. 1976. A simple method for computing the inverse of a numerator relationship matrix used in prediction of breeding values. Biometrics 32:69–83.
- Hill, W. G. 2014. Applications of population genetics to animal breeding, from Wright, Fisher and Lush to genomic prediction. Genetics 196:1–16.
- Hoffmann, A. A. 2000. Laboratory and field heritabilities: lessons from *Drosophila*. In T. Mousseau, B. Sinervo, and J. Endler, eds.), Adaptive genetic variation in the wild. Oxford Univ. Press, New York, Oxford.
- Houle, D. 1992. Comparing evolvability and variability of quantitative traits. Genetics 130:195–204.
- Keller, L. F., and A. J. Van Noordwijk. 1993. A method to isolate environmental effects on nestling growth, illustrated with examples from the Great Tit (Parsus major). Funct. Ecol. 7:493–502.
- Keller, L. F., P. R. Grant, B. R. Grant, and K. Petren. 2001. Heritability of morphological traits in Darwin's Finches: misidentified paternity and maternal effects. Heredity 87:325–336.
- Kruuk, L. E. B. 2004. Estimating genetic parameters in natural populations using the 'animal model'. Philos. Trans. R. Soc. B Biol. Sci. 359:873– 890.
- Kruuk, L. E. B., and J. D. Hadfield. 2007. How to separate genetic and environmental causes of similarity between relatives. J. Evol. Biol. 20:1890– 1903.
- Küchenhoff, H., S. M. Mwalili, and E. Lesaffre. 2006. A general method for dealing with misclassification in regression: the misclassification SIMEX. Biometrics 62:85–96.
- Lande, R., and S. J. Arnold. 1983. The measurement of selection on correlated characters. Evolution 37:1210–1226.

Lush, J. L. 1937. Animal breeding plans. Iowa State College Press, Ames, IA.

- Lynch, M., and B. Walsh. 1998. Genetics and analysis of quantitative traits. Sinauer Associates, Sunderland, MA.
- Macgregor, S., B. K. Cornes, N. G. Martin, and P. M. Visscher. 2006. Bias, precision and heritability of self-reported and clinically measured height in Australian twins. Hum. Genet. 120:571–580.
- Magder, L. S., and J. P. Hughes. 1997. Logistic regression when the outcome is measured with uncertainty. Am. J. Epidemiol. 146:195–203.
- Meffert, L. M., S. K. Hicks, and J. L. Regan. 2002. Nonadditive genetic effects in animal behavior. Am. Nat. 160(Suppl 6):S198–S213.
- Meuwissen, T. H. E., B. J. Hayes, and M. E. Goddard. 2001. Prediction of total genetic value using genome-wide dense marker maps. Genetics 157:1819–1829.
- Mitchell-Olds, T., and R. G. Shaw. 1987. Regression analysis of natural selection: statistical inference and biological interpretation. Evolution 41:1149–1161.
- Møller, A., and M. D. Jennions. 2002. How much variance can be explained by ecologists and evolutionary biologists? Oecologia 132(4):492– 500.
- Morrissey, M. B., and I. B. J. Goudie. 2016. Analytical results for directional and quadratic selection gradients for log-linear models

of fitness functions. *bioRxiv*. https://www.biorxiv.org/content/early/2016/02/22/040618.

- Morrissey, M. B., and K. Sakrejda. 2013. Unification of regression-based methods for the analysis of natural selection. Evolution 67:2094–2100.
- Morrissey, M. B., L. E. B. Kruuk, and A. J. Wilson. 2010. The danger of applying the Breeder's equation in observational studies of natural populations. J. Evol. Biol. 23:2277–2288.
- Morrissey, M. B., D. J. Parker, P. Korsten, J. M. Pemberton, L. E. B. Kruuk, and A. J. Wilson. 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.
- Muff, S., A. Riebler, L. Held, H. Rue, and P. Saner. 2015. Bayesian analysis of measurement error models using integrated nested Laplace approximations. J. R. Stat. Soc. Appl. Stat. C 64:231–252.
- Muff, S., M. A. Puhan, and L. Held. 2018. Bias away from the Null due to miscounted outcomes? A case study on the TORCH trial. Stat. Methods Med. Res.: *In press*.
- Nakagawa, S., and H. Schielzeth. 2010. Repeatability for Gaussian and non-Gaussian data: a practical guide for biologists. Biol. Rev. Camb. Philos. Soc. 85:935–956.
- Palucci, V., L. R. Schaeffer, F. Miglior, and V. Osborne. 2007. Non-additive genetic effects for fertility traits in Canadian Holstein cattle. Genet. Sel. Evol. 39:181–193.
- Peek, M. S., A. J. Leffler, S. D. Flint, and R. J. Ryel. 2003. How much variance is explained by ecologists? Additional perspectives. Oecologia 137:161–170.
- Price, G. R. 1970. Selection and covariance. Nature 227:520-521.
- Price, T. D., and P. T. Boag. 1987. Selection in natural populations of birds. Pp. 257–287 in F. Cooke, and P. Buckley, eds., Avian genetics. Academic Press, Cambridge, MA.
- Price, T. D., and P. R. Grant. 1984. Life history traits and natural selection for small body size in a population of Darwin's Finches. Evolution 38:483–494.
- Richardson, S., and W. R. Gilks. 1993. Conditional independence models for epidemiological studies with covariate measurement error. Stat. Med. 12:1703–1722.
- Robertson, A. 1966. A mathematical model of the culling process in dairy cattle. Anim. Sci. 8:95–108.

- Roff, D. A. 2007. A centennial celebration for quantitative genetics. Evolution 61:1017–1032.
- Rue, H., S. Martino, and N. Chopin. 2009. Approximate Bayesian inference for latent Gaussian models by using integrated nested Laplace approximations (with discussion). J. R. Stat. Soc. B (Stat. Methodol.) 71: 319–392.
- Senneke, S. L., M. D. MacNeil, and L. D. Van Vleck. 2004. Effects of sire misidentification on estimates of genetic parameters for birth and weaning weights in Hereford cattle. J. Anim. Sci. 82:2307–2312.
- Smith, J. N. M., P. Arcese, and D. Schulter. 1986. Song sparrows grow and shrink with age. AUK 103:210–212.
- Steinsland, I., C. T. Larsen, A. Roulin, and H. Jensen. 2014. Quantitative genetic modeling and inference in the presence of nonignorable missing data. Evolution 68:1735–1747.
- Stephens, D. A., and P. Dellaportas. 1992. Bayesian analysis of generalised linear models with covariate measurement error. In J. M. Bernardo, J. O. Berger, A. P. Dawid, and A. F. M. Smith, eds. Bayesian statistics. Vol. 4. Oxford Univ. Press, Oxford, U.K.
- van der Sluis, S., M. Verhage, D. Posthuma, and C. V. Dolan. 2010. Phenotypic complexity, measurement bias, and poor phenotypic resolution contribute to the missing heritability problem in genetic association studies. PLOS One 5:e13929.
- Wilson, A. J. 2008. Why  $h^2$  does not always equal VA/VP? J. Evol. Biol. 21:647–650.
- Wilson, A. J., D. Réale, M. N. Clements, M. B. Morrissey, E. Postma, C. A. Walling, L. E. B. Kruuk, and D. H. Nussey. 2010. An ecologist's guide to the animal model. J. Anim. Ecol. 79:13–26.
- Wolak, M. E., and J. M. Reid. 2017. Accounting for genetic differences among unknown parents in microevolutionary studies: how to include genetic groups in quantitative genetic animal models. J. Anim. Ecol. 86:7– 20.
- Zheng, Y., R. Plomin, and S. von Stumm. 2016. Heritability of intraindividual mean and variability of positive and negative affect: genetic analysis of daily affect ratings over a month. Psychol. Sci. 27:1611–1619.

# Associate Editor: M. Morrissey Handling Editor: P. Tiffin

# Supporting Information

Additional supporting information may be found online in the Supporting Information section at the end of the article.

Appendix 1. Supplementary text and figures (pdf).

- Appendix 2. Supplementary text and figures for simulation study (pdf).
- Appendix 3. R script for the simulation and analysis of pedigree data.
- Appendix 4. R script for heritability in snow voles.

Appendix 5. R script for selection in snow voles.

Appendix 6. R script for response to selection in snow voles.